

Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES

TECHNICAL FRUITS OF NUCLEAR RESEARCH

by W. J. OOSTERKAMP *).

621.039

The article below is based on the address delivered on 26th January 1962 by Dr. W. J. Oosterkamp upon his inauguration as Professor Extraordinary at the Technische Hogeschool, Eindhoven. Some of the lantern slides shown on that occasion are reproduced here together with a number of additional illustrations. Bibliographical references have also been added.

The extensive nuclear research undertaken in numerous laboratories during the past thirty years has produced many results that were useful for practical applications. For example:

(1) It has led to a clear insight into the operation of particle accelerators and enabled considerable experience to be gained in their construction. These devices are capable of imparting large amounts of kinetic energy to electrons and other particles without requiring correspondingly high voltages. This is turned to advantage in the design of devices for generating X-rays of great penetrating power, which are used in the non-destructive examination of thick objects and for the irradiation of malignant tumours. High-energy electrons are also employed for the latter purpose.

(2) It has made it possible to produce on an industrial scale numerous artificial radioactive substances, which find application in chemistry, biology, medicine, agriculture and industry.

(3) It has opened up a new source of energy to help meet the world's growing power requirements.

We shall deal briefly with these three points in turn, but will touch first on the hazards involved in nuclear engineering.

Hazards of nuclear engineering

The hazards are twofold: the radioactive radiation and the risk of explosions. On the latter point I can be brief. The hazard of explosions exists in many technical installations. Electric power stations, for example, nowadays use steam pressures of 175 atm

and temperatures of almost 600 °C; and hydrogen gas at 50 atm and 500 °C is commonplace in the production of petrol. To me, as a layman in these fields, such pressures and temperatures seem most alarming; and yet, as far as I know, their use causes no catastrophes.

I would like to go into more detail about the radiation hazard. When energetic charged particles, such as alpha and beta or other electron rays, pass through matter they ionize atoms and molecules in their path or bring them into an excited state, thereby successively losing energy until they come to rest. Depending on its initial energy, one particle can give rise in this way to tens of thousands of ions or excited atoms. Rays of non-charged particles, such as neutrons and X-ray or gamma photons, can cause in their interaction with matter the liberation of energetic charged particles, as for instance electrons or recoil nuclei, which produce similar ionizations in their turn. Ionization and excitation may bring about a chemical change in the molecule struck, and if that molecule forms part of a living cell, may produce a biological effect such as to damage or destroy the cell. Exposure to large doses of radiation can result in serious bodily injury, e.g. the X-ray burns that were a fairly frequent occurrence in the period when the use of X-rays was pioneered. Recovery from somatic damage of this nature is possible to a considerable extent, the recovery being more complete the lower the dose of radiation received. Nowadays the maximum permissible doses — which do not of course apply to therapeutic irradiations — are much lower than those that cause appreciable somatic damage.

*) Philips Research Laboratories, Eindhoven.

Experiments on animals have shown that ionizing radiation received in the gonads can cause genetic mutations which are cumulative and permanent. The animal itself notices nothing, but the consequences may appear in later generations, and it is assumed they will invariably be deleterious. Investigations by Sobels at Leyden ¹⁾ may well point to a means of effectively protecting the transmitters of heredity in man. It is not certain whether any threshold value exists for these hereditary changes, i.e. a dose below which no effect is produced. Many experts believe that a threshold value is unlikely and that even the smallest doses can cause a hereditary effect, although with correspondingly less probability.

A human being can be exposed not only to radiation from outside but also to radiation from radioactive substances inside his own body, which are either present naturally, or have been inhaled or taken in with food and drink, or which have entered in some other way. The natural radioactive potassium isotope K-40 and the radioactive carbon isotope C-14 alone are responsible for producing in a human body weighing 70 kg about 7000 radioactive disintegrations per second; in other words, every second 7000 projectiles are fired, each of which gives rise on an average to about 10 000 ionizations. This has now been going on from generation to generation for the last half a million years, without however causing the destruction of the human race. Added to this is a roughly equal dose due to cosmic radiation, and a dose twice as large from the radioactivity in the immediate environment.

Recommendations concerning the maximum permissible radiation doses and the maximum permissible concentrations of radioactive substances in the air, in the drinking water and in the human body have been agreed on by an international commission ²⁾. These recommendations, so far as they relate to a genetic effect, are based partly on experiments with insects and animals, in particular fruit flies and mice, and partly on the level of natural radiation and on the fact that this can vary fairly considerably from one place to another. In most places at sea level the natural dose is about 100 millirads per year ³⁾, which implies an energy ab-

sorption in the human body of 30 picowatts per kilogram. There are also fairly densely populated areas where the natural dose is substantially higher: in Bolivia the dose is 30 per cent higher owing to the altitude of the country (3000-4000 m); in a part of France numbering about 7 million inhabitants the dose is twice as high, because of the radioactive constituents of the granite rocks; and in a part of the state of Kerala in India, which has 100 000 inhabitants, the local monazite sand results in a dose ten times as high. In the latter region the World Health Organization has started investigating the hereditary characteristics of the population ⁴⁾.

In connection with the expected increase in the uses of radioactive substances and nuclear energy, and the consequent higher incidence of radioactive waste products in various countries, an investigation was instituted some ten years ago into the exposure of man to ionizing radiation caused by man-made sources. The unexpected conclusion was that medical diagnostic procedures using X-rays were the greatest source of exposure ⁵⁾. In some countries the genetically important dose due to X-ray examinations is roughly of the same magnitude as the natural dose, whereas all other artificial sources together, including the radioactive fall-out from atom-bomb tests, give rise to a dose which is a mere few per cent of that due to natural radiation. When this fact was established, the radiologists — in the Netherlands on the initiative of the Netherlands Radiological Society and of the Board of Health — and the designers of medical X-ray equipment considered ways and means of improving the situation. The result is that authoritative members of the medical profession now believe that the medical dose can be considerably reduced without detracting in any way from the value of the diagnostic examination.

As regards protection against the radiation hazard, it can be said in general that the present methods and instruments of measurement are adequate for determining whether the radiation dose and the concentrations of radioactive material remain within the prescribed safety limits. Moreover it is considered to be technically feasible to meet the recommended standards. A necessary condition, however, is that those working with radiation should exercise discipline. Lower dose levels than those internationally recommended as the maximum permissible are also to be achieved. Their

¹⁾ Personal communication from Professor F. M. Sobels, Laboratory of Radiation Genetics, Leyden University.

²⁾ Recommendations of the International Commission on Radiological Protection (adopted September 9, 1958), Pergamon Press, London 1959.
Recommendations of the International Commission on Radiological Protection; Report of Committee II: Permissible Dose for Internal Radiation, Pergamon Press, London 1959.

³⁾ The definitions of the rad and other dose units will be found in: J. Hesselink and K. Reinsma, Dosimeters for X-radiation, Philips tech. Rev. 23, 55-66, 1961/62 (No. 2).

⁴⁾ Personal communication.

⁵⁾ Exposure of man to ionizing radiation arising from medical procedures, Report of the ICRP and ICRU, Physics in Medicine and Biology 2, 107-151, 1957.

realization entails sharply increasing costs, and this raises a serious problem of policy. It is necessary to weigh the benefits of these technical advances and their costs, including those involved in protection, against the risks. In doing so one should make a comparison with other risks to which man is exposed, as for instance road traffic, which, in the Netherlands, claims about 2000 fatalities every year.

As far as the present situation is concerned, the health and casualty statistics of industries and research centres in the field of nuclear energy compare favourably with those of other industries.

I shall now turn to a review of some of the above-mentioned fruits of nuclear research. First the accelerators.

Particle accelerators for industrial and medical use

We shall discuss here only electron accelerators, which are used among other things for generating X-rays of great penetrating power. X-rays are as a rule produced by causing fast electrons to strike a target. The higher the energy of the electrons the more penetrating are the X-rays generated. This is important in the medical irradiation of deep-seated tumours. For that purpose an X-ray tube operated at 800 kV had already been installed in the Antoni van Leeuwenhoekhuis (cancer hospital) at Amsterdam before 1940⁶⁾. To prevent electrical breakdown, this tube was an immovable assembly of complicated glass structures, insulators and electrodes. The tube was energized by an open high-tension generator. After 1945 the application of microwave techniques, which had been highly developed during the war for radar purposes, led to the construction of a much more elegant device: the linear accelerator⁷⁾. In this device it is possible, without requiring very high voltages, to accelerate the electrons to extremely great energies, which are in principle unlimited if only the accelerator can be made long enough and supplied with sufficient high-frequency power. In the United States (Stanford University) a project exists for an accelerator nearly 2 miles long, with a final electron energy of 20 000 million electron-volts (ultimately to be raised to 45 000 MeV). For the applications to be discussed, much lower energies are sufficient, i.e. about 5 MeV.

The first linear accelerators built in Great Brit-

⁶⁾ J. H. van der Tuuk, A million volt X-ray tube, *Philips tech. Rev.* **4**, 153-161, 1939.

⁷⁾ D. W. Fry, The linear electron accelerator, *Philips tech. Rev.* **14**, 1-12, 1952/53.

C. F. Bareford and M. G. Kelliher, The 15 million electron-volt linear electron accelerator for Harwell, *Philips tech. Rev.* **15**, 1-26, 1953/54.

ain were used as neutron sources for nuclear experiments⁸⁾. If the target is encased in beryllium, the absorption of X-rays in beryllium nuclei gives rise to two α particles plus one neutron per nucleus. For higher electron energies a uranium target is used, and neutrons are produced by photonuclear reaction in the target itself.

A 4 MeV linear accelerator for making radiographs of very thick metal work-pieces, e.g. 10-30 cm steel, can be a compact apparatus and thus easily handled. An accelerator of this type is at present being used for examining the welds of the reactor pressure vessels of the 500 MW nuclear power station under construction at Trawsfynydd in North Wales⁹⁾ (*fig. 1*). The efficiency of a power station equipped



Fig. 1. 4 MeV industrial radiography installation, equipped with a linear accelerator (Mullard), in operation for examining the welds of one of the pressure vessels under construction for the nuclear power station at Trawsfynydd, North Wales⁹⁾. (Photograph by courtesy of Atomic Power Constructions Ltd.)

⁸⁾ J. D. Cockcroft, High-energy electron accelerators as pulsed neutron sources, *Nature* **163**, 869, 1949.

⁹⁾ R. F. Hanstock, Electron accelerators for site radiography, *Nuclear Power*, February 1961.

The accelerator used for the investigation represented in *fig. 1* is described by T. R. Chippendale in: A 4 MeV industrial radiography installation, *Philips tech. Rev.* **23**, 197-215, 1961/62 (No. 7).

with gas-cooled nuclear reactors increases appreciably with the pressure of the cooling gas. A limit is set to the pressure by the wall thickness of the steel pressure vessel, which has to be welded on site. Advances in welding technique having made it possible to weld thicker steel plates, the radiographic unit had to be adapted accordingly. The radiation from normal X-ray apparatus, operated at 300 kV — in some cases 400 kV — was found to be no longer penetrating enough, so that for the latest pressure vessels, which have a wall thickness of almost 12 cm, use has been made of the gamma radiation from radioactive cobalt-60 and also, for about a year now, of the 4 MeV linear accelerator, the radiation intensity from which is roughly 30 times greater than that from a cobalt source of one kilocurie. There are indications that neither is the best method of examining the welds of such large wall thicknesses but that better results can be achieved with ultrasonic radiation. This doubt as to the usefulness of the linear accelerator does not, however, apply to the examination of large castings.

In Great Britain in particular, wide use is made of the linear accelerator for X-ray therapy (*fig. 2*). At the 1961 annual British Radiological Congress the linear accelerator was called the "work-horse of the large treatment centre"¹⁰). A comparison between a 4 MeV linear accelerator and the above-mentioned 800 kV glass X-ray tube reveals the following progress: greater compactness, a five-fold increase in electron energy, dose rate 25 times higher, greater mobility and no exposed high tension.

If the energy of the electrons used for generating X-rays exceeds 5 MeV, the X-rays do not become correspondingly more penetrating. With electron rays, however, this is indeed the case. The depth of their penetration is roughly proportional to the energy, about 0.3 cm/MeV. In radiation therapy the dose distribution in the irradiated body may sometimes be more favourable with electrons than with X-rays. For this reason electron irradiation has gained ground. To reach deep-seated tumours however, this calls for electron energies of 20 MeV and even higher. A linear accelerator for such energies would be unmanageably long (the length is roughly speaking 2 metres per 10 MeV).

Other constructions then enter into consideration, as for instance accelerators in which the electrons describe circular paths, such as *betatrons*. The electrons are made to revolve in a circular orbit by means of a magnetic field normal to their velocity. Acceleration using a travelling wave is then not

¹⁰) C. W. Miller, Brit. J. Radiology 35, 182, 1962.

readily possible, but the electrons can be accelerated by the electric field which is induced by a *change* in the magnetic field in which the electrons move. The energy gain per revolution is low, being of the order of magnitude of 10 eV, so that to acquire an energy of some tens of MeV the electron must complete a few million orbits, thereby covering a distance of some thousands of kilometres.

A betatron can be fed with alternating current of 50 c/s; this, among other things, makes it very simple in construction. To function properly, however, it must satisfy three conditions, which it took about twenty years to recognize.

The first ideas of the cyclical acceleration of electrons by means of magnetic induction were described by Slepian in 1922 in an American patent. As far as is known, he never put them into effect. At about the same time Wideröe, who was studying electrical engineering at the Karlsruhe Institute of Technology, was thinking along similar lines. His ideas led to an extensive experimental investigation, on which he based a thesis earning him his Doctor's degree at Aachen¹¹). He derived the condition which the magnetic field has to fulfil in order for the electrons to continue describing the same orbit. His experiments, however, yielded no positive result. He had more success with other experiments described in his thesis, which may be regarded as the precursors of the linear accelerator, and with which energies of about 2 MeV were achieved.

In about 1935 Steenbeck, working in the Siemens laboratories, derived the criterion for the stability of the electron orbit against small disturbances¹²). This enabled him to obtain favourable experimental results, but the electron beam produced was very weak. His ideas too went no further than the patent in which they were described. In about 1940 the stability criterion was also derived by Kerst at the University of Illinois. His calculations further showed that the maximum current obtainable is proportional to the initial energy of the electrons and that a high initial energy is therefore required¹³).

The complete foundation had now been laid for the betatrons of today. Wideröe was not yet finished with the betatron idea. He took up the thread again in 1944, and the result was a variable-energy betatron, for energies of 5 to 35 MeV, adapted to generating both

¹¹) R. Wideröe, Über ein neues Prinzip zur Herstellung hoher Spannungen, Arch. Elektrotechn. 21, 387-406, 1928. — For the principles of the betatron see A. Bierman and H. A. Oele, Betatrons with and without iron yoke, Philips tech. Rev. 11, 65-78, 1949/50.

¹²) M. Steenbeck, Beschleunigung von Elektronen durch elektrische Wirbelfelder, Naturwiss. 31, 234-235, 1943.

¹³) D. W. Kerst and R. Serber, Electronic orbits in the induction accelerator, Phys. Rev. 60, 53-58, 1941.

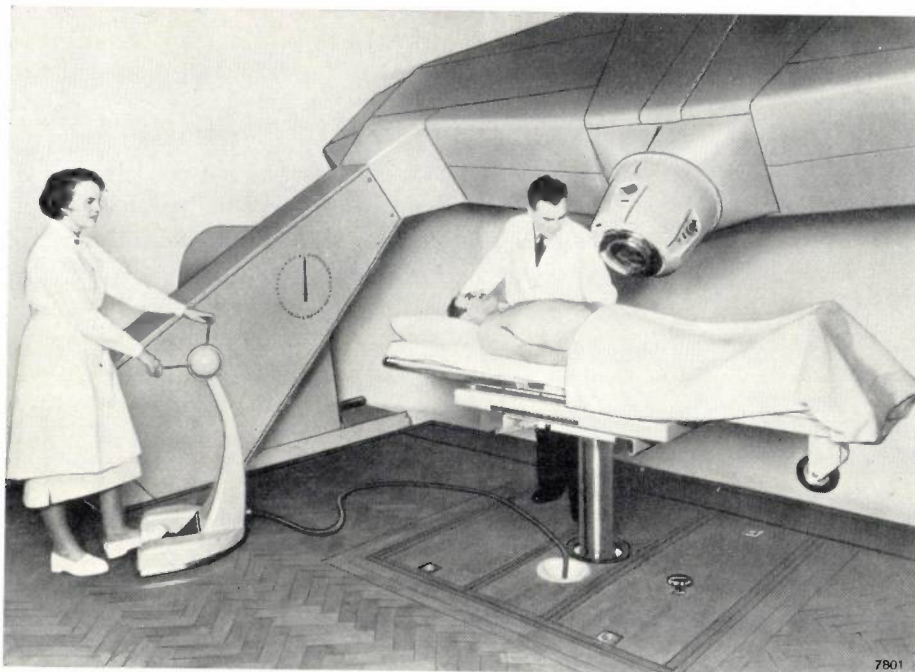


Fig. 2. X-ray therapy installation, equipped with a 4 MeV linear accelerator (Mullard).

X-rays and fast electrons ¹⁴⁾ (fig. 3). In 1942 Gund and others took up Steenbeck's work where he had left off. This led to a highly compact and mobile apparatus, for energies up to 18 MeV. This too can generate both X-rays and electron beams ¹⁵⁾ (fig. 4).

¹⁴⁾ R. Wideröe. Das Betatron, Z. angew. Physik 5, 187-200, 1953.

¹⁵⁾ K. Gund and H. Berger, Die 15-MeV-Elektronenschleuder für medizinische Anwendung der Siemens-Reiniger-Werke, Strahlentherapie 92, 489-505, 1953.

Owing to the great length of time the electrons remain in the accelerating tube, the concentration of electric charge is very high, and partly for that reason the average current in a betatron is relatively low, being less than one per cent of that in a linear accelerator.

Despite the fact that the linear accelerator and the betatron are already in fairly wide use in the medical world, it is still debatable in how far therapy

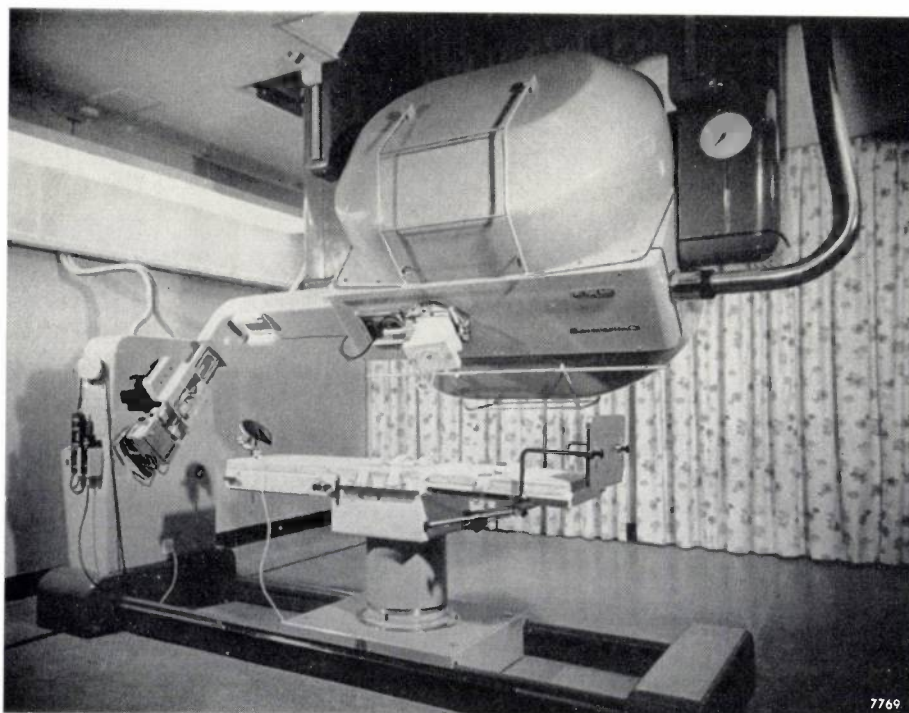


Fig. 3. Betatron for radiotherapy with X-rays or electrons of 5 to 35 MeV. (Photograph by courtesy of Brown-Boveri & Co.)

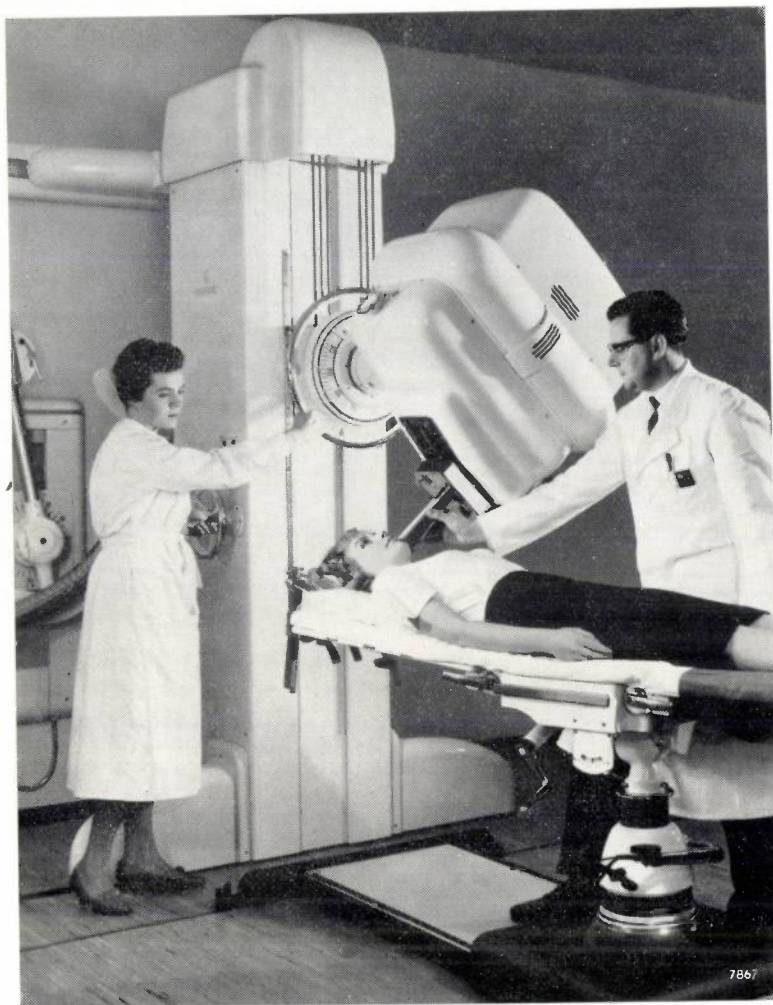


Fig. 4. Betatron for therapy with X-rays or electrons up to 18 MeV. (Photograph by courtesy of Siemens-Reiniger-Werke.)

using electron rays of extremely high energy is really an advance on the conventional methods of radiation therapy.

Applications of radioactive isotopes

Artificial radioactive substances are made by the irradiation of a stable isotope with neutrons in a nuclear reactor, by chemical separation from the fission products of a reactor, and by bombardment of certain elements with light or heavy hydrogen ions

Fig. 5. Investigation of sand movements due to sea currents, with the aid of a radio-isotope. The radioactive material (scandium-46, contained in zeolite, a granular material resembling North-Sea sand, and mixed with a quantity of sand) is supplied in plastic sacks and emptied into a metal container. The photograph shows a sack being opened with a knife on the end of a long stick. An empty sack is seen in the foreground. The closed container is let down on to the sea bed and there discharged. This method has meanwhile been replaced by a less hazardous one in which the radio-isotope is placed in the container by the manufacturer. (Photograph by courtesy of the Delta Authority of the Dutch National Water Board.)

accelerated in a cyclotron¹⁶). The fact that these substances emit radioactive rays is put to good use in several ways.

One of their applications is as "tracers". Radioactive substances owing to their radiation reveal their location and concentration, while chemically they behave in exactly the same way as the stable isotopes of the same element. By substituting a radioactive isotope for a certain atom in molecules involved in a chemical or biological process, one can follow the behaviour of these molecules during that process¹⁷.

Radioactive tracers are also employed in hydraulic engineering. They are playing a useful role for example in the Netherlands Delta project, in which it is very important to know the changes caused in the

¹⁶) A. H. W. Aten Jr. and J. Halberstadt, The production of radio-isotopes, Philips tech. Rev. 16, 1-12, 1954/55.

¹⁷) For a general survey see A. H. W. Aten Jr. and F. A. Heyn, The use of isotopes as tracers; The technique of investigations with radioactive and stable isotopes, Philips tech. Rev. 8, 296-303 and 330-336, 1946.



transport of sand by sea currents when the hydrographic situation around the coast is radically altered by damming and other interventions. In this investigation a sand-like substance containing a known amount of radioactive scandium-46 is dumped at a particular point on the sea bed (*fig. 5*). The radioactivity of the sea bed is then measured over a certain period in the area around the dumping site, making it possible to chart the movements of sand during that period¹⁸).

In medicine radioactive iodine-131 is frequently used for investigating the function of the thyroid gland. When iodine is orally administered to healthy people, in the form for example of sodium iodide, after absorption in the small intestine an average of 35% settles in the thyroid gland. The take-up of too much or too little iodine indicates a morbid condition. After administering radioactive iodine, the take-up can be ascertained by measuring the radioactive radiation from the thyroid gland at specific times with a detector placed above the area of interest (*fig. 6*). The sensitive part of the detector is surrounded with a lead shield having a frontal opening designed to ensure that only the gamma radiation from the thyroid gland is picked up. Using a radiation limiter with a very small aperture, it is possible to scan the area of the thyroid point by point, thus obtaining a two-dimensional activity diagram¹⁹). From this it can be seen whether



Fig. 6. Investigating the activity of the thyroid gland by means of a radioactive tracer, iodine-131, orally administered to the patient some time beforehand. The radiation pick-up (a scintillation detector) above the patient is shielded in such a way that only the gamma rays from the thyroid are measured. (Photograph by courtesy of St. Annadal Hospital, Maastricht.)

¹⁸) J. J. Arlman, J. N. Svašek and B. Verkerk, The use of radioactive isotopes for the study of littoral drift, Philips tech. Rev. 21, 157-166, 1959/60.

¹⁹) This scanning method, applied to various organs and using refinements in observation technique, is described in five articles published in Amer. J. Roentgenology 87, 128-170, 1962 (No. 1).

the shape of the functional part of the thyroid is normal and deviations from the normal shape, if any, can be recognized (*fig. 7a, b*).

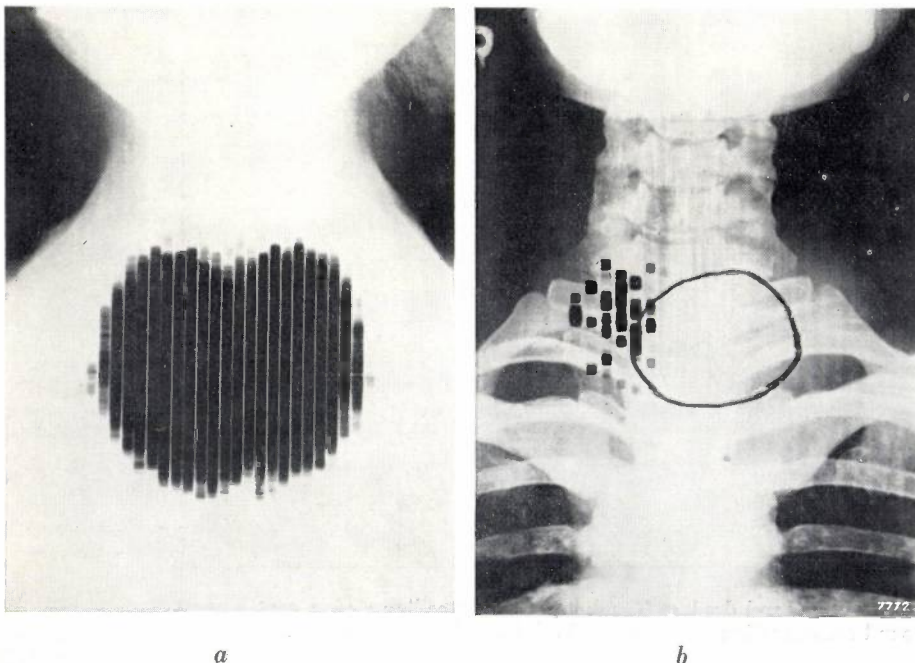


Fig. 7. By fitting the scintillation detector, in the arrangement shown in *fig. 6*, with a collimator having a very small aperture, the thyroid area can be scanned point by point and a 2-dimensional "scintigram" obtained of the thyroid gland. From this it can be seen whether the active part of the thyroid gland (the part that takes up iodine) is normal in shape, and if not, what the abnormalities are. The scintigrams shown here of a healthy (*a*) and a diseased subject (*b*) are each projected on to a radiograph of the same person. In case (*b*) the thyroid is active only at the periphery, top left; a tumour has made the other (circled) part inactive. (Photographs by courtesy of Dr. Coenegracht, St. Annadal Hospital, Maastricht.)

Radioactive substances can also be used for measuring the thickness of materials, the concentration of mixtures, and the height of liquid levels in closed vessels. As an illustration, I shall describe a sand-concentration meter which is at present used on about ten suction dredgers in the Netherlands. If the sand content in the mixture of sand and

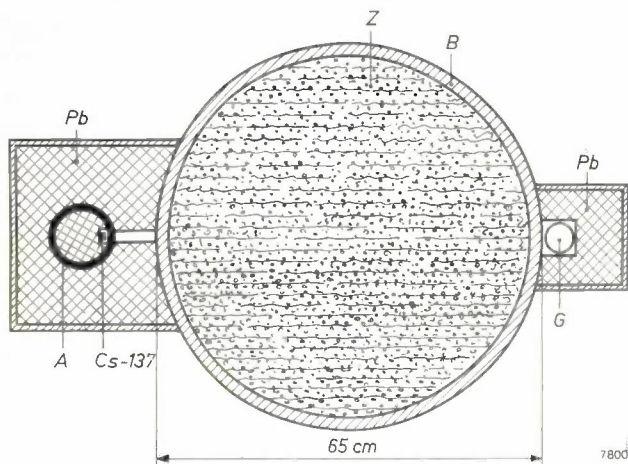


Fig. 8. Use of a radioactive isotope (caesium-137) for continuously measuring the sand concentration in a mixture of sand and water (*Z*) extracted through the outlet pipe of a sand dredger. Attached to one side of the pipe *B* is a capsule containing the radioactive material. The capsule is located in an opening on the periphery of a lead cylinder *A*, which can be turned about its axis inside a lead shield *Pb*. Mounted at the other side of the pipe *B* is a radiation detector *G*, also enclosed in a lead shield *Pb* to make the radiation passing through the sand and water mixture harmless. The intensity of the radiation that reaches *G* depends on the concentration of the sand. The radiation source can be "switched off" by slightly rotating the cylinder *A*.

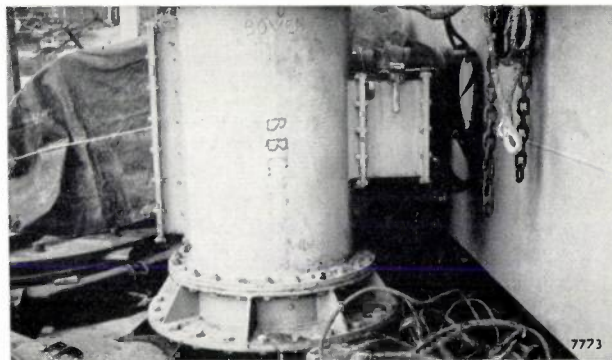


Fig. 9. Vertical part of the outlet pipe of the dredge to which, as shown in fig. 8, a large holder with the radioactive source and a smaller holder with the radiation detector are fitted.

water extracted through the dredger pipe is too low, then a needlessly large amount of water is being pumped. If the sand content is too high, the drag is excessive and there is a risk of a stoppage. To measure the concentration, a capsule containing caesium-137 is fitted to one side of a vertical part of the outlet pipe (*figs 8 and 9*). The capsule is enclosed in a lead shield, which absorbs the radiation emitted in undesired directions. The lead is safeguarded against mechanical damage by an iron housing. The sensitive part of a radiation detector, equally well protected, is mounted diametrically opposite. Water does not absorb the gamma rays emitted by the caesium-137 to the same extent as sand. The more sand the mixture contains, the smaller is the intensity of the radiation reaching the detector and hence the lower the reading on the meter. The meter is located in the control cabin of



Fig. 10. Control cabin of the sand dredger. The meter near the ceiling gives a continuous indication of the sand concentration in the pumped mixture of sand and water.



Fig. 11. Application of the gamma rays from a radio-isotope (iridium-192) for radiographing the welds in a pipeline. The holder with the radiation source is introduced into the pipe on a kind of trolley and pushed to the site of the weld, where an X-ray film is laid around the pipe. The source is operated and the film exposed by remote control. (Photograph by courtesy of "Röntgen-Technische Dienst", Rotterdam.)

the dredger (*fig. 10*), so that the operator can take appropriate action as soon as a reading shows any deviation.

In principle the gamma rays from radioactive isotopes could replace X-rays for industrial radiography. For the same size of emissive area, however, the radiation intensity from an X-ray unit may be as much as a hundred times higher, making it possible to use more economical exposure times. For this application, therefore, the use of radioactive substances is only justified when the radiation source has to be set up at a place which is not readily accessible (if at all) for an X-ray tube or where no electricity is available, or where the radia-

tion from ordinary X-ray units is not penetrating enough and radioactive substances are available which emit sufficiently hard gamma rays.

A good example is the use of radioactive iridium-192 in the examination of welds in pipelines. The holder containing the radioactive source is placed either on the outside of the pipe, diametrically opposite the film (so that the radiation passes through two parts of the wall successively), or in the centre of the pipe so that a radiograph can be taken of a complete weld on a film laid around the pipe. The photograph in *fig. 11* shows a source-holder being introduced into a pipe on a kind of trolley. After the film has been placed in position, the radiation source can be remotely uncovered to expose the film. This method was employed, for example, in the laying of the oil pipeline from Rotterdam to Cologne.

For the same reasons as explained in connection with the X-ray examination of materials, in clinical radiography radio-

active isotopes can replace X-ray units only in exceptional cases, particularly since in medical work it is often necessary to radiograph moving organs, which calls for very short exposure times.

In X-ray therapy the size of the emissive area is of less importance, and for this reason a gamma-ray source, in particular cobalt-60, is sometimes used instead of the conventional 200-300 kV X-ray units (*fig. 12*). The radiation from cobalt-60 is about as penetrating as the X-radiation from a 4 MeV linear accelerator. The choice between them is generally decided by secondary factors.

In analytical chemistry increasing use is being made of neutron activation analysis. The sensitivity

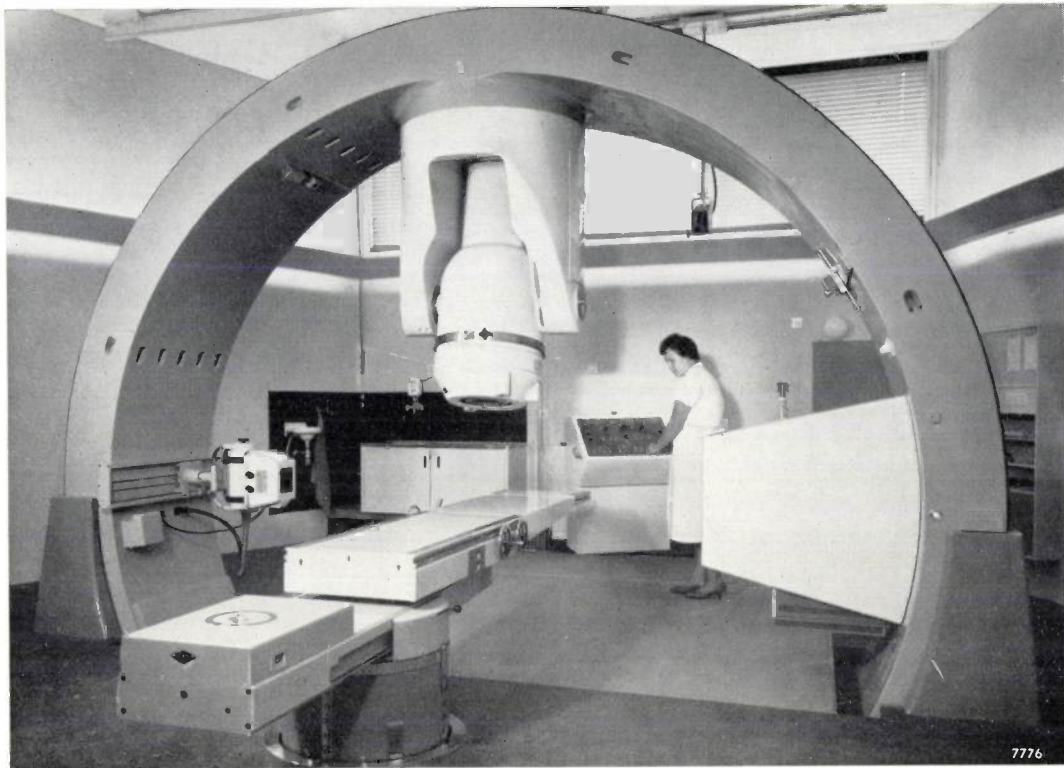


Fig. 12. Radiotherapy installation using cobalt-60 as radiation source. The source itself consists of a number of gilded discs of metallic cobalt, each with an activity of about 200 curies. Although the cobalt discs together weigh no more than a few dozen grammes, about 1500 kg of lead are needed to shield off the radiation in undesired directions. In the photograph the "cobalt radiator" is seen in the middle, suspended from an enormous ring stand which continues under the floor, and which enables the patient to be irradiated from any direction. On the left can be seen an X-ray diagnostic unit, and on the right, diametrically opposite, an X-ray image intensifier with television camera, making it possible to collimate the rays and direct them accurately onto the part of the patient to be treated. (Photograph by courtesy of Smit Röntgen N.V., Leyden.)

of this method is in some cases greater than that of "wet" or spectrochemical analysis. The substance to be analysed is bombarded by neutrons, either by leaving it for some time in a nuclear reactor or by exposing it to the radiation from a transportable neutron source²⁰). Depending on the composition, radioactive isotopes are thereby formed from various elements present in the sample. Each radioactive isotope is characterized by its half life and by the nature and energy of the radiation emitted. Analysis of the radiation emitted by the bombarded sample (*fig. 13*) reveals the presence of certain elements and also, in most cases, their concentration. Sometimes the activation analysis has to be preceded by chemical separation.

²⁰) O. Reifenschweiler and K. Nienhuis, The neutron tube, a simple and compact neutron source, Philips tech. Rev. 23, 325-337, 1961/62 (No. 11).

Use of nuclear energy

There are two ways of releasing the energy locked up in the atomic nucleus: by the fission of very heavy nuclei, such as uranium-235, or by the fusion of very light nuclei. The sun has been radiating energy to the earth for some six thousand million years now; this solar energy is entirely due to nuclear fusion, and the same applies to its conserved form in the fossilized fuels coal, oil and so on. The amount of solar energy reaching the earth daily is about 100 000 times greater than the daily consumption of fuels throughout the world. Even if the earth's population were doubled and the energy consumption per head of population were equal to twice the present consumption in the Netherlands, solar radiation would still be 4000 times greater than our requirements. Are we good stewards of the energy so lavishly given to us? And if we wish to make better use of it, can we?

In principle the latter would seem possible, but

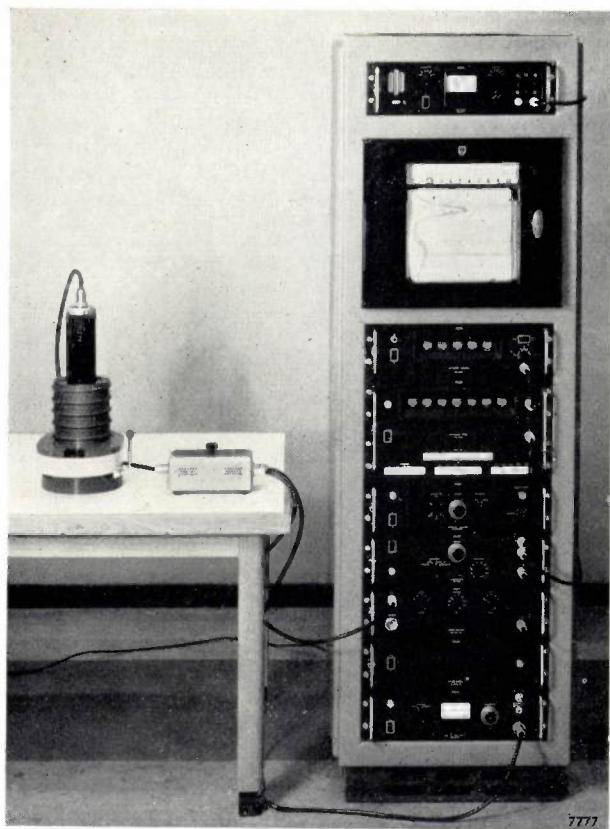


Fig. 13. Gamma-ray spectrometer for neutron activation analysis. On the table, left, is a scintillation detector enclosed in a lead shield, at the bottom of which is placed the sample to be chemically analysed. Beside it on the table is a pre-amplifier for the detector signal. The remaining electronic equipment is housed in the rack on the right.

only along biological lines. It will be a long time yet before the oceans are covered with floating plantations and all the deserts planted with trees. Another possibility would be to increase the efficiency of natural photosynthesis, that is the conversion of solar radiation into chemical energy through the agency of chlorophyll. There is, however, little prospect of realizing this on a large scale either at present. Some other means than the direct use of the sun's energy will therefore have to meet the world's growing energy requirements.

Compared with nuclear fission, nuclear fusion has the attractive aspect of giving rise to much less radioactive waste. The necessity of safely storing for many hundreds of years the large quantities of radioactive fission products from nuclear power stations is an awkward complication. Although it is thought that these problems are technically surmountable, they call for elaborate safeguards and moreover involve a great deal of expense. Intensive investigations in the field of controlled nuclear fusion have not yet, as far as is known, shown much prospect of possible technical application. This is not particularly surprising, for according to present

knowledge, temperatures of about 100 million °C would be needed to ensure success. Even on the sun the nuclear fusion processes are far from stable.

The progress of nuclear fission has been entirely different. The nuclear reactors based on this principle have nearly all, right from the beginning, worked in accordance with the physicists' predictions. In the application of nuclear power for the generation of electricity, Great Britain is in the lead. In the United States, where there is no pressing need for a new source of energy, it was preferred first of all to try out various systems on their merits. There too a number of nuclear power stations are now in operation, one of which has the largest working reactor in the world. This is the nuclear power station at Dresden, 50 miles south-east of Chicago, where one reactor delivers an electrical power of 180 MW. The reactor under construction at Indian Point *) will have an even higher output, 275 MW, and is expected to enter into operation in the course of 1962. This plant has some interesting aspects. The fuel is not only enriched with uranium-235, which is usual in this type of reactor with water under high pressure, but also contains thorium-232. This is not in itself fissionable, but during the operation of the reactor it is converted, by the capture of a neutron, into uranium-233, which is very readily fissionable in this type of reactor. In this way the fuel is continuously supplemented. One of the factors that have hitherto adversely affected the efficiency of a nuclear power station is that the maximum steam temperatures obtainable are lower than those in conventional power stations. In the Indian Point station the temperature of the steam delivered by the reactor at 230 °C is raised to 540 °C with an *oil-fired* superheater. This makes the overall thermal efficiency a great deal better than if nuclear energy alone were to be used; the nuclear energy is responsible for 60% of the power output. Finally, it is worth mentioning that the authorities have had sufficient confidence to erect this nuclear power station at a site on the river Hudson only 25 miles from the densely populated city of New York.

In Great Britain there are now two nuclear power stations in operation, each with four reactors. One is at Calder Hall (*fig. 14*) and the other at Chapelcross; each delivers a maximum electrical power output of 45 MW per reactor, i.e. 180 MW per station. Two further stations, for 275 MW at Bradwell, and for 300 MW at Berkeley (*fig. 15*) are expected to be put into full operation in

*) This reactor is now critical, and is being worked up to full power. *Ed.*

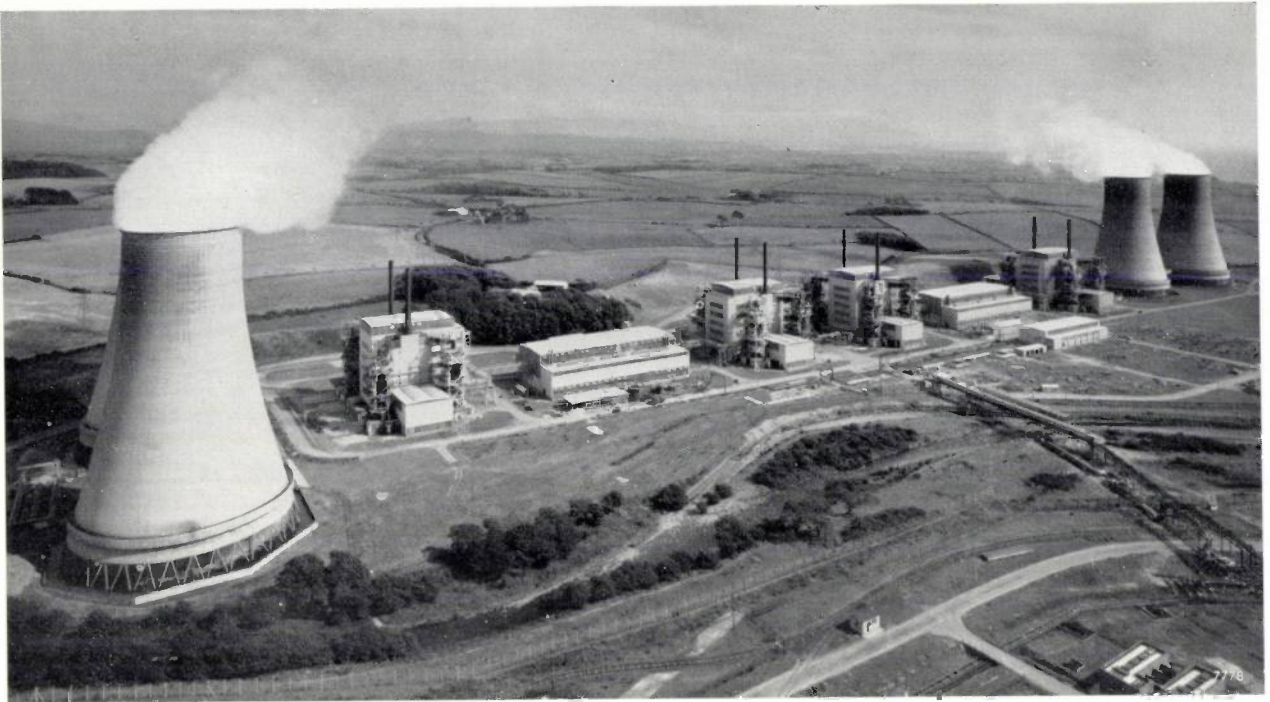


Fig. 14. View of the Calder Hall nuclear power station in Cumberland, England. The last of the four reactors was started up in April 1959 for supplying electricity to the national grid. (Photograph by courtesy of the U.K. Atomic Energy Authority).

a few months time. Several others are under construction or are in the planning stage. The nuclear power programme, to be completed in about six years, envisages a total capacity of 5000 MW. This is 14 per cent of the present capacity

in Great Britain and more than the present capacity in the Netherlands. Although the first nuclear plant, at Calder Hall, has an efficiency of less than 20 per cent, this figure will rise to nearly 30 per cent in the plants now under construction.

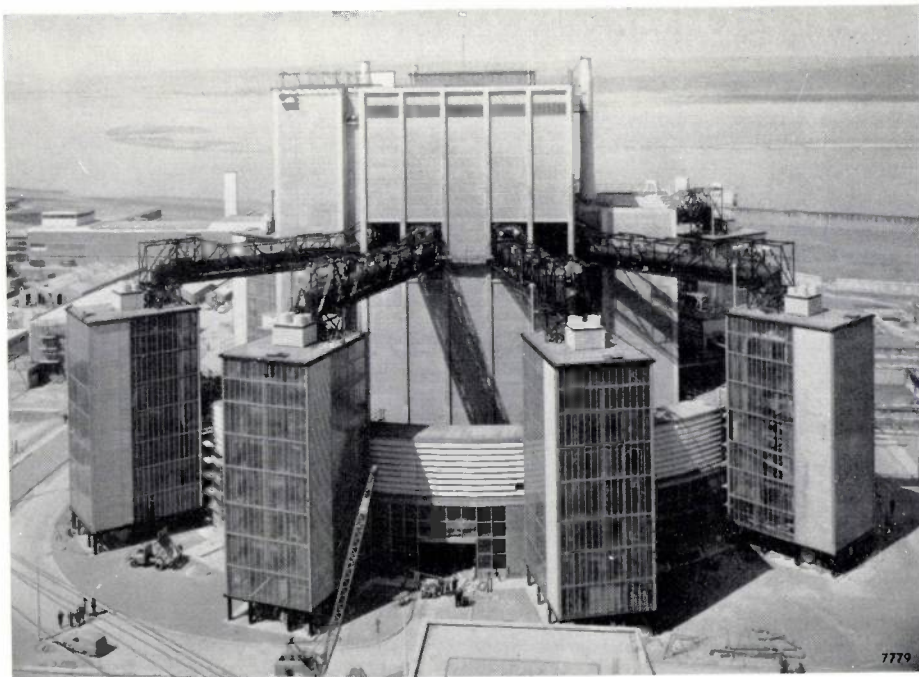


Fig. 15. A reactor plant nearing completion at the Berkeley nuclear power station in Gloucestershire, England, at the mouth of the Severn. (Photograph by courtesy of the U.K. Atomic Energy Authority.)

Fig. 1. Machine for fully automatic CO₂ welding. The carriage *W* travels at an accurately adjustable speed over a flat plate or on rails, so that the welding head with nozzle *E* moves over the joint to be welded. The welding head (shown separately on the right) can be adjusted to any desired position. With this machine welding currents up to 700 A and rates of travel up to 200 cm/min are possible.

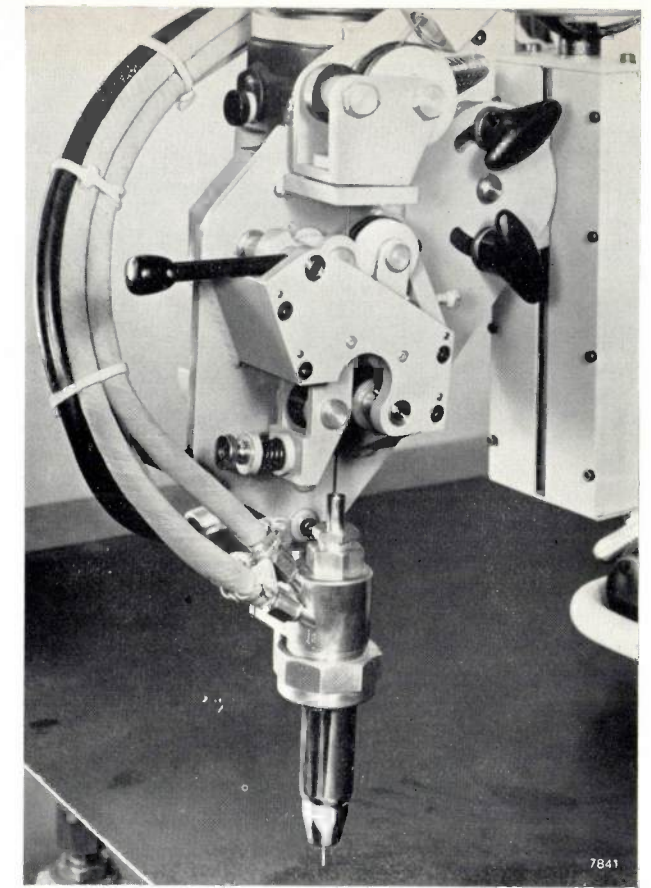
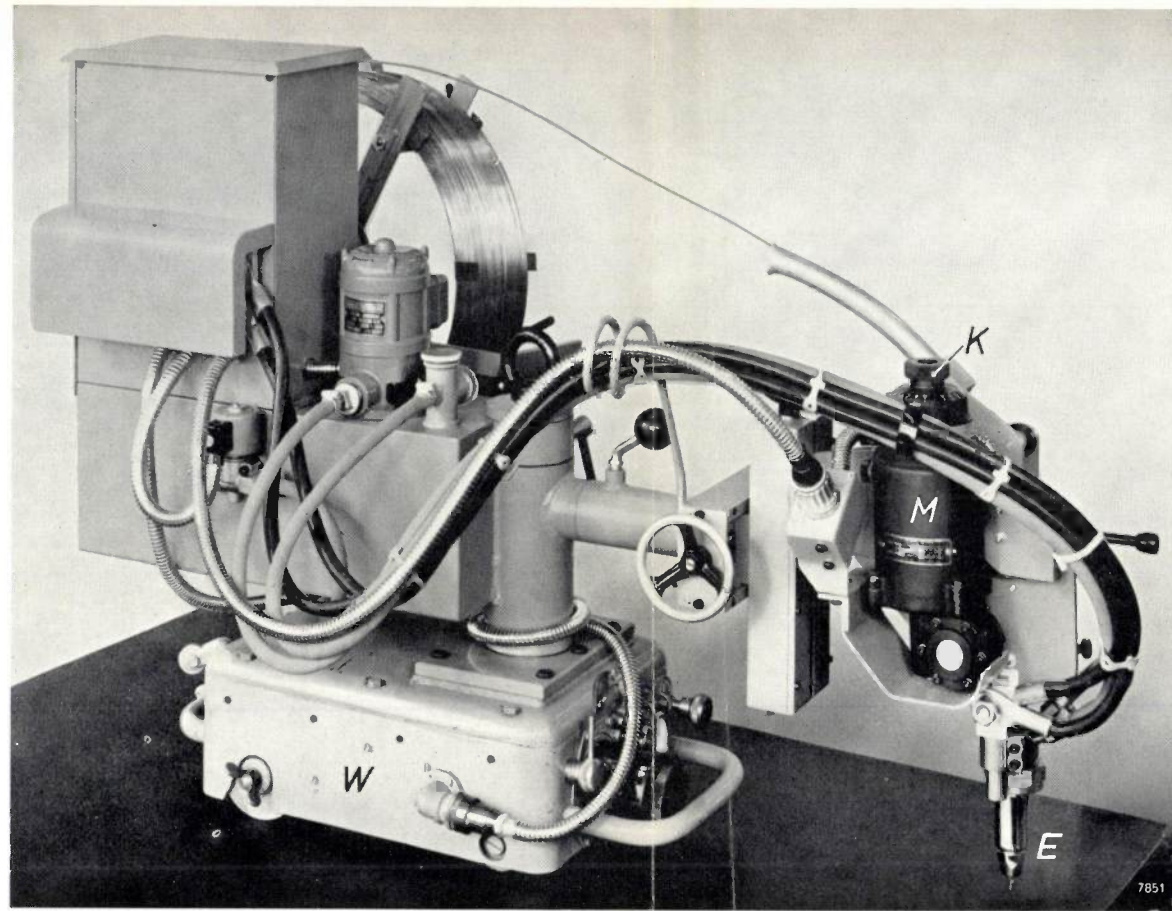
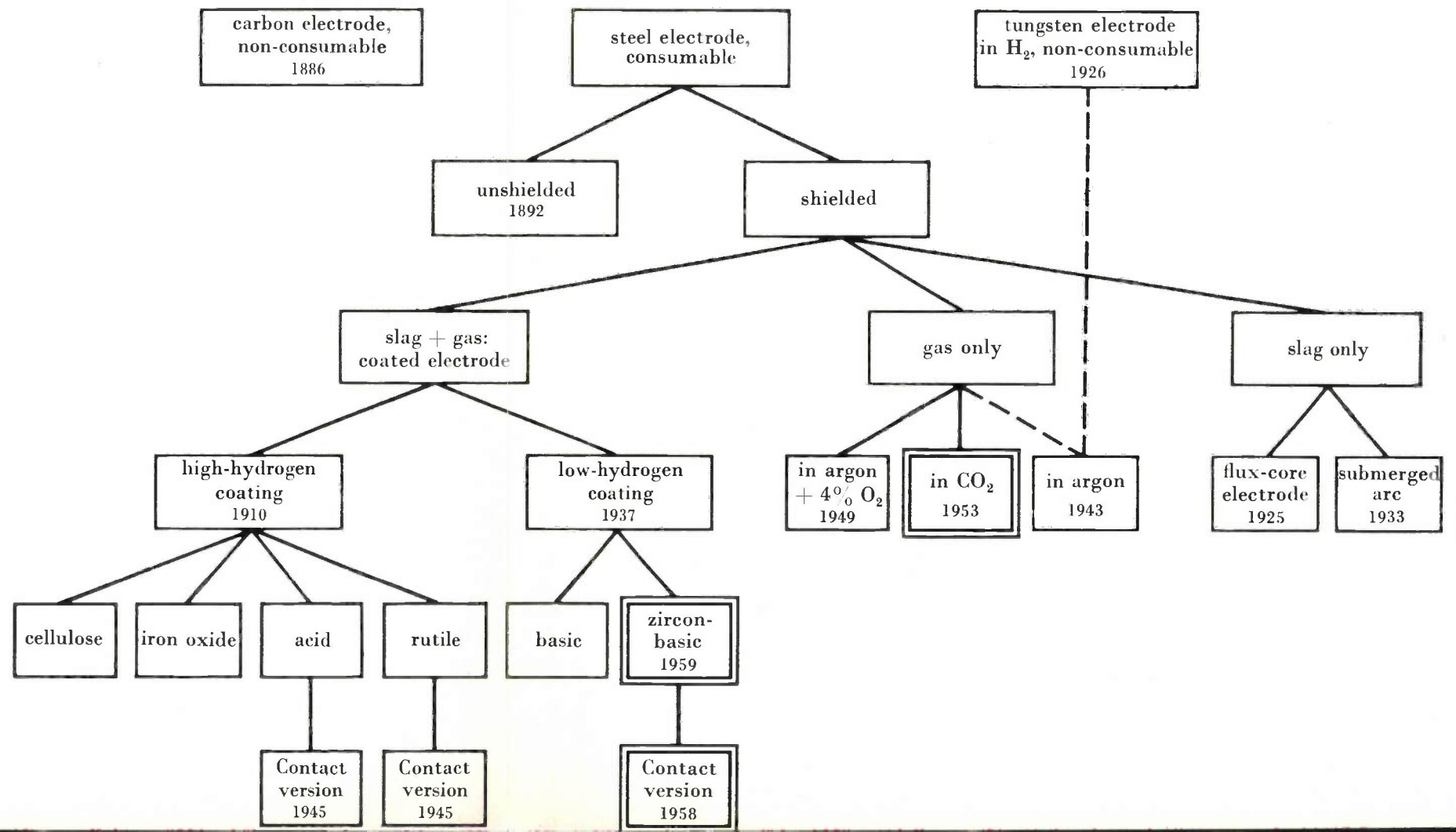


Fig. 2. Diagram of the development of various methods of arc-welding steel, classified according to the manner in which the molten metal is shielded from the atmosphere. The methods framed by double lines are discussed further in the article. The coated electrodes provide shielding by slag and gas combined. Numerous methods, not shown in this diagram, have been proposed in recent years for combined gas and slag shielding.



SOME MODERN METHODS OF ARC-WELDING STEEL

by P. C. van der WILLIGEN *).

621.791.8

Electric arc welding is nowadays the most widely used method of obtaining a mechanically strong joint between metal parts. The term arc welding covers a very wide variety of techniques. In this article the author confines himself to the arc welding of steel, including the various mild (unalloyed) steels as well as low-alloy types. Three Philips developments are discussed: Contact electrodes, zircon-basic electrodes and CO₂-shielded welding.

The extensive use made of the arc welding of steel can be judged from the annual consumption of electrodes. In the Netherlands some 275 million electrodes were consumed in 1955¹⁾. The Netherlands consumption of steel in that year was 2.5 million metric tons, so that per ton of steel 110 electrodes were used. The average weight of an electrode at that time being 40 grams, this means an electrode consumption of 4.4 kg per ton of steel.

The corresponding figure in most other countries is somewhat lower, a relatively large proportion of Dutch steel consumption being attributable to ship building, where arc welding is booming. In the United States, on the other hand, a big steel consumer is the automobile industry, where more use is made of resistance welding; nevertheless, U.S. consumption of welding electrodes in 1955 was no less than 2.9 kg per metric ton of steel.

In view of the enormous scale on which the arc welding of steel is finding application, it is not surprising that numerous investigations have been undertaken in this field and a wide variety of welding methods developed. We shall discuss here three developments that have taken place in the Philips Research Laboratories at Eindhoven: Contact electrodes, zircon-basic electrodes (both types essentially overlap each other, as we shall see) and CO₂-shielded welding. A machine for automatic CO₂ welding can be seen in *fig. 1*; particulars are given in the relevant section of this article.

To place these developments in the right setting, we shall first give a survey of welding methods classified in accordance with one particular point of view, without however attempting to be complete and without being concerned too much with chronological order. We shall return to our classification later when comparing the various methods in regard to such points as economy, quality of the weld, and the possibility of automatic operation.

Survey of arc-welding methods

In the first years after the invention of arc welding (carbon electrodes in 1886, consumable metal electrodes in 1892) it was not yet realized that to obtain a mechanically sound weld it is essential to exclude the air from the vicinity of the arc; if the droplets of the weld metal, which are heated in the arc to a temperature several hundred degrees above the melting point, are exposed to the atmosphere they absorb some tenths of a percent of nitrogen and oxygen, resulting in a brittle and porous weld.

All the welding methods to be mentioned here are therefore designed to shield the weld metal from the atmosphere. The methods are shown in *fig. 2*, classified according to the means of shielding adopted.

An effective way of protecting the molten metal against the atmosphere is to weld in vacuum. This is in fact done in some cases, the heat for welding being generated by an electron beam. This method cannot of course be counted among the forms of arc welding, and it is obviously impracticable for the great majority of arc-welding applications.

A generally useful solution to the problem of shielding the weld metal was found in about 1910 with the development of the *coated* (or *covered*) *electrode*. During the melting process the coating produces *slag*, which envelopes the droplets with a molten shell, and it also evolves *gas*, which excludes the air from the vicinity of the arc. Initially, up to about 1937, the only known coated electrodes produced an acid slag and a gas mixture consisting mainly of carbon monoxide and hydrogen²⁾. These electrodes are eminently satisfactory for welding mild (unalloyed) steels, but in the case of low-alloy steels, which are nowadays gaining increasing

*) Philips Research Laboratories, Eindhoven.

¹⁾ W. Gerritsen, Smit Mededelingen 14, 125, 1959 (in Dutch).

²⁾ Particulars of the numerous kinds of acid electrodes will be found, for example, in the *Welding Handbook*, published by the American Welding Society, New York 1960, 4th Edn. This work also contains detailed information on many arc-welding questions that can only be touched upon here in passing.

Even so, the electric power thus produced is still more expensive than that delivered by conventional electrical power stations.

The first British industrial nuclear reactors, built at Windscale, were solely intended for the production of plutonium for atomic bombs. The reactors built subsequently at Calder Hall and Chapelcross were primarily intended for plutonium production, the electric power being a by-product. The large reactors now under construction are specifically designed for the generation of electricity ²¹).

An objection sometimes made against nuclear power stations is that the radioactivity of their waste products may be a danger to public health. It is known with certainty, however, that coal-mining underground is a frequent cause of disease, in spite of all the advances made in mining hygiene in the last decades. The mining of uranium ore can also be dangerous to the health of the miners, the escape

of radioactive radon being a particular hazard. But by adopting appropriate measures, so far as they have not already been introduced, uranium mining can probably be raised to a higher hygienic level than coal mining. If this can indeed be done, the nuclear power station will represent a social advance on the coal-fired type.

Summary. Inaugural address presented at the Technische Hogeschool, Eindhoven. In the past thirty years nuclear research has led to the development of numerous techniques, and the author discusses some of their applications in various fields. Electrons of very high energy, e.g. from 4 to 35 million electron-volts, can be used for medical irradiation, or for generating very hard X-rays, which are used for the same purpose and also for industrial radiography. The electrons are obtained from accelerators, e.g. linear accelerators or betatrons. In nuclear reactors, or in cyclotrons, radioactive isotopes are produced which can be used as "tracers" in various investigations (examples: for the study of littoral drift, the absorption of iodine in the thyroid gland), for measurements based on the absorption of radiation (example: sand concentration in the outlet pipe of a sand dredger), for industrial radiography (welds in pipelines) or for medical therapy (cobalt irradiator). Mention is also made of neutron activation analysis. To illustrate the use of nuclear energy for the generation of electricity, some data are given on British and American nuclear power stations. The survey is preceded by a consideration of the hazards involved in nuclear engineering, and the author returns to this point at the end. With appropriate safeguards he considers the dangers to be no greater or even less than are encountered elsewhere in engineering or daily life.

²¹) For a survey of the nuclear power stations under construction or planned in Great Britain, see: U.K. power reactors building progress, Nucl. Engng. 7, 108-109, 1962 (No. 70). Particulars of the further development of reactors for electricity production are given by W. Vinke and P. J. Kreyger, Power reactor experiments, Atoomenergie 4, 33-39, 1962 (No. 2).

importance, they tend in unfavourable circumstances to cause cracks in the welded joints and other difficulties. This led to the development of *basic* (low-hydrogen) *electrodes*, the coating of which yields a basic slag and a gas which is largely CO and contains very little hydrogen. An independent investigation in about 1943 at Eindhoven led to the development of the *Contact* (iron-powder) *electrodes*, which were marketed in 1947 as a special version of the acid-coated electrodes³⁾. From these, with a view to applying the same principle to the basic electrodes, the *zircon-basic coating* was developed⁴⁾.

Meanwhile, however, other means were found of shielding the weld metal from the air: whereas the coated electrode does this by producing slag and gas, it proved possible to protect the weld either by slag or gas alone (fig. 2).

Before going into this point, we shall examine the question whether gas shielding alone is fundamentally sound. For if the deposited metal is shielded without slag, we sacrifice several additional functions which slag fulfils or is supposed to fulfil: 1) it provides for metallurgical purification of the weld metal; 2) since some of its constituents readily emit electrons, it has a stabilizing effect on the arc, thereby facilitating welding with alternating current; 3) it protects the weld, as the arc moves away, against too rapid cooling by the air; 4) it continues to protect the cooling metal against the action of atmospheric oxygen and nitrogen.

The first function mentioned may indeed be important on occasion, but it can be dispensed with if the deposited metal has the same (or higher) purity as the parent metal, i.e. the steel to be welded. The second function can be fulfilled by arc-stabilizing substances, added in small quantities to the steel itself. The other functions, at least as far as steel is concerned, have no real significance: cooling by the air is immaterial, in view of the very high heat dissipation through the workpiece itself; and experience shows that the air does no harm to the cooling weld in the case of steel⁵⁾.

As long ago as 1925 an electrode was brought out, designed to provide shielding by the slag alone: this was the "flux-core electrode", in which substances that form slag and also contribute to the stabilization of the arc are contained in an axial cavity of the electrode wire. This did not, however, provide adequate protection of the weld metal, and these electrodes have therefore virtually dropped out of use.

Good results, on the other hand, were obtained with another method of shielding by slag alone; this was called "submerged-arc welding", and came out in about 1933 in the United States⁶⁾. The arc in this case burns in a cavity inside a coarse-grained layer of powder (flux). The arc is therefore not visible, but this is not a serious objection as the method was in fact primarily developed for automatic welding, and has so far remained the most widely employed for this purpose. The flux is made by melting a mixture of slag-forming substances and, after solidification, grinding it into grains of a particular size. The mixture may be sintered instead of melted; in both operations, the temperature is so high as to remove all substances that give off gas, including possible sources of hydrogen. As we shall see, the absence of hydrogen is an important characteristic.

In the welding methods with gas shielding only, the absence of hydrogen is also essential. One would not think so, however, judging from the first method of this kind, *atomic-hydrogen welding*, dating from 1926, which uses an arc between two tungsten electrodes in a hydrogen atmosphere. The hydrogen atoms formed by dissociation in the white-hot arc recombine outside the arc and thus, together with the radiation from the arc, melt the workpiece and the added weld metal. The method is still used here and there for metals and alloys other than steel. Where steel is concerned, however, the appropriate method in this category is argon-arc welding, developed about 20 years ago in the United States. The weld metal is shielded here by an inert gas — which in America may be helium as well as argon, but in Europe is always argon — which is directed around the arc from a gas cylinder. In its original form, in which the arc is struck between a non-consumable tungsten electrode and the parts to be welded, the method is still employed on a wide scale for special alloy steels and other metals. Since 1949 a method derived from inert-gas welding has gained popularity for welding normal steel, the tungsten electrode being replaced by a consumable bare steel wire, which is fed from a motor-driven reel. In the classification given in fig. 2, this method is denoted as "argon + 4% O₂", the addition of 3 to 5% oxygen to the inert gas having proved desirable to avoid porous welds.

The latter method was extended by the process of CO₂ welding, which, in 1953, was conceived and

³⁾ P. C. van der Willigen, *Welding J.* 25, 313S, 1946 and *Philips tech. Rev.* 8, 161 and 304, 1946.

⁴⁾ P. C. van der Willigen, *Welding News* No. 89 (March 1958), No. 123 (July 1961) and No. 124 (Sept. 1961).

⁵⁾ This is not the case in the welding of metals that have a great affinity for oxygen and nitrogen, e.g. aluminium or titanium. When such metals are welded by a gas-shielded method, extra gas shielding may be necessary during cooling.

⁶⁾ See the book mentioned in reference³⁾; further L. Wolff in "Werkstoff und Schweissung" (Ed. F. Erdmann-Jesnitzer), Akademie-Verlag, Berlin 1951, p. 458.

developed in the Philips Research Laboratories⁷⁾: the argon, which is rather costly, is replaced here by carbon dioxide, which is cheap. Although this is not an inert gas, it provides complete protection for the weld metal when properly used. The method also has many other attractive features, and has therefore gained considerable ground in the course of a few years⁸⁾.

The classification in fig. 2 is, of course, far from complete. In particular, combined shielding by slag and gas is found not only in the above-mentioned coated electrode, but also in numerous welding methods subsequently developed, which were in fact based on shielding by gas or by slag alone. A method of CO₂ welding has been introduced, for example, in which a little powdered flux is added just in front of the arc by means of a magnetic field. There are also various techniques for special cases, such as stud welding and spot arc welding, which might also be included in this diagram. Passing these over, we shall now turn to the three above-mentioned Philips developments, prefacing our discussion with a short account of the basic electrodes.

Basic electrodes

With *acid*-coated electrodes, which were the only covered electrodes known up to 1937 (and are still used in enormous numbers), the arc atmosphere is roughly 40% hydrogen, the rest consisting of water vapour, carbon dioxide and carbon monoxide. The latter two gases originate from carbonates and carbohydrates in the coating. H₂ and H₂O come from substances such as water-glass and kaolin, added to the coating as binders and moulding compounds. For ordinary mild steel these electrodes are very suitable, but when they were used for free-cutting steel (which contains a fairly high percentage of sulphur) they gave porous welds. They proved equally unreliable for low-alloy steel, the welds sometimes showing microcracks. Other difficulties were also encountered, such as under-bead cracking, reduced tensile strength of the deposited metal (noticeable in the occurrence of "fish-eyes" in tensile test specimens taken from this metal) and flaking of enamel from the weld in subsequently enamelled parts.

⁷⁾ Netherlands patent application No. 176 664, March 1953, and U.S. patent No. 2 824 948 in the name of P. C. van der Willigen and H. Bienfait. See also P. C. van der Willigen and L. F. Defize, *Schweissen und Schneiden* 9, 50, 1957. An English translation of the latter article has appeared in *Welding News*, No. 83, 2-14, 1957.

⁸⁾ For the original form of argon-arc welding the recommended international designation is TIG (tungsten inert gas), for the variant with consumable bare wire MIG (metal inert gas) and for CO₂ welding MAG (metal active gas).

At a time when no answer had yet been found to these undesired effects, efforts were being made to make electrodes that would yield a *basic* type of slag (with CaCO₃ and CaF₂ as the characteristic constituents of the coating). At first the results were negative, and porous welds were obtained. Sound welds were not produced in this way until a start was made, in France, with "baking" the coated electrodes at high temperatures (about 450 °C). The basic electrode then quickly came into use. In the Netherlands round about 1942 these electrodes were referred to as "St 52" electrodes, since they were found particularly useful for a low-alloy steel of that name. As mentioned, this high-quality structural steel could not always be satisfactorily welded using acid-type electrodes. After 1942 it gradually became clear that the above-mentioned difficulties were attributable to hydrogen⁹⁾, and only then was it understood that it was necessary to bake the electrodes at high temperature in order to remove residual hydrogen sources such as water. For this reason, basic electrodes are still frequently called low-hydrogen electrodes to indicate their essential feature. This, incidentally, is a typical example of a successful technique being applied long before its underlying principles had properly been understood.

Basic electrodes, with their excellent mechanical properties, fell short in one important respect: they were not easy to weld with. The slag was not easily removable, the surface of the bead was not smooth, the arc was not very stable (particularly when welding with alternating current and low open-circuit voltage), and so on. The latter difficulty was due to the fluoride in the coating; because of its high electron affinity, the fluorine accelerated the deionization of the arc and thus hampered its periodic re-striking. This snag was overcome by adding to the coating magnesium powder, which acts as an arc stabilizer, or, in the case of another basic electrode, by applying two coatings one on top of the other, the inner coating containing no fluorine. The other difficulties, however, remained.

The chance of further improving the welding properties came unexpectedly in 1943 from an invention made in another connection.

Contact electrodes

In that year we were working on experiments which it was hoped would improve the thermal

⁹⁾ G. L. Hopkins, *Trans. Inst. Welding* 7, 76, 1944. A. E. Flanigan, *Welding J.* 26, 193S, 1947. See also J. D. Fast, *Causes of porosity in welds*, *Philips tech. Rev.* 11, 101-110, 1949/50; J. D. Fast, *Low-hydrogen welding rods*, *Philips tech. Rev.* 14, 96-101, 1952/53.

efficiency of welding electrodes, i.e. reduce the heat generated in proportion to the deposited quantity of metal. Possible means to that end seemed to be to reduce the depth of penetration (and hence waste less heat on needless melting of the material of the workpiece) by not subjecting *all* the metal to the passage of the current. With this in mind we decided, after various experiments, to distribute uniformly in the coating a quantity of iron in the form of iron filings (later iron powder). This resulted in a very thick coating with a high iron content. When we started welding with these electrodes, we noticed that after striking an arc it was not necessary to hold the electrode at a certain distance from the workpiece, but that it could be rested upon it, i.e. the electrode could be touch-welded. During welding with the heavily coated electrode a deep cup forms at the end of the coating which automatically keeps the arc length constant (*fig. 3*). Moreover, we

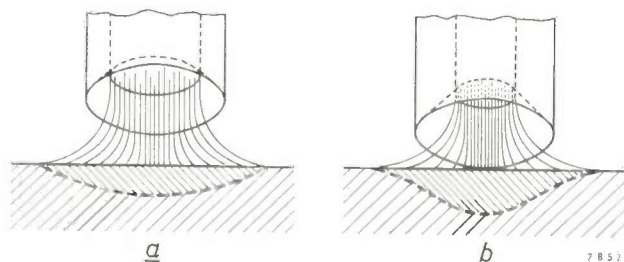


Fig. 3. *a*) Welding with a normal coated electrode and free arc. *b*) Welding with a Contact electrode. The coating being thicker, a deeper cup is formed than in (*a*). The rim of the cup can thus rest on the surface of the workpiece. The arc here is concentrated on a thinner core wire, so that the penetration in the middle is somewhat deeper than in (*a*).

found that given a suitable percentage of iron powder with the proper grain size, the coating has sufficient electric conductivity to restart the arc automatically by contact with the workpiece. Since these electrodes are in mechanical as well as in electrical contact with the workpiece via the coating we gave them the name "Contact electrodes"³⁾.

Further experiments showed that existing types of electrodes could be adapted to the Contact principle, by transferring part of the core iron in a finely powdered form to the coating, in which it is intimately mixed with the other constituents. The iron-to-slag ratio is then unchanged; the slag has the same composition and the weld is given roughly the same mechanical properties as with the original electrode. The advantages are not merely greatly increased ease of welding, but also improved shielding (due to the deep cup) against the nitrogen in the

atmosphere, and above all greater welding speed. Because of the shielding effect of the cup, the heat radiation by the arc is used more effectively, so that for a given power more metal is deposited per unit time — the effect aimed at in the original experiments, although in a different direction — and less of the molten metal is lost by spatter. As a result (and because the arc voltage is somewhat higher) Contact electrodes give an average rate of deposition of 0.24 grams of metal per ampere-minute compared with 0.17 grams with conventional coated electrodes (with a free arc).

In about 1947 Philips put on the market Contact versions of *acid* electrodes. In the following years similar types were produced by manufacturers in other countries, e.g. Russia. In Great Britain and the United States they first came into use in 1952¹⁰⁾: the welders there were initially opposed to the introduction of these electrodes which made welding easier and faster. The iron content in the coating of present-day Contact electrodes generally weighs roughly half of the iron core. The "deposition efficiency" of the electrode, i.e. the ratio between the weight of the deposited metal and the consumed core, is consequently about 150%.

Zircon-basic electrodes

Soon after 1943 we attempted to make a Contact version of the basic type of electrode, whose welding properties were in greatest need of improvement. Our efforts at the time were frustrated by the inadequate *viscosity* of the slag formed.

It will be evident that the viscosity of the slag is of considerable importance from the point of view of welding in various positions, e.g. the welding of flat grooves, vertical upward and overhead welding, and standing-fillet welding. In vertical upward welding, for example, a fairly thin, fluid slag is permissible or even desirable, as it easily runs off downwards; for standing-fillet welds, on the other hand, a thin, fluid slag is unsuitable because it collects at the lower edge of the bead, displacing the metal there, and thus weakening the final joint (*fig. 4*).

At the melting point of steel, which is about 1500 °C, it is very difficult to measure the viscosity of the slag exactly with the available instruments because the instruments are attacked. Data on this subject are therefore scarce and unreliable. Nevertheless it can be said that the viscosity of acid slag at 1500 °C is in the region of 10 poise, and that of basic slag is

¹⁰⁾ In English-speaking countries Contact electrodes are often referred to as iron-powder electrodes. In Germany the term "Hochleistungselektroden" is used.

roughly half this value. (The viscosity of both is thus much higher than that of the molten steel itself, which is 0.025 poise.) Every welder knows that basic slag is much less viscous than acid slag: he would say that the latter has the consistency of treacle and the former that of water. If now, in order

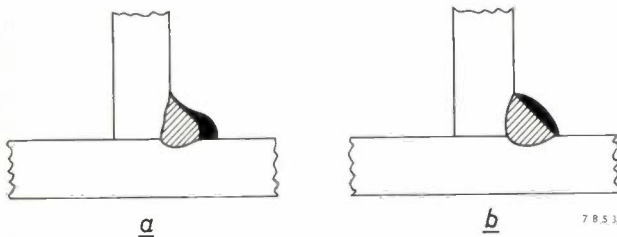


Fig. 4. For depositing a standing-fillet weld an electrode is required that forms a viscous slag. If the slag is too fluid, it pushes the molten metal at the lower side of the bead to one side; this can be seen in (a), while (b) shows the desired shape.

to make a Contact version, iron powder is added to the coating of a basic electrode, the already thin slag becomes even more fluid. Attempts to thicken the slag by adding, for example, a high percentage of rutile (TiO_2) to the coating, yielded the desired result but reduced the basicity of the slag¹¹), the consequence being poorer mechanical properties and a greater tendency for the weld to show porosity.

After many years of research, a good basic Contact electrode was at last produced in 1957⁴) by adding to the coating about 20% zirconium oxide (ZrO_2). It proved impracticable to measure the viscosity of the slag from this coating with any reliability, not only because of the difficulty mentioned but also because of the added complication that the ZrO_2 particles in the "zircon-basic" slags tend to settle out during the measurement (the specific gravity of ZrO_2 is 5.6 g/cm³, that of the basic slag about 3 g/cm³). This very fact, however, indicates that the slag must have a higher viscosity, for it is evident that the ZrO_2 crystals, whose melting point is 2700 °C, have not reacted with the other constituents and are present as such in the fluid slag. A suspension of solid particles in a liquid has a viscosity η_s which is greater than the viscosity η_0 of the liquid itself. This is expressed by the Einstein viscosity equation:

$$\eta_s = \eta_0 \left(1 + \frac{5}{2} \varphi \right),$$

where φ is the total volume fraction of the dispersed particles.

¹¹) The basicity is given by the molecular ratio of basic and acid oxides; in our case these are mainly CaO and SiO_2 respectively, with CaO : $\text{SiO}_2 \approx 1.5$.

The viscosity-enhancing effect of the ZrO_2 proved useful both for making a Contact version of the basic electrode and also for the normal versions. With the new "zircon-basic" electrode, type Ph 86, standing-fillet welds can easily be made. The Contact version of this electrode, type C 6, which has a deposition efficiency of about 160% and is now marketed alongside the free-arc version, is particularly suitable for the rapid filling of Vee joints. Owing to the greater viscosity of the slag, zircon-basic electrodes cannot however be used for vertical upward welding.

For welding in various positions an important point besides the viscosity is the *melting range* of the slag (i.e. the difference in temperature between the softening point and the melting point, which lies between 1200 and 1400 °C for most slags)¹²). This melting range can be measured by heating pieces of slag on a platinum sheet in a furnace and observing the change of shape through an optical pyrometer. The value found for various acid slags is 100-200 °C, and only about 40 °C for normal basic slags. The zircon-basic slag is found to have a much longer melting range, more like that of acid slags.

Returning for a moment to the ZrO_2 particles present in suspension, a side-effect of their presence is that the slag is light grey in colour where it has not been exposed to the air: this is, as it were, an identification mark of the zircon-basic electrode. See fig. 5, which shows photographs of the underside of two kinds of slag.

The presence of ZrO_2 crystals in the slag has also been demonstrated by X-ray diffraction diagrams. At the normal basicity of 1.5, ZrO_2 evidently behaves neutrally; at a higher basicity it appears that the compound CaZrO_3 is formed.

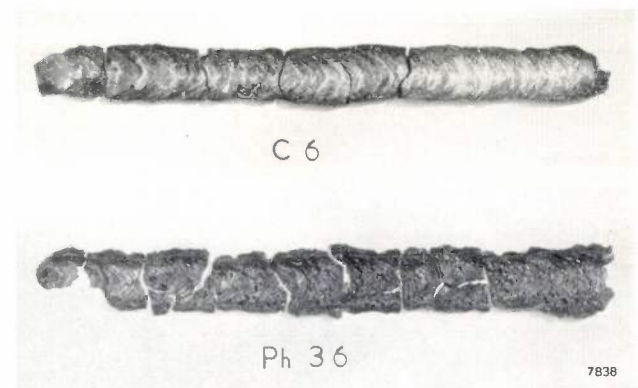


Fig. 5. Slag formed by a zircon-basic and a basic electrode (C 6 and Ph 36, respectively). The photo shows the underside of the slag, which was in contact with the bead. The slag of C 6 appears light grey in the middle where it was not exposed to the air; this is due to the dispersed ZrO_2 particles (which are also concentrated near the surface, due to sedimentation of the particles). If the bead is deposited not on a flat plate, as here, but in a groove, the air has less access to the slag and the light-grey colour is perceptible up to the edges.

¹²) G. J. Pogodin-Alexejew, *Theorie der Schweißprozesse*, Verl. Technik, Berlin 1953, p. 188.

Notable features of the C 6 and Ph 86 electrodes are the ease with which the slag is removable even from narrow grooves, and the smoothness of the weld. The surface of the weld metal is less ridged, and the transition between weld and parent metal is smoother than with the basic electrode (fig. 6). It is due to this that fatigue tests of welded joints in the

electrodes give appreciably better values; see the graphs in fig. 7a and b, which also give sketches of the type of impact bars used and some data on the welds. The plates measured $400 \times 400 \times 76$ mm. The total number of passes in each welded joint was 24. After each pass we reversed the workpiece and waited until the temperature had dropped below 100°C , after which it was cooled with water and dried. This method gave a better approximation to the situation during the welding of very thick plates in practice.

In many modern methods of welding the welder must take precautions to avoid inhaling too much of the fumes. This also applies to the basic

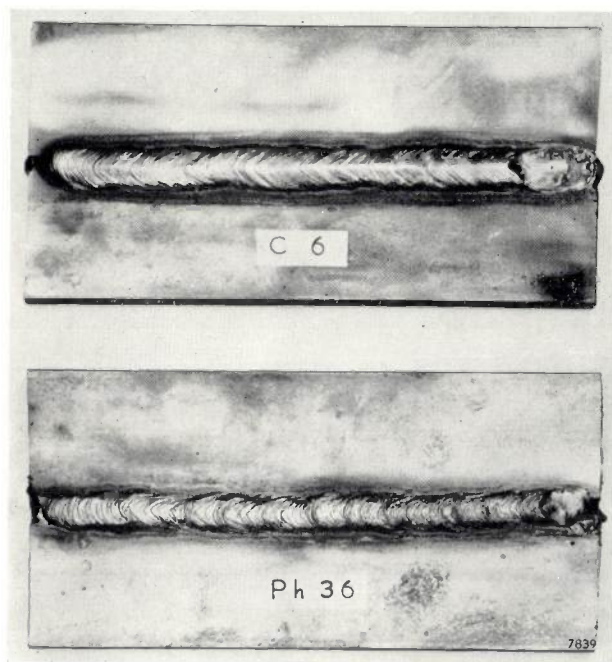


Fig. 6. Appearance of an unprepared butt weld in 13 mm plate with 2 mm gap; one made with a zircon-basic Contact electrode C 6 and the other with a normal basic electrode Ph 36. The surface of the latter weld is uneven and the weld metal does not spread out very well.

as-welded condition (i.e. not smooth-ground or stress-relieved by heat treatment) have shown values of fatigue strength about 30% higher than similar welds made with normal basic electrodes.

As mentioned at the beginning, the introduction of basic electrodes was stimulated by the increasing use of low-alloy steel, which makes lighter constructions possible but at the same time imposes high demands on the mechanical properties of the welds. Another tendency in modern engineering is the use of thicker and thicker plates, e.g. for reactor pressure vessels, for very large ships and in the chemical industry. For welding such plates, too, basic electrodes have proved equally successful. We have carried out extensive tests to determine the notch toughness of welded joints in very thick plate (76 mm) made with type C 6 and Ph 86 electrodes. At room temperature these electrodes were found to give practically the same high notch-toughness value as the normal basic electrode Ph 36 S. At low temperatures, on the other hand, the zircon-basic

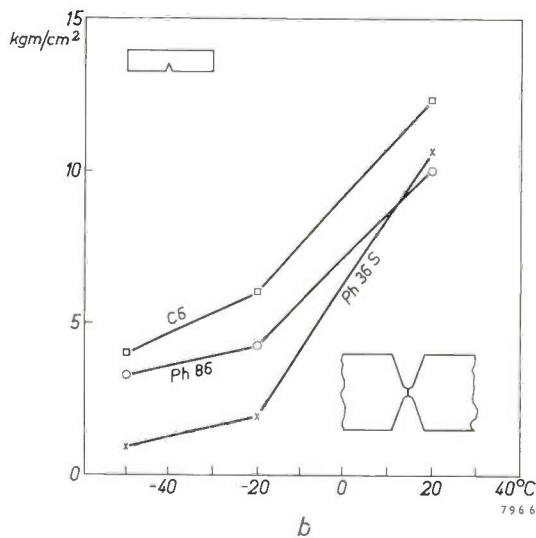
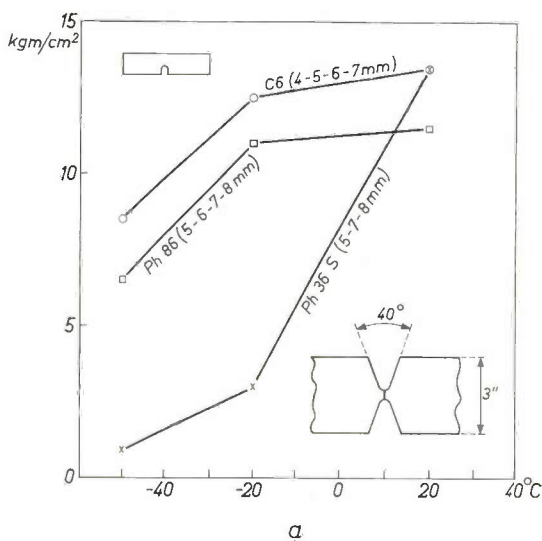


Fig. 7. a) Average notch toughness of a welded joint made with three types of electrode, as a function of temperature. The notch toughness was measured on Charpy test specimens 7×10 mm² of the shape shown top left (U notch), taken from a welded joint in 3 inch thick plate, the preparation being as shown bottom right. b) As (a) using Charpy test specimens 8×10 mm² with a Vee notch as shown top left.

electrodes, the vapours from which contain a small amount of fluoride. *Table I* shows the total amount of fluorine per 100 grams of deposited metal measured in the welding vapours from basic and zircon-basic electrodes. The quantities found for zircon-basic electrodes are considerably smaller, which was to be expected since the coating of zircon-basic electrodes contains much less fluoride.

Table I. Measured quantity of fluorine in welding vapour from various electrodes.

Electrode, type and thickness	mg F in welding vapour per 100 g deposited metal
C 6 (4 mm)	55-56
Ph 86 (5 mm)	75-77
Ph 36 S (5 mm)	97-139

CO₂-shielded welding

The development now to be described constitutes an entirely different approach to the problems of arc welding.

In 1952 we started experimental welding in mixtures of CO and CO₂, based upon our knowledge of the composition of the arc atmosphere in welding with basic electrodes. We soon stopped adding CO, which is relatively expensive and is moreover poisonous. Welding with an arc atmosphere consisting entirely of CO₂ had already been tried in 1926, but the result at that time was a brittle and porous weld. We shall describe below how this problem was finally solved, but first it will be useful to outline the whole method as at present employed.

Main features of CO₂ welding

As mentioned at the beginning of this article, CO₂ welding is based on the use of consumable bare wire which is unwound from a reel and fed by a motor-drive to the arc. This method makes it possible to weld long joints uninterruptedly, and can readily be made automatic or semi-automatic.

The principal difference between welding with separate electrodes and welding with consumable bare wire consists in the way in which the current is supplied. The current is fed to the wire close to the tip by means of a feed tube as sketched in *fig. 8*. This makes it possible to work with much higher current densities, and hence at greater speed, than when using electrodes where the current is normally applied to the rear end, i.e. at a distance of 35 or 45 cm from the arc. An electrode of 2 mm core-wire diameter, for example, is welded with a current of at the most 70 A; a higher current would make

the electrode too hot, perhaps even red-hot, giving rise to decomposition of the coating and other difficulties. In CO₂ welding with consumable wire of the same thickness, the welding currents can be as high as 700 A.

A machine for fully automatic CO₂ welding was shown in *fig. 1*. The motor (*M* in this figure) feeds the wire from the reel to the welding arc through two hardened-steel rollers. The wire is clamped with a variable pressure between the rollers, and the speed of the motor can be continuously varied by a manual control (*K*). The nozzle (*E*), in the axis of which is located the feed tube and out of which

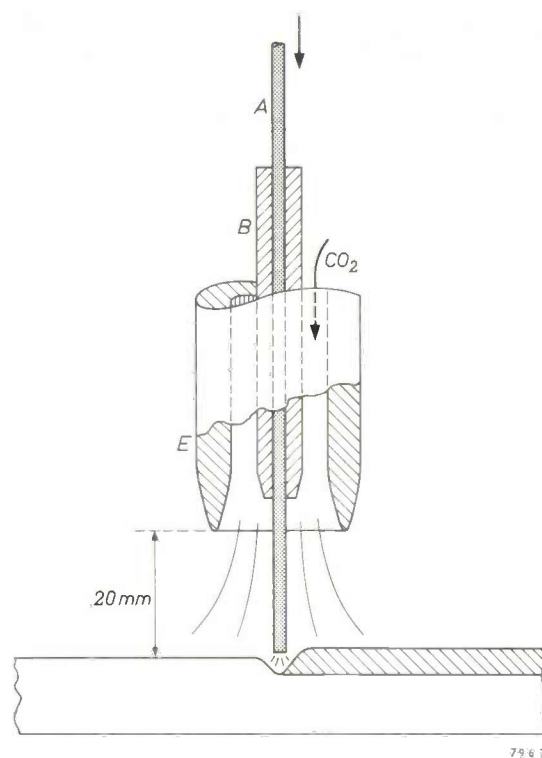


Fig. 8. Sketch of the bottom section of the welding head. *E* is the water-cooled nozzle from which the CO₂ is supplied. The consumable bare wire *A* is fed through tube *B*. The current is supplied to the wire via the same tube, the mouth of which is only about 25 mm from the wire tip. To ensure good contact it is essential not to straighten the wire beforehand.

The wire projecting from the tube *B* is pre-heated by the current it carries. This effect, which is extremely important with the high current used and the small diameter of the wire, can largely be controlled by varying the length of the wire protruding (here about 25 mm). This effect was recognized in 1943 by T. Hehenkamp and described in U.S. patent No. 2 475 835.

the CO₂ gas flows through an annular opening, is double-walled and is fitted with closed-circuit water cooling.

Complete equipment for semi-automatic CO₂ welding can be seen in *fig. 9*.

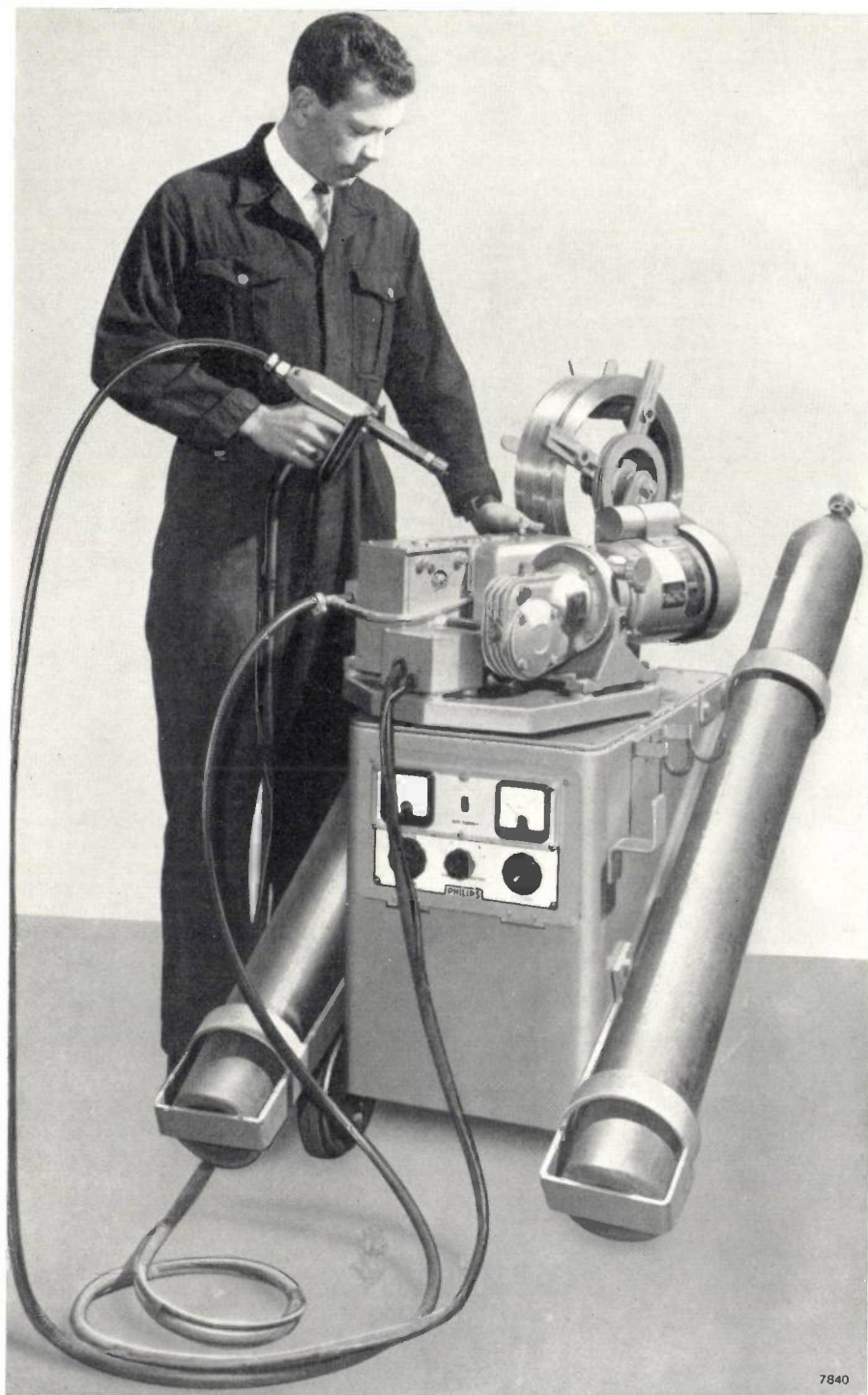


Fig. 9. Complete equipment for semi-automatic CO₂ welding. An electric motor feeds the consumable bare wire from the reel through the long flexible pipe to the welding gun, which the operator passes along the joint to be welded. The current (DC up to 300 A in this machine) is supplied near to the wire tip; the carbon dioxide (from one of the two gas cylinders) flows through a nozzle in the gun and envelopes the wire electrode and the arc.

The heat generated by the current is very efficiently used in all gas-shielded methods of welding, since no slag needs to be melted. Consequently the deposition rate (at a given current) is considerably higher in these methods than, for example, when using the electrodes discussed above: in CO₂ welding

the deposition rate is between 0.30 and 0.40 grams per ampere-minute. A further advantage of CO₂ welding is the particularly favourable shape of the penetration: whereas argon-shielded welding produces a V-shaped penetration (fig. 10a), CO₂-shielded welding produces a U-shaped penetration (fig. 10b) which is moreover, using the same current, several millimetres deeper than with argon. The upshot is that CO₂ welding, generally speaking, permits less bevelling of the parts to be welded, and thus narrower grooves that can be filled in fewer passes.

All this makes CO₂ welding a very fast process. The gain in welding speed compared with methods employing separate electrodes is of course enhanced by the fact that it is no longer necessary to stop every minute or so to insert a new electrode. Compared with the continuous method of submerged-arc welding (described at the end of this article) a favourable feature of CO₂ welding is that no time is wasted in removing slag from the bead.

Apart from the welding speed it is important to consider the mechanical properties of the welds obtained. We shall presently examine in some detail

the manner in which the metal is shielded in CO₂ welding, but it may be mentioned at this stage that the deposited metal has roughly the same mechanical properties as obtained when using basic electrodes. As to the properties of the weld, if good use is made of the deep penetration which is a feature of

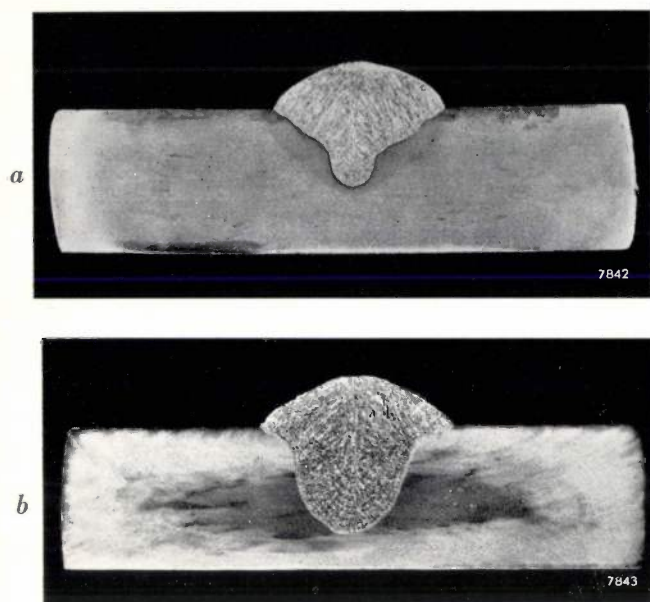


Fig. 10. Etched cross-section of weld beads, deposited with 2 mm thick bare wire on a flat plate 13 mm thick, at a current of 480 A and a rate of travel of 40 cm per minute: (a) welded in argon + 5% oxygen, (b) welded in CO_2 . In (a) the penetration is V-shaped, 7 mm deep; in (b) it is U-shaped, 10 mm deep.

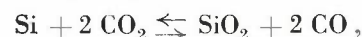
CO_2 welding, the deposited metal will be mixed with a high percentage of molten parent metal. This is one of the main points governing the mechanical properties of the weld. Also of importance is the number of layers: several layers entail repeated annealing of the previous beads, and this as a rule improves the structure and increases the toughness of the weld. The deep penetration in CO_2 welding does, however, call for extra precautions in some cases. If the width of the penetration is less than the depth, which may happen if the arc is very short and the current high, internal shrinkage cracks (hot cracking) may occur. This phenomenon is also known in other welding methods. The remedy in all cases is to increase the arc length and/or to widen the groove to be welded.

Shielding the metal in CO_2

When the molten metal in the arc is enveloped by CO_2 , it is completely shielded against the nitrogen in the atmosphere. Although the metal droplets transferred in the arc are also thereby shielded against oxygen in the air, they are nevertheless highly susceptible to oxidation: the CO_2 itself dissociates at a relatively low temperature into CO and oxygen, and at the temperature of the arc there is likely to be almost complete dissociation, i.e. a high partial oxygen pressure (see fig. 13).

The oxidizing action of the dissociated CO_2 was evidently the reason for the porosity of the welds which caused the failure of the experiments in 1926. The experiments begun in 1952 were based on the

idea that the oxidizing action could be compensated by adding deoxidizing agents to the steel wire, i.e. alloying elements such as titanium, silicon, manganese, etc., the oxides of which possess a high heat of formation and which therefore largely bind the "available" oxygen by, e.g., the following reactions:



In view of the high partial oxygen pressure expected (even higher than 50% in the case of complete dissociation in the arc) we thought at first that fairly high percentages of additions would be necessary. The experiments soon showed, however, that the CO_2 had little oxidizing action, less even than the air (partial oxygen pressure 20%) — a surprising fact to which we shall return below. Under favourable conditions the addition of 0.3% Si + 0.3% Mn was found to be sufficient to prevent porosity in the weld metal⁷). The results of a representative experiment, illustrating the relatively low oxidation, are summarized in Table II. In this experiment a particular kind of weld was made successively in three different gas atmospheres: CO_2 , air, and

Table II. Burn-off of silicon from consumable steel wire, thickness 1.6 mm, welded in three different gas atmospheres. In all three cases the wire contained 0.9% Si, together with 0.1% C, 1.6% Mn, 0.02% P and 0.02% S. The metal was deposited in a wide Vee joint in 16 mm thick steel plate, using direct current of 340 A (wire positive) at 32-33 V arc voltage and a travel speed of 40 cm/min. The joint was filled in eight passes, with interpass cooling to 100 °C. Gas feed at nozzle 25 litres per minute.

Shielding gas	Si in deposited metal %	Si burn-off in consumed wire %
CO_2	0.46	48
Air	0.30	67
Argon + 20% O_2	0.27	70

argon + 20% O_2 . In the latter atmosphere oxygen was present with roughly the same partial pressure as in air, so that this experiment could show whether the nitrogen in the air had anything to do with the oxidation. The welding wire had roughly the same composition as that now used in the CO_2 welding process; it contained 0.9% Si and 1.6% Mn, together with various other elements. The tabulated results of the analysis of the deposited metal show that in CO_2 -shielded welding relatively little Si disappears (the "burn-off" of Si is by far the smallest). Further experiments demonstrated that, as regards the burn-off of Si and Mn, the CO_2 behaves like argon plus about 9% O_2 .

This favourable behaviour makes it possible to use relatively low percentages of Si and Mn, which

can readily be added as alloying elements to the steel wire. Of course, one cannot expect the added Si and Mn to give a 100% shielding of the iron (for reasons of reaction kinetics), even though the Si, for example, is by no means completely consumed (see Table II): a small quantity of Fe is oxidized, and all reaction products together form silicates, which are found in the form of a thin film of slag on the weld bead. If a second bead is to be applied to the first, this film is in many cases not brushed off before welding.

In regard to the composition of the steel of the weld wire, it is also necessary to take account of the carbon which is always present. A very low carbon content in the deposited metal cannot be realized — even if it were desirable — since, when welding with a wire containing very little carbon (e.g. 0.03%), some carbonizing of the deposited metal always occurs (up to about 0.05%). The explanation is that the transferred droplets are enveloped in a blanket of almost pure CO, from which it is possible for iron to take up carbon.

In view of its toxicity, we cannot be indifferent to the further fate of this CO, produced by dissociation of CO₂ and also by the above-mentioned reactions with Si and Mn etc. One might be inclined to suppose that the CO would be oxidized by the oxygen in the surrounding air. Measurements of the quantity of CO present, with and without the admission of air to the surroundings of the CO₂ arc, have shown⁷⁾ that this is only partly the case: because of the marked dilution of CO with CO₂, some of the CO does not undergo combustion and remains in the welding vapours. At a supply of 20 litres of CO₂ per minute, the fumes from the welding site are found to contain 0.5 litre of CO per minute. Air conditioning is therefore just as essential with CO₂ welding as with most other welding processes, at least in small workshops.

Droplet transfer in CO₂ welding

When developing a welding process it is desirable to have a clear idea of the way in which the molten metal is transferred in the arc to the workpiece. Right from the beginning of our work on CO₂ welding we therefore studied the droplet transfer with the aid of a high-speed 16 mm cine camera¹³⁾.

¹³⁾ See the articles mentioned in reference 7). The same set-up was used for studying the droplet transfer of coated electrodes: P. C. van der Willigen and L. F. Defize, Philips tech. Rev. 15, 122, 1953/54. See also: L. F. Defize and P. C. van der Willigen, Droplet transfer during arc welding in various shielding gases, Brit. Welding J. 7, 297-305, 1960. (The Sir William J. Larke medal of the Institute of Welding for 1961 was awarded to Messrs. Van der Willigen and Defize for the latter article. — Ed.)

As a result we were able to give a reasonably satisfactory explanation for the relatively slight oxidizing action discussed in the foregoing.

The camera used can take up to 3000 exposures per second; projecting the film at normal speed (24 frames per second), the welding process is thus seen in slow motion, being slowed down more than 100 times. Fig. 11 shows the set-up employed. A weld bead is deposited on a flat steel plate. The arc remains stationary and the plate is moved, so that the cine camera and the illumination can be kept in a fixed position. When filming we used a long arc,

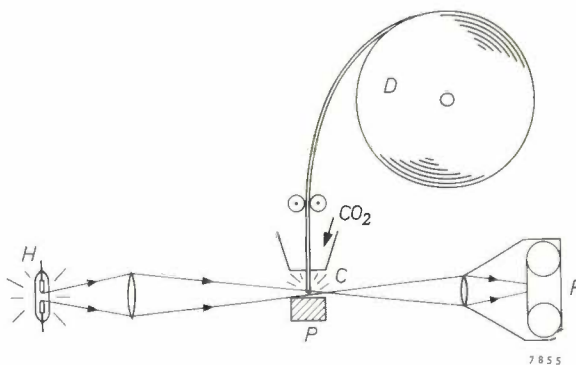


Fig. 11. Set-up for filming droplet transfer in CO₂ welding. C arc with tip of consumable wire automatically fed from the reel D. The workpiece P moves in the direction of the observer. F 16 mm cine camera capable of 3000 exposures per second. H water-cooled high-pressure mercury lamp, which, with the aid of a lens, provides a sufficiently brilliant background for the arc during filming.

the melting tip of the wire being about 7 mm above the surface of the plate, with the object of observing as much as possible of what was happening in the arc. In actual welding practice a short arc is always preferable, in order among other things to reduce spatter loss (this applies to all arc-welding methods). The tip of the wire is then located about 2 mm or even less above the surface of the workpiece. This is possible owing to the fact that, at the enormous current density used in CO₂ welding (about 20 000 A/cm²), the "arc pressure" blows, as it were, a cavity into the parent metal immediately after it has been melted, thus creating its own arc length. The arc then burns largely in the cavity, thereby reducing spatter loss and also energy losses due to radiation.

The slow-motion films show that the electrode material is transferred in a CO₂ atmosphere in the form of *very coarse droplets*, the diameter of which is between 3 and 4 mm, independent of the diameter of the wire used; see Table III. The droplets grow asymmetrically on the tip of the wire towards the side of the bead already deposited. This oblique melting of the wire appears to be bound up with the

Table III. Average size of transferred metal droplets in CO₂ welding with wire of various thicknesses.

Wire diameter mm	Current A	Current density A/cm ²	Droplet diameter mm
2.5	490	10 000	3.7
2.0	315	10 000	4.0
1.2	110	10 000	3.6

fact that at the fairly rapid relative displacement of arc and workpiece (e.g. 110 cm/min) the arc is in contact with the hottest part of the workpiece, i.e. the metal last melted, and that the tip of the wire receives radiation mainly from that side. At a given moment the growing droplet separates, the arc jumps over from the droplet to the wire, and the droplet falls spinning (this can clearly be seen in the film) along a curved path into the weld pool. Since the arc creates its own length, as we have seen above, the coarse droplets do not give rise to short-circuiting. Fig. 12 gives some frames from a colour film taken with the camera pointed obliquely downwards, showing the cavity blown by the arc, and the oblique melting of the wire.

It was known from a previous investigation¹⁴) (and confirmed by our experiments) that the droplet transfer in argon is quite different, particularly when a small percentage of oxygen has been added to the argon: *small* droplets are formed at the tip of the melting wire and split off in the direction of the wire. Increasing the current reduces the size of the droplets, until at last they are no longer separately distinguishable¹⁵).

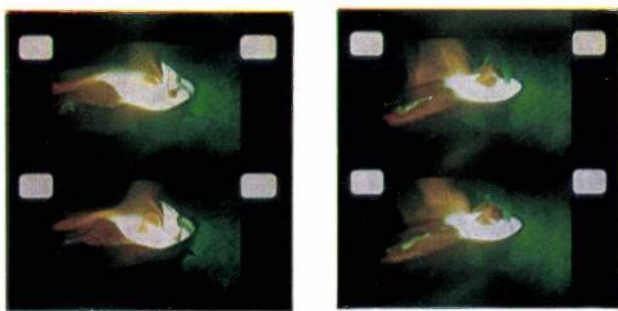


Fig. 12. Frames from a 16 mm colour film of CO₂ welding, using a short arc (as in practice). The film speed was 2000 frames per second, the camera being directed obliquely downwards, giving a view of the cavity blown by the arc in the workpiece. On the left can be seen the growing end of the weld bead; right, above the weld pool, the molten tip of the wire and the droplet passing in the arc. (This is particularly clear in the photos on the right, which also show under the weld pool the reflection of the electrode wire in the gleaming surface of the workpiece.)

¹⁴) H. T. Herbst and T. McElrath, *Welding J.* **30**, 1084, 1951.

¹⁵) Our investigation related to helium as well as argon: as far as droplet transfer is concerned, He resembles A at high currents, and CO₂ at low currents.

The slow-motion pictures reveal another remarkable difference between argon and CO₂: with argon the arc plasma envelopes the tip of the wire, including the hanging droplet, but with CO₂ only part of the droplet is covered, no more than half its surface area. Without going too much into the details, some of which have still not been explained, we should mention here that this difference in plasma coverage is probably attributable to a difference in the contraction of the plasma¹⁶). Both with argon and CO₂ there is the familiar contraction due to magnetic attraction of the current paths ("pinch effect"), but the contraction is increased by local cooling of the gas¹⁷), especially in the neighbourhood of the relatively cold wire tip, and this effect will be most marked in gases with the highest thermal conductivity — a property which is particularly noticeable in gases which dissociate at the high arc temperatures, such as H₂ and N₂ and above all CO₂ (see fig. 13).

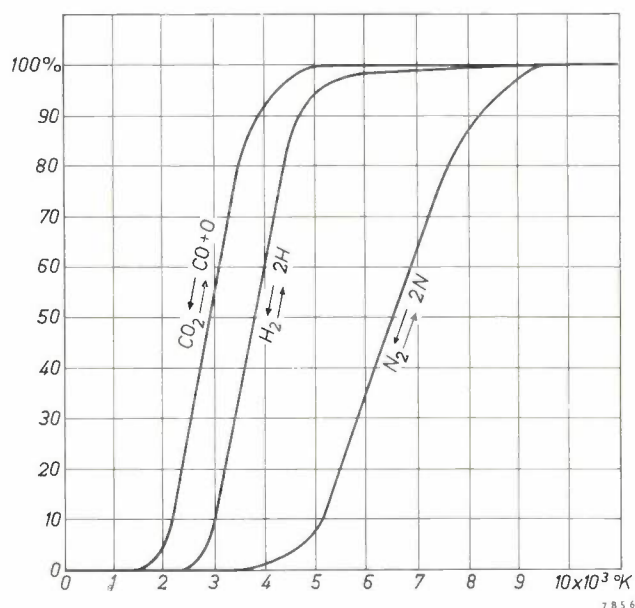


Fig. 13. Dissociation of various poly-atomic gases, as a function of temperature, at a pressure of 1 atm. (After: H. C. Ludwig, *Welding J.* **38**, 296S, 1959.)

The differences found between argon and CO₂ in regard to droplet size and plasma contraction are undoubtedly related to one another. If the cross-section of the contact surface between plasma and droplet is smaller than the cross-section of the neck of the droplet, as in CO₂, the result is an electro-

¹⁶) See last article of¹³).

¹⁷) In the "plasma burner", which is now in use for producing exceptionally high temperatures (e.g. 10000 °C), a water-cooled copper wall is employed for constricting the arc plasma. See, for example, *Welding and Metal Fabrication* **27**, 287, 1959.

dynamic force directed upwards, which is increased by the pressure of the gas in the arc. This enables the droplet to grow to quite a size until the pinch effect in the neck, together with the force of gravity and surface tension, cause the droplet to separate and the arc to jump back to the wire (see above). If the contact surface is larger than the cross-section of the droplet neck, however, as in argon (fig. 14), the resultant force is directed downwards. The droplets then break off and fall into the weld pool without having a chance to grow large.

We have dwelt on these effects at some length because, in our opinion, they make the relatively slight oxidation of the steel in the CO₂ atmosphere to some extent understandable. For owing to the large volume of each droplet transferred, and also to the fact that the droplet is only partially in contact with the extremely hot plasma, the surface exposed to oxidation is relatively small. It is probable, however, that other factors have a bearing on the oxidation, as for example the very short time in which the droplet forms and passes in the arc, and the violent perturbation of the molten material. Our hypothesis is not yet capable of providing a quantitative explanation of the processes involved, especially the measured burn-off of Si and Mn.

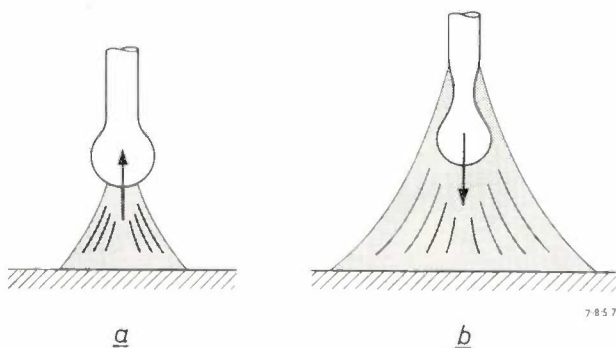


Fig. 14. Droplet formation on the wire tip, *a*) CO₂ welding, *b*) welding in argon + 5% oxygen. In (*a*) the arc plasma (shaded area) covers only a relatively small part of the droplet surface. Electrodynamical and other forces produce a resultant upward force, enabling the droplet to grow to an appreciable size; in (*b*) the resultant force is in the downward direction, and the droplets remain small.

Comparison of CO₂ welding with other methods

In the engineering industry there are many machining operations that resemble one another to some extent but which exist in their own right side by side in fields that often overlap, as for example boring, drilling, broaching, punching, milling, reaming, etc. One cannot say that any one of these methods is "the best", and the same applies to the many and various methods of welding. Any comparison between them should rather serve as an attempt to delineate their useful scope.

If we first compare CO₂ welding with the methods using basic and zircon-basic electrodes, we note that while they are all suitable for mild steel and low-alloy steel, CO₂ welding entails higher investment costs and therefore will not be considered where the volume of welding work is small. Moreover, zircon-basic electrodes give a smoother weld and will accordingly be preferred where a particularly trim finish is required. Another limitation of CO₂ welding, at least in the present state of the art, is the need to use direct current, which again calls for more expensive welding equipment. Alternating current is less suitable because CO₂ has a high ionization potential, which means that a very high open-circuit voltage would be required from the welding transformer to ensure reliable re-striking. (This drawback can be eliminated by using special wire containing as additive or coating a substance having a high thermionic emission, e.g. caesium.)

In the case of mass production, however, or where very long joints have to be welded, the investment costs are outweighed by the greater speed of CO₂ welding. This process can be used to best advantage for automatic or semi-automatic work, to which it is ideally adapted. In fact, where fully automatic work is concerned, there is only one other method that can compare with CO₂ welding, and that is submerged-arc welding, which has hitherto been the most widely adopted automatic welding process.

Submerged-arc welding was described very briefly at the beginning of this article. Fig. 15 gives a clearer idea of its operation. Bare wire is continuously fed from a reel and melted in an arc which is blanketed by the uninterrupted deposition of a

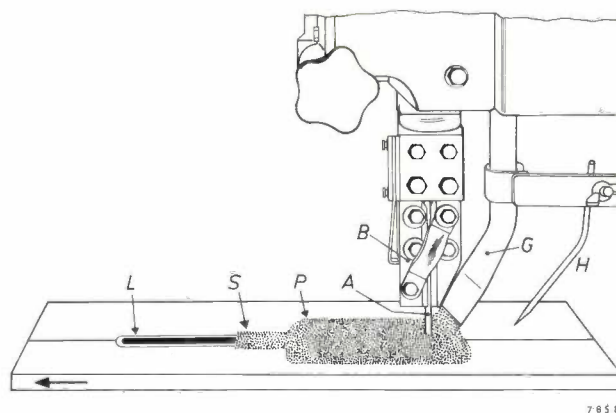


Fig. 15. Submerged-arc welding. *A* consumable bare wire. *B* current supply. *G* tube through which the slag-forming granular flux *P* is continuously deposited over the arc. The arc itself is therefore invisible and burns in a cavity inside the layer of flux. *S* weld bead with layer of slag. *L* bead from which slag has been removed. *H* adjustable pointer for directing the arc to the groove to be welded.

granular slag-forming flux. The fact that the arc is thus submerged largely governs the characteristics of the process: 1) very high currents can be used since the intense radiation is not troublesome; 2) relatively good use is made of the arc heat, which is retained inside the fusible material; 3) the method lends itself well only to fully automatic welding, for if the operator wanted to control the welding head manually he would have to work by feel, being unable to see the arc. This is the price to be paid for the advantage of not having to shield the eyes. As regards point (1) currents of 1000 A and higher were initially used, permitting very high welding speeds and often making it possible to produce thick welded joints in a single pass. It was later established that the structure of such a single-pass weld was too coarse, to the detriment of the mechanical properties. Extremely high currents are therefore no longer the rule, but the welding speeds are still very high.

In comparing CO₂ welding and submerged-arc welding, then, it must be stated that neither method differs much as far as welding speed is concerned. Preference for one method or the other will therefore depend rather on secondary factors, some of which we will mention here.

In the CO₂-shielded method the nozzle through which the gas streams gradually becomes constricted by spatter during welding, and therefore has to be regularly cleaned. Submerged-arc welding, being free from spatter, does not have this drawback. On the other hand, time is wasted with this method in removing the surplus flux and the slag formed around the weld, particularly in the case of narrow joints, where the slag has to be chipped away. This has the further disagreeable consequence of making it difficult to keep the working area clean. One of the virtues of gas-shielded welding is precisely the fact that the gas, after having served its purpose, disappears of its own accord or can easily be continuously exhausted.

We have already pointed out that the melting of the slag costs heat. This reduces the rate of deposition in grams per ampere per minute. *Table IV*, which includes some deposition rates mentioned earlier in this article, shows that the highest values are obtained in CO₂ welding.

The profile of the weld bead in CO₂ welding is often quite convex. This is especially noticeable in square butt welds; in submerged-arc welding the surface of the weld is smoother.

Submerged-arc welding sometimes gives trouble in a humid atmosphere, because the flux absorbs moisture. CO₂ welding, on the other hand, may be

Table IV. Deposition rate of metal in grams per ampere per minute (approximate values) in various welding methods.

Coated electrodes with free arc	0.17
Contact electrodes	0.24
Submerged-arc welding	0.24-0.30
CO ₂ -shielded welding	0.30-0.40

hampered by side winds, making it necessary to fit up a screen or curtain, or to increase the gas supply.

Finally, as the flux in submerged-arc welding is held in place by gravity, vertical upwards and overhead welding are of course impossible. These positions present no problems in CO₂ welding.

Although the CO₂ method of welding has made headway all over the world in a remarkably short time, and is still gaining ground, it will be clear from the survey presented here that we are not predicting that it will eventually supersede submerged-arc welding. It is rather to be assumed that engineering will have room for both methods, just as there will continue to be scope for coated welding electrodes and for the numerous variants of processes in which the weld metal is shielded by slag, by gas or by combinations of both.

Summary. The many existing methods of arc-welding steel can be classified according to the way in which the molten metal is shielded from the atmosphere. This can be done by slag (as in submerged-arc welding), by gas (as in argon-arc or CO₂ welding), or by slag and gas combined (e.g. using coated electrodes). Three developments are discussed in the latter two categories, due to Philips Research Laboratories in Eindhoven: Contact electrodes, zircon-basic electrodes and CO₂ welding.

Contact (iron-powder) versions of acid-coated electrodes have been widely used since 1947 for their ease and speed of welding and other useful properties. A Contact version of basic (low-hydrogen) electrodes, which are especially suitable for low-alloy steels, was difficult to make because of the low viscosity of the slag. The addition of ZrO₂ to the coating brought an improvement, and the "zircon-basic" electrodes developed on this basis (type Ph 86, and the Contact type, C 6) produce excellent results in down-hand welding, both as regards the appearance of the weld and its mechanical properties.

In CO₂ welding, consumable bare wire, containing about 0.9% Si and 1.6% Mn as deoxidizing additives, is fed by an electric motor from a reel to the arc, which is shielded by carbon dioxide gas issuing from a nozzle around the tip of the wire. Although a very high partial oxygen pressure is to be expected, due to dissociation of the CO₂ in the arc, the above percentage, or even lower, of Si and other additives is surprisingly enough capable of almost entirely preventing oxidation of the deposited metal. This is attributed to the large size of the droplets transferred in the arc, as observed with a high-speed cine camera, and which is probably bound up with the marked contraction of the arc plasma (in argon, for example, the contraction is much less marked and the droplets are therefore smaller).

Like submerged-arc welding, CO₂ welding lends itself very well to automatic operation. In a comparison of the various methods of arc welding a rough indication is given of their appropriate fields of application.



A TRANSISTOR CAR RADIO WITH PUSH-BUTTON TUNING

by D. PASMA *) and G. SPAKMAN *).

621.396.62:621.382.3:629.113

Two important characteristics desirable in car radios are small dimensions and ease of operation. As regards the first point, the transistor made a considerable advance possible. This article describes one of the most recently developed sets (type N5X04T), which is entirely equipped with transistors. The receiver is tuned to the principal stations by push-buttons; the mechanism employed and the measures taken to ensure accurate tuning are discussed at some length.

Differences between car radios and other receivers

Certain characteristics are required of car radios which are not so rigorously imposed, if at all, on portable sets or domestic receivers. We shall begin with a brief review of these characteristics.

The first point to be considered is the *mechanical construction*. Although a car radio should be small enough for it to be mounted in or under the dashboard of a normal car, it should be sufficiently sturdy to withstand the shocks and vibrations to which it is constantly subjected in a moving vehicle.

To screen the set against electrical interference from the engine, it must be completely enclosed in a metal housing, and this in its turn calls for special measures to ensure adequate heat removal. Further-

more the construction should make it possible to carry out repairs quickly and easily. With the latter point in mind, the mounting of the set in the car should be a simple operation.

Electrically, too, special demands are made on car radios. The power has to be supplied by the car battery, which may have a voltage of 6, 12 or perhaps 24 V, and which in some cases may be earthed at the positive pole and in other cases at the negative pole. The set should be adaptable to each of these situations. To limit the extra drain on the car battery, the current consumption of the receiver should be as low as possible.

As regards circuitry, the demands made on a car radio do not differ much from those imposed on a good domestic receiver. A high sensitivity is required from both types of set. In the fairly small

*) Radio, Television and Record-playing Apparatus Division, Philips, Eindhoven.

aerials used on cars, the signal strength is in many cases not much more than $10 \mu\text{V}$, and reasonable reception at such a low aerial voltage should still be possible. For example, at a modulation depth of 30%, a power of at least 0.5 W should be available at the output. This fairly high power is necessary in a car in view of the fairly noisy conditions usually present. For this reason, too, the *maximum* output should be fairly high; at the present time it is usual to be able to supply 5 or 6 W to the loudspeaker without appreciable distortion.

Since considerable variations in signal strength occur when the car is in motion, a very effective automatic gain control is necessary. In this respect, too, the demands are no lower than those made on a high-quality domestic receiver.

Some of the above-mentioned requirements also apply to other transportable receivers. Small size, for example, is always important. Since in most cases, however, a smaller output is sufficient, the power taken from the source and hence the heat generated are considerably lower. This makes it easier to solve the problems of heat dissipation arising from the reduced dimensions.

In some respects the demands made on the receiver part of mobile VHF equipment resemble those applicable to a car radio. Such equipment does not, however, like the car radio, have a continuously variable tuning frequency, but operates on various fixed frequency bands. This considerably simplifies certain problems of circuitry. On the other hand, a mobile VHF unit, as professional equipment, has to meet even more rigorous demands in mechanical respects than a car radio.

Development of the car radio

The building of car radios was started as long ago as in the thirties. These sets were designed for reception in the long and medium wavebands. The first receivers were "straight sets", i.e. they had no frequency changer, but a switch was very soon made to superheterodyne receivers.

The valves and other components then available were rather bulky, and the early sets, with built-in loudspeaker, had a volume of 8.5 litres. Sets of this size necessarily had to be fitted outside the reach of the driver; usually they were fitted to the bulkhead (fire wall) between engine space and car interior. The driver operated the set by Bowden cables from a control box ¹⁾. With the development of smaller valves and components, the size of subsequent series of car radios was gradually reduced, a trend which was helped by not building-in the loud-

speaker. Nevertheless, the sets were still too big to be mounted near the driver, and the use of Bowden cables still remained necessary. Because of the drawbacks attached to these cables (stiff running, backlash), efforts continued to be made to reduce the dimensions sufficiently to enable the set to be mounted in the dashboard. These efforts reached fruition in 1945, with the aid of a new range of small valves ²⁾.

The new set (type NX570V) consisted of two parts. One part, which could be mounted in the dashboard, contained the receiver proper; the power pack and the loudspeaker were contained in a separate box which could be fitted to a suitable point on the bulkhead. The two parts were connected by a multicore screened cable.

The volumes of the two parts were 1.6 l and 4.7 l respectively, giving a total of 6.3 litres. In this respect, then, not much progress had yet been made. An important advance, however, was that the manual controls were now on the receiver itself, thus dispensing with the need for Bowden cables.

A measure that contributed a great deal to the further reduction of volume was the change from tuning by means of ganged variable capacitors to *permeability tuning*, using inductors in which the position of the ferrite core can be varied. This considerably reduces the size of the tuned circuits. A further advantage of permeability-tuned inductors is that they are much less sensitive than capacitors, with their large vanes, to mechanical vibrations, and thus cause less microphony. Yet another advantage of these inductors is that, given a well designed mechanical system, they are much more readily adaptable to push-button tuning. We shall return to this point presently. Finally, it is worth mentioning that using a capacitive aerial and a tuning circuit having a constant capacitance and a variable inductance, the voltage multiplication Q is virtually independent of the tuning frequency. (This means that the sensitivity of the receiver is much less dependent on the tuning frequency.)

In the earlier sets the power pack accounted for a substantial part of the total volume. To obtain the anode voltage for the valves from the car battery, use was made of a vibrator. This converted the battery voltage into an alternating voltage, which was stepped up and rectified to produce the required DC supply of about 220 V ³⁾. Rectification was

¹⁾ J. W. Alexander, A car radio, Philips tech. Rev. 3, 112-118, 1938.

²⁾ G. Alma and F. Prakke, A new series of small radio valves, Philips tech. Rev. 8, 289-295, 1946.

³⁾ J. Kuperus, On the construction of vibrators for radio sets, Philips tech Rev. 6, 342-346, 1941, and also the article in reference ¹⁾.

effected either by using a separate contact on the vibrator or a rectifying valve. Since the unit also contained a transformer, together with smoothing capacitors and choke, it always took up a fair amount of space. The obvious line of development was towards a receiver in which the battery voltage could be used directly for feeding the valves. It was possible to make valves that could operate on an anode voltage of about 6 V, with the exception, however, of the output valve; no valve could deliver a power of 2 to 5 W on such a low voltage as this. For this reason the construction of such a receiver only became possible after the advent of the transistor. One of the first transistors capable of delivering a sufficiently high power (type OC 16) was used in the output stage of a car radio, the other stages of which were equipped with valves for low-voltage operation. The output stage contained one OC 16 transistor, or two of them in push-pull. In these sets, called "hybrid receivers", or VT (valve-transistor) receivers, which were first sold in 1957, a vibrator was thus no longer needed.

Since this made it possible at the same time to dispense with the transformer and smoothing components (choke and capacitors) a considerable reduction in volume was achieved. The radio-frequency and intermediate-frequency circuits were contained in a housing of 1.3 l volume, which could be mounted in the dashboard, whilst the audio-frequency section was mounted, without the loudspeaker, in a separate housing of only 0.9 l. The volume was thus 2.2 l. Owing to various sizes of loudspeaker being used, the *total* volume varied 3 and 4.2 l. This considerable size reduction was possible despite the fact that the set could deliver 6 W, while only 2 W could be obtained from the earlier car radios.

Although substantial space saving was achieved by eliminating the vibrator, transformer, rectifier and smoothing components, the gain was to some extent offset by the need to introduce extra suppressor filters into the supply lines. The likelihood of ignition interference was considerably greater now that the valves were directly connected to the battery.

An intractable problem in the development of VT sets was the design of an effective automatic gain control. The main difficulty was that at the low anode voltage of 6 V, even a very low negative grid bias reduced the anode current of the valves to zero. As a result, severe distortion occurred when the sets came into the proximity of strong transmitters. Although this trouble could be avoided by careful design, the normal spread in the data of the

valves and other components made it impossible to guarantee the absence of this undesirable effect in all sets.

The difficulty described is much less serious at an anode voltage of 12 V. In countries where the great majority of cars are equipped with 12 V batteries, these VT sets are therefore still being made.

The mechanical and electrical construction of the first VT sets was identical with that of other receivers; the components, for example, were interconnected in the conventional way by copper wire. In such very compact units there was a considerable risk of short-circuits, and in fact these occurred repeatedly after installation in the vehicle. This difficulty was overcome by the use of printed wiring in car radio receivers, which not only substantially reduced the risk of wiring faults, but also allowed a further reduction of dimensions.

When transistors were sufficiently far advanced to enable them to be used for radio-frequency and intermediate-frequency amplification, it became possible to build receivers which operated entirely on transistors. All-transistor car radios are now also being made, which are both smaller and consume less current. The entire set (excluding the loudspeaker) can now be contained in a housing of 1.7 litres, the dimensions of which allow it to be mounted in the dashboard. A receiver of this kind (type No. N5X04T) will be described in this article.

The circuit

The set is designed for reception in the long-wave and medium-wave bands. It is equipped with ten transistors and three germanium diodes, the various functions of which are illustrated in figs 1 and 2. These figures show the basic circuits, omitting non-essential parts. We shall consider some particulars of these circuits.

Fig. 1 represents the radio-frequency, intermediate-frequency and detector sections. *I* is the radio-frequency stage, which uses a transistor (Tr_1) of the type OC 170. K_1 and K_2 are the two RF tuning circuits, using permeability-tuned inductors. The way in which the RF tuning circuits and the oscillator circuit are switched over from the long-wave to the medium-wave band is not shown in this diagram. The aerial is connected to a capacitive tapping of K_1 . The aerial lead incorporates a choke, L_1 , for suppressing ignition interference from the engine. To enable the set to be connected to aerials of different sizes, a variable capacitor C_1 is included in the screened aerial lead (see also *fig. 5b*).

The RF stage *I* is followed by the self-oscillating

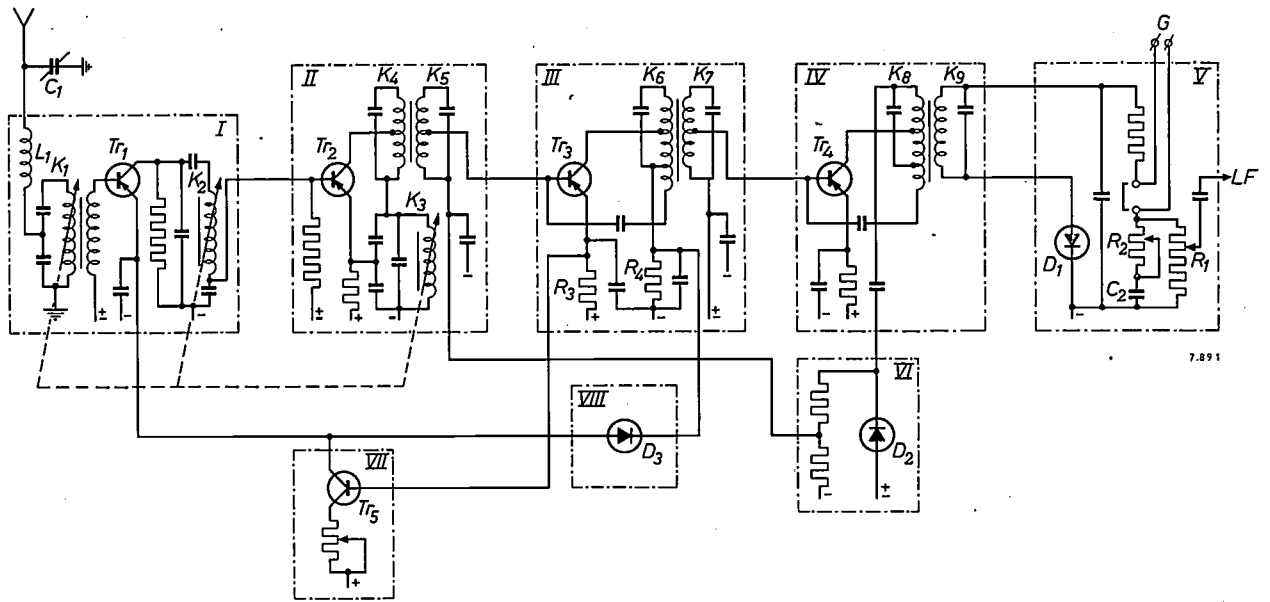


Fig. 1. Simplified diagram of the radio-frequency, intermediate-frequency and detector circuits of the N5X04T car radio. Among the parts omitted are the circuit elements required for switching from the medium-wave to the long-wave band. The symbols +, - and ± indicate that the relevant points are connected to the positive or negative supply cable or to a voltage divider between these cables. The various parts of the diagram that may be regarded more or less as separate units are surrounded by dot-dash rectangles. I RF stage. II mixer

stage. III and IV IF stages. V detector stage. VI, VII and VIII parts of circuit involved in automatic gain control. C₁ is a variable capacitor incorporated in the aerial lead (see also fig. 5b). The maximum capacitance variation is 60 pF, which makes the set suitable for connection to aerials whose capacitance (including that of the screened aerial cable) is between 45 and 105 pF. An "Automignon" record player can be connected to the terminals marked G. In that case the connection in the receiver between two points has to be broken.

mixer II, which has a transistor (*Tr*₂) of the type OC 44. *K*₃ is the oscillator circuit, which is tuned, like *K*₁ and *K*₂, by varying the position of a ferro-cube core in the inductors. A tuned transformer, consisting of the tuned circuits *K*₄ and *K*₅, forms the first intermediate-frequency band filter.

The mixer is followed by two intermediate-frequency stages III and IV, containing transistors *Tr*₃ and *Tr*₄ of the type OC 45. The second and third IF band filters are formed respectively by circuits *K*₆, *K*₇ and *K*₈, *K*₉. The inductors of the two primary circuits *K*₆ and *K*₈ are provided with a few extra turns (shown at the bottom of the coils) which are connected via capacitors to the input end of the relevant transistors. This largely eliminates the feedback in the transistors. (This technique is known as neutralization.)

The detector stage V is a normal diode detector using a germanium diode *D*₁ of the type OA 79. The potentiometer *R*₁ is the volume control, and potentiometer *R*₂ with capacitor *C*₂ form the tone control.

The parts of the circuit denoted by VI, VII and VIII form the automatic gain control.

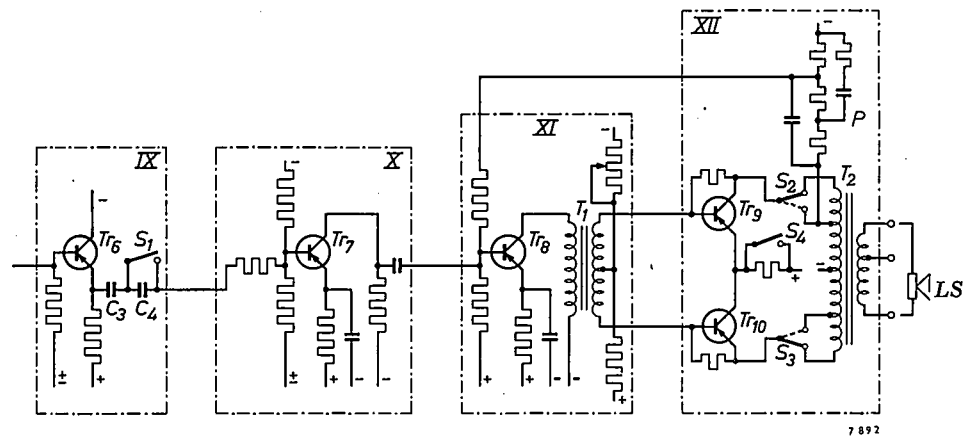


Fig. 2. Simplified diagram of the audio-frequency section. IX emitter-follower stage. X and XI audio-frequency amplifier stages. XII output stage. LS loudspeaker. P negative-feedback network. S₁ speech/music switch. If the set is operated from a 6 V battery, optimum matching between the transistors *Tr*₉, *Tr*₁₀ and the loudspeaker is obtained by using only part of the primary of transformer *T*₂. For this purpose switches *S*₂ and *S*₃ are used (shown here in the position for 12 V). If the set is connected to a 12 V battery, switch *S*₄ is used to introduce a resistor into the common emitter lead of *Tr*₉ and *Tr*₁₀; this resistor is short-circuited in the case of a 6 V supply.

The automatic gain control works as follows. The signal voltage on the primary of the last intermediate-frequency transformer, K_8 , is rectified by the diode D_2 (type OA 79) and the rectified signal is fed via a filter to the base of the first IF transistor Tr_3 . Consequently, if the signal strength increases the emitter-collector current and the gain of Tr_3 decrease. The change in the voltage across the emitter resistance R_3 of Tr_3 now controls, through the intermediary of transistor Tr_5 , the gain of the radio-frequency transistor Tr_1 . To allow the reception of very strong signals (e.g. 1 V or higher) without overloading the IF amplifier, a diode D_1 (also of type OA 79) is connected between the emitter of Tr_1 and the top of the resistor R_4 in the collector lead of Tr_1 . During the reception of small signals the emitter of Tr_1 is at a negative potential with respect to the top of R_4 . As a result D_1 is cut off. Large signals, however, make this diode conduct, chiefly because of the drop in the voltage across R_4 . Between the base and the emitter of Tr_1 there now appears a positive voltage of about 1 V, cutting off this transistor completely and allowing the very strong aerial signal to reach the mixer transistor Tr_2 only via internal capacitances, thus undergoing a strong attenuation.

Fig. 2 shows a simplified circuit diagram of the audio-frequency section. The first audio-frequency stage, IX, consists of a transistor Tr_6 of the type OC 75, in the common-collector (emitter-follower) configuration. This stage provides no voltage amplification; its presence in the circuit is to enable a crystal pick-up (of an "Automignon" record-player) to be connected to the terminals G (see fig. 1). As a crystal pick-up has a high internal resistance and thus calls for a large load resistance, it was necessary to give the first audio-frequency stage a large input resistance. One means to this end is to use a transistor in the common-collector connection, the input resistance of which is much higher than that of transistors in other arrangements. A consequence of this high input resistance is that the detector stage V also has a high input resistance. This made it possible to connect the detector in parallel with the last IF circuit K_9 , and no tap on the coil of K_9 was necessary for this connection.

Stages X and XI are audio-frequency amplifiers. The coupling between IX and X is effected by the two capacitors C_3 and C_4 in series. The capacitor with the lower capacitance, C_4 , can be bypassed by the switch S_1 . When S_1 is open the low notes are not so strongly reproduced, which is an advantage for the reception of speech (speech/music switch). Transistors Tr_7 and Tr_8 are types OC 71 and OC 79 respectively. The latter transistor is coupled via a transformer T_1 to the push-pull output stage XII. This contains two transistors of type OC 26 (Tr_9 and Tr_{10}). The loudspeaker LS is connected across the output transformer T_2 . If required, two loudspeakers can be connected in parallel. To ensure proper matching in that case,

a tap is provided on the secondary of T_2 . Negative feedback is provided between stage XI and stage XII via the network P in the figure. The output stage can deliver a power of 6 W, at which the distortion is 10%.

The diagram in fig. 3 again simplified, represents the power supply circuit. The battery is connected to the terminals a and b. The chokes L_2 and L_3 , together with capacitors C_5 and C_6 , form the interference-suppressor filters. The on/off switch S_5 is combined in the conventional way with the volume control R_1 (see fig. 1).

To adapt the set for operation on a 6 V or 12 V battery, according to requirements, various connections have to be changed. The switches used for the purpose are S_2, S_3, S_4, S_7 and S_8 in figs 2 and 3;

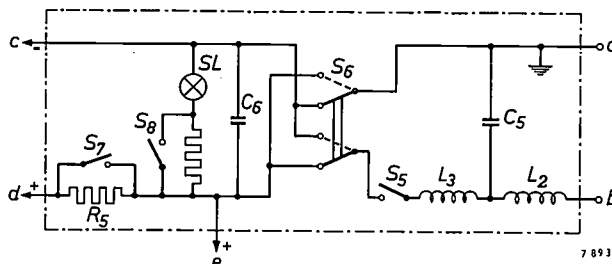


Fig. 3. Simplified circuit diagram of the power pack. S_5 on/off switch. The switches marked S_7 and S_8 serve for adapting the receiver to the car battery voltage (6 V or 12 V). The battery is connected between points a and b. The supply voltage for the output stage XII is obtained from the point marked e, use thus always being made of the entire voltage available (6 or 12 V). All other stages, however, are fed from point d. At both battery voltages the supply voltage in the set is kept at 6 V by opening or closing switch S_7 , as the case may be. The dial illumination bulb SL is connected directly to the 6-V battery (S_8 closed). If the battery voltage is 12 V, a resistance is connected in series with the lamp by opening S_8 . The double-throw switch S_6 enables the set to be operated irrespective of whether the positive or negative battery terminal is connected to earth: turning the switch to one or the other position ensures that points c, d and e always have the indicated polarity.

further particulars are given in the captions. The double-throw switch S_6 in fig. 3 enables the set to be used in cars in which the positive pole of the battery is connected to earth as well as in those where the negative pole is earthed.

In the actual receiver it is a particularly simple matter to switch from one supply voltage to another. The leads to the switches S_2, S_3, S_4, S_7 and S_8 terminate at contact points which are all disposed on a small panel. A multipole plug effects the necessary connections. The contacts are so arranged that, in order to change over from a 6 V to a 12 V supply, or vice versa, the plug simply has to be reversed. The polarity switch S_6 also takes the form of a reversible plug. Both plugs can be seen in fig. 6.

Construction

The problems involved in the construction of a car radio are not only concerned with the efforts to reduce dimensions, mentioned at the outset. Small dimensions are also important in portable receivers, but in their case the construction is simplified in as much as the manual controls and the tuning dial can be disposed if necessary on different sides of the set. Since a car radio has to be mounted in the dashboard, the controls and dial all have to be on the same side, and indeed on one of the smallest sides in view of the relatively small space available on the dashboard. An added difficulty is, that it should be possible to operate the set wearing fairly thick gloves, and therefore the controls should not be too small. The receiver type N5X04T measures $180 \times 174 \times 54$ mm. The space thus available for a tuning dial and the controls was 54×180 mm. This side also had to accommodate the means of securing the set in the dashboard, to which we shall return presently.

The receiver can be tuned by means of five *push-buttons*, which select three stations in the medium-wave band and two in the long-wave band. A control knob is provided for tuning to other stations in the normal way. The facility for tuning to a station by pressing a push-button is of considerable importance in a car radio, since it means that the driver's attention is not so distracted as when tuning by ear with a knob. If he knows the sequence of the five stations to which the push-buttons are preset, he can operate the set by touch alone.

The construction of the set was substantially influenced by the choice of push-button tuning. We have already mentioned that permeability-tuned inductors are preferable for this purpose to variable capacitors. The rectilinear movement of the inductor cores in the tuning process, and the small mass of the moving parts, make it easier in this case to design a simple and accurately functioning tuning system than if conventional variable capacitors were to be used.

In a receiver without push-button tuning, the RF tuning section is not tied to a particular part of the housing. The mechanical coupling between the tuning control, the variable inductors and the dial cursor can almost invariably be satisfactorily effected by a system of cords or other coupling elements, particularly since there need be no fixed relation between the tuned frequency and the position of the tuning control: the station is after all tuned in by ear. In the case of push-button tuning, however, such a mechanism is unsuitable owing to the considerable precision required to adjust the tuning device by means of each push-button. The motion of the push-buttons must be transmitted as directly as possible to the ferrite cores. The inductors should therefore be brought forward as much as practicable, roughly in the middle of the set, so that they form a single assembly together with the push-buttons. To achieve the desired compactness the remainder of the circuit should be grouped around this assembly so as to waste as little space as possible. *Fig. 4* illustrates how this is done in the receiver under discussion. *F* is the front plate of the set (seen from

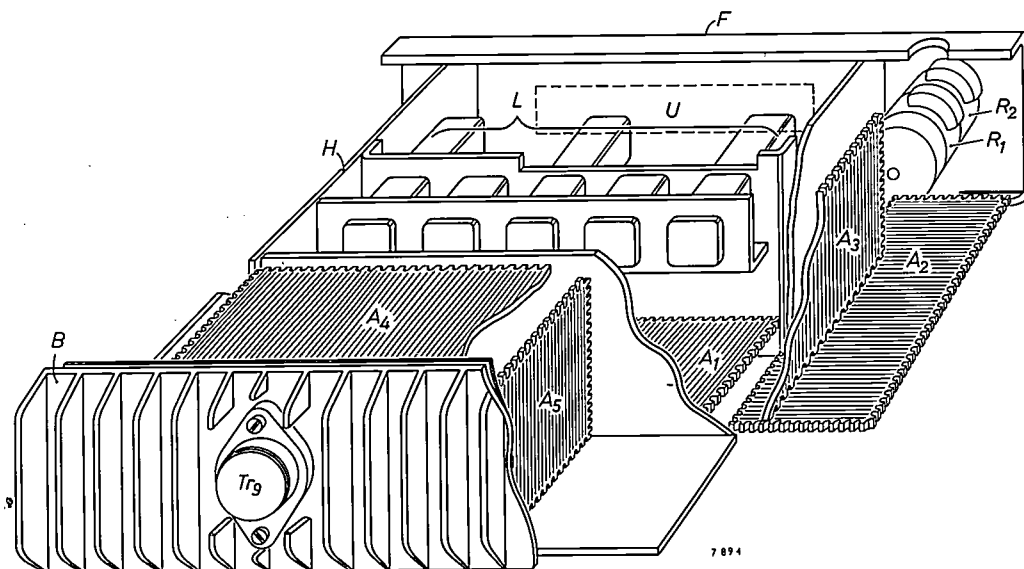
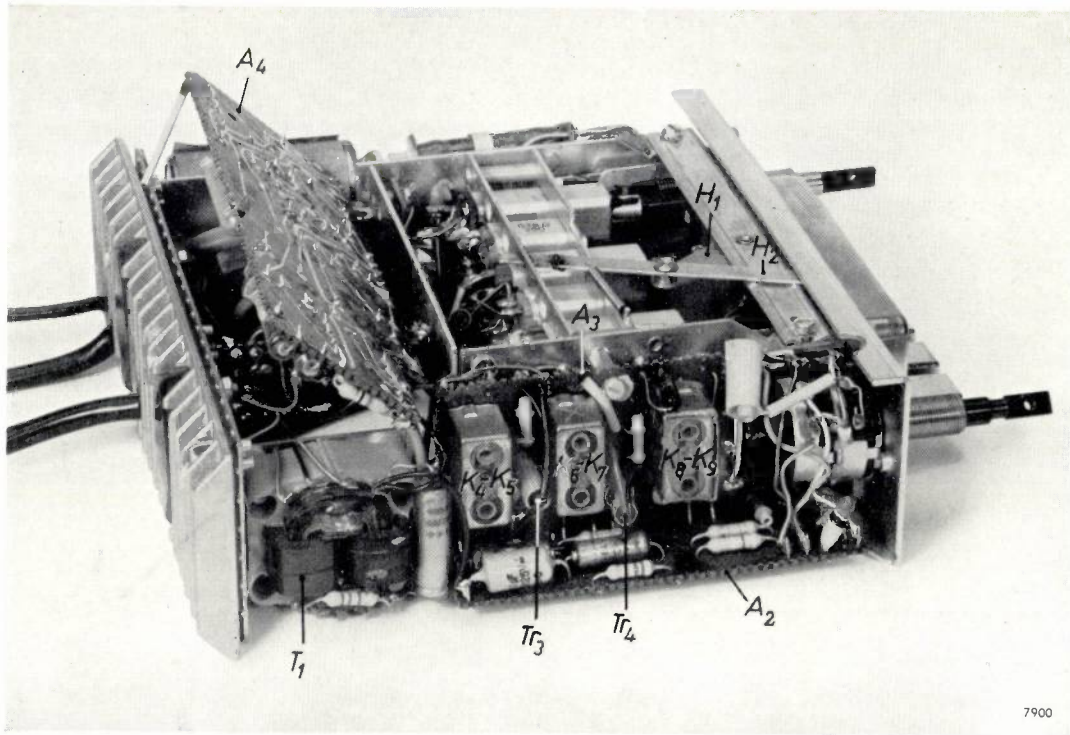
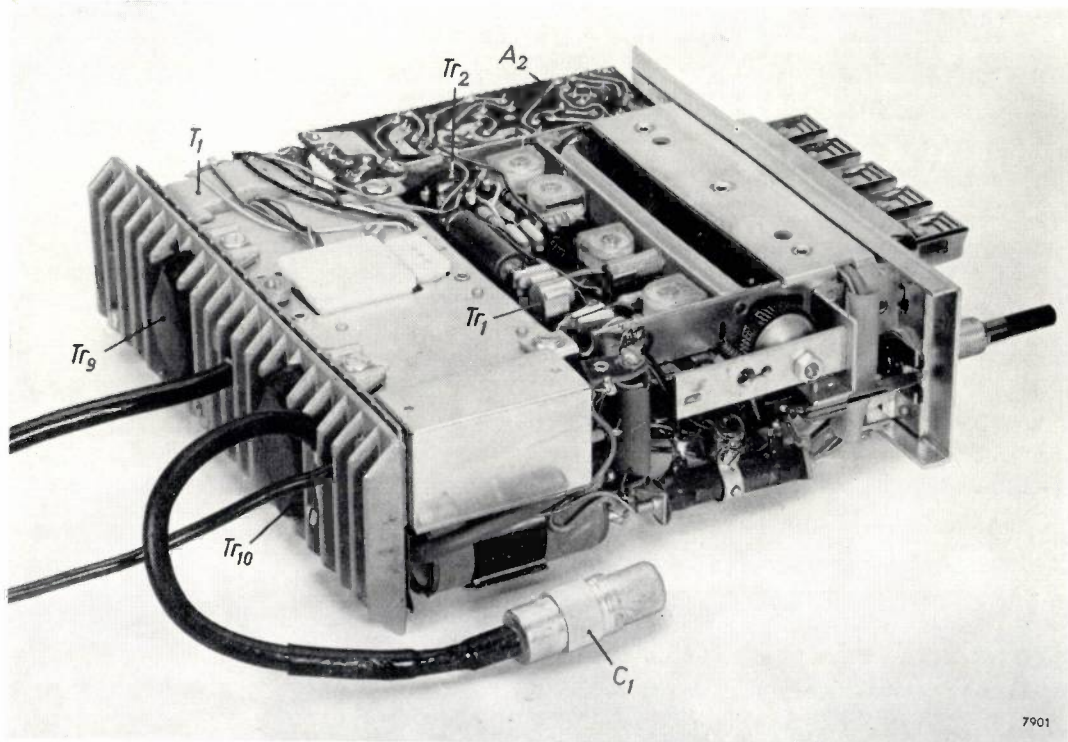


Fig. 4. Sketch showing the lay-out of the principal parts of the N5X04T car radio. *F* front plate. *U* dial. *A*₁ to *A*₅ printed-wiring panels. *B* cooling plate. *H* chassis plate. *L* inductors. *Tr*₁ one of the two output transistors. *R*₁ volume control. *R*₂ tone control.



a



b

Fig. 5. The receiver with the housing removed.
 a) Right way up. Panel A_4 is raised to allow access to the components beneath it.
 b) Upside down. In the foreground can be seen the variable capacitor C_1 in the aerial lead (see fig. 1), enabling the set to be matched to aerials of different capacitance. The letters denoting various components correspond to those used in figs 1 to 4.

the rear), on the outside of which is the dial U . Immediately behind are the inductors L of the radio-frequency and oscillator circuits. The five plates marked A_1 to A_5 are printed-wiring panels⁴⁾. On these panels are mounted most of the smaller components, such as resistors and capacitors. Panel A_1 carries most of the components of the radio-frequency section, panels A_2 and A_3 carry the intermediate-frequency and detector components, and A_4 and A_5 carry the audio-frequency circuits, with the exception of some heavy components like the transformers T_1 and T_2 (see fig. 2). The latter are mounted directly on the metal chassis.

Although much less heat is generated in transistor sets than in receivers fitted with valves, in this case it was nevertheless necessary to pay attention to cooling. The most heat is generated in the two output transistors, Tr_9 and Tr_{10} ; these are therefore fixed to the aluminium plate, B , which forms the back cover of the set. This plate is provided with cooling ribs and the transistors are mounted on the outside, so that they are in fact outside the actual set. (In fig. 4 only one of the output transistors is visible; see also fig. 5.)

Opposite panel A_3 are the volume control R_1 and the tone control R_2 (see also fig. 1), which are operated with the aid of two concentric shafts, and thus take up little space on the front plate. The on/off switch (S_5 , see fig. 3) is as usual coupled to R_1 . Mounted on the other side of the front plate is the mechanism (not visible in fig. 4) for tuning to stations that cannot be selected with the push-buttons, and also the speech/music switch (S_1 in fig. 2). The operating shafts for the tuning control and S_1 are also concentric. The part of the set to the left of the compartment marked H in fig. 4 contains the power supply components (fig. 3).

The connection points for the printed-wiring panels are mostly arranged around the periphery, enabling each panel to be easily disconnected and replaced if necessary. Panel A_4 is connected by means of flexible wires, and is held in place between lugs. It can thus be pulled up, as shown in fig. 5a, without breaking any of the connections, thus allowing easy access to the components beneath it.

Various other measures have been taken to ensure ready access to the majority of components. The chassis is so designed that after disconnecting only one lead and removing two screws, the rear section can be turned through 90° with respect to

the front section. It can be fixed in that position with two screws.

Fig. 6 shows a receiver opened up in this way, which is not only useful for making repairs but is also used in assembly.

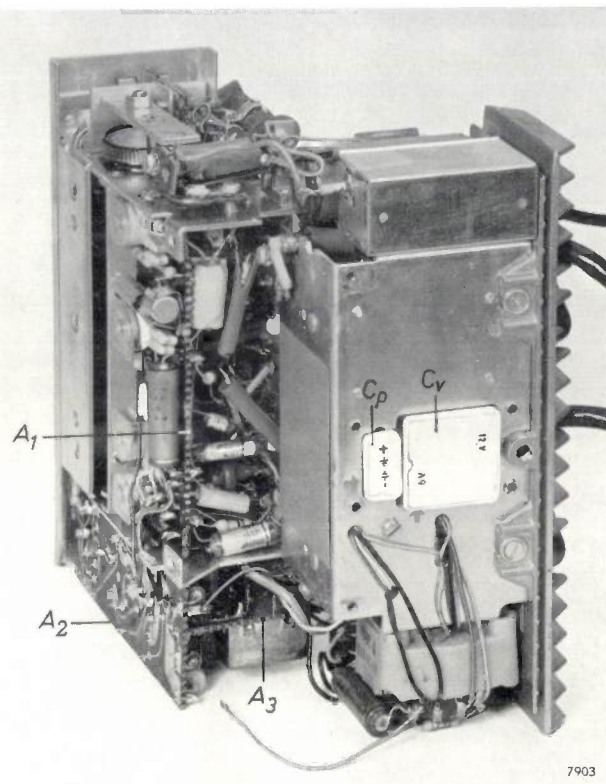


Fig. 6. For assembly or repair work, the rear portion of the receiver can be turned through 90° with respect to the front part. C_v and C_p are reversible plugs for adapting the set to batteries of different voltage (6 V or 12 V) and different polarity. A_1 , A_2 , A_3 are printed-wiring panels.

The tuning mechanism

We shall now consider in more detail the mechanism of tuning by means of the five push-buttons. Each push-button ensures that the ferrite cores in the inductors of both RF circuits and in the oscillator circuit are moved into pre-set positions. The mechanism employed is shown in a very simplified form in fig. 7. The cradle N pivots around the line AA' and is coupled to the bar J_1 . The movement of this bar is transmitted by rods T to the ferrite cores in the inductors L . Only one of the five push-buttons, D , is represented in the drawing. Mounted on the push rod E of this button is a semi-circular segment S . When the button is depressed, the cradle N takes up the position corresponding to the position of the segment, and this in turn brings the bar J_1 into the position corresponding to the tuning of the receiver to a

⁴⁾ R. van Beek and W. W. Boelens, Printed wiring in radio sets, Philips tech. Rev. 20, 113-121, 1958/59.

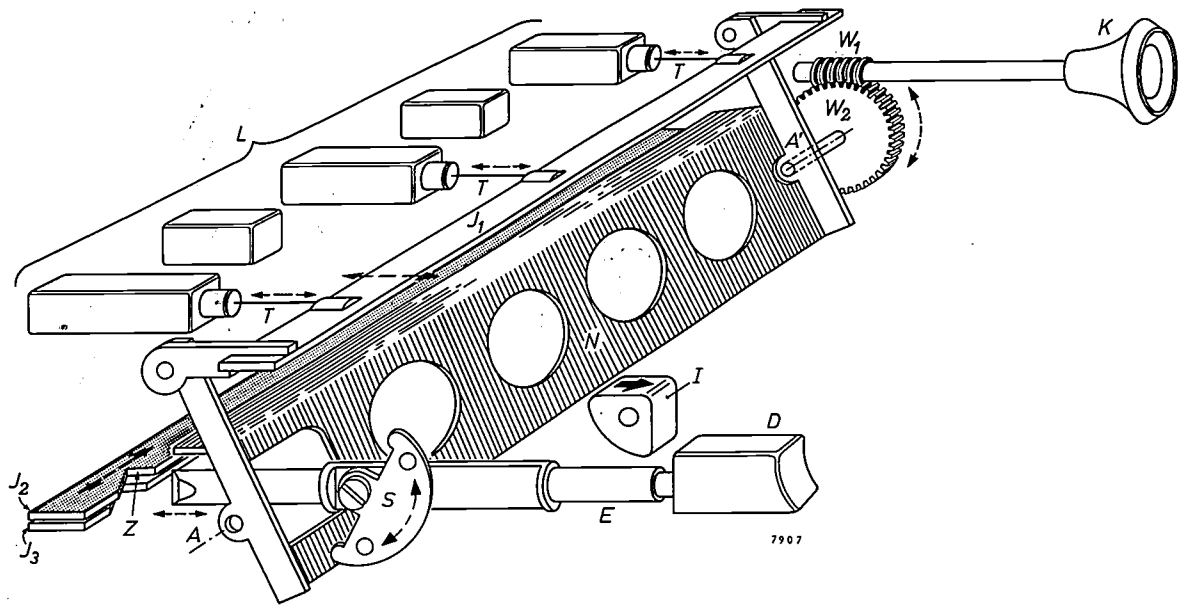


Fig. 7. Simplified sketch of tuning mechanism. *D* push-button. *E* push rod. *S* segment. *N* cradle which turns about the line *AA'*. *W*₁ worm. *W*₂ wormwheel. *K* tuning knob. *J*₁, *J*₂ and *J*₃ moving bars. *T* tie-rods. *L* inductors. *I* indicator block.

particular station. Since the segments *S* of the five push-buttons all have different positions, the operation of each button tunes the set to one of five different stations. When none of the buttons are depressed, the cradle can be moved by means of the control knob *K*, which is coupled to the shaft of the cradle *N* by a gear system consisting of a worm *W*₁, a wormwheel *W*₂ and a friction clutch (not shown).

Each segment *S* in its normal position is locked to its push rod *E* so that, whenever a push-button is depressed, the cradle *N* takes up the same position. The position of *S*, however, can be altered very simply. A mechanism, not drawn in fig. 7, loosens segment *S* when the push-button is pulled outwards. By now tuning with *K* to a certain station and then pushing *D* in again, *S* takes up the position corresponding to that of *N*. When the push-button is next depressed, *S* is again locked to *E*. In this way, then, a station once tuned in with the tuning control *K* can repeatedly be obtained by depressing one of the push-buttons.

The push-buttons also operate the waveband switch. The sliding contacts of this switch are mounted on a strip of insulating material which is connected to the sliding bar *J*₂ shown in fig. 7. The latter is provided with notches, one of which (*Z*) can be seen in the figure. Whenever a button is depressed, the chamfered end of the relevant rod *E* pushes against the bevelled edge of one of these notches, causing *J*₂ to take up one of the two extreme posi-

tions, and thus setting the waveband switch to one of the two wavebands. The notches in *J*₂ are so arranged that the long-wave band is switched in with two of the push-buttons and the medium-wave band with the three others.

Another sliding bar *J*₃ is moved at the same time as *J*₂; this serves to disconnect the friction clutch between the wormwheel and the cradle shaft when one of the push-buttons is depressed. (We shall return to this point later.)

The use of a sliding switch instead of the conventional rotary type made it possible to place the contacts immediately under the inductors, thus giving short connections and a very convenient layout.

Directly behind each push-button there is a tilting block *I* marked with an arrow. When a push-button is depressed, the block is tipped into a position in which the arrow becomes visible in front of an opening in a screen (not shown in the drawing). When the button is released, the block remains locked in that position. When another push-button is depressed the block is unlocked and springs back into its original position, an arrow then appearing above the other push-button. In this way it can always be seen which button was last depressed, allowing rapid identification of the station selected.

When the receiver is tuned either by push-button or tuning control, the station indicator must move along the dial. Its movement must be perpendicular to that of the sliding bar *J*₁. The most obvious method of coupling the tuning mechanism to the cursor

is to use a cord passing over pulleys. A cord drive system, however, could not be used here because the friction involved would detract from the high precision required of the tuning mechanism. The simple linkage mechanism adopted, and shown in *fig. 8*, caused much less friction.

The link H_1 pivots about a fixed point 1 and carries at end 3 a spindle which rests in a slot in the bar J_1 (see *fig. 7*). The end 2 is hinged to a strip H_2 , which can slide in a slot in a fixed plate Q (The spring 4 keeps H_2 pressed against one side of this slot.) The other end of H_2 carries the cursor H_3 , which is easily detachable. U denotes the position of the dial. It can easily be shown that the pivoting of H_1 (caused by the spindle 3 following the movement of J_1) results in a practically linear movement of H_3 ; with the dimensions of the components used in this set, the maximum deviation from a straight line is only 0.3 mm. The cursor thus travels straight along the dial with sufficient accuracy.

Tuning precision

As mentioned in the foregoing, a very high degree of precision is required of the tuning mechanism. Some figures will illustrate this. In the inductors used in the N5X04T set the total variation of the self-inductance is achieved by displacing the ferrite core 15 mm. This is sufficient to cover a frequency

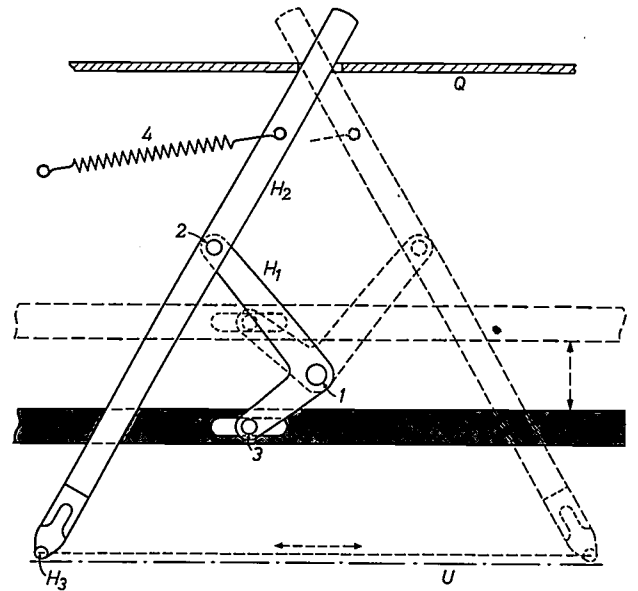


Fig. 8. Sketch of dial cursor mechanism. H_1 link with fixed pivoting point 1 ; the link carries at 3 a spindle which slides in a slot in J_1 (see *fig. 7*). H_2 strip connected with H_1 at point 2 . Q chassis plate. H_1 cursor. U dial. 4 spring.

illustrated in *fig. 9*. The push rod, which was shown in *fig. 7* for simplicity as a single rod, consists in reality of two parts capable of moving relative to one another, and denoted by E_1 and E_2 in *fig. 9*. The part E_1 , which carries the push-button D , can

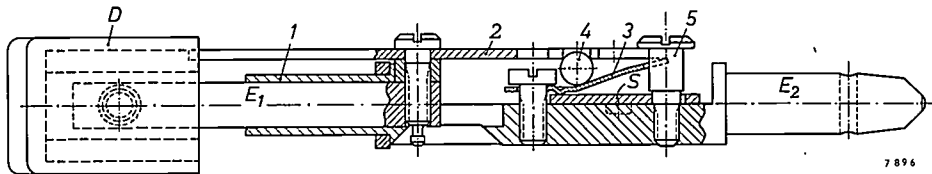


Fig. 9. Cross-section of push-button. E_1 and E_2 two parts of the push rod; E_1 slides inside the sleeve 1 of E_2 . 2 and 3 leaf springs. S segment. 4 steel ball. 5 screw. D push-button.

range of 1100 kc/s in the medium waveband. If we now specify that, when one of the push-buttons is depressed, the set should be tuned to the required station with a deviation of no more than 0.5 kc/s, then the average precision with which the cores must be displaced is $(0.5/1100) \times 15$ mm ≈ 7 μ m. This maximum permissible deviation is the sum of the deviations that may be due to misalignment of the segments S (see *fig. 7*) and the play in, and deformation of, the drive system. The components of this system must therefore be light and very rigid, and springs must be used to ensure that any play in the connections has no effect on the adjustment of the cores. Once the segment S has been fixed in position, any movement in relation to the push rod E must be prevented by very rigid clamping. The construction used for this purpose is

slide in a bushing 1 which is attached to E_2 . The segment S is pivoted on the screw 5 . *Fig. 9* represents the situation in which S can move freely. When E_1 is now moved to the right in relation to E_2 , a leaf spring 2 presses a hardened steel ball 4 to the right and causes a leaf spring 3 to clamp the edge of the segment S . Very considerable leverage is exerted by 3 as a result of its special profile. The force with which S is clamped is roughly 40 kg, which is sufficient to prevent any slip in the segment, even in frequent use.

To prevent displacement of the segment S in the pivoting point, the hole through which the screw 5 protrudes is not round but slightly V-shaped (see *fig. 10*). This ensures that there is no play at this position when the segment S is pressed against the cradle N (see *fig. 7*).

The precision with which the cradle returns to the same position when one of the push-buttons is depressed depends, among other things, on the torque required to move this component with the bar J_1 , the tuning cores and the cursor mechanism. This torque is determined to a large extent by the friction clutch by means of which the cradle shaft A' is connected to the wormwheel W_2 (see fig. 7). Apart from increasing the torque necessary to turn the cradle when the wormwheel is stationary, the friction clutch also causes the cradle to rebound slightly whenever a push-button is depressed, largely because the depression of the button produces slight torsion in the cradle. There may even be some rebound in the friction clutch itself. The consequence is an additional error in the reproducibility of the push-button tuning.

Operation of the tuning control K may also cause some rebound in the friction clutch, primarily due to vibrations to which the set is subjected after the tuning.

The degree of precision required in the alignment of the cradle appears from the following figures. The angle through which the cradle can turn is roughly 40° . In the medium-wave band this corresponds to a frequency coverage of 1100 kc/s. If we specify 0.2 kc/s as the maximum permissible detuning due to rebound, then the maximum angle over which the cradle may rebound is $(0.2/1100) \times 40^\circ = 26''$.

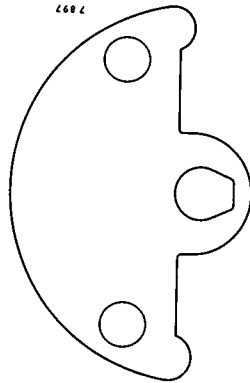


Fig. 10. Segment used in pushbutton mechanism (denoted by S in figs. 7 and 9). The central opening is slightly V-shaped to reduce play.

To minimize the chance of cradle rebound, a mechanism is used by means of which, whenever a push-button is depressed, the friction clutch between the wormwheel and the cradle shaft is put out of operation. For ease of illustration, this mechanism was omitted in fig. 7, but is drawn separately in fig. 11, together with the friction clutch in cross-section. The wormwheel W_2 can be turned on the shaft A' of the cradle N ; fixed to the same shaft is a diaphragm M of beryllium copper, which is pressed into a tapered recess in the wormwheel by a retainer 1 . The pressure is applied by the spring 4 via a lever 3 and a screw 2 . When one of the push-buttons is depressed, the sliding bar J_3 moves to the right (see fig. 7), as a result of which the lever 3 also moves to the right against the action of the

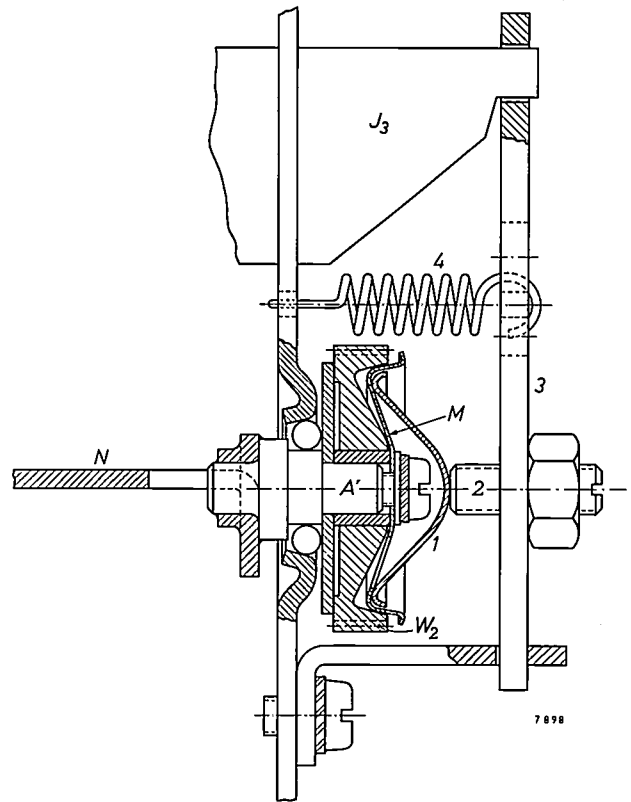


Fig. 11. Cross-section of the mechanism which disconnects the friction coupling between the tuning control and the cradle when one of the push-buttons is depressed. N cradle with shaft A' . W_2 wormwheel. M diaphragm. 1 retainer. 2 screw. 3 lever. 4 spring. J_3 sliding bar.

spring 4 . The pressure of the retainer 1 on the diaphragm M is then removed, thereby disconnecting the coupling between the wormwheel W_2 and the shaft A' . When the push-button is released, the coupling is again restored by the spring 4 .

Another factor influencing the precision of the tuning mechanism is that, when a push-button is depressed, the cradle is not only turned, but slightly bent. The bending produced by the force F_d exerted by a push-button is illustrated in fig. 12, exaggerated

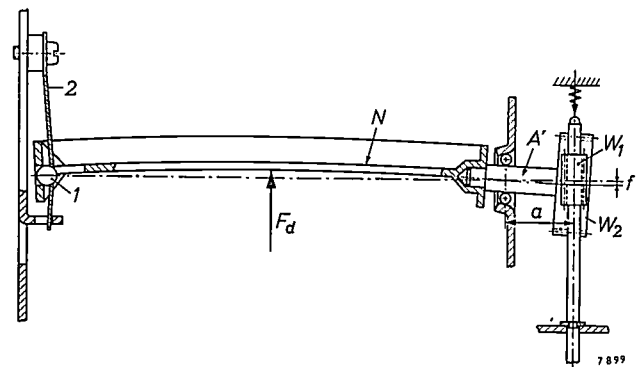


Fig. 12. Simplified sketch of the cradle N , bent under the force F_d exerted by one of the push-buttons. W_1 worm. W_2 wormwheel. 1 steel ball. 2 leaf spring. a distance from centre of wormwheel to centre of ball bearing. f displacement of wormwheel due to bending of the cradle.

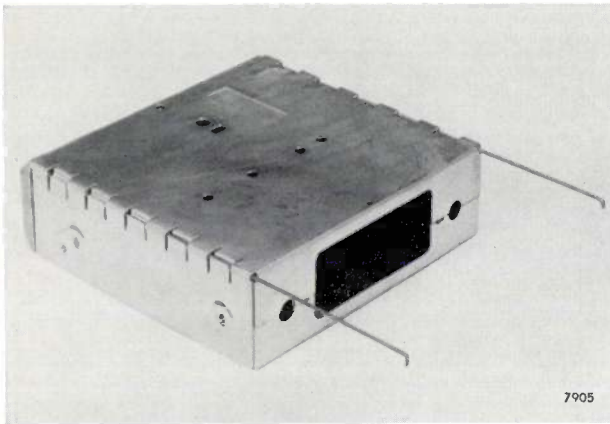


Fig. 13. The two parts of the metal housing are provided with lugs which are held together by two stainless-steel pins, thus ensuring good electrical screening.

for the sake of clarity. The effect disappears when the push-button is released. Since the wormwheel W_2 , which meshes with the stationary worm W_1 , also moves upon the bending and rebound of the cradle, the resultant displacement of wormwheel and cradle may cause an impermissible tuning error if suitable measures are not taken to prevent it. To give some idea of the amount of bending permissible, it may be mentioned that a displacement f of the wormwheel over a distance of $5 \mu\text{m}$ can cause a detuning of 0.6 kc/s .

The wormwheel displacement produced by a given degree of bending is less the smaller is the distance a between the centre of the right-hand ball bearing and the centre of the wormwheel. Among the measures therefore taken to reduce this distance to as little as 3 mm was to mount the right ball-bearing on the outside of the chassis and to provide the left-hand side of the cradle with a ball thrust bearing, in which the ball I is forced outwards by a leaf spring 2. The cradle was also given such a profile as to make it very rigid. These and similar measures made it possible to reduce the detuning in question to a permissible value.

The receiver housing

In designing the receiver housing, two important factors had to be taken into account. Firstly, the housing had to prevent the penetration of electrical interference from the engine. Secondly, provision had to be made for opening and closing the set quickly and easily. These requirements would be incompatible if one were to build a metal housing, the various parts of which were screwed together, for effective screening in that case would call for the use of numerous screws to ensure good electrical contact at many points over the parts of the housing. This, of course, would make it impossible to

open and close the set quickly. A satisfactory solution was found in a housing consisting of two parts each fitted with lugs in which two pins of stainless steel can be inserted (see fig. 13). The lugs act as contact springs, thus ensuring effective electrical screening.

The cooling block B (see fig. 4) is slid on to the bottom part of the housing at the rear, and is clamped between the two parts when the housing is closed. It is connected to the bottom part of the housing by only one screw, which in most cases is the same screw used for securing the set in the car. In such cases, the receiver housing is in fact "screwless" after removal from the dashboard.

The receiver is mounted in the dashboard by means of two heavy threaded bushings, fitted concentrically with the operating shafts. If the opening in the dashboard does not correspond to the dimensions of the set, a special ornamental plate can be used. If necessary the set can be additionally supported at the rear.

Fig. 14 shows a photograph of the complete receiver, and in the title photograph it can be seen mounted in the dashboard of a car. The set can also be mounted in the car in other ways; for example it can be mounted in a cover suspended under the dashboard.



Fig. 14. The complete car radio receiver N5X04T.

As the loudspeaker is not incorporated in the set, it can be set up in any appropriate part of the car. If required, two loudspeakers can be used, connected in parallel to the tapping available for that purpose on the output transformer.

Summary. After briefly considering the requirements to be met by a car radio, the authors review the development of these sets in the last twenty years. A description follows of one of the latest car radio receivers, type N5X04T, which is entirely equipped with semiconductor devices, viz. ten transistors and three germanium diodes. The circuitry and construction are discussed, with special emphasis on the tuning mechanism using push-buttons, which calls for very high precision to ensure reproducible tuning.

ABSTRACTS OF RECENT SCIENTIFIC PUBLICATIONS BY THE STAFF OF N.V. PHILIPS' GLOEILAMPENFABRIEKEN

Reprints of those papers not marked with an asterisk * can be obtained free of charge upon application to the Philips Research Laboratories, Eindhoven, Netherlands, where a limited number of reprints are available for distribution.

- 2954:** B. Jansen: A rapid and accurate method for measuring the thickness of diffused layers in silicon and germanium (Solid-state electronics 2, 14-17, 1961, No. 1).

Surface-layer thicknesses can be measured very practically by making use of the brittleness of silicon and germanium, thus eliminating cumbersome, precise grinding and polishing procedures. The germanium or silicon slice is broken in a special manner producing a rather flat cleavage surface, in which any disturbing fracture line is more or less perpendicular to the very sharp, long edges. In this surface the exposed *P-N* junctions are marked as very thin lines by one of the well-known methods (electrolytical, chemical). The distances between these lines or to the edge of the cleavage surface are measured under a high-power metallographic microscope with an eyepiece micrometer. Magnification factors mostly used are 400 or 600 times. Very thin layers of about 1 μm are still measurable. The method is used for diffused layers as well as for alloyed contacts or combinations as in the case of alloy diffusion.

- 2955:** J. Ubbink: Moderne ruisarme versterkers, I. Masers (Ingenieur 73, O1-O6, 1961, No. 3). (Modern low-noise amplifiers, I. Masers; in Dutch.)

An introduction to the maser (which gives Microwave Amplification by Stimulated Emission of Radiation), a new sort of low-noise amplifier. Three types of masers are discussed. The ammonia maser is not suitable for use as an amplifier, but it is an excellent frequency standard. Solid-state masers, which make use of the magnetic properties of e.g. ruby (Al_2O_3 containing a small percentage of Cr^{3+} ions), are better amplifiers. The ruby can be used in two ways: in a cavity resonator or in a waveguide. The first method has several disadvantages; the second is much better, if a system is used in which the group velocity of the microwaves is much smaller than in free space (a "slow-wave" structure). A practical execution of such a maser is described.

- 2956:** H. Mooijweer: Moderne ruisarme versterkers, II. Parametrische versterkers (Ingenieur 73,

O23-O32, 1961, No. 7). (Modern low-noise amplifiers, II. Parametric amplifiers; in Dutch.)

Discussion of the principle of operation of parametric amplifiers. The effective input noise temperature is higher than that of a maser (see 2955), but a parametric amplifier is simpler to make. Some actual amplifiers and some possible designs are described. The article contains extensive references to the literature.

- 2957:** M. J. Sparnaay and J. van Ruler: The adsorption of oxygen gas on germanium and surface conductivity (Physica 27, 153-162, 1961, No. 2).

A new method is described for studying electrical surface properties of semiconductors. The (very simple) method consists of resistivity measurements, at oxygen pressures ranging from 10^{-9} to 10^{-2} mm Hg, of thin germanium single crystals, the diameter of one sample varying from 1 mm to 10^{-2} mm. The crystals can be given the desired shape by "burning off" the germanium at 700 °C in an oxygen gas pressure of 10^{-2} mm Hg. By this method the roles played by surface conductivity and bulk conductivity can easily be separated. Results are given for intrinsic germanium and are in qualitative agreement with results obtained by Handler *et al.*

- 2958:** J. L. Meijering: On the thermodynamics of the Au-Pt system (Phys. Chem. Solids 18, 267-268, 1961, No. 2/3).

Comment on a publication by Weiss and Tauer, who claimed to be able to explain the asymmetry of the phase diagram of the system Au-Pt from measurements of the specific heat only. It is shown that the phase diagram calculated by the method of Weiss and Tauer in fact differs in two essential points from that found experimentally.

- 2959:** C. Haas: Vibratiespectra van kristallen (Ned. T. Natuurk. 27, 105-118, 1961, No. 4). (Vibration spectra of crystals; in Dutch.)

Data about vibrations in crystals obtained by measurement of specific heat and thermal conductivity are difficult to interpret. This is because all the

vibrations in the crystal contribute to these quantities. This article discusses two methods whereby more selective information about the lattice vibration can be obtained: infrared spectroscopy and Raman-spectroscopy.

2960: A. Schmitz: De tunneldiode (Ned. T. Natuurk. 27, 133-142, 1961, No. 4). (The tunnel diode; in Dutch.)

Esaki published his article on the tunnel diode in Jan. 1958. The most important property of this diode is a part of the I - V characteristic with a negative slope. Various workers have studied the properties and applications of this diode. This article gives a survey of the present state of knowledge in this field.

2961: M. J. Sparnaay: Elektrische Doppelschichten (Original lectures IIIrd Int. Congr. Surf. Act., Cologne, Sept. 1960, Vol. II, pp. 232-253, Universitätsdruckerei Mainz GmbH). (Electrical double layers; in German.)

An electrical double layer is often formed at the boundary between two phases which contain a sufficient number of free charge carriers. This phenomenon is very important in colloidal systems, but also in semiconductors, as regards both surface effects and P - N junctions. In this article a survey is given of the theory of these double layers for various cases; the points of similarity and dissimilarity between the different cases are pointed out.

2962: J. Bergsma, J. A. Goedkoop and J. H. N. van Vucht: Neutron diffraction investigation of solid solutions $AlTh_2D_n$ (Acta crystallogr. 14, 223-228, 1961, No. 3).

$AlTh_2$ absorbs hydrogen readily. This article describes an investigation, carried out with the aid of neutron diffraction, of the manner of incorporation of the hydrogen in the $AlTh_2$ lattice. In fact, deuterium was used instead of hydrogen, in connection with the demands made by neutron-diffraction techniques. The investigation was restricted to solid solutions of composition $AlTh_2D_n$, with $n = 0, 2, 3$ and 4 . It was found that the hydrogen is taken up in tetrahedral interstices between thorium atoms, just as in thorium which contains no aluminium. If $n = 4$, all the available tetrahedral interstices are filled. At lower values of n , i.e. when the filling is incomplete, no order could be found in the distribution of the deuterium atoms over the available sites, even at a temperature of $82^\circ K$. See also Philips tech. Rev. 23, 69, 1961/62 (No. 3).

2963: P. Massini: Self-absorption correction for isotopes emitting weak beta rays (Science 133, 877-878, 1961, No. 3456).

If it is desired to determine the radiation intensity per mg of a radioactive sample from measurements of the radiation actually emitted, a correction must be applied for the absorption by the sample itself. The correction factor is often given as a function of the mass of the sample. There is disagreement in the literature as to the form of this correction curve for β -emitters. The author shows that this disagreement may be due to differences in the geometry of the sample and the measuring set-up.

2964: H. C. Hamaker: Examples of designed experiments (Industr. Qual. Control 17, No. 9, 16-20, 1961).

Textbooks on the application of statistics to experimental investigations usually give examples where the design of the experiment is taken as given, and the experimental results are successfully processed with the aid of standard statistical methods such as variance analysis. In practice, the situation is often different. This publication describes two cases from industrial practice where the help rendered by the statistician consisted in analysing the problem and designing a suitable experiment or series of experiments on the basis of this analysis. The results of these experiments spoke for themselves: there was no need to subject them to any form of statistical analysis.

This aspect of the task of a statistician is often neglected in textbooks of statistics. See also Philips tech. Rev. 22, 105, 1960/61.

2965: P. C. van der Willigen: Booglasmethodes voor staal en de rol die de waterstof daarbij speelt (Chem. Weekbl. 57, 170-176, 1961, No. 14). (Methods of arc welding of steel, and the role of hydrogen in such methods; in Dutch.)

A short survey of the historical development of arc-welding methods for steel. Welding electrodes with an "acid" coating are still used on an immense scale. A characteristic of the acid coating is that it produces much hydrogen during welding. Hydrogen has come to be regarded as an enemy of good welding, for reasons which are stated in this article. Hydrogen-free welding methods (submerged-arc welding, use of electrodes with basic coatings, welding with a bare wire in a protective atmosphere of argon or CO_2) are therefore gaining more and more ground. Hydrogen in steel is also one of the causes of the cracking of enamel.

Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES

MECHANIZED MOUNTING OF COMPONENTS ON PRINTED-WIRING PANELS

by R. van BEEK *) and A. J. HALBMEYER *).

621.397.62.049.75.002.72

In recent years increasing use has been made of printed-wiring panels in electronic equipment. Initially the components were mounted (inserted and fixed) on these panels largely by hand. The mounting of many components has now been entirely mechanized. At Philips in Eindhoven a machine has been in use for some time which mounts components automatically on printed-wiring panels for use in television sets.

I. GENERAL CONSIDERATIONS ON MECHANICAL MOUNTING

Introduction

In the printed-wiring technique boards of insulating material are used, on one side of which the required pattern of conducting connections is formed from a thin layer of copper. On the other side of these *panels*, as they are called, are placed the circuit components (resistors, capacitors, inductors, valve holders, etc.) with their terminal wires or tags inserted into holes in the panel, and brought into good electrical contact with the printed wiring by dip soldering. The way in which this was done in the first Philips radio sets provided with printed wiring was the subject of an earlier article in this journal ¹⁾.

The mounting of the components, although at that time already mechanized to some extent, was largely manual. Insertion of the components in a flat panel, however, lends itself very well to mechanization, for each part of the flat surface is readily accessible to insertion tools. In the Eindhoven factories of Philips, components are now mounted entirely automatically on panels for television sets. The method will be described in this article.

The mounting procedure in itself is simple: it may consist in taking a component from a magazine,

bending the terminal wires into the appropriate shape, inserting them into specific holes in the panel, and cutting and bending the wires under the panel. The time needed for these operations varies from about $\frac{1}{2}$ to $1\frac{1}{2}$ seconds. In some cases the process can be even simpler, for example if the component has terminal pins or tags that fit into a particular pattern of holes; in that case the pre-bending and cutting operations are dispensed with.

Before the mounting procedure can begin, the panel must be properly aligned in relation to the insertion tool, and after the component has been fitted the assembly of panel and component has to be carried away. Together, these operations also take very little time, roughly the same as that mentioned above.

Two methods of mechanized mounting

As long as no components have yet been mounted, the panels can simply be stacked one on top of the other, and their conveyance to the assembly point is easily mechanized. Once the first component has been fixed to the panel, however, the situation changes entirely: the previously simple panel is now a vulnerable and not easily stacked assembly, which requires careful handling if damage is to be avoided. For this reason once the insertion of components has started it is best to continue until all mechanically

*) Radio, Television and Record-playing Apparatus Division, Philips, Eindhoven.

¹⁾ R. van Beek and W. W. Boelens, Printed wiring in radio sets, Philips tech. Rev. 20, 113-121, 1958/59.

insertable components have been secured to the panel, thus avoiding the formation of stocks of partly assembled panels. This can be done in two ways:

- 1) After the first insertion tool has fitted the first component, the panel is immediately passed to a second insertion tool that fits the next component, after which it passes to the third insertion tool, and so on until all components are in position.
- 2) After the insertion tool has fitted the first component, the same tool fits the succeeding components. For this purpose it is necessary to shift and/or turn the panel in such a way that each component is fitted in the right position.

The second procedure has the following drawbacks compared with the first:

- a) For each insertion operation the panel has to be brought into correct alignment in accordance with three data: two in respect of position and one in respect of direction.
- b) Either the designer must limit his choice to components of one particular size, or the insertion tools must be reset after each operation (i.e. adapted to components of different sizes).

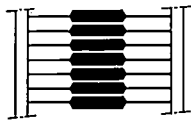

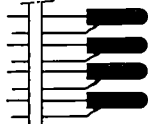


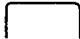






this is the form which we have chosen. As opposed to assembly on a turntable, this system can easily be extended, which has proved very useful in practice.

Subdivision of components according to shape

Before dealing with the fully mechanical process of insertion and attachment, it is necessary to consider the *form* of the components to be mounted. Not all available components are suitable for mechanical insertion. The main forms of mechanically insertable components are represented in *Table I*. Machine tools (insertion heads) have been developed for the automatic insertion and fixing of all these main forms of components. The construction of these heads is dealt with in Part II and photographs of them can be seen in figures 5, 7, 8 and 9.

Components with axial terminal wires (column AW in Table I), such as resistors, are inserted by an AW head. They are fed to the head in the form of a tape with the ends of the wires stuck to two strips. Bridging wires (column BW in Table I) serve for interconnecting certain points of the printed wiring. A special BW head cuts them from a roll of bare wire and fits them to the panel. (Since bridging wires

Table I. Main forms of mechanically insertable components.

Name	Components with axial wires	Bridging wires (jumpers)	Pin-up components	Valve holders
Symbol	AW	BW	PU	VH
As received from supplier				
Prepared for insertion by machine				
Fitted to panel				

7942

If the procedure under (1) is adopted, the above-mentioned alignment and resetting are necessary only once in a large series of insertion operations. The installation can therefore be much simpler than in procedure (2).

Conveyance of the panels to the insertion tools can be effected either on a band or chain conveyor, or by means of a turntable around which the tools are arranged. The first case results in an *assembly line*, and

can be regarded as an extreme form of AW component, we later decided in certain cases to feed them in the same way as AW components and insert them with an AW head; for the sake of simplicity, in Part I of this article we shall treat bridging wires as AW components.)

An entirely different head construction is needed for inserting standing or "pin-up" components, such as the ceramic capacitor shown in column PU in

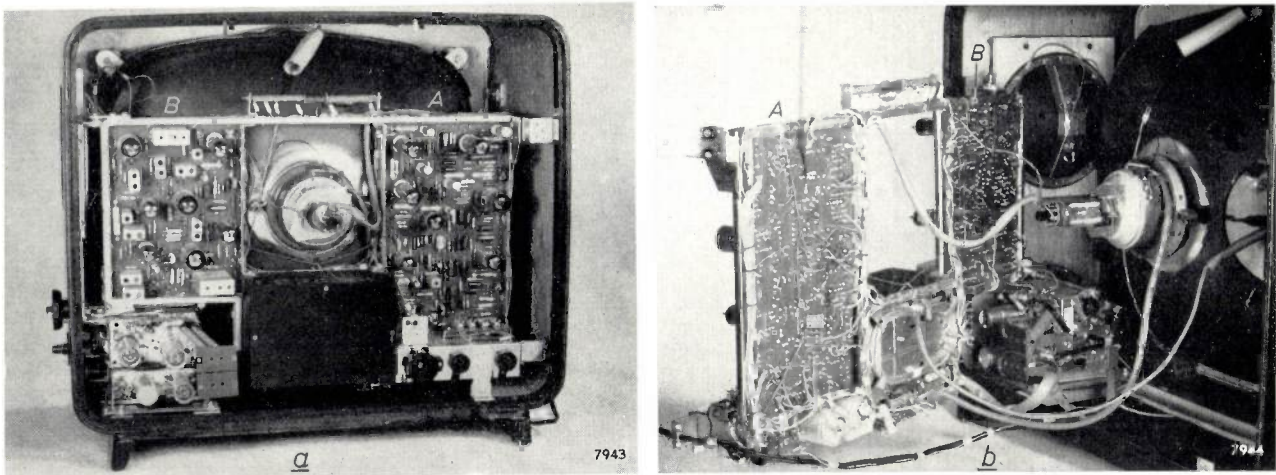


Fig. 1. Television set type 21 TX 310A. It contains two printed wiring panels, *A* and *B*, on which most components have been mounted mechanically. The component side of the panels can be seen in (a), and in (b) the side with the printed wiring and soldered joints.

Table I. The same applies to the head developed for mounting valve holders (column VH in Table I).

Distribution of the panel components

We shall take as an example a particular television set. This set contains two printed wiring panels (fig. 1). Table II gives a survey of the numbers of the various components used, divided into those which are and those which are not suitable for mechanical insertion. The AW components are subdivided here into non-raised and raised components (i.e. mounted in contact with the panel or at some distance from the surface).

If an assembly line were set up for each of the two types of panel, with a separate insertion head for each component, each line would require about 100 assembly stations (102 for panels *A* and 92 for panels *B*; see Table II). As will presently be shown, assembly lines of this length cannot be kept in operation, owing to frequent malfunctioning. Even if this were not the case, this method would still have the drawback of a much too great a production capacity; each of the lines would be able to fit about 1000 panels an hour with components, so that the panel

production capacity would correspond to the production of about 1000 television sets an hour — a number which, even in the largest plants, is too high to make the installation an economic proposition.

To reduce this enormous production capacity, one could assemble on a single line series of *A* and *B* panels alternately. The panel production would then correspond to about $1000/2 = 500$ television sets an hour, which is still an exceptionally high number. Moreover, the line would then have even more assembly stations than before, i.e. about 130 (89 AW + 32 PU + 6 VH stations, see Table II) and would therefore be impracticable.

In order to shorten the assembly lines one might consider subdividing the panels, e.g. panel *A* into three and panel *B* into two smaller panels with the mechanically mountable components divided, for instance, as shown in Table III. The production capacity would then be sufficient for producing about $1000/(3 + 2) = 200$ television sets an hour. This figure would still be excessive for a small factory, but could make the installation profitable in a large factory.

Table II. Division of components of two printed-wiring panels for the television set type 21 TX 310A.

Type of component	Panel A					Panel B					$\Sigma A + \Sigma B$
	AW	AW raised	PU	VH	ΣA	AW	AW raised	PU	VH	ΣB	
Mechanically insertable	89	0	8	5	102	54	0	32	6	92	194
Not mechanically insertable	38	14	11	3	66	22	16	9	1	48	114
Total	127	14	19	8	168	76	16	41	7	140	308

Although all five types of sub-panels can now be assembled on a single assembly line with about 50 stations (30 AW + 16 PU + 3 VH stations, see Table III), which is a lot less than just mentioned, the number is still too high. Furthermore, it must be remembered that subdivision of a panel is costly, for not only does it multiply numerous operations (sawing-off, hole-punching, dip-soldering, setting-up in mounting and inspection machines) but in addition connections have to be provided between the sub-panels themselves, and each sub-panel demands certain measures for fixing it in the set. The costs which all this involves outweigh the economies obtained from mechanical mounting.

Table III. Division of panels A and B of Table I into three and two sub-panels respectively, and the resultant distribution of components.

	Type of component		
	AW	PU	VH
$\frac{1}{3}$ panel A	$\frac{89}{3} \approx 30$	$\frac{8}{3} \approx 3$	$\frac{5}{3} \approx 2$
$\frac{1}{2}$ panel B	$\frac{54}{2} = 27$	$\frac{32}{2} = 16$	$\frac{6}{2} = 3$
Number of assembly positions	30	16	3

To enable components to be mounted on a single panel by means of an assembly line which is not too long and has no excessive production capacity and in the numbers given in Table II, we have adopted the following procedure: *the panel is passed several times along the same, relatively short assembly line, the insertion heads of which are reset before mounting each successive set of components.*

In the construction of a fully mechanized assembly line on this principle, we must take into account the fact already mentioned that partly assembled panels are vulnerable and awkward to handle. This difficulty is met by enclosing each panel, before assembly, in a protective holder (a frame structure — see fig. 10 — which will be discussed in Part II). This also prevents the laminated panels from becoming worn after repeated setting-up in successive assembly stations (more than a hundred times), which would make accurate centring impossible. A consequence of the system using holders is that between each two successive passes the holders have to be returned from the end of the assembly line to the beginning and positioned there in a suitable way.

Optimum arrangement of assembly line

Fig. 11 will give the reader some idea of the arrangement of the assembly line in use at Philips.

The figure shows the band conveyor 1 which carries along the panels held in the holders 2 to under the insertion heads 3. The band conveyor 5 returns the holders with uncompleted (in this case blank) panels to the beginning of the line. At 7 the fully assembled panels are replaced by blank ones.

The number of assembly stations the line should have and the number of holders in circulation, in order to keep the invested capital as low as possible in proportion to the production capacity, can be determined from:

- 1) the cost of the assembly line and the holders,
- 2) the time taken to reset the insertion heads for each successive set of components,
- 3) the cycle time *c* taken for fitting and transporting one component.

It will be useful to consider the production capacity first.

Production capacity

The theoretical production capacity (the theoretical maximum number of assembly operations per hour) of each assembly station is denoted by *p_{th}*. If the cycle time for one assembly operation, including the transport of the panel to the next station, is *c* seconds, we should have *p_{th}* = 3600/*c* operations per hour and per station, if it were not for the fact that the machine between each two passes of the series of *m* panels is unproductive for a certain time (*s* seconds). Since a pass of *m* panels takes *mc* seconds, we have

$$p_{th} = \frac{mc}{mc + s} \times \frac{3600}{c} = \frac{3600}{c + \frac{s}{m}} \dots (1)$$

The real or effective production capacity per assembly station, *p_{eff}*, is smaller than the theoretical value:

$$p_{eff} = \eta_b \eta_s p_{th}, \dots (2)$$

where the factors η_b and η_s , which are both less than unity, are of different origin as is shown in the following:

- 1) The number of components mounted during a single pass may be smaller than the number of stations available for that type of component. Some insertion heads are therefore idle during the pass. This is represented by the factor η_b , which thus depends on the number of assembly stations in the line and their function. The value of η_b can be made to approach unity by making the fullest possible use of the capacity of the equipment; consequently, various components that would be suitable for mechanical mounting must then be

mounted manually. A formula for η_b will be derived presently.

- 2) The factor η_s represents the production losses due to failures in operation or malfunctioning of assembly stations. If one or more of the n assembly stations is out of operation for a fraction a of the time, and thus holds up the remainder of the assembly line, then

$$\eta_s = (1 - a)^n \dots \dots \dots (3)$$

Table IV gives the values of η_s (in %) for various values of a (in %) and n . It can be seen that the

Table IV. The factor η_s , in %, calculated from eq. (3), as a function of the number of assembly stations n and of the average fraction a (in %) of the time that an insertion head holds up production.

a %	Number of assembly stations n							
	1	2	4	8	16	32	64	128
0.0	100	100	100	100	100	100	100	100
0.5	99.5	99.0	98.1	96.1	92.3	85.2	72.5	52.6
1.0	99.0	98.1	96.1	92.3	85.2	72.5	52.6	27.7
1.5	98.5	97.0	94.1	88.6	78.5	61.7	38.0	14.5
2.0	98.0	96.0	92.2	85.1	72.4	52.2	27.2	7.4
2.5	97.5	95.1	90.4	81.7	66.7	44.5	19.8	3.9

result of a disturbance lasting, say, 1.5% of the time is more serious the greater is the number n of assembly stations. Long assembly lines are thus not suitable for practical purposes, as mentioned earlier.

From (1) and (2) it follows that

$$p_{\text{eff}} = \eta_b \eta_s \frac{3600}{c + \frac{s}{m}} \dots \dots \dots (4)$$

The assembly line in operation at Philips has a cycle time c of 3 seconds, an unproductive time s per pass of 60 seconds, and the number of circulating panel holders m is 100. It therefore follows from (4) that in the ideal case ($\eta_b = 1$), best use is made of capacity when:

$$p_{\text{eff}} = \eta_s \frac{3600}{3.60} = 1000 \eta_s \text{ insertions}$$

per hour and per assembly station; for this assembly line, then, p_{eff} is numerically identical with $10 \times \eta_s$ in % and can thus be read directly from Table IV.

Choice of the number of holders and assembly stations

Let us now return to our object, which is to choose the number of circulating holders m and the number of assembly stations n in such a way as to minimize the cost of the assembly line in relation to the production capacity.

An assembly line is composed of the following elements:

- 1) A number n (to be decided) of assembly stations, coupled together by a band conveyor. The cost of one station, together with the appertaining support for the band conveyor, will be denoted by N .
 - 2) An installation designed to return the holders with incomplete or blank panels to the beginning of the line and also to keep a reserve of holders (this installation comprises — see fig. 11 — a slide 4, a band conveyor 5 and a lift mechanism 6), and a device for removing the fully assembled panels from the holders and replacing them by blank ones (see 7 in fig. 11). The cost of these two installations together will be called K .
 - 3) A number of holders m (to be decided). The cost of each holder will be called M .
- The total cost I of the assembly line is thus

$$I = nN + K + mM.$$

The cost i_s per assembly station, in relation to the effective production capacity per assembly station, is therefore:

$$i_s = \frac{I}{np_{\text{eff}}} = \frac{nN + K + mM}{n} \frac{1}{\eta_b \eta_s} \frac{c + \frac{s}{m}}{3600} \dots \dots \dots (5)$$

By calculating $\partial i_s / \partial m$ from (5) and equating to zero, we find the value m_0 of m at which i_s is lowest. This value further depends on n :

$$m_0 = \sqrt{\frac{s}{c} \frac{nN + K}{M}} \dots \dots \dots (6)$$

In practice it is found that N , K and M are in the ratio of 100 : 150 : 2.2. Using these figures, and putting $s = 60$ seconds and $c = 3$ seconds, eq. (6) can be written

$$m_0 = 30.2 \sqrt{n + 1.5} \dots \dots \dots (7)$$

Fig. 2 shows a plot of i_s versus m , where $\eta_b = 1$ and $\eta_s = 1$, for various numbers n of assembly stations. The minimum of the curves is so flat that it is of little consequence as far as i_s is concerned if, for practical reasons, we choose a round number for m and thus deviate slightly from the optimum value m_0 .

Fig. 3 shows i_s as a function of the number of assembly stations n for various values of the number of holders m and of the stoppage percentage a (i.e. the time loss due to a faulty head), in the case where $\eta_b = 1$. The marked influence of the stoppage percentage is evident. *The feasibility of mechanized assembly therefore depends upon the possibility of keeping the stoppage percentage very low.* The extent to which

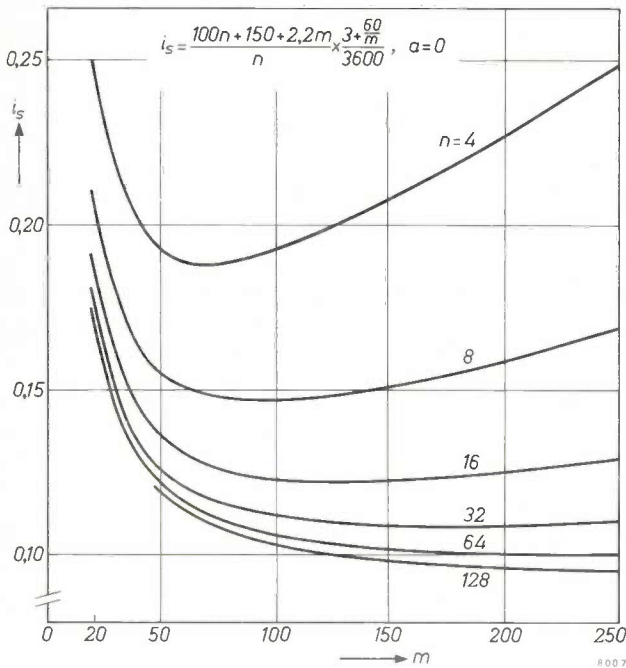


Fig. 2. The price i_s of an assembly line per assembly station, in relation to the effective production capacity per station, as a function of the number of circulating holders m , with the number of stations n as parameter, for the case $\eta_s = 1$ ($a = 0$) and $\eta_b = 1$.

this can be done depends largely on the extent of the deviations of the components from the prescribed shape; bent or eccentrically positioned terminal wires, for example, have an adverse influence. If the stoppage percentage can be reduced to 1%, an assembly line with 16 stations and 100 holders is a good proposition. The point corresponding to this in fig. 3 is marked with a circle ²⁾.

For a stoppage percentage of 2% and $n = 16$, we see from Table IV that η_s is equal to 72.4 %, so that (for $\eta_b = 1$) $p_{\text{eff}} = 724$. Given these conditions an assembly line with 16 stations can thus mount $16 \times 724 = 11\,584$ components per hour.

We shall now consider the case where the line is equipped for mounting AW components, of which there are 89 on panel A in the above-mentioned television set and 54 on panel B (Table II). Optimum use is made of the capacity of the line if respectively 80 and 48 of these are mechanically mounted (the nearest multiples of $n = 16$ less than 89 and 54); for each television set we thus have $80 + 48 = 128$. The production of the line is then sufficient for an average

of $11\,584/128 \approx 90$ sets per hour. If more components suitable for mechanical mounting become available, the number of mechanically mountable components can be increased, e.g. to 128 on the one panel and to 80 on the other. The line will then be fully occupied producing panels for an average of $11\,584/(128 + 80) \approx 55$ sets per hour, and will thus also be suitable for smaller factories. This demonstrates the great importance of mechanically mounting as many components as possible.

Increasing the number of mechanically mounted components to the maximum, however, means that the factor η_b , which was assumed to be unity, drops to less than 1. In the passes made by a panel along the line there will always be one in which not all insertion heads mount components (unless the number of components to be mounted per panel is an exact multiple of n). In the case just considered, with $n = 16$ and 89 AW components on one panel, all 16 heads are in operation during 5 passes; in one pass, however, only $89 - 5 \times 16 = 9$ heads are in use, and thus 7 are idle. The number of heads that remain idle during this pass, averaged over a sufficient number of different panels, is $\frac{1}{2}(n - 1)$. If an average of \bar{C} components per panel are mechanically mounted, the number of passes being equal, the theoretical mounting capacity is thus $\bar{C} + \frac{1}{2}(n - 1)$, where \bar{C} is the useful part and $\frac{1}{2}(n - 1)$ the non-useful part. It follows then that

$$\eta_b = \frac{\bar{C}}{\bar{C} + \frac{1}{2}(n - 1)} = \frac{1}{1 + \frac{n - 1}{2\bar{C}}} \dots (8)$$

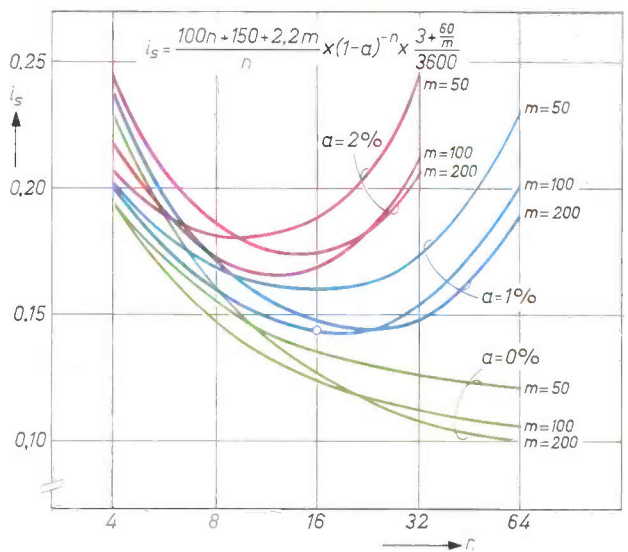


Fig. 3. Plot of i_s versus n , with m as parameter, for a percentage time loss a due to malfunctioning: $a = 0$ (green curves), $a = 1\%$ (blue curves) and $a = 2\%$ (red curves).

²⁾ The number of stations was finally brought to 24, as described in Part II. If a were in fact 1%, one line with 24 stations would be rather more advantageous than two lines with 12 stations. Practice has since proved that the estimate of $a = 1\%$ was on the optimistic side. At $a = 1.5\%$ the line with $n = 24$ has the same production capacity as two lines, one with $n = 9$ and the other with $n = 10$. The latter combination also has the advantage of greater flexibility.

(N.B. This formula only applies if there are a sufficient number of different panels.)

Fig. 4 shows a plot of η_b versus the number of stations n for various values of the average number of mechanically mounted components \bar{C} per panel.

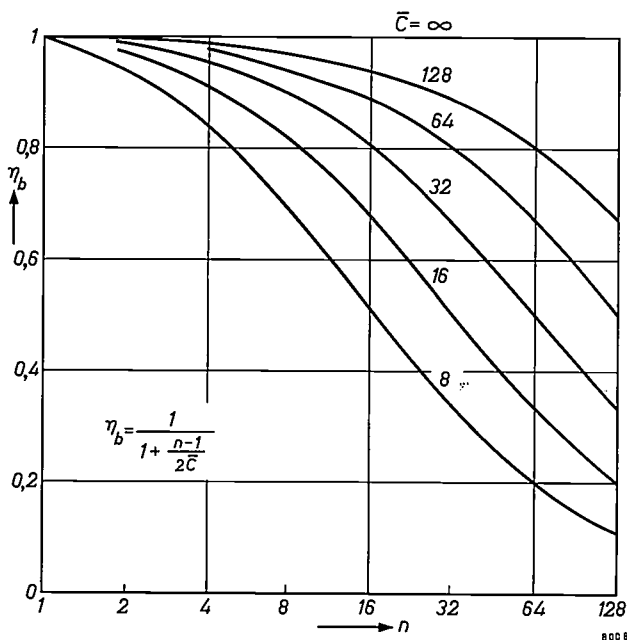


Fig. 4. Machine occupation factor η_b as a function of the number of assembly stations n , the parameter being the average number of mechanically mountable components \bar{C} per panel.

If all mechanically mountable components in the above television set are in fact so mounted — i.e. 89 + 54 components in 6 + 4 = 10 passes, instead of 80 + 48 components in 5 + 3 = 8 passes — we find

$$\eta_b = \frac{89 + 54}{16 \times 10} = 0.89.$$

The average production would now be reduced to

$$0.89 \times \frac{11584}{89 + 54} \approx 72 \text{ sets per hour.}$$

Manual assistance of assembly line

Operators are needed to attend the assembly line. Their jobs are:

- 1) to supply the machine with the components to be mounted,
- 2) to programme the machine,
- 3) to reset the machine periodically,
- 4) to remove fully assembled panels and introduce blank ones,
- 5) to repair faults.

The total time needed for these jobs determines the number of operators. The individual times can best be studied by considering the time taken for a pass of m holders along the line.

Job 1. For every pass of m holders a new batch of m components must be fed to each insertion head. This can be done while the machine is in operation. The time t_1 needed for this operation is constant for every insertion head, e.g. 10 seconds.

Job 2. The machine has to be repeatedly programmed with the coordinates of the place where each insertion head has to fit its component in the next pass of the panels. This programming takes a time t_2 for every pass of m panels, varying from 20 to 28 seconds for each insertion head. It is possible to make t_2 (like t_1) fall entirely within the operating period of the machine, so that it does not contribute to the above-mentioned time s during which the machine is unproductive between successive passes of m panels.

Job 3. The resetting of an insertion head takes a certain amount of time. For each pass of m holders these times together constitute the resetting period t_3 of the machine; t_3 varies from 2.3 to 12 seconds. If only one operator were available for resetting each head, t_3 would be equal to the unproductive time s .

Job 4. The time needed to replace each fully assembled panel by a blank one is constant and adapted to the rhythm of the machine; the cycle time c for removing a completed panel, plus an equal length of time c for introducing a blank panel. If a panel, before completion, has to pass along the line r times, the replacement time t_4 per pass of m panels is $2mc/r$. When $m = 100$, $c = 3$ seconds, and $r = 10$, then $t_4 = 60$ seconds.

Job 5. The operators responsible for the above four jobs are also required to repair normal faults. The time taken for this is obtained by multiplying the sum of the above-mentioned times by $1/\eta_s$.

We thus arrive at a total time for the assistance of the line per pass of m panels, which is given by $(t_1 + t_2 + t_3 + t_4)/\eta_s$. Table IV has shown that η_s decreases faster with increasing stoppage percentage a the greater is n and thus the longer is the line. The desire to use as few operators as possible is therefore another argument for keeping the assembly lines short. A line with 24 stations and 100 circulating holders requires no more than 5 operators.

In view of the large number of components C per panel, it is evident that even if only a very small percentage of the mounting operations are failures, the result will still be that a substantial proportion of the completed panels will contain faults. Let u be the fraction of failures in the total number of mounting operations, then the fraction G of the faultlessly assembled panels is:

$$G = (1-u)^C.$$

If, for example, $C = 40$ and 80 , and $u = 0.5\%$ and 1% , we obtain the following figures:

	$u = 0.5\%$	$u = 1\%$
$C = 40$	$G = 82\%$	$G = 67\%$
$C = 80$	$G = 67\%$	$G = 45\%$

These figures indicate that it is not sufficient to take random samples when inspecting the panels; each completed panel must be separately inspected and, if necessary, repaired.

Conclusions

The following conclusions may be drawn from the above discussion.

Construction of the assembly line

- 1) To keep down investment costs it is of the utmost importance that the insertion and resetting times should be short.
- 2) To keep the resetting time short the manipulations on the machine between operations should be effected as quickly as possible. These manipulations are: changing magazines with components, and deciding the insertion points for the component on the panel.
- 3) Stoppages cause the production capacity per

assembly station to decrease as the number of stations n increases. The cost per station is a minimum at a specific value of n .

- 4) The fewer assembly positions there are, the better the use that can be made of capacity in a varying programme.

The component batch

- 1) To limit stoppages due to faults, the components to be used should have a shape and quality such that they cannot damage or soil the machine (e.g. due to paint scrapings settling in the form of dust, or caked together with oil or wax, between moving parts).
- 2) To make good use of available capacity, as many as possible of the various kinds of components should be of the same shape.
- 3) To reduce manual operations the component should be systematically packed and supplied in a form suitable for mechanical handling.

Design of the panel

- 1) As many as possible of the components to be mounted on one panel should be suitable for mechanical fitting.
- 2) The components chosen by the designer should as far as possible be geometrically similar (or better still, congruent).

II. MECHANICAL CONSTRUCTION

Requirements

The requirements imposed on the construction of a mechanical assembly line have been touched upon in Part I. Examining these requirements in more detail, we see that the design engineer is faced with the following tasks:

- 1) To design universal insertion heads for mounting different components.
- 2) To design a universal holder capable of holding panels of different sizes, carrying them to the appropriate place under the insertion heads and fixing them there, and protecting them from damage during conveyance.
- 3) To design a conveyor mechanism which keeps the holders in circulation.
- 4) To devise a system which minimizes the time required for successively resetting the insertion heads.

- 5) To develop a drive mechanism which safeguards the machine against damage that might be caused by the failure of one or more insertion heads.

In the following paragraphs we shall examine the manner in which the design engineer tackled these assignments and the solutions adopted for various problems. First of all we shall give a general description of the layout and operation of the assembly line in use at Eindhoven for the mechanical assembly of panels for television sets.

Operation of assembly line

A simplified diagram of the assembly line is given in fig. 11. Above a band conveyor 1 which carries the panel holders 2, there are a number of insertion heads 3 which mount the components on the panels. At present there are 24 insertion heads, viz:

18 AW, 2 BW, 2 VH and 2 PU heads.

Since a panel contains a lot more than 24 components, each panel must pass several times through the machine. For this purpose the assembly line is organized as follows.

100 holders circulate constantly through the machine. They are carried under the insertion heads horizontally (lying flat) by a twin-belt conveyor 1. At the end of this first conveyor the holders are conducted via a slide 4 to a lower steel band conveyor 5. The holders are thereby turned to the vertical position and returned by the conveyor 5 to the beginning of the assembly line. Here there is a lift device 6, which lifts the holders from conveyor 5, turns them back to the horizontal position and deposits them on the band conveyor 1.

Each holder contains the same kind of panel, and these 100 panels are now fitted with the first series of 24 components. This being done the insertion heads are quickly reset, a new supply of components (100) is introduced, and the second series of 24 components are mounted. When all mechanically mountable components are in place, the panels are replaced by 100 blank ones.

The reasons why there are fewer insertion heads than there are components per panel have been given in Part I: if there were an insertion head for each component the line would be much too long and a fault could then have serious consequences; moreover the production (e.g. 2 million panels a year) would by far exceed the production required.

For replacing the completely assembled panels by blank ones, an unloading station is fitted at the end of the assembly line (at 7 in fig. 11). The completed panels are removed from the holders automatically, but the blank ones are inserted by hand. Manual insertion is necessary to ensure that the notches in the reference face (which we shall deal with presently) fit over the appropriate pins. If this were not so, then the holder would be unable to bring the panel into proper alignment under the insertion mechanism of the heads, resulting in faulty mounting of the relevant component. Provisions have been made to make the time needed to fit a panel in an empty holder equal to the time in which a head mounts components. In our assembly line this time is 2 sec.

The time taken to mount a component and then convey the holder to the next head is 3 seconds (the cycle time c). The machine thus works automatically for 100×3 seconds = 5 minutes, and as soon as the 100th holder has left the first insertion head, it switches that head off, enabling the operators to start resetting the heads for the next pass of the holder. As soon as a head has been reset, it is again switched on. Once all heads have been successively reset and

the last one switched on again, the assembly line resumes its automatic operation.

The assembly line components

The insertion heads

In Part I it was mentioned that four kinds of insertion head have been developed. i.e. for components with axial terminal wires (AW), for bridging wires (BW), for valve holders (VH) and for pin-up components (PU).

The *AW* head is shown in fig. 5. The operation of the insertion mechanism is represented schemat-

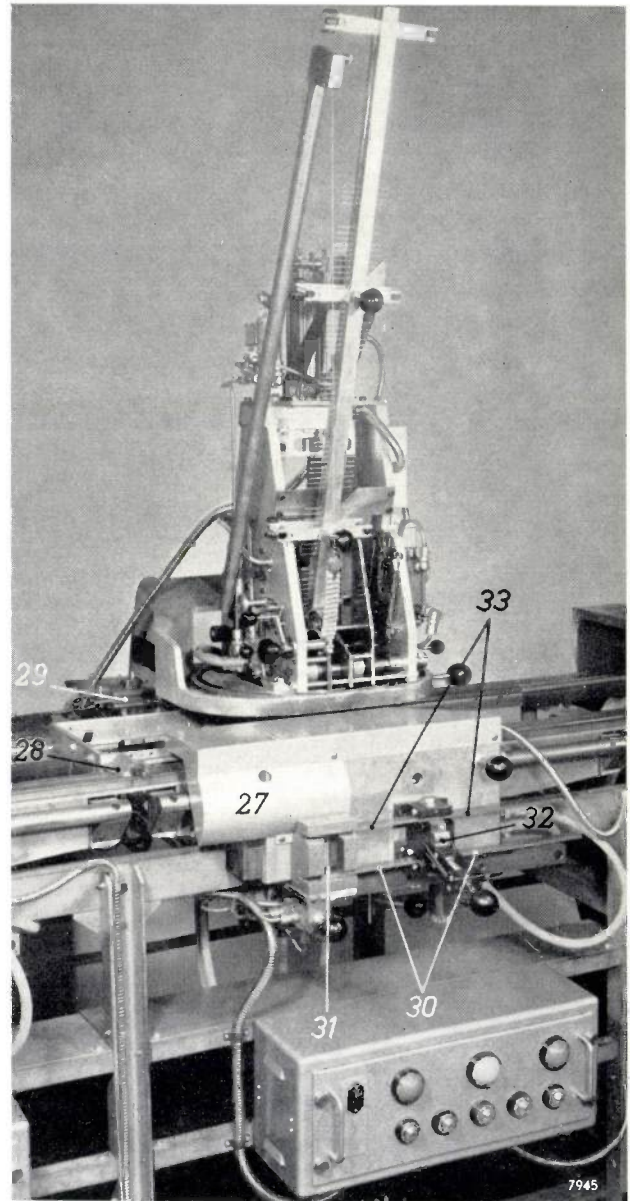


Fig. 5. Insertion head for components with axial terminal wires (AW head). 27 support. 28 pin on holder. The bar 29 forces the holder with its studs against two pins in the support. 30 and 33 interchangeable setting bars with holes that determine the x coordinate of the component to be inserted in the present and in the next pass respectively. 31 fixed pin. 32 locking device.

ically in *fig. 6a-h*. The component *8* is lifted from the component magazine *9* by two mechanical fingers and laid on two horizontally movable anvils *10* (*fig. 6a*). Two bars *11* (*fig. 6b*) now descend on the outside of the anvils and bend the terminal wires of the component as shown in *fig. 6c*. The moment these bars touch the panel *12* the anvils are withdrawn and two other bars *11a* descend (*fig. 6d*). These first push the wires through two holes in the panel and then through holes in a cutting and bending jig *13* (*fig. 6e*), underneath the panel. This device cuts off the wires to the appropriate length and bends them over to one side under the panel (*fig. 6f* and *g*). When this is done, the insertion mechanism returns to its starting position, whereupon the next component (*8'*) is lifted from the magazine and placed on the anvils (*fig. 6h*).

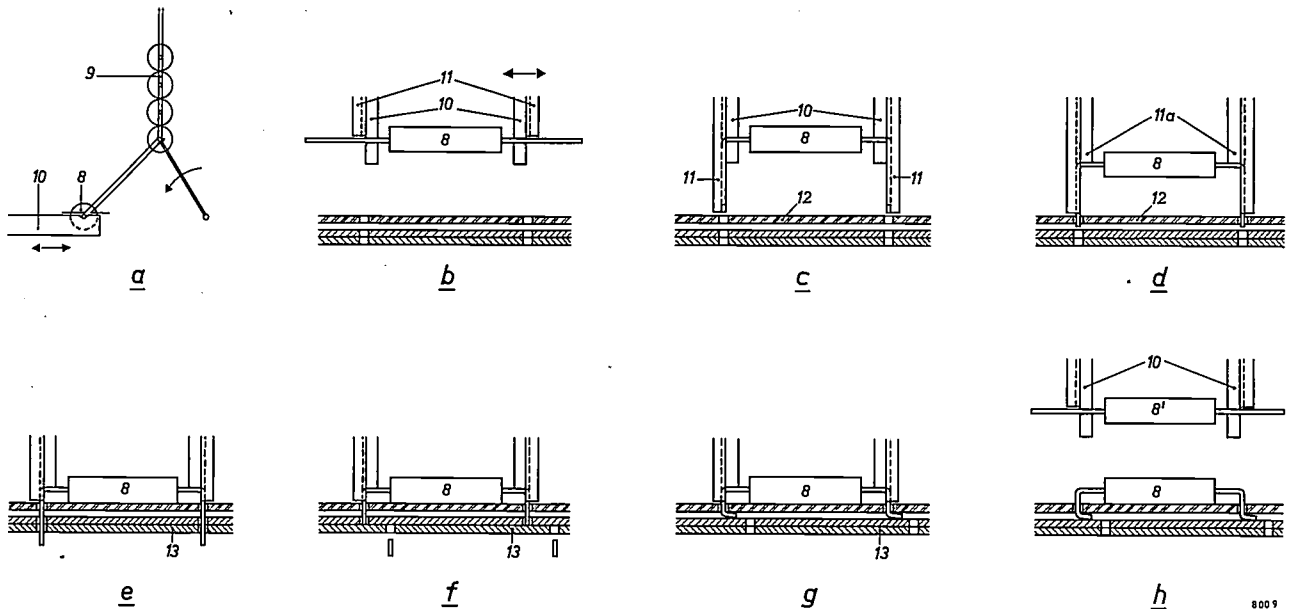


Fig. 6. Operation of the insertion mechanism of the AW head.

- a) The component *8* is supplied from magazine *9* to the anvils *10*.
- b) and c) Bars *11* descend to bend the terminal wires ready for insertion.
- d) Bars *11a* push the terminal wires into specific holes of the panel *12*.
- e) The terminal wires have now passed through the holes in the cutting and bending jig *13*, which cuts them off (*f*) and bends them over (*g*).
- h) The next component (*8'*) is placed on the anvils *10*.

The insertion mechanism is symmetrical in relation to the middle of the component and is divided into a stationary and a sliding part (*fig. 6b*). The point where the terminal wires have to be bent is preset with the aid of a steel strip. To enable components to be mounted parallel either to the long or the short side of the panel, the insertion mechanism can be rotated through 90° about the vertical axis of the stationary part; for this purpose it is mounted on a turntable, which is fixed to the U-shaped frame of

the head. The cutting and bending jig turns at the same time but otherwise needs no resetting, as the cutting and bending plate contains a sufficient number of holes. The distance between these holes is equal to the basic grid spacing $e = 2.54 \text{ mm}$ ($\approx 0.10 \text{ in.}$) between the holes in the panel (see below) or to a whole multiple of e .

The components, stuck to strips, are fed to the AW heads in the form of a tape (Table I).

An *insertion head for bridging wires* can be seen in *fig. 7*. The construction is much the same as that of the AW head. It differs only in that the component magazine is replaced by a reel of bare wire (*14*) and that the BW head contains a mechanism which supplies and cuts off pieces of wire of the required length. The length is set by adjusting the sliding part of the insertion mechanism. A bridging wire is

mounted on the panel in the same way as an AW component.

The *valve-holder insertion head* (*fig. 8*) is entirely adapted to the manner in which the valve holders are supplied. In cooperation with the suppliers of these components a form was adopted that lent itself readily to mechanical insertion. This made it possible to keep the insertion mechanism of the VH head very simple.

The valve holders are delivered in series of fifty,

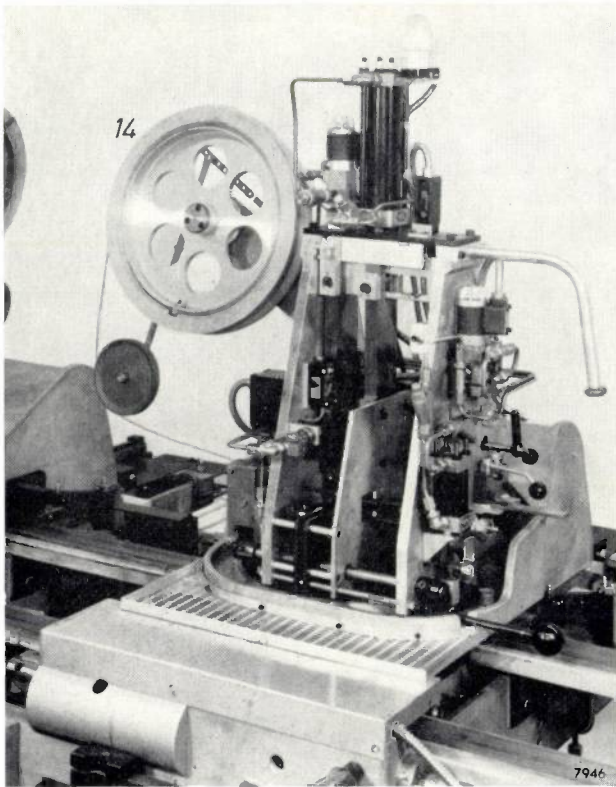


Fig. 7. Insertion head for bridging wires (BW head), largely identical with the AW head (figs 5 and 6). The supply magazine here is a reel of bare wire 14.

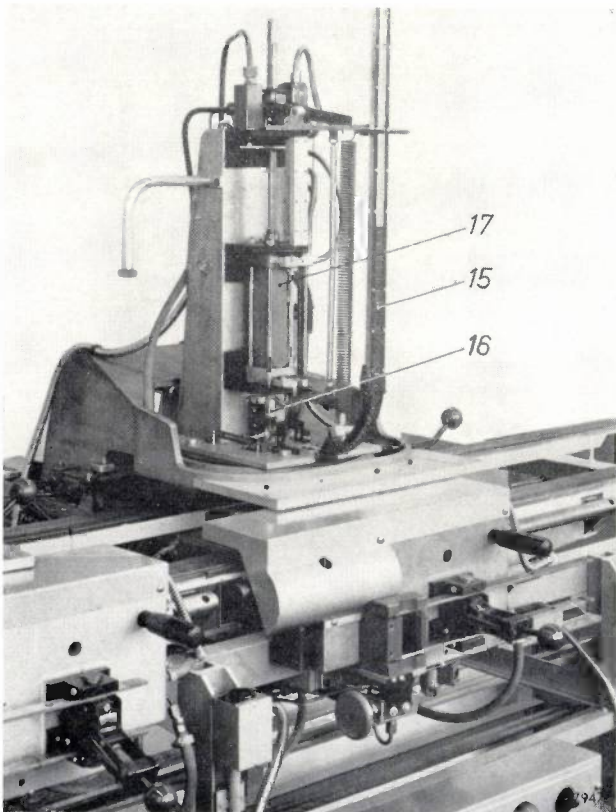


Fig. 8. Insertion head for valve holders (VH head). 15 magazine containing valve holders, which are attached to a strip of resin-bonded paper. 16 punch which successively stamps a valve holder from the strip and mounts it on the panel. 17 air cylinder which drives the mechanism.

fixed to a strip of resin-bonded paper 1 metre long (15 in fig. 8). A punch 16 successively stamps each valve holder from the strip and mounts it in a single continuous movement on the panel. At the same time the terminal tags of the valve holder are bent underneath the panel.

The insertion head for pin-up components (ceramic capacitors) is shown in fig. 9. The distance between the terminal wires of the PU components is three times the standard spacing of 2.54 mm. The PU components are supplied attached to a strip

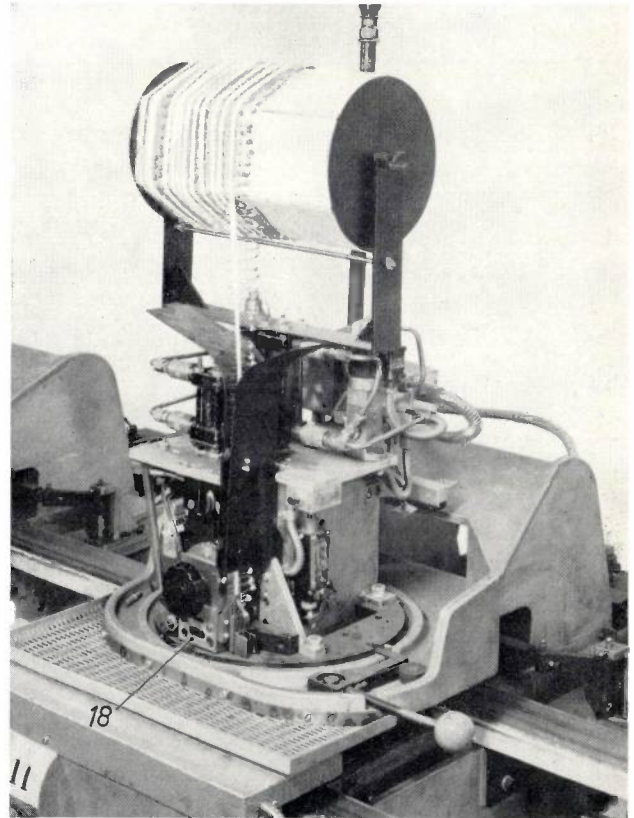


Fig. 9. Insertion head for pin-up components (PU head). The strip of components is wound on the hexagonal drum. The component-feed mechanism is marked 18.

(Table I). A mechanism 18 feeds in one component at a time. Pincers in the insertion mechanism cut the component from the strip and insert the wires into specific holes in the panel. A cutting and bending device underneath the panel cuts the wires to the required length and bends them against the panel.

The BW, VH, and PU heads are all fitted with the same turntable construction as the AW heads, thus offering the same choice of two mounting directions differing by 90° .

The drive system for the insertion mechanisms will be discussed at the end of this article.

The holder

The holder consist of a rectangular frame (*fig. 10*) with four lugs by which it rests on the band conveyor which carries it to the insertion heads. The holders are large enough to take a panel of the maximum size (150 by 300 mm). In view of the large number required (100), the holders are made by casting. Light metal was chosen to keep the weight down.

A panel is secured in the holder by a movable strip 20 which presses two recesses in the reference face of the panel (discussed below) against two of the five pins in the side 19 of the holder; which two depends on the size of the panel. At the other side of the holder there are two studs 21, one of which is V-shaped. When the holder arrives under an insertion head, a roller 22 presses these studs against two stop pins in the insertion head, which fixes the holder's position. The holder is thereby raised 1 mm so that it is clear of the continuously moving band conveyor.

The two rear lugs are fitted with pins 23. After leaving the last insertion head the holder is sus-

pending vertically from these pins, and is returned in this position to the beginning of the assembly line.

The conveyor mechanism

The conveyor mechanism consists of two twin band conveyors in continuous motion, situated one above the other and travelling in opposite directions (1 and 5 in *fig. 11*). A device at each end transfers the holders from one conveyor to the other (*fig. 10* and *fig. 12*). The upper band conveyor (24 in *fig. 12*) carries the horizontal holders to the insertion heads and the lower conveyor (25) returns the holders, now vertical, to the beginning of the assembly line.

The upper conveyor consists of two leather belts side by side, covered by a wear-resistant plastic layer. The high coefficient of friction of this material ensures that the holders are carried along once they have left an insertion head. The surface of the lower conveyor is of steel, so as to reduce the friction between the continuously moving conveyor and the pins of the intermittently stationary holders.

The supply of 100 holders is distributed in groups over the length of the steel conveyor. This is done in order to distribute the weight evenly along the conveyor and to avoid excessive driving forces at the beginning of the assembly line.

For each group of holders there is a pneumatically operated lock system which, as soon as a new holder enters the machine via the upper conveyor, lets one holder pass, which is then carried by the lower conveyor to the next group. In this way, holders are supplied and removed synchronously, and thus cannot accumulate.

At the beginning of the line is the lift installation referred to (*fig. 10*), which raises the holders from the lower to the upper conveyor and returns them from the vertical to the horizontal position. At the end of the line is the slide which has also been referred to (26 in *fig. 12*); this transports the holders from the upper to the lower conveyor and turns them from the horizontal into the vertical position. The conveyor system for the holders thus constitutes a closed loop.

Resetting the insertion heads

One of the design engineer's tasks was to devise a system for resetting the insertion heads in the shortest possible time. Before dealing with the system adopted, we shall take a closer look at the panel on which the components are to be mounted.

The panel contains holes through which the terminal wires of the components have to be inserted in order to join them to the wiring pattern on the other side. The holes are arranged in accordance

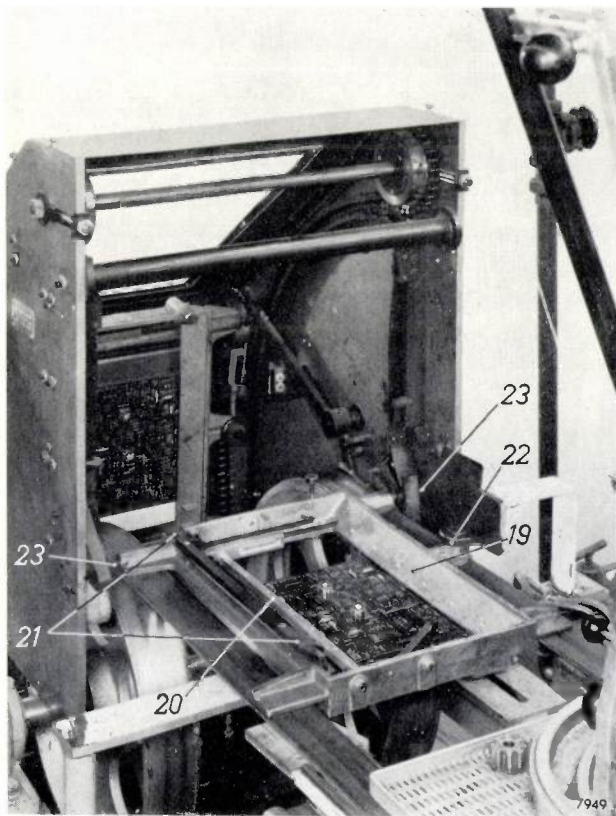


Fig. 10. Beginning of the assembly line. In the foreground a holder containing an almost blank panel. Strip 20 presses the panel with its reference face against pins in the holder (these pins, not visible here, are located at 19). The studs 21 ensure that the holder is raised above the conveyor upon arrival under an insertion head. 22 pressure roller. 23 pins by which the holder is suspended when it is returned to the beginning of the line.

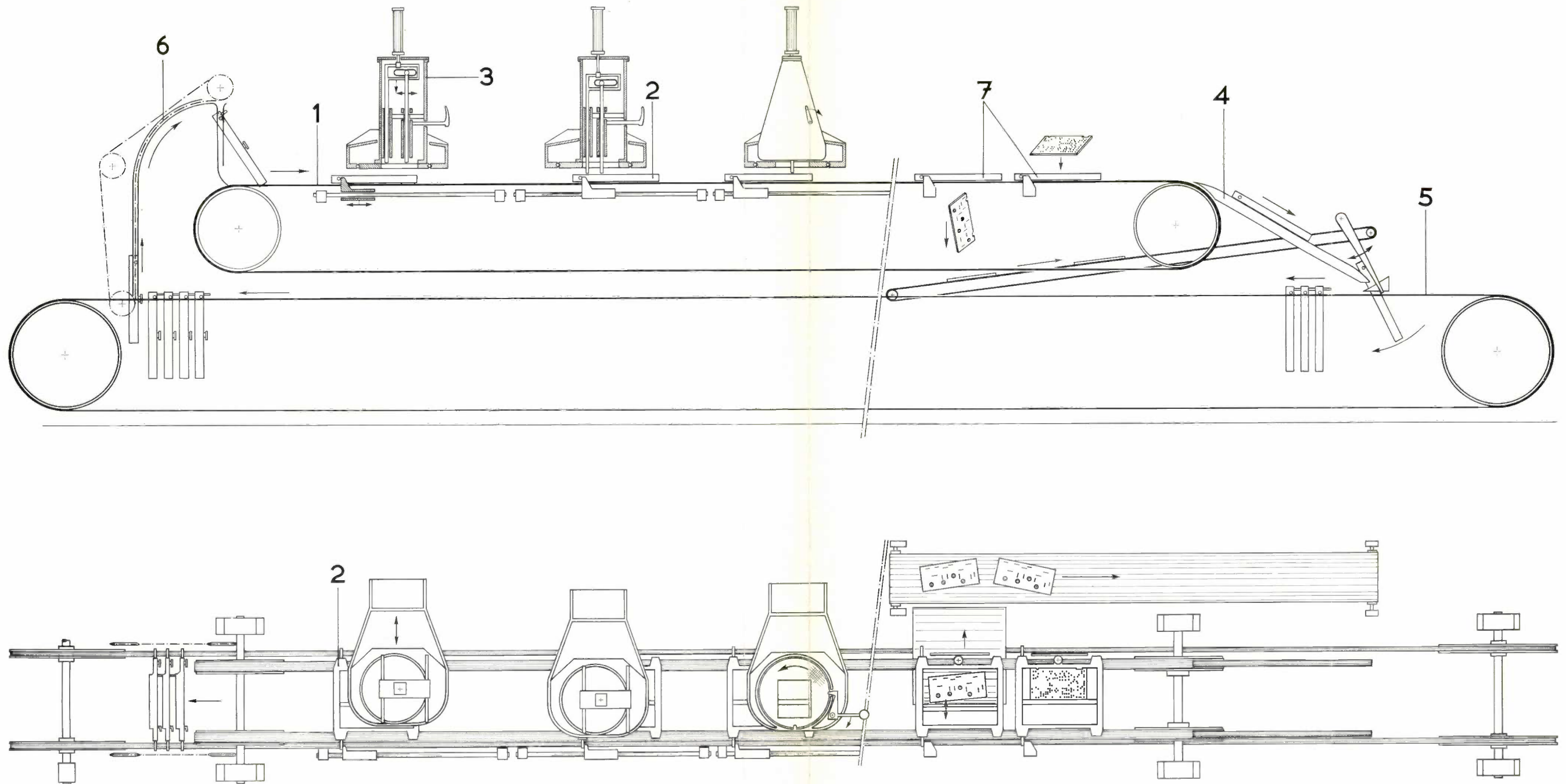


Fig. 11. Simplified sketch of the assembly line mechanically mounting components on printed-wiring panels for television sets. 1 twin leather belt-conveyor, which carries the panel holders 2 from left to right under the insertion heads 3. The panel holders arrive via a slide 4 in the vertical position on a steel band conveyor 5, which returns them to the beginning of the line where they are lifted by an elevator device 6, and deposited horizontally on the upper conveyor 1. Fully assembled panels are removed from the holders, carried off to the right, and replaced by blank panels at station 7.

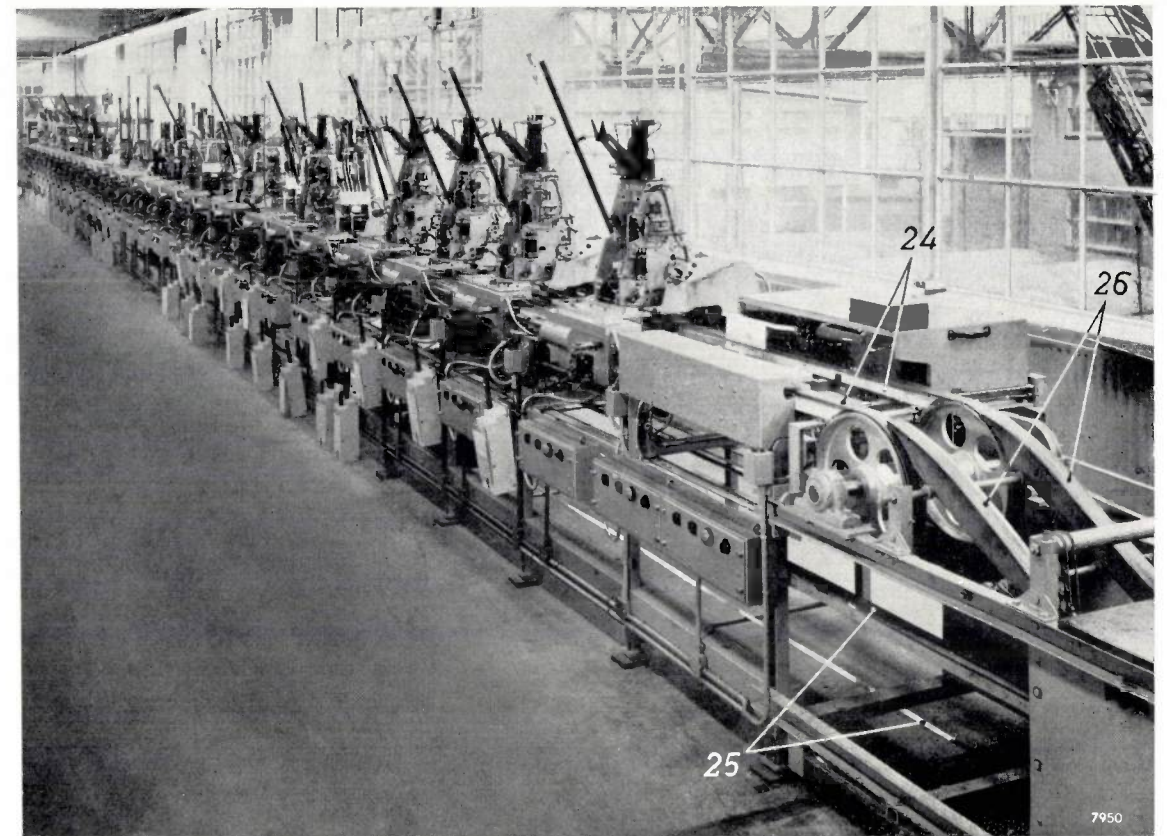


Fig. 12. View of the assembly line in operation at Philips, at present with 24 assembly stations. In the foreground is the end of the line. At 24, where the unloading and reloading station is located, can be seen the twin band conveyor. 26 slide that carries the holders to the lower conveyor, 25, which returns them in the vertical position to the beginning of the line.

with the standard system laid down in publication 97 of the International Electrical Commission (I.E.C.).

The rectangular panel is thought of as divided into a grid of squares with sides $e = 2.54$ mm, parallel to the sides of the panel. On the largest panel the squares cover a rectangle measuring $120e$ by $62e$ (fig. 13). The position of each hole, which must be

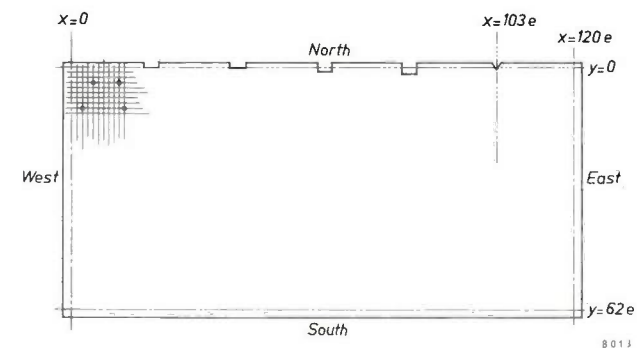


Fig. 13. The I.E.C. has laid down that the holes in printed wiring panels may only be applied at the corners of a standard grid of squares with sides $e = 2.54$ mm (≈ 0.10 in.). In panels of the largest type the grid covers a rectangle with sides $120e$ and $62e$.

One long side of the panel shown represents the reference face standardized at Philips, with a V-notch at $x = 103e$ and four rectangular recesses of differing depths (less than four on smaller panels).

at one of the corners of the squares, is determined by coordinates x and y with respect to axes parallel with the long and short sides of a panel respectively. The diameter of the holes is standardized at 1.3 mm, the thickness of the panels at 1.6 mm.

At Philips it is also standardized that one long side of the panel should form the *reference face* mentioned earlier. For that purpose this side is provided with four rectangular recesses and one V-shaped notch, as illustrated in fig. 13. The panel lies with the notch and the far left recess against two of the five pins in the holder. The three recesses in between are slightly deeper and are thus clear of the other pins, so that the panel is supported at two points only.

Fig. 13 relates to the largest type of panel. Smaller types are curtailed from the left in fig. 13, and thus lack one or more of the rectangular recesses. The panel is then supported in the holder by the V-notch and by the remaining far left recess (on one of the three pins that are not used for the largest panel).

With a view to simple and rapid resetting of the insertion heads, it is necessary to take into account that one component and the next may differ in the following respects:

- 1) the location of the component on the panel,
- 2) the direction in which it is to be mounted: "east-west" or "north-south" (fig. 13),
- 3) the length over which the terminal wires have to be bent, and
- 4) the diameter.

Furthermore, the magazines must be designed so that when empty, they can be rapidly replaced by full ones. We shall examine these points in turn.

1) The first question concerns the method of correctly positioning an insertion head over the panel.

At the front of the insertion head a support 27 (fig. 5) is fitted, which can be moved in the x direction (coinciding with the direction in which the panel is conveyed). The holder arrives under the insertion head at a speed of 1 metre per second, where it is braked by means of an air cylinder against which a pin 28 in the holder comes to rest. A bar 29 clamps the holder by pressing the studs on the holder against the round pins fixed in the support. The position of the support is determined by the lower of two metal setting bars (30). This bar contains two holes. A rigid pin 31, fixed to the frame, fits into the left hole of the setting bar, and a second pin, located in a locking device 32, fits into the right hole. In this way

then, the position of the support is fixed, and at the same time the x coordinate is set, the spacing between the two holes in the strip being chosen so as to correspond to the specified value of x .

The other setting bar in fig. 5, marked 33, serves for setting the next position of the support. During the operation of the machine this bar can be replaced by another which supplies new information. This system of pre-selection enables a new x coordinate to be set within 2 seconds, as all the operator has to do is to lift the handle of the locking device 32 and move the support to the left or right (depending on the location of the hole in the setting bar) until the spring-loaded pin of the locking device snaps into the hole in the upper setting bar (33).

To set the y coordinate the U-shaped frame 34 (fig. 14), to which the insertion mechanism 35 is fixed, can be moved perpendicular to the direction in which the holders are conveyed. The position of this frame is determined by means of setting bars 36 in the same way as the position of the support, using a locking device (37), two bars (36) and a fixed pin (38). The operation, however, is different, the frame being set pneumatically in the new position. When changing the setting bar not in use, the personnel operate

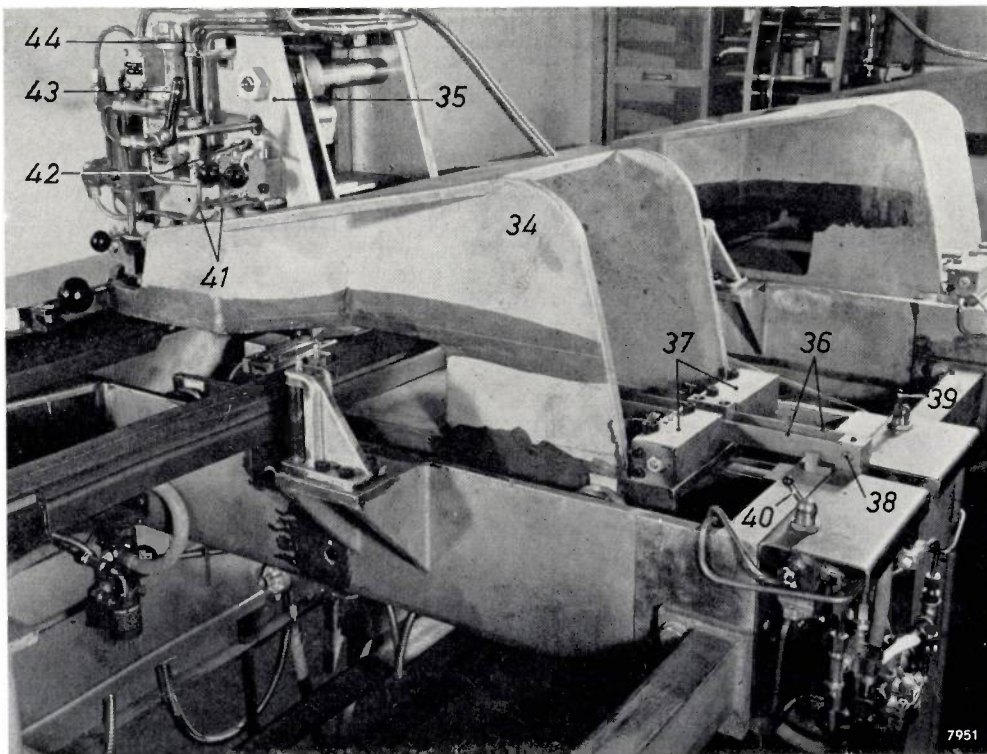


Fig. 14. Mechanism for setting the y coordinate of the point where the insertion head has to insert a component. 34 U-shaped frame and 35 insertion mechanism. 36 two setting bars, each of which fixes a y coordinate. 37 locking device. 38 fixed pin. 39 and 40 manual controls for operating air valves. 41 strip which fixes the length over which the components terminal wires are to be bent (leg spacing). 42 fixed pin. 43 handle used for resetting the leg-spacing for different components. 44 exchangeable strip to compensate for differences in component diameters.

the valves 39 and 40 at the rear of the machine. The setting of valve 39 determines which strip will presently have to be scanned for resetting the head, and the setting of valve 40 indicates whether the frame will have to be moved forward or backward to find the hole in the new strip. These indications are taken up in a pneumatic system, which is started by a push-button on the front of the machine. The setting of a new y coordinate takes 0.3 second.

Since, as mentioned, the setting bars can be changed in order to preselect a new x or y coordinate while the machine is in operation, the actual loss in time is only 2.3 seconds.

When the above operations are being carried out the centre line of the fixed part of the insertion mechanism is brought into position, as discussed under the heading *The insertion heads*.

2) Apart from the location, an important consideration in resetting the insertion heads is the direction in which the components have to face. To make the most efficient use of the space on the panel, the components face only in the x or in the y direction, i.e. "east-west" or "north-south". We have seen that for this reason the insertion mechanism (35 in fig. 14) can be rotated through 90°.

3) The third point concerns the setting of the length through which the terminal wires (e.g. of AW components) have to be bent (leg-spacing). As discussed in the foregoing, the insertion mechanism for AW components consists of a fixed and a sliding part. The spacing between these parts is fixed by a setting bar 41 (fig. 14) with two holes. A rigid pin 42, fixed to the frame, fits into one hole. The sliding part carries a locking device, the pin on which snaps into the other hole. Here too, the construction is duplicated, similar to that for setting the x and y coordinates, and preselection is possible.

The adjustment is made by means of the handle 43, the new setting bar being scanned while the handle is pushed in or pulled out (depending on the position of the hole in the new setting bar). This setting takes 2.3 seconds.

4) If no account were taken of the differences in diameter that may be met in successive components, there would be a risk of terminal wires being sheared during insertion. To avoid this it is necessary to reset the insertion mechanism every time a component with a different diameter comes along.

This resetting is very simple. The tool pressing on the component descends until a stud on the tool meets a stop. All that is necessary is to place between the stud and the stop a metal strip (44 in fig. 14), the thickness of which is equal to half the diameter of the component. This strip limits the stroke

of the press tool so as to prevent the terminal wires from slipping. This resetting takes 3 seconds.

5) Finally, a word about the construction of the magazines. These are designed to contain exactly as many components as there are holders in circulation, in our case 100. The magazines are thus empty after the 100 holders have passed down the line. They are mounted on two plates provided with open slots into which fit pins attached to the magazine. The time taken to replace an empty magazine by a full one is about 19 seconds.

If we add the various times needed for resetting the insertion head, we come to a total of 31 seconds. The actual time is at present even shorter, for since components have been supplied in tape form it has been possible to simplify the feed mechanism and thus shorten the resetting times. Moreover it will hardly ever be necessary to make all the resetting adjustments mentioned to every insertion head. The number of resetting adjustments are arranged so as to be a minimum. Of course, new x and y coordinates will regularly have to be set, and empty magazines replaced by full ones. It has been proved that with good organisation the total average time needed for resetting an AW head can be reduced to about 24 seconds.

Safety measures concerning the drive mechanism

The insertion mechanisms of the assembly line are all driven by air cylinders (one or two per insertion head) which are controlled by electromagnetically operated valves³⁾. Fig. 8 shows the air cylinder 17 of a VH head.

Since the mounting of a component consists of a series of operations done step by step, it is possible in each case to check whether in fact one operation has been completed before it is the turn of the next one. Use is made of this in our assembly line.

If a fault occurs for some reason or another, the machine stops automatically. This prevents damage being done in the event of the failure of an insertion head. A red lamp lights up on the head responsible for stopping the machine.

Each insertion head is protected by the following safeguards:

- a) The holder clamping mechanism can only work if the holder is properly aligned in front of the support.
- b) The insertion mechanism cannot start inserting a component until the holder is properly clamped to the support.

³⁾ See Philips tech. Rev. 23, 378, 1961/62 (No. 12).

- c) The holders are not released until all insertion heads have completed the prescribed operations and have returned to their starting position.

Final remarks

In the foregoing we have described how the design engineer met the requirements summarized at the beginning of Part II. The result is a production machine which — within the design limitations of the printed-wiring panel as a carrier of electronic components — is universally usable. Since the insertion heads can be reset quickly, the machine is equally capable of efficiently fitting small series of panels with components.

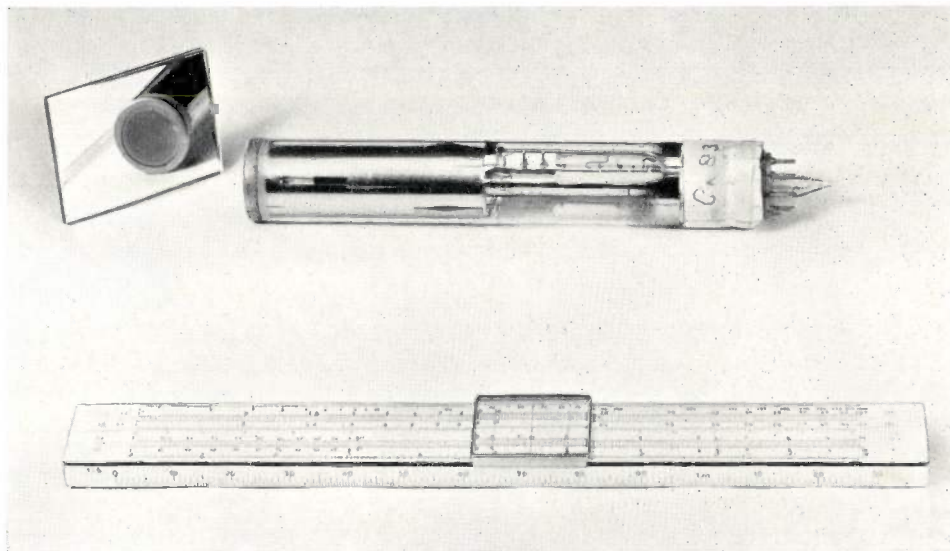
The increasing adoption of mechanical mounting methods will influence the development of electronic components and building elements (combinations of several components). It is to be expected that efforts

will be made right at the beginning of this development to design the components in such a way as to keep the construction of insertion machines as simple as possible.

Summary. This article is devoted to a machine, in the form of an assembly line, which is in use at Philips in Eindhoven for automatically mounting components on printed-wiring panels for use in television sets.

Part I is a general discussion of the mechanical mounting of components. Among the points dealt with are two methods of mechanized insertion, the grouping of mechanically mountable components according to their shape and their distribution over the panels, production capacity, the choice of the number of panel holders and the number of assembly stations so as to minimize the cost of the machine in relation to the production capacity, and the manual maintenance of the machine. Considerable significance is attached to the stoppage percentage (the fraction of the time lost owing to malfunctioning of the machine).

Part II deals with the mechanical construction. A discussion of the operation of the assembly line is followed by a review of its principal components.



THE "PLUMBICON", A NEW TELEVISION CAMERA TUBE

621.385.832.564.4

The most commonly used types of television camera tubes are the image orthicon, the image iconoscope and the vidicon. The latter type, in which the storage target is a photoconductive layer consisting of, e.g. antimony trisulphide Sb_2S_3 , has the virtues of simple construction and easy operation. The vidicons available up to now, however, only supply good pictures at high levels of illumination: at low levels the dark current of the light-sensitive surface

becomes very troublesome, and moreover the response of the tubes is then too slow for the demands of broadcast television.

The Philips Research Laboratories, after development work covering several years, have now produced a new camera tube which is based on the vidicon principle, and thus has the same advantages of simplicity and easy operation, but uses a different photoconductive material. The photoconductive

- c) The holders are not released until all insertion heads have completed the prescribed operations and have returned to their starting position.

Final remarks

In the foregoing we have described how the design engineer met the requirements summarized at the beginning of Part II. The result is a production machine which — within the design limitations of the printed-wiring panel as a carrier of electronic components — is universally usable. Since the insertion heads can be reset quickly, the machine is equally capable of efficiently fitting small series of panels with components.

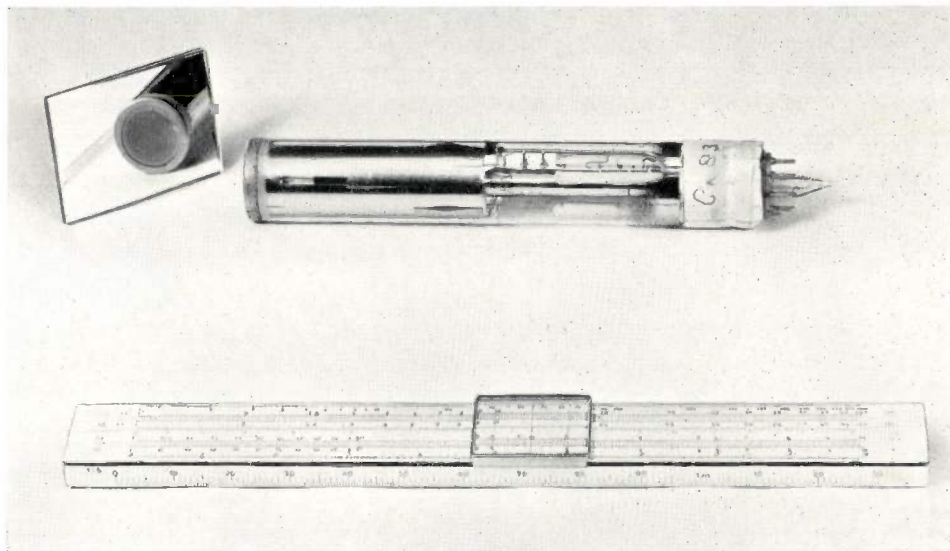
The increasing adoption of mechanical mounting methods will influence the development of electronic components and building elements (combinations of several components). It is to be expected that efforts

will be made right at the beginning of this development to design the components in such a way as to keep the construction of insertion machines as simple as possible.

Summary. This article is devoted to a machine, in the form of an assembly line, which is in use at Philips in Eindhoven for automatically mounting components on printed-wiring panels for use in television sets.

Part I is a general discussion of the mechanical mounting of components. Among the points dealt with are two methods of mechanized insertion, the grouping of mechanically mountable components according to their shape and their distribution over the panels, production capacity, the choice of the number of panel holders and the number of assembly stations so as to minimize the cost of the machine in relation to the production capacity, and the manual maintenance of the machine. Considerable significance is attached to the stoppage percentage (the fraction of the time lost owing to malfunctioning of the machine).

Part II deals with the mechanical construction. A discussion of the operation of the assembly line is followed by a review of its principal components.



THE "PLUMBICON", A NEW TELEVISION CAMERA TUBE

621.385.832.564.4

The most commonly used types of television camera tubes are the image orthicon, the image iconoscope and the vidicon. The latter type, in which the storage target is a photoconductive layer consisting of, e.g. antimony trisulphide Sb_2S_3 , has the virtues of simple construction and easy operation. The vidicons available up to now, however, only supply good pictures at high levels of illumination: at low levels the dark current of the light-sensitive surface

becomes very troublesome, and moreover the response of the tubes is then too slow for the demands of broadcast television.

The Philips Research Laboratories, after development work covering several years, have now produced a new camera tube which is based on the vidicon principle, and thus has the same advantages of simplicity and easy operation, but uses a different photoconductive material. The photoconductive

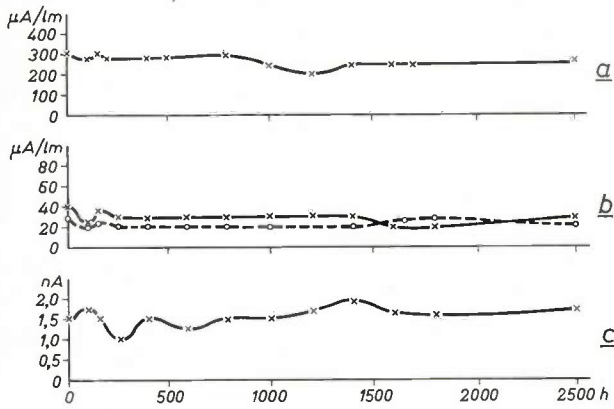


Fig. 1. Some characteristics of laboratory versions of the "Plumbicon" camera tube life tests: a) sensitivity to white light (tungsten light, colour temperature 2870 °K); b) sensitivity to red light (solid line) and blue light (dashed line); this test showed that the spectral sensitivity of the tube does not vary appreciably during its useful life; c) dark current in nanoamperes.

layer consists of vapour-deposited lead-monoxide, PbO, and we have therefore called this tube the "Plumbicon"¹⁾. By vapour-depositing the PbO under accurately controlled conditions (atmosphere, temperature, etc.), it has proved possible to obtain a layer that has a very low dark current ($< 5 \times 10^{-9}$ A); the result is an exceptionally uniform picture. The

A great deal of attention was paid in the work on this new tube to its speed of response. Slowness of response is sometimes due to centres that also affect the other properties of the camera tube. It was found possible, by carefully choosing the deposition conditions, to reduce the time lag in the response to less than 0.2 second, which is scarcely perceptible to the eye, and in particular to make it independent of the level of illumination, without losing the favourable properties mentioned. The useful life of the tube was of course also the subject of extensive investigations. The useful life now achieved is longer than that expected of television camera tubes of studio quality. Certain characteristics of the "Plumbicon" were found to remain more or less unchanged even after operation for several thousands of hours; see fig. 1.

Fig. 2 illustrates the speed of response and uniformity of the "Plumbicon" picture compared with the pictures photographed simultaneously by an Sb₂S₃ vidicon and an image orthicon. The laboratory version of the tube itself is shown in the title photograph.

A problem more or less distinct from those mentioned above was to overcome the notorious white

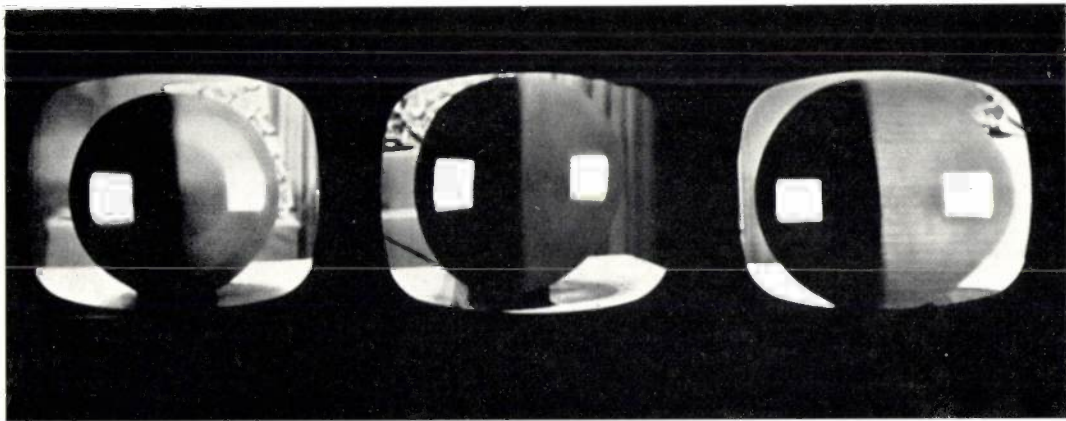


Fig. 2. Photographs (exposure $\frac{1}{25}$ sec) of three monitors connected, from left to right, to an Sb₂S₃ vidicon, a "Plumbicon" and an image orthicon; the three camera tubes — each optimally adjusted — were all directed at the same object, a disc rotating at 40 r.p.m. The Sb₂S₃ vidicon has a relatively slow response, which can be judged from the smears behind the rotating white square. The dark halo inherent in the image orthicon can be seen here around the white square, and there is also some lack of uniformity in the background of the image.

sensitivity of the layer is also very high (> 150 μ A/lumen). In this respect the "Plumbicon" camera tube rivals the most sensitive camera tubes at present in use for broadcast television — image orthicons — whilst at the same time it is free of the spurious signals inherent in the latter cameras (e.g. dark halos).

¹⁾ An earlier communication concerning a camera tube using a sensitive layer of lead monoxide was published some years ago in this journal: L. Heijne, P. Schagen and H. Bruining, Philips tech. Rev. 16, 23, 1954/55.

spots — strongly localized extra-high contributions to the dark current — which are also encountered in other kinds of camera tube. This effect has been almost entirely suppressed in the new tube.

It is intended in due course to publish in this journal a detailed discussion of this tube and its characteristics.

E. F. de HAAN *).

*) Philips Research Laboratories, Eindhoven.

WIRELESS-POWERED TOYS

by J. F. van OORT *) and W. BAKKER *).

688.727-83:621.398

*Although toys are not usually produced by Philips, and research in this field is not a normal subject of our journal, we believe our readers will be interested in the idea presented below. It relates to a problem long familiar to traction engineers **) and the solution offered, though limited, may well be adaptable to more serious applications.*

The idea of remotely controlling toy models of ships, cars and aeroplanes by radio is nearly as old as radio engineering itself. By present standards the principle is simple: the "controller" has a small radio transmitter with which he can send out coded signals; the remotely controlled model contains a receiver tuned to the transmitter, electro-mechanical devices which decode the signals and use them to control the mechanism, and electric batteries or some other source of power for the mechanism and the receiver. It is well known that the application of this principle is nowadays not confined to toys,

but extends to the control of "drones" (pilotless aircraft targets), rockets and artificial satellites.

Details are given in this article concerning an entirely different system of remote control, in which a wireless method is used not only to *steer* the model but also to supply the *energy* needed for all movements and for decoding the signals. The most obvious advantage of this is, of course, that the remotely controlled model need carry no power source of its own, such as accumulators or dry batteries. Another advantage is that it makes it possible to simplify the actual control. The principle will be briefly described in its application to a small model of a fork-lift truck (*fig. 1*).

The surface on which the model travels is surrounded by an inductive loop (partly visible in *fig. 1*) through which a signal generator passes an

*) Philips General Advertising Division, Exhibition Department, Eindhoven.

**) See, for example, the description of a "non-contact traction system" using capacitively transmitted radio-frequency power, by G. I. Babat, J. tech. Phys. U.S.S.R. 16, 555, 1946 (also in Engrs Digest (London) 8, 78, 1947).

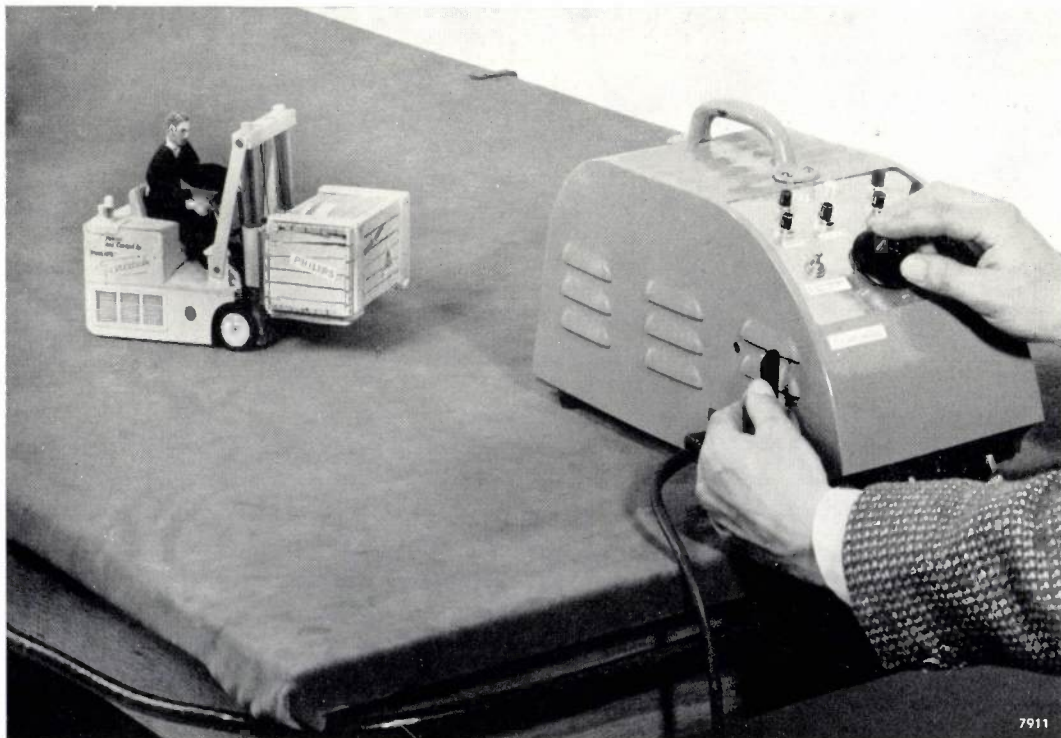


Fig. 1. Model of a fork-lift truck, wireless-powered and steered by the control apparatus on the right.

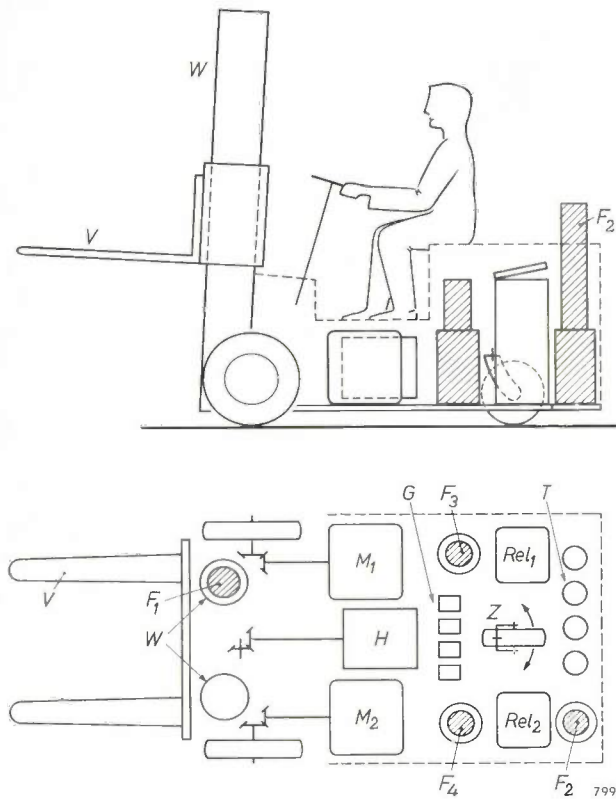


Fig. 2. Layout of the components in the model. F_1 to F_4 wire-wound antenna rods of ferroxcube (F_1 is contained in one of the two guides W for the fork V). T trimming capacitors for the four tuned circuits. G rectifiers. Antenna rods F_1 and F_2 supply the power for the DC motors M_1 and M_2 respectively, each of which drives one of the front wheels through its own gearbox and bevel-gear system. Z free-swivelling castor. H motor for the fork-lift mechanism.

To switch from forward drive into reverse, a current pulse of 15 kc/s is generated in the ring cable by means of a push-button on the signal generator (which connects a small capacitor in parallel with the tuning capacitor). The antenna rod F_3 , tuned to 15 kc/s, then energizes the relay Rel_1 with holding contact, which reverses the current direction in both motors M_1 and M_2 . A second identical current pulse reverses the situation again. With antenna rod F_4 , which is tuned to 30 kc/s, and relay Rel_2 the two main circuits are switched over in a similar manner from motors M_1 , M_2 to the fork-lift motor H (the two circuits are then in parallel). The relay Rel_1 , which reverses the direction of travel, now does the same for the lifting mechanism.

alternating current of say 20 kc/s. The model is equipped with a vertical antenna rod of ferroxcube wound with a coil which, together with a small trimming capacitor, forms a resonant circuit tuned to the relevant frequency. Because of the high frequency and the fact that the ferroxcube attracts a substantial part of the magnetic flux inside the inductive loop, a considerable e.m.f. is generated in this circuit. If in addition the antenna is provided with a coupling coil, then, given suitable matching, the power picked up in this way can be used to light an electric bulb or, after rectification of the induced voltage, to drive a small DC motor.

Actually our small model contains not one but four antenna rods, each with its own tuned circuit.

Two of the circuits, with antennae F_1 and F_2 in fig. 2, are tuned respectively to 19 and 21 kc/s; each of them, via two rectifiers, feeds one of two 3 W motors, one of which drives the right front wheel and the other the left front wheel of the truck. The signal-generator frequency can be continuously varied between the two frequencies, as a result of which the power in one tuned circuit is altered relative to the other, thus causing one motor to run faster than the other, or both at the same speed. In this way the power transmission is very simply combined with the steering control. The speed of travel too can be very simply controlled by varying the current through the inductive loop. The third and fourth antenna rods serve respectively for switching the truck into reverse, and for operating the fork-lift. The circuits of these antennae are tuned to 15 and 30 kc/s respectively. The caption to fig. 2 explains the switching mechanism and fig. 3 shows most of the components of the model.

The signal generator comprises mainly two EL 34 pentodes and draws a maximum of 150 W from the mains. The inductive loop, which may have an area, for example, of 1×2 metres, forms part of the tuning inductance. At 50 W generator power the field strength in the middle of the loop is about 4×10^{-5} Wb/m² (0.4 gauss). The power radiated by the system is negligible.

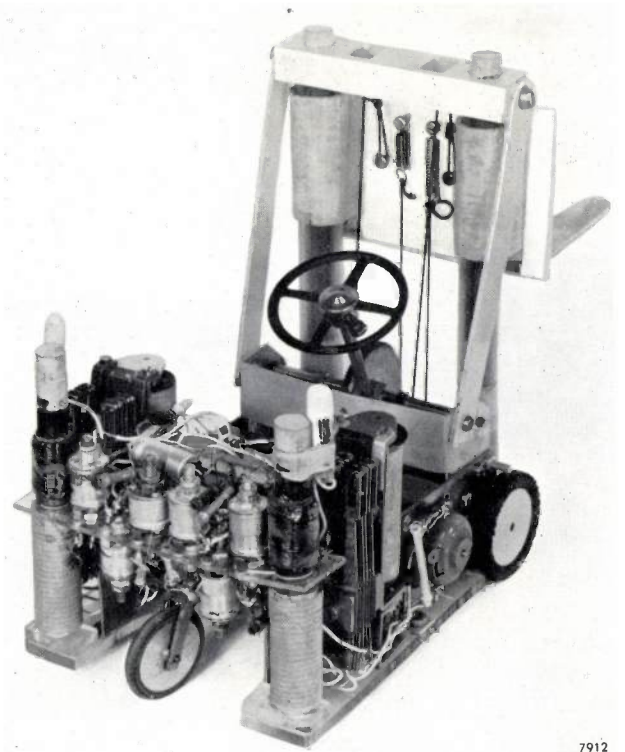


Fig. 3. The model opened up.

The method of power transmission employed here may be compared with that of a transformer: the primary is represented by the loop system, the iron core by the ferroxcube rod, and the secondary by the winding around the rod. The fact that this transformer works with reasonable efficiency in spite of the extremely inefficient "magnetic circuit", is due to the high frequency used and to the properties of the ferroxcube, which combines at such frequencies a high permeability with very low eddy current and other losses.

Models like the one described here have also been made for more serious purposes, in particular for instruction and practice in driving lessons. The same principle could also be applied to the operation of equipment in hermetically sealed rooms or in places in laboratories and factories that are not readily accessible owing to radiation or high-tension hazards. We have already mentioned the advantages compared with more conventional systems, i.e. the elimination of batteries which have to be periodically changed or recharged, and the simplicity of the

equipment. To enable this particular model to perform all its functions with equal flexibility using the normal method of radio control ¹⁾, quite a number of valves or transistors and associated circuit elements would have been required. These conventional systems are of course still necessary where the objects to be controlled have to be more freely manoeuvrable and have to perform more intricate operations, which also calls for higher power.

¹⁾ A. H. Bruinsma, Radio-controlled models, Philips tech. Rev. 15, 281-285, 1953/54.

Summary. A description of a model of a fork-lift truck, in which the driving motors and lifting mechanism draw their power from the magnetic field of an inductive loop enclosing a manoeuvring surface of e.g. 1×2 m. To this end the model carries wire-wound antenna rods of ferroxcube, tuned to the frequency of an alternating current (about 20 kc/s) passed through the inductive loop with the aid of a signal generator. Very simple means are used not only to supply the power but also to control the movements of the model. — This principle of transmitting power to a moving object may perhaps find application in special cases in laboratories.

ETCH PITS ON A ZINC-SULPHIDE CRYSTAL

546.47'221:548.572

The cubic modification of zinc sulphide, also called zinc blende or sphalerite, has a structure resembling that of diamond: all the zinc atoms are tetrahedrally surrounded by sulphur atoms, and vice versa.

Closer examination reveals that the crystal can be regarded as built up of alternate layers of zinc and sulphur atoms, stacked one above the other. In this structure a layer of zinc and a layer of sulphur atoms are one above the other; at a somewhat greater distance in this stacking direction, which is a crystallographic $\langle 111 \rangle$ axis, there again follow a layer of zinc and a layer of sulphur atoms close together, and so on (see *fig. 1*). The zinc and sulphur atom layers after each large spacing are in the opposite sequence if the stacking is continued in the opposite direction. This means that the $\langle 111 \rangle$ axis is a polar axis.

This polarity appears in a surprisingly beautiful way in the macroscopic-chemical behaviour of such a crystal, i.e. when etched. The photographs in *figs 2* and *3* illustrate this: they show the opposite end faces perpendicular to the polar axis (i.e. a (111) plane and a $(\bar{1}\bar{1}\bar{1})$ plane) of the same cubic zinc-sulphide crystal, after etching with hydrochloric acid. It can be seen that one of these crystal faces is covered with etch pits, bounded by planes ($\{100\}$ planes),

from which the triaxial symmetry of the $\langle 111 \rangle$ axis can easily be recognized. On the opposite end face the etch pits are rounded, although here too the triaxial symmetry is still recognizable.

Similar phenomena have been observed by other investigators, both on zinc sulphide and on substances with the same crystal structure, in particu-

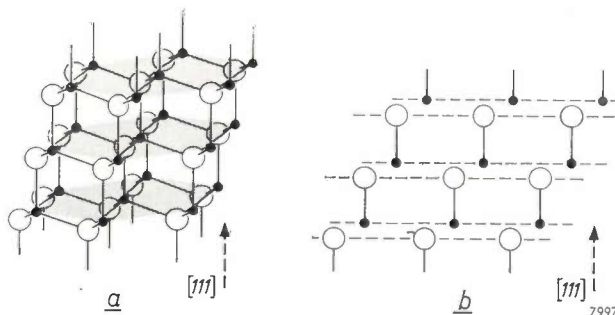


Fig. 1. Crystal structure of cubic ZnS (sphalerite).

- Perspective sketch. Successive (111) planes are occupied alternately by Zn and by S atoms (denoted by spots and circles respectively). The Zn (111) planes are represented as shaded areas. The $[111]$ axis, marked by the arrow, is a polar axis.
- Projection of the lattice onto a (110) plane, i.e. normal to the polar axis. Note the alternate large and small spacings between the successive planes.

The method of power transmission employed here may be compared with that of a transformer: the primary is represented by the loop system, the iron core by the ferroxcube rod, and the secondary by the winding around the rod. The fact that this transformer works with reasonable efficiency in spite of the extremely inefficient "magnetic circuit", is due to the high frequency used and to the properties of the ferroxcube, which combines at such frequencies a high permeability with very low eddy current and other losses.

Models like the one described here have also been made for more serious purposes, in particular for instruction and practice in driving lessons. The same principle could also be applied to the operation of equipment in hermetically sealed rooms or in places in laboratories and factories that are not readily accessible owing to radiation or high-tension hazards. We have already mentioned the advantages compared with more conventional systems, i.e. the elimination of batteries which have to be periodically changed or recharged, and the simplicity of the

equipment. To enable this particular model to perform all its functions with equal flexibility using the normal method of radio control ¹⁾, quite a number of valves or transistors and associated circuit elements would have been required. These conventional systems are of course still necessary where the objects to be controlled have to be more freely manoeuvrable and have to perform more intricate operations, which also calls for higher power.

¹⁾ A. H. Bruinsma, Radio-controlled models, Philips tech. Rev. 15, 281-285, 1953/54.

Summary. A description of a model of a fork-lift truck, in which the driving motors and lifting mechanism draw their power from the magnetic field of an inductive loop enclosing a manoeuvring surface of e.g. 1×2 m. To this end the model carries wire-wound antenna rods of ferroxcube, tuned to the frequency of an alternating current (about 20 kc/s) passed through the inductive loop with the aid of a signal generator. Very simple means are used not only to supply the power but also to control the movements of the model. — This principle of transmitting power to a moving object may perhaps find application in special cases in laboratories.

ETCH PITS ON A ZINC-SULPHIDE CRYSTAL

546.47'221:548.572

The cubic modification of zinc sulphide, also called zinc blende or sphalerite, has a structure resembling that of diamond: all the zinc atoms are tetrahedrally surrounded by sulphur atoms, and vice versa.

Closer examination reveals that the crystal can be regarded as built up of alternate layers of zinc and sulphur atoms, stacked one above the other. In this structure a layer of zinc and a layer of sulphur atoms are one above the other; at a somewhat greater distance in this stacking direction, which is a crystallographic $\langle 111 \rangle$ axis, there again follow a layer of zinc and a layer of sulphur atoms close together, and so on (see *fig. 1*). The zinc and sulphur atom layers after each large spacing are in the opposite sequence if the stacking is continued in the opposite direction. This means that the $\langle 111 \rangle$ axis is a polar axis.

This polarity appears in a surprisingly beautiful way in the macroscopic-chemical behaviour of such a crystal, i.e. when etched. The photographs in *figs 2* and *3* illustrate this: they show the opposite end faces perpendicular to the polar axis (i.e. a (111) plane and a $(\bar{1}\bar{1}\bar{1})$ plane) of the same cubic zinc-sulphide crystal, after etching with hydrochloric acid. It can be seen that one of these crystal faces is covered with etch pits, bounded by planes ($\{100\}$ planes),

from which the triaxial symmetry of the $\langle 111 \rangle$ axis can easily be recognized. On the opposite end face the etch pits are rounded, although here too the triaxial symmetry is still recognizable.

Similar phenomena have been observed by other investigators, both on zinc sulphide and on substances with the same crystal structure, in particu-

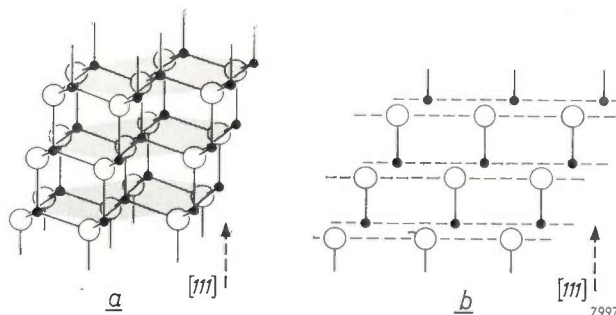


Fig. 1. Crystal structure of cubic ZnS (sphalerite).

- Perspective sketch. Successive (111) planes are occupied alternately by Zn and by S atoms (denoted by spots and circles respectively). The Zn (111) planes are represented as shaded areas. The $[111]$ axis, marked by the arrow, is a polar axis.
- Projection of the lattice onto a (110) plane, i.e. normal to the polar axis. Note the alternate large and small spacings between the successive planes.

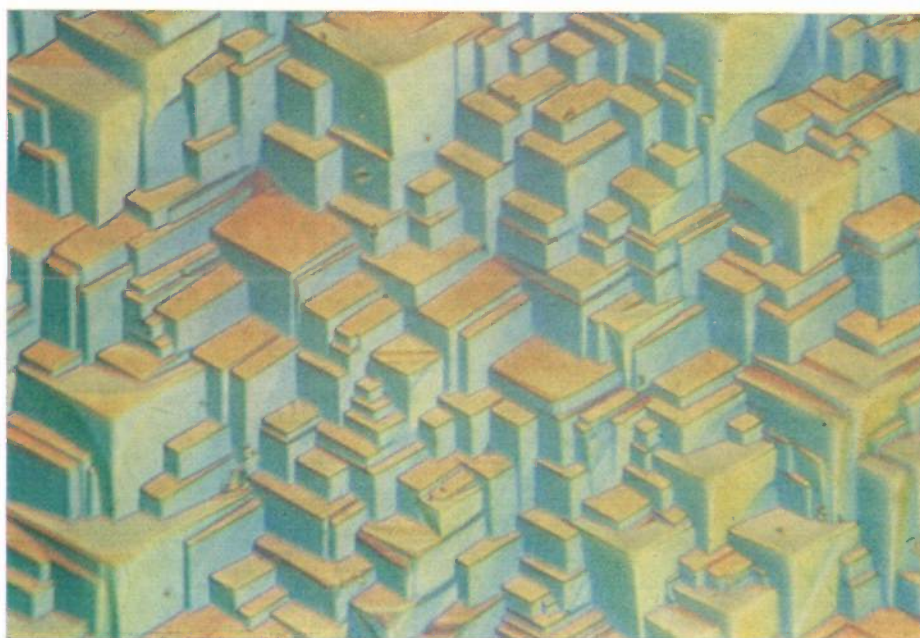


Fig. 2. Photomicrograph of an end face — (111) plane — of a cubic ZnS crystal, etched with hydrochloric acid. Taken by Nomarski's interference-contrast method. Magnification approx. 600 times.

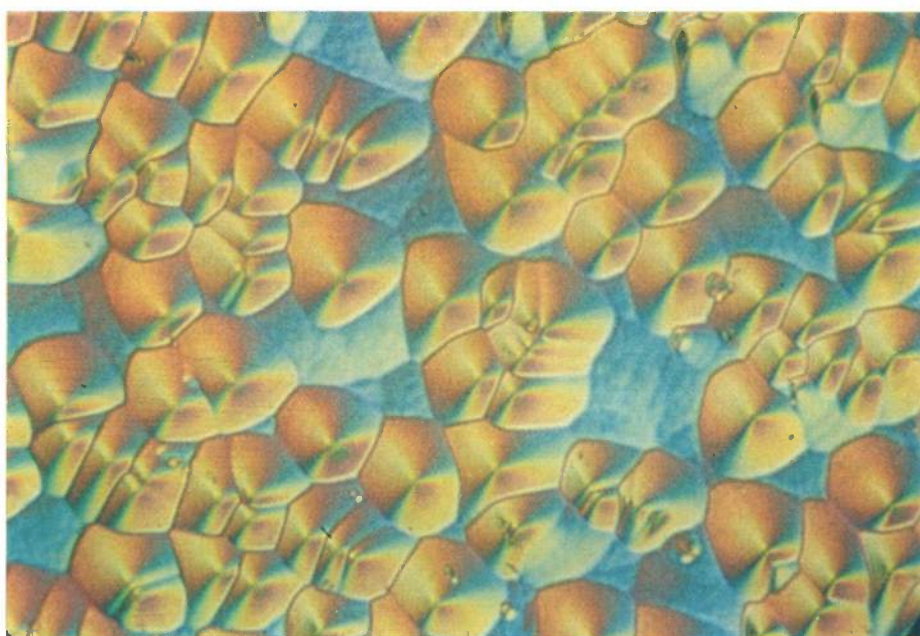


Fig. 3. Photomicrograph of the opposite face — $(\bar{1}\bar{1}\bar{1})$ plane — of the same crystal as in fig. 2, with the same magnification.

lar gallium arsenide ¹⁾ ²⁾ ³⁾. It was found possible in that case to determine at which ends of the polar axis the one or the other etch pattern appears. A

¹⁾ H. A. Schell, *Ätzversuche an Galliumarsenid*, Z. Metallk. 48, 158-161, 1957.

²⁾ J. G. White and W. C. Roth, *Polarity of gallium arsenide single crystals*, J. appl. Phys. 30, 946-947, 1959.

³⁾ J. Woods, *Etch pits and dislocations in cadmium sulphide crystals*, Brit. J. appl. Phys. 11, 296-302, 1960.

crystal like ZnS or GaAs which is somewhat ionic in character and has a structure consisting of two kinds of atoms stacked in alternate layers, must contain equal numbers of both types of atoms in order to comply with the condition of electrical neutrality. Disregarding crystal defects (which will not affect the broad lines of the picture) a crystal of this type is thus bounded at one end of a polar axis

by a plane containing *one* kind of atom, and on the other by a plane containing the *other* kind. Coster, Knol and Prins ⁴⁾ have reported a method of determining the kind of atom in the outer layer by means of X-ray analysis (scattering of X-rays having a wavelength near an absorption edge of one of the atoms of the crystal). Using this method White and Roth ²⁾ have found that in GaAs the "angular" etch pits occur on the Ga side, i.e. where the metal ions

are, and the "rounded" etch pits on the As side. By analogy with this, it seems probable that in our case fig. 2 is attributable to the Zn side, and fig. 3 to the S side of the crystal investigated.

The photographs reproduced here were taken with a Reichert microscope fitted with a Nomarski interference-contrast device ⁵⁾, from which the colours originate.

A. J. ELAND *).

⁴⁾ D. Coster, K. S. Knol and J. A. Prins, Unterschiede in der Intensität der Röntgenstrahlenreflexion an den beiden 111-Flächen der Zinkblende, Z. Phys. **63**, 345-369, 1930.

⁵⁾ G. Nomarski and P. R. Weill, Revue Métall. **52**, 121, 1955.
*) Philips Research Laboratories, Eindhoven.

ABSTRACTS OF RECENT SCIENTIFIC PUBLICATIONS BY THE STAFF OF N.V. PHILIPS' GLOEILAMPENFABRIEKEN

Reprints of those papers not marked with an asterisk * can be obtained free of charge upon application to the Philips Research Laboratories, Eindhoven, Netherlands, where a limited number of reprints are available for distribution.

2966: J. Hornstra: The role of grain boundary motion in the last stage of sintering (Physica **27**, 342-350, 1961, No. 3).

The last stage of the sintering process, when the pores are no longer connected with one another, is limited by diffusion of vacancies from the pores to grain boundaries. It has also been suggested that plastic deformation is a possible fundamental mechanism for this process. It is shown in this article that plastic deformation can indeed play a role, since such deformation and the diffusion of vacancies can interact. The diffusion is, however, the basic mechanism.

2967: E. E. Havinga: The temperature dependence of dielectric constants (Phys. Chem. Solids **18**, 253-255, 1961, No. 2/3).

It is shown that the change in the dielectric constant of an isotropic or cubic material when the temperature is raised can be split into three parts: 1) a decrease due to a fall in the number of polarizable particles per unit volume as a result of the thermal expansion; 2) an increase, because a particle is more easily polarizable when it has more room; 3) a change due to the variation of the polarizability at constant volume with the temperature. A simple method of deriving these three contributions from measurements is discussed. The three contributions are cal-

culated for some of the substances for which the necessary data are available (LiF, KCl, NaCl and BaTiO₃).

2968: J. Jacobs and J. N. Walop: Determination of the concentration of myxoviruses and of some specific antisera, based on neuraminidase activity (Nature **189**, 334-336, 1961, No. 4761).

Description of a method for the determination of the concentration of myxo-viruses, based on the measurement of the enzyme (neuraminidase) activity of the virus. The method has been used for a number of strains of influenza virus (A, A₁, A₂ and B), and for the virus of Newcastle disease (a poultry disease). It has been found that the homologous antibodies against in particular the strains of Asian influenza can be estimated by determining to what extent they inhibit the neuraminidase activity of the virus. These chemical methods of estimation have the advantage that no virological experience is necessary.

2969: H. W. van den Meerendonk and G. G. J. Bos: Assortimentsbepaling van seizoenprodukten (Statistica neerl. **15**, 177-187, 1961, No. 2). (Stock planning for seasonal articles; in Dutch.)

Treatment of the problem of how many of the

various lines of seasonal articles (e.g. seasonal delicacies at Christmas and Easter, winter and summer clothing) a dealer should buy. Two cases are discussed, viz where the distribution of the demand is discrete and where it is continuous.

2970: U. Enz: Magnetization process of a helical spin configuration (J. appl. Phys. **32**, suppl. to No. 3, 22S-26S, 1961).

The magnetic properties of dysprosium can be explained in terms of a model in which the magnetization directions of successive basal planes of the hexagonal structure differ by a constant angle, so that a helical spin distribution arises. The behaviour of such a spin distribution in a magnetic field is calculated, and the results are compared with experiment. A similar spin configuration is proposed for the hexagonal oxide $\text{Sr}_2\text{Zn}_2\text{Fe}_{12}\text{O}_{22}$.

2971: H. Zijlstra: Magnetic annealing of "Ticonal" G magnet steel (J. appl. Phys. **32**, suppl. to No. 3, 194S-196S, 1961).

Short description of an investigation which has been fully dealt with in the author's thesis (2853).

2972: P. Massini: Translocation of 3-amino-1,2,4-triazole in plants, II. Inhibition study (Acta bot. neerl. **10**, 99-104, 1961, No. 1).

Continuation of an investigation whose earlier results have been described in publications 2644 and 2853f. It has been found that the way in which 3-amino-1,2,4-triazole (AT) spreads from the leaves of tomato or bean plants through the rest of the plant depends on the illumination, and that the transport is inhibited by HCN. The transport of AT by the transpiration stream after application to the stem is however not inhibited by HCN. The concentration of HCN needed to block the transport of AT from the leaves was also enough to inhibit the respiration of stem sections.

2973: H. Koopman: U.V. spectra of derivatives of 1,3,5-triazine (Rec. Trav. chim. Pays-Bas **80**, 158-172, 1961, No. 2).

It is known that derivatives of 1,3,5-triazine have a herbicidal action. This publication describes the synthesis of four new 2,4-dichloro-6-(p-substituted phenyl)-1,3,5-triazines. The ultraviolet spectra of these compounds, and also of other compounds described previously (see 2586, 2815 and 2834) are discussed.

2974: L. A. Æ. Sluyterman: A simple method of establishing disulphide interchange in pro-

tein constituents (Biochim. biophys. Acta **48**, 429-436, 1961, No. 3).

Two molecules, both of which consist of groups of atoms connected by two sulphur atoms (disulphide bridges), can under some circumstances exchange groups of atoms, producing new substances; for example, a mixture of two symmetrical disulphides A-SS-A and B-SS-B can give the asymmetrical disulphide: $\text{A-SS-A} + \text{B-SS-B} \rightleftharpoons 2\text{A-SS-B}$. It is sometimes useful to know in connection with the investigation of proteins whether certain conditions have an effect on such exchange reactions. This has been investigated with the aid of paper electrophoresis for the case of oxidized glutathione and its mono- and di-acetyl derivatives. It has been found that substances which contain an SH group act as catalysts for the exchange reaction. Na_2SO_3 and KCN, which have the property of liberating SH groups from disulphides, also act as catalysts; KCN less so than Na_2SO_3 . It was shown with the aid of amperometric titrations that traces of SH groups in oxidized glutathione cause a slow spontaneous exchange; compounds which fix SH groups suppress this spontaneous reaction.

2975: A. C. van Dorsten: Statistical effects in electron microscope image recording in relation to optimal electron-optical magnification and picture quality (Proc. Eur. reg. Conf. on electron microscopy, Delft 1960, Vol. 1, pp. 64-68, publ. Ned. Ver. voor Electronenmicroscopie, Delft).

An investigation of the effect of the electron density on the picture quality in electron microscopy. The author starts from the assumption that the quality of a photograph obtained with an electron microscope, and thus the information which can be obtained from the photograph without undue strain, depends mainly on the contrast. It was found from experiments with series of photographs in which the number of contrast steps gradually increased that nearly all the people involved in the experiments called a photograph "optimal" when it contained 4-6 contrast steps. The number of electrons which must fall on the most heavily exposed image elements in order to ensure a given number of contrast steps follows from a simple statistical calculation. These results hold for all methods of making the image visible: photographic, xerographic, etc. The author then calculates the electron-optical magnification needed with various photographic emulsions in connection with the demands made on the electron density. In practice the magnification must always be greater than this

because of the halo effect. The electron density is thus not a limiting factor for photographic methods of electron microscopy. The calculations given in this publication may however be of value for the investigation of other methods of making the image visible; for example, the author comes to the conclusion that an image intensifier — while it may sometimes be of considerable use — cannot substantially improve the perception of structural details in an electron micrograph, and may even cause a deterioration in the perception of detail.

2976: A. C. van Dorsten and H. F. Premsele: Low-voltage electron microscopy (as 2975, pp. 101-104).

Theory predicts that the quality of electron micrographs of very thin samples of low contrast can be substantially improved by using a low beam voltage (below 15 kV): this increases the contrast considerably. An experimental investigation carried out with a normal electron microscope intended for use at 50 kV but used at 13 kV showed that the expected improvement is in fact found. The choice of film is very important if the best results are to be obtained. Certain modifications had to be made to the electron microscope for work at this low voltage: an insert was used in the anode to keep the field strength at the tip of the cathode up to the desired level, and special measures had to be taken to prevent the formation of surface charges, which have more effect on the electron trajectories at lower voltages.

2977: C. Berghout: Über die Suszeptibilität und den elektrischen Widerstand homogener und inhomogener Kupfer-Eisen-Legierungen (Z. Metallk. 52, 179-186, 1961, No. 3). (The susceptibility and electrical resistance of homogeneous and inhomogeneous copper-iron alloys; in German.)

This publication consists mainly of a report of measurements of the electrical resistance and the magnetic susceptibility of copper alloys containing a very fine iron precipitate. This precipitate is formed by heating the originally homogeneous alloys, containing at most 1% Fe, at temperatures round about 400 °C. In this condition, the face-centred cubic iron, which is normally non-magnetic, is found to possess a ferro-magnetic moment, which however disappears if the sample is heated for too long. If the sample is then subjected to intensive plastic deformation, the aggregates of iron atoms are made smaller, and the precipitate becomes magnetic again. The plastic deformation also gives rise to magnetic anisotropy.

2978: R. van Strik: A method of estimating relative potency and its precision in the case of semi-quantitative responses (Proc. Symp. on quantitative methods in pharmacology, Leiden, May 1960, editor H. de Jonge, pp. 88-100, North-Holland Publ. Co., Amsterdam 1961).

The potency of biologically active substances, e.g. drugs, is usually determined in the following way. The preparation to be tested and a qualitatively similar standard preparation of known activity are administered to groups of test objects, e.g. animals, each preparation in a series of increasing doses, and the (specific) reaction observed after administration is measured. The relationship between the intensity of this reaction and the dose (the "dose-effect curve") is determined separately for each preparation from the experimental data. If the dose is plotted on a logarithmic scale, two parallel curves should be obtained, provided the two preparations are qualitatively similar. The horizontal distance between the two curves represents the potency ratio of the two preparations, from which the activity of the unknown can be calculated. In practice, the log-dose-effect curve can often be represented by a straight line over a sufficiently wide range. It is then possible to calculate the potency ratio with its confidence limits in a rather simple way.

Quantitative measurement of a continuously variable reaction is not always possible however. In such cases, the reactions are often classified in a series of ordered categories as e.g. —, +, ++, +++, etc. This publication describes a method which allows the computation of the potency ratio and its confidence limits under such circumstances. The method consists in arranging all the observed reactions together in order of increasing intensity, instead of classifying them into categories. The potency ratio with its confidence limits can be calculated from the rank numbers thus obtained, by applying standard bio-assay computations. The reliability of the method, which is illustrated by an example, and its limitations are discussed at length, together with some of the difficulties that may arise in practical application.

2979: T. Kralt, H. D. Moed, V. Claassen, Th. W. J. Hendriksen, A. Lindner, H. Selzer, F. Brücke, G. Hertting and G. Gogolak: Reserpine analogues (Nature 188, 1108-1109, 1960, No. 4756).

Reserpine is a naturally occurring alkaloid with a sedative effect; it is also useful as a remedy for

high blood pressure. This publication is a preliminary report of the synthesis of a number of compounds whose chemical structure is similar to that of reserpine, and of an investigation of their possible medical applications. Some of the compounds are indeed very promising for use as drugs.

2980: J. H. Stuy: Radiation inactivation of intracellular transforming deoxyribonucleic acid (thesis Utrecht, May 1961).

The hereditary properties of a living cell are determined by the structure of the various parts of very long molecules of deoxyribonucleic acid (DNA), which are present in the nucleus of the cell (in the chromosomes, if there are any). These DNA molecules have a structure rather like a ladder: the sides, which spiral round each other, are made of alternate sugar (deoxyribose) and phosphate groups, while the rungs are provided by organic bases. DNA can be extracted from bacteria, but the molecule probably breaks into pieces in the process. If such extracted DNA is brought into contact with other bacteria of the same species, it is possible under certain conditions for the bacteria to exchange parts of their DNA molecules with corresponding fragments of the added DNA which has been taken up by the cell; this phenomenon is known as bacterial transformation. The transformation can be quantitatively investigated by giving the extracted DNA one or more properties ("markers") which the receiving bacteria lack (e.g. resistance to a given antibiotic).

Use is made of bacterial transformation in the investigation of the lethal damage done to bacteria by radiation described in this thesis. It is found that both ultraviolet (non-ionizing) radiation and X-rays (which do ionize) damage DNA in the living cell to the anticipated extent. There are indications that the cell is able to repair the damage done by ultraviolet radiation, but an important result of exposure to X-rays is considerable breakdown of DNA after irradiation. This is probably due to the activation of an enzyme, deoxyribonuclease (DNase), which splits DNA and which is present in an inactive form in each cell. Further measurements have shown that DNA produced in cells which have lost the ability of unlimited division (such cells are said to be inactivated or dead) is biologically normal.

The effect of various kinds of radiation has been investigated with the aid of a system of three coupled markers (i.e. markers which are all on the same fragment of DNA). It has been found that the effects of electron beams are very similar to those of the enzyme DNase. It is concluded from this that the

main effect of electron beams (and also of X-rays) is the breakage of ester links between phosphate and deoxyribose. The principal effect of ultraviolet radiation is a degradation of the organic bases.

2981: W. Albers: Diffusion of arsenic in germanium from the vapour phase (Solid-state electronics 2, 85-95, 1961, No. 2/3).

Investigation of the manner in which arsenic penetrates a piece of germanium placed in arsenic vapour. Part of a quartz vessel was maintained at a temperature of 200 °C; a small piece of arsenic was placed in this part. A piece of germanium was placed in another part of the vessel, which was at a higher temperature (750, 800, 850 and 900 °C in successive experiments). After a certain time (usually 2¼ hours) the arsenic concentration in the germanium was determined as a function of the distance below the surface. The diffusion constant D was determined as a function of the temperature from the results of these measurements. The variation of D with temperature T obeys the relationship $D = D_0 \exp(E/kT)$, $D_0 = 3 \text{ cm}^2\text{sec}^{-1}$ and E (the activation energy of the diffusion) = 56 kcal mole⁻¹. The actual value of the arsenic concentration, e.g. at the surface of the germanium, was found to depend considerably on the experimental conditions. If the vessel contained only arsenic vapour, it made a great difference if the two parts of the vessel which were maintained at the different temperatures were connected by a capillary or by a wide pipe. A narrow capillary hinders the arsenic vapour near the germanium in interacting appreciably with the walls of the cold part of the vessel. The arsenic molecules in the vapour are then much more strongly dissociated than when the connecting part is a wide pipe, which aids this interaction. As a result of this, the arsenic concentration in the germanium is much lower (by a factor of more than 10) when a wide connecting pipe is used than with a capillary. If the vessel contains an inert gas, this hinders the interaction with the cold wall even when a wide connecting pipe is used, so that in this case no lowering of the arsenic concentration is found.

The means used to nullify the disturbing effect of "thermal conversion" are described.

2982: J. A. Kok and C. E. G. M. M. van Vroonhoven: Aspects of electrical breakdown of liquid insulating material, II (Appl. sci. Res. B 9, 125-132, 1961, No. 2).

Electrical breakdown of insulating liquids can be prevented by the addition of soaps or resins to the liquid in question. These are adsorbed by impurities

floating in the liquid, and prevent these impurities from flocculating, when they may give rise to breakdown. It has been found that if soap and resin are added together in certain proportions (the soap may also be formed by gradual saponification of insulating oil) immediate flocculation (and thus rapid breakdown) may result. The operation of aromatic stabilizers and inhibitors, most of which were developed by dye research, is explained. See also 2641 and 2554.

2983: G. Krijl and J. L. Melse: Calorimetric determination of metal coating thicknesses on small objects (Trans. Inst. Metal Finishing 38, 22-26, 1961, part I).

Many small metal-plated articles are mass-produced nowadays. It is desirable to be able to measure the average thickness of the metal coating on such objects rapidly and accurately during the manufacturing process. The authors have developed a method for determining this by measuring the amount of heat evolved when the coating is dissolved off in a suitable solvent. Use is made of a simple calorimeter specially developed for the purpose. The calorimeter is calibrated for a given object by determining the heat evolved when a metal coating of known thickness is dissolved. A great advantage of this method is that no balance is required. It has been used with success for zinc and cadmium on steel, nickel on stainless steel, and silver or phosphor bronze.

2984: J. A. W. van Laar: Die Unterrostung von lackiertem Stahl—Entstehung, Bestimmung, Bekämpfung (Dtsch. Farben-Z. 15, 56-67 and 104-116, 1961, Nos 2 and 3). (Under-rusting of painted steel — causes, determination and prevention; in German.)

If steel is painted properly, it will not rust. If however the paint is removed locally, e.g. by scratching, the rusting process can continue under the paint, and may also cause damage to the paint at the same time. This phenomenon is called "under-rusting". This publication describes an extensive investigation of under-rusting under various conditions. Since the normal standards for estimating damage due to rust are not applicable here, new ones were developed for this purpose, based on the rate at which under-rusting proceeds at right angles to a scratch produced in the paint in a specified way.

The theoretical and practical aspects of under-rusting are compared with those of two other important modes of attack, viz blistering and loss of adhesion of the paint film due to the influence of water and of brine.

2985: J. G. van Santen and G. Diemer: Photorectifier based on a combination of a photoconductor and an electret (Solid-state electronics 2, 149-156, 1961, No. 2/3).

See Philips tech. Rev. 23, 310-315, 1961/62 (No. 10).

2986: M. T. Vlaardingerbroek and U. Weimer: Wiselwerking van plasma en elektronenbundel (Ned. T. Natuurk. 27, 207-212, 1961, No. 6). (Interaction of a plasma with an electron beam; in Dutch.)

If an electron beam modulated by a high-frequency signal is passed through a gas-discharge plasma, the modulation may increase under certain circumstances, which offers a possibility of amplification. This publication gives a discussion of the mechanism of the growth of the waves in such a system, and also describes an experiment by which the growth of the waves is demonstrated.

2987: W. K. Hofker: Stralingsdetectie met *P-N*-overgangen in halfgeleiders (Ned. T. Natuurk. 27, 213-217, 1961, No. 6). (Radiation detection with *P-N* junctions in semiconductors; in Dutch.)

Short description of the principle of operation and the properties of a new type of radiation detector, which can be used with especial advantage for the detection of heavy particles (fission products, α -particles, deuterons).

2988: W. van Gool: Fluorescence centres in ZnS (thesis Amsterdam, Jan. 1961).

In this thesis are discussed the defect chemistry of ZnS and the relationship between the fluorescence of this substance and its lattice imperfections. As an introduction to the experimental part, some well-known aspects of fluorescent ZnS are mentioned.

The experimental part describes first the preparation of fluorescent ZnS and photoconductive CdS. Special attention is paid to phosphors which contain only a coactivator (Al, Cl), and to phosphors with equal amounts of activator (Cu, Ag or Au) and coactivator (Al, Sc, Ga or In). The spectral distribution of the fluorescent light is measured and the reproducibility of the measurements is discussed. The fluorescence and the glow curves of a series of phosphors containing increasing amounts of dope is measured at various temperatures.

The theoretical part begins with a discussion of the defect chemistry of ZnS. This is followed by a thermodynamic description of the interaction of oxygen with ZnS and CdS. The calculation is split up into three parts relating to: the macroscopic

phase diagram, the composition of the atmosphere and the defect chemistry of the sulphides. The conditions under which one can make powder with a constant concentration of oxygen in the crystal lattice and a varying concentration of sulphur vacancies, and *vice versa*, are then deduced.

The rest of the thesis is concerned with the interpretation of the experimental results in terms of the defect chemistry and the positions of the energy levels. If the concentrations of activator and coactivator are not equal, there are many conceivable ways in which these impurities could be taken up in the crystal lattice. The experimental data at present available are not enough to allow us to choose between the various possibilities. The main reason for this is the fact that our knowledge of vacancies and interstitial atoms is still very slight.

If the concentrations of activator and coactivator are equal, more positive conclusions can be drawn. It is assumed in the interpretation of these results that the added impurities are associated in the crystal lattice. The model suggested is sufficient for a qualitative explanation of the measured fluorescence, its variation with the temperature, and the glow curves. Other interpretations which have been given in the literature are also discussed. The last chapter suggests which further investigations would yield most information about the fluorescence and the defect chemistry of ZnS.

2989: C. M. Hargreaves: On the growth of sapphire microcrystals (J. appl. Phys. **32**, 936-938, 1961, No. 5).

Sapphire (α -Al₂O₃) microcrystals can be made by oxidizing aluminium in a stream of hydrogen which contains a little water vapour. Some of the crystals produced in this way have the shape of extremely thin needles, or "whiskers", while others are very thin platelets. In many substances, crystals of these shapes grow at low supersaturations by the deposition of material on the steps formed by screw dislocations on the surface. It is not however certain whether this growth mechanism applies to the pure Al₂O₃ crystals grown in this way. The optical and electron-microscopic investigations mentioned in this paper provide no evidence of spiral growth patterns associated with screw dislocations. It follows from thermodynamic calculations that — in contrast to a previous statement in the literature — the transport of Al via the vapour phase during the growth of the sapphire crystals most probably occurs as Al₂O and not as AlO.

2990: W. Hondius Bolding: Quality and choice of Potter-Bucky grids, IV. Focus-grid distance limits, V. The contrast improvement factor (Acta radiol. **55**, 225-235, 1961, No. 3).

A discussion of the use of Potter-Bucky grids for improving the quality of X-ray photographs; continued from **2793**. In Part IV, the limits between which the distance from the tube focus to the grid can be varied is discussed. As a rough measure of the improvement which can be obtained with such a grid, the lead content of the grid may be used (see Part III, **2793**). The contrast-improvement factor defined in Part V is however a better measure, though more difficult to determine.

2991: G. D. Rieck and H. A. C. M. Bruning: Thermal diffusion of oxygen and nitrogen in zirconium (Nature **190**, 1181-1182, 1961, No. 4782).

Preliminary report of an investigation of the behaviour of oxygen and nitrogen dissolved in zirconium. Such solutions are found to exhibit the Ludwig-Soret effect, i.e. the dissolved gas moves towards the colder part of the metal if a temperature gradient is present.

2992: A. R. W. Muyen: Optimum lot-size policy if tools break down frequently (Operat. Res. Quart. **12**, 41-53, 1961, No. 1).

A well-known formula gives the optimum length of a production run, so that the cost of interrupting production is balanced against warehouse costs. This formula does not consider the possibility that the tools used in the production may break. Let us suppose that there is an appreciable chance that tool breakage will in fact occur before the production run has reached the optimum length mentioned above. It is shown in the present publication that in that case one can determine two numbers Q_1 and Q_2 ($Q_1 < Q_2$) such that if a tool breaks down after more than Q_1 but less than Q_2 articles have been produced, it is best to stop the run. If a tool breaks before Q_1 , it is best to replace the tool and carry on. The run should in any case be stopped after the production of Q_2 articles. The values of Q_1 and Q_2 are derived, and the results expressed in graphical form. In this derivation it is assumed that there are an infinite number of tools in reserve. A Monte Carlo calculation shows however that the values of Q_1 and Q_2 found in this way also apply for a finite stock of tools.

Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES

A NEW VASODILATING PHARMACEUTIC

by H. D. MOED *).

615.717

*In vascular disorders like arteriosclerosis, vital functions of the human organism are endangered because the blood is unable to circulate freely through certain parts of the body. There is accordingly a demand from the medical side for a drug that is capable of restoring freedom of circulation where this has been impaired. This has given impetus to a whole series of investigations in various countries; in the Philips research laboratories, Netherlands, it has led to the development of "Duvadilan" **), which has quickly been adopted for large-scale medical employment as a vasodilating and spasmolytic agent. The following article seeks to give an account of these investigations, which may be regarded as being of pharmacological and chemical as well as medical interest.*

Introduction

The body contains a great variety of substances that influence the activity of organs and other functional systems. It is one of the functions of the pharmacologist to ascertain the actions of such substances and the mechanisms underlying their effects. A common procedure in pharmacological research is to modify the molecule of the active substance in various ways and to investigate the associated changes in its pharmacological properties. This knowledge may be of value when a drug with a given action is being sought. Investigations of this kind, carried out on adrenaline, have yielded results that have an important bearing on the treatment of certain vascular disorders such as arteriosclerosis in the head and limbs, and B rger's and Raynaud's diseases.

Adrenaline is a hormone secreted by the adrenal gland and possesses a particularly complex set of actions in the body, to which we will start by devoting some attention.

- 1) Adrenaline has an effect on the blood-vessels, *dilating* the vessels in certain kinds of tissue (skeletal muscle, for example) and *constricting* those in other kinds (the skin and the intestines).
- 2) It also affects the heart, causing tachycardia (speeding up the heart rate) and strengthening the contractions.

- 3) It has a relaxant effect on the smooth musculature of the bronchial passages and of the intestine.
- 4) It plays a part in sugar metabolism, the effect being to raise the concentration of sugar in the blood.

A certain pattern can be seen in these many effects if adrenaline is regarded as the hormone that makes the body capable of sudden great exertion. Adrenaline has in fact been called the "escape hormone". The skeletal muscles used in a fast escape require an extra generous supply of blood; hence the dilation of their blood-vessels and the reinforcement of heart action. Other parts of the body that are of secondary importance for the time being have to make do with a smaller supply; hence the vasoconstriction in the region round the intestine and the relaxation of muscles in the intestinal wall. Dilation of the bronchial passages, resulting in increased oxygen uptake, is another desideratum for flight.

Important as it may be for self-preservation in man and animals, the mere fact of this complex set of actions, together with the difficulty of controlling them properly, means that adrenaline is not entirely suitable for therapeutic purposes. Its vasodilatory effect could be very useful in the treatment of the diseases named above, in which the free circulation of the blood is impeded by the spastic state of the vessels, for example, or by thickening of their walls; but the accompanying effects on the heart (palpitation) would be troublesome and sometimes even harmful. Another disadvantage is that adrenaline is ineffective when given orally.

*) N.V. Philips-Duphar, Weesp, Netherlands.

***) Registered trade-mark by N.V. Philips-Duphar and sold under this trade-mark in most countries. In some countries "Duvadilan" is sold under a different trade-mark, e.g. "Vasodilan", "Dilavase" and "Cardilan".

Extensive research at the Central Laboratory of N.V. Philips-Duphar, Weesp (Netherlands), has resulted in the synthesis of isoxsuprine*), a derivative of adrenaline which is marketed under the name of "Duvadilan". This preparation has proved to be a particularly useful one in that it induces marked vasodilation when given orally (or by other routes) while having little effect on the heart. It is accordingly of considerable value in medicine.

In the section which follows, the actions and molecular structures of adrenaline and isoxsuprine will be compared with those of some other adrenaline derivatives in order to show more clearly the specific properties of isoxsuprine. Attention will then be given to the synthesis of the drug, and to the mechanism of the relevant chemical reaction. Finally, clinical data and the results of animal experiments will be given to illustrate the pharmacological actions of isoxsuprine.

Adrenaline derivatives

Relationships between adrenaline, noradrenaline and the sympathetic nervous system

As we have learned from the above, one sometimes wants to retain one effect of a substance while eliminating another. The belief that this is possible implies the assumption that each of the given effects is the result of action on an appropriate *receptor*, and that the structural elements of the substance essential for its action on one kind of receptor are different from those essential for its action on another kind. Receptors can be thought of as regions of tissue with a specialized structure on which the drug takes hold. There is experimental evidence in support of the idea that different kinds of receptors exist, and that they work independently. Firstly, altering the chemical structure of a substance may modify its individual effects *in different ways*. Secondly, it is possible for *some* of the effects to be suppressed by a second substance, the other effects persisting.

Convincing evidence of the first kind can be obtained by comparing the effects of adrenaline with those of noradrenaline and isopropylnoradrenaline, which are substances chemically related to adrenaline and are classed as *sympathomimetic* (this term will be explained later). Noradrenaline is another hormone secreted by the adrenal gland, and isopropylnoradrenaline is a compound derived from it. The effects of these three compounds are compared in *Table I*; from this it is clear that noradrenaline only *partly* shares the effects of adrenaline, and that isopropyl-

Table I. Some effects of noradrenaline, adrenaline and isopropylnoradrenaline.

	noradrenaline	adrenaline	isopropylnoradrenaline
Vasodilation (skeletal muscles)	0	+	+
Vasoconstriction (skin, intestines)	+	+	0
Bronchodilation	0	+	+
Tachycardia	0	+	+

noradrenaline produces precisely those effects of adrenaline which are not shared by noradrenaline. On the basis of such observations Ahlquist¹⁾ postulated the existence of two types of sympathomimetic receptors which he called α and β .

In order to understand the meaning of the name "sympathomimetic", let us look into the links between adrenaline, noradrenaline and the function of the nervous system. A number of vital functions, including metabolism and blood flow, are to some extent under the control of the autonomic nervous system. It is a dual form of control exercised by the two parts, sympathetic and parasympathetic, into which the autonomic system can be divided. In general, stimulation of the sympathetic system results in enhanced activity, whereas stimulation of the parasympathetic system results in relaxation and recuperation.

In both these systems stimuli are transmitted from nerve cell to nerve cell and finally from nerve cell to tissue by means of chemical substances. In both it is acetylcholine that passes on the stimulus from one nerve cell to the next, and in the parasympathetic nervous system the same chemical is responsible for transmission to the tissues; but in the sympathetic system this latter responsibility is reserved for noradrenaline and, in a lesser degree, for adrenaline.

Substances whose action on the tissues produce effects like those of stimulation via the autonomic nervous system are classed as "sympathomimetic" or "parasympathomimetic", depending on the stimulated systems ("mimetic" meaning "mimicking"). The opposite categories of "sympatholytic" and "parasympatholytic" include substances with the

¹⁾ R. P. Ahlquist, Amer. J. Physiol. 153, 586, 1948; see also, for example, E. J. Ariëns, Modern concepts in relationships between structure and biological activity, Symposium VII, 1st International Pharmacological Meeting, Stockholm 1961.

*) Generic name registered at the World Health Organization.

ability of *suppressing* the effects in question ("lytic" meaning "loosing" or "releasing from").

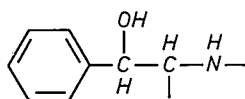
Let us now return to the postulation of Ahlquist regarding the existence of two types of sympathomimetic receptors. Vasoconstriction in the skin is thought to be the result of sympathomimetic substances acting on the α -receptors and is therefore classed as α -sympathomimetic; β -sympathomimetic effects include vasodilation in skeletal muscles and dilation of the bronchial tubes. Noradrenaline is believed to act on α -receptors only, adrenaline on both types, and isopropylnoradrenaline on β -receptors only.

For brevity, we shall refer to the above as α -mimetic and β -mimetic substances. The terms " α -lytic" and " β -lytic" are also used with reference to the opposite categories — see above. The theory is that these lytic substances attach themselves to α -receptors or β -receptors without giving rise to the appropriate effect; by engaging with or occupying the receptors, it is thought, they prevent adrenaline and noradrenaline exercising their particular sympathomimetic effects.

Now that this background has been sketched in, it can be seen that a big advance will have been made in the direction of a "specific" vasodilating agent if a substance can be produced which combines a β -mimetic with a α -lytic action. The line followed in our investigations will be easier to understand if this is kept in mind. Other factors that had an important bearing on these investigations are mentioned in the next section, in which the structure of isoxsuprine is described.

The chemical structure of isoxsuprine as compared with the structures of adrenaline, noradrenaline and other related substances

Fig. 1 shows the structural formulae of adrenaline, noradrenaline and isoxsuprine. As can be seen, they all contain the group



This group appears to be essential if a sympathomimetic effect is to be obtained. In this connection *ephedrine* (fig. 1) is of some interest. This is a naturally occurring alkaloid prepared from the herb known to the Chinese as Ma Huang. In 1923 Chen and Schmidt²⁾ showed that ephedrine has much the

same properties as adrenaline. But there is an important difference: adrenaline is readily broken down, and for that reason is ineffective when given orally, whereas ephedrine is relatively free from this disadvantage. This is probably due to two things, the first being the methyl group which has been substituted on to carbon atom β (see fig. 1) and which serves in some degree to protect the neighbouring nitrogen from enzymatic action (which is one of the ways to breakdown the molecule). Secondly, the elimination of hydroxyl groups attached to the benzene nucleus also reduces the risk of breakdown, though in a different way. At the same time, unfortunately, the loss of hydroxyl groups greatly reduces the sympathomimetic activity of the molecule.

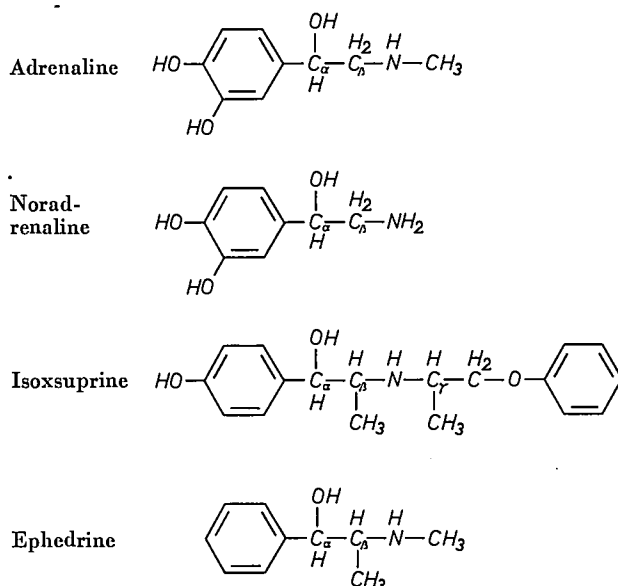


Fig. 1. The structural formulae of adrenaline, noradrenaline, isoxsuprine and ephedrine.

Mention has already been made of the striking differences which exist between the effects of adrenaline, noradrenaline and isopropylnoradrenaline. These differences indicate that the group to be substituted on to the nitrogen atom has an important influence on the sympathomimetic activity of the resulting compound. The question has been carefully investigated in experiments in which the methyl group in the adrenaline molecule was replaced by various other alkyl groups. In general this was found to result in a decrease in α -mimetic potency *without any loss of β -mimetic action*. In some cases there was even reversal of the α -mimetic action to an α -lytic one. In these investigations³⁾ the compound that attracted the most attention was the previously mentioned isopropylnoradrenaline. This compound

²⁾ K. K. Chen and C. F. Schmidt, J. pharmacol. exper. Ther. 24, 339, 1924.

³⁾ H. Konzett, Arch. exper. Pathol. 197, 41, 1940/41.

has a bronchodilating effect three times as great as that of adrenaline. It is also a potent vasodilator. It is entirely free from the α -mimetic effects of adrenaline on the nerve cells, but still has the undesirable effect of stimulating the heart, particularly when given orally. It affects the heart to a much less extent when inhaled and in this form it is very useful for relieving asthma.

The search for other groups suitable for attachment to the nitrogen was influenced by the thought that it might be possible to derive a substance which would relax vascular spasm in virtue of properties other than and additional to β -mimetic activity. K \ddot{u} lz⁴⁾ was the first to substitute phenylalkyl groups on to the nitrogen. Fig. 2 shows the chemical structure of the substance synthesized by K \ddot{u} lz, together with that of papaverine. The latter is a spasmolytic drug possessing a direct action on muscle. Although K \ddot{u} lz states that he was interested in a spasmolytic action, it is not certain that the analogy with papaverine was uppermost in his mind. Be that as it may, the molecule he obtained can be regarded as an opened-out papaverine molecule. The molecule obtained also has a powerful α -lytic action.

In the synthesis of isoxsuprine, instead of K \ddot{u} lz' phenylalkyl a *phenoxyalkyl* group was attached to the nitrogen, a similar structural change having produced an increase in the α -lytic activity of related substances.

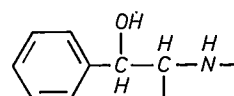
An isopropyl was chosen to give additional protection to the nitrogen atom by means of a second methyl group.

So far no mention has been made of the fact that isoxsuprine is a *racemic mixture (racemate)*, i.e. a

mixture of two kinds of molecules whose structures are mirror images of each other (see below). Nor has it been mentioned that three other racemic mixtures exist whose components have the same structural formula as those of isoxsuprine. This last fact especially has an important bearing on the preparation of an agent having the desired therapeutic properties. As in the preparation of many related substances, these substances tend to occur together. In this case, since they exhibit big differences in their value as drugs, it was important to separate them. Isoxsuprine proved to be the most powerful vasodilator and to have the least effect on the heart. The difficulties involved in the preparation form the subject of the next section.

We end this section with a summary of the characteristics of isoxsuprine as follows:

a) It has a β -sympathomimetic action due to the group



b) It has acquired an α -sympatholytic action in virtue of the substituted phenoxyalkyl group.

c) It has a spasmolytic action attributable to built-in analogies with the structure of papaverine.

d) It is effective when given orally owing to the elimination of a hydroxyl group and the protection of the nitrogen by the two methyl groups.

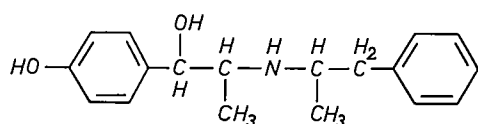
e) Its effect on the heart is slight.

The following section is mainly concerned with the chemical aspect of the investigations; in the final section of the article, we shall return to the pharmacological aspects.

Preparation of isoxsuprine by asymmetric synthesis

Sometimes the chemist is faced with the problem of synthesizing a substance without at the same time producing other substances with the same chemical composition (isomers). The same problem arises in an acute form in cases where the substances in question have the same structure as well as the same composition, and differ only in the spatial arrangement or configuration of the component atoms within the molecule (stereoisomers).

To take a simple example, let us consider the reduction of butanon-2 to butanol-2. Figs 3a, b and c represent molecules of these substances. The atom shown occupying a central position in the structure of butanol-2 is known as an asymmetric carbon atom; it is possible for the surrounding atoms



Phenylisopropyl-noroxypheдрine

Papaverine

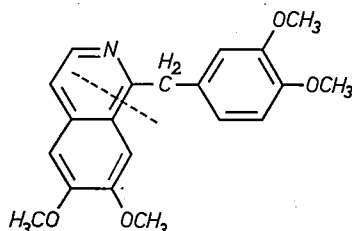


Fig. 2. The structural formulae of phenylisopropyl-noroxypheдрine and papaverine. The close similarity between the two structures will be clearer if the papaverine molecule is imagined to be cut at the point indicated by the dotted line, and then stretched out.

⁴⁾ F. K \ddot{u} lz and M. Schneider, *Klin. Wochenschrift* **28**, 535, 1950.

obtained, and these could be separated by fractional crystallization. However, the isoxsuprine yield would be much smaller.

*Mechanism of the stereospecific reduction in isoxsuprine production*⁶⁾

The asymmetric or stereospecific reduction in the production of isoxsuprine is of theoretical as well as of practical interest; but first it will be as well to give names to the four possible racemates. All have the same chemical structure. We shall take the symbol I, the initial letter of isoxsuprine, to denote this structure, adding some qualification to show which of the various configurations around carbon atoms α , β and γ is being referred to. The accepted prefixes erythro- and threo-, indicating configurations similar to those of erythrose and threose, two well-known sugars, would be suitable for the structures arising around α and β , the two carbon atoms that are close together. The atoms of isoxsuprine, it has been found, are arranged around α and β in the same way as the component atoms of erythrose, and we can therefore give isoxsuprine the name erythro-I. There is a second racemate with the same configuration around α and β as erythro-I; the only difference lies in its opposite configuration around carbon atom γ . We shall refer to this second racemate as alloerythro-I. Accordingly, we can apply the names threo-I and allothreo-I to the two other racemates, which have "threo" configurations around carbon atoms α and β . A diagrammatic explanation is given in fig. 4.

Some processes of synthesis whereby these configurations can be obtained are shown schematically in fig. 5. It can be seen that, with the exception of IV, each of these processes yields either "erythro" or "threo" configurations. This is only to be expected if carbon atom β does, in fact, have a certain effect on the steric course of the reduction process, in the manner described above. As the carbon atom β is unable to influence the configuration of the more distant carbon atom γ , the "allo" forms are always present as a sort of by-product.

It will be clear from the diagram that the course of the synthesis depends on (1) the number of groups attached to the nitrogen atom in the molecule to be reduced, and (2) on the way reduction takes place, i.e. on whether it is reduced by LiAlH_4 or, alternatively, by hydrogen, palladium and carbon. In the order in which the reactions are shown in fig. 5, their products show a complete reversal (from "erythro"

to "threo" configurations). For the formation of "threo" configurations it appears to be necessary for the nitrogen atom to have two large groups attached to it; moreover, the reducing agent must be LiAlH_4 . With a view to explaining this reversal let us try first of all to visualize how carbon atom β exercises its influence under conditions of reduction with LiAlH_4 .

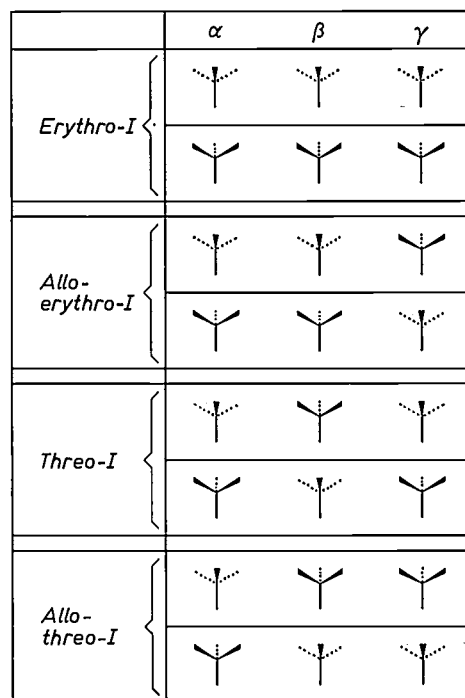
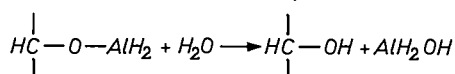
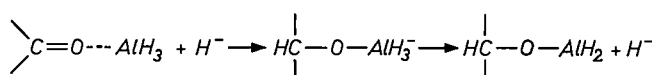
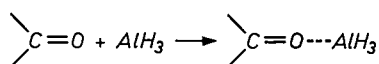
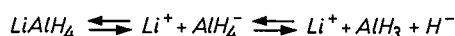


Fig. 4. The isomers that are of interest in the synthesis of isoxsuprine (here called erythro-I). Special symbols have been used to represent configurations which are mirror images. The configurations in question arise around the three asymmetric carbon atoms denoted by α , β and γ . $2^3 = 8$ combinations are possible; these eight stereoisomers constitute four pairs of enantiomers, each pair occurring in a racemate. The names "erythro-I" and "threo-I" have been given to the pairs arising out of combinations between α and β which show structural analogies with erythrose and threose. The prefix "allo" serves to distinguish the two pairs arising out of further combinations with γ .

The reduction reaction of a carbonyl (CO) group with LiAlH_4 is fairly well known. A hydride (H^-) ion combines with the carbon atom, and an aluminium complex attaches itself to the oxygen. The reaction may take place in the following stages:



⁶⁾ J. van Dijk and H. D. Moed, Rec. Trav. chim. Pays-Bas 78, 22, 1959.

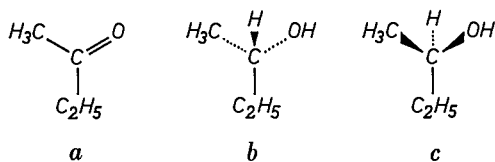
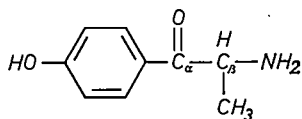


Fig. 3. Diagram a) shows the structure of butanon-2; b) and c) represent the two possible three-dimensional configurations of butanol-2. Atoms and radicals joined to the central carbon atom by dotted lines are to be imagined as occupying positions under the plane of the drawing. The wedge-shaped linking symbols indicate that the atom or radical in question is above the plane of the drawing.

to take up either of two spatial arrangements which are mirror images of each other (and which are optical antipodes or enantiomers, this meaning that they rotate the plane of polarisation of transmitted light through the same angle in opposite sense). For the butanon-2 molecule the possibility of alternative configurations does not arise; the opportunity of assuming one configuration or the other occurs during the reduction process, when the hydrogen atom is attaching itself to the carbon atom shown in a central position. The stereoisomer of butanol-2 formed depends on which side the hydrogen atom approaches the butanon molecule. It will be fairly obvious that nothing can be done to "guide" the hydrogen atom; the direction from which it comes is purely a matter of chance. This means that carried out on a large scale, the reduction process yields equal numbers of the molecules of the two stereoisomers — a racemic mixture, in other words. The conditions for *asymmetric synthesis* — a process which will yield only one of the possible isomers — must be such that attachment to the molecule is made more difficult for a hydrogen atom approaching from one direction than for a hydrogen atom approaching from the other. The only agency capable of impeding the approach of the hydrogen atom from one of the two directions is the introduction of a *second asymmetric atom in the vicinity*. The preparation of isoxsuprine will serve as an illustration of this principle of asymmetric synthesis⁵⁾.

A substance with the structure

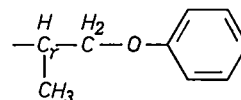


can be taken as starting material for isoxsuprine. A possible first step would be to reduce the carbonyl group. Carbon atom α will become asymmetric owing to the reduction process, as does the central

carbon atom in butanon-2. Carbon atom β is already asymmetric and it appears to decide the outcome of the reduction process, determining which of the two possible configurations will arise around the new asymmetric carbon atom α .

A small complication is added by the fact that the starting substance is itself a racemate, whose two components differ in their spatial arrangement around carbon atom β . The two configurations influence the reduction process in opposite ways, and consequently this process yields two structures which are mirror images — another racemate, in fact, consisting of two optical antipodes. Thus, in this more complicated case, the asymmetric synthesis decides the outcome of one special racemate instead of the two possible ones.

The next stage in the synthesis of isoxsuprine is to attach the group



to the nitrogen atom; this involves the addition of another asymmetric atom to the two already present. The third asymmetric atom is too far away from the others to allow the principle of asymmetric synthesis to be applied, therefore at this stage of the synthesis ordinary methods are employed. The result is a mixture of two different racemates, which are now separated by fractional crystallization. Separation by this method is possible owing to differences in solubility between the two racemic mixtures — differences which do not exist between the optical antipodes of which each racemic mixture is composed. The question of solubility is linked to that of structure: the separable molecules are *not* mirror images of one another, whereas the components of each mixture *are*. (The former are accordingly classified as diastereoisomers, the latter as enantiomers.)

To sum up, the molecule to be synthesized contains three asymmetric atoms α , β and γ ; the possible number of stereoisomers is therefore eight. Since these are paired, each pair consisting of mirror-image isomers or enantiomers, four different racemates can be assumed to exist. Isoxsuprine, which is one of the racemates, has been prepared on the one hand by an asymmetric reduction process which excludes two of the possible racemic mixtures, and on the other hand by fractional crystallization, by means of which the two remaining ones were separated.

As can now be inferred, it is possible to synthesize isoxsuprine without using asymmetric reactions. A mixture of all four racemic mixtures would be

⁵⁾ For a more detailed treatment of asymmetric synthesis, see for example V. Prelog *et al.*, *Helv. chim. Acta* 36, 308, 320, 325, 1953, and D. J. Cram and F. A. Abd Elhafez, *J. Amer. Chem. Soc.* 74, 5828, 1952.

Opinions are still divided as to the nature of the aluminium complex, but there is no doubt that it is a bulky one. By making a model of the molecule that undergoes reduction in reaction VI some idea of the effect of attachment of a complex of this kind can be obtained. To an extent dependent on their distance apart, the aluminium complex and the groups attached to the nitrogen atoms will obstruct each other; this effect is known as steric hindrance. The most favourable spatial arrangement is that arising when the aluminium complex is situated between the smallest groups (CH_3 and H — see fig. 6a) attached to the adjacent asymmetric atom, therefore being as far away from the nitrogen atom as possible. The course of the reduction process can be predicted for such a fixed molecule. A hydride ion approaching from under the plane of the drawing meets a *methyl group* belonging to the neighbouring atom; a *hydrogen atom* is the only obstacle in the path of a hydride ion approaching from above. The "energy hill" that the ion has to pass over is much lower in the latter than in the former case. The configuration arising in this latter case is in fact the "threo" configuration obtained in practice.

Let us now consider what happens when there is one group less attached to the nitrogen atom. As before, the most favourable position for the aluminium complex is between the smallest groups of the neighbouring asymmetric atom. However, the molecular model now shows that the "effective bulk" of the nitrogen atom and the single group attached to it is roughly the same as that of the methyl group. There is therefore no way of predicting the preferred molecular structure from the size of the groups. Yet there is a second factor which may be decisive: the place of the missing group has been taken by a hydrogen atom which has the tendency to form a *hydrogen bridge* with an electron donor, available in this case in the form of the oxygen belonging to the carbonyl group. If this hydrogen bond is formed, the position will be as shown in fig. 6b. It will be seen that from the viewpoint of the approaching hydride ion, precisely the opposite situation now prevails. On approaching the carbon atom from below, the hydride ion finds a hydrogen atom in its path; coming from above, it is obstructed by a methyl group. This time the approach from below will be the least difficult, and we may therefore expect the resulting molecule to have the "erythro" configuration.

It may be added that reduction with palladium and carbon rather than with LiAlH_4 is also to be regarded as a factor favouring the "erythro" configuration. Reduction takes place here in an acid environment, in which the nitrogen atom is able to

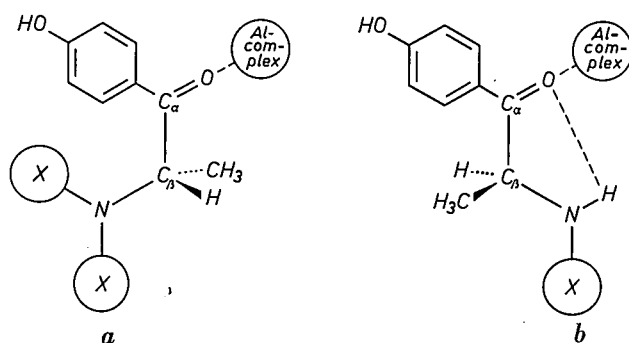


Fig. 6. Molecular models to explain the course of the reaction when LiAlH_4 is used as reducing agent. The method of representing the third dimension is the same as in fig. 3. Diagram a) shows a molecule in which two large groups X have been substituted on to the nitrogen atom. In the course of the reaction the nitrogen atom takes up a position as remote as possible from the Al complex attached to the carbonyl group — a result of "steric hindrance". In diagram b) only one X cluster has been substituted on to the nitrogen and hence the steric hindrance effect is weaker. It is also possible for a hydrogen bond (the dashed line) to form. The reduction process, which consists of the attachment of a hydride ion to carbon atom a, depends on the state of the molecule, i.e. that shown in a) or that shown in b). The substituents attached to the neighbouring carbon atom β hinder the approach of the hydride to some extent, the methyl group offering more of an obstacle than the hydrogen atom. In structural state a), the hydride is more likely to attach itself to the side of the molecule facing the reader, so giving rise to the "threo" configuration; in state b), it is more likely to attach itself to the other side, giving rise to the "erythro" configuration.

collect an extra proton. Having acquired a positive charge, the nitrogen atom is certain to take up a position close to the negatively polarized oxygen atom in the carbonyl group.

This gives a satisfactory explanation of the stereospecific nature of the reduction reaction and the reversal referred to above. The only thing that still has to be explained is why reaction IV, which yields *both* configurations, departs from the above described pattern. This can easily be explained. Apart from accelerating the reduction of the carbonyl group, the palladium-carbon catalyst helps to detach the benzyl group. For the molecules where the benzyl is detached before the carbonyl is reduced, the course of the reaction will then be the same as in II and the end-product will have the "erythro" configuration. For the molecules in which the sequence of these events is reversed, "threo" configurations may form.

Now that the synthesis of isoxsuprine has been described and an explanation has been offered for the stereospecific character of the reduction reaction, some experimental evidence for the pharmacological action of isoxsuprine will be given. From now on the drug will be referred to under its trade-name of "Duvadilan" ⁷⁾. The evidence in question includes findings both from animal experiments and from clinical trials.

⁷⁾ For pharmaceutical purposes "Duvadilan" is prepared as the hydrochloride of isoxsuprine.

The pharmacology of "Duvadilan"

Animal experiments

Many methods of investigation have been used to determine the effect of "Duvadilan" on *blood-vessels*. Firstly, its effect on blood flow through an ear of a rabbit has been investigated: the vessels were constricted by administering adrenaline; then "Duvadilan" was given and was found to reverse the constriction effect, the blood flow being restored⁸⁾. Secondly, the blood flow was measured in the artery in a rear leg of a dog; *fig. 7* shows the percentage increases obtained in the rate of flow as a function of dose⁹⁾. Thirdly, the effect of "Duvadilan" was

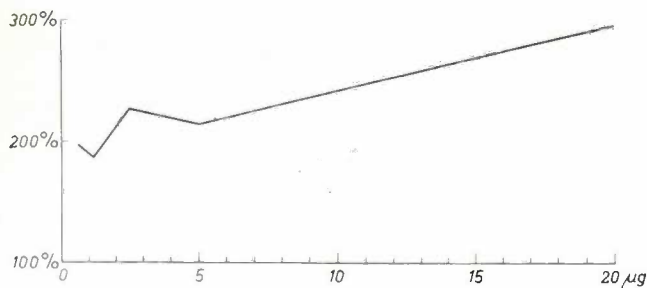


Fig. 7. Arterial blood flow in a rear leg of a dog, as a function of the dose of "Duvadilan" injected into the artery. Blood flow values are expressed as a percentage of the normal value. (After Clark⁹⁾.)

⁸⁾ F. Brücke, Wiener Klin. Wochenschrift 68, 183, 1956.

⁹⁾ The graph has been plotted from data provided by B. B. Clark, and published in F. Kaindl, S. S. Samuels, D. Selman and H. Shaftel, *Angiology* 10, 185, 1959.

compared with that of a physiological salt solution in experiments on the artery in a leg of a rabbit; the observed changes in blood flow have been plotted in *fig. 8*¹⁰⁾.

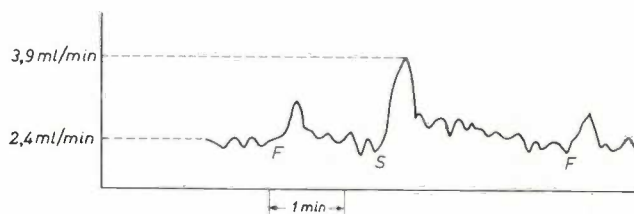


Fig. 8. Record of arterial blood flow in a leg of a rabbit. The effect of "Duvadilan", which was injected in the artery at instant S, was compared with that of a physiological salt solution, administered at instants F. (After Hyman and Winsor¹⁰⁾.)

In addition, experiments have been performed with "Duvadilan" on less accessible regions such as the heart and the brain. The effect on the heart muscle was investigated as follows. Characteristic changes appear in the electrocardiogram of a dog after treatment with a hormone known as "Pitressin" (Parke Davis trade-mark). The changes are caused by a shortage of oxygen: the coronary arteries, which supply the heart, go into a spasm as a result of the administration of "Pitressin". *Fig. 9a* is an electrocardiogram showing this effect. *Fig. 9b*

¹⁰⁾ C. Hyman and T. Winsor, *Acta pharmacol. et toxicol.* 17, 59, 1960.

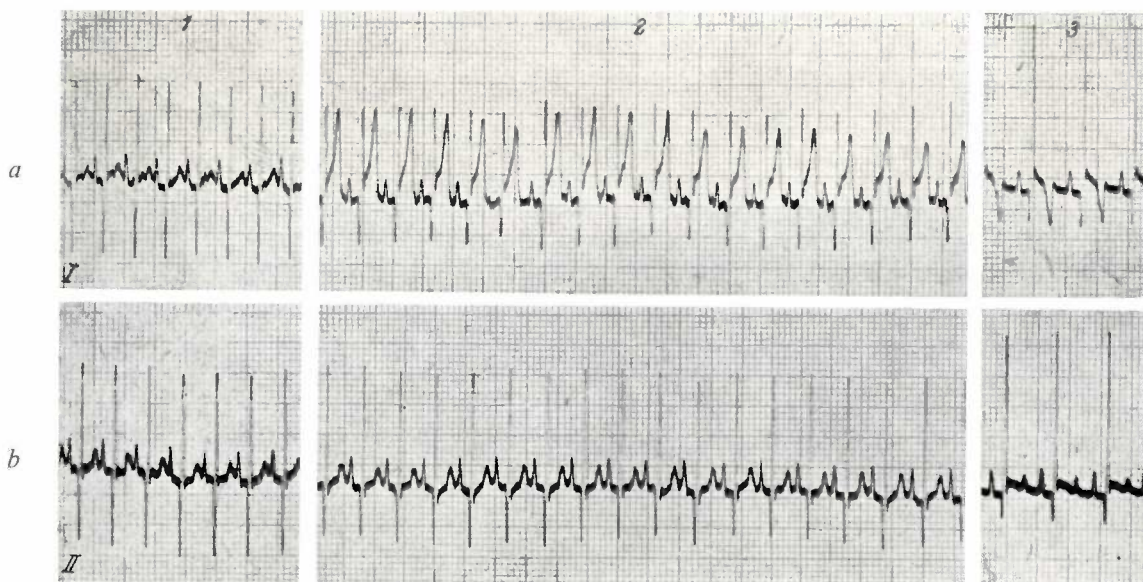


Fig. 9. Two electrocardiograms, recorded at an interval of eight days, showing the response of a dog to treatment with a) "Pitressin", which induced coronary spasm, and b) "Pitressin" and "Duvadilan". The traces to the left of the first vertical division were obtained before administering "Pitressin"; the traces between the vertical divisions were recorded immediately after, and those on the right five minutes after "Pitressin" was administered. On occasion b) the dog received a subcutaneous injection of 0.25 mg of "Duvadilan" per kg of body weight 15 minutes before the "Pitressin" was administered. It is clear that "Duvadilan" suppresses the effect of "Pitressin". (After Brücke⁸⁾.)

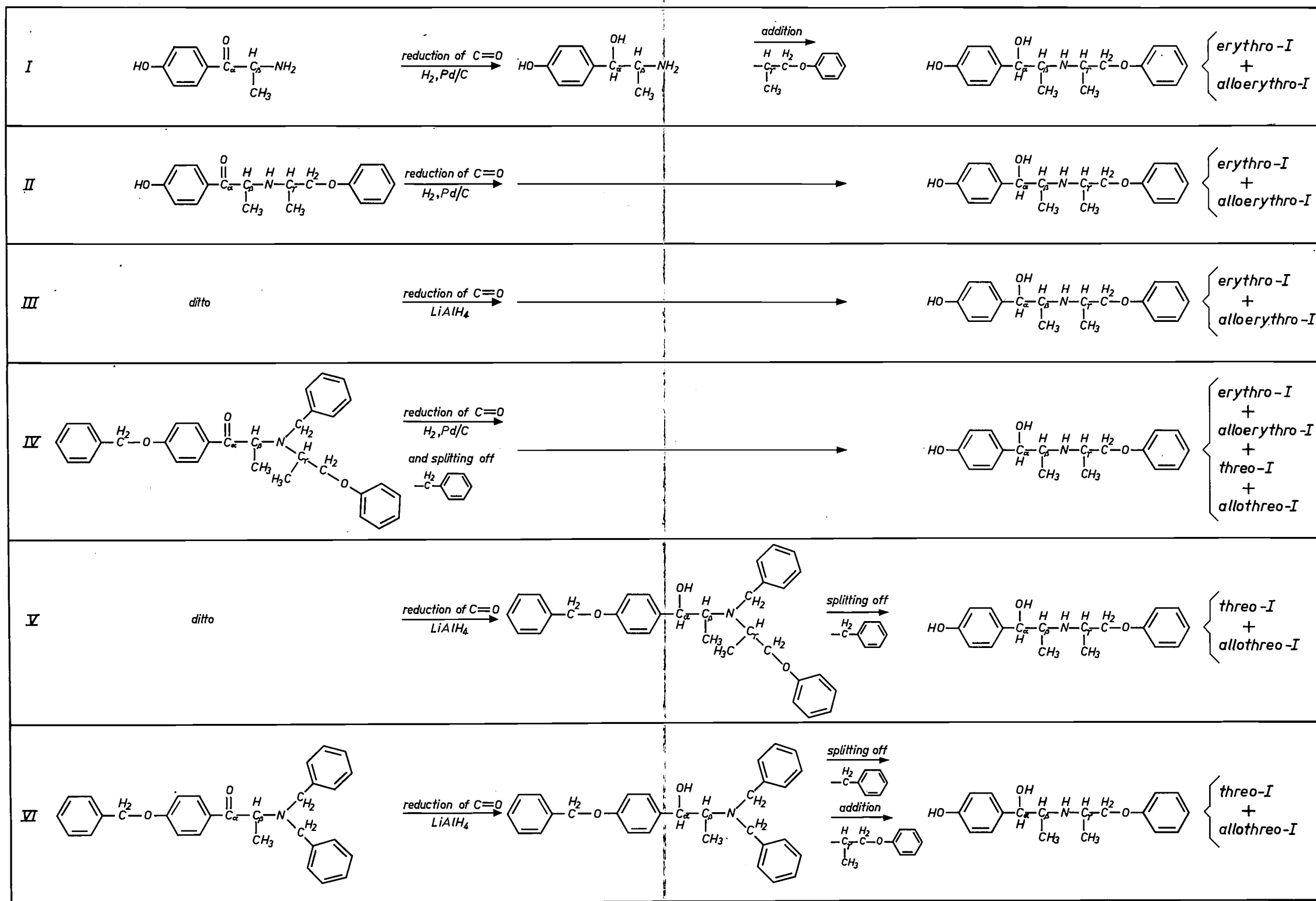


Fig. 5. Reactions yielding isoxsuprine (erythro-I) and the isomeric racemates alloerythro-I, threo-I and allothreo-I. All the starting substances can be produced by conventional methods.

was recorded a few days later; on this occasion the same dog was given a subcutaneous injection of "Duvadilan" before being treated with "Pitressin". In this case the latter drug appears to have had no effect⁸⁾.

To illustrate the effect on the blood supply to the brain, experiments have been performed in which the cerebral blood-vessels of dogs and cats were photographed via a window in the top of the skull¹¹⁾¹²⁾. Afterwards the diameters of the vessels appearing in the photographs (fig. 10) were measured. Table II gives the results of an experiment of this kind in which the animal was injected with 0.1 mg of "Duvadilan" per kg of body weight. It is evident from this that "Duvadilan" dilates the blood-vessels of the brain.

The effect of "Duvadilan" upon the heart was measured from the heart rate (number of beats per minute) as determined from an electrocardiogram. Following intravenous administration of 0.2 mg of "Duvadilan" per kg of body weight, the heart rate of non-anaesthetized dogs was found to increase rapidly at the outset, reaching twice the starting value and thereafter slowing down again; 70 minutes after medication it had fallen to 24% above the starting value⁸⁾.

It should be borne in mind that the speed of the heart can alter for two reasons: it may change owing to direct action on the heart, and also in response to a change in blood pressure. This latter possibility must be excluded if it is desired to record changes in heart rate that are a true measure of direct action on the heart. A heart-lung preparation is used for this purpose, reflex acceleration being eliminated by cutting the nerve connections that make it possible. In a preparation of this kind, the heart rate was

Table II. Diameters (in microns) of arteries in the cerebrum of a dog, as determined from photographs taken before and after "Duvadilan" was injected into a vein in a rear leg. (After Jourdan and Faucon¹²⁾.)

Time	Artery A	Artery B	Artery C	Artery D	Artery E
16.46 h	429	143	71	209	126
16.49 h	429	143	77	209	110
16.50 h	injection of 0.1 mg "Duvadilan" per kg body weight				
16.51 h	434	159	82	236	132
16.53 h	440	165	93	236	121
16.56 h	445	159	93	242	110
17.00 h	445	165	88	253	121
17.02 h	429	165	99	247	121
17.18 h	434	176	99	335	159
17.19 h	434	165	88	330	165

¹¹⁾ H. Wahrenbourg, P. Pruvot, A. Sueur and J. Lekieffre, *Thérapie* 16, 314, 1961.

¹²⁾ F. Jourdan and G. Faucon, *Thérapie* 14, 1075, 1959.



a



b

Fig. 10. Photographs taken through a window in the skull of a cat, showing cerebral blood-vessels. Exposure a) was made at the same instant as 0.1 mg of "Duvadilan" per kg of body weight was injected into a blood-vessel in a rear leg; exposure b) was made 90 seconds afterwards. Of the two blood-vessels in the middle of the photographs, the artery on the left was greatly dilated while the vein on the right underwent no change. (After Wahrenbourg, Pruvot, Sueur and Lekieffre¹¹⁾.)

found to increase by between 20% and 50% following administration of 5 mg of "Duvadilan". The corresponding increase in the rate at which blood was pumped around the body (the "heart minute volume") was about 20%⁸⁾. On the basis of these findings "Duvadilan" may be said to have a comparatively weak direct action on the heart.

Clinical evidence

Numerous medical practitioners have reported the results of "Duvadilan" treatment of their patients.

A check on whether blood flow has really improved or not is provided by a technique known as *plethysmography* which allows the blood stream in a given part of the body to be measured from the outside.

There is a pressure difference of about 100 mm Hg between the blood entering and the blood leaving a part of the body such as a finger, foot, arm or leg. This difference is that between the arterial and venous pressures. By applying an inflatable cuff to these parts of the body it is possible to stop the flow of blood through the vein while not interfering, or

hardly interfering, with that through the artery. The cuff is inflated to a pressure higher than the venous but lower than the arterial pressure. The result is that the volume of the affected part increases. The change in volume per unit time *immediately after blockage* is equal to the rate of arterial blood flow.

One way of measuring the change in volume is by means of an inelastic container into which the limb or part is inserted and then sealed off from the exterior, the space thus enclosed being connected to a volumometer. Blood flow in the limb or part can be determined from a curve showing changes in the enclosed volume, as traced by the volumometer.

The tracings shown in *fig. 11* were obtained in this way. They record blood flow in the calf of a patient

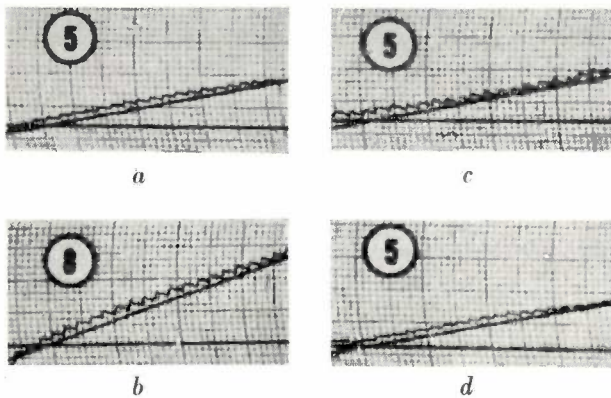


Fig. 11. Plethysmograms *a*) and *b*) showing the volume increases undergone by a patient's right calf before and after injection of 10 mg of "Duvadilan" into a leg artery. As a control, increases in the volume of the left calf were recorded at the same time (plethysmograms *c* and *d*). The slope of the trace is a measure of the blood flow. In the right calf this increased by 80% in consequence of the "Duvadilan" injection. The effect of the drug when administered in this way does not extend to the other limb.

The heartbeat shows up in the fluctuations of the plethysmographic trace, which can accordingly serve also as a record of heart rate. (After Hyman and Winsor¹⁰.)

before and after "Duvadilan" was injected into the artery supplying the calf. It will be seen that the slope of the curve increased by 80% after the injection. As a control the same kind of measurement was performed simultaneously on the untreated calf; there is no evidence of any increase in blood flow through this calf, indicating that the drug possesses a local action when administered in this manner¹⁰).

The method just described reveals *overall* changes in blood flow through a given part of the body. It is also desirable to know to what extent these changes benefit the various tissues through which the blood circulates (muscle, skin), and assist the various functions the blood performs in these tissues (heat regulation, metabolism). A method specially designed

for ascertaining effects on metabolism in a given kind of tissue consists of measuring the rate at which the blood carries away a known quantity of radioactive sodium (²⁴Na) that has been injected into the tissue. It was found in this way that after the injection of 10 mg of "Duvadilan" into a muscle, the "²⁴Na clearance rate" increased by about 50%¹⁰).

Other data relating to the effects of "Duvadilan" have been obtained by using the drug to treat vascular disorders involving "intermittent claudication" (i.e. occasional limping). Sufferers from such disorders are prevented from walking long distances by the onset of spasm in the blood-vessels. The average distance the patient is able to walk without experiencing pain is an inverse measure of the seriousness of the condition. "Duvadilan" treatment has been found to increase the average walking distance.

"Duvadilan" has also been used to treat necrotic skin diseases in which, owing to an inadequate blood supply, patches of skin die. Administration of "Duvadilan" has been found to promote the healing of these necrotic skin areas.

An impression of a drug's *action on the heart* can be obtained not only by taking the pulse or inspecting an electrocardiogram but also by measuring variations in blood pressure. It should be explained that the pressure of blood in the circulatory system alternates in the rhythm of the heartbeat between two values, the *systolic* and the *diastolic* pressure. The amplitude of the pressure change (i.e. the systolic minus diastolic pressure, otherwise known as "pulse pressure") is related to the amount of blood forced into circulation with each contraction of the heart. This amount is called the *stroke volume*. *Fig. 12* shows blood pressure values before and after a

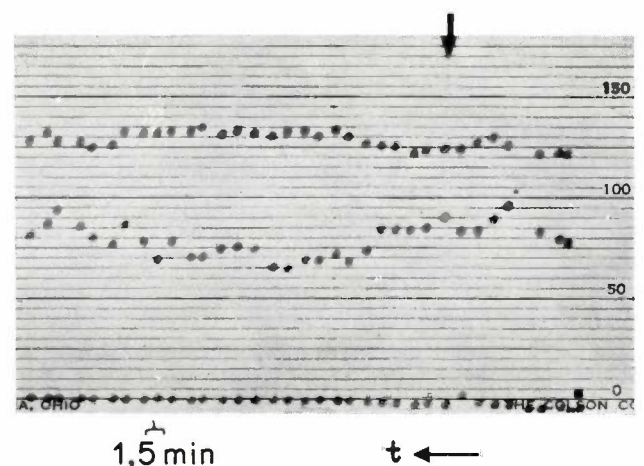


Fig. 12. Record of a patient's blood pressures. The upper chain of dots represents systolic, the lower chain diastolic pressures. At the instant marked by an arrow, a 10 mg dose of "Duvadilan" was given in the form of an intramuscular injection. Time is represented from right to left. (After Hyman and Winsor¹⁰.)

"Duvadilan" injection. It reveals a slight increase in pulse pressure after the injection, which can be taken as an indication that the stroke volume had increased to some extent¹⁰). Further evidence in this direction is provided by radiographs showing the outline of the heart wall. The displacement of the heart wall as recorded by this technique, known as X-ray kymography, can be taken as a basis for calculating the stroke volume. On this basis the stroke volume has been found to increase by about 7% after a slow intravenous injection of 10 mg of the vasodilator¹⁰).

The effect of "Duvadilan" on the *heart minute volume* has also been investigated. As already indicated, the minute volume is the output of blood from the heart over a period of one minute; it is therefore the product of stroke volume and heart rate. The minute volume can be measured by injecting a small quantity of radioactive serum into the bloodstream at some point, and then recording, by means of a detector placed above the heart, for example, the rate at which radioactive material is passing this second point. From results of such measurements it has been calculated that the minute volume increases by about 30% after a slow intravenous injection of 10 mg of "Duvadilan". Evaluation of electrocardiographic and other direct evidence of heart rate changes shows an increase of about 20% in the number of beats per minute following an intramuscular injection of 10 mg of "Duvadilan". This figure agrees more or less with those for stroke and minute volume¹⁰). All in all, these data reveal that "Duvadilan" has a relatively small effect on the heart, in human beings as well as in laboratory animals.

In conclusion it may be mentioned that the drug is being employed on an increasing scale in obstetrics. "Duvadilan" appears to be very effective for relieving uterine spasm, in virtue of its β -mimetic action and probably of its papaverine-like spasmolytic action as well. Indeed, many of the reports on "Duvadilan" have come from obstetric departments, particularly during the last few years. Thus, having first gained recognition as a vasodilator, "Duvadilan" is proving to be of value as a spasmolytic agent.

As to further developments in this field, it will be recalled that "Duvadilan" is composed of two optical antipodes; investigations now in progress are concerned with the individual contributions of the component enantiomers to the effects of "Duvadilan". These may be expected to provide new knowledge that will certainly be of scientific interest, and may possibly be of medical value.

Summary. "Duvadilan" (isoxsuprine), which was introduced a few years ago, is a drug with a specific vasodilating action and is of value for treating certain vascular and other disorders. It is suitable for oral administration as well as by other routes. Its pharmacological action can be regarded as comprising a β -sympathomimetic, an α -sympatholytic and a spasmolytic component. These three properties of isoxsuprine can be explained with reference to the properties of substances to some extent analogous with it, namely adrenaline, noradrenaline, ephedrine, papaverine, etc. Isoxsuprine is a racemic mixture (racemate) which, during synthesis, has to be separated from three other racemates, all the substances concerned being isomers of each other. One stage of separation is effected by exploiting stereospecific reduction (using LiAlH_4 , or alternatively H_2 , Pd and C). A mechanism is suggested to explain the relevant reaction. The concluding section contains data on the pharmacology of "Duvadilan", derived from animal experiments and clinical practice.

ASSEMBLY OF 5-kW COMMUNICATION TRANSMITTERS

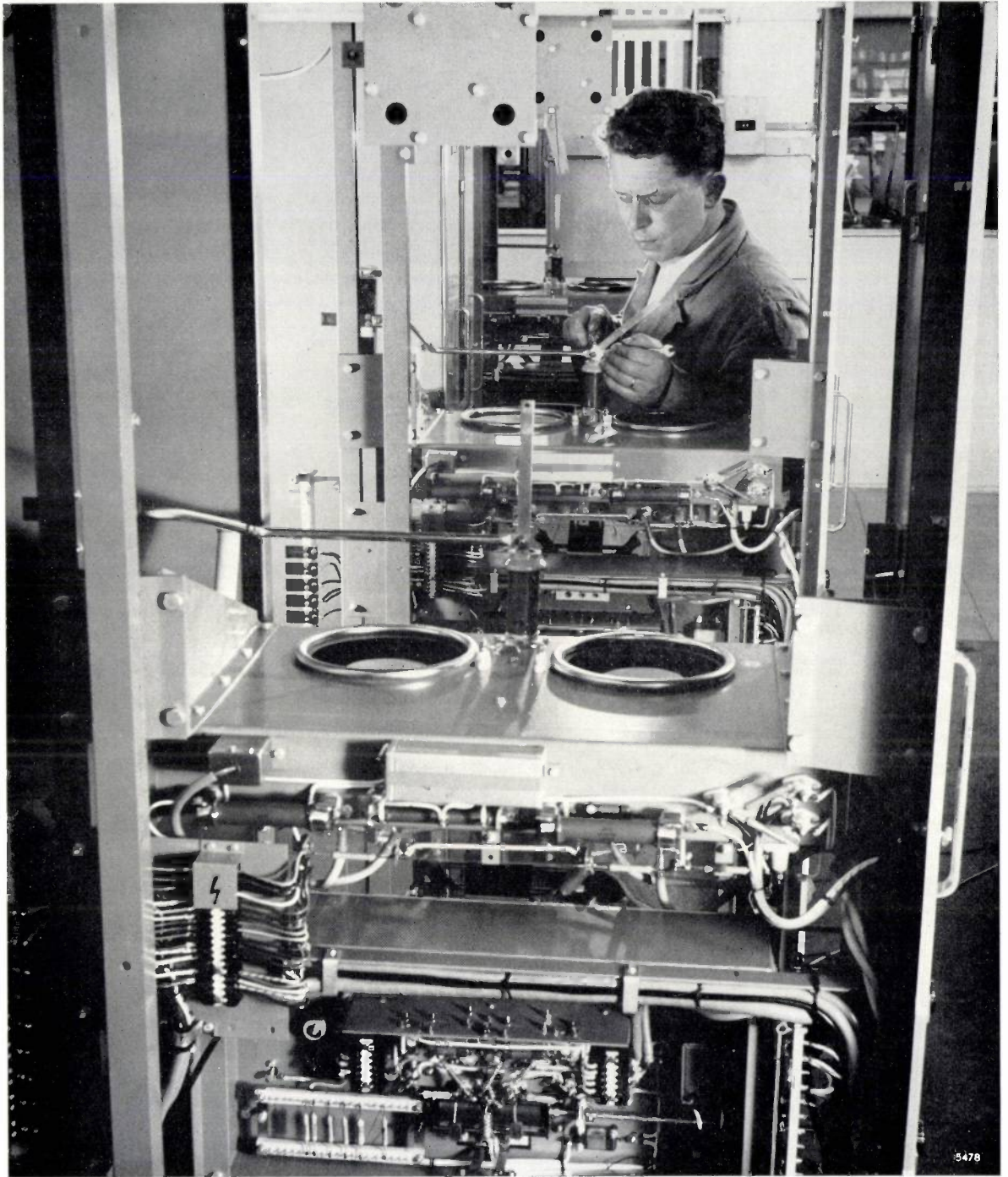


Photo Maurice Broomfield

The photograph shows 5-kW communication transmitters being assembled in the Huizen (Netherlands) factory of N.V. Philips' Telecommunicatie-Industrie. Transmitters of this type consist of a power-supply cabinet, one or more radio-frequency cabinets, and possibly a modulator cabinet. Each RF cabinet is tuned to a fixed frequency. When propagation conditions make it necessary to change the frequency, two manual operations are all that is needed to switch from the cabinet in use to another.

The mechanic is seen here behind the second of three frames slid out of their cabinets. He is fitting the HT cable that connects the decoupling capacitor to the output stage. The latter, which is not yet in place, is mounted on a separate chassis to make it interchangeable.

Above the drive unit in the front cabinet can be seen the airducts for the cooling of the two output valves, which are in push-pull.

AN EXPERIMENTAL APPARATUS FOR RECORDING TELEVISION SIGNALS ON MAGNETIC TAPE

621.397.3:621.395.625.3

Video signals (frequencies up to 5 Mc/s) cannot be recorded on magnetic tape by the normal method employed in sound recording because the tape would have to run at far too great a speed (more than 15 metres per second). For television recording, therefore, systems have been developed that use *moving* recording heads, which write narrow tracks more or less obliquely across the magnetic tape. One possible method, as in the Ampex system, for example, works with several magnetic heads operating in succession, and another system uses only a single head. For some years a system of the latter kind has been in operation in the Philips Research Laboratories at Eindhoven. An apparatus based on this system now gives an image quality that meets the high requirements of studio work. Various

devices have still to be developed and fitted, however, before the apparatus can be used for practical television transmissions.

The principle is illustrated in *fig. 1*. A magnetic tape 2.5 cm wide is moved at a speed of 38 cm/s over a helical track which covers almost the entire periphery (about 353°) of a stationary drum. The pitch of the helix is slightly less than the tape width. Around the periphery of the drum, of diameter 305 mm, is a slit in which runs a recording head (video head, *E*) fixed to the edge of a disc that rotates around the drum axis at high speed (50 r.p.s.); see *fig. 2*. In one revolution, then, the rotating head records a track about 1 m long set obliquely across the tape, and owing to the slight displacement of the tape during this revolution the next track lies

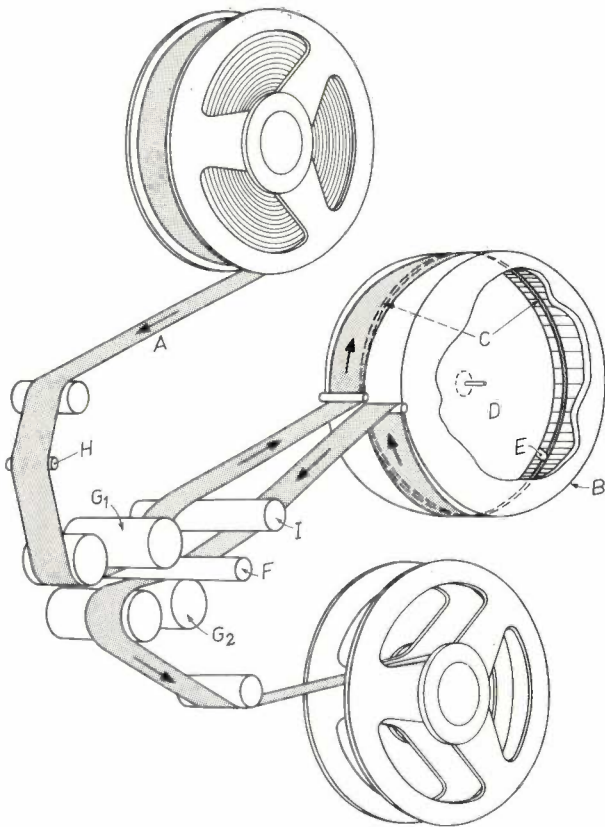


Fig. 1. Principle of the system. *A* magnetic tape. *B* stationary drum with slit *C*. *D* rotating disc with video head *E*. *F* drive shaft, turning at constant speed, with two pressure rollers *G*₁ and *G*₂. *H* erasing head. *I* guide roller.

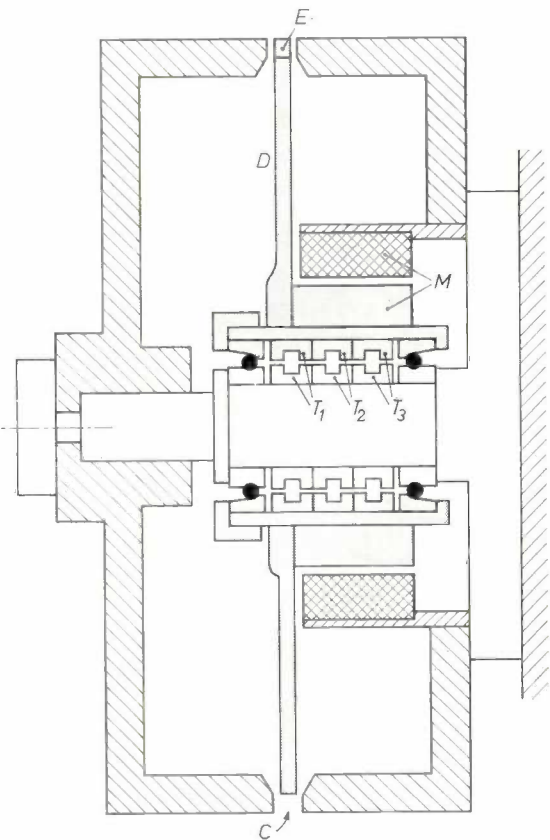


Fig. 2. Cross-section of the drum. *C* slit. *D* disc with video head *E*. The disc is driven by a motor *M* contained inside the drum. *T*₁, *T*₂, *T*₃ transformers. The rotating assembly is indicated by grey shading.

beside the first, and so on. In this way, except for two narrow edges the entire tape is filled with narrow recorded tracks; see *fig. 3*. The video head — a ferroxcube head with a gap of $2\ \mu\text{m}$ — is only $180\ \mu\text{m}$ wide, and this is therefore the track width.

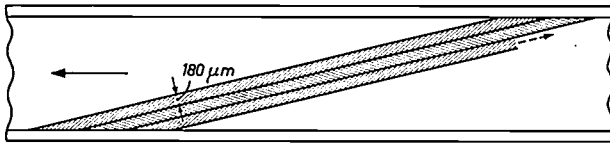


Fig. 3. Sketch representing the tracks recorded on the relatively slowly moving magnetic tape by the rotating head during successive revolutions. In reality the slope of the tracks is much less steep in relation to the length of the tape; each track is about 1 m long.

Since the disc rotates at a speed of 50 r.p.s., each track contains exactly one complete scanning raster, i.e. one picture frame. During recording, the rotation of the disc is synchronized with the picture frames, the phase being set so that the flyback coincides exactly with the moment at which the recording head moves from the end of one track to the beginning of the next.

The same head is used for the *playback* of the picture, the tape and the head moving in the same way as before. To ensure that the head then runs exactly on the centre line of the recorded track (and not, for example, on the boundary between two tracks), the tape must be made to run during playback with the correct "phase" in relation to the revolving video head. For this purpose an auxiliary head, which is stationary, impresses on one of the narrow edges of the tape the normal frame synchronizing signals during the recording of a picture. During playback these signals control the transport of the tape; in other words they act as a kind of magnetic film perforation. The other narrow edge carries the sound. The synchronizing head and the audio head are accommodated in the guide *I* (see *fig. 1*).

For the purpose of recording, the video signal is converted into a signal with frequency modulation, in such a way that the frequencies 5 and 7 Mc/s correspond respectively to the black and white levels. The signals are not applied to the rotating video head via slip rings, but much more reliably via a transformer (with ferroxcube cores), the secondary of which rotates together with the disc. This transformer is contained inside the stationary drum (T_1 in *fig. 2*).

Also contained in the drum are two similar transformers (T_2 and T_3), one of which is intended for feeding a transistor pre-amplifier which rotates with the disc and amplifies the

signal from the video head for playback. The other transformer is intended to be used, in the future, for taking the signal from a monitor head mounted diametrically opposite to the actual video head on the rotating disc. The idea is to use this head, delayed by half a revolution of the disc ($1/100$ s), for monitoring the picture while it is being recorded.

The fact that a complete picture frame is recorded on one track made it possible to render crosstalk between neighbouring tracks virtually harmless. With the modulation system adopted a signal suffers very little interference if the spurious signal (from another track) picked up simultaneously by the head has nearly the same frequency, i.e. if the signal causing the interference corresponds to roughly the same brightness. In the apparatus described this can in fact be ensured by arranging that the more than 300 line-synchronizing pulses in each track are situated exactly next to the pulses in the neighbouring tracks (*fig. 4b*). Plainly, adjacent points of neighbouring tracks will always represent practically the same part of a picture and hence the same brightness level. The picture played back is then virtually free of crosstalk interference,

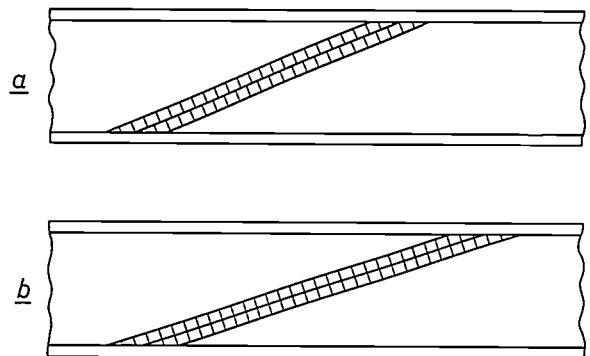


Fig. 4. Illustrating the relative positions of the series of line-synchronizing pulses in two neighbouring tracks: *a*) not aligned, *b*) aligned to minimize crosstalk.

as can be seen in *fig. 5b*. If, however, the line pulses of one track are slightly displaced with respect to those of the neighbouring tracks (*fig. 4a*), particularly troublesome interference occurs when there is crosstalk from a neighbouring line pulse during the scanning of a white part of the picture. This effect can be seen in *fig. 5a*.

To ensure that the line pulses have the required alignment (following a pattern as in *fig. 4b*), it is necessary for the drum diameter to be equal to one of a few discrete values, for a given tape speed. The diameter of 305 mm which we chose is one of the series of the values corresponding to a tape speed of 38 cm/s.

Since one track (length about 1 m) contains exactly one frame, the apparatus described also



Fig. 5. Television picture recorded and played back with the apparatus described. In *a*) the line-synchronizing pulses were out of alignment (fig. 4*a*); in *b*) the situation was as in fig. 4*b*. In both cases deliberate steps were taken to produce strong crosstalk, i.e. a ratio of 3 : 1 between the desired and the interference signal in the playback head.

has the remarkable facility of playing back the instantaneous picture *when the tape is stationary*. In that case the beginning and end points of a scanning track of the head are displaced by one track width (because there is now no contribution from the tape movement to the relative movement of tape and head), so that for a while the head scans two neighbouring tracks simultaneously. Surprisingly little is noticed of this, however, owing to the above-mentioned behaviour in respect to crosstalk. It proved to be possible to scan such a stationary pic-

ture for quite a considerable time, e.g. for several minutes, without the magnetic layer suffering any appreciable wear from the video head, which in this case passes over the same piece of tape 3000 times a minute.

In due course we intend to publish in this journal a more detailed description of the mechanical and electronic parts of the apparatus¹⁾. It may be mentioned here, however, that a problem to be solved — in fact the main problem of all video recording methods — was the accurate fixing of the recorded video signals in time and space. For black-and-white television the maximum permissible fluctuation amounts to one image point ($\approx 1/600$ picture width, i.e. 10^{-7} s, which is about $5 \mu\text{m}$ on our tape). A servo system in our equipment ensures that the synchronization between the tape drive and the disc drive remains within this close tolerance. Variations in the tension of the tape, which would cause noticeable elongations and hence errors in signal location, are prevented by transporting the tape, both to and from the drum, by pressure rollers (G_1 and G_2 in fig. 1, which keep the tape pressed against the driving shaft F , rotating with a constant speed). Thermal expansion of the tape is compensated by reducing the tensioning; this adjustment is at present done by hand.

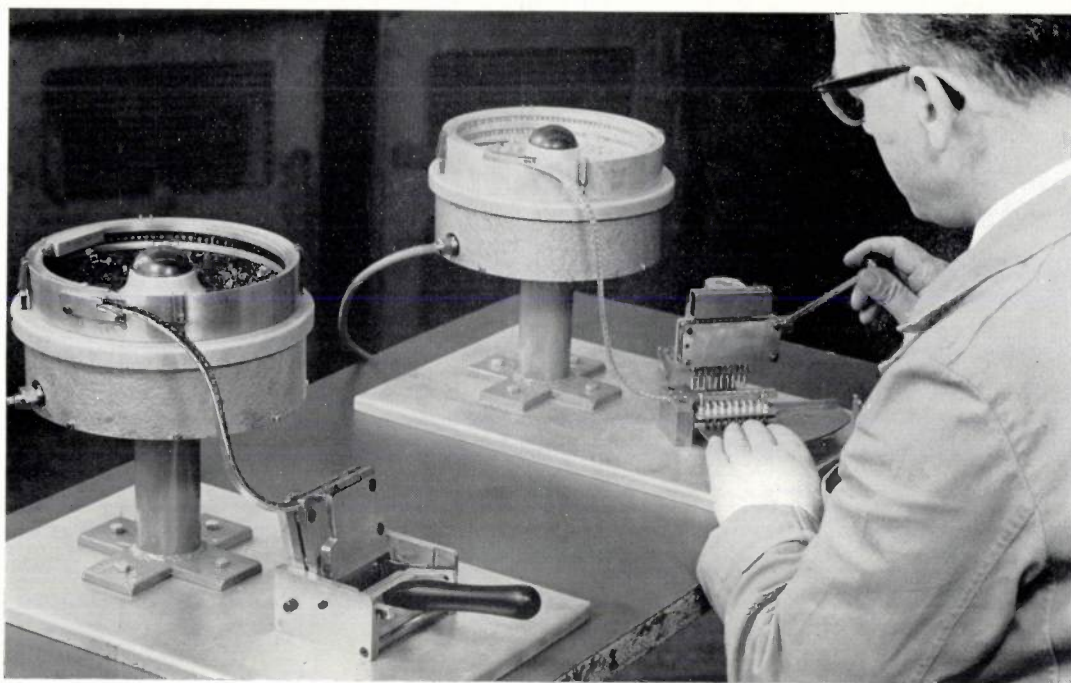
For recording colour-television pictures the permissible fluctuations in the location of the video signals are at least 20 times smaller than for black-and-white television. For this purpose

the synchronization system referred to is inadequate. With a view to future experiments on the recording of colour-television signals a special synchronization system is being developed, one feature of which is an electronically controlled delay circuit.

F. T. BACKERS *) and J. H. WESSELS *).

¹⁾ The present experimental form is described at somewhat greater length in an article shortly to be published in the *Nachrichtentechnische Zeitschrift*.

*) Philips Research Laboratories, Eindhoven.



VIBRATORY FEEDERS

by H. G. de COCK *).

621.867.52

Vibratory feeders are conveyor devices that automatically feed a regular supply of small components to a production process. Literature on the operation of vibratory feeders is scarce. The article below summarizes a study of the subject made by the author for the purpose of production mechanization in the Philips factories.

Introduction

The vibratory feeder occupies an important place among the mechanisms used for supplying small components in an ordered manner to production processes. This is largely due to the fact that in vibratory feeders each component to be fed is separately set in motion, so that the components do not have to push each other forward. This means that nearly all kinds of components, whatever their shape and size, can be fed in any orientation by a vibratory feeder. Examples are carbon resistors, small capacitors, ceramic sleeves, filaments, rivets, nuts and screws, washers etc.

A vibratory feeder of the type described in this article consists of a round drum 1 (fig. 1), which holds the components to be fed. The drum is mounted on inclined leaf springs 2, which are fixed to a counterweight 3 resting on rubber pads 4. The

drum is vibrated by an electromagnetic or pneumatic drive system 5. The vibration consists of a rotating component in the horizontal plane, and a vertical component. The inside wall of the drum is fitted with a helical track 6, up which the components 7 move from the batch at the bottom. In this process the components can be easily sorted (i.e. broken components or different types that have inadvertently slipped in between are removed) and oriented (set in the right position).

The title photo shows a set-up in which two modern, pneumatically driven vibratory feeders are supplying two different kinds of components of a coil to one point of use, where they are processed nine at a time. In addition to fulfilling the functions of feeding, orienting and sorting, this set-up offers a considerable saving of labour by making multiple handling possible.

The most important feature of a vibratory feeder is the maximum speed at which it can transport the

*) Philips Radio, Television and Record-playing Apparatus Division, Eindhoven.

components. This depends on several variables, and in this article we shall consider which these are and how they should be chosen to obtain the optimum result. Consideration will also be given to the maximum spring tension to be expected when the vibratory feeder is in operation.

After treating the theoretical behaviour of a part on a vibrating inclined plane, various experiments based on this theory will be described ¹⁾. This will be followed by a nomogram that can be used for designing vibratory feeders.

The results obtained apply equally well to straight vibratory feeders ²⁾. Compared with the latter, however, the drum form has the advantage that it can contain a larger quantity of components and that components that drop out during sorting and orienting fall back into the drum.

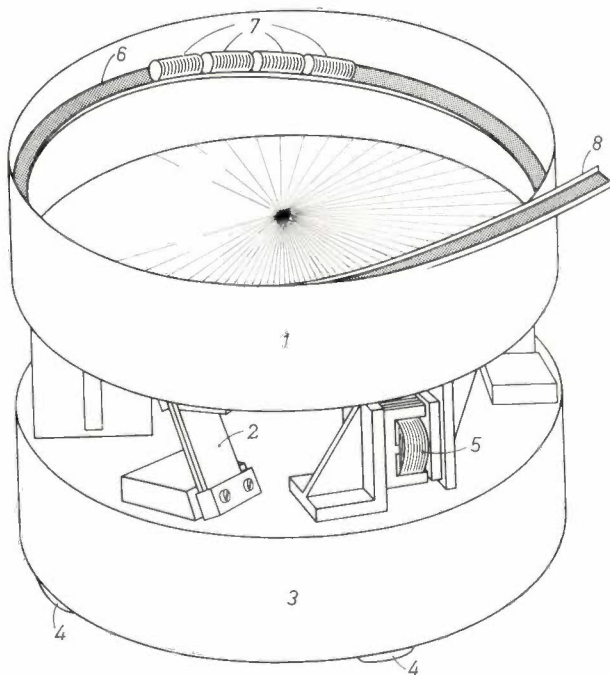


Fig. 1. Sketch of a vibratory feeder. The drum I is mounted on inclined leaf springs 2 attached to a counterweight 3, which rests on rubber pads 4. An AC-actuated electromagnet 5 keeps the drum vibrating in relation to the counterweight. The inside wall of the drum is fitted with a helical track 6 along which the components 7 move upwards to the exit 8.

Theory of the vibratory feeder

The transport of parts along a flat track (in general inclined) is possible if the parts either slide along the track or — by being repeatedly thrown upwards — move along the track in a series of hops. In a vibratory feeder the parts are given either one of these motions or both, in that the track is kept in harmonic vibration the direction of which makes a certain angle β with the track (fig. 2). This is the purpose of the inclined leaf springs on which the drum is mounted (fig. 1).

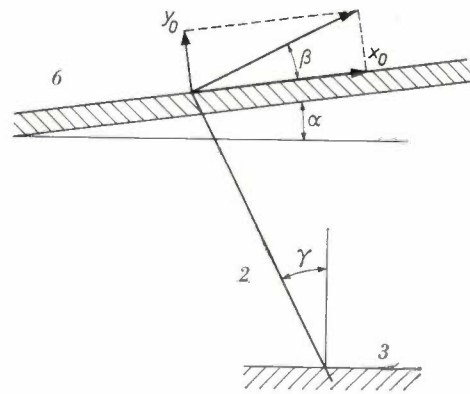


Fig. 2. The vibrating track 6 of a vibratory feeder, with leaf spring 2 and counterweight 3 (cf. fig. 1). The angle of inclination of the track is α . The track vibrates in a direction which is perpendicular to the plane of the spring and makes an angle β (vibration angle) with the track. In the neutral position the leaf spring makes an angle γ (spring angle) with the vertical; when the leaf spring and the track are at equal distances from the centre line of the drum, γ is equal to $\alpha + \beta$.

If the amplitude of vibration of a vibratory feeder is gradually increased from zero, the parts will successively describe the following motions:

- At first the parts follow the movement of the vibrating track, i.e. there is no relative movement and therefore no transport.
- The parts start to slide periodically up the track and are thus transported.
- A stage is then reached in which the parts slide alternately upwards and downwards.
- Finally the parts periodically leave the track for a moment.

The case that arises depends on the acceleration of the track. More important than the amplitude of vibration x_0 (in the direction of the track) is, therefore, the amplitude of acceleration $x_0\omega^2$, where ω is the angular frequency of the vibration.

The values of $x_0\omega^2$ at which the above-mentioned cases (a), (b), (c) and (d) occur can be derived theoretically, as shown in *Appendix I*. They are found to be a function of the following variables (fig. 2): the angle of inclination α of the track, the

¹⁾ As far as we know the only noteworthy literature on vibratory feeders consists of two recent articles in the Russian journal *Stanki i instrument*. Reference may be made to the English publication of the journal: V. A. Povidaylo, Design calculations and construction of vibrating hoppers, *Machines and tooling* 30, No. 2, 5-9, 1959, and V. A. Povidaylo, Optimum vibratory feeder operating conditions, *Machines and tooling* 31, No. 5, 2-6, 1960. These articles only came to our notice after we had completed our own investigations. The author's results agree in broad lines with ours.

²⁾ See also S. Böttcher, Beitrag zur Klärung der Gutbewegung auf Schwingrinnen, *Fördern und Heben* 8, 127-131, 235-240 and 307-315, 1958.

vibration angle β relative to the track, and the static coefficient of friction μ between the parts and the track.

The calculation shows the following.

If
$$\frac{x_0\omega^2}{g} \geq \frac{\mu + \alpha}{1 + \mu \tan\beta}, \dots (1)$$

where g is the acceleration due to gravity, the parts periodically slide upwards.

If
$$\frac{x_0\omega^2}{g} \geq \frac{\mu - \alpha}{1 - \mu \tan\beta}, \dots (2)$$

the parts will periodically slide downwards.

If
$$\frac{x_0\omega^2}{g} \geq \cot\beta, \dots (3)$$

the particles will periodically leave the track.

In expressions (1), (2) and (3), $\tan\beta$ predominates. Fig. 3 shows the various limiting values of $x_0\omega^2/g$ for the various forms of motion as a function of the vibration angle β . It is assumed here that $\alpha = 2^\circ (= 0.035 \text{ radians})$ and $\mu = 0.4$.

As shown in Appendix I, a condition for the proper operation of a vibratory feeder is:

$$\frac{\mu + \alpha}{1 + \mu \tan\beta} < \frac{\mu - \alpha}{1 - \mu \tan\beta} < \cot\beta, \dots (4)$$

indicating that β must be greater than the abscissa of the point of intersection P in fig. 3 and smaller than the abscissa of the point of intersection Q .

If the amplitude of vibration and thus also the amplitude of acceleration of such a vibratory feeder is increased, the parts will therefore:

a) follow the movement of the track (i.e. not be transported) as long as

$$\frac{x_0\omega^2}{g} < \frac{\mu + \alpha}{1 + \mu \tan\beta},$$

b) slide periodically upwards if

$$\frac{\mu + \alpha}{1 + \mu \tan\beta} \leq \frac{x_0\omega^2}{g} < \frac{\mu - \alpha}{1 - \mu \tan\beta},$$

c) slide alternately upwards and downwards if

$$\frac{\mu - \alpha}{1 - \mu \tan\beta} \leq \frac{x_0\omega^2}{g} < \cot\beta,$$

since both (1) and (2) are satisfied (the upward sliding motion predominates), and

d) periodically jump clear of the track when

$$\frac{x_0\omega^2}{g} \geq \cot\beta$$

(in this case the parts may also slide in both directions, as both (1) and (2) are still satisfied).

If, instead of (4), the following condition were to apply:

$$\frac{\mu + \alpha}{1 + \mu \tan\beta} > \frac{\mu - \alpha}{1 - \mu \tan\beta},$$

so that β lay in the region left of P in fig. 3, the downward sliding would start at a smaller amplitude, in other words earlier, than the upward sliding. A vibratory feeder of such a kind would be incapable of upward feeding. The boundary case, given by point P in fig. 3, is determined by

$$\frac{\mu + \alpha}{1 + \mu \tan\beta_{cr}} = \frac{\mu - \alpha}{1 - \mu \tan\beta_{cr}},$$

where β_{cr} is the critical vibration angle. From the latter equation it follows that:

$$\tan\beta_{cr} = \frac{\alpha}{\mu^2}. \dots (5)$$

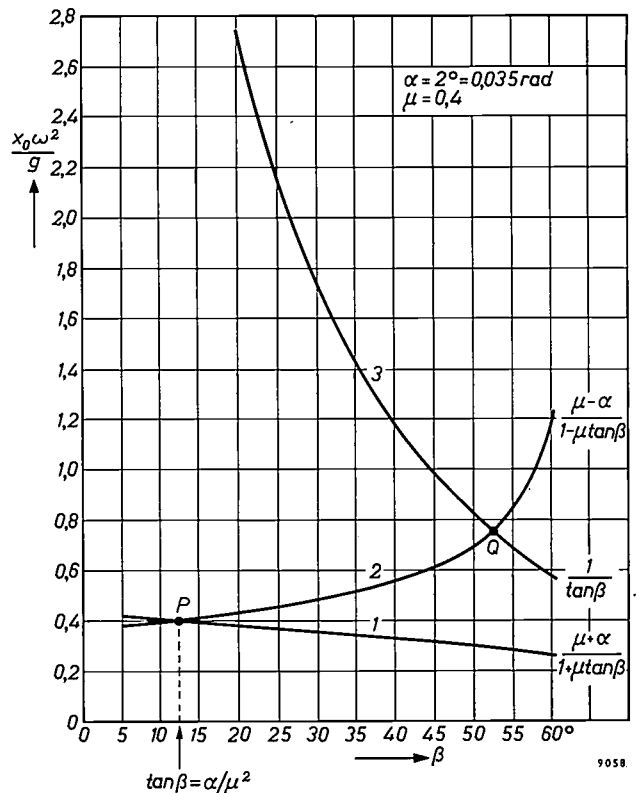


Fig. 3. As a function of the vibration angle β , with $\alpha = 2^\circ$ and $\mu = 0.4$, curve 1 gives a plot of $(\mu + \alpha)/(1 + \mu \tan\beta)$, curve 2 a plot of $(\mu - \alpha)/(1 - \mu \tan\beta)$ and curve 3 a plot of $\cot\beta$. Below curve 1 lies the region in which the parts follow the movement of the track (and in which they are thus not transported), above curve 1 begins the region in which they slide periodically upwards. Curve 2 separates the region of periodic upward sliding from the region of alternate upward and downward sliding, Curve 3 marks the beginning of the region in which the parts periodically leave the track. For good operation of the vibratory feeder, β should lie between the projections onto the abscissa of the points of intersection P and Q .

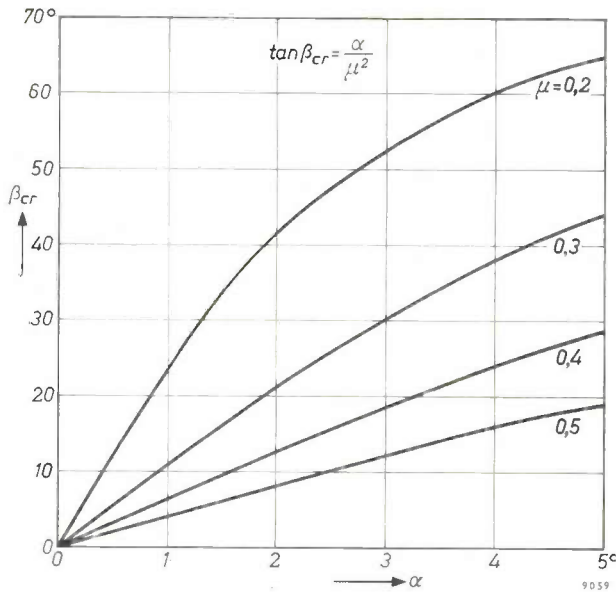


Fig. 4. The critical value β_{cr} of the vibration angle ($\tan \beta_{cr} = \alpha/\mu^2$) as a function of the angle of inclination α , for various values of the coefficient of friction μ .

Thus, for upward transport along a vibrating inclined track the vibration angle β must be greater the steeper and/or smoother is the track.

Fig. 4 gives a plot of the critical vibration angle β_{cr} versus α for various values of μ .

Practice has shown that optimum operation of a vibratory feeder occurs when the vertical component of the acceleration has an amplitude between $2g$ and $2.5g$, so that for most of the time the parts are clear of the track. Although we can then hardly speak of the parts sliding, and we have just shown the significance of the critical vibration angle by considering sliding phenomena, the critical vibration angle is nevertheless of great importance in these conditions too. To be certain of upward transport at all times it is therefore necessary to ensure that $\tan \beta$ is greater than α/μ^2 . This has been confirmed by experiments.

Results of experiments

Verification of theory by experiments on a model

The effects so far discussed are, as we have seen, governed by the acceleration of the track. This makes it possible to choose a frequency low enough, and at the same time the vibration amplitude large enough to allow the effects to be observed with the naked eye. For this purpose a model was made on which the variables α , β , μ and the frequency $f = \omega/2\pi$ could easily be altered. This model, in which the track is straight, is shown in fig. 5.

The results of the measurements on this model are represented by the dashed curves in fig. 6. As can be

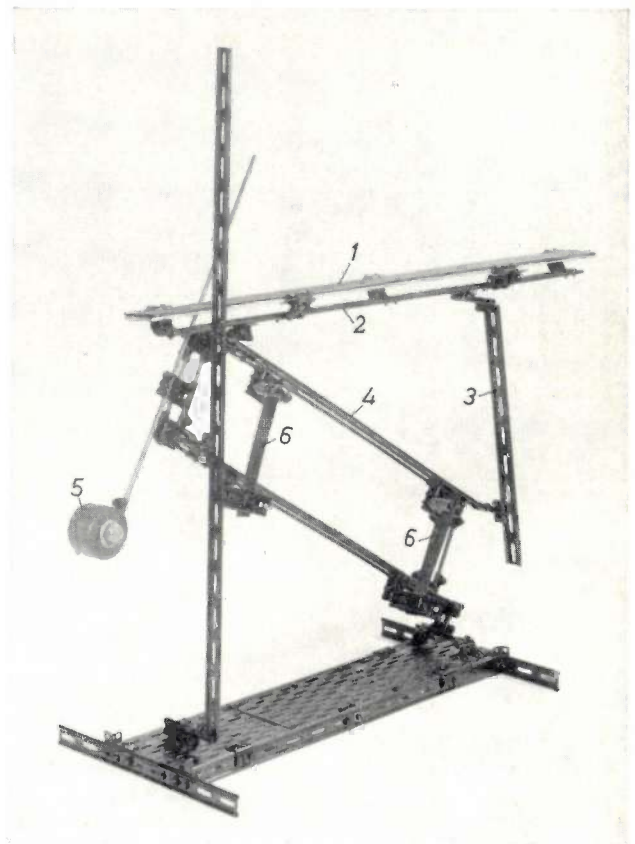


Fig. 5. Model used for experiments. 1 vibrating track (in this case straight) forming a rigid assembly with the beams 2, 3 and 4 and an adjustable weight 5. The track is carried by leaf springs 6, attached to the base. The track is set in vibration by hand and vibrates at its natural frequency, which can be varied by shifting the weight 5. Low frequencies and large amplitudes are used so that the effects can be followed with the naked eye. The track can be fitted with surfaces of differing roughness.

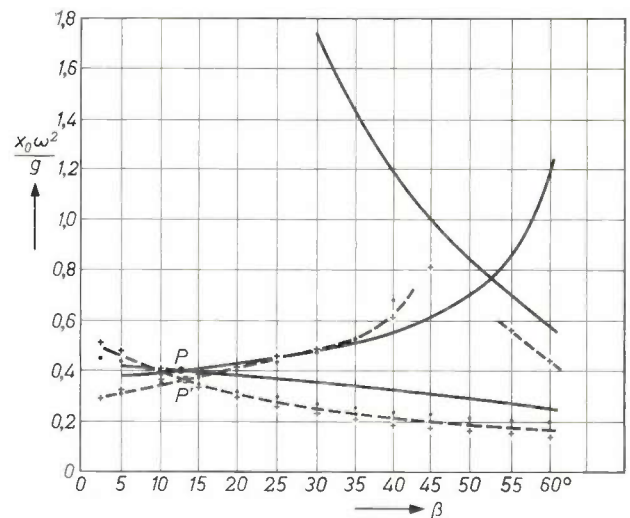


Fig. 6. The three limit curves for the various forms of motion. $x_0\omega^2/g$ is plotted versus the vibration angle β , for $\alpha = 2^\circ$ and $\mu = 0.4$. The solid lines are the theoretical curves (identical with those in fig. 3), the dashed lines the measured curves. The agreement is good, as can be seen from the close proximity of the points of intersection P and P' (with abscissa $\beta_{cr} = \tan^{-1}(\alpha/\mu^2)$).

seen, they are in reasonable agreement with the theoretical (solid) curves, particularly as regards the proximity of P' to P , with respect to the abscissa $\beta_{cr} = \tan^{-1}(a/\mu^2)$.

The experiments showed that the static coefficient of friction μ for a wide variety of components and track surfaces always had a value between 0.35 and 0.45. The spread in the values for one and the same combination of materials was often found to be greater than the difference between different combinations.

Measuring the transport speed of a vibratory feeder

As mentioned in the introduction, the most important feature of a vibratory feeder is the maximum speed at which it can transport the components. The foregoing theory gives the conditions under which each of the various forms of motion occur, so it is possible, in principle, to find the transport speed by calculating the distance over which the components move in one period. This calculation is only feasible, however, if certain simplifying assumptions are made, for example regarding the collisions between components and the track and between the components themselves. It is difficult to ascertain the validity of such assumptions. A simpler course is to measure the speed of the components directly at various values of the parameters.

A series of experiments was therefore carried out using a vibratory feeder built up from interchange-

able parts, thus permitting simple alteration of the parameters. The vibration amplitude x_0 was measured with a vibration pick-up, and the transport speed of the components with a stopwatch. The experimental set-up is shown in *fig. 7*. The vibration pick-up and ancillary equipment were used to measure, in addition to the amplitude, the velocity and the acceleration of the track (in this set-up the horizontal component of these quantities), and their variation was displayed on an oscilloscope.

It was found that, in spite of the electromagnetic drive, the vibration was almost purely harmonic. A second vibration pick-up, perpendicular to the first, revealed that the direction of the vibration was purely perpendicular to the leaf springs, proving that we were in fact dealing with a linear harmonic vibration, as assumed in the theory.

The leaf springs were situated at a distance of 100 mm from the vertical axis of the drum, a distance equal to the radius of the drum, so that the spring angle γ (*fig. 2*) was equal to $a + \beta$. The speed experiments were performed with frequencies from 50 to 100 c/s and the following values of γ and a :

$$\begin{aligned} \gamma &= 45^\circ, 35^\circ, 25^\circ, 15^\circ \text{ and } 7^\circ, \\ a &= 1^\circ, 1.5^\circ \text{ and } 2^\circ. \end{aligned}$$

The components used in the experiments were: bare carbon resistors with brass terminal caps, length 11 mm, diameter 3.4 mm; silver-plated contact springs for a switch ($12 \times 3.5 \times 3$ mm) and

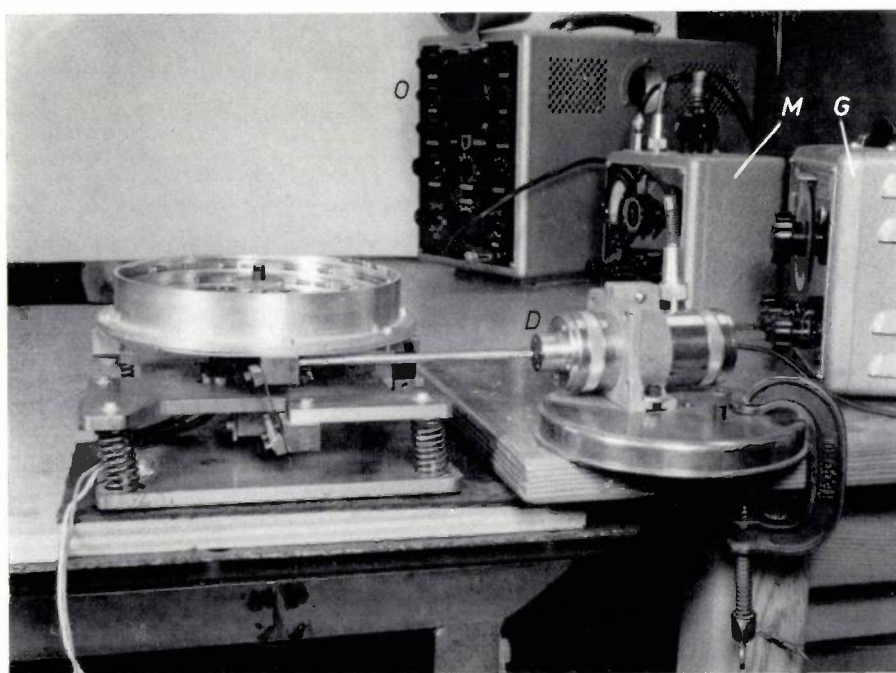


Fig. 7. A vibratory feeder (left) in a measuring set-up. The vibratory feeder is driven at variable frequencies by a signal generator G via an amplifier (not visible). By means of a vibration pick-up D (type PR 9261) and an amplitude measuring apparatus M (type PR 9250/01) the vibration is displayed on an oscilloscope O . The parts on the feeder track are carbon resistors.

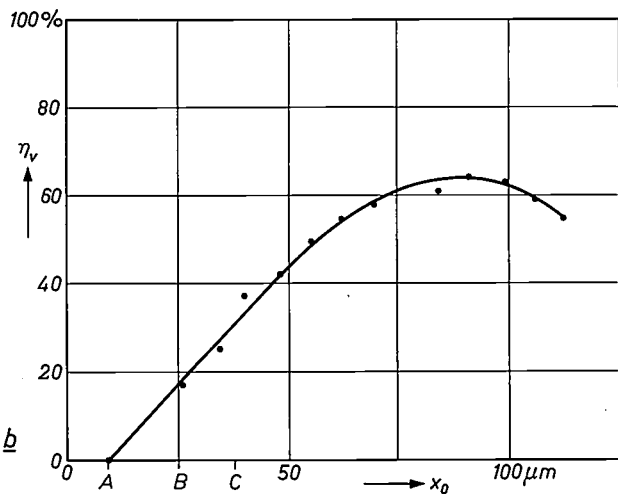
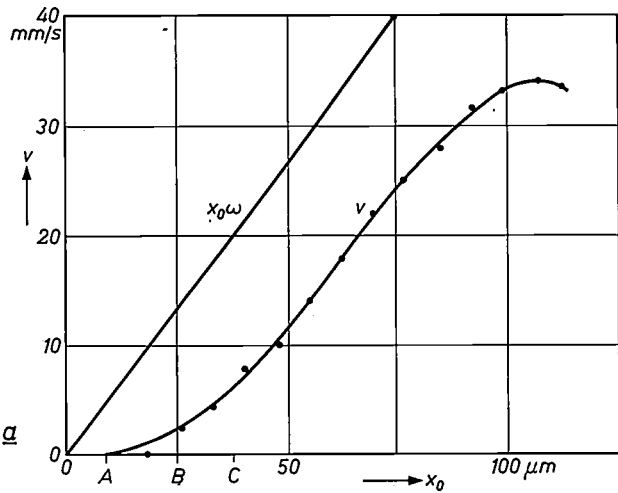


Fig. 8. a) Curve v : measured values of the speed v at which given components are transported by a particular vibratory feeder, as a function of the amplitude x_0 of the track; $\alpha = 1^\circ 27'$, $\beta = 43^\circ 33'$, $\alpha + \beta = 45^\circ$, $f = 86$ c/s. b) For the same case a plot of the "speed efficiency" $\eta_v = v/x_0\omega$ versus x_0 . At the values A , B and C of x_0 , the components, as x_0 is increased, begin respectively to slide upwards, to slide alternately upwards and downwards, and to leave the track.

carbon brushes ($9 \times 4 \times 2.5$ mm).

The transport speed was measured without the components being restricted in their movements, i.e. without any piling up at the bottom of the drum and without obstruction from sorting and orienting devices. The measured speed v was plotted versus the horizontal amplitude of vibration $x_0 \cos \alpha$ ($\approx x_0$), found with the vibration pick-up. As an example, fig. 8a shows the results of one of these experiments. Using the results of the model experiments, this graph also indicates the values of x_0 at which the components in succession started periodically sliding upwards (at $x_0 = A$), started alternately sliding upwards and downwards (at $x_0 = B$) and started to leave the track (at $x_0 = C$).

In fig. 8b we have plotted, versus x_0 , the "speed efficiency" η_v . This is defined as the ratio of the transport speed v to the amplitude $x_0\omega$ (also plotted in fig. 8a) of the track velocity in the forward direction:

$$\eta_v = \frac{v}{x_0\omega}$$

Obviously the transport speed can never exceed the velocity amplitude $x_0\omega$.

Finally, fig. 9 gives a survey of a number of speed experiments. This time the transport speed v is plotted as a function of the amplitude $x_0\omega^2$ of the track acceleration.

From our numerous series of experiments the following conclusions can be drawn.

- 1) A vibratory feeder works optimally when the vertical component of the acceleration has an amplitude between $2g$ and $2.5g$; see the values shown by dots on the curves in fig. 9. At these values of $x_0\omega^2$ the maximum useful transport speed is achieved at the respective value of γ . Although a large amplitude may result in a somewhat higher transport speed, the components are then bounced so violently that there can no longer be any question of sorting and orientation.
- 2) Provided that the acceleration amplitude ($x_0\omega^2$) is kept constant, variation of the frequency does not significantly affect the maximum useful transport speed.

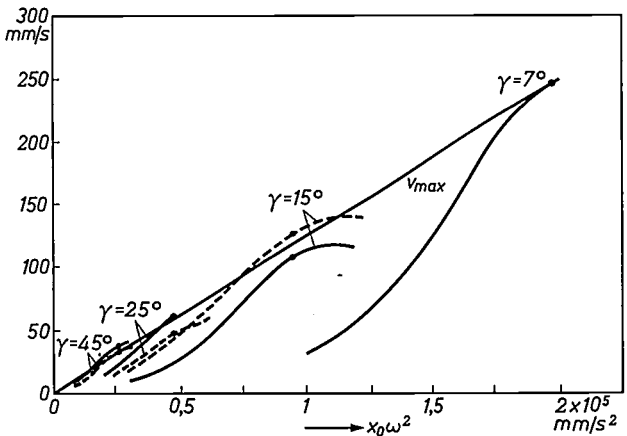


Fig. 9. Survey of the results of various speed experiments on actual vibratory feeders: v as a function of $x_0\omega^2$ for $\alpha = 1^\circ 27'$ and the indicated values of the spring angle γ . At each value of γ use is made of one or two frequencies f :

$\gamma = 45^\circ$...	$f = 66$ and 86 c/s,
$\gamma = 25^\circ$...	$f = 58$ and 81 c/s,
$\gamma = 15^\circ$...	$f = 86$ and 104 c/s,
$\gamma = 7^\circ$...	$f = 78$ c/s.

The dashed curves relate to the lower of the two frequencies. The values of the maximum useful transport speed, denoted by a dot on each curve, are found to lie approximately on a straight line (v_{max}).

- 3) For given vibration angle β , variation of the angle of inclination α between the values of 1° and 3° , which occur in practice, does not cause any significant change in the maximum useful transport speed.
- 4) The same applies to variation of the coefficient of friction μ .
- 5) The mass of each of the transported components has no influence, nor, within very wide limits, has the shape.

The experiments from which these conclusions were drawn were performed with the above-mentioned vibratory feeder, the drum diameter of which was 200 mm. The mass of each component was a few grammes or less. There is no reason to assume that the results obtained would not apply to other drum diameters and heavier components. This has been confirmed by some incidental observations on both smaller and larger vibratory feeders (diameters 80 and 500 mm).

Maximum useful speed

As just mentioned, for optimum operation the amplitude of the vertical component of the acceleration must have a value of between $2g$ and $2.5g$:

$$x_0\omega^2 \tan (\alpha + \beta) = (2 \text{ to } 2.5)g. \quad (6)$$

It appears that when this condition is fulfilled the speed efficiency also reaches its maximum value $\eta_{v \max}$. The maximum useful transport speed is thus:

$$v_{\max} = \eta_{v \max} x_0\omega = \frac{\eta_{v \max}}{\omega} x_0\omega^2 = \frac{\eta_{v \max}}{\omega} \times \text{constant}. \quad (7)$$

If η_v were to reach the same maximum value at every vibration frequency, the maximum transport speed would apparently be inversely proportional to the frequency. This is not the case, however: the maximum speed efficiency is found to be roughly directly proportional to the frequency. Within the frequency range from 50 to 100 c/s it has been found empirically that:

$$\eta_{v \max} = 0.785 \frac{f}{100}, \quad (8)$$

with f in c/s. At higher frequencies this relation obviously no longer applies, since η_v cannot be greater than 1; above 130 c/s the maximum useful transport speed would indeed be inversely proportional to the frequency.

From (7) and (8) it follows that:

$$v_{\max} = 1.25 \times 10^{-3} x_0\omega^2 \text{ mm/s}. \quad (9)$$

In combination with (6), and given $g \approx 10^4 \text{ mm/s}^2$, this gives:

$$v_{\max} = (25 \text{ to } 31) \cot (\alpha + \beta) \text{ mm/s}. \quad (10)$$

High transport speeds can thus be obtained if $\alpha + \beta$ is small.

Loading of the springs

Although we have seen that frequencies in the normal range have no influence on the transport speed, it is as well to take into account the relation between the frequency and the loading of the springs when the vibratory feeder is working optimally.

For a leaf spring rigidly clamped at both ends the maximum energy per unit volume that may be stored upon deflection is given by:

$$\frac{W}{V} = \frac{\sigma_b^2}{18E}.$$

Here W is the energy stored by the spring, V the volume of the spring, σ_b the maximum permissible amplitude of the alternating bending stress in the spring, and E the modulus of elasticity of the spring material. From (9) it was seen that the maximum useful transport speed v_{\max} is directly proportional to the acceleration amplitude, so that at v_{\max} , the maximum deflection e of the spring will be inversely proportional to the square of the chosen frequency f :

$$e \propto f^{-2}.$$

The vibrating mass and the springs form a system having a natural frequency f_0 , the square of which is proportional to the stiffness C of the springs:

$$C \propto f_0^2.$$

If we make the frequencies f and f_0 equal, then the energy $W = \frac{1}{2}Ce^2$ which is stored in the deflected springs is:

$$W = \frac{1}{2}Ce^2 \propto f^2 f^{-4} = f^{-2}.$$

A higher frequency means, therefore, that at the same alternating stress one can use springs of smaller volume, or with the same volume the bending stresses in the springs are lower.

The influence of the angle of inclination and the coefficient of friction

Although the angle of inclination α and the coefficient of friction μ have no direct influence on the maximum useful transport speed, they are indirectly of importance so far as they concern the critical angle of vibration β_{cr} , for according to (5), $\tan \beta_{cr} = \alpha/\mu^2$. A small α or a large μ gives a small β , and hence a small value of $\alpha + \beta$, thus giving a high maximum

useful transport speed (see eq. 10). In practice, however, it is difficult to vary μ to any appreciable extent: as mentioned above, the value of μ found in very divergent cases was always between 0.35 and 0.45.

The results of the various experiments now make it possible to analyse more exactly the behaviour of the individual components in the optimum setting. This analysis will be found in *Appendix II*.

Nomogram for the design calculation of vibratory feeders

Let us recapitulate the conditions which a properly functioning vibratory feeder is required to satisfy.

Upward transport is ensured when

$$\tan \beta \geq \frac{a}{\mu^2} \dots \dots (5)$$

The maximum useful transport speed is obtained when the vertical component of the amplitude of acceleration satisfies

$$x_0 \omega^2 \tan (a + \beta) = (2 \text{ to } 2.5)g \dots \dots (6)$$

This speed is then

$$v_{\max} = 1.25 \times 10^{-3} x_0 \omega^2 = 0.05 x_0 f^2 \text{ mm/s. } (9)$$

With the aid of these relations a nomogram has been drawn (*fig. 10*), from which the design of a vibratory feeder can easily be calculated.

From the origin radial lines are drawn which with the horizontal axis $I-I'$ enclose angles $a + \beta$. The

parameter $x_0 f^2 (= x_0 \omega^2 / 4\pi^2)$ is plotted along the axis $I-I'$. This parameter is a measure of the amplitude of the acceleration along the track. In view of the small value of a , this amplitude may be assumed to be equal to the amplitude of the horizontal component of the acceleration. Along the vertical axis one should thus be able to read on the same scale (not shown in the figure) the value of the amplitude of the vertical acceleration component for various angles $a + \beta$. The lines $II-II'$ and $III-III'$ correspond to a vertical acceleration of $2g$ and $2.5g$ respectively, and thus form the boundaries of the region in which the operation of the vibratory feeder is optimum. The line $III-IV'$ represents graphically the equation $v_{\max} = 0.05 x_0 f^2$.

The maximum useful transport speed at a given angle $a + \beta$ (irrespective of the values of a, μ and f) can be found by vertically projecting on the line $III-IV'$ the section $C'C$ which the lines $II-II'$ and $III-III'$ cut off from the line IC drawn from the origin at the chosen angle $a + \beta$. This gives the points D' and D . The ordinates of D' and D can be read on the left in mm/s and give the limits within which v_{\max} will lie. The interval between these limits is due to the difference between $2g$ and $2.5g$; in this way account is taken of the experimental spread and of the approximations introduced when interpreting the results, without either complicating the nomogram or making it too inaccurate.

We now have to take the critical vibration angle into account. The three curves in *fig. 10* are based on

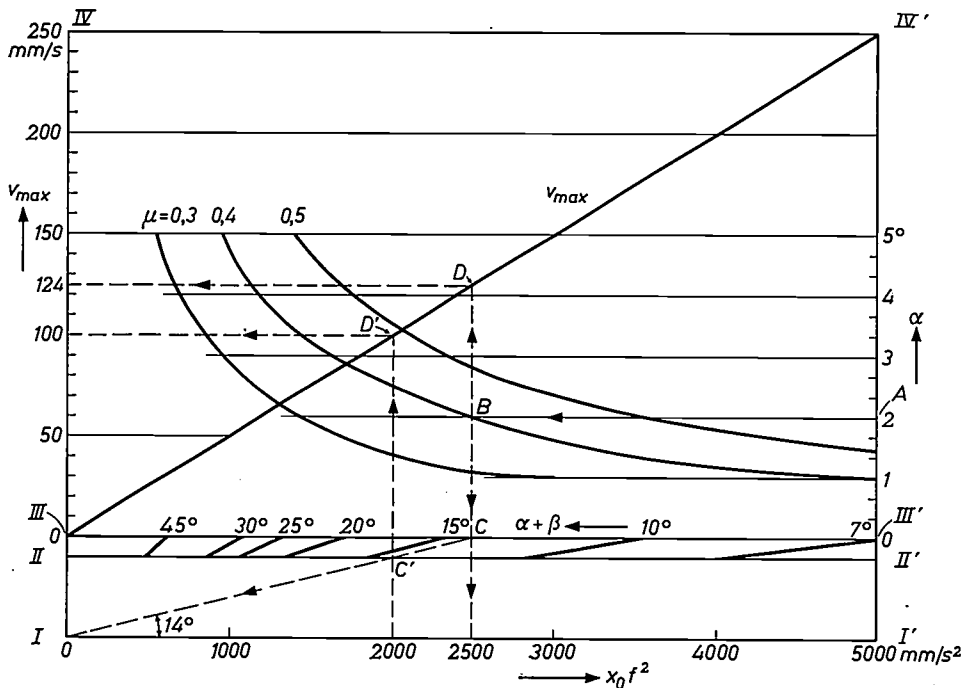


Fig. 10. Nomogram for design calculations of vibratory feeders.

the relationship $\tan \beta_{cr} = a/\mu^2$ for $\mu = 0.3, 0.4$ and 0.5 , the a values (see the right-hand scale) following from

$$\alpha + \beta = \tan^{-1} \frac{\alpha(1 + \mu^2)}{\mu^2 - \alpha^2}$$

being plotted here as a function of the values $\alpha + \beta$ on the abscissa $III-III'$. (With the simplification $\tan \alpha \approx \alpha$ — for small α — this expression follows directly from $\tan \beta = a/\mu^2$.) The minimum required angle $\alpha + \beta$ can now be read from the horizontal axis $III-III'$, perpendicularly below the point where the curve for the relevant value of μ intersects the horizontal line indicating α . If the distance R_s from the springs to the drum axis is equal to the radius R_t of the track (which is often the case), the angle read from the nomogram is also the minimum required spring angle γ ; in general, however,

$$\tan \gamma = \frac{R_t}{R_s} \tan (\alpha + \beta).$$

We can also read from the nomogram the vibration amplitude of the springs at a given spring angle γ , this amplitude (in mm) being equal to the length IC , measured on the scale of $x_0 f^2$ and divided by f^2 . Of course, this is only the deflection of the end of the spring at the drum side. To obtain the total deflection of the spring we must add the deflection at the other end attached to the counterweight.

The nomogram thus makes it possible to calculate the maximum bending stress in the springs. This stress σ_b is:

$$\sigma_b = \frac{3Ehe}{l^2},$$

where h is the thickness, l the length and e the total deflection of the spring.

As an illustration we shall give a worked-out example.

Example

For constructional reasons (drum diameter and pitch of the track, determined among other things by the dimensions of the components to be conveyed) we have, say, $\alpha \geq 2^\circ$.

Let μ be 0.4 (μ can be measured by determining the tangent of the angle of inclination at which a component just begins to slide with the track stationary).

The minimum spring angle required at the smallest useful angle of inclination ($\alpha = 2^\circ$) is found as follows. Through point A for $\alpha = 2^\circ$ (fig. 10) we draw a horizontal line. This cuts the curve for $\mu =$

0.4 at B . From B we draw a line BC perpendicular to the horizontal axis $III-III'$ and join C to I . The angle which IC makes with the axis $I-I'$ is the required minimum angle $\alpha + \beta$, here 14° . The critical vibration angle is therefore $14^\circ - 2^\circ = 12^\circ$.

Depending on the required transport speed we make the angle $\alpha + \beta$ equal to or greater than 14° . For $\alpha + \beta = 14^\circ$ the expected transport speed, as can be read from the vertical axis $III-IV'$ via D' and D , lies between 100 and 124 mm/s.

If this speed is too low we must reduce the critical vibration angle, that is, in accordance with (5), make α smaller and/or make μ larger (the latter, however, is difficult). Given $\alpha = 1^\circ$ we can reduce $\alpha + \beta$ to 7° , which brings v_{max} between 204 and 250 mm/s. The effect of reducing α , however, might be to make the track pitch smaller than the lower limit imposed on the pitch by the size of the components. To avoid this the drum diameter can be increased.

In the case where $R_s = R_t$ and $\alpha + \beta = 14^\circ$, the length IC corresponds to an acceleration of $2500/\cos 14^\circ = 2527$ mm/s². The vibration amplitude of the drum end of the spring, perpendicular to the plane of the spring, is then

$$\begin{aligned} \text{at } f = 100 \text{ c/s} & \dots \dots 0.26 \text{ mm} \\ \text{and at } f = 50 \text{ c/s} & \dots \dots 1.03 \text{ mm.} \end{aligned}$$

If the counterweight has, for example, the same mass and the same moment of inertia as the drum, the total deflection of the springs is twice as large.

Appendix I: Calculation of the conditions for the various forms of motion of the parts in a vibratory feeder

Consider a part of mass m situated on an inclined track (angle of inclination a) which describes a harmonic vibration of frequency $f = \omega/2\pi$ at an angle β with the plane of the track.

The vibration can be resolved into a component $x = x_0 \sin \omega t$ along the track and a component $y = y_0 \sin \omega t$ perpendicular to the track (fig. 2). As long as the part follows the motion of the track, it obeys these equations of motion:

$$\begin{aligned} \text{along the track} & \dots \quad F_f - mg \sin a = m\ddot{x} = -mx_0\omega^2 \sin \omega t, \\ \text{perpendicular to track} & \quad F_n - mg \cos a = m\ddot{y} = -my_0\omega^2 \sin \omega t. \end{aligned}$$

Here F_f represents the frictional force and F_n the component of the reaction which the track exerts on the part; $mg \sin a$ and $mg \cos a$ are the components of the force of gravity along the track and perpendicular to it.

As long as $|F_f| < \mu F_n$, the part follows the motion of the track. To discover whether this condition is satisfied, we have plotted in fig. 11a, as a function of the phase ωt of the vibration, the quantities

$$\left. \begin{aligned} |F_f| &= |mg \sin a - mx_0\omega^2 \sin \omega t| \\ \mu F_n &= \mu mg \cos a - \mu my_0\omega^2 \sin \omega t. \end{aligned} \right\} \dots (11)$$

Fig. 11b also gives the curves of the deflection x , the velocity \dot{x} and the acceleration \ddot{x} of the track in the x direction. At the values of the parameters chosen here the condition $|F_f| < \mu F_n$ is evidently constantly fulfilled; the part thus follows the motion of the track.

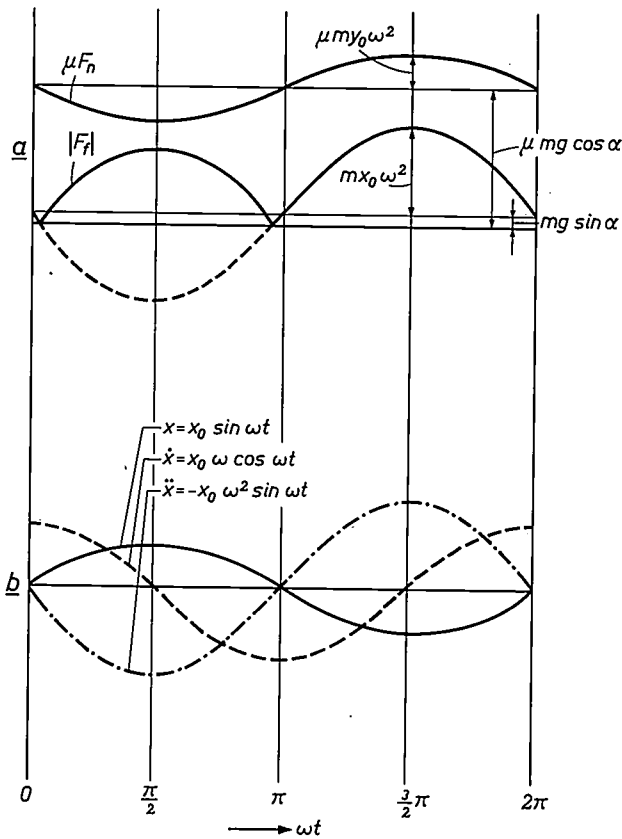


Fig. 11. Variation with ωt of (a) forces and (b) displacement x , velocity \dot{x} and acceleration \ddot{x} in the x direction of a vibratory feeder, given $\alpha = 3^\circ$, $\beta = 35^\circ$, $f = 100$ c/s, $\mu = 0.45$, $x_0 = 7.6 \mu\text{m}$. In (a) $|F_f|$ is the absolute value of the frictional force and F_n the normal component of the reaction which the track exerts on the part.

If, however, $|F_f|$ becomes equal to μF_n — e.g. because the amplitude is increased — the part will start sliding. It then depends on the phase at which this happens whether the part slides upwards or downwards. And if F_n becomes zero, the part leaves the surface of the track.

The curves in fig. 12a-d represent $|F_f|$ and μF_n at successively larger amplitudes. Fig. 12a corresponds to fig. 11a and is included for the sake of completeness. In fig. 12b the condition $|F_f| = \mu F_n$ is satisfied at $\omega t = \frac{1}{2}\pi$ (point A). At this instant the vibrating track reaches its topmost position and the part therefore starts periodically sliding upwards. If at $\omega t = \frac{3}{2}\pi$, in the descending stroke, $|F_f|$ becomes equal to μF_n (the point B in fig. 12c) the part will alternately slide upwards and downwards. With increasing amplitude a point is finally reached, at $\omega t = \frac{1}{2}\pi$, where the condition $\mu F_n = 0$ is also satisfied (point C in fig. 12d); from that instant the part will leave the surface of the track momentarily once in every period.

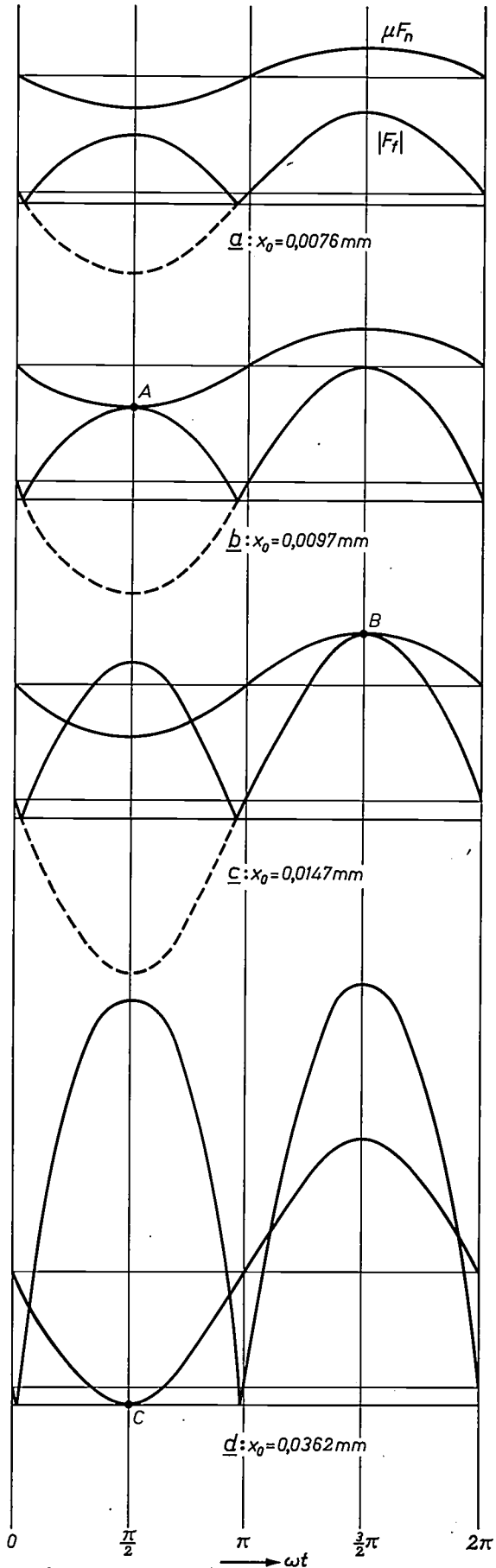
Fig. 12. The forces $|F_f|$ and μF_n as functions of ωt , for an amplitude x_0 increasing from (a) to (d).

Case (a) (identical with that of fig. 11a): $x_0 = 7.6 \mu\text{m}$. The force $|F_f|$ is constantly smaller than μF_n ; the part follows the movement of the track.

Case (b): $x_0 = 9.7 \mu\text{m}$. At point A, where $\omega t = \frac{1}{2}\pi$, the friction $|F_f| = \mu F_n$. The part begins to slide upwards.

Case (c): $x_0 = 14.7 \mu\text{m}$. At point B, where $\omega t = \frac{3}{2}\pi$, the friction $|F_f| = \mu F_n$. The part starts sliding alternately upwards and downwards.

Case (d): $x_0 = 36.2 \mu\text{m}$. At point C, where $\omega t = \frac{1}{2}\pi$, the friction $F_f = 0$. Once in every period the part leaves the track.



Since the curves in fig. 12 are simply graphical representations of the equations of motion of a part that follows the movement of the track, they give no indication of the motion of a part which is sliding or which is separated from the track. In other words, the curves lose their significance beyond point *A*, *B* or *C*. The figures do, however, give the sequence in which the cases *A*, *B* and *C* will occur as the amplitude increases. It is also clear that, if the upward sliding begins at a smaller amplitude than the downward sliding, the former will always have a greater initial speed than the latter; moreover, under the condition mentioned, the friction resisting the sliding motion is in fact less in the upward than in the downward direction. In all respects, then, the upward sliding predominates.

It also appears that the mass of the part has no influence on the occurrence of sliding and bouncing, for equations (11) can be divided by *m* without altering the relation between *F_f* and *μF_n*. This is bound up with the fact that the whole process is an interplay of forces; the important quantity is not the vibration amplitude *x₀* but the acceleration amplitude *x₀ω²*.

The sequence in which the upward and downward sliding and leaving the track occur is governed by the parameters *α*, *β* and *μ*. The smallest acceleration amplitude *x₀ω²* at which once in every period, during the upward stroke of the vibration, the condition $|F_f| = \mu F_n$ is just fulfilled (point *A* in fig. 12*b*) follows from:

$$\mu mg \cos \alpha - \mu m y_0 \omega^2 = m x_0 \omega^2 - mg \sin \alpha.$$

Given $y_0 = x_0 \tan \beta$ and for small values of *α* ($\sin \alpha \approx \alpha$ and $\cos \alpha \approx 1$) this changes to:

$$\frac{x_0 \omega^2}{g} = \frac{\mu \cos \alpha + \sin \alpha}{1 + \mu \tan \beta} \approx \frac{\mu + \alpha}{1 + \mu \tan \beta}.$$

The part thus starts periodically sliding when:

$$\frac{x_0 \omega^2}{g} \geq \frac{\mu + \alpha}{1 + \mu \tan \beta}.$$

In a similar way it can be deduced that the boundary case of point *B* in fig. 12*c* is determined by:

$$\frac{x_0 \omega^2}{g} = \frac{\mu - \alpha}{1 - \mu \tan \beta}.$$

Thus it follows that the part will periodically slide downwards when

$$\frac{x_0 \omega^2}{g} \geq \frac{\mu - \alpha}{1 - \mu \tan \beta}.$$

The smallest value of *x₀ω²* at which *F_n* just becomes zero once in every period (point *C* in fig. 12*d*) follows from:

$$\mu mg \cos \alpha - \mu m y_0 \omega^2 = 0$$

or
$$\frac{x_0 \omega^2}{g} = \frac{\cos \alpha}{\tan \beta} \approx \cot \beta.$$

In this way we have derived the conditions under which the various forms of motion occur.

Appendix II: Behaviour of an individual part on an optimally adjusted vibrating track

The experimental results described on pages 89-90 permit a more detailed analysis of the behaviour of an individual part on a vibrating track with optimum setting. The distance covered per period is:

$$\eta_v x_0 \omega f^{-1} = 2\pi \eta_v x_0.$$

In practice η_v lies between 40 and 80%, so that per period the part covers a distance along the track of about 2.5 to 5 times the vibration amplitude *x₀*.

The motion in the *y* direction, perpendicular to the track, is represented for the idealized case in fig. 13. The vertical acceleration component — which, owing to the small value of *α*, is virtually equal to the acceleration in the *y* direction — has an

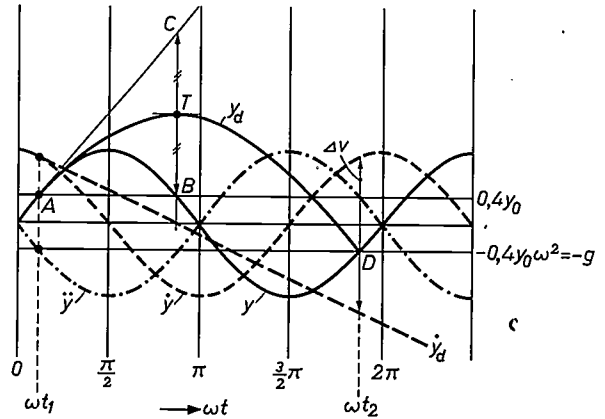


Fig. 13. Analysis of the manner in which a separate part of a vibrating track moves in the *y* direction, on the assumption that $y_0 \omega^2 = 2.5g$. The values plotted are the displacement *y*, the velocity \dot{y} and the acceleration \ddot{y} of the track. At ωt_1 the part leaves the track and describes a parabola *ATD*; \dot{y}_d shows the vertical velocity component of the part. At ωt_2 (point *D*) the part lands again on the track; the consequences of the impact (velocity difference Δv) are not considered.

amplitude of, say, 2.5*g*. The phase ωt_1 at which the part leaves the track follows from the fact that the vertical acceleration of the part at that moment, $-2.5g \sin \omega t_1$, is equal to $-g$, so that $\sin \omega t_1 = 1/2.5 = 0.4$ and $\omega t_1 = 23.6^\circ$. At this moment the deflection of the track in the *y* direction is $0.4y_0$ and the velocity in that direction is $y_0 \omega \cos \omega t_1 = 0.916y_0 \omega$. These values are also the initial conditions for the now following free motion of the part in the *y* direction. In relation to the neutral position of the track, the equation of this motion is:

$$y_d = y_0 \sin \omega t_1 + (y_0 \cos \omega t_1)(\omega t - \omega t_1) - \frac{g}{2\omega^2}(\omega t - \omega t_1)^2.$$

The moment *t₂* at which the part, after describing part of a parabola, again joins the track can therefore be found from:

$$y_0 \sin \omega t_2 = y_0 \sin \omega t_1 + (y_0 \cos \omega t_1)(\omega t_2 - \omega t_1) - \frac{g}{2\omega^2}(\omega t_2 - \omega t_1)^2,$$

or, since $y_0 \omega^2 = 2.5g$:

$$\sin \omega t_2 = \sin \omega t_1 + (\cos \omega t_1)(\omega t_2 - \omega t_1) - \frac{1}{5}(\omega t_2 - \omega t_1)^2. \quad (12)$$

This equation gives the relation between ωt_2 and ωt_1 . Even without solving this to find ωt_2 , it is plain that ωt_2 is independent of frequency, for the equation only contains phase angles.

The parabola described by the part is easily constructed. The parabola begins at point *A* (fig. 13) of the sine curve *y* at ωt_1 (where $\ddot{y} = -g$) and at this point has a common tangent *AC* with the sine curve. From this the rest follows. The top *T* of the parabola lies halfway between *B* and *C* directly above the point where the vertical velocity \dot{y}_d of the part is zero. The right half of the parabola meets the sine curve *y* again at point *D*, which is thus the point where the part again lands on the track. It is evident that the abscissa ωt_2 of *D* cannot be very different from 2π .

This abscissa can be calculated approximately by expanding $\sin \omega_2 t$ into a series and neglecting all terms in that series except the first one. Eq. (12) then becomes a quadratic equation in $\omega_2 t$. Only one of the two roots is of interest here, which is $\omega_2 t = 331.3^\circ$ (accurate up to about 1°). As found above, $\omega_2 t$ is independent of frequency.

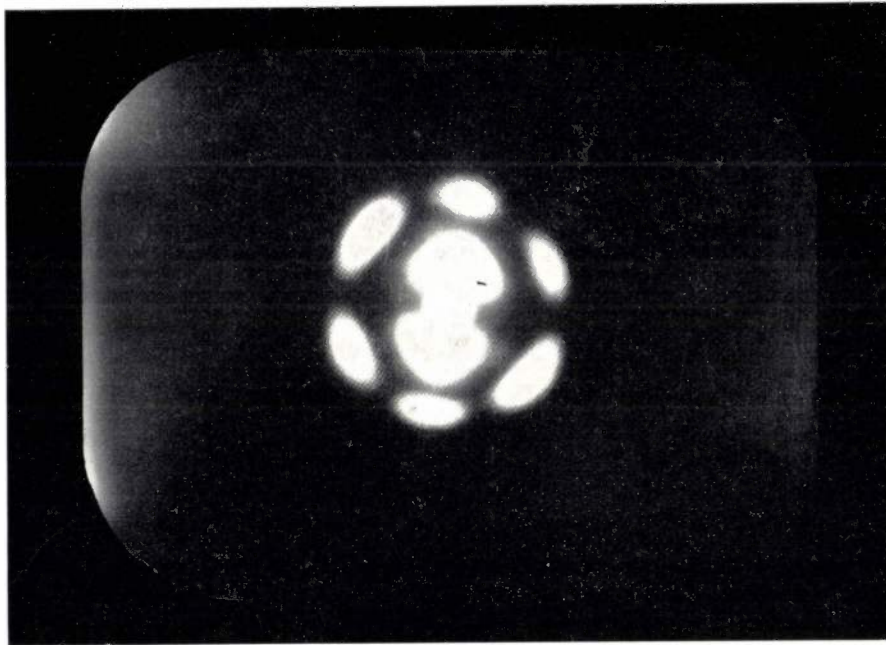
It should be noted at this point that the parabola corresponds only to an idealized representation, for we disregard the fact that at point *D* the part and the track collide with each other with a velocity difference Δv (see fig. 13). The velocity

difference Δv is proportional to the velocity amplitude $y_0 \omega$ of the track. Since, in an optimally adjusted vibratory feeder, this amplitude is inversely proportional to the frequency, this accounts to a large extent for the drop in the speed efficiency $\eta_{v \max}$ at lower frequencies.

The above considerations show that the vibratory feeder gives optimum operation when the part in principle just touches the track once in every period, and then only when the track just reaches its maximum velocity in the upward direction.

Summary. Vibratory feeders are widely used for conveying small components in an ordered sequence to production processes. The components are transported by means of a vibrating inclined track. Parts on such a track can be made to behave in one of the following ways: they can follow the movement of track, and are thus not transported; they can slide periodically either upwards or downwards; they can slide alternately upwards and downwards; or they can periodically leave the track. In a theoretical treatment the author first considers which parameters govern each of these motions, and in partic-

ular analyses the conditions in which upward or downward transport is to be expected. The results of experiments using a model are then described. These show satisfactory agreement with the theory. The theory, which was supplemented by systematic velocity experiments on actual vibratory feeders, supplies the basis for a nomogram useful for design calculations. The application of the nomogram is explained with the aid of a worked-out example. Consideration is also given to the loading of the springs on which the vibratory feeder is mounted.



A SMALL, STABLE GAS LASER

537.525:535.339

The electromagnetic radiation emitted by an unmodulated *radio transmitter* has the form of a single, practically continuous sine wave. It is thus coherent and monochromatic. The radiation emitted by a normal *light source*, e.g. a sodium lamp, on the other hand, consists (even if we only consider one single spectral line) of innumerable waves, emitted with random phase by the atoms involved. It is thus incoherent, and the spectral lines are relatively broad.

Now some years ago a new type of light source was introduced which, like a radio transmitter, emits a coherent and monochromatic radiation:

the LASER. This name is formed from the initials of "Light Amplification by Stimulated Emission of Radiation". (One also sometimes speaks of an "optical MASER", where M stands for Microwave.) The above-mentioned properties, and yet others due to the special design of this light source (e.g. a sharp beaming of the emitted power), make the laser of great scientific interest and open wide perspectives for its practical application¹⁾.

¹⁾ A. L. Schawlow and C. H. Townes, *Infrared and optical masers*, *Phys. Rev.* **112**, 1940-1949, 1958. A good review, with extensive references, is to be found in W. Kaiser, *Der optische Maser*, *Physica status solidi* **2**, 1117-1145, 1962 (No. 9).

This abscissa can be calculated approximately by expanding $\sin \omega_2 t$ into a series and neglecting all terms in that series except the first one. Eq. (12) then becomes a quadratic equation in $\omega_2 t$. Only one of the two roots is of interest here, which is $\omega_2 t = 331.3^\circ$ (accurate up to about 1°). As found above, $\omega_2 t$ is independent of frequency.

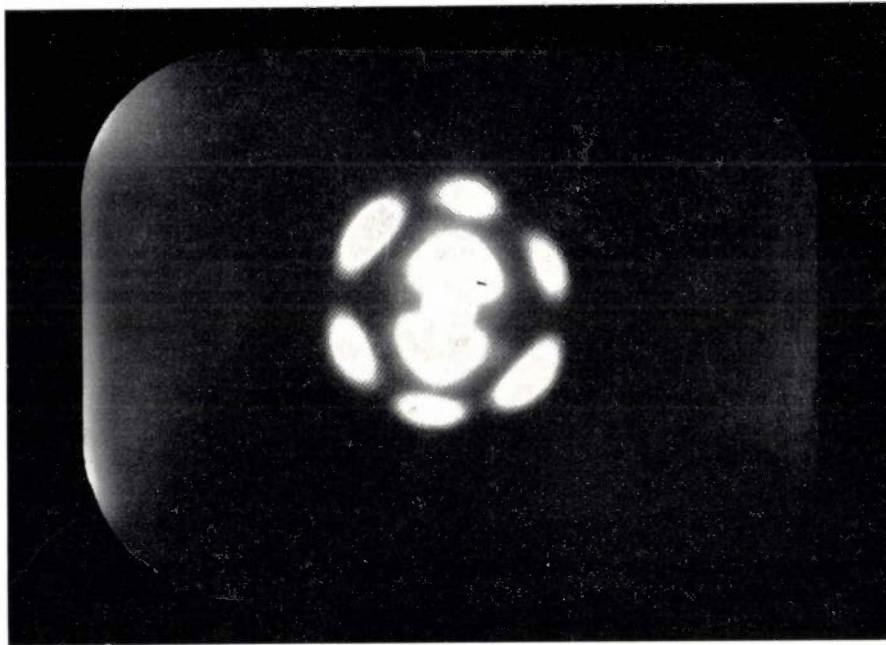
It should be noted at this point that the parabola corresponds only to an idealized representation, for we disregard the fact that at point *D* the part and the track collide with each other with a velocity difference Δv (see fig. 13). The velocity

difference Δv is proportional to the velocity amplitude $y_0 \omega$ of the track. Since, in an optimally adjusted vibratory feeder, this amplitude is inversely proportional to the frequency, this accounts to a large extent for the drop in the speed efficiency $\eta_{v \max}$ at lower frequencies.

The above considerations show that the vibratory feeder gives optimum operation when the part in principle just touches the track once in every period, and then only when the track just reaches its maximum velocity in the upward direction.

Summary. Vibratory feeders are widely used for conveying small components in an ordered sequence to production processes. The components are transported by means of a vibrating inclined track. Parts on such a track can be made to behave in one of the following ways: they can follow the movement of track, and are thus not transported; they can slide periodically either upwards or downwards; they can slide alternately upwards and downwards; or they can periodically leave the track. In a theoretical treatment the author first considers which parameters govern each of these motions, and in partic-

ular analyses the conditions in which upward or downward transport is to be expected. The results of experiments using a model are then described. These show satisfactory agreement with the theory. The theory, which was supplemented by systematic velocity experiments on actual vibratory feeders, supplies the basis for a nomogram useful for design calculations. The application of the nomogram is explained with the aid of a worked-out example. Consideration is also given to the loading of the springs on which the vibratory feeder is mounted.



A SMALL, STABLE GAS LASER

537.525:535.339

The electromagnetic radiation emitted by an unmodulated *radio transmitter* has the form of a single, practically continuous sine wave. It is thus coherent and monochromatic. The radiation emitted by a normal *light source*, e.g. a sodium lamp, on the other hand, consists (even if we only consider one single spectral line) of innumerable waves, emitted with random phase by the atoms involved. It is thus incoherent, and the spectral lines are relatively broad.

Now some years ago a new type of light source was introduced which, like a radio transmitter, emits a coherent and monochromatic radiation:

the LASER. This name is formed from the initials of "Light Amplification by Stimulated Emission of Radiation". (One also sometimes speaks of an "optical MASER", where M stands for Microwave.) The above-mentioned properties, and yet others due to the special design of this light source (e.g. a sharp beaming of the emitted power), make the laser of great scientific interest and open wide perspectives for its practical application¹⁾.

¹⁾ A. L. Schawlow and C. H. Townes, *Infrared and optical masers*, *Phys. Rev.* **112**, 1940-1949, 1958. A good review, with extensive references, is to be found in W. Kaiser, *Der optische Maser*, *Physica status solidi* **2**, 1117-1145, 1962 (No. 9).

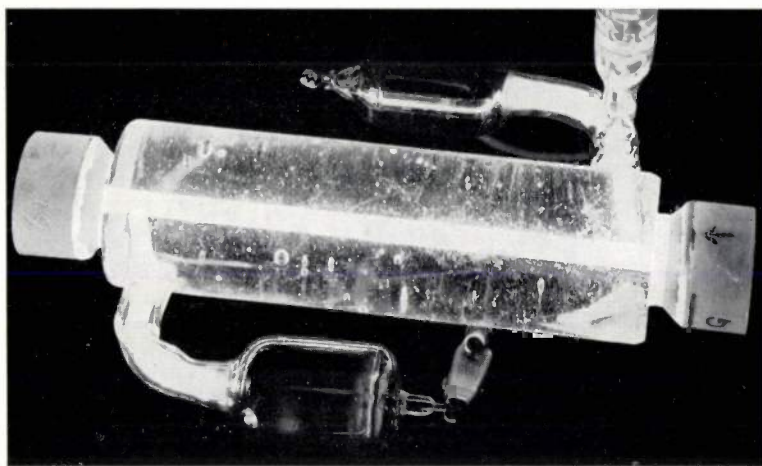


Fig. 1. The small gas laser. It consists of a heavy cylindrical block of fused-quartz 12 cm long, whose ends are polished flat to within $0.1 \mu\text{m}$, with a hole only 3 mm in diameter bored along the axis. Two fused-quartz blocks which are polished flat to within even closer limits and are provided with infrared-reflecting interference mirrors (reflection coefficient 99%) are placed against the ends of the cylinder, where they give a vacuum-tight seal by adhesion alone. The mirrors thus do not need to be sealed on by the application of heat. The great thickness of the cylinder (35 mm) is necessary to give the system sufficient mechanical rigidity (no deformation due to e.g. gravity) and to ensure that the temperature distribution in the "resonant cavity" of the laser has very good cylindrical symmetry.

The electrical discharge in the gas mixture with which the tube is filled (85% He + 15% Ne, pressure 3 torr) is produced by applying a DC voltage between two electrodes which are sealed into side tubes which connect with transverse borings in the fused-quartz cylinder. The length of the active gas column between the two transverse borings is 10 cm.

Many research teams are therefore at present engaged in investigations on and development of the laser.

Our own work in this direction has recently led to a result which we would like to describe here briefly²⁾. We hope to make it the subject of a more extensive article in due course.

The work in question concerns the *gas laser* (as opposed to the solid-state laser, which we shall not consider here). The stimulated-emission effect, on which the operation of the laser is based, consists in the fact that the probability of an atom making

the transition from an excited state to a state of lower energy, and at the same time emitting radiation of the corresponding wavelength, is increased if the atom in question is situated in a field of radiation of the same wavelength; moreover, the transition then occurs at such a moment that the emitted radiation is *in phase* with the stimulating radiation. The latter is thus *amplified*. The amplification which can be obtained in a gas is generally small, e.g. a few percent per metre path length of the radiation in the gas. Very large amplifications however can be obtained in a relatively short gas column if parallel plane *mirrors* are placed at each end, and the distance between the mirrors is adjusted to exactly a whole number of half wavelengths of the radiation: the radiation can then pass to and fro many times through the gas, thus giving a kind of positive feedback exactly in phase; the resulting amplification may even be so great that self-excitation of the radiation becomes possible. We then

have not a light *amplifier* (analogous to the maser for microwaves³⁾) but an *oscillator* for light waves, i.e. a light *source* which can continuously emit radiation.

The stimulated-emission effect can only be perceived if there are enough atoms in the required excited state. This is ensured by means of a (usually fairly complicated) excitation mechanism with its own energy source ("optical pump"). In the gas laser in question, a mixture of neon and helium is used as the medium. The neon atoms give stimulated emission at a wavelength of $1.153 \mu\text{m}$, i.e. in the near infrared.

³⁾ See e.g. R. W. DeGrasse, D. C. Hogg, E. A. Ohm and H. E. D. Scovil, Ultra-low-noise measurements using a horn reflector antenna and a traveling-wave maser, *J. appl. Phys.* **30**, 2013, 1959.

²⁾ See also: H. G. van Bueren, J. Haisma and H. de Lang, *Physics Letters (Amsterdam)* **2**, 340, 1962 (No. 7).

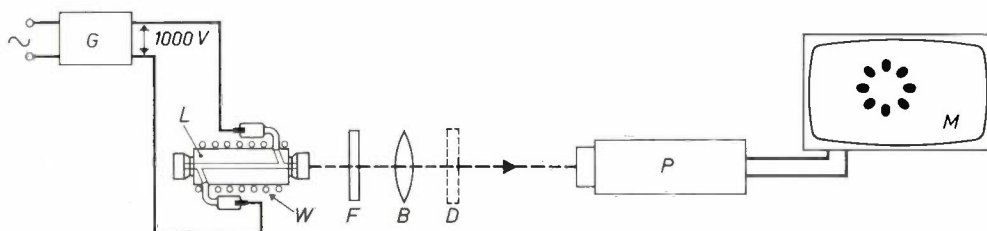


Fig. 2. The whole set-up. *L* laser tube, with voltage supply *G*. The laser radiation emerging from one of the end windows is focused on the photocathode of the infrared-sensitive lead-oxide vidicon *P* by means of the lens *B*. The filter *F* cuts off the visible radiation of the gas discharge, which is not due to laser action. *D* is a rotatable polarizer which is used to investigate the state of polarization of the emitted radiation. *M* television monitor on which the spatial distribution of the emitted laser radiation is made visible. *W* heater wire.

The helium atoms can take care of the pump effect if they in their turn are excited by an electric field. A mixture of neon and helium of suitable composition (other necessary conditions being fulfilled) can therefore be made to exhibit laser action if it is simply placed between parallel mirrors as described above and a gas discharge is produced in it⁴).

Originally, gas columns 50 to 100 cm in length were used in order to obtain sufficient amplification. There is much to be said for using shorter columns: the set-up then becomes much more manageable and less sensitive to mechanical vibrations, and it is easier to satisfy the condition that the mirrors should be accurately parallel. We have now succeeded in making a laser with a gas column only 10 cm long, owing among other things to the use of a narrow channel for the gas discharge, and to the adoption of a type of mirror with reflection and other losses of less than 1% and with deviations from flatness of not more than about $0.03 \mu\text{m}$ and of not more than 1 second of arc from the proper alignment. This tube is shown in *fig. 1*; some constructional details are given in the caption. The construction is very simple, and no difficult adjustments are required to obtain the laser action. The accurate adjustment of the distance between the two mirrors (which can also be regarded as bringing the electromagnetic resonant cavity, formed by the laser tube, into resonance) is carried out by means of a heating wire wound round the tube, which gives an adjustable expansion. The continuously radiated power amounts to some tenths of a milliwatt.

Fig. 2 shows the whole experimental set-up. The DC supply for the laser is very simple, as it only has to deliver about 10 watts at a voltage of about 1000 volts. The whole set-up is therefore easily transportable. The infrared laser radiation obtained is made visible by means of closed-circuit television using a lead-oxide vidicon ("Plumbicon") which has been developed in these laboratories and whose photoconductive layer is sensitive to infrared⁵). This layer is placed in the focal plane of a lens which is situated in line with the laser. An image is then produced on the screen of the TV picture tube e.g. like that shown in the title photograph. The regular pattern of spots of light seen here may be ascribed to narrow beams of light leaving the laser tube at very small angles, and which probably correspond to various standing waves in the tube (various resonance states due to different modes of vibration in the resonant cavity). The different beams of light also differ in their polarization direction. A detailed interpretation of these phenomena cannot be given here. The strong dependence of the laser effect on the length of the tube appears most clearly if one observes the image on the screen for some time during the warming-up or cooling-down of the tube: the pattern will then be seen to change gradually, owing to the appearance and disappearance of the various modes. This has been recorded on film, and we show a few frames of such a film in *fig. 3*.

J. HAISMA *),
S. J. van HOPPE *),
H. de LANG *),
J. van der WAL *).

⁴) A. Javan, W. R. Bennett, Jr. and D. R. Herriott, Population inversion and continuous optical maser oscillation in a gas discharge containing a He-Ne mixture, *Phys. Rev. Letters* **6**, 106-110, 1961 (No. 3).

⁵) E. F. de Haan *et al.*, not yet published.

*) Philips Research Laboratories, Eindhoven.

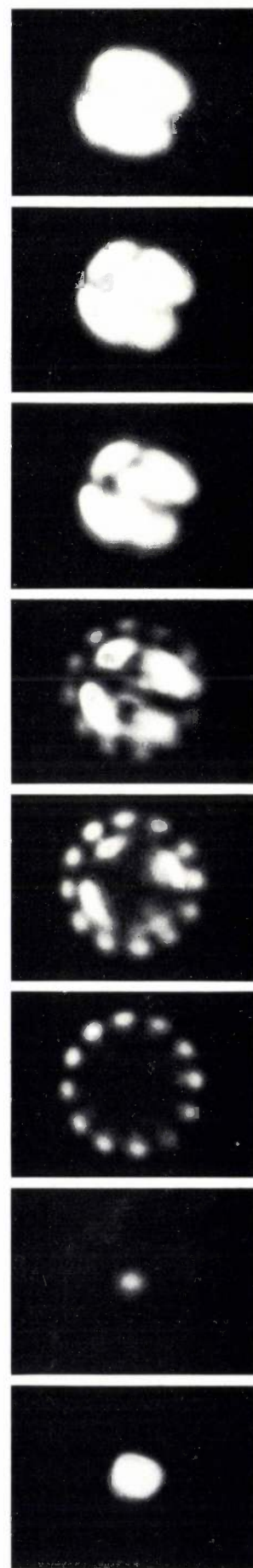


Fig. 3. Series of 16-mm film frames, as recorded from the screen of the TV tube. Time interval between successive frames 3 seconds.

ABSTRACTS OF RECENT SCIENTIFIC PUBLICATIONS BY THE STAFF OF N.V. PHILIPS' GLOEILAMPENFABRIEKEN

Reprints of these papers not marked with an asterisk * can be obtained free of charge upon application to the Philips Research Laboratories, Eindhoven, Netherlands, where a limited number of reprints are available for distribution.

- 2993:** A. Venema: Vacuum techniques in Europe (Advances in electron tube techniques, Proc. 5th U.S. Nat. Conf., New York, Sept. 1960, editor D. Slater, pp. 166-171, Pergamon, London 1961).

A brief survey of the present state of vacuum techniques in Europe, with special reference to pumping systems (including those for ultra-high vacuum) and gas analysis.

- 2994:** P. Massini: Movement of 2,6-dichlorobenzonitrile in soils and in plants in relation to its physical properties (Weed Res. **1**, 142-146, 1961, No. 2).

2,6-dichlorobenzonitrile (code name H133) is a weed-killer. It is of importance to know how this substance can best be applied. The movement of H133 in the soil and in plants has therefore been investigated, using samples labelled with ^{14}C in the nitrile group. It has been found that the movement of H133 in the soil and in plants is mainly determined by the following properties of the H133: its relatively high volatility, its low solubility in water and its strong adsorption on lignin, humus and fats. Indications of a transformation of H133 in plants were also found.

- 2995:** T. Kralt, W. J. Asma, H. H. Haecck and H. D. Moed: Reserpine analogues, I. β -indolyethylamine derivatives (Rec. Trav. chim. Pays-Bas **80**, 313-324, 1961, No. 4).

- 2996:** T. Kralt, W. J. Asma and H. D. Moed: Reserpine analogues, II. β -phenylethylamine derivatives (Rec. Trav. chim. Pays-Bas **80**, 330-357, 1961, No. 4).

- 2997:** T. Kralt, W. J. Asma and H. D. Moed: Reserpine analogues, III. Alkoxy- β -phenylethylamine derivatives (Rec. Trav. chim. Pays-Bas **80**, 431-445, 1961, No. 5).

Reserpine, an alkaloid which lowers the blood pressure and also has a sedative effect, contains three chemical groups in its molecule which together determine its important pharmacological properties: the β -indolyethylamine group, the tertiary nitro-

gen atom and the alcohol group esterified by trimethoxybenzoic acid. In an attempt to discover the relationship between the chemical structure of this compound and its physiological effects, a number of reserpine analogues containing one or more of the above-mentioned groups have been synthesized to see whether they also have the same physiological effects. These publications describe the synthesis and chemical investigation of the analogues. The pharmacological investigation is published elsewhere.

- 2998:** H. Koelmans and H. G. Grimmeiss: The incorporation of Cu in CdIn_2S_4 photoconductors (Physica **27**, 606-608, 1961, No. 6).

The authors describe the further progress of their investigations of the photoconductive properties of CdIn_2S_4 (see also **2819**).

This compound is made photoconductive by activation with Cu. This is done by adding Cu_2S and In_2S_3 (in the ratio 1:5) to the CdIn_2S_4 . The amounts added are so small that normal X-ray methods can give no information about the mode of incorporation of the Cu in the spinel lattice of the CdIn_2S_4 . In order to investigate this point, the Cu_2S and In_2S_3 were added in the same ratio but in much larger concentrations. Investigation of the resulting compounds has shown that the spinel structure is maintained, while some of the Cd atoms are replaced by Cu, and an equal number by In. The compound CuIn_5S_8 , formed by replacing all the Cd atoms, still has the spinel structure.

This same method of "exaggerating the activation" has been used to see how Cu is incorporated when only Cu_2S is added. It has been found that the same mechanism is followed, accompanied by the formation of a second phase.

As the result of another stage of this investigation, monocrystals of activated CdIn_2S_4 have been prepared. The absorption coefficient of such monocrystals has been determined as a function of the photon energy.

- 2999:** J. B. de Boer: The application of sodium lamps to public lighting (Illum. Engng. **56**, 293-312, 1961, No. 4).

Sodium lamps have a number of advantages for the lighting of traffic routes. Perceptibility is increased in sodium light, and the glare nuisance is reduced compared to that of other light sources. This is shown from the results of a number of investigations. The use of sodium lamps also means lower capital costs, operating costs and maintenance costs in many cases. These economic advantages are now reinforced by the fact that the luminous efficiency of sodium lamps has been increased to 112 lumens per watt. The characteristic colour of sodium light can be made use of for leading traffic through or round a town. See also Philips tech. Rev. 23, 258-272, 1961/62 (No. 8/9).

This publication is the text of a lecture which the author delivered in America, where sodium lamps are much less used than in Europe. A discussion with a number of American lighting experts following on this lecture is also reproduced.

3000: C. Z. van Doorn: Method for heating alkali halides and other solids in vapors of controlled pressure (Rev. sci. Instr. 32, 755-756, 1961, No. 6).

If KCl is heated in potassium vapour of controlled pressure, slight deviations from the stoichiometric composition are produced, leading to characteristic colours. The study of the nature of this effect is of considerable theoretical interest. The pressure of the unsaturated metal vapour is controlled by making use of a principle which somewhat resembles that of a mercury-diffusion pump. The apparatus described can be used for controlling the vapour pressure of all substances which have a low vapour pressure near the triple point, e.g. Na, K, Rb, Cs, Hg, S, Se, Te, P (yellow). It has already also been used for heating CdS crystals in an atmosphere of Cd or S.

3001: N. W. H. Addink and L. J. P. Frank: Het zink- (koolzuuranhydrase-) gehalte van bloed van lijders aan neoplastische ziekten (10de Jaarboek van Kankeronderzoek en Kankerbestrijding in Nederland, 1960, pp. 11-21). (The zinc (carbonic anhydrase) content of the blood of sufferers from neoplastic diseases; in Dutch.)

In a comparative analytical (spectrochemical) investigation of the mineral components of the blood of healthy persons and of cancer patients, the greatest differences were found with zinc. The zinc content of healthy blood increases with age, but in patients with progressive neoplastic diseases the zinc level falls continuously, with a few excep-

tions (bone and lung cancer, leukemia). The course of the illness as diagnosed by the doctor is in agreement with conclusions drawn from the zinc level of the blood in 80-90% of the cases. Zinc occurs in the carbonic anhydrase present in the erythrocytes. Some properties of this metallo-protein are discussed. The authors of this article hope that this protein will be tested for anti-carcinogenic action, although it is not yet certain whether the observed deficiency is to be regarded as a primary cause of neoplastic diseases or as an accompanying effect. See also 2726 and 2832.

3002: H. J. Heijn and J. C. Selman: The Philips computer PASCAL (IRE Trans. on electronic computers EC-10, 175-183, 1961, No. 2).

See Philips tech. Rev. 23, 1-18, 1961/62 (No. 1).

3003: H. Zijlstra: Device for the rapid measurement of magnetic anisotropy at elevated temperatures (Rev. sci. Instr. 32, 634-638, 1961, No. 6).

In this method for the rapid measurement of magnetic anisotropy, the magnetic sample is suspended in a self-exciting torsion pendulum which hangs in a magnetic field. A relationship can be derived between the resonance frequency of the pendulum and the anisotropy energy (which is a measure of the magnetic anisotropy). The resonance frequency lies between 10 and 100 c/s, so that the time required for these measurements is less than in other methods. It is thus possible to follow a rapidly varying anisotropy in this way, e.g. during the heat treatment of the sample in a magnetic field. This method has been used to measure the anisotropy of a monocrystal of "Ticonal G" magnet steel during cooling from 900 to 400 °C in a magnetic field.

3004: J. J. Engelsman and A. M. J. M. Claassens: Anodic stripping using a rotating mercury drop (Nature 191, 240-241, 1961, No. 4785).

The polarographic determination of extremely small quantities of metals can be carried out accurately by using as electrode a platinum wire rotating about its axis, with a tiny drop of mercury (weight about 0.5 mg) on its tip. This refinement of an established method has already given good results with zinc, cadmium, lead and copper.

3005: G. H. Jonker and S. van Houten: Semiconducting properties of transition metal oxides (Halbleiterprobleme Vol. 6, editor F. Sauter, pp. 118-151, Vieweg, Brunswick, Germany, 1961).

A review article dealing with the semiconducting properties of oxides of the iron group of transition metals. These properties are determined by the electrons from the d shell, which also cause the much more deeply studied magnetic properties of these compounds. The authors give a simple phenomenological description of the semiconducting properties, which can be used as a basis for the study of the individual oxides. To avoid complications, the discussion of the individual compounds is restricted to those with a simple crystal structure, e.g. the rocksalt, corundum, rutile, spinel or perovskite structure.

3006: J. L. Ouweltjes: The specification of colour rendering properties of fluorescent lamps (*Die Farbe* **9**, 207-246, 1960, No. 4/6).

Two methods have been proposed for investigating the colour-rendering properties of a light source. In the first method, a number of standard colours are compared in the light under investigation and in the light of a standard light source. In the second method, the spectrum is divided into a number of bands, and the energy radiated by the light source in question in each band is compared with the energy radiated by the standard light source in the same band. The first part of this article describes briefly how the spectral energy distribution of a "TL" lamp can be calculated from the energy distribution and quantum yield of the phosphors used. In the second part it is shown that eight test colours are more than enough to determine the colour rendering of a "TL" lamp. The third part contains an attempt to compare the results of the two above-mentioned methods, which was however unsuccessful. For the moment, the method using test colours seems to be the most reliable. See also *Philips tech. Rev.* **13**, 249-260, 1956/57.

R 421: S. A. Wytzes: Theoretical considerations on the collimators of an X-ray spectrograph (*Philips Res. Repts* **16**, 201-224, 1961, No. 3).

An investigation of the way in which the profile and the intensity of the lines measured with an X-ray spectrograph depend upon the dimensions of the collimators and the properties of the reflecting crystal. Simple graphical methods are derived for determining the line profile for two cases: the case where the angle between the normal to the lattice planes and the bisector between the incident and the dif-

fracted rays exhibits a certain spread, but the angle of diffraction is exactly 2θ , and the case where both angles exhibit a spread.

R 422: A. Baelde: The influence of non-uniform base width on the noise of transistors (*Philips Res. Repts* **16**, 225-236, 1961, No. 3).

The experimentally determined and calculated values of the noise resistance of a transistor do not agree with each other. It is shown that this must be ascribed to variations in the thickness of the base. The effective base resistance, which mainly determines the noise resistance, is then larger than the value derived from admittance measurements.

R 423: P. A. H. Hart and F. L. van der Vinne: The measurement of the noise quantities of the EC 56 at 1400 Mc/s (*Philips Res. Repts* **16**, 237-244, 1961, No. 3).

Noise measurements have been carried out on the microwave tube EC 56 at a frequency of 4000 Mc/s (see **R 393**). The authors have carried out similar measurements at 1400 Mc/s, at which frequency the experimental technique is much more complicated. The four measured characteristic noise quantities of the above-mentioned tube can be used to calculate the minimum noise factor of a travelling-wave tube which uses the same electron beam. The results of the measurements are in agreement with the theory given by Vlaardingerbroek (**R 393**).

R 424: J. H. N. van Vucht: Kinetic study of the reaction of Th_2Al with H_2 (*Philips Res. Repts* **16**, 245-265, 1961, No. 3).

The absorption of hydrogen by a getter powder occurs by means of an autocatalytic reaction. This reaction is initially slow, since the grains of metal are completely covered with a film of impurities, e.g. oxide, through which hydrogen can only diffuse slowly. When, however, the metal has taken up enough hydrogen to split the film up by expansion effects, the absorption can take place much more quickly.

R 425: M. Koedam: Cathode sputtering by rare-gas ions of low energy (*Philips Res. Repts* **16**, 266-300, 1961, No. 3).

Second part of a thesis, which has already been summarized in full under **R 416**.

Philips Technical Review

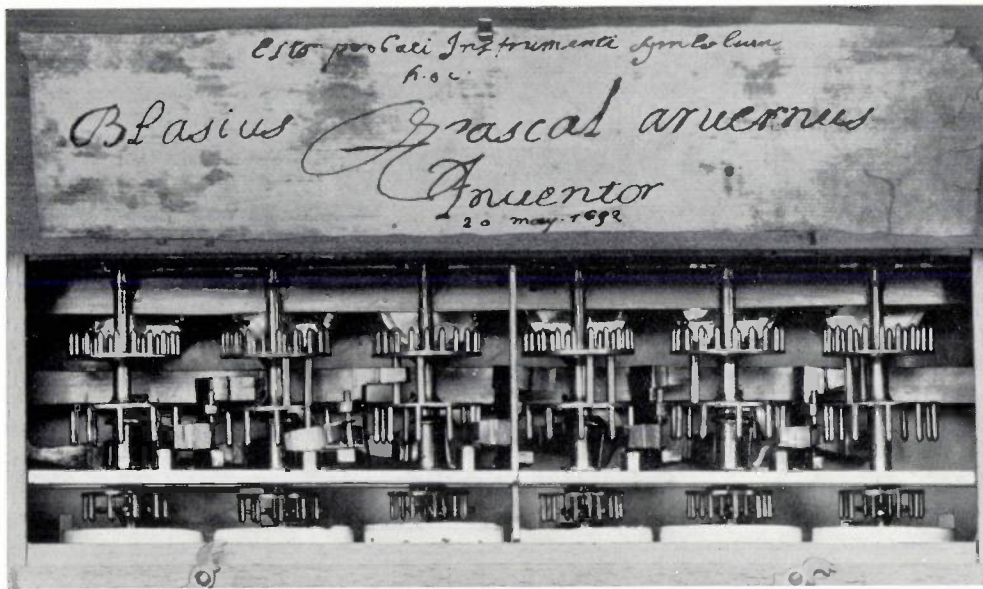
DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES



In this issue of our journal we present a number of articles which, though drawn from different departments within Philips, have a single subject in common: they illustrate the function which the electronic computer fulfills as a tool of science, technology and economics, for research and routine work. It need hardly be said that this function has already assumed considerable importance; and it is continually gaining further ground at such a speed that a selection of articles as brought together here may well have lost its point in a few years time, when the electronic computer will perhaps have become just as commonplace a tool as the slide rule and the desk calculator are today.

Since the creation of Philips Computing Centre — one department of which is shown in the above photograph — hundreds of problems have been tackled there

with electronic computers. Our aim in this number is to give some idea of the diversity of the problems concerned: the selection comprises the development of a clover-leaf cyclotron, the processing of orders for corrugated cardboard, the calculation of potential fields and electron trajectories, the vacuum deposition of resistors, and the design of magnets for loudspeakers. Moreover, the problems chosen demonstrate the variety of detailed mathematical tasks entrusted to the computer. Three short articles deal with less commonly considered aspects of the electronic computer: one gives a glimpse of the history of automatic computation; one discusses the solution of a puzzle — a light-hearted but instructive example of a non-arithmetical problem; and the last brings the machine unexpectedly close to the reader with the aid of a gramophone record.



British Crown Copyright, Science Museum, London

THE CALCULATING MACHINE OF BLAISE PASCAL

681.14(091)

The digital electronic computer PASCAL in Philips Computing Centre has been given this name in honour of the French mathematician and philosopher Blaise Pascal¹⁾, who in 1642, at the age of eighteen, designed a calculating machine at Rouen. His object with this machine, which became known as the Pascaline, was to ease the burden on his father who, as a tax official, had a great deal of figure work to do. Although in the later years of his short life (Pascal died in 1662, almost exactly three hundred years ago) he was mainly concerned with other matters, he nevertheless had more than 50 models made of his machine²⁾, each being an improvement on the preceding ones. In 1645 he presented one to Chancellor Pierre Séguier, through whose good offices he obtained in 1649 a royal privilege on his invention; in 1647 he showed one to Descartes; in 1652 he finally arrived at a form that satisfied him; he sent one machine to Queen Christina of Sweden, and another he demonstrated personally to a distinguished gathering in Paris — successfully, to judge from the poetic effusion of a contemporary³⁾.

One model of the Pascaline dating from 1652 has been well preserved and is to be seen at the Conservatoire des Arts et Métiers, Paris. The title photograph is of a replica in the Science Museum in London. The Paris Conservatoire has three other machines; all four bear the arms of the Pascal family (see fig. 1).

Various mechanical aids to arithmetical work, such as the time-honoured abacus and the graduated rods invented by Napier in 1617 (Napier's "bones"), were already in use before Pascal's machine. But Pascal went an essential step further, in that his machine contained a discontinuous mech-

³⁾ Muse historique, Loret, of 14th April, 1652 (see the book quoted in footnote ²⁾, page 57).

"Je me rencontrai l'autre jour
Dedans le petit Luxembourg,
Au quel beau lieu que Dieu bénie
Se trouva grande compagnie,
Tant duchesses que cordons bleus,
Pour voir les effets merveilleux
D'un ouvrage d'arithmétique,
Autrement de mathématique,
Où, par un talent sans égal
Un auteur qu'on nomme Pascal,
Fit voir une spéculative
Si claire et si persuasive,
Touchant le calcul et le jet,
Qu'on admira le grand projet.
Il fit encor sur les fontaines
Des démonstrations si pleines
D'esprit et de subtilité,
Que l'on vit bien, en vérité,
Qu'un très beau génie il possède
Et qu'on le traita d'Archimède."

¹⁾ W. Nijenhuis, The PASCAL, a fast digital electronic computer for the Philips Computing Centre, Philips tech. Rev. 23, 1-18, 1961/62 (No. 1). — It should be mentioned that according to some, the name is an acronym derived from Philips Automatic Sequence Calculator.

²⁾ P. Humbert, L'oeuvre scientifique de Blaise Pascal, Albin Michel, Paris 1947, p. 56.

anism — the “sautoir” — for automatically carrying over tens, etc., in adding operations (*fig. 2*). This is the basis of all digital techniques and the logical consequence of the digital or positional system of writing numbers.

Pascal's contemporaries were aware of the potentialities of his idea. Speaking of the “machine arithmétique” his sister Gilberte expressed it thus ⁴): “This accomplishment has been regarded as something new in nature, to have reduced to a machine a science that belongs entirely to the mind, and to have found the means of performing all operations with complete certainty, without the need for reasoning”. People felt a kind of uneasiness or amazement about the Pascaline, such as many of

us feel today about automation, which seems capable, through the use of electronic computers, of taking over our whole function of logical thought. Gilberte Pascal, incidentally, went on to say: “This effort tired him very much, not because of the brainwork or of the mechanism, which he found without any trouble, but because of the difficulty of making the workers understand all these things”. Indeed one may assume that the realization of Pascal's invention was seriously hampered by the fact that the schooling and probably the equipment of the instrument makers at that time were inadequate for making such intricate devices with the necessary precision. No model of the Pascaline seems to have worked for long without faults, and the manufacture of calculating machines on a commercial scale had to await the perfecting of the mechanism and a general improvement in the standard of engineering.

⁴) Pascal, *Pensées et Opuscules*, éd. par L. Brunschvicg, Hachette, Paris 1945 (p. 10).

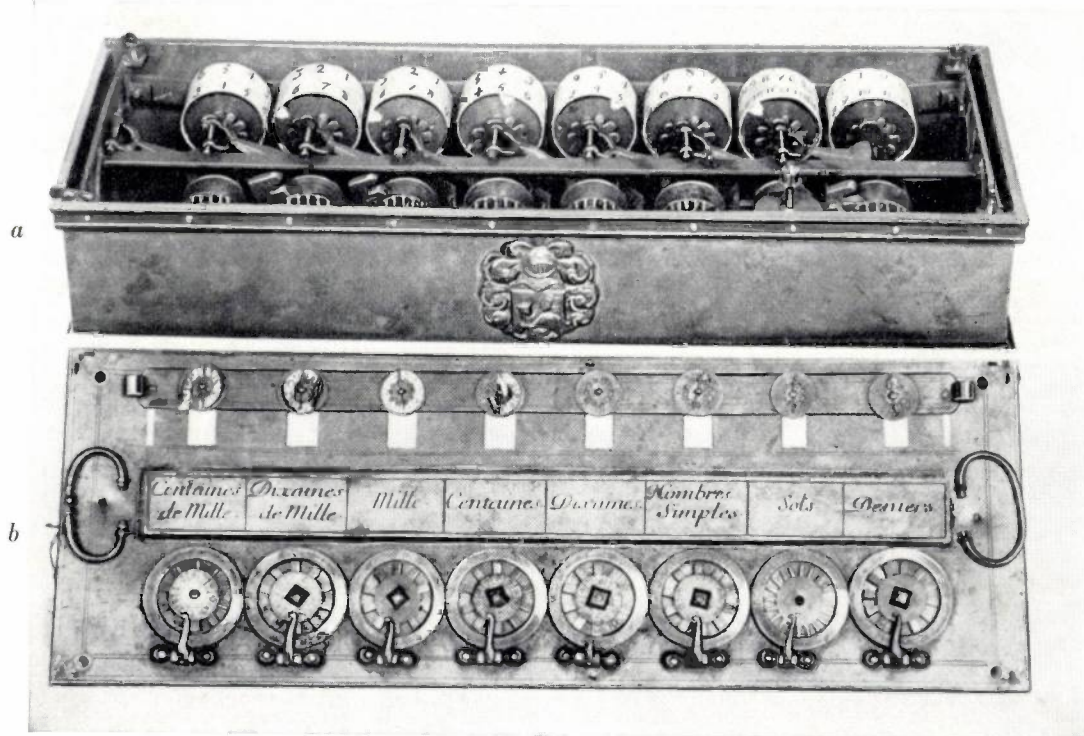


Photo Conservatoire des Arts et Métiers, Paris

Fig. 1. One of the four models of Pascal's calculating machine preserved in the Conservatoire des Arts et Métiers, Paris. This model, like various others, was designed for the addition of money up to 1 million livres. For this purpose six decimal places are available, plus a seventh place with 20 units for the sous and an eighth with 12 units for the deniers. (The same division, into pounds, shillings and pence, has persisted in Great Britain up to the present day.) The divisions can clearly be seen, on the removed cover (b), on the eight selector discs which serve for setting the digits to an amount to be added.

It is worth noting that even the earliest calculating machines demonstrate in this way that digital computing is not tied to the decimal system. The binary system ¹) employed in electronic computers is just another variant.

The machine is operated by inserting a peg in each of the eight selector dials and turning the dial through successive stops. The number thus set, and the result of the addition when the next number is set, appear in the sight holes in the cover, below which rotate the figure wheels seen in (a). The carry-over of the tens (and the twelves and twenties) is automatic.

Subtraction is done by pushing down the bar above the sight holes which carries the eight “register” wheels, thus exposing the top halves of the sight holes, in which there now appear the figures in the reverse sequence (the complements respectively of 10, 12 or 20). A number is set and subtracted by turning the selector dial in the same direction as for addition. Pascal devised this method because his automatic carrying device, the “sautoir”, worked only in one direction (see *fig. 2*).

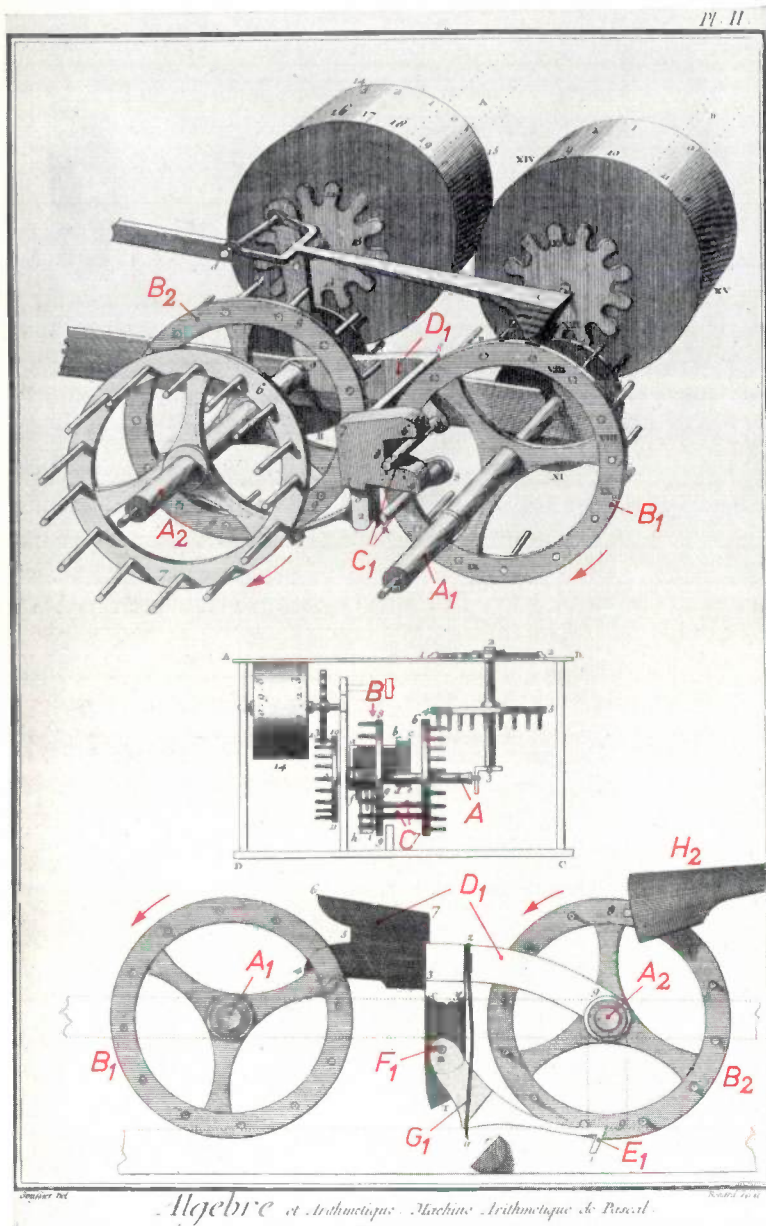


Photo Science Museum, London

Fig. 2. Mechanism for automatic carry-over in Pascal's calculating machine. The drawing is reproduced from the Diderot and d'Alembert Encyclopedia, Paris 1752-1777. We have added the letters printed in red, for denoting components. The mechanism (the "sautoir") works roughly as follows. The drawing in the middle shows all components for one digit; the selector dial is seen on the cover, at the right, and the figure wheel or drum is on the left, under the cover. The movement of the dial is transmitted by two pairs of gear wheels (pin wheels) via the shaft A to the figure wheel. The middle pin wheel B on this shaft serves for carrying over the tens. In the top drawing can be seen the pin wheel B_1 for one digit and the pin wheel B_2 for the next higher digit. Towards the end of a full revolution of B_1 , the two pins C_1 engage the two teeth of the doubly-bent lever D_1 turning about the spindle A_2 and lift the lever. When B_1 has completed a full revolution (i.e. completing a ten) the pins C_1 release the teeth of D_1 , the lever drops and a pawl E_1 on the lever pushes the pin-wheel B_2 one step further. This is made clear by the bottom drawing, which shows the components from the other side. The arm of pawl E_1 hinges on the spindle F_1 and is lifted by the leaf spring G_1 , so that when D_1 drops, the pawl can engage a pin on B_2 whereas during the lifting of D_1 (and also when B_2 is turned independently) the pins are free to slide off along the arm of E_1 . A catch H_2 prevents B_2 from being dragged in the wrong direction when D_1 is lifted.

Some stages in this further development of calculating machines may usefully be mentioned. In Britain, in 1666, Morland built a calculating machine (two examples of which are preserved in London) which, compared with the Pascaline, represented a step backwards. The machine worked on the same principle — the adding of figures by successive rotations of a kind of selector dial through discrete angles — but there was no automatic carrying device. In 1672 Leibniz began work in Hanover, and later in Paris, on a calculating machine based on a new idea, the "stepped gear", which could also perform multiplication and division. He worked on this for many years helped by various instrument makers. It was not until 1694 that his first machine was completed, and even then seems never to have been reliable in operation. This machine is still at Hanover, and a replica is in the Deutsches Museum in Munich. In Padua in 1709 Poleni utilized the same principle as Leibniz and conceived a novel, highly effective mechanism for the automatic carry-over — virtually the same construction is still used in mechanical counting mechanisms today, such as mileometers, gas and electricity meters, etc. A wooden model of his machine so disappointed Poleni, however, that he destroyed it. A machine built in 1727 by Antonius Braun fared better. Embodying a device similar to that used by Leibniz and Poleni⁵⁾, this machine (fig. 3) was put together with great care and precision — Braun was apparently both inventor and craftsman — and gives the impression of having worked well although it does not appear to have been easy to operate. After numerous other intermediate stages the first calculating machine to be manufactured on a commercial scale appeared in 1820; this machine was designed by Charles Xavier Thomas of Colmar and remained on the market, with few modifications, for almost 100 years.

⁵⁾ J. Nagler, Beschreibung der Rechenmaschine des Antonius Braun, Blätter für Technikgeschichte, No. 22, pp. 81-87, Springer, Vienna 1960.

Fig. 3. Calculating machine made by Antonius Braun in 1727. It was intended as an aid to surveying work, and could add, subtract, multiply and divide. Whether it worked satisfactorily is not known. The photo shows the machine without the cylindrical side panel that protects the mechanism from dust. The top plate, with the setting levers and figure dials, bears a Latin inscription in which the maker ("Opticus Et Mathematicus") humbly dedicates the instrument to the Emperor Charles VI. The instrument can be seen in the Technisches Museum für Industrie und Gewerbe at Vienna⁵⁾.

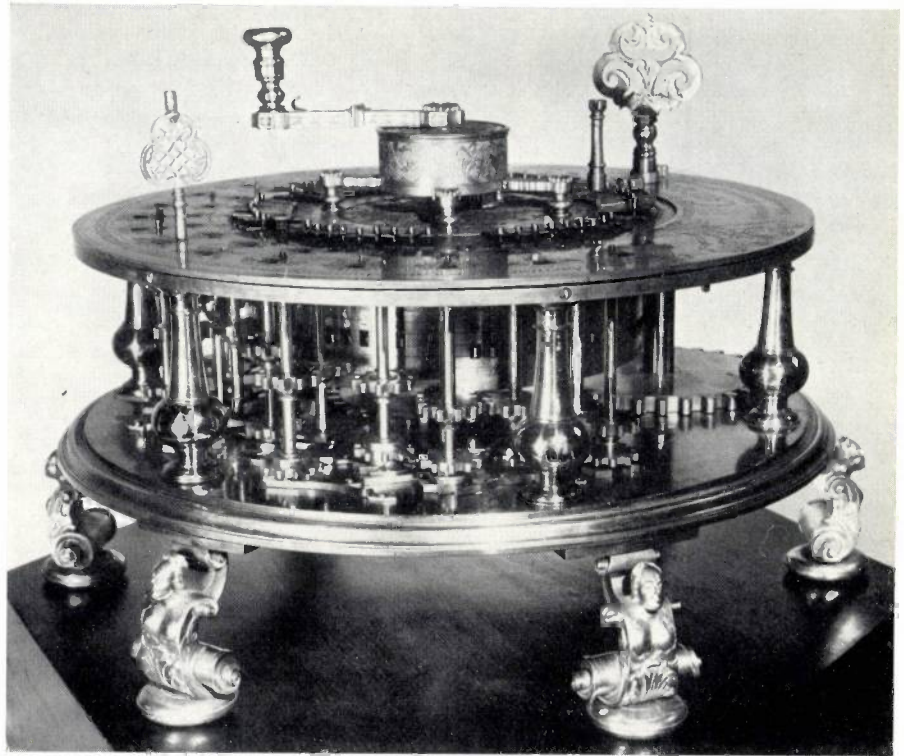


Photo Technisches Museum, Vienna

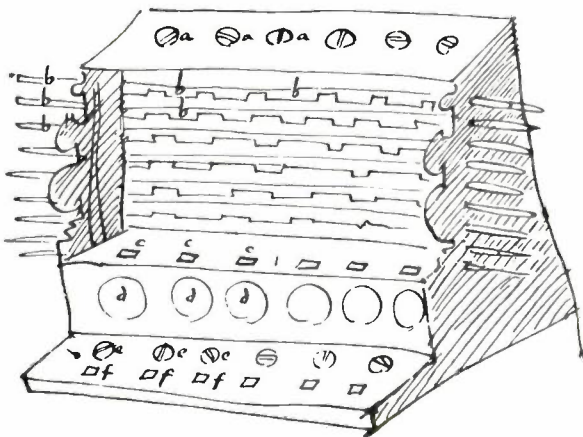


Fig. 4. Sketch of W. Schickard's calculating machine, taken from his letter to Kepler of 25.2.1624. The text referring to the machine reads (translated from the Latin): "I shall outline the arithmetic apparatus in more detail another time; being in haste the following must suffice: *aaa* are the top ends of vertical cylinders, on which are written the multiplications of the figures, and those [multiplications] which are necessary can be seen through the sliding windows *bbb*. Fixed on the inside to *ddd* are gear wheels with 10 teeth, that mesh with one another such that if any wheel on the right turns round ten times, the wheel to the left of it turns round once; or if the first-mentioned wheel makes a hundred turns, the third wheel turns once, etc. To wit, [they all do this] in the same direction, for which purpose an identical intermediate wheel *h* was necessary. Any given intermediate wheel sets all to the left of it in motion, in the requisite proportion; but none to the right of it, which called for special measures. The number on these wheels is visible through the holes *ccc* in the centre ledge. Finally, the letters *e* on the bottom ledge denote rotary knobs and the letters *f* are again holes through which figures used when working can be seen."

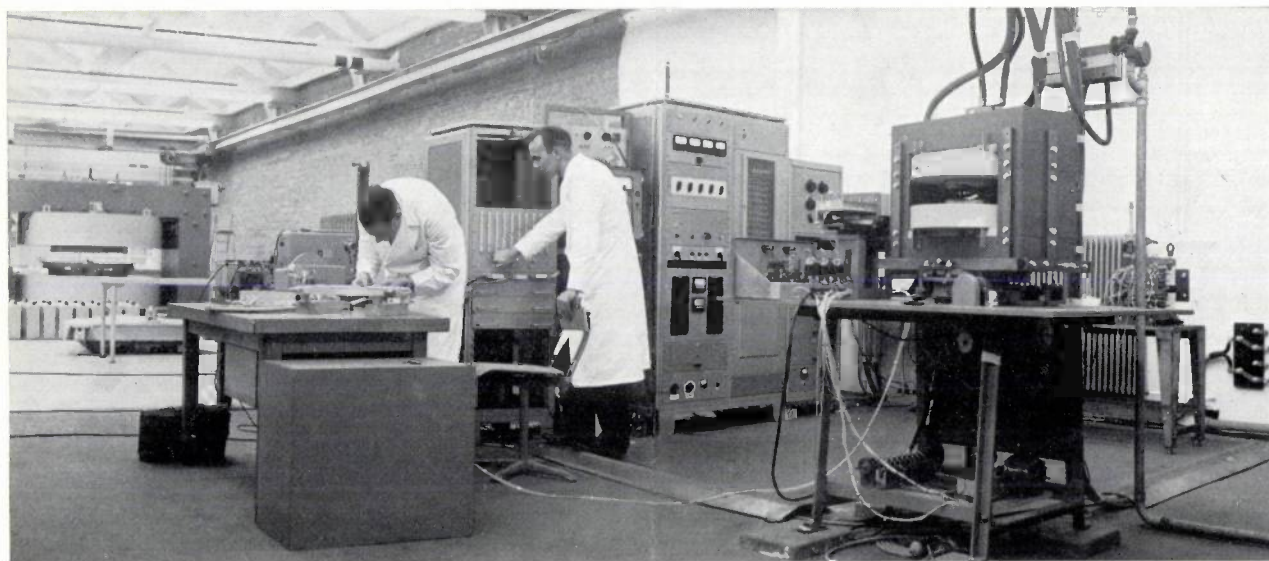


This account of the earliest history of the calculating machine cannot be closed without mentioning that a few years ago Hammer and v. Freytag-Löringhoff⁶⁾ discovered a predecessor of the Pascaline: the hebraist, astronomer and mathematician Wilhelm Schickard at Tübingen had already constructed in 1623, i.e. 20 years before Pascal, a calculating machine with automatic carry-over of tens which could apparently add and subtract (even alternately, which was not possible with the Pascaline) and which, moreover, had a device that facilitated multiplication and division. In two letters to Kepler, dated 20.9.1623 and 25.2.1624, Schickard reports and describes his invention (fig. 4), and on the basis of this description and a sketch found in the papers Schickard left behind, a reconstruction has been made of the machine. Unlike Pascal however, Schickard evidently did not arouse the interest of his contemporaries in his machine. A second, improved model which he had designed was destroyed by fire before completion, and probably also because of the war at the time and his death soon after (he and his whole family died of the plague in 1635), his invention was immediately forgotten.

S. GRADSTEIN *).

⁶⁾ B. v. Freytag-Löringhoff, Wiederddeckung und Rekonstruktion der ältesten neuzeitlichen Rechenmaschine, VDI-Nachrichten 14, No. 39, 21st December 1960 (p. 4).

*) Philips Research Laboratories, Eindhoven.



INVESTIGATION OF THE MAGNETIC FIELD OF AN ISOCHRONOUS CYCLOTRON

by N. F. VERSTER *) and H. L. HAGEDOORN *).

621.317.42:621.384.611.2

The invention of the cyclotron by Lawrence was based on the fact that a homogeneous magnetic field forces ions of a given mass to describe circular orbits with a *constant* period, independent of their velocity. The ions can therefore be given greater and greater velocities with the aid of an HF alternating electric field (applied across the slit between the two "Dees") (fig. 1). The factor limiting the velocity which can be attained in this way has already been mentioned several times in this journal¹⁾: as the energy of the ions increases, the relativistic increase of their mass begins to be appreciable; for protons, this increase amounts to 1% per 10 MeV. In order to keep the period of revolution of the ions constant despite this fact, the magnetic field would have to increase slightly with increasing distance from the centre. It is found in practice however that the magnetic field must actually *decrease* slightly with increasing distance from the centre. This is because the ions can oscillate vertically and horizontally about their ideal (horizontal) orbit, and the magnetic field should be slightly barrel-shaped in order to stabilize the vertical oscillation, which implies that the field should decrease with increasing distance from the centre. Since the period of revolution of the ions is

thus not constant, the particles do not always pass the slit at the moment when the electric field has its maximum accelerating value (phase 0°), but are gradually retarded. When the phase lag reaches 90° , the ions are no longer accelerated at all. This effect can be reduced to a certain extent by using a high

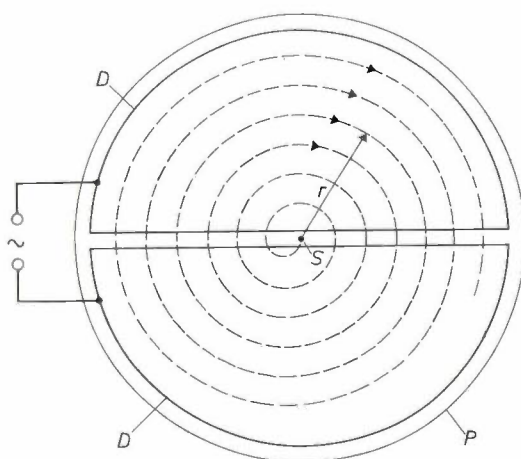


Fig. 1. Principle of the acceleration of ions in a cyclotron. The ions supplied by the source *S* are accelerated in the electric field between the two D-shaped electrodes *D* (the "dees"). They describe circular orbits under the influence of a magnetic field (perpendicular to the plane of the paper), so that they keep on passing the slit between the dees. The voltage applied between the dees is AC of such a frequency that the field between the dees just reverses in the time which it takes the ions to describe half a circle. *P* is the edge of one of the poles of the magnet.

*) Philips Research Laboratories, Eindhoven.

¹⁾ For example, W. de Groot, Cyclotron and synchrocyclotron, Philips tech. Rev. 12, 65-72, 1950/51.

dee voltage (100 kV or more), since the particles then reach a given energy in fewer revolutions; but even so, one cannot get beyond 10 or 20 MeV with protons in this way.

In the *synchrocyclotron*, this problem is solved by varying the frequency of the accelerating voltage periodically ¹⁾²⁾. The ion current is then no longer continuous, but consists of a series of short pulses. The final energy can be very high in machines of this type: the biggest synchrocyclotron in the world produces protons of not less than 730 MeV.

As long ago as 1938, Thomas ³⁾ had shown theoretically that it would be possible to get both the desired stabilization in the vertical direction and a constant period of revolution by giving up the idea that the magnetic field had to possess rotational symmetry, i.e. that the ideal orbits of the ions had to be circles. It appears that a strong azimuthal variation in the magnetic induction can contribute so much to the vertical stabilization that the average induction over one complete orbit (which determines the period in this case too) can even be allowed to *increase* with increasing distance from the centre. However, the realization of the magnetic field needed for this purpose seemed so difficult that for some time no one attempted to put this principle into practice. The interest in this idea, however, was renewed after the development of the principle of "strong focusing", which somewhat resembles Thomas's idea; this principle was originally used in proton synchrotrons ⁴⁾ and was further developed by the Midwestern University Research Association ⁵⁾. In 1958 the first cyclotron to operate on Thomas's principle was completed in Delft, Netherlands ⁶⁾. Since then, another five of these *isochronous* cyclotrons have been completed. The name "isochronous" refers to the constant period of revolution; they are also sometimes called AVF cyclotrons (from Azimuthally Varying Field).

A cyclotron of this type is under construction in the Philips Laboratories in Eindhoven (Geldrop section), under the general direction of A.C. van Dorsten. Like various other isochronous cyclotrons which have already been completed, this machine will be able to accelerate *different* sorts of ions to a *continuously variable* final energy. Protons will be able to be accelerated to a maximum of 25 MeV ⁷⁾.

In the development of this cyclotron, in particular for the accurate realization of the desired magnetic field, considerable use was made of the electronic computer PASCAL. In this article we shall describe how this was done. We shall begin by discussing in somewhat more detail the theory of the motion of the ions in the complicated magnetic field and the consequences of the *variable* energy. We shall then describe how the magnetic field is measured and how the data thus obtained are processed with the aid of the PASCAL, and finally how the behaviour of the particles in the magnetic field was calculated and how this led to the accurate determination of the field ⁸⁾.

The motion of the ions

We shall introduce cylindrical coordinates r, θ, z , with the z axis along the axis of the poles of the magnet and the plane $z = 0$ in the median plane between the two poles.

In the *rotationally symmetrical* field of the classical cyclotron, with an induction $B(r)$ in the median plane, ions of mass m , charge e and velocity v describe a circular orbit in the median plane, whose radius r_0 is given by

$$r_0 B(r_0) = mv/e. \dots \dots (1)$$

The angular velocity is

$$\omega = v/r_0 = eB(r_0)/m, \dots \dots (2)$$

and it will be seen that this angular velocity (and thus the required frequency f of the accelerating field) is indeed constant as long as the relativistic increase of m is neglected. If we take this increase into account, it is found that for the angular frequency to remain constant the following relationship must be satisfied:

$$B(r) = \frac{m\omega}{e} = \frac{m_0\omega}{e\sqrt{1 - (v/c)^2}} \approx \frac{m_0\omega}{e} [1 + \frac{1}{2}r^2(\omega/c)^2]. (3)$$

²⁾ Some of the synchrocyclotrons which have been partly or wholly built by Philips are described in: F. A. Heyn, Philips tech. Rev. **12**, 241, 247 and 349, 1950/51 and F. A. Heyn and J. J. Burgerjon, *ibid.* **14**, 263, 1952/53 (these articles deal with the cyclotron in Amsterdam); W. Gentner *et al.*, *ibid.* **22**, 141, 1960/61 (three articles concerning the CERN 600-MeV synchrocyclotron for which Philips made the HF installation); G. T. de Kruiff and N. F. Verster, *ibid.* **23**, 381, 1961/62 (No. 12) (concerning the 160-MeV synchrocyclotron at Orsay).
³⁾ L. H. Thomas, Phys. Rev. **54**, 580, 1938.
⁴⁾ E. D. Courant, M. S. Livingston and H. S. Snyder, Phys. Rev. **88**, 1190, 1952 and **91**, 202, 1953.
⁵⁾ K. R. Symon, D. W. Kerst, L. W. Jones, L. J. Laslett and K. M. Terwilliger, Phys. Rev. **103**, 1837, 1959.
⁶⁾ F. A. Heyn and Khoe Kong Tat, Rev. sci. Instr. **29**, 662, 1958. The same authors give a more detailed description in: Sector-focused cyclotrons (Proc. Conf. Sea Island, Georgia, Feb. 1959), pp. 29-39. See also Khoe Kong Tat, The isochronous cyclotron, thesis Delft, 1960.

⁷⁾ Further details may be found in N. F. Verster *et al.*, Some design features of the Philips AVF prototype, Nucl. Instr. Meth. **18/19**, 88-92, 1962.
⁸⁾ See also N. F. Verster and H. L. Hagedoorn, Computer programs for an AVF cyclotron, Nucl. Instr. Meth. **18/19**, 327-335, 1962.

This is known as the isochronous field form, and we shall denote it by $B_{is}(r)$ from now on. The frequency f_z of the vertical oscillations of the ions about the circle of radius r_0 (the equilibrium orbit) is given by ⁹⁾

$$(f_z/f)^2 = -k, \quad \dots \dots (4)$$

where k is the "field index", defined by the equation

$$k = \frac{r}{B} \frac{dB}{dr} = \frac{d \log B}{d \log r}.$$

For stability in the vertical direction, it is necessary that f_z should be real, i.e. that k should be negative. This makes it impossible to satisfy equation (3), which is thus the reason why the rotationally symmetrical magnetic field is useless when the relativistic increase of mass becomes appreciable.

The frequency f_r of the radial oscillations is given by:

$$(f_r/f)^2 = 1 + k. \quad \dots \dots (5)$$

Radial stability is thus achieved as long as $k > -1$.

If we now abandon the idea of a rotationally symmetrical magnetic field, as Thomas suggested, then the form of the magnetic field in the median plane $z = 0$ can be described quite generally by means of a Fourier series:

$$B(r, \theta) = \bar{B}(r) \left[1 + \sum_{n=1}^{\infty} A_n(r) \cos n\theta + \sum_{n=1}^{\infty} B_n(r) \sin n\theta \right],$$

or

$$B(r, \theta) = \bar{B}(r) \left[1 + \sum_{n=1}^{\infty} C_n(r) \cos n\{ \theta - \varphi_n(r) \} \right]. \quad (6)$$

We cannot discuss the theory of the motion of the ions in such a field in great detail here. In the first place, this theory states that radial stability is only possible if the first two harmonics are absent or at least very weak ($C_1 \ll 10^{-3}$ and $C_2 \ll 10^{-2}$). In order to satisfy this condition, the magnetic field is given three-fold or four-fold symmetry about the z axis. Equation (6) then only contains terms in which n is a multiple of 3 (thus 3, 6, 9, ...) or 4 (thus 4, 8, 12, ...). Three-fold symmetry was chosen for our cyclotron, and we shall therefore restrict ourselves to this case, where $n = 3$ is the lowest occurring term (and the most important).

In order to produce a magnetic field $B(r, \theta)$ of three-fold symmetry, three steel "hills" each about 60° wide (and thus with "valleys" in between) are built on to each pole of the magnet, which thus looks something like a clover leaf (see fig. 2a); such cyclotrons are therefore also sometimes called *clover-leaf*

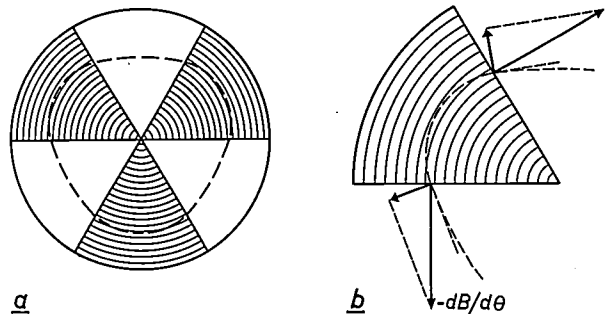


Fig. 2. a) Sketch of a magnetic field with 3-fold symmetry. The field is stronger in the shaded areas; the orbits of the ions are therefore more strongly curved here, so that the orbits assume a somewhat triangular shape. b) Since the ions do not cross the boundaries between the various sectors at right angles, at these boundaries they encounter a field which decreases with increasing distance from the centre. The larger arrows represent the (negative) field gradients, and the smaller ones the components of these at right angles to the ion orbit shown here.

cyclotrons. Instead of the circular orbit of the classical cyclotron (eq. 1), an ion now describes a somewhat *triangular* orbit, since the curvature of the orbit is greater near the hills than near the valleys. It will be seen that the triangular orbit does not cut the radial boundaries between the hills and the valleys at right angles. In these six boundary regions, the magnetic field at right angles to the orbit always decreases with increasing distance from the centre — which is exactly the situation required for vertical stability, as we have seen for the case of the classical cyclotron. It is thus plausible that the desired stability can be obtained in an AVF cyclotron, despite the fact that the average field actually increases with increasing distance from the centre (fig. 2b).

The stabilizing effect is increased even further if the boundaries between the hills and the valleys are made slightly *spiral* instead of the simple radial shape (fig. 3a); this makes the angle between the orbit and the edge of the hill alternately large and small. It is this form which has been chosen for our cyclotron. The phase angle φ of each Fourier term will then indeed be a function of r , as indicated in equation (6) (which is not the case with the simple radial hills); the relationship is given by $r d\varphi/dr = \tan \gamma$, where γ is the spiral angle as defined in fig. 3b.

The equations of motion of an ion in a complicated

⁹⁾ This formula follows directly from the equations of motion of the ion. See e.g. ¹⁾ or p. 386 of the last article quoted in ²⁾. We may mention that in these articles the definition of the field index n differs from that used here: $n = -k$.

field, as described by equation (6), cannot be solved exactly. This is one of the greatest problems encountered by the designer of an isochronous cyclotron. Here too, the electronic computer provides the solution to the problem: the orbits were determined by numerical analysis, with the aid of the PASCAL. The functions $\bar{B}(r)$, $A_3(r)$, $B_3(r)$, $A_6(r)$, $B_6(r)$, . . . , which are needed for this calculation, were in their turn calculated by the PASCAL from measurements of the magnetic field (Fourier analysis). Because of the three-fold symmetry, it is only necessary to make measurements in one 120° sector, along a number of arcs of circles of different radii. (A number of check measurements were made through 360° to make sure that the amplitudes of the first and second harmonics were indeed as small as specified above.) The required isochronous field distribution $\bar{B}_{is}(r)$, the corresponding angular velocity ω , and the oscillation frequencies f_z and f_r were then determined from these calculated orbits, again with the aid of the PASCAL.

When we speak here of "the" magnetic field which is measured and for which the ion orbits are calculated, we mean in the first place an actual field obtained using clover-leaf poles. The results of the calculations are naturally just what is needed to determine the *corrections* which must be made to the field to ensure exact isochrony. We shall discuss this in more detail in the next section.

The calculations in question, which we shall also discuss further below, are rather complicated and time-consuming, even with the PASCAL. We have therefore also paid some attention to the possibility of designing a cyclotron using an *approximate* analytical solution of the equations of motion in the clover-leaf magnetic field. The simplifying assumptions which are necessary in order to obtain a solution in this way are however so drastic that it seemed doubtful whether the solution thus obtained would be of any practical use. Moreover, the expressions

for the isochronous field, the angular velocity and f_r are still very complicated¹⁰). We will give here only the most important terms in these expressions:

$$\bar{B}_{is}(r) \approx \frac{m_0 \omega}{e} \left[1 + \frac{1}{2} r^2 \left(\frac{\omega}{c} \right)^2 - \frac{1}{16} C_3^2 - \frac{1}{32} \left(r \frac{dC_3^2}{dr} \right) \right], \quad (7a)$$

$$\omega \approx \frac{e \bar{B}}{m} \left[1 + \frac{1}{16} C_3^2 + \frac{1}{32} \left(r \frac{dC_3^2}{dr} \right) \right], \quad (7b)$$

$$(f_z/f)^2 = -k + \left(\frac{1}{2} + \tan^2 \gamma \right) C_3^2. \quad \dots \quad (8)$$

We have (again with the aid of the PASCAL) computed the value of the complete expressions in all cases, using the above-mentioned functions A_3 , B_3 , A_6 , B_6 , . . . obtained from the Fourier analysis of the magnetic field. The values thus found agree very well with the values previously calculated without making any simplifying assumptions, so it would seem to be justified to use the above-mentioned expressions in further projects.

The stabilizing effect of the clover-leaf field, which we postulated on the basis of general arguments above, is given quantitative expression in equation (8): even if the field index is positive (field increasing with increasing distance from the centre), vertical stability is still possible as long as the coefficient C_3 is large enough, and the spiral shape of the sector boundaries ($\gamma > 0$) also helps in this direction.

In order to give some idea of the values found in practice, we may mention that the maximum field index of our clover-leaf cyclotron is +0.05. Further, $C_3 \approx 0.3$ and $\gamma = 35^\circ$, from which $(f_z/f)^2 = +0.04$.

Variation of the final energy of the ions

If it is desired to bring the beam of high-energy ions produced in a cyclotron out of the accelerating space (which is the case in most nuclear physical experiments), then a beam-extraction system must be introduced into the magnetic field¹¹). Ions with lower energies than the maximum are in principle available in orbits with smaller radii, but in practice it is not possible to extract them by shifting the extraction system towards the middle of the cyclotron. In order to vary the energy of the ions, which is desirable for many investigations, one must therefore extract the ions from the same orbit and vary the induction B , i.e. alter the current flowing through the energizing coils of the magnet.

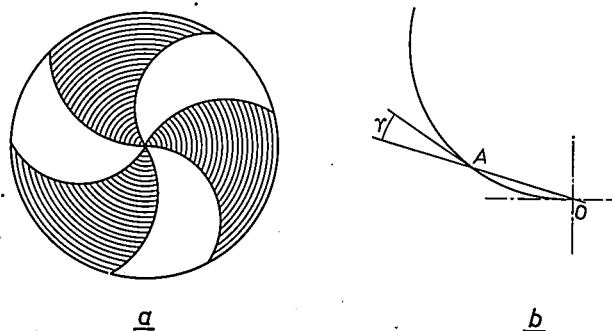


Fig. 3. Sketch of an azimuthally varying field of 3-fold symmetry, with sectors bounded by spirals. a) Overall view of the field. b) The angle γ of the spiral at the point A of the spiral sector boundary is the angle between the radius OA and the tangent at A .

¹⁰) The equations of motion and the approximate analytical solutions referred to here may be found in H. L. Hagedoorn and N. F. Verster, *Orbits in an AVF cyclotron*, Nucl. Instr. Meth. 18/19, 201-228, 1962.

¹¹) The design of a magnetic beam-extraction system for a synchrocyclotron has recently been described in this Review in the last article quoted in ²).

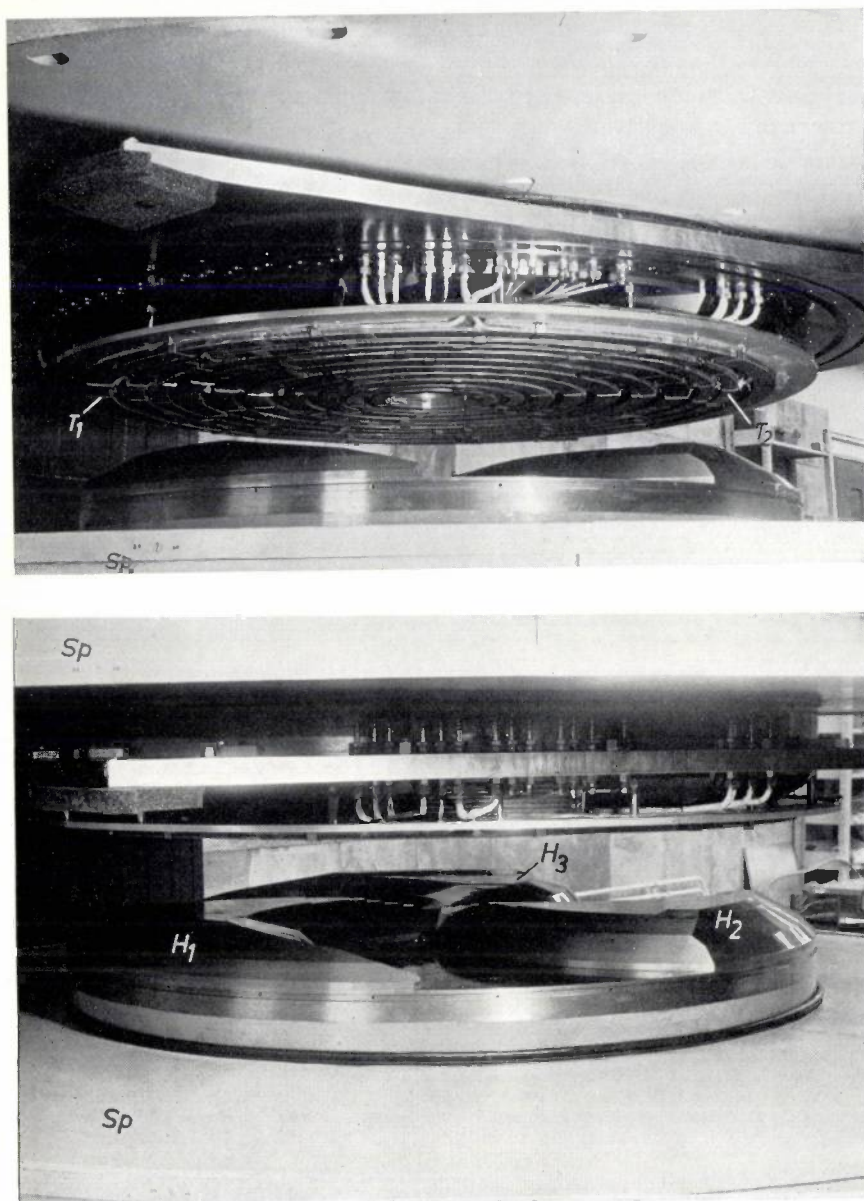


Fig. 4. In the isochronous cyclotron under construction in the Philips Research Laboratories, the desired form of the magnetic field is obtained with the aid of "hills" and "valleys" on the pole pieces and ten pairs of concentric trimming coils. When these photos were made, the trimming coils had as yet only been fitted on one pole. The lower photo shows the shape of the bottom pole of the magnet, and the upper one the ten concentric trimming coils on the top pole. *Sp* main coil of the magnet. H_1 , H_2 and H_3 hills, separated by the three valleys. T_1 and T_2 the two outermost trimming coils.

This method, which is used in some classical cyclotrons and in most isochronous ones (including ours), has important consequences.

In the first place, the angular velocity ω is directly proportional to B (see eq. 2); the HF oscillator which supplies the dee voltage must thus be tunable. We shall not discuss this point in any further detail.

Secondly, when the current in the excitation coils is altered, the field distribution $B(r)$ in the air gap does not change quite uniformly, owing to local magnetic saturation of the steel of the pole pieces. Although the edges of the hills on the poles have been

bevelled so as to prevent strong field concentration at the edges and thus to limit saturation as much as possible, this effect is still noticeable. However, this does not matter as much as might be thought, since for other reasons the given field distribution $B(r)$ can no longer be used when another final energy is required. This has nothing to do with the stability of the orbits — the Fourier coefficients $A_n(r)$ and $B_n(r)$ usually remain large enough to ensure vertical stability, even if they have been altered by the saturation effect — but with the isochronism.

This is because when the ion energy on the fixed outer circle is altered, the field distribution needed for isochronism changes since the relativistic increase of mass is now a different function of the radius. (This is the third and most important consequence.) The same applies if we want to accelerate different kinds of particles (e.g. deuterons and α particles as well as protons). In our cyclotron, we can accelerate protons to 25 MeV with an induction on the axis of $B_0 = 1.4 \text{ Wb/m}^2$; the mean induction at the periphery must then be 0.037

Wb/m^2 more, because of the increase of mass of 2.5%. With the same value of B_0 , deuterons are accelerated to 12.5 MeV, the increase of mass is then about 0.6%, and the mean induction at the periphery need now only be about 0.009 Wb/m^2 greater than B_0 . In order to realize the desired isochronous field for each ionic species and each final energy, it is thus not enough simply to alter the current flowing through the excitation coils of the magnet; extra measures must be taken.

For this purpose we have fitted ten pairs of concentric trimming coils on the two poles of the magnet (see fig. 4). Excitation of one of these pairs mainly

effects $\bar{B}(r)$ for values of r lower than the radius of the coil. The largest coil causes B_0 to change by about 0.025 Wb/m^2 . Each of these ten pairs of coils (and the main coil itself) must now be activated in an appropriate manner for each desired ionic species and each final energy. This kills two birds with one stone, since changes in the magnetic saturation of the magnet steel are also compensated for. The determination of the excitation currents is in fact the main problem treated in this article.

It should be mentioned in this connection that a very great accuracy is demanded for this purpose. Even though we have chosen the relatively high value of 50 kV for the dee voltage, a proton still has to make 250 revolutions to reach an energy of 25 MeV. Let us suppose that we have chosen $\bar{B}(r)$ only 0.1% less than the required isochronous field throughout; the ion would then undergo a phase lag of 0.36° with respect to the dee voltage each revolution, and after 250 revolutions the phase lag would be 90° . It is thus clear that the isochronous field must be realized with an accuracy of considerably better than 1 part in a 1000.

This means that the magnetic field must be measured at a relatively large number of points. Since the effect of each of the ten pairs of coils must be investigated for a number of values (e.g. ten) of the current through the main excitation coil, the number of measurements involved is enormous (more than 100 000). We therefore had to develop a special technique for measuring, recording and processing the data in question. We shall now discuss this technique in some detail.

The measurement of the magnetic field

The magnetic field of the isochronous cyclotron under construction in the Philips Research Laboratories has been measured — one could perhaps better say “mapped” — with the aid of automatic measuring equipment whose real active element is a *Hall probe*. This probe is placed in a holder which can be moved along a rotating arm mounted on a vertical axis which is situated at the centre of the magnetic field (fig. 5). The displacement of the probe along the arm can occur in steps of 2 mm or multiples thereof, and the arm rotates in steps of 2° (or 4° or 6°).

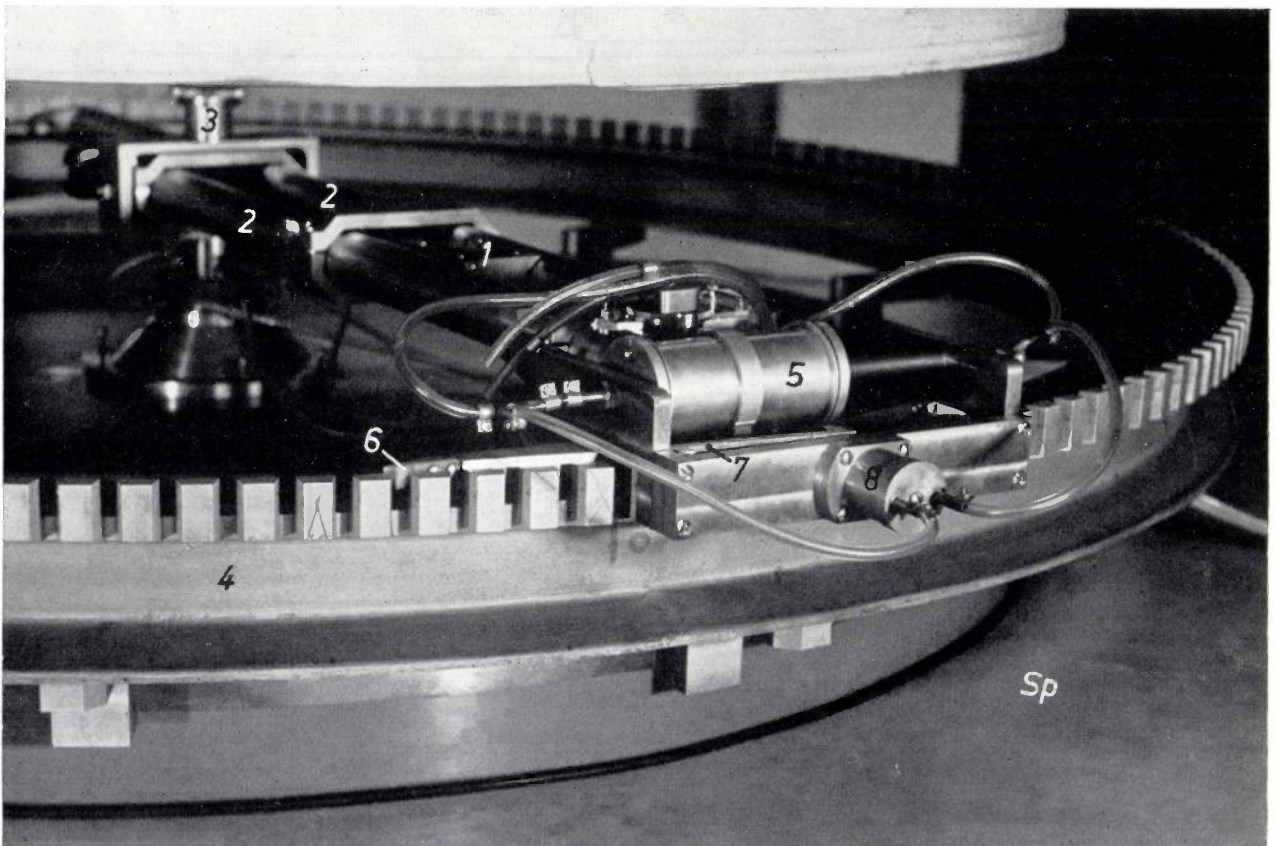


Fig. 5. The automatic equipment for measuring the magnetic field. The measuring element, a Hall probe, is contained in the holder 1 which can be displaced radially between the bars 2 which together form the arm. The arm, which can rotate about the vertical axle 3 situated at the centre of the magnetic field, is supported at its outer end by the toothed aluminium ring 4.

Each tooth corresponds to a rotation of 2° . The motion of the arm and the holder is operated pneumatically. The rotation is actuated by the cylinder 5 and the catch 6, while the exact positioning of the arm is controlled by a stud on the bar 7, which is operated by the cylinder 8. Sp is the top of the main coil surrounding the lower pole of the magnet.

Since we wanted to measure the field with an accuracy of about 1×10^{-4} Wb/m² (1 gauss) and as the field gradient at the edges of the sectors is very large, the position of the Hall probe had to be accurate to within about 0.05 mm. The azimuthal position of the arm is accurate to within 12 seconds of arc owing to the use of a toothed ring round the edge of the magnet, on which the arm rests. The required precision of the displacement of the holder along the arm was obtained with the aid of a perforated bronze strip.

The displacement of the Hall probe is produced by means of compressed air. This method lends itself both to automation and to operation in a magnetic field. The equipment allows a measurement to be carried out every 5 seconds. The course of each 5-second cycle is determined by a rotary switch. It is possible to carry out measurements along the circumference of a circle ($r = \text{constant}$), and also along a given radius ($\theta = \text{constant}$). A similar, but smaller and simpler, device was constructed for measurements on a model magnet (scaled down 5 times) to determine the best shape of the hills and valleys in the poles of the magnet.

The Hall probe was fed with a current of 18 mA. The sensitivity at this current is about 3×10^{-2} Vm²/Wb (= 3 μ V per gauss). The DC voltage thus obtained was amplified 300 times using a chopper amplifier, with a very constant amplification factor ($1 : 10^4$), owing to strong negative feedback ($5000 \times$). The noise in the output signal corresponds to about 1 μ V, so a difference in induction of 1×10^{-4} Wb/m² is easily detected.

Also in the interests of precision, the current supply of the Hall probe is kept very constant (better than 1 in 3×10^4), and the probe itself is kept in a small thermostat.

The Hall probe was calibrated by placing it, together with the measuring probe of a proton-resonance set-up, in the field of an auxiliary magnet¹²). This field was then so adjusted that nuclear magnetic resonance occurred at 4 Mc/s and at all harmonics of it up to the 21st (i.e. up to 84 Mc/s). The inductions corresponding to the lowest and highest frequencies are 0.0940 and 1.9738 Wb/m² respectively. By working with a series of harmonics, we were independent of the quality of the standard oscillator as regards linearity. This procedure gave 22 calibration points (including that for zero field).

The relationship between the induction and the output voltage of the Hall probe is linear to within

¹²) The measurement of magnetic fields by means of the nuclear magnetic resonance of protons is described in e.g. Philips tech. Rev. 15, 55, 1953/54.

about 1%, but because of the great precision required it was necessary to represent it by a seventh-degree polynomial. Regular checks of the output voltage at the highest resonance frequency (84 Mc/s) have shown that the apparatus is very stable: only small and infrequent corrections were necessary. This correction is made by varying the current through the Hall probe until the output voltage returns to its original value. The deviation in the output voltage at lower frequencies then corresponds to less than 1×10^{-4} Wb/m². Fig. 6 gives an idea of the variation in the measured values due to random fluctuations ("noise").

The amplified output signal of the Hall probe is fed to a digital voltmeter with a capacity of four digits and a sensitivity of 1 mV per unit. The equipment is so adjusted that this meter reads -9995 in the absence of a magnetic field and +9995 at an induction of 2.0000 Wb/m²: one unit thus corresponds very nearly to 1×10^{-4} Wb/m². The reading of this meter is automatically typed out, and also punched in 5-channel telex tape, the latter form being suitable for feeding into a computer. Fig. 7 gives an impression of the equipment during measurements on the model magnet.

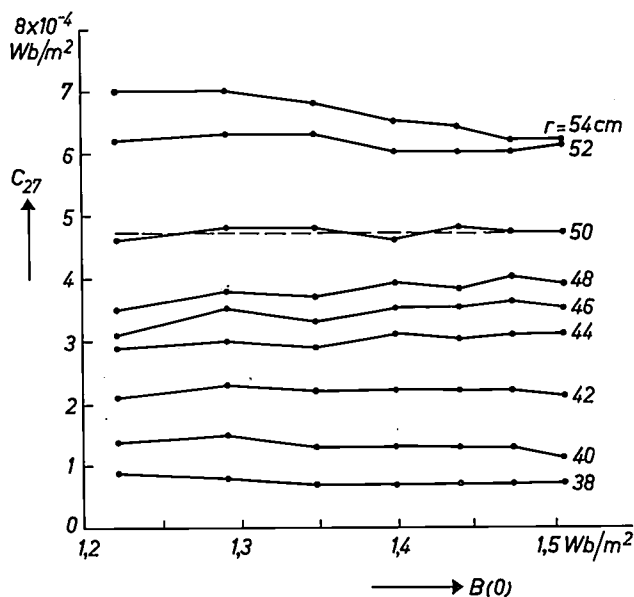


Fig. 6. The very slight spread in the measurements may be seen from e.g. this plot of the amplitude C_{27} of the 27th harmonic of the induction around circles of various radii, as a function of the excitation of the main coil (which is expressed as the induction $B(0)$ at the centre of the field). The spread in the points plotted here is of the order of 10^{-5} Wb/m². (This can be seen clearly from the points for $r = 50$ cm, where a smooth curve — in this case a straight line — has been drawn in.) Since each point represents the mean of 60 measurements, the spread of the individual measurements is about $10^{-5} \times \sqrt{60} \approx 0.8 \times 10^{-4}$ Wb/m². (It should be remembered in this connection that the mean values are not obtained from continuously variable values, but from readings of the digital voltmeter, which are "quantized" in steps of near enough 1×10^{-4} Wb/m².)

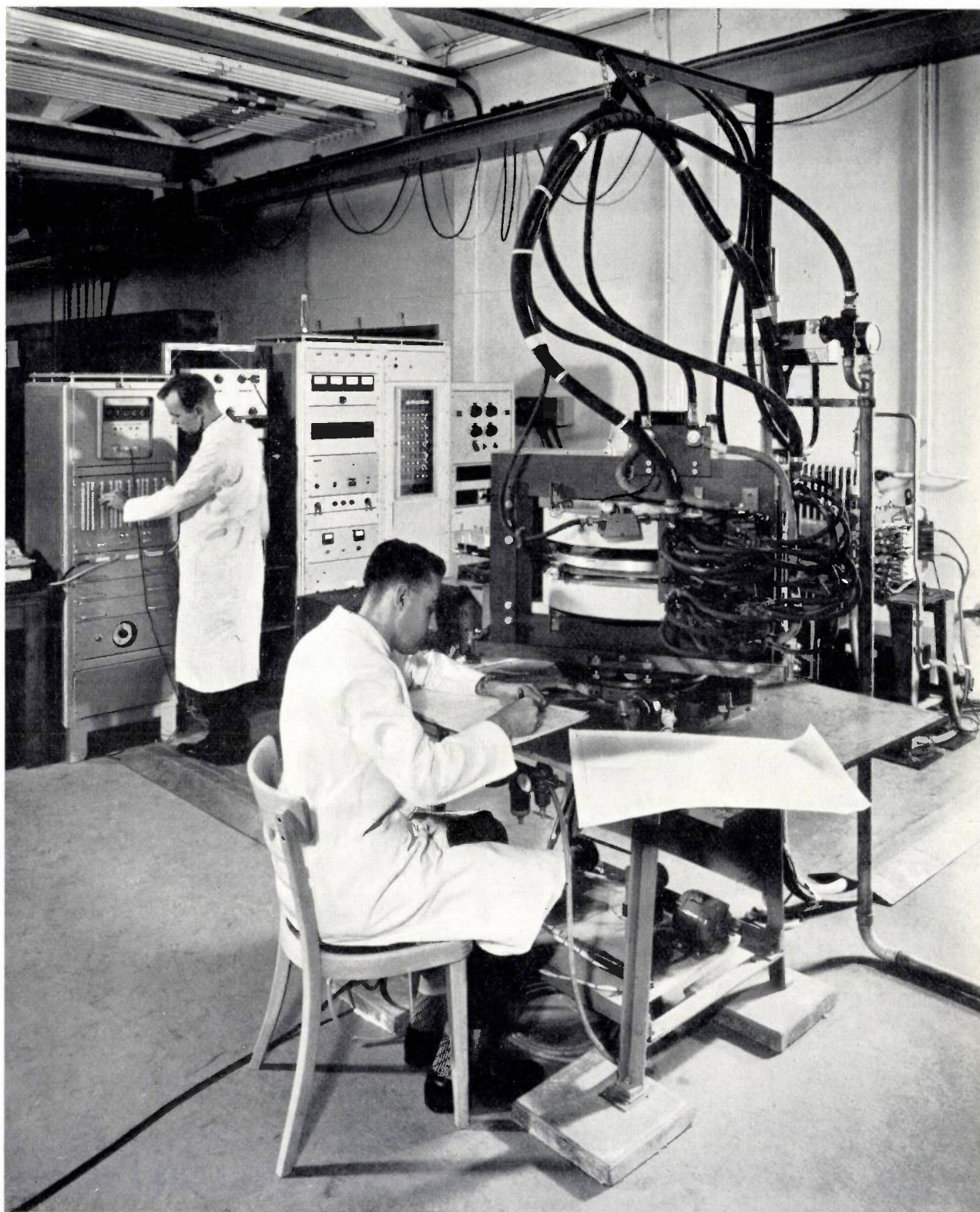


Photo Adolf Morath

Fig. 7. The above photo shows field measurements being carried out on the test magnet (middle foreground), in order to determine the most suitable form of the hills and the valleys. In this case, the Hall probe was not situated on a rotating arm, but in a slit in the disc which can be seen between the poles of the magnet. At the top of the left-hand rack of measuring equipment can be seen the digital voltmeter, and underneath it the keyboard for setting the heading which enables the series of measurements punched on a given tape to be identified. The apparatus in the other racks is mainly for the stabilization of the magnet current. On the right (only partly visible) is the proton-resonance equipment. The rotary switch which controls the field-measurement equipment is under the table which carries the magnet (largely not visible).

The investigation of this model magnet (scaled down 5 times) was necessary because it is impossible to calculate the required form of the hills and valleys exactly. The form which gives the best results must be found by trial and error, by finding how different variables change with different shapes of the pole pieces and then trying to interpolate to the best shape. Fairly rough models can be used to begin with, but the final ones must have a precision of at least 1 in 10^3 . Naturally, the quality of the steel used for the model magnet, and the induction, must correspond to that intended for use in the cyclotron itself. The power consumption of the coil of our model magnet was therefore no less than 25 kW.

The identification of a measurement and the detection of errors

In order to be able to identify the series of measurements punched on a given tape, we provided each tape with a "heading" of ten digits. The measuring equipment was therefore provided with a key-board containing ten columns of ten keys, to enable the entire heading to be set by hand. After it has been completely set up, it is punched into the tape by pressing a button.

The first digit of the heading indicates the type of measurement (e.g. 2 = azimuthal series, 5 = radial series, 9 = calibration), the next three contain the serial number of the field investigated, the next three the position (θ in a radial series and r in an azimuthal series), the next two e.g. the interval at which measurements were made and the last one the calibration formula to be used. An example of a piece of tape with such a heading is shown in fig. 8. To control the "input-check programme"

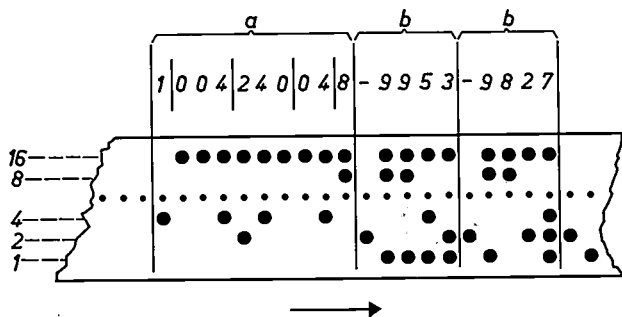


Fig. 8. The beginning of a series of measurements as punched on 5-channel tape. Each column contains one number (character). The tape is read, one character at a time, in the direction of the arrow. The first ten characters form the "heading" a , and characterize the nature of the measurement, the measuring interval, etc. Each subsequent set of five characters b forms a four-figure number (one measured value), preceded by a plus or minus sign. The measured values are thus not completely coded in the binary system, but punched per (decimal) digit. The coding is as follows. A hole in the fifth channel means that the rest of this character is a number. If this hole is absent, the rest is a code sign; for example, a hole in the second channel signifies a minus sign, and a hole in the first channel a plus sign. In order to facilitate the reading of the tape, the first digit of the heading is coded differently (no hole in channel 5; number punched is digit in question plus 3). The small dots represent the sprocket holes.

(discussed below), we use a heading which starts with the digit 0, while the second digit indicates the instruction. For example, 08 means "reject the last measurement", $05 x_1 \dots x_8$ means "read $x_1 \dots x_8$ as a positive number and regard it as a measurement".

The punched tapes on which calibration measurements are recorded are also provided with a heading, which contains the date and an identification number.

The automatic processing of the data makes it possible to handle the vast number of measurements, but it also provides the investigator with a new problem. When data are processed by hand, there is always some sort of check on the *reliability* of the measurements — for example, if the results are plotted in a graph it can be seen instantly whether one value deviates by an improbable amount from the expected value. In automatic data processing, there is absolutely no check unless special measures are taken.

We have therefore provided both programmes, according to which the computer processes the results of the measurements, with a part which ensures that each value is compared with the other values in the same series so that no point deviates too much from a smooth curve imagined drawn through the plotted results. Every value which deviates from the expected value by more than 2×10^{-3} Wb/m² (20 gauss) is detected and replaced by a value obtained by interpolation from neighbouring values. These programmes also provide an indication of the spread of the accepted values.

The various computer programmes

We shall now survey the computer programmes according to which the data obtained with the equipment described above are processed. For the sake of simplicity, we shall restrict ourselves here to a short description of what the programmes achieve and their significance. The mathematical basis of some of the programmes is briefly described in the final section of this article¹³). All the programmes involved are shown in the block diagram of fig. 9. In this diagram, the blocks M represent different types of measurements and the blocks P different programmes carried out by the computer PASCAL. The connecting lines T represent the temporary storage of the results of a measurement or a calculation on punched tape, and the feeding of these results into a later programme.

¹³) Most of these programmes were prepared by the staff of the Philips Computing Centre. In particular, H. Q. J. Meershoek prepared all the programmes in which the orbit equations had to be solved.

The three programmes *P1*, *P2* and *P4*, which are used for processing the punched tapes obtained from the different types of measurement of the field (*M1* to *M3*), all begin with the "input-check programme". This programme ensures that the machine decodes the characters on the tape, and investigates them to make sure that the tape it has been given really contains a series of field measurements; it does this by checking the *sequence* of the characters. As we have mentioned above, on the tape must be either a "heading" or a four-figure number coming from the voltmeter, preceded by a plus or minus sign. If this is not so, the series in question is rejected. As we have said, the instructions (05, 08, etc., see above) determine which special part of the input-check programme should be carried out in a given case.

The checking and identification of a tape (the input-check programme) may be divided into the following steps:

- 1) A character is "read", i.e. as soon as the sprocket hole passes the reading station of the machine, the character is transferred as a binary number to the arithmetical unit; the machine then checks whether the number is greater or less than 16.
- 2) If it is less than 16 (i.e. if there is no hole in channel 5), it is either a plus or a minus sign or the start of a heading; the machine investigates this further.
 - 2.1) If it is a 1 or a 2, it represents a plus or minus sign, and it should be followed by four digits (i.e. channel 5 should be punched four times); the machine checks whether this is so.
 - 2.2) If it is between 3 and 12, it is the first character of a heading, and should be followed by nine other digits (i.e. channel 5 should be punched nine times); the machine checks this too.
 - 2.3) If it is 13, 14 or 15 the machine rejects the series, since these numbers do not occur in the code.
- 3) If it is between 16 and 25 (i.e. if there is a hole in channel 5), it represents one of the digits from 0 to 9; the cipher in question is found by subtracting 16 from the number punched in the tape: the result *V* of a measurement can thus be calculated as follows from the four characters in question, k_1 to k_4 :

$$V = 10 [10 \{10(k_1 - 16) + k_2 - 16\} + k_3 - 16] + k_4 - 16.$$
- 4) If the number lies between 26 and 31 (larger numbers cannot be punched in the binary notation on 5-channel tape), the machine again rejects the series. After a series is rejected, all punchings are ignored until one between 4 and 12, i.e. until the first digit of the heading of the next series is met.

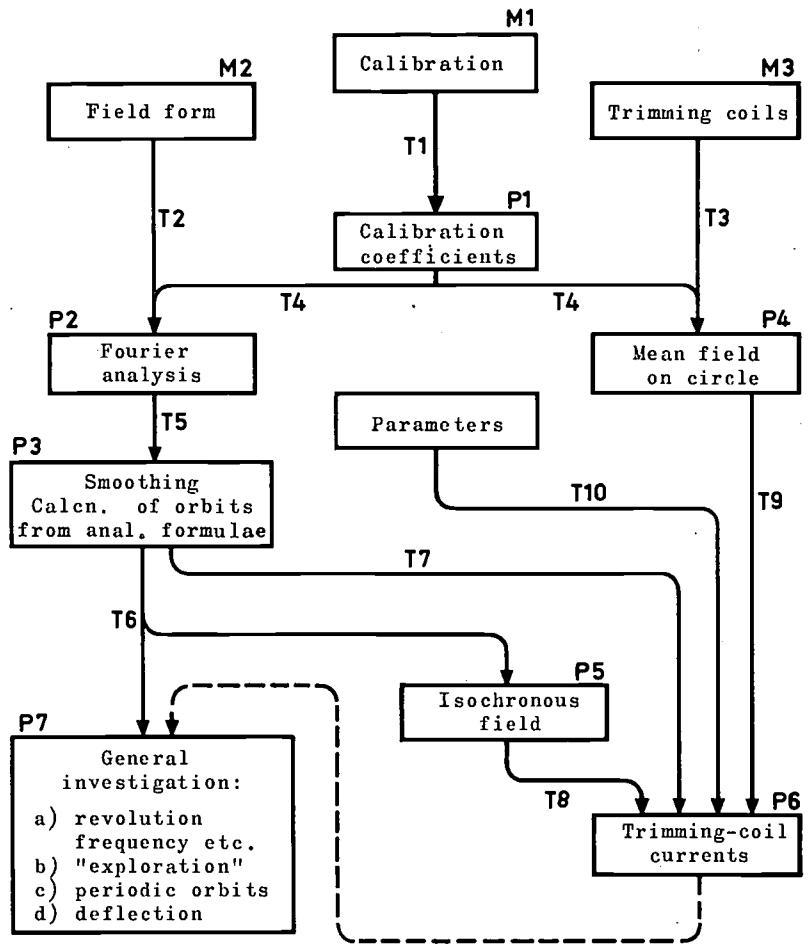


Fig. 9. A schematic view of the nature and sequence of the calculations which are carried out. The blocks *M* represent measurements, the blocks *P* are computer programmes, and the lines *T* symbolize the storage of the results of a given programme on tape, which is later fed into a subsequent programme.

By punching the results of the measurements per (decimal) digit, the coding of large numbers in the binary notation is avoided, but it naturally uses somewhat more binary digits (bits); the complete binary coding of the numbers takes place in the machine. The machine takes $\frac{1}{4}$ second to treat a series of 60 measurements, plus heading, in the manner described above (reading, checking, decoding of the numbers punched in the tape and systematic storing). This time is determined by the speed at which the tape passes the reading station; the time needed for the actual calculation is about half of this. As we have mentioned above, the time taken to make the measurements themselves is 60×5 seconds, i.e. 5 minutes.

The programmes P1, P2 and P4

As mentioned above, the curve through the 22 experimental points involved in a calibration is represented by a seventh-degree polynomial:

$$B = \sum_{n=0}^7 c_n x^n, \dots \dots \dots (9)$$

where *x* is a given reading of the digital voltmeter and *B* is the corresponding induction. In the calibration, *B* is punched in the tape as the frequency

of the proton resonance, and x as *three* readings of the voltmeter, i.e. each point on the calibration curve is obtained from three measurements. The eight coefficients c_n are calculated with the aid of programme *P1* by the method of least squares¹⁴).

The programmes *P2* and *P4*, which process the results of measurements *M2* and *M3* respectively, both start (after the input-check programme) with a part in which the measured values are converted into values of B by use of eq. (9).

The rest of the programme *P2* carries out a Fourier analysis on all azimuthal series of measurements (which extend, as we have mentioned, over one 120° sector). Since the sectors are assumed to be identical, only those coefficients appear whose serial number is a multiple of three. These calculations give in the first place the quantities \bar{B} , a_{3m} and b_{3m} which occur in the expression

$$B(\theta) = \bar{B} + \sum_{m=1}^{\infty} a_{3m} \cos 3m\theta + b_{3m} \sin 3m\theta. \quad (10)$$

It will be seen that this is nothing more than equation (6), expressed in a more suitable form for our purposes.

The machine continues the Fourier analysis as far as the value M of m at which a and b are of the same order of magnitude as the noise. (The values of a and b decrease rapidly with increasing m .)

The machine is then instructed by *P2* to check whether there are any experimental values of B which differ by more than 2×10^{-3} Wb/m² from the values calculated from equation (10), using the M pairs of Fourier coefficients which have just been found (Fourier synthesis). If such experimental values are found, they are replaced by a value found by interpolation. If the maximum deviation is still too large after three such corrections, the whole series of measurements is rejected. The punched tape *T5* which the machine delivers on completion of *P2* contains \bar{B} and the coefficients a_{3m} and b_{3m} for $m = 1$ to 6; higher harmonics are neglected in the further calculations. The printed paper which the machine delivers at the same time as *T5* gives, moreover, the maximum value of the difference between the measured values and the Fourier synthesis, as a measure of the noise involved in the measurements.

The programme *P4* processes the results of the

measurements on the ten pairs of trimming coils. In this case the main field is first measured, and then the field obtained when current flows through each pair of trimming coils in turn. Ten radial series covering an entire 120° sector are measured each time. Fluctuations of more than 2×10^{-3} Wb/m² are again detected and replaced by values obtained by *radial* interpolation. The $\Delta \bar{B}(r)$ values for each pair of coils are punched on tape (*T9*). The heading of this tape contains among other things the serial number of the main field and that of the pair of trimming coils involved.

The programmes *P3*, *P5* and *P6*

With the aid of the programme *P3* all results obtained with the aid of *P2* which refer to one particular field (and which thus have the same field serial number in their heading) are combined to give tables of $\bar{B}(r)$, $a(r)$ and $b(r)$. These series are also "smoothed" and any numbers which may be missing are filled in by interpolation.

The following quantities are then calculated, using the above-mentioned approximate analytical solution of the equations of motion:

- 1) The frequencies of the horizontal and vertical oscillations of an ion about its equilibrium orbit (neglecting the effect of the relativistic increase in mass).
- 2) The variation of \bar{B} needed for isochronism (ditto).
- 3) Two parameters which characterize the horizontal stability¹⁵).

When these calculations are completed, the machine punches *two* tapes. One (*T6*) contains the 13 tables which give $B(r)$, $a_{3m}(r)$ and $b_{3m}(r)$ for $r = 0, 2, 4, \dots, 70$ cm. (One table of $\bar{B}(r)$ and twelve of the Fourier coefficients.) The other tape (*T7*) contains only $\bar{B}(r)$. The tape *T6* can be used in the programmes in which the orbit equations have to be integrated, while *T7* is used in the programme *P6*, for calculating the currents which must flow through the trimming coils in order to make the mean field satisfy the conditions for isochrony.

The programme *P5* contains parts in which the orbit equations are integrated numerically, so tape *T6* is made use of here. By means of an iteration process, which we shall not describe further, the equilibrium orbits are calculated for radii increasing from 4 to 54 cm, in steps of 2 cm. (We shall denote such a series by $r = 4(2)54$ cm.)

P5 then calculates the form \bar{B}_{is} of $\bar{B}(r)$ needed for isochronism, also by an iteration process. First of all, the angular velocity ω is calculated for each of the

¹⁴ Descriptions of this method can be found both in chapters on regression analysis in text-books of statistics (e.g. P. G. Hoel, Introduction to mathematical statistics, Wiley, New York 1954 or M. J. Moroney, Facts from figures, Penguin books 1951) and in books on experimental techniques in physics (e.g. F. Kohlrausch, Praktische Physik, 20th edn., Vol. 1, chap. 1.22, Teubner, Stuttgart 1955).

¹⁵ The parameters D_1 and D_2 of the article quoted in ¹⁰).

equilibrium orbits which have just been calculated. As in *P3*, the contribution of the relativistic increase of mass is neglected; this is later inserted for each ionic species (each value of e/m) separately. Since ω is roughly proportional to \bar{B} , the corrections ΔB which have to be applied to achieve isochronism can be calculated with good accuracy from the corresponding values of $\Delta\omega$. Since this process converges rapidly, it only needs to be repeated a few times to obtain a very good approximation.

Finally, f_r and f_z are calculated by integration of the orbit equations, also for non-relativistic particles.

On the basis of the above results an estimate is made of the desired form of $\bar{B}_{is}(r)$, taking the relativistic increase of mass into account, for ions with various values of e/m . Using this estimate as a first approximation, the above-mentioned iteration process is used to calculate $\bar{B}_{is}(r)$ for protons, deuterons, α particles, ${}^3\text{He}^{++}$ ions and ${}^4\text{He}^+$ ions. These have the e/m values 1, $\frac{1}{2}$, $\frac{1}{2}$, $\frac{2}{3}$ and $\frac{1}{4}$ (taking that of the proton as 1).

Programme *P6* then calculates the currents I_j which should flow through the various trimming coils in order to give the best possible approximation to the desired field in a given case (i.e. for a given ionic species and a given final energy). These calculations make use of the values of $\bar{B}(r)$ for the main coil (measurement *M2*, punched in tape *T7*), the isochronous (mean) field $\bar{B}_{is}(r)$ calculated with *P5* (tape *T8*), the known influence of each pair of trimming coils on the variation of the mean field (tape *T9*) and finally a group of parameters, punched in tape *T10*, which define what is understood by "the best possible".

Expressed in words, the condition "the best possible" entails that:

- 1) The phase difference between the dee voltage and the revolution of the ions must be as small as possible.
- 2) In order to prevent vertical instability, the difference between the field index of the actual field and that of the isochronous field must not exceed a certain small value. (Since the ideal field distribution is approximated to with a finite number of trimming coils, $\bar{B}(r)$ fluctuates about the theoretical required value. Care must be taken that the corresponding fluctuations in k are not too great.)
- 3) In order to be able to extract the beam with different excitations of the magnet, the fringing field must be given a certain form; the actual form must be as close as possible to this.

The currents in question are then calculated by the method of least squares on the basis of these

three conditions together with the subsidiary condition that the current should not exceed a certain maximum value. The three conditions are weighted to different degrees. The above-mentioned group of parameters includes, among others, these weights, the maximum value of r associated with conditions 1) and 2), the lower limit of r associated with the third condition (the different regions must naturally overlap somewhat), and the maximum permissible values of the currents.

The programme *P7*

The block *P7* in fig. 9, which is marked "general investigation", represents a group of programmes, each of which contains one part in which the orbit equations are solved (and where the tape *T6* must be used, as in *P5*), and a specific part which organizes the further course of the calculation. In some cases, a disturbance in the field can also be introduced, e.g. a first harmonic $c_1(r) \cos \Theta$.

The programme *P7a* ("revolution frequency etc.") resembles *P5* but is simpler. It calculates the equilibrium orbits for each kind of particle with the measured field distribution, for $\bar{r} = 4(2)54$ cm, and the momentum, energy, oscillation frequencies and revolution frequency for each of these orbits.

The programme *P7b* ("exploration") investigates the behaviour of ions which oscillate in the horizontal plane with a large amplitude. For this purpose, each ion is followed for ten revolutions. The state of the ion at a given moment is expressed by the radius r and the slope $dr/d\Theta$ of the orbit at that moment with respect to the circle passing through that point; this slope will be denoted by r' and expressed in cm/rad. In order to distinguish free oscillations from the forced oscillations which the ions exhibit as a result of the azimuthal variation of the field, the ions are observed "stroboscopically", i.e. we only consider r and r' at $\Theta = 0^\circ, 120^\circ$ and 240° . The forced oscillations are thus eliminated. By allowing the machine to carry out these calculations for a large number of initial values an idea is soon obtained of the radial motion of the ion. The calculated values of r and r' can be automatically plotted in a graph by a data plotter. An example of such a radial "phase plot" is shown in fig. 10. The crosses in this figure represent invariant points, i.e. they refer to orbits in which the same values of r and r' are found after each revolution.

The values of r and r' for these "periodic orbits" can be found by iteration, using programme *P7c*, starting from an r - r' pair lying close enough to one of the invariant points. The periodic orbits which one finds in theory (there are at least seven of them,

including the equilibrium orbit of course) include both stable and unstable ones. The stability of the equilibrium orbit is greater if the crosses representing the three neighbouring unstable (or metastable) orbits in the phase plot are further away from the cross representing the equilibrium orbit itself.

By repeating this process for somewhat greater or smaller values of the momentum, it can be found how the points corresponding to the stable periodic orbits move in the phase plot as a function of the momentum.

Finally, the group *P7* contains a subgroup *P7d* of rather complicated programmes ("deflection"), which provide information about the possibility of beam extraction (both magnetic and electrostatic).

Some of the mathematical methods used

Fourier analysis (P2)

Let us suppose that we have measured a function $B(\Phi)$ at N points Φ_i distributed evenly round the circumference of a circle (the angle $\Phi_i = (i-1) 2\pi/N$, where $i = 1, \dots, N$) and that we want to calculate the coefficients of the series

$$B(\Phi_i) = \bar{B} + \sum_{k=1}^{N/2} a_k \cos k\Phi_i + b_k \sin k\Phi_i \dots \quad (11)$$

The functions $\cos k\Phi$ and $\sin k\Phi$ in which $B(\Phi)$ is expanded are orthogonal, i.e. the sum over i of their products is equal to zero:

$$\sum_{i=1}^N \sin k\Phi_i \sin l\Phi_i = 0$$

if $k \neq l$. Further:

$$\sum_{i=1}^N \sin^2 k\Phi_i = \sum_{i=1}^N \cos^2 k\Phi_i = N/2.$$

With the aid of these relationships we find

$$\bar{B} = \frac{1}{N} \sum B(\Phi_i),$$

$$a_k = \frac{2}{N} \sum B(\Phi_i) \cos k\Phi_i,$$

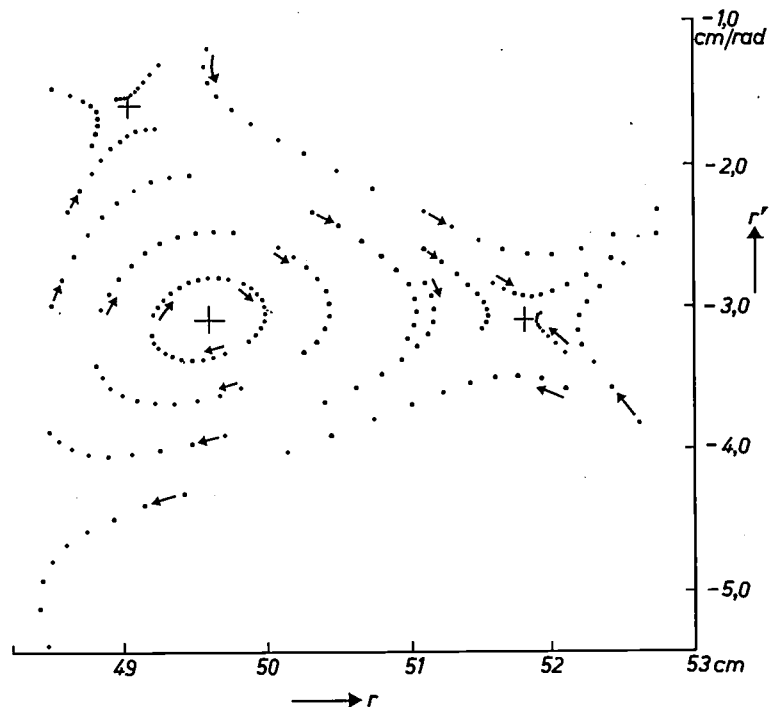


Fig. 10. The programme "exploration" (*P7b*) follows for ten revolutions the radial behaviour of ions which started from a point outside the equilibrium orbit, and at a certain angle to this orbit (but in the median plane). The situation at a given azimuthal angle is characterized by the radius r and the derivative $dr/d\theta$ (denoted by r'). This graph will be seen to consist of groups of ten points, corresponding to the successive passages of a single ion. The crosses represent orbits on which the ions have the same values of r and r' after each revolution. Such "periodic orbits" may be stable (e.g. the equilibrium orbit itself, represented here by the cross at $r = 49.6$ cm) or unstable. The average radius of the equilibrium orbit in this case was 48 cm. Even for this orbit r' is not in general zero since the orbit is not circular.

$$b_k = \frac{2}{N} \sum B(\Phi_i) \sin k\Phi_i.$$

Owing to the above-mentioned orthogonality, each shorter series

$$B_M^*(\Phi_i) = \bar{B} + \sum_{k=1}^M a_k \cos k\Phi_i + b_k \sin k\Phi_i, \quad (12)$$

where $M < N/2$, satisfies the condition laid down by the method of least squares:

$$\sum_i (B_i - B_{M,i}^*)^2 = \text{minimum} \dots \quad (13)$$

Here $B_{M,i}$ is a shortened way of writing $B_M(\Phi_i)$. (Because of the orthogonality, the series has only to be made longer to get a better approximation; there is no need to change the existing terms.) The quantity $B_i - B_{M,i}^*$, which is equal to the sum of the higher terms of the series which have been left out (the residue) will be denoted here by $d_{M,i}$:

$$d_{M,i} = \sum_{k=M+1}^{N/2} a_k \cos k\Phi_i + b_k \sin k\Phi_i. \quad (14)$$

The value of $d_{M,i}$ can be found by subtracting the right-hand side of (12) from the measured values of B .

During the calculation of the Fourier coefficients, the machine constantly keeps track of $d_{M,i}$. The first value, $d_{0,i}$, which applies as long as no coefficients have yet been calculated, is equal to $B - \bar{B}$ and is calculated as soon as \bar{B} has been determined. After the calculation of each successive pair of coefficients, the current value of $d_{M,i}$ is replaced by the next one, which is one Fourier term smaller:

$$d_{M,i} = d_{(M-1),i} - a_M \cos M\Phi_i - b_M \sin M\Phi_i.$$

After each stage of the calculation, the autocorrelation coefficient C_F of the residue d_i (we shall leave out the subscript M from now on) is calculated from the formula

$$C_F = \frac{\sum_{i=1}^{N-1} d_i d_{i+1} + d_N d_1}{\sum_{i=1}^N d_i^2}.$$

If the residue still contains periodic variations which are larger than the random fluctuations (noise) and whose Fourier coefficients have an index k which is less than $N/6$, it can be shown that C_F is positive. If the residue consists mainly of noise, C_F is negative. The machine continues the Fourier analysis until C_F becomes negative, so that no useless Fourier coefficients are determined.

It is not so important to finish the Fourier analysis at the right moment in order to save time — the time saved is not so much — but because in this way we split up the measured series $B(\Phi_i)$ as well as possible into an essential part B^* and a noise part d (cf. the discussion of $P2$ above).

The calculation of the calibration coefficients ($P1$)

The calibration curve contains 22 points $B_i = B(x_i)$, where x_i denotes the readings of the digital voltmeter. The method of least squares demands that we choose the coefficients of the approximation:

$$B_i^* = \sum_{n=0}^M c_n x^n$$

in such a way that B_i^* satisfies an equation similar to (11).

This calculation is based on the same principle as that of the Fourier coefficients. We start by transforming x_i^n in a group of orthogonal functions $X_{n,i}$ according to the procedure used for the construction of Legendre polynomials:

$$X_n = x^n - \sum_{l=0}^{n-1} \alpha_{n,l} x^l,$$

where the coefficients $\alpha_{n,l}$ are so chosen that

$$\sum_i X_{n,i} X_{l,i} = 0 \text{ if } n \neq l.$$

Here, too, a shorter expansion of $B(x_i)$ in powers of $X_{n,i}$ is a least-squares approximation. We found by investigation of the residue that this becomes smaller than 10^{-4} Wb/m² if M is chosen greater than or equal to 7. This is the reason for the statement that the calibration curve must be represented by a seventh-degree polynomial.

The smoothing of the Fourier coefficients ($P3$)

The tables of $\bar{B}(r)$, $a_{3m}(r)$ and $b_{3m}(r)$ based on measurements of B at the equidistant points $r_i = 0(2)64$ cm are smoothed in the following manner.

Let a function $y_i = y(x_i)$ be given at the equidistant points x_i . We determine a smooth function \bar{y}_i with the aid of the equation

$$\sum_i [\delta^4(\bar{y}_i)]^2 + K(\bar{y}_i - y_i)^2 = \text{minimum}, \quad (15)$$

where δ^4 is the fourth-order difference:

$$\delta^4(y_i) = y_{i+2} - 4y_{i+1} + 6y_i - 4y_{i-1} + y_{i-2},$$

and K is a parameter. By means of condition (15) we demand both that the function \bar{y}_i should be smooth (in the sense that the fourth-order differences are small) and that the deviations $d_i = \bar{y}_i - y_i$ should be small. The function is smoothed more strongly as K is chosen smaller. Naturally, K must not be chosen so small that the function y_i is mutilated and information is lost. The most suitable value of K is found during the calculation by using successively $K = 100, 30, 10, 3$, etc. in the determination of \bar{y} . After each calculation, the autocorrelation coefficient C_S of the deviation d is calculated from the formula

$$C_S = \frac{\sum_{i=1}^{N-1} d_i d_{i+1}}{\frac{1}{2}d_1^2 + \frac{1}{2}d_N^2 + \sum_{i=2}^{N-1} d_i^2}.$$

The difference between the definition of C_S and that of the autocorrelation coefficient C_F which was used in the Fourier analysis is due to the fact that in the Fourier analysis we were dealing with a cyclic series of values, i.e. a series in which the correlation between the last and the first point also plays a role.

For large K , C_S is negative; there is then no correlation between the deviations d_i and d_{i+1} , which is an indication that these consist entirely of noise. If K is so small that information is "planed off", there is naturally a correlation between the corrections applied, and C_S is positive. The calculation is therefore ended at the value of K at which C_S just becomes positive.

The processing of the measurements on the trimming coils (P4)

We have measurements $B(r_i, \theta_j)$ for $r_i = 0(2)64$ cm and $\theta_j = 0(12)108^\circ$ of the main field and of the field obtained when one pair of coils is activated in addition to the main coil. We obtain the contribution $\Delta B(r_i, \theta_j)$ of the trimming coil by taking the difference between the two values measured at each point.

It can be shown from equation (10) for the magnetic field that

$$\frac{1}{10} \sum_{j=0}^9 \Delta B(r_i, \theta_j) = \Delta \bar{B}(r_i) + \Delta a_{30}(r_i).$$

In words: the mean value of the difference ΔB measured at ten equidistant points distributed over a 120° arc is equal to the difference between the mean values of B plus the difference between the amplitudes of the 30th harmonic ($\cos 30 \theta$ is $+1$ at all the points in question; the other Fourier terms cancel out). Since a_{30} is only a few times 10^{-4} Wb/m², the term Δa_{30} may be neglected and we may write:

$$\Delta \bar{B}(r_i) = \frac{1}{10} \sum_{j=0}^9 \Delta B(r_i, \theta_j).$$

The measurements are checked by means of a relationship which is also derived from equation (10) by taking the sum $D(r_i)$ of the alternating series $+\Delta B(r_i, \theta_0) - \Delta B(r_i, \theta_1)$, etc. This gives:

$$D(r_i) = \sum_{j=0}^9 (-1)^j \Delta B(r_i, \theta_j) = 10 \Delta a_{15}(r_i).$$

Now a_{15} is less than 5×10^{-3} Wb/m² and Δa_{15} is less than 10^{-4} Wb/m². If everything is in order, therefore, $D(r_i)$ should be at the most about 10^{-3} Wb/m². The machine calculates the value of $|D|$ for each case, and checks whether this is less than 2×10^{-3} Wb/m². If this is not so, the fourth difference δ^4 is

calculated for each radial series, and the series in which the largest value of δ^4 occurs is corrected by radial interpolation. If $|D|$ is still too large, the series of measurements is rejected.

Summarizing, we may say that the programmes described above allow us to do the following things: to process a large amount of data without allowing improbable values to slip past unnoticed; to find the random fluctuations in all series of measurements, i.e. check whether the desired precision is indeed reached and maintained; to calculate the currents in the trimming coils needed to give different ionic species the desired final energy; to get an idea of the stability of the orbits; and finally, to obtain enough insight into the behaviour of the ions to make possible further development of this type of cyclotron (e.g. with respect to the beam extraction).

Summary. An isochronous cyclotron under construction in Philips Research Laboratories will enable different kinds of ions to be accelerated to a variable final energy (protons up to a maximum of 25 MeV). In an isochronous cyclotron the influence of the relativistic increase of mass on the revolution frequency is compensated by making the induction $\bar{B}(r)$ of the magnetic field increase with the radius r . The vertical stability of the orbits is ensured by making B vary periodically in the azimuthal direction by means of sector-shaped hills and valleys on the pole pieces. The Philips cyclotron is also fitted with ten pairs of trimming coils, which allow the form of $\bar{B}(r)$ to be varied to meet the demands imposed by the different ionic species and the different final energies. This article describes how this complicated field is measured — for ten different excitations of the main coil, in connection with the desired variation of the final energy — and how the more than 10^6 measured values are processed with a computer. The calculations include 1) construction of the calibration curve of the field-measuring equipment, 2) Fourier analysis of the field along circular contours, 3) smoothing of the tabulated values of $\bar{B}(r)$ and the Fourier coefficients, 4) the calculation of the effect of the trimming coils, 5) finding the field distribution $\bar{B}(r)$ needed for isochrony, 6) calculating the current through the trimming coils to give the best approximation to this field distribution, and finally various checks, e.g. for detecting improbable experimental results.

TRIM-LOSSES IN THE MANUFACTURE OF CORRUGATED CARDBOARD

by H. W. van den MEERENDONK *) and J. H. SCHOUTEN **).

65.012.122:676.76

The enormous amounts of packing material which the Philips factories require for dispatching their products are largely supplied by Philips' own corrugated-cardboard factory, which in turn draws its major raw material for this purpose — paper — from the Company's own paper factory. *Fig. 1* may give some idea of the large quantities involved.

*) Technical Efficiency and Organization Department, Philips, Eindhoven.

***) Philips Computing Centre, Eindhoven. Drs. Schouten died suddenly in November 1961.

Two large machines are installed in the corrugated-cardboard factory. One of them produces a sheet of corrugated cardboard about 2170 mm wide, the other a sheet about 1525 mm wide. In both machines the continuous cardboard sheet (see *figs 2* and *3*) is fed to a section in which it is cut into rectangular pieces. This is done by two types of cutting mechanism. The first type is a rotating cutting wheel which slits the sheet lengthwise; the second type consists of a set of "cut-off shears" which slice the



Fig. 1. The roll stores of Philips' corrugated-cardboard factory at Eindhoven. Every year these stores handle 36 000 metric tons of paper, which are processed into corrugated cardboard.

sheet transversely (see *fig. 4*). A considerable number of the first kind of cutting mechanism can be used in each of the machines so as to produce a number of narrow strips, each of adjustable width. Each machine has only two cut-off mechanisms, mounted one above the other at the end of the sheet and each adjustable to the required cut-off length. One or more of the strips produced by slitting the sheet lengthwise can be fed to the upper cut-off mechanism, and the others to the lower one.

In many cases the cut pieces of corrugated cardboard are passed for further processing to a department where they may, for example, be printed and made into boxes.

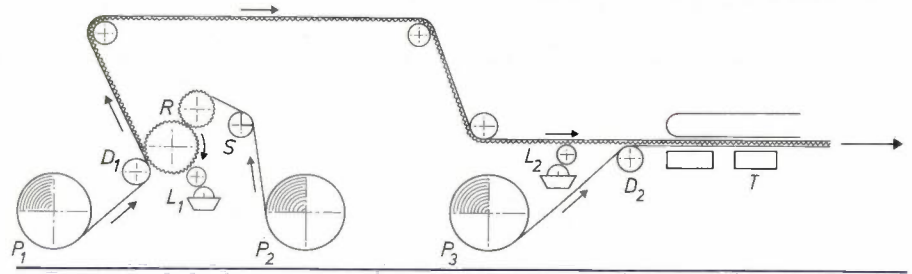


Fig. 2. Diagram illustrating the manufacture of corrugated cardboard. The corrugated cardboard is made from three sheets of paper fed from rolls P_1 , P_2 and P_3 . On the roller S the inside sheet is moistened and heated by steam. This sheet is then passed through the heated "corrugator roller" R , which impresses the required corrugation. Immediately afterwards the glue is spread over the tops of the corrugations of this inside sheet (glue pan and "rollers" L_1) and one outside sheet is then pressed against this side by a pressure roller D_1 . Further on, glue is spread on the other side of the inside sheet (glue pan and rollers L_2) and the other outside sheet pressed against it. D_2 is a guide roller. The sheet is then heated by steam plates T to harden the glue and dry the sheet. The finished sheet emerging at the right-hand end is then passed to the mechanisms which cut it into rectangular pieces of the required size (see *fig. 4*).

"Twin-board" (double-wall) corrugated cardboard is made from five rolls of paper, the second and fourth sheets being corrugated.

Let us now examine the production process. Every week the factory has to deal with a series of orders received from the various Philips factories (and from outside customers). In view of the wide diversity of

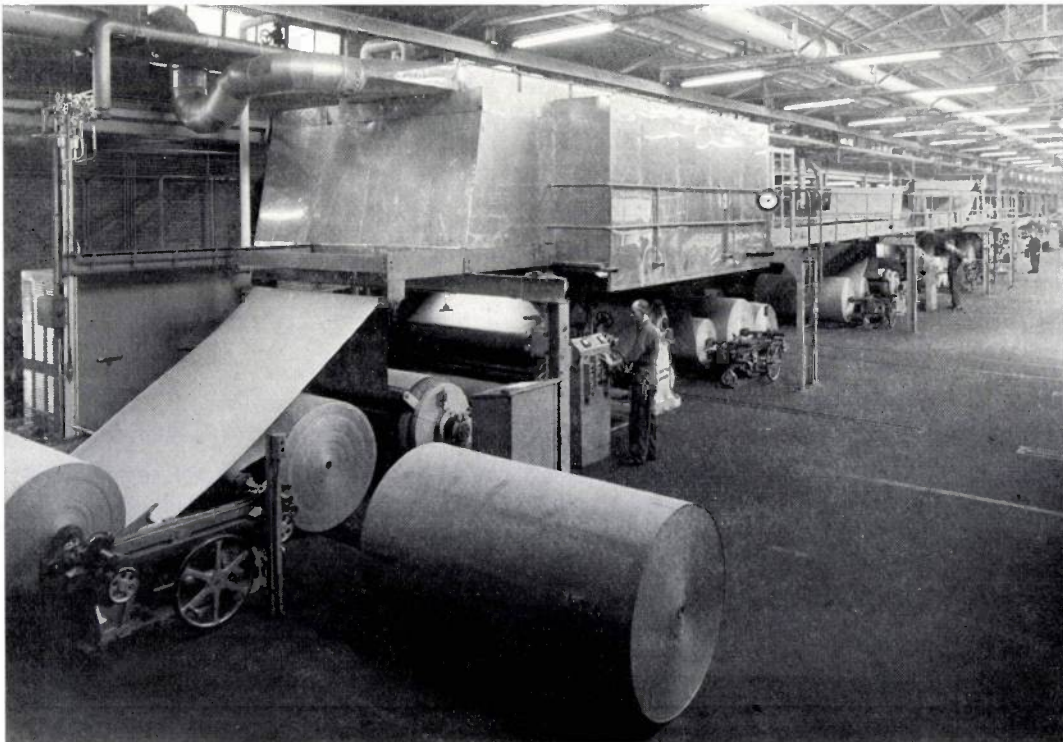


Fig. 3. Machine which produces corrugated "twin-board" sheet, 2170 mm wide. The production process begins on the left. The raw material — five sheets of paper — is fed into the machine from large rolls. Beside every roll of paper two other rolls stand by to take over immediately when it is used up, the machine thus being continuously fed with paper. The first stand-by for the roll on the far left (P_1 *fig. 2*) can be seen on the right of it in the machine, the second stand-by lies next to the latter in the foreground, outside the machine. Both reserves for the other working rolls are already in the machine. Above the operator at the control panel can be seen a large ventilation hood. The speed of the moving sheet can be read from the meter at the right end of the hood.

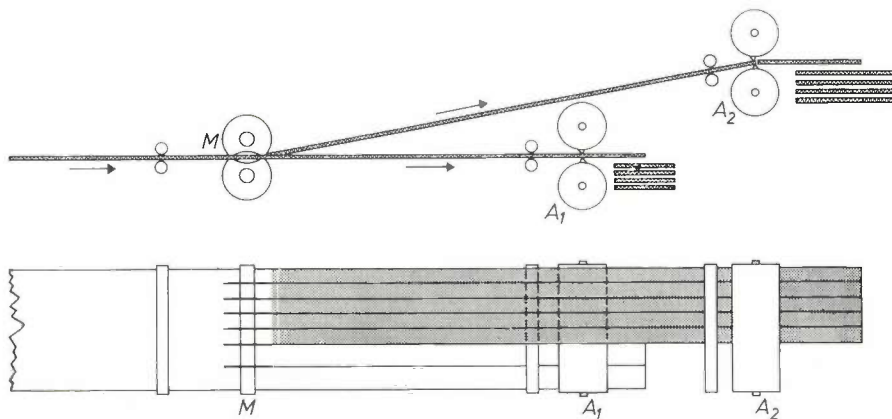


Fig. 4. Sketch of the two kinds of cutting mechanism. The corrugated-cardboard sheet arriving from the left is cut lengthwise by a number of cutting wheels M , the distances between which are adjustable, into a number of strips of the required widths.

Some of the strips then move straight on to "cut-off shears" A_1 which apply the transverse cuts, and the others travel upwards (shown shaded in the plan view) to another set of cut-off shears A_2 . These shears are moved up and down by an eccentric mechanism at an adjustable speed so as to cut off pieces of the required lengths. The reciprocating mechanism is so designed that the shears during the actual cutting process move at the exact speed of the travelling sheet, thus avoiding stoppages and tearing of the sheet.

the products to be packed, it is obvious that the orders differ considerably in regard to the size of pieces required. The problem every week is to consider how the two continuous sheets of corrugated cardboard produced are to be cut into pieces of the dimensions ordered. For in general the width of the pieces in an order from a factory will not divide a whole number of times into the fixed width of the uncut cardboard sheet, so that a strip of a certain width will be left over (the "trim"). By arranging pieces pertaining to various orders side by side in the uncut sheet, an attempt can be made to make fuller use of the available width (fig. 5); but then it must be ensured that an even broader strip will not be wasted when combining the remaining orders. The object must be to plan the combination of orders so as to minimize the total trim-losses of a week's production.

The orders can also differ as to the "quality" required, e.g. as to type of paper and depth of corrugation. This does not essentially affect the problem, but means that it must be considered separately for every group of orders of the same quality.

If there are not many orders, it is not too difficult to calculate which combination will give the lowest trim-losses. As the number of orders increases, the number of possible combinations increases so quickly that it soon becomes impracticable to compute "by hand" the minimum trim-losses for all combinations. This meant that reliance had to be placed on the experience and "flair" of the production planner.

The advent of the electronic computer has changed this situation. In cooperation with the production office of the corrugated-cardboard factory we have

analysed the problem of combining orders and programmed it for the PASCAL. For some time now, by way of experiment, the PASCAL has been computing at the beginning of every week the most favourable distribution of certain groups of orders (of different quality) from the total batch of orders. It does this in a matter of minutes. At the same time the production office still draws up its schedule in the old way for the same groups of orders, so that by comparing the results and analysing various methods

adopted in practice for increasing the efficiency of the factory, the means can be found to refine and suitably adapt the computing programme.

The manner in which the computations are performed using the PASCAL is the subject of this article.

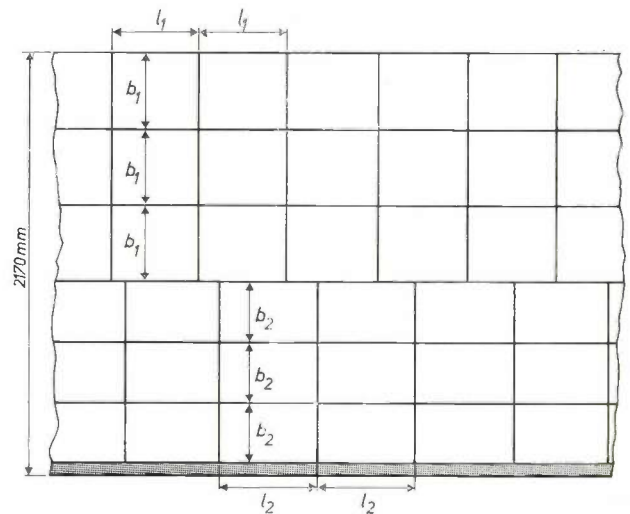


Fig. 5. The machine produces a continuous sheet of a fixed width, e.g. 2170 mm. The individual pieces of different sizes for the various orders have to be cut in such a way as to produce minimum waste of the available width (the trim is shown shaded). For this purpose a number of pieces from two different orders are cut side by side out of the width of the sheet (only from two orders, because the machine can only cut two different lengths from the sheet).

Linear programming

The corrugated-cardboard problem described here belongs to a large category of problems that can all be treated by the same method. A feature common to all these problems is that a total "yield" from

various sources has to be maximized, or a total sum in costs (or a loss) on various "expenditures" has to be minimized, while the sources or expenditures are not independent but interrelated through one or more subsidiary conditions. It is easily seen that this description fits the corrugated-cardboard problem. Problems of the same nature are found in widely diverse fields: e.g. in agriculture, where various crops of different prices and annual yield are cultivated in a given area; in stock breeding, where the feed has to be planned with specified proportions of starch, protein and fat, from a number of types of fodder of varying composition and price; in an oil refinery, where the crude oil is processed into a variety of cracking products of different values; further, in the distribution of shipping cargoes, in the control of traffic, and so on. In practice the above-mentioned subsidiary conditions frequently consist of *linear* relations between the variables. The task of determining the optimum production or expenditure programme in such cases has therefore been given the name "*linear programming*" — a somewhat confusing name, since it has nothing to do with the programming of the computer, which is generally employed for carrying out the calculations.

Methods of linear programming have only been developed in the last 15 years or so — in connection with the advent of electronic computers — and there have been numerous publications dealing with the mathematical theory and the practical applications of these methods¹⁾. We cannot of course go deeply into the theory here, but we shall try to explain, with the aid of a numerical example, the reasoning behind the procedure which we have adopted for the corrugated-cardboard problem. This procedure is called the Simplex method.

Numerical example (Simplex method)

For simplicity we assume that there is only one corrugated-cardboard machine, producing a sheet 2150 mm wide, and for the present we shall disregard the limitation of having only two different cut-off lengths. We consider a batch of only three orders, as specified in *Table I*.

Table I. Data of the batch of three orders taken as an example.

	Width	Length	Number of pieces
Order I	600 mm	1000 mm	1800
Order II	500 mm	800 mm	2400
Order III	230 mm	300 mm	4500

¹⁾ S. I. Gass, *Linear programming*, McGraw-Hill, New York 1958. W. W. Garvin, *Introduction to linear programming*, McGraw-Hill, New York 1960. S. Vajda, *Readings in linear programming*, Pitman, London 1958.

The respective widths of the pieces ordered are 600, 500 and 230 mm. We can thus cut 3 pieces for order I from the available width, or 4 pieces for order II, or 9 pieces for order III. These three modes of cutting ("combinations") are denoted D_1 , D_2 , D_3 ; see *Table II*. The use of these combinations is the most primitive manner of executing the orders.

Table II. The three most primitive "combinations" for cutting pieces for the different orders from the corrugated-cardboard sheet (width 2150 mm).

	D_1	D_2	D_3
Order I	3	—	—
Order II	—	4	—
Order III	—	—	9
Width of the trim	350 mm	150 mm	80 mm

We must then cut from the sheet of cardboard x_1 metres according to combination D_1 , then x_2 metres according to combination D_2 , and x_3 metres according to D_3 . The lengths x_1 , x_2 and x_3 must satisfy the equations:

$$\left. \begin{aligned} 3x_1 + 0x_2 + 0x_3 &= 1800, \\ 0x_1 + 4x_2 + 0x_3 &= 1920, \\ 0x_1 + 0x_2 + 9x_3 &= 1350. \end{aligned} \right\} \dots (1)$$

On the right side of the first equation is the total length (in metres) obtained when all pieces produced for order I are laid end to end. On the left are the contributions to this length of the pieces cut in accordance with D_1 , D_2 and D_3 . Similarly, the second and third equations apply to orders II and III respectively. It follows from (1) that

$$\begin{aligned} x_1 &= 600 \text{ m}, \\ x_2 &= 480 \text{ m}, \\ x_3 &= 150 \text{ m}. \end{aligned}$$

Instead of speaking of the width and length of the trim, it is obviously sufficient to take the total length of sheet used, $x_1 + x_2 + x_3$: in the solution with the lowest trim-losses this will be a minimum. With the primitive method of cutting in accordance with D_1 , D_2 and D_3 , we have $x_1 + x_2 + x_3 = 1230$ m.

This method is certainly not the most economical, for it can immediately be seen that in addition to three pieces for order I the working width can accommodate a further piece for order III. We shall now survey all the possible methods of cutting. This survey is presented in *Table III*, and contains, in addition to D_2 and D_3 , ten combinations C_1 - C_{10} . The system by which these combinations are found is self-explanatory, as also is the fact that D_1 in this connection can at once be replaced by C_1 .

Table III. All possible combinations for the given batch of orders. By way of illustration, the width of the trim is given for each combination, although these values are not used in the computations.

	(D ₁)	C ₁	C ₂	C ₃	C ₄	C ₅	C ₆	C ₇	D ₂	C ₈	C ₉	C ₁₀	D ₃
Order I	(3)	3	2	2	1	1	1	1	—	—	—	—	—
Order II	(—)	—	1	—	3	2	1	—	4	3	2	1	—
Order III	(—)	1	1	4	—	2	4	6	—	2	5	7	9
Width of the trim in mm	(350)	120	220	30	50	90	130	170	150	190	0	40	80

We shall now discuss a formal procedure, from which in addition mathematical evidence is obtained that C₁ is more economical than D₁, although this is visible at a glance.

We cut x₁' , x₂' and x₃' metres in accordance with D₁, D₂ and D₃, and in addition p₁ metres in accordance with C₁. To satisfy order I we should then have:

$$3x_1' + 0x_2' + 0x_3' + 3p_1 = 1800,$$

and similarly for orders II and III. The three linear equations, which are a generalization of system (1), we write symbolically:

$$x_1' \begin{Bmatrix} 3 \\ 0 \\ 0 \end{Bmatrix} + x_2' \begin{Bmatrix} 0 \\ 4 \\ 0 \end{Bmatrix} + x_3' \begin{Bmatrix} 0 \\ 0 \\ 9 \end{Bmatrix} + p_1 \begin{Bmatrix} 3 \\ 0 \\ 1 \end{Bmatrix} = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix},$$

or more concisely:

$$x_1' \{D_1\} + x_2' \{D_2\} + x_3' \{D_3\} + p_1 \{C_1\} = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix} \dots (2)$$

The newly introduced combination C₁ can be expressed in terms of the original combinations D₁, D₂, D₃. Three relations are applicable (one for each order), which we can write in the same way as above:

$$\{C_1\} = a_1 \{D_1\} + a_2 \{D_2\} + a_3 \{D_3\} \dots (3)$$

From Table III it is easily verified that in our case the coefficients a₁, a₂, a₃ are:

$$a_1 = 1, \quad a_2 = 0, \quad a_3 = \frac{1}{9}.$$

Substitution of (3) in (2) with these coefficients gives:

$$(x_1' + p_1) \{D_1\} + x_2' \{D_2\} + (x_3' + \frac{1}{9} p_1) \{D_3\} = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix} \dots (4)$$

Since this is the same system of equations as (1), with different variables, it follows immediately that:

$$\left. \begin{aligned} x_1' + p_1 &= x_1 = 600 \text{ m,} \\ x_2' &= x_2 = 480 \text{ m,} \\ x_3' + \frac{1}{9} p_1 &= x_3 = 150 \text{ m.} \end{aligned} \right\} \dots (5)$$

We have now cut a total of x₁' + x₂' + x₃' + p₁ metres, and this, according to (5), is equal to x₁ - p₁ + x₂ + x₃ - $\frac{1}{9} p_1$ + p₁ = 1230 - $\frac{1}{9} p_1$. The total length of sheet used is thus indeed smaller if we use the combination C₁. How far can we go with this, i.e. how large can we make p₁? This is governed by the condition that x₁' , x₂' , x₃' must always be ≥ 0; obviously we cannot cut negative lengths. By increasing p₁ we reduce, according to (5), the lengths x₁' and x₃' , and x₁' is the first of these to become zero, viz. at p₁ = 600. The sheet length used is then 1230 - 600/9 = 1163 $\frac{1}{3}$ m, a gain of 66 $\frac{2}{3}$ m, and the cutting pattern now obtained — in which D₁ is completely replaced by C₁ — can be written

$$480 \{D_2\} + 83 \frac{1}{3} \{D_3\} + 600 \{C_1\} = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix} \dots (6)$$

In the same way we can examine whether there is any further gain to be obtained by introducing combination C₂, with a given sheet length p₂. For this purpose we consider the equations:

$$x_1'' \{C_1\} + x_2'' \{D_2\} + x_3'' \{D_3\} + p_2 \{C_2\} = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix} \dots (7)$$

Here too {C₂} can be expressed in the combinations used earlier (this is generally possible provided the starting combinations are linearly independent). It is easily verified that:

$$\{C_2\} = \frac{2}{3} \{C_1\} + \frac{1}{4} \{D_2\} + \frac{1}{27} \{D_3\} \dots (8)$$

Substitution in (7) gives:

$$(x_1'' + \frac{2}{3} p_2) \{C_1\} + (x_2'' + \frac{1}{4} p_2) \{D_2\} + (x_3'' + \frac{1}{27} p_2) \{D_3\} = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix} \dots (9)$$

which, by comparison with system (6), yields:

$$\left. \begin{aligned} x_1'' + \frac{2}{3} p_2 &= 600, \\ x_2'' + \frac{1}{4} p_2 &= 480, \\ x_3'' + \frac{1}{27} p_2 &= 83 \frac{1}{3}. \end{aligned} \right\} \dots (10)$$

The total sheet length cut is now:

$$x_1'' + x_2'' + x_3'' + p_2 = 600 - \frac{2}{3}p_2 + 480 - \frac{1}{4}p_2 + \\ + 83\frac{1}{3} - \frac{1}{27}p_2 + p_2 = 1163\frac{1}{3} + \frac{5}{108}p_2.$$

It is seen that the introduction of C_2 gives a poorer result and is therefore pointless.

Incidentally, owing to the limitation of our corrugated-cardboard machine, mentioned at the beginning, we anyhow cannot use C_2 because it is only possible to cut parallel strips of the sheet in two lengths. We can therefore only take combinations of two orders, so that C_5 and C_6 must be scrapped as well as C_2 .

Let us take the procedure one step further. We introduce p_3 metres of C_3 :

$$x_1''' \{C_1\} + x_2''' \{D_2\} + x_3''' \{D_3\} + p_3 \{C_3\} = \\ = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix}, \dots \quad (11)$$

in which we can substitute for C_3 :

$$\{C_3\} = \frac{2}{3}\{C_1\} + 0\{D_2\} + \frac{10}{27}\{D_3\}. \dots \quad (12)$$

We find for the cut length:

$$x_1''' + x_2''' + x_3''' + p_3 = 1163\frac{1}{3} + p_3 - \frac{2}{3}p_3 - \frac{10}{27}p_3 = \\ = 1163\frac{1}{3} - \frac{1}{27}p_3.$$

Here again, therefore, something is to be gained; as a maximum we can choose $p_3 = 225$ m, making $x_3''' = 0$. Thus the combination D_3 is now superseded by C_3 and the solution becomes:

$$450 \{C_1\} + 480 \{D_2\} + 225 \{C_3\} = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix},$$

with a total consumed length of $1163\frac{1}{3} - \frac{225}{27} = 1155$ m.

Continuing in this manner we find, after trying all combinations:

$$194 \{C_1\} + 289 \{C_3\} + 640 \{C_4\} = \begin{Bmatrix} 1800 \\ 1920 \\ 1350 \end{Bmatrix}, \dots \quad (13)$$

with a cut length of 1123 metres.

We must now see whether a further improvement is possible by re-introducing in the last solution one of the other combinations eliminated or rejected in the previous steps. In the present example it is found that this is not the case with any combinations, and therefore solution (13) offers the optimum method of executing the three orders.

In the foregoing we have at various points skimmed over certain questions that have an important place in the theory of linear programming. In particular, it is not self-evident that a "basic solution" can

always be given directly, as in our example using D_1 , D_2 , D_3 . Theory shows that for any batch of m orders such a basic solution making use of not more than m combinations can always be indicated. The theory also answers the question as to how this changes if the m orders have to be carried out not with one but with two machines: in that case the starting basic solution should, in general, consist of $m + 1$ combinations.

We cannot deal here with these theorems and the theoretical foundation of linear programming in general, nor with the graphic interpretation that can be given of the algebraic method outlined above. In this connection reference may be made to the literature ¹⁾.

Let us now take another look at the numerical example. The computations in each attempt to replace one of the tried combinations by another — an iteration — is a simple operation. It is clear, however, that the amount of computing work will accumulate enormously if a large batch of orders has to be dealt with: dealing with m orders involves in each step the handling of systems of m equations with m unknowns. In spite of the methods of simplifying and speeding up the computation, which will presently be dealt with, it is impracticable even for batches of only say 20 orders to calculate the exact optimum solution without the aid of an electronic computer.

The problem in practice; complications

We have seen that the first step in computing for a batch of m orders consists in finding all possible combinations in which two types of pieces (i.e. having two different cut-off lengths) fall within the width of the corrugated-cardboard sheet. It can, however, be immediately decided on mathematical grounds that certain combinations cannot yield an advantage compared with others out of the series. The work of computation can be greatly simplified if we remove these combinations from the series beforehand. We are then left with a number in the order of m^2 combinations.

Moreover, we can speed up the computing procedure considerably if, instead of trying all combinations one after the other, as we did in our numerical example, we choose for the iteration the "best" new combination. This is possible if we express not only the "next" cutting combination of the series in the set of combinations last used, but all other cutting combinations. For this purpose we need the coefficients a_1, a_2, \dots, a_m (see eq. 3) for all these cutting combinations, but these can fairly easily be determined by a manipulation derived from the theory of

linear equations ²⁾. From the numerical example the reader can easily verify (see, for example, eq. 10) that for a cutting combination to yield a gain, the condition must be fulfilled:

$$a_1 + a_2 + \dots + a_m > 1.$$

Now, the extent to which the sum exceeds the value 1 is a measure of the speed at which the trim-losses will decrease if we introduce an increasing length (p) for that cutting combination, and we therefore choose for the next iteration the cutting combination for which $\sum a_i$ is maximum. It is found that on adopting this procedure the optimum solution is reached after a number of iterations in the order of 1.5m.

The computation procedure is accelerated in this way, but opposing this are various complications in the handling of the corrugated-cardboard problem in practice, partly due to extra restrictions and partly to extra tolerances. The most important of these are mentioned below.

The first restriction is bound up with the fact, mentioned in the introduction, that the factory uses two production machines. Normally the aim will be to keep both machines fully occupied. The week's orders must therefore be allocated so that the total lengths of sheet produced by both machines are roughly the same. The permissible difference may be put at, say, 5%.

Another limitation is the need to take into account the time taken for resetting the machines. The change to a different order combination (C), for which a specified length x of sheet is to be produced, makes it necessary to reset all cutting mechanisms, during which process the machine is out of operation. This not only means a loss of production but also increases scrap, owing to a certain length of sheet remaining too long in the drying section of the machine. It is therefore only worth resetting the machines if the length x is not unduly small. Consequently, the strict theoretical solution must also take account of whether such a combination with small x exists, and that combination will have to be eliminated at the cost of a somewhat greater theoretical trim-loss.

It is normally required that the corrugation in the cardboard should run parallel with a particular side of the rectangular pieces to be delivered. This side is then called the "width", and the side at right-

angles to it must be placed in the longitudinal direction of the uncut sheet. In some cases, however, (e.g. for cover sheets) the direction of the corrugation is not important. The length and width of all or part of these pieces in such an order may therefore, if convenient, be interchanged. In the procedure described this means that the number of combinations to be tried will be even larger, and that the appertaining coefficients, i.e. in our example the figures in Table III, may contain *fractions*. Suppose, for instance, that the length and width h are permitted to be exchanged for order *II* in Table I. We then have a new combination, C_{11} , in which the working sheet width (2150 mm) can contain 1 piece of length 500 mm for this order plus 2 pieces for order *I*. As regards its contribution to the number of pieces for order *II* per linear metre of uncut sheet, the "exchanged" piece is equivalent to 800/500 "non-exchanged" pieces for order *II*, so that C_{11} consists of 2 pieces for *I* and $1\frac{3}{5}$ pieces for *II*. Furthermore we obtain a new combination, C_{12} , because 2 exchanged pieces for order *II* can be laid beside 1 non-exchanged piece for the same order; C_{12} thus consists of $4\frac{1}{5}$ pieces for order *II*. Similarly we arrive at combinations C_{13} and C_{14} ; see Table IIIa.

Table IIIa. Supplementary to Table III for the case where the length and width of the pieces for order *II* can be exchanged.

	C_{11}	C_{12}	C_{13}	C_{14}
Order <i>I</i>	2	—	—	—
Order <i>II</i>	$1\frac{3}{5}$	$4\frac{1}{5}$	$3\frac{1}{5}$	$1\frac{3}{5}$
Order <i>III</i>	—	—	2	5
Width of the trim	150 mm	50 mm	90 mm	200 mm

A tolerance of considerable importance in planning concerns the fixing of the required *number* of pieces in each order. An arrangement can be made with customers to allow the factory a margin of say $\pm 5\%$ in the number to be delivered. This allows a refinement of the linear programming procedure, known as an "upper bound technique". This makes it possible, within the prescribed limits, to fix the number of pieces for each order in such a way as to minimize the ultimate total percentage of scrap. In the computing work this means carrying out some extra iterations.

The programme for the computer

We will not go any deeper into the computing procedure, and will conclude by touching briefly on the composition of the computer programme.

The orders are supplied to the machine in code form on punched tape. The code contains data on the width, the length and number of pieces, whether

²⁾ The manipulation consists in first forming an inverse matrix from the numbers of the basis combinations (cf. Table II). The coefficients a_1, a_2, \dots for each cutting combination are then found by simply multiplying the numbers of that combination by the rows of the inverse matrix.

or not the width and length are interchangeable, and whether there is any permissible margin in the number of pieces, and the machine provides each order with a serial number. If the pieces in two orders happen to be the same size, the machine at once places them together under one serial number.

The first part of the computer programme now consists in producing the list of possible combinations (pairs) of orders that can fit into the available sheet width. (In some cases, where the same lengths are required for different orders, *three* or more widths can be combined; the machine can be programmed to produce these combinations also.) In orders where the dimensions are interchangeable, the lengths are also used for the formation of pairs. Combinations that on mathematical grounds will not occur in the final solution, are immediately removed from the list by the machine, and so are combinations that are undesirable for technical reasons (e.g. with a view to subsequent processing). A list of combinations is prepared for the machine which delivers the broad sheet, and another list for the machine delivering the narrow sheet.

The second part of the computing programme consists in producing a basic solution, similar to $D_1-D_2-D_3$ in the example discussed. This is followed by the "iterations", in each of which the variation of the introduced lengths p is subjected to the condition that none of the lengths x should be negative, but also to the condition that the total lengths of the two uncut sheets should not differ by more than 5%. After an optimum solution has been found, the part

of the programme involving the "upper bound technique" is effected, which determines the final optimum, with permissible deviations in the numbers of pieces per order.

The last part of the programme is concerned with eliminating the possible combinations with cut lengths *lower* than a certain limit. Mathematically speaking, this step upsets everything that has been done before, since the final solution will as a rule no longer be optimum. An exact solution, however, would take too long to compute, and it may safely be assumed that for all practical purposes the discrepancy will not be of much consequence.

Finally, to provide some idea of what the results of a computation by PASCAL look like, we give in *Table IV* an example of a batch of nine orders 1-9, including four "interchangeable" orders, and in *Table V* the distribution as computed and printed out by the PASCAL. In orders 3 and 8 the length

Table IV. Example of a batch of nine orders supplied to the PASCAL for processing.

Order number	Width mm	Length mm	Number of pieces required	Inter-changeable (1 = yes)
1	1580	1414	365	—
2	1500	1290	4080	1
3	1475	1400	638	1
4	735	1500	3367	—
5	540	1455	850	—
6	530	1350	1275	—
7	500	1520	6734	—
8	359	900	5360	1
9	182	1440	3535	1

Table V. Result of the computation for the orders in Table IV as printed out by the PASCAL. The first column shows that seven combinations are cut from the wide sheet and three from the narrow sheet. The third column indicates for each combination how many pieces of the respective orders are laid side by side in the uncut sheet. This number, times the number of "cuts", gives the number of pieces. In the last column "fin" is printed if the pieces obtained from the combination complete the total number required in an order.

Machine, combination number	Order number	Number of pieces in combination	Number of pieces	Length mm	Number of cuts	Inter-changed (1 = yes)	Width mm	Number of metres	Finished
Wide 1	2	1	2007	1500	2007	1	1290	3011	
	4	1	2007	1500	2007	—	735		
Wide 2	5	4	852	1455	213	—	2160	310	fin
	1	1	365	1414	365	—	1580	517	fin
6	1	383	1350	383	—	530			
Wide 4	6	1	261	1350	261	—	530	352	
	9	9	2205	1440	245	—	1638		
Wide 5	6	4	632	1350	158	—	2120	214	fin
Wide 6	8	2	5364	359	2682	1	1800	963	fin
	9	2	1336	1440	668	—	364		fin
Wide 7	3	1	638	1475	638	1	1400	941	fin
	4	1	627	1500	627	—	735		
Narrow 1	2	1	2073	1290	2073	—	1500	2675	fin
Narrow 2	7	3	6735	1520	2245	—	1500	3412	fin
Narrow 3	4	2	734	1500	367	—	1470	550	fin

should be taken as the "width", and in order 2 some of the pieces should be produced with the length and breadth unchanged and some interchanged. The total length of sheet produced by the wide machine is seen to be 6310 m, that of the narrow machine 6638 m (difference about 5%). The total area produced is 23 818 m². The number of pieces produced differs very little from the number ordered (the greatest discrepancy is 6 pieces, a consequence

of rounding-off in the calculation); they cover altogether a total area of 23 103 m². The trim-losses thus come to 715 m², i.e. 3.0%.

The actual time taken by the PASCAL to compute this batch of orders was about 11 seconds. Ordinarily the batches of orders are larger. The computing time for a normal batch of, say, 30 orders is about 60 seconds. Feeding-in the computing programme and batch data takes about 25 seconds.

Summary. A machine in a corrugated-cardboard factory delivers a continuous sheet of cardboard with a fixed width of e.g. 1.5 or 2 metres. When this is cut into individual pieces, of different sizes to comply with different orders, losses result from the fact that full use is not made of the whole width of the sheet. An attempt can be made to fit several pieces into the available width in such a way as to minimize the total trim-losses. Finding the optimum cutting arrangement for a given batch of orders is a weekly recurring problem in the corrugated-cardboard factory. The problem belongs to a large category of problems found in widely diverse fields of technology and economics. The mathematical technique devised for handling such problems — linear programming — is explained in this article with the aid of a numerical example. If the batch of orders is at

all large, e.g. 20 or more orders, the computing work is feasible only with the aid of an electronic computer. As an aid to the production office of Philips' corrugated-cardboard factory, the computations have been programmed for the PASCAL. At present every week, by way of experiment, PASCAL computes several batches of orders, each for a specified quality of cardboard. The problem is complicated by various extra requirements, e.g. that the capacity of the two production machines available in the factory should be fairly uniformly utilized, and by the fact that there are some tolerances, in particular that the width and length of pieces in some orders are interchangeable, and that some margin is allowed in the numbers of pieces supplied in a given order. These refinements are included in the programme.

CALCULATION OF POTENTIAL FIELDS AND ELECTRON TRAJECTORIES USING AN ELECTRONIC COMPUTER

by C. WEBER *).

518.5:537.213

Introduction

In electron optics it is important to know the position of the electrons in the space between electrodes of given configuration and given potentials as a function of time. The problem can be divided into two parts: determination of the potential field between the electrodes, and the solution of the equation of motion of the electrons in this field. The solutions can be found by various methods.

- 1) *Purely mathematical calculation of the potential field and the trajectories.* This method is seldom applied; for there are not many electrode systems in which the potential field can be calculated without using numerical methods, and only in a few of these systems are the electron trajectories amenable to computation.
- 2) *The rubber membrane.* A horizontally stretched rubber sheet, deformed to follow the shape and potentials of the electrodes and over which steel balls are rolled, is a mechanical analogue of the electrical problem for "two-dimensional" cases ¹⁾. The method is very convenient but not highly accurate. The influence of space charge can only be determined with considerable difficulty ²⁾.
- 3) *Conductive paper.* The potential of "two-dimensional" fields can also be determined using a sheet of conductive graphite paper ³⁾, on which the outlines of the electrodes are painted with conductive paint. Here again, it is difficult to take into account the space charge effects.
- 4) *The electrolytic tank.* This is another analogue system for determining potential fields, and can be used with a second analogue system for automatically recording the electron trajectories ⁴⁾. The electrolytic tank is a valuable tool for electron

optics, as it can be employed for determining electron trajectories in a wide variety of electrode systems. The tank method shares the drawback of the rubber membrane that space charge effects can only be calculated by a cumbersome procedure.

Another disadvantage of the tank method is that it takes a lot of time to make the individual models of the electrodes. The process is quicker using a "wedge tank" ⁵⁾, simulating such a small sector of the electrodes — which have to be rotational-symmetric — that the models can be made flat in one direction. The trajectories near to the axis (paraxial trajectories) cannot be determined accurately in a tank, and it is precisely these trajectories that are often particularly important.

- 5) *The resistance network.* This is yet another analogue system ⁶⁾, and along with the electrolytic tank is an important method of determining potentials. There are networks for solving two-dimensional problems and others for rotational-symmetric three-dimensional problems. The method is both quick and reasonably accurate. Space charge effects can be determined by applying appropriate currents to the junctions of the network. To obtain sufficient accuracy this should be done at a fairly large number of junctions, which of course reduces the convenience of the method.

Starting from the potentials on the axis of the network with rotational symmetry, the paraxial electron trajectories can be determined with the aid of an analogue computer ⁷⁾.

- 6) *Numerical methods.* Until recently the numerical determination of potential distributions was not widely used owing to its relatively slow convergence to a solution. It is only since the advent of fast electronic computers that this method has

*) Philips Research Laboratories, Eindhoven.

¹⁾ P. H. J. A. Kleijnen, The motion of an electron in two-dimensional electrostatic fields, Philips tech. Rev. 2, 338-345, 1937. See also the photo in Philips tech. Rev. 13, 16, 1951/52, and figs. 5 and 6 in Philips tech. Rev. 14, 122, 1952/53.

²⁾ G. Alma, G. Diemer and H. Groendijk, A rubber membrane model for tracing electron paths in space charge fields, Philips tech. Rev. 14, 336-344, 1952/53.

³⁾ W. Clausnitzer and H. Heumann, Ausmessung elektrischer Felder mit Hilfe von halbleitenden Schichten, Z. angew. Phys. 2, 443-446, 1950.

⁴⁾ J. L. Verster, An apparatus for automatically plotting electron trajectories, Philips tech. Rev. 22, 245-259, 1960/61. See also the literature referred to in note ¹⁾ of that article.

⁵⁾ M. Bowman-Manifold and F. H. Nicoll, Electrolytic field plotting trough for circularly symmetric systems, Nature 142, 39, 1938.

⁶⁾ J. C. Francken, The resistance network, a simple and accurate aid to the solution of potential problems, Philips tech. Rev. 21, 10-23, 1959/60.

⁷⁾ A. J. F. de Beer, H. Groendijk and J. L. Verster, The plotting of electron trajectories with the aid of a resistance network and an analogue computer, Philips tech. Rev. 23, 352-362, 1961/62 (No. 11).

become feasible. A great advantage is that it can be used without much extra trouble for solving Poisson's equation in electrostatics, which relates to cases involving space charge.

The method adopted in a given case depends largely on the nature of the problem. The choice is influenced by considerations such as the accuracy required, whether or not space charge can be disregarded, the speed at which the final result can be obtained and the shape of the electrodes. If an electronic computer is available the numerical method is as a rule preferable for cases involving space charge.

Calculating space-charge effects entails solving a more complicated equation and in addition presents

be used for calculating the space charge from the electron trajectories, in which case the entire calculation is performed mechanically.

Of this extensive problem: the calculation of (a) potential fields, (b) electron trajectories and (c) space charge, in this article we shall consider only the first two parts, (a) and (b), on which work has now been completed. The calculations were carried out from the middle of 1959 to the middle of 1961 in Philips Computing Centre⁸⁾ on an electronic computer, type 650, supplied by the International Business Machines Corporation (*fig. 1*). A wide range of problems in which the space charge is known or may be assumed to be zero can be solved by the method about to be described.



Fig. 1. Two IBM 650 electronic computers in Philips Computing Centre at Eindhoven. A_1 and A_2 are the actual computers, B_1 and B_2 the punched-card machines, and C_1 and C_2 the power units. The line printer D serves for both installations.

the following difficulty. The space charge occurring in Poisson's equation cannot be calculated until the trajectories are known, but to calculate the trajectories it is necessary to solve Poisson's equation, and thus to know the space charge. To get out of this impasse we adopt a *method of iteration*. We begin by estimating the space charge and then use the estimated values for solving Poisson's equation; from the potentials found we calculate the trajectories, and using these correct the space charge; we then solve Poisson's equation once more, and so on. This iterative process converges to a solution fairly quickly²⁾. Given an electronic computer, it can also

In the meantime the electronic computer PASCAL, developed at Philips, came into use⁹⁾, and a start was recently made in the Computing Centre on programming parts (a) and (b) and also (c) on this much faster machine. It will soon be possible to begin calculating potential fields, electron trajectories and space charge distributions with the PASCAL computer.

⁸⁾ A. J. W. Duijvestijn advised on the mathematics involved, and the programming was carried out by H. B. Nota, both of Philips Computing Centre.

⁹⁾ W. Nijenhuis, The PASCAL, a fast digital electronic computer for the Philips Computing Centre, Philips tech. Rev. 23, 1-18, 1961/62 (No. 1).

There are cases in which a magnetic as well as an electric field has to be taken into account. This is so, for example, in travelling-wave tubes, for which the electron trajectories have also been calculated¹⁰. The potential field in this case is found in the same manner as described below, but the process of calculating the trajectories and the space charge is entirely different.

Conversion of Poisson's differential equation into a difference equation

The electric potential V of a system containing a space charge ρ , which is a position function, is given by Poisson's differential equation:

$$\frac{\partial^2 V}{\partial x^2} + \frac{\partial^2 V}{\partial y^2} + \frac{\partial^2 V}{\partial z^2} = -\frac{\rho}{\epsilon_0}, \quad \dots \quad (1)$$

where ϵ_0 is the dielectric constant of free space.

Since we are concerned in this article solely with cases possessing rotational symmetry, we make the rectangular coordinates x, y and z into cylindrical coordinates r, φ and z (fig. 2), taking the axis of rotation as the z axis. V is then dependent on φ , so that the derivatives with respect to φ are zero. Eq. (1) is now:

$$\frac{\partial^2 V}{\partial r^2} + \frac{1}{r} \frac{\partial V}{\partial r} + \frac{\partial^2 V}{\partial z^2} = -\frac{\rho}{\epsilon_0}. \quad \dots \quad (2)$$

For solving this kind of problem numerically, use can best be made of *difference equations*. For this reason we derive from the differential equation (2), which gives a relation between certain derivatives of V , a difference equation which describes the relation between the potential at closely adjacent points. For this purpose we express the derivatives of V at a point P_0 in terms of the potential V_0 of P_0 and in terms of the potentials of nearby points whose distances h to P_0 are so small that we can permissibly

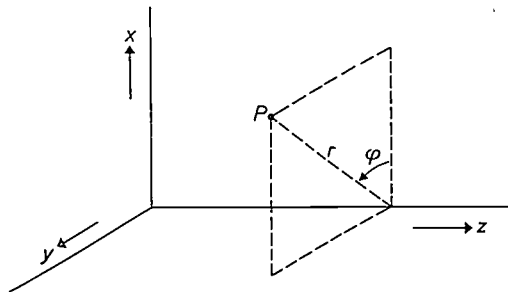


Fig. 2. Point P with rectangular coordinates x, y and z , and with cylindrical coordinates r, φ and z .

neglect the terms to the third and higher powers of h in the series presently to be given.

We distinguish between two cases: P_0 not on the z axis, and P_0 on the z axis.

Points not on the z axis

For points P_0 not on the z axis it is sufficient to take four neighbouring points, P_1, \dots, P_4 , situated at distances h_1, \dots, h_4 from P_0 (fig. 3). As will be seen later, it is possible in many cases to choose these

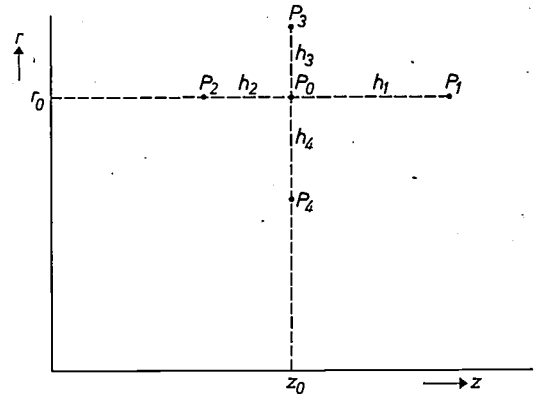


Fig. 3. Point P_0 at a distance r_0 from the z axis (the axis of rotational symmetry) surrounded by points P_1, P_2, P_3 and P_4 at distances h_1, \dots, h_4 from P_0 .

four distances equal to one another; this is not possible in all cases however so we shall derive here the difference equation for the more general case of dissimilar distances h .

We expand V into a Taylor series around the point P_0 , of which the coordinates are r_0 and z_0 :

$$V(r,z) = V(r_0, z_0) + \frac{1}{1!} \left[\left(\frac{\partial V}{\partial r} \right)_0 (r-r_0) + \left(\frac{\partial V}{\partial z} \right)_0 (z-z_0) \right] + \frac{1}{2!} \left[\left(\frac{\partial^2 V}{\partial r^2} \right)_0 (r-r_0)^2 + 2 \left(\frac{\partial^2 V}{\partial r \partial z} \right)_0 (r-r_0)(z-z_0) + \left(\frac{\partial^2 V}{\partial z^2} \right)_0 (z-z_0)^2 \right] + \dots \quad (3)$$

The coordinates of P_1, \dots, P_4 are respectively $(r_0, z_0+h_1), (r_0, z_0-h_2), (r_0+h_3, z_0)$ and (r_0-h_4, z_0) . From (3) we find for the potentials V_1, \dots, V_4 at these points:

$$\left. \begin{aligned} V_1 &= V_0 + \left(\frac{\partial V}{\partial z} \right)_0 h_1 + \frac{1}{2} \left(\frac{\partial^2 V}{\partial z^2} \right)_0 h_1^2 + \dots \\ V_2 &= V_0 - \left(\frac{\partial V}{\partial z} \right)_0 h_2 + \frac{1}{2} \left(\frac{\partial^2 V}{\partial z^2} \right)_0 h_2^2 - \dots \end{aligned} \right\} \quad (4)$$

and similarly for V_3 and V_4 .

If the points P_1, \dots, P_4 are so close to P_0 that the terms containing third and higher powers of h_1, \dots, h_4 may be neglected, the first and second

¹⁰ By E. Deimel, Electron Tubes Division, Philips Eindhoven. See: Mikrowellenröhren, Vorträge der Internationalen Tagung in München 7-11 June 1960, pp. 493-507; Vieweg, Brunswick 1961.

partial derivatives of V to z and r are found from (4) to be:

$$\left. \begin{aligned} \left(\frac{\partial V}{\partial z}\right)_0 &= \frac{h_2}{h_1(h_1+h_2)}V_1 - \frac{h_1}{h_2(h_1+h_2)}V_2 + \frac{h_1-h_2}{h_1h_2}V_0, \\ \left(\frac{\partial^2 V}{\partial z^2}\right)_0 &= \frac{2}{h_1(h_1+h_2)}V_1 + \frac{2}{h_2(h_1+h_2)}V_2 - \frac{2}{h_1h_2}V_0, \\ \left(\frac{\partial V}{\partial r}\right)_0 &= \frac{h_4}{h_3(h_3+h_4)}V_3 - \frac{h_3}{h_4(h_3+h_4)}V_4 + \frac{h_3-h_4}{h_3h_4}V_0, \\ \left(\frac{\partial^2 V}{\partial r^2}\right)_0 &= \frac{2}{h_3(h_3+h_4)}V_3 + \frac{2}{h_4(h_3+h_4)}V_4 - \frac{2}{h_3h_4}V_0. \end{aligned} \right\} \dots (5)$$

In this way we have expressed the derivatives of V at point P_0 in the potential of P_0 and in the potentials of the neighbouring point P_1, \dots, P_4 . Substitution of (5) in the differential equation (2) changes the latter into the following difference equation:

$$\frac{2}{h_1(h_1+h_2)}V_1 + \frac{2}{h_2(h_1+h_2)}V_2 + \frac{2r_0+h_4}{h_3(h_3+h_4)r_0}V_3 + \frac{2r_0-h_3}{h_4(h_3+h_4)r_0}V_4 = -\frac{\rho_0}{\epsilon_0} + \left[\frac{2}{h_1h_2} + \frac{2r_0+h_4-h_3}{h_3h_0r_0} \right] V_0, \dots (6)$$

where ρ_0 is the space charge at the point P_0 .

If points P_1, \dots, P_4 can be taken as being at equal distances from P_0 ($h_1 = h_2 = h_3 = h_4 = h$), we can simplify (6) to

$$V_1 + V_2 + \left(1 + \frac{h}{2r_0}\right)V_3 + \left(1 - \frac{h}{2r_0}\right)V_4 = -\frac{\rho_0 h^2}{\epsilon_0} + 4V_0. \dots (7)$$

Points on the z axis

Equations (6) and (7) do not apply on the z axis, where r_0 is zero and thus the terms with $1/r_0$ are infinitely large. To derive a difference equation that does hold for points on the z axis we return to Poisson's differential equation in its original form, (1). We now dispose six points around P_0 as indicated in fig. 4. Having regard to the rotational symmetry we set the points P_3, \dots, P_6 at equal distances (h_4) from P_0 , so that the potentials of these four points are identical ($V_3 = V_4 = V_5 = V_6$). By a procedure similar to that described in the previous section we then arrive at the following difference equation:

$$\frac{2}{h_1(h_1+h_2)}V_1 + \frac{2}{h_2(h_1+h_2)}V_2 + \frac{4}{h_3^2}V_3 = -\frac{\rho_0}{\epsilon_0} + \left(\frac{4}{h_3^2} + \frac{2}{h_1h_2} \right) V_0. \dots (8)$$

In cases where h_1 and h_2 are equal to h_3 ($= h$), equation (8) simplifies to:

$$V_0 + V_2 + 4V_3 = -\frac{\rho_0 h^2}{\epsilon_0} + 6V_0. \dots (9)$$

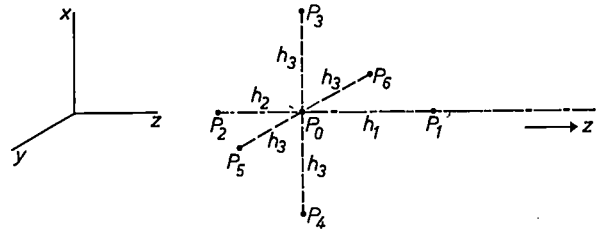


Fig. 4. Point P_0 on the z axis, surrounded by the six points P_1, \dots, P_6 at distances h_1, h_2 and h_3 from P_0 (because of the rotational symmetry P_4, P_5 and P_6 are at the same distance h_3 from P_0 as P_3).

Boundary conditions

In order to determine the potential field we must also know the boundary conditions, that is we must know the potential at a closed "boundary" (the surface of the electrodes). To calculate the potential field between, for example, two concentric conducting spheres, it is necessary to know the potentials both of the inner and the outer sphere.

The method of iteration

For the purpose of calculating a potential field with the aid of the difference equations we introduce in the $r-z$ plane, parallel with the coordinate axes, a close-meshed grid consisting of squares of mesh width h . As a simple example, fig. 5 shows the situation for the case just mentioned, where the potential field is to be calculated between two concentric spheres, i.e. the potential at each of the intersection points of the grid between the boundaries (between the circles on which the spheres cut the $r-z$ plane; in view of the symmetry it is sufficient here to take semicircles). We number the total of N intersection points that lie between (not on) these boundaries from 1 to N , as indicated for a number of these points in fig. 5. The intersection points that happen to lie exactly on a boundary are marked with a cross.

Some of the points $i = 1, \dots, N$ lie in the r or the z direction at a distance less than h from the boundary (e.g. points 111, 228, 260, 315). As far as these points are concerned, not all distances to the surrounding points P_1, \dots, P_4 can be equal, and we must therefore use equation (6) or (8) applicable to dissimilar distances h . By far the most points $1, \dots, N$, however, are so far removed from the boundary that the distances h can be taken as iden-

tical, thus permitting the use of the simplified formulae (7) or (9).

We begin by filling in arbitrary values for the potentials at points 1, . . . , N of the grid, although of course these values will not generally satisfy the difference equation. For the potential at point 1 we

large number of points in the manner described. It is therefore reasonable to look for a method to make the process converge faster, thus reducing the number of cycles.

Consider point i of the grid. We denote the potential found for this point after k cycles by $V_i^{(k)}$, and

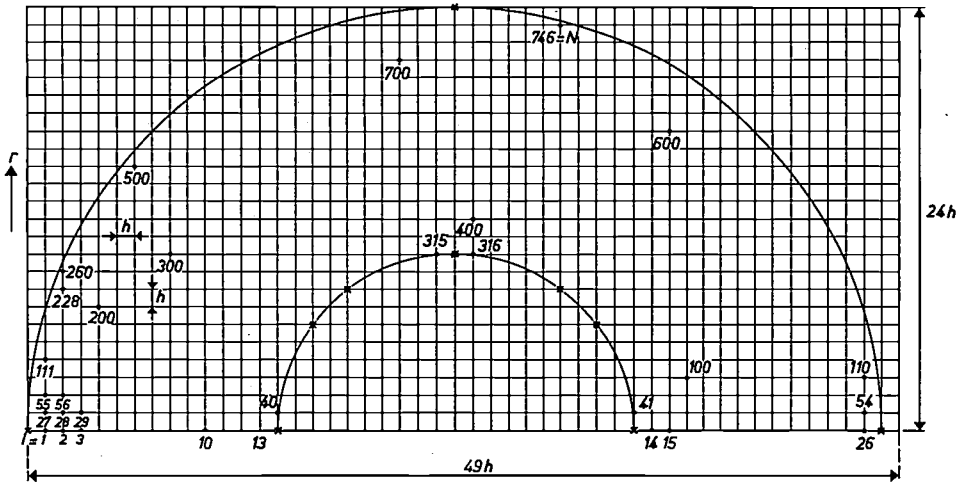


Fig. 5. The r - z plane is provided with a grid of squares of mesh width h . For computing the potential at the grid points the cross-section of the electrodes — here two concentric spheres — is drawn in the r - z plane (semicircles only, in view of rotational symmetry). The grid points between the semicircles are numbered in order of the horizontal rows. The grid points that lie on one of the semicircles are marked with a cross.

can now find a value that does satisfy the difference equation, by substituting for the arbitrary value the one that can be calculated by means of the difference equation from the potentials at the surrounding points. Having thus found a better value for the potential at point 1, we make the potential of point 2 satisfy the difference equation in the same manner. This means, it is true, that the potential at point 1 no longer satisfies the equation, but this will presently be remedied. We proceed in the same manner for points 3 up to N .

After completing this cycle we have a situation in which most of the potentials found do not yet satisfy the difference equation. We therefore begin, again at point 1, on a second cycle, then on a third, and so on. The more often we do this the better will the potentials meet the equation, and thus the smaller will be the corrections to be applied. It can be proved that with each cycle the values converge closer towards a final solution, and that this solution is independent of the initial values chosen and satisfies the difference equation at all points.

Speeding up the convergence by successive over-relaxation

Even with an electronic computer it takes a considerable time to calculate the potential at a suitably

that found for $k + 1$ cycles by $V_i^{(k+1)}$. As k increases, these potentials approach the end value V_i (fig. 6). At a given stage of the $(k + 1)$ st cycle we have calculated $V_i^{(k+1)}$, while the previous value, $V_i^{(k)}$, is still noted for the point i . If, for example, $V_i^{(k+1)}$ is larger than $V_i^{(k)}$ it is reasonable to assume that the still unknown value $V_i^{(k+2)}$ will likewise be larger than $V_i^{(k+1)}$. Accordingly — anticipating the

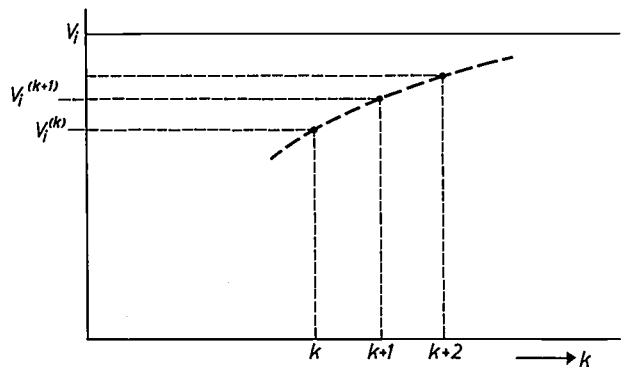


Fig. 6. The potential computed for grid point i approaches the value V_i after an increasing number of iterative cycles. After k cycles the value $V_i^{(k)}$ is found. If, after $k + 1$ cycles, for example, the value $V_i^{(k+1)}$ was greater than $V_i^{(k)}$, it is probable that $V_i^{(k+2)}$ would be greater still after $k + 2$ cycles. The convergence is accelerated by substituting for $V_i^{(k)}$ not $V_i^{(k+1)}$ but a somewhat greater value, in accordance with eq. (10) (successive over-relaxation).

development — instead of substituting $V_i^{(k+1)}$ for $V_i^{(k)}$ we at once take a somewhat larger value, viz:

$$V_i^{(k)} + \omega \left[V_i^{(k+1)} - V_i^{(k)} \right], \dots \quad (10)$$

where the factor ω is greater than 1. This method is referred to as successive over-relaxation.

The value of ω must be carefully chosen. With $\omega = 1$ we again have the situation where the value $V_i^{(k+1)}$ is inserted. If we make ω too large, the next value will be farther away from V_i instead of being nearer to it, and the speed of convergence will be lower instead of higher. There is a specific optimum value of ω at which the process converges fastest. It can be demonstrated that this optimum value depends only on the configuration of the electrodes, and not on the potentials initially adopted. It can also be proved that the best value of ω is independent of the cycle number k ¹¹⁾.

For certain simple configurations the optimum value of ω can be calculated (see appendix). From this we can estimate the optimum value for less simple configurations.

Successive overrelaxation can speed up the iteration process considerably. If, for example, the boundary is a cylinder of length $49h$ and radius $24h$, the iteration process, using the optimum value of ω , is shortened by a factor of 24. When we mention that even then the computer took more than $1\frac{1}{2}$ hours to solve our problem, it will be evident just how substantial this saving is.

Calculation of potential fields using the IBM 650

The IBM 650 computer, in Philips Computing Centre, which was used for calculating potential fields and electron trajectories, has a single magnetic drum type memory for 2000 "words", capable of storing both numbers and instructions¹²⁾. The computer works on the decimal system. A word consists of ten digits.

In accordance with the requirements of the computer for the kind of problem involved we take a fixed grid of 50 points (49 meshes) in the z direction and 25 points (24 meshes) in the r direction, giving a total of $50 \times 25 = 1250$ points. The boundary with the given potentials must lie within the rectangle of 49×24 meshes (fig. 5). Each point of the grid is given a corresponding word in the memory. The information contained in each word has to comprise the potential V and the space charge ρ at the relevant

point, together with an instruction I , to which we shall return presently.

The calculation of V is performed with an accuracy of four digits, ρ is also given in four digits, while two digits are required for I . In the word (ten digits) the first two places are used for I , the next four for V , and the last four for ρ (fig. 7).

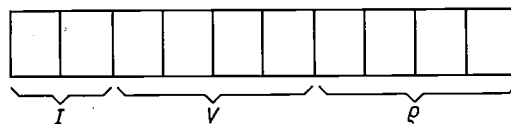


Fig. 7. A "word" of the IBM 650 computer consists of ten digits. The first two digit positions are used for the instruction I , the next four for the potential V , and the last four for the space charge ρ .

The machine has to distinguish between the following five kinds of points:

- a) Points with a fixed potential; these are boundary points and points outside the boundary.
- b) Points on the z axis to which eq. (9) applies.
- c) Points on the z axis to which eq. (8) applies, i.e. where at least one of the distances h_1, h_2, h_3 differs from h .
- d) Points not on the z axis to which eq. (7) is applicable.
- e) Points not on the z axis to which eq. (6) applies, i.e. where at least one of the distances h_1, \dots, h_4 differs from h .

This distinction is made using the number stored in the instruction part (98, 88 or 89), in combination with the values of r_0 ($= 0$ or > 0):

	$r_0 = 0$	$r_0 > 0$
$I = 98$	kind (a)	kind (a)
$I = 88$	kind (b)	kind (d)
$I = 89$	kind (c)	kind (e)

Before starting the machine on the computations we ensure that all initial data have been fed in, i.e. that the instruction part and the space charge part of each word contain the correct values. On the boundaries the potential must have the prescribed values; arbitrary figures are inserted for the potentials in the other points.

The computation procedure is represented by the block diagram in fig. 8. The machine successively scans all points of the grid, beginning bottom left, then passing along the z axis to the right, then from left to right along all points one row higher, and so on. If it is the turn of point j , for example¹³⁾,

¹¹⁾ For practical reasons a fixed number ω is used in all cases which is the same for all points (ω is thus also independent of i).

¹²⁾ An explanation of these terms will be found in the article quoted in ⁹⁾.

¹³⁾ The letter i used in the foregoing denoted one of the N points of the grid where the potential had to be computed, and thus excluded points on and outside the boundary (fig. 5). The computer, however, examines successively all 1250 points of the grid; to make a distinction here the letter j is used.

block 3 is interrogated to discover to which of the five kinds this point belongs. Depending on the result of the interrogation, one of the blocks 4a, . . . , 4e is selected (the letters a, . . . , e correspond to the above division into five kinds). For points with a fixed potential, block 4a does not change the potential at the memory location of point j . Blocks 4b, . . . , 4e, on the other hand, replace the potential at the memory location j by another one, calculated from the relevant difference equation using the method of successive over-relaxation in accordance with eq. (10). The arithmetic unit can carry out simple mathematical operations such as adding, subtracting, multiplying and dividing, and numbers can also be transported from the memory to the arithmetic unit

and vice versa. Each operation, e.g. addition or transport, has its own number code. The programme consists of a succession of code numbers which establish the correct order of the operations.

Let us now examine the programme section, represented by block 4b. First of all, starting from the number j , it is necessary to compute, by adding and subtracting operations, the memory locations of the potentials V_1, V_2, V_3 and V_0 which occur in eq. (9). Next, the potential part has to be separated from the words contained in these memory locations. The space-charge part pertaining to point j must also be split off. From these potentials and from the space charge the new potential is now computed. This is done by additions and multiplications, carried out

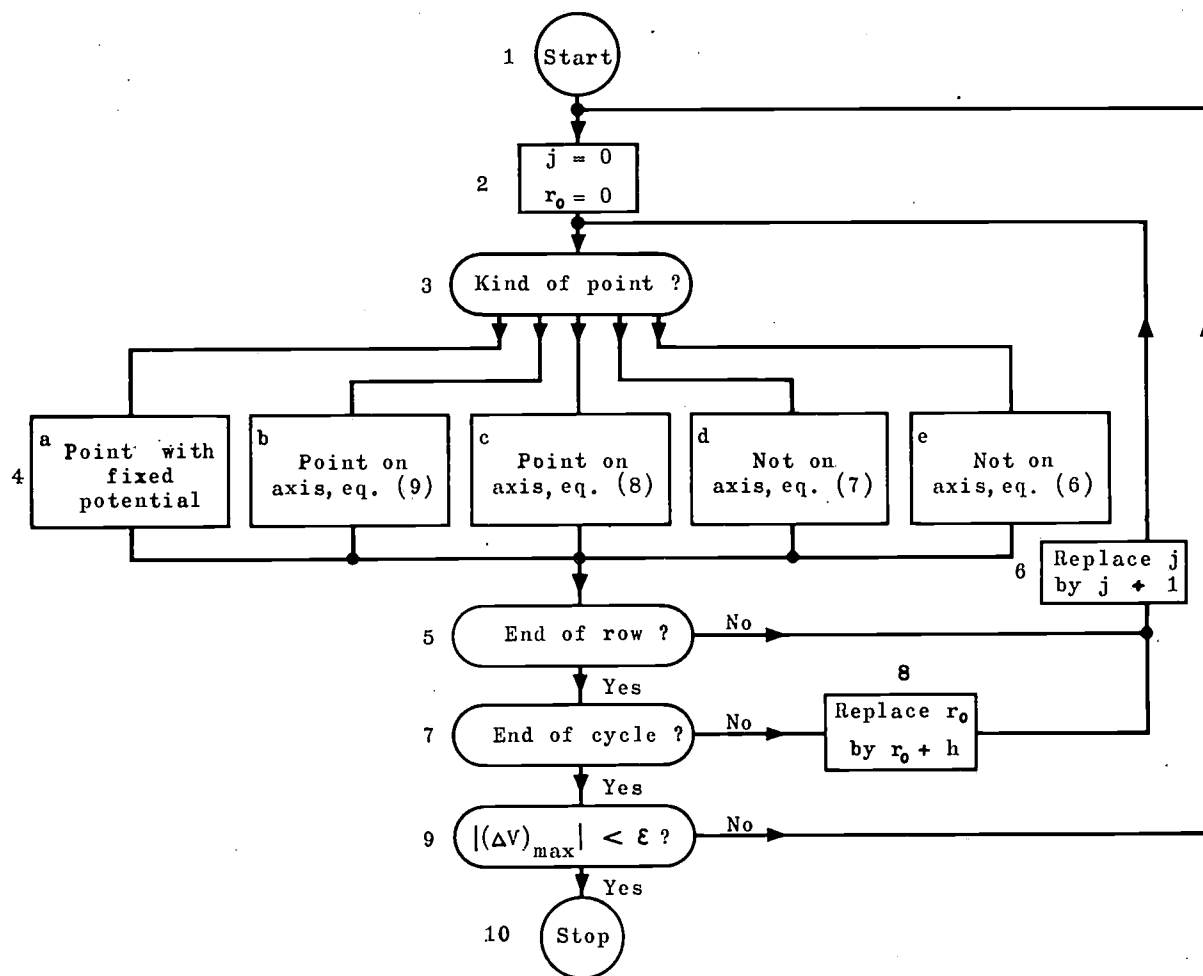


Fig. 8. Block diagram representing the calculation of a potential field on the IBM 650. 1 start. In block 2, quantities j and r_0 are made equal to their initial values. Block 3 ascertains to which of the five kinds a, . . . , e point j belongs. Depending on the result of this investigation, block 4a retains the potential of the relevant point, or one of the blocks 4b, . . . , 4e modifies the potential in agreement with the appropriate formula. Next, block 5 examines whether the point is at the end of a row. If it is not, a start is made on the next point, for which purpose block 6 receives the instruction: "replace j by $j + 1$ ". If the point was the last one on a row, block 7 investigates whether the end of a cycle has been reached. If it has not, the next row (block 8) is started; if it has, block 9 investigates whether $|(\Delta V)_{\max}|$ is now smaller than ϵ . On this depends whether the machine should proceed with the computation or should stop (block 10).

step by step one after the other. Finally the potential part of the word pertaining to point j must be replaced by the new potential, and the resulting new word stored in the relevant memory location. Once this has been done the computer proceeds with the part of the programme formed by block 5. At point $j + 1$ the same part of the programme is repeated, with the potentials and space charge pertaining to this new point.

Block 5 investigates each point to see whether it lies at the end of a horizontal row. If this is not so in the case of point j , the same procedure is applied at point $j + 1$ of the row; block 6 ensures that j is replaced by $j + 1$. If, however, j lies at the end of a row, the calculation follows a different course in the block diagram: block 7 then ascertains whether the end of the cycle has been reached. If this is not the case, the next row is started upon, for which purpose r_0 in the difference equation must first be increased by one unit h (block 8). The next row is then started, and so on.

If the end of the cycle had in fact been reached, the machine then starts on the following cycle, beginning with the point at the bottom left corner (fig. 5), after j and r_0 have been made equal to zero. In this way, however, the machine would never stop computing. To have some criterion as to whether the iterative process has gone far enough we let the machine calculate for every point not merely the new potential but also the difference between the old and the new one. This difference $(\Delta V)_j$ is computed in blocks 4b, . . . , 4e (for the sake of clarity we have omitted this in discussing the procedure in block 4b). The absolute value of $(\Delta V)_j$ is stored in a permanent location in the memory. After this the difference between the old and the new potential is computed for the next point, the machine determines whether the absolute value $|(\Delta V)_{j+1}|$ of this difference is larger or smaller than $|(\Delta V)_j|$. The larger of the two amounts $|(\Delta V)_j|$ and $|(\Delta V)_{j+1}|$ now occupies the location in the memory where $|(\Delta V)_j|$ was stored. At the end of the cycle the memory contains the largest potential correction, $|(\Delta V)_{\max}|$, found in this cycle. After a large number of cycles the value of $|(\Delta V)_{\max}|$ will finally approach zero. At the end of each cycle a predetermined amount ϵ is deducted in block 9 from $|(\Delta V)_{\max}|$. If the result is positive or zero, the machine then starts on a new cycle; if the result is negative, the machine stops.

If we make ϵ equal to 1 or 2 units of the last digit location, we find that $|(\Delta V)_{\max}|$ does not usually reach such a small value. This is a consequence of the over-relaxation, for when the iterative process is almost completed the rounding-off errors become

much more significant. Successive over-relaxation is applied to these errors, with the result that the process no longer converges. For this reason we choose to make ϵ equal to about 5 units of the last digit location. When $|(\Delta V)_{\max}|$ has reached this greater value ϵ , we set the computer in such a way as to make the factor ω in (10) equal to 1. The computer now works without over-relaxation, and we let it go on computing for a while, thus improving the accuracy of the end result.

It is important that the computation should be fast especially in those blocks that are frequently repeated, such as block 4d, which computes the potential in the most numerous kind of points (d).

Another thing to be considered is the scale adopted. As mentioned, four digital locations are available for the potential, and it is important that all these four locations should be used. This implies that the unit of potential should not be so ineptly chosen that all numbers begin with one or more zeros. On the other hand no digits should lie to the left of the four available locations. The unit of potential can best be chosen such that the highest potential V_{\max} occurring in the problem is given the numerical value 9999. In this way the most effective possible use is made of the available locations.

If we denote the numerical values of the potential given by the machine as $V = 9999 V/V_{\max}$, we can reduce eq. (7) to:

$$v_1 + v_2 + \left(1 + \frac{1}{2n}\right)v_3 + \left(1 - \frac{1}{2n}\right)v_4 = -\frac{9999 e_0 h^2}{\epsilon_0 V_{\max}} + 4v_0.$$

Here $n = r_0/h$ is a dimensionless quantity and is always a whole number. It can be seen from the equation that the most suitable unit of space charge is $\epsilon_0 V_{\max}/9999h^2$.

The difference equation of the trajectories

Now that the potential field is known we can proceed to calculate the trajectories described by the electrons in this field. We again use the cylindrical coordinates r , φ and z .

In the z direction the electric field has the component $-\partial V/\partial z$. The force acting on an electron (charge $-e$, mass m) in the z direction is thus $e\partial V/\partial z$, and the acceleration in the z direction is

$$\frac{d^2z}{dt^2} = \frac{e}{m} \frac{\partial V}{\partial z}, \text{ set } = L(r, z). \quad \dots \quad (11)$$

The acceleration in the r direction, d^2r/dt^2 , consists in the first place of an analogous term: $(e/m)\partial V/\partial r$, and in addition a term $r(d\varphi/dt)^2$ derived from the centrifugal force. The last term can be written

differently. Because of the rotational symmetry no other force acts in the φ direction, so that according to the law of conservation of momentum, $r^2 d\varphi/dt = r_0^2 (d\varphi/dt)_0$ (the subscript 0 refers to the initial state at the time $t = 0$). The contribution of the centrifugal force to the acceleration in the r direction, $r(d\varphi/dt)^2$, can therefore also be written: $(r_0^4/r^3)(d\varphi/dt)_0^2$. For the total acceleration in the r direction we then find:

$$\frac{d^2r}{dt^2} = \frac{e}{m} \frac{\partial V}{\partial r} + \left[r_0^2 \left(\frac{d\varphi}{dt} \right)_0 \right]^2 \frac{1}{r^3}, \text{ set} = K(r, z). \quad (12)$$

Together with the initial conditions, the equations (11) and (12) govern the motion of the electron, and thus the position of the electron at any given moment.

To solve the problem numerically we transform the differential equations (11) and (12) into difference equations. The distance of the electron from the z axis, r , is a function of the time t . We expand this function around $t = t_0$ in a Taylor series:

$$r(t + \Delta t) = r(t_0) + \left(\frac{dr}{dt} \right)_{t_0} \Delta t + \frac{1}{2} \left(\frac{d^2r}{dt^2} \right)_{t_0} \Delta t^2 + \frac{1}{6} \left(\frac{d^3r}{dt^3} \right)_{t_0} \Delta t^3 + \frac{1}{24} \left(\frac{d^4r}{dt^4} \right)_{t_0} \Delta t^4 + \frac{1}{120} \left(\frac{d^5r}{dt^5} \right)_{t_0} \Delta t^5 + \dots$$

We consider r at the moments $t_0 - \tau$ and $t_0 + \tau$, which differ from t_0 by a small amount τ . In the above series we substitute first τ and then $-\tau$ for Δt , and add the two results; we make τ small enough to allow terms with powers of τ higher than 5 to be neglected. The result is:

$$r(t_0 + \tau) + r(t_0 - \tau) = 2r(t_0) + \tau^2 \left(\frac{d^2r}{dt^2} \right)_{t_0} + \frac{\tau^4}{12} \left(\frac{d^4r}{dt^4} \right)_{t_0}. \quad (13)$$

In a similar way d^2r/dt^2 can be expanded in a series, which in the same way leads to

$$\tau^2 \left(\frac{d^2r}{dt^2} \right)_{t_0 + \tau} + \tau^2 \left(\frac{d^2r}{dt^2} \right)_{t_0 - \tau} = 2\tau^2 \left(\frac{d^2r}{dt^2} \right)_{t_0} + \tau^4 \left(\frac{d^4r}{dt^4} \right)_{t_0}. \quad (14)$$

We now make use of the differential equation (12). If the values of r and z at the time t are known, the value of the function K at the time t can be determined; we denote this value by K_t . By eliminating (d^4r/dt^4) from (13) and (14) and using (12) we find the required difference equation:

$$r(t_0 + \tau) = 2r(t_0) - r(t_0 - \tau) + \frac{1}{12} \tau^2 K_{t_0 - \tau} + \frac{10}{12} \tau^2 K_{t_0} + \frac{1}{12} \tau^2 K_{t_0 + \tau}. \quad (15)$$

In the same way we find for the z direction:

$$z(t_0 + \tau) = 2z(t_0) - z(t_0 - \tau) + \frac{1}{12} \tau^2 L_{t_0 - \tau} + \frac{10}{12} \tau^2 L_{t_0} + \frac{1}{12} \tau^2 L_{t_0 + \tau}. \quad (16)$$

If r and z at the times $t_0 - \tau$ and t_0 are known, equations (15) and (16) enable us to calculate r and z at the time $t_0 + \tau$, but not before a further difficulty has been overcome. K and L can be calculated at the times $t_0 - \tau$ and t_0 because we know r and z at those times. But this does not hold for K and L at the time $t_0 + \tau$. We must therefore again adopt an iterative method. Starting from the values of K and L at the times $t_0 - \tau$ and t_0 we make by extrapolation an estimate of K and L at the time $t_0 + \tau$; with the estimated values we compute r and z from (15) and (16), and use the values thus found for r and z for again computing K and L . If the new values differ too much from the original estimates, we compute r and z once again with the new values. This process is repeated until K and L stay constant. We have then also found the correct values of r and z . As a rule the process converges so quickly that the final result is reached after only one correction.

The difference equation of the beginning of the trajectory

We shall now consider two cases: (a) where the first and second points of the trajectory are given (which are passed with a time difference τ); (b) where the starting position and the initial velocity are known.

- a) If the first and the second point are known we calculate the third point by the method just described, then the fourth point from the second and third, and so on. In this way we find points of the trajectory which are passed by the electron with equal time differences.
- b) If the starting position and the initial velocity are given, the problem then is how to arrive at the second point so that we can proceed by the method mentioned under (a). This is done by twice integrating equations (11) and (12) for the initial state with respect to time.

Suppose that K is known in the first four points ($t = 0, \tau, 2\tau$ and 3τ). We approximate to K by a polynomial. The necessary precision requires one of the third degree:

$$K_t = a_0 + a_1 t + a_2 t^2 + a_3 t^3,$$

where the coefficients a_0, \dots, a_3 are functions of the given quantities K_0, \dots, K_3 such that the polynomial for the first four points acquires the respective values K_0, K_1, K_2, K_3 . We thus find not only the second point of the trajectory but at the same time the third and the fourth.

Integrating the polynomial twice with respect to time we obtain:

$$r(t) = r_0 + \left(\frac{dr}{dt} \right)_0 t + \frac{1}{2} a_0 t^2 + \frac{1}{6} a_1 t^3 + \frac{1}{12} a_2 t^4 + \frac{1}{20} a_3 t^5. \quad (17)$$

The integration constants are chosen so that (17) for $t = 0$ satisfies the initial conditions.

If we express the coefficients a in the given K quantities we can calculate r_1, r_2 and r_3 from (17):

$$\begin{aligned}
 r_1 &= r_0 + \tau \left(\frac{dr}{dt} \right)_0 + \frac{97}{360} \tau^2 K_0 + \frac{19}{60} \tau^2 K_1 - \frac{13}{120} \tau^2 K_2 + \\
 &\qquad\qquad\qquad + \frac{1}{45} \tau^2 K_3, \\
 r_2 &= r_0 + 2\tau \left(\frac{dr}{dt} \right)_0 + \frac{23}{45} \tau^2 K_0 + \frac{23}{15} \tau^2 K_1 - \frac{2}{15} \tau^2 K_2 + \\
 &\qquad\qquad\qquad + \frac{2}{45} \tau^2 K_3, \\
 r_3 &= r_0 + 3\tau \left(\frac{dr}{dt} \right)_0 + \frac{39}{40} \tau^2 K_0 + \frac{27}{10} \tau^2 K_1 + \frac{27}{40} \tau^2 K_2 + \\
 &\qquad\qquad\qquad + \frac{3}{20} \tau^2 K_3.
 \end{aligned}
 \tag{18}$$

Similar formulae apply to z ; all that is necessary is to replace r by z and K by L in (18).

To find the first four values of r and z we again use an iterative method; for the first estimate we make the K 's and L 's equal to K_0 and L_0 respectively; we then calculate r and z for the first four points using (18) and the corresponding equations for the z 's, and from these r and z values find better values for K and L . This process also converges quickly.

When the first four points have been found we compute the other points using (15) and (16).

Calculation of K and L

To determine the quantities K and L defined by (11) and (12) we need the partial derivatives of the potential V with respect to r and z at a given point (r, z) . For this purpose we first compute the derivatives at the nearest grid point (r_0, z_0) with the aid of the following formulae (see fig. 9):

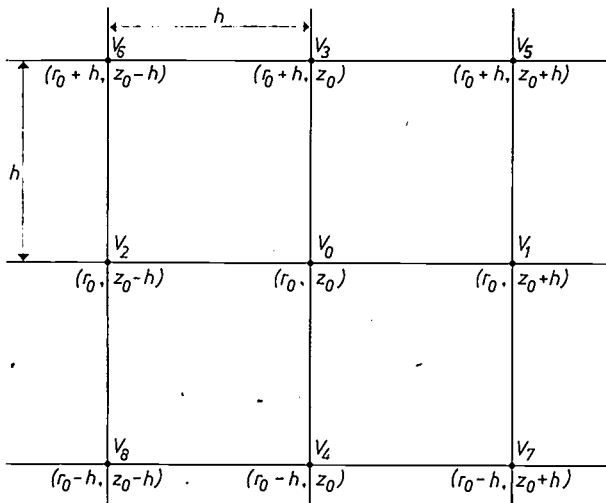


Fig. 9. On the basis of this figure, equations (19) were obtained.

$$\left. \begin{aligned}
 \left(\frac{\partial V}{\partial z} \right)_{r_0, z_0} &= \frac{V_1 - V_2}{2h}, & \left(\frac{\partial V}{\partial r} \right)_{r_0, z_0} &= \frac{V_3 - V_4}{2h}, \\
 \left(\frac{\partial^2 V}{\partial z^2} \right)_{r_0, z_0} &= \frac{\frac{V_1 - V_0}{h} - \frac{V_0 - V_2}{h}}{h} = \frac{V_1 + V_2 - 2V_0}{h^2}, \\
 \left(\frac{\partial^2 V}{\partial r^2} \right)_{r_0, z_0} &= \frac{\frac{V_3 - V_0}{h} - \frac{V_0 - V_4}{h}}{h} = \frac{V_3 + V_4 - 2V_0}{h^2}, \\
 \left(\frac{\partial^2 V}{\partial r \partial z} \right)_{r_0, z_0} &= \frac{\frac{V_5 - V_6}{2h} - \frac{V_7 - V_8}{2h}}{2h} = \frac{V_5 + V_8 - V_6 - V_7}{4h^2}.
 \end{aligned} \right\} \dots \tag{19}$$

The derivatives at point (r, z) can now be calculated by means of a "two-dimensional" series expansion around the grid point (r_0, z_0) . Using these derivatives we find K and L from equations (11) and (12).

Computing the trajectories with the IBM 650

The block diagram in fig. 10 shows the way in which the computer determines the electron trajectories. To start with, the initial conditions $r_0, z_0, (dr/dt)_0, (dz/dt)_0$ and $(d\varphi/dt)_0$ are fed to the machine by means of a punched card. The first step for the computer is to determine iteratively from (18) the first four points of the trajectory, the first estimate being, as mentioned above, $K = K_0$ and $L = L_0$.

When the first four points have been computed in this way, the other points are found using equations (15) and (16). The only difference between the programme for a succeeding point and that for the point just dealt with is that $t_0 + \tau$ is substituted in all cases for t_0 . The machine now computes by extrapolation the values of K and L at the time $t_0 + \tau$.

In the r - z plane there is a given region G within which the trajectories have to be computed. If, after some time, the machine finds a point (r, z) outside the region G , this means that the trajectory has been computed far enough and that the computation of the next trajectory can be started with the new initial conditions. When all punched cards with initial conditions have been used, all trajectories have been computed and the machine stops.

As regards the scale, we note that the unit of length can most usefully be taken as the greatest length and the unit of velocity as the highest velocity occurring in the problem. These two units fix the unit of time. The use of these units ensure that the order of magnitude of the numbers obtained is such that the accuracy of the calculation is maximum.

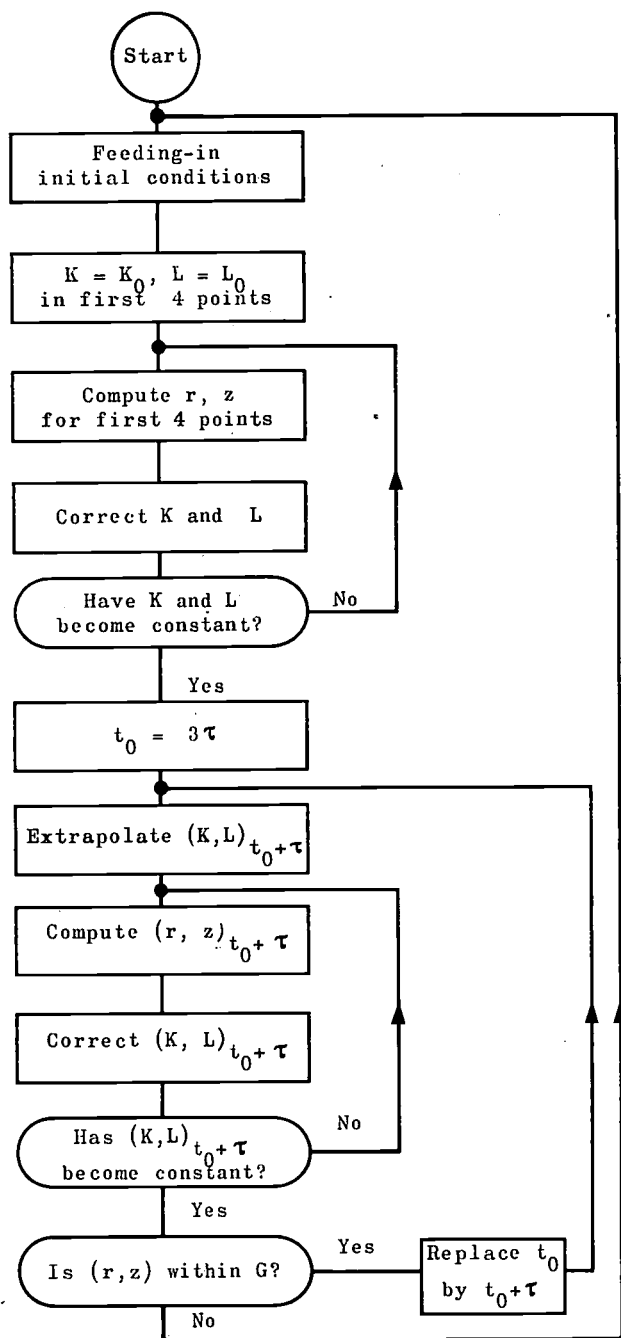


Fig. 10. Block diagram illustrating the computation of electron trajectories with the IBM 650.

Examples

Two concentric spheres

As our first example we again consider a case with two concentric conducting spheres. We now let the inner sphere function as cathode. Assuming that the electrons have no initial velocity on leaving the cathode, the space charge and the potential inside

such a system are capable of exact calculation¹⁴⁾, the results of which can be used to check the results obtained from the computer.

For the given space-charge distribution, the potential field between the spheres was calculated in the IBM 650 by the iterative method (radius of inner sphere $7h$, potential 0000; radius of outer sphere $21h$, potential 9999; these radii thus differ from those in fig. 5).

Since the space charge at the cathode is infinitely large, the numerical equations (6) and (7) are no longer valid in the neighbourhood of the cathode. However, the law connecting space charge and cathode distance is known¹⁴⁾; a correction can therefore be applied in the computation, and this was done in the present example.

Altogether 60 iteration cycles were required, in the course of which $|(\Delta V)_{\max}|$ dropped to 1. Each cycle lasted about 100 seconds, so that the total time was about 100 minutes. Table I gives the value of $|(\Delta V)_{\max}|$ after each cycle. The greatest deviation of the iteratively calculated potential from the one exactly calculated amounted to three units in the last digit.

Table I. Values of $|(\Delta V)_{\max}|$ after each of the 60 cycles which, in calculating the potential field between two concentric spheres, were required to reduce $|(\Delta V)_{\max}|$ to one unit in the last digit. The 60 cycles took about 100 minutes to complete.

Cycle number	$ (\Delta V)_{\max} $	Cycle number	$ (\Delta V)_{\max} $	Cycle number	$ (\Delta V)_{\max} $
1	9999	21	0814	41	0034
	7422		1030		0032
	7025		1060		0024
	7324		0855		0018
5	7605	25	0547	45	0021
	7384		0463		0014
	7042		0444		0008
	6683		0464		0010
	5913		0487		0009
10	4933	30	0423	50	0007
	4413		0281		0007
	3971		0241		0007
	3565		0154		0006
	3293		0072		0005
15	2544	35	0081	55	0002
	2698		0100		0002
	1787		0082		0001
	1982		0059		0001
	1757		0058		0001
20	2377	40	0043	60	0001

In a spherically symmetric configuration, electrons leaving the cathode with a radial initial velocity obviously describe straight, radial trajectories. Seven such trajectories are represented in fig. 11.

¹⁴⁾ I. Langmuir and K. B. Blodgett, Currents limited by space charge between concentric spheres, Phys. Rev. 24, 49-59, 1924. — The tables in this article were not accurate enough for our purpose, and we therefore recalculated them with greater accuracy.

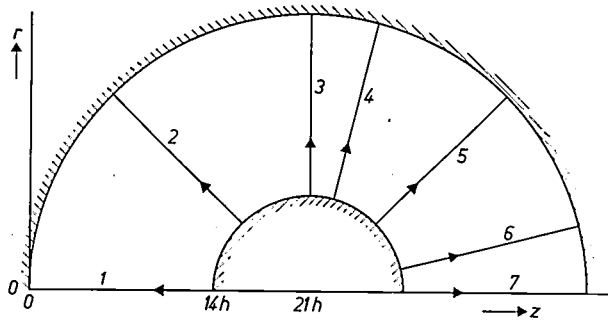


Fig. 11. Seven electron trajectories (1, . . . , 7) between two concentric spheres with radii $7h$ and $21h$. The small sphere acts as the cathode. The initial velocity of the electrons is assumed to be zero.

We have also calculated these trajectories from the potential field determined by iteration; the maximum deviation from the theoretical trajectories was only 0.03 mesh width.

Electron gun of a cathode-ray tube

The second example relates to the potential field of an electron gun of a cathode-ray (picture) tube. A sketch of the gun in cross-section is given in fig. 12. In the calculation the space charge was assumed to be zero.

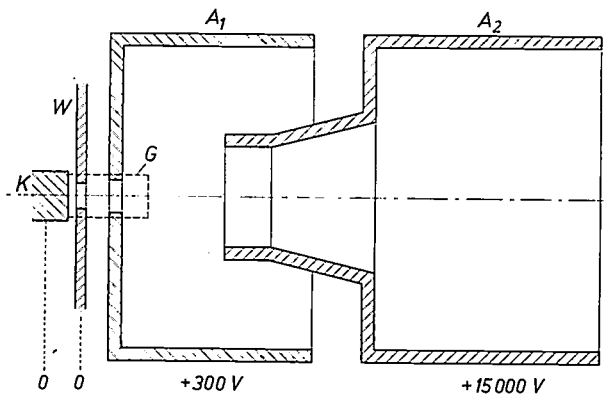


Fig. 12. Axial cross-section of the electron gun of a cathode-ray tube. *K* cathode. *W* Wehnelt cylinder. *A*₁ first anode. *A*₂ second anode. The form of the potential along the dotted rectangle *G* was predetermined by means of a resistance network; the values found were used as boundary conditions for the iterative computation of the potential field using the IBM 650.

With the aid of a resistance network ⁶⁾ an estimate on a reduced scale was made of the potentials along the dotted rectangle *G* in fig. 12. Then, using this potential distribution as the boundary condition, the potential field inside this rectangle was computed by the iterative method. The result can be seen in fig. 13, which also shows the form in which the line printer of the computer delivers the solution. The potential scale was chosen so that the maximum potential (bottom left in fig. 13) was represented by

about 9999. (In this case, for reasons which are not relevant here, 0400 was added to every potential value, hence the fact that the potential of the cathode and that of the Wehnelt cylinder had the value 0400.)

Finally the trajectories which are described in this potential field by the electrons leaving the cathode with zero velocity were also computed. Some of these trajectories are indicated in fig. 13.

Appendix: Speed of convergence of iterative methods

To give some idea of the speed of convergence in iteration problems of the kind discussed, we shall consider first the following simple iteration problem:

$$V_0^{(k+1)} = \frac{1}{4}V_1^{(k)} + \frac{1}{4}V_2^{(k)} + \frac{1}{4}V_3^{(k)} + \frac{1}{4}V_4^{(k)} + \frac{1}{2n} [V_3^{(k)} - V_4^{(k)}], \dots \quad (20)$$

where $n = r_0/h$. The potential value of the $(k+1)$ st cycle is thus computed entirely from the old potential values of the k th cycle, without applying successive over-relaxation. The method we shall adopt is that of Garabedian ¹⁵⁾. For this purpose we write (20) in this form:

$$\frac{V_1^{(k)} + V_2^{(k)} - 2V_0^{(k)}}{h^2} + \frac{V_3^{(k)} + V_4^{(k)} - 2V_0^{(k)}}{h^2} + \frac{V_3^{(k)} - V_4^{(k)}}{2nh^2} = \frac{4}{h} \frac{V_0^{(k+1)} - V_0^{(k)}}{h}$$

We now replace the superscript k , which increases in steps of 1, by a superscript t which increases continuously, in such a way that every time k increases by 1, t increases by the amount h . As a result of this substitution the difference expression for $(4/h)(\partial V_0/\partial t)$ appears in the right-hand side of the last equation. Returning to the differential form we now obtain the following differential equation:

$$\frac{\partial^2 V}{\partial z^2} + \frac{\partial^2 V}{\partial r^2} + \frac{1}{r} \frac{\partial V}{\partial r} = \frac{4}{h} \frac{\partial V}{\partial t}$$

The method of separating the variables yields the solution:

$$V = V^{(\infty)}(r,z) + \sum_{a=1}^{\infty} b_a \exp(-h\lambda_a t/4) V_a(r,z), \dots \quad (21)$$

where V_a and λ_a are given by the solution of the following eigenvalue equation:

$$\frac{\partial^2 V_a}{\partial z^2} + \frac{\partial^2 V_a}{\partial r^2} + \frac{1}{r} \frac{\partial V_a}{\partial r} + \lambda_a V_a = 0, \dots \quad (22)$$

where the boundary values are zero.

With the given boundary conditions, $V^{(\infty)}(r,z)$ is the solution of the problem. If we begin with arbitrary initial values of the potential, the error can be resolved into the eigenfunctions V_a ; the coefficients are then the constants b_a from (21). The manner in which each term of this series decreases after a number of iteration cycles is given by $\exp(-h\lambda_a t/4)$. After each iteration cycle the term thus decreases by the factor $\exp(-h^2\lambda_a/4)$. After a large number of cycles the speed of

¹⁵⁾ P. R. Garabedian, Estimation of the relaxation factor for small mesh size. Mathematical tables and other aids to computation 10, 183-185, 1956.

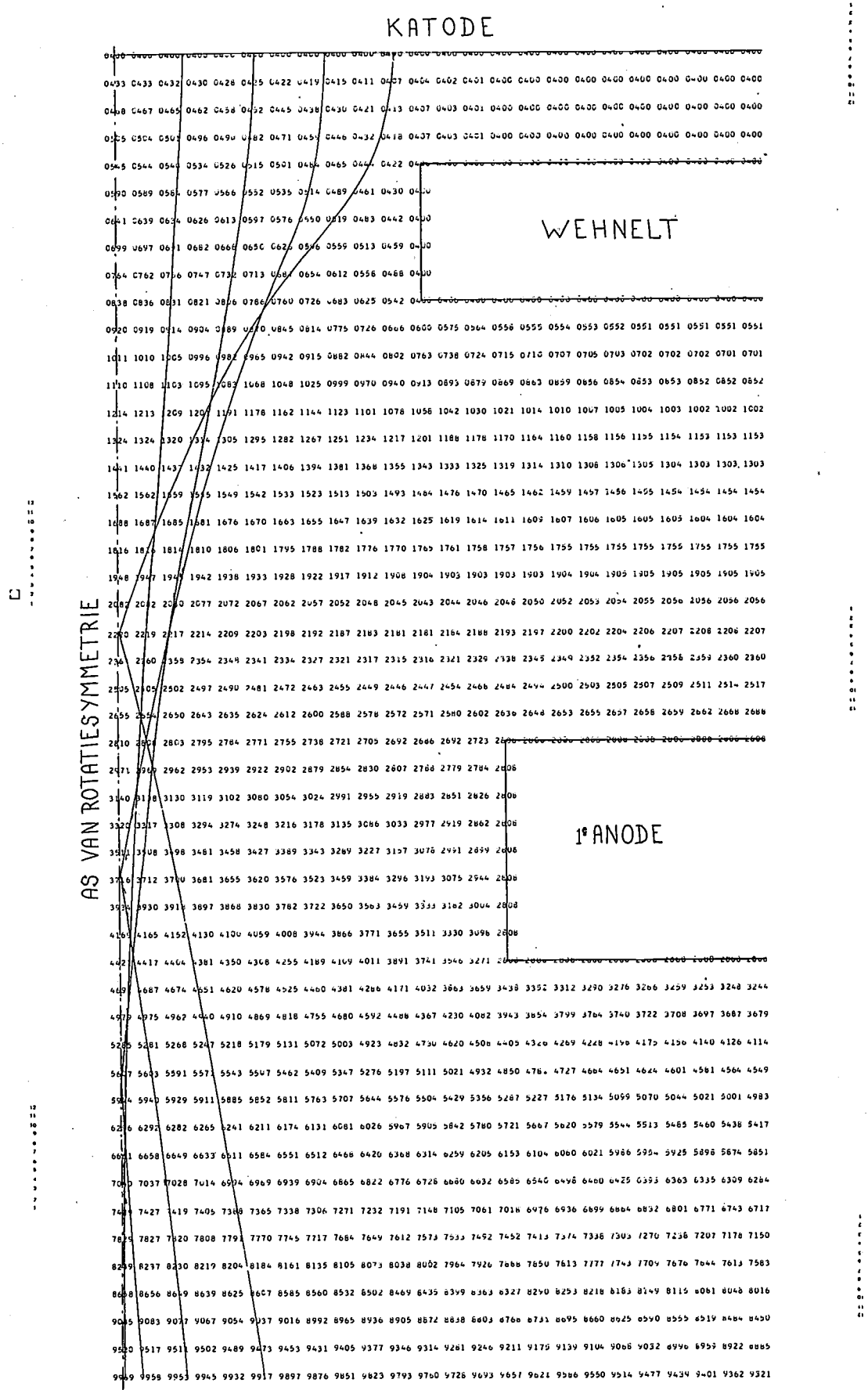


Fig. 13. The form in which the IBM 650 delivers the solution of a problem, in this case the potential field within the rectangle G in the electron gun represented in fig. 12. (All values have to be reduced by 0400.) The figure shows the axis of rotational symmetry, the periphery of the electrodes and some electron trajectories with initial velocity zero.

convergence is determined by the smallest eigenvalue λ_a ; the other terms have then become negligibly small.

If the boundary of the area has the form of a cylinder of revolution, with length ph and radius qh (fig. 14), the eigenvalues of (22) can be computed by the method of the separation of variables, and we find for the smallest eigenvalue:

$$\lambda_{\min} = \frac{\pi^2}{p^2h^2} + \frac{\mu_1^2}{q^2h^2},$$

where μ_1 is the first zero point of the bessel function of zero order.

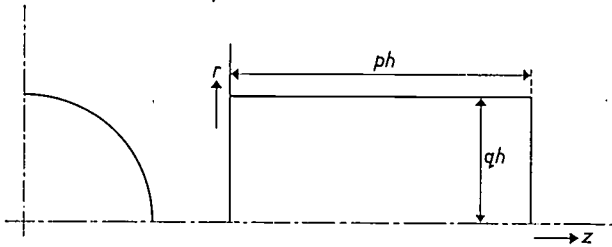


Fig. 14. Cylinder of revolution with the z axis as the axis of symmetry and bounded by flat planes perpendicular to the z axis. Length ph , radius qh . In this simple case λ_{\min} can be computed. The cylinder, with properly chosen p and q , is a sufficiently accurate approximation for many more complicated cases.

After a reasonably large number of iteration cycles have been computed, the speed of convergence will, in practice, be governed by this λ_{\min} . After each succeeding iteration the error will thus decrease by the factor

$$\exp \left[- \left(\frac{\pi^2}{4p^2} + \frac{\mu_1^2}{4q^2} \right) \right] \approx 1 - \frac{\pi^2}{4p^2} - \frac{\mu_1^2}{4q^2}.$$

(The approximation holds if p^2 and q^2 are large compared with unity.)

Similar considerations to those discussed for this simple case

are also applicable when use is made of the method of successive over-relaxation. The over-relaxation factor ω then occurs in the factors by which the eigenfunctions V_a decrease. The value of ω can be found at which the absolute value of these factors drops to minimum; this value of ω gives the fastest convergence¹⁶⁾.

If the boundary of the area is not a cylinder of revolution, equation (22) is not so easy to solve. In such cases the best course is to approximate to the boundary by a suitably chosen cylinder, and to apply the optimum value of ω thus found to the case with the original boundary¹⁷⁾.

¹⁶⁾ Another method of finding this optimum value of ω is to substitute the value found for λ_{\min} in certain formulae derived from a theory put forward by Young. See: D. Young, Iterative methods for solving partial difference equations of elliptic type, Trans. Amer. Math. Soc. 76, 92-111, 1954.

¹⁷⁾ D. Young, ORDVAC solutions of the Dirichlet problem, J. Ass. Computing Machinery 2, 137-161, 1955.

Summary. An important problem in electron optics is to determine the trajectories described by the electrons in the space between electrodes of given configuration, each with a known potential. This article deals with the numerical method of solving this problem. A great advantage of the numerical method is that it is also applicable to cases involving space charge (Poisson's equation). The method calls for a great deal of computation, and has only become feasible since the advent of electronic computers.

After discussing the transformation of Poisson's differential equation into a difference equation, the article deals with the iteration process and the method of successive over-relaxation to accelerate convergence. The difference equation for the trajectories is then derived. It is explained how the potential fields and orbits are computed using an IBM 650 machine of Philips Computing Centre (the programming of the fast machine PASCAL has only recently been completed).

The subject is illustrated with two examples relating to the space between two concentric spheres and the space within the electron gun of a cathode-ray tube.

An appendix deals with the speed of convergence of iterative methods.

VACUUM DEPOSITION OF RESISTORS

by P. HUIJER *), W. T. LANGENDAM *) and J. A. LELY *).

621.316.849:686.49

Vacuum deposition of metals and alloys

Certain parts of electronic circuits such as resistors, capacitors and the electrical connections can be made by depositing metal films on a base or substrate, such as a thin glass plate. Other elements like coils and transistors can easily be soldered to the film. It is possible in this way to obtain very compact circuits (fig. 1). Now that electronic systems for data

processing, automatic control, etc., are becoming more and more complicated, this technique, which is called microminiaturization, is beginning to assume great importance.

The procedure for making *electrical connections* by vacuum deposition is as follows. A charge of metal with a high conductivity is heated to above its melting point inside an evacuated bell-jar. Some of the vapour given off by the metal condenses on the cold substrate in the bell-jar. Evaporation is continued until a layer of the required thickness has formed; for example, a $0.1 \mu\text{m}$ film will form within a few minutes. A film with a given pattern can be obtained either by covering the substrate with a mask during deposition, or by etching away unwanted parts of the film afterwards¹⁾.

The pure metals used for vacuum-deposited wiring are not suitable for making *resistors* by the same technique. The resistivity of common metals is so small that only a film a few atoms thick would offer the required electrical resistance, and films of this thickness could not be prepared within sufficiently close tolerances. Thus, for the deposition of resistors, alloys such as nickel-chromium must be used whose specific resistance ($100 \mu\Omega\text{cm}$) is many times as great as that of pure metals. Many alloys have the additional advantage of possessing a much smaller temperature coefficient (less than $400 \times 10^{-6}/^\circ\text{C}$ in the case of 80% nickel and 20% chromium).

Fig. 2 is a (purely schematic) phase diagram showing the equilibrium between the liquid and vapour phases of nickel-chromium alloys at a constant temperature. From this diagram it is possible to deduce what happens when an alloy is evaporated out of the liquid phase in the same way as an unalloyed metal. Consider a liquid phase whose Cr concentration is c_1 and whose state is represented by point *A*; the vapour above the liquid and in equilibrium with it, corresponding to point *A'*, has a much higher Cr concentration c_1' . The reason for this is that each component evaporates independently of the other, at its own particular rate, Cr evaporating much faster than Ni. The system is continuously losing Cr, which condenses on the cold

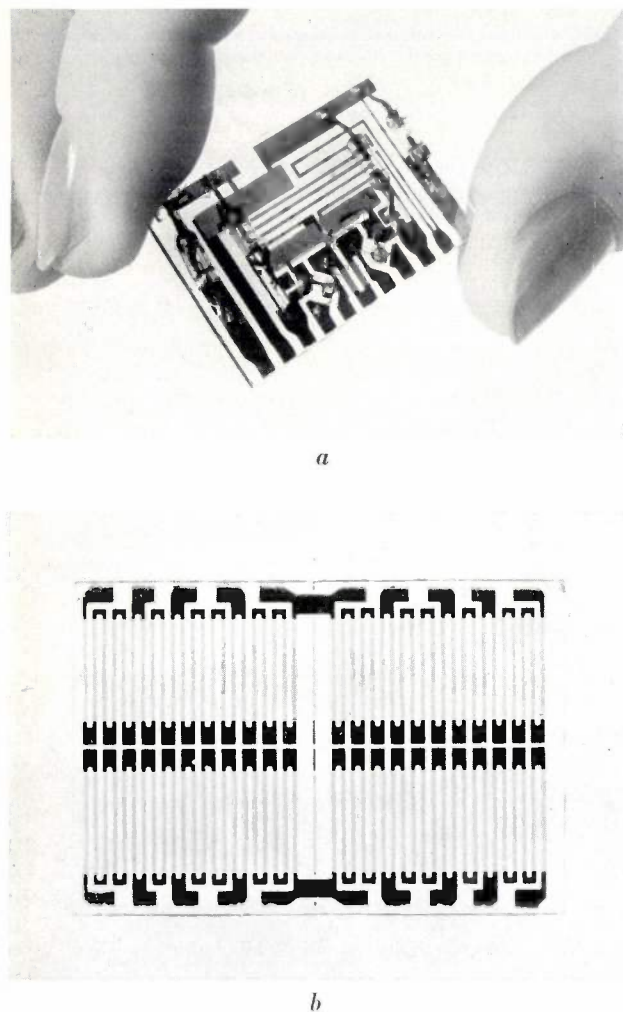


Fig. 1. a) Electronic circuit consisting in part of vacuum-deposited elements.

b) Series array of vacuum-deposited resistors, to which tapplings may be made at various points. The dark areas are conductors and the light-coloured strips are resistors. Each of the resistors lying between two conductive areas has a value of $5.6 \text{ k}\Omega$. The substrate dimensions are $20 \times 30 \text{ mm}$.

*) Industrial Components and Materials Division, Philips, Eindhoven.

¹⁾ For information on vacuum deposition in general, see for example L. Holland, *Vacuum deposition of thin films*, Chapman & Hall, London 1956.

substrates and the walls of the bell-jar, the higher Cr concentration in the vapour than in the liquid being obtained by drawing on the liquid phase. The consequent decrease in the Cr concentration of the liquid phase affects the whole volume of liquid, not just a surface layer, because of the stirring action of the convection currents. In this way the Cr concentration in the liquid drops to c_2 and that in the vapour to c_2' (corresponding to points B and B' in the diagram). The decrease continues until all the Cr has been used up (point E in the diagram). The process is in all respects analogous to the fractional distillation of an alcohol and water mixture, for example.

In general, the properties of an alloy depend to a large extent on its composition. Vapour produced in this way cannot be used to deposit films with uniform electrical properties because its Cr concentration is constantly changing. If a large number of exactly similar resistors is required, some means must be found for controlling the composition of the vapour. We shall briefly describe three appropriate methods.

In the first method, known as "flash evaporation"²⁾, the starting material is a correctly pro-

portioned mixture of the components of the alloy in powder form. Small amounts of the mixture are dropped at intervals onto an evaporator, a hot tantalum strip for example, the temperature of which is so high that the whole amount evaporates almost instantaneously. Consequently the vapour in the bell-jar has the right composition at all times, and the same applies to the condensed film. The method works quite well even when there is a big difference between the evaporation rates of the components.

In the second method the individual components are evaporated continuously but each from its *own source*³⁾. That is to say, mixing takes place in the vapour phase. The success of the method depends on accurate adjustment of the individual evaporation rates. It has proved possible in practice to control evaporation of the components with the required accuracy, and so to achieve efficient resistor production by this method.

The third method is based on the evaporation from the solid state of a charge of alloy prepared in advance, and not out of the liquid state which would lead to fractional distillation. This third method also gives good results⁴⁾, but the situation in this case is more complicated. While studying the method in detail we came up against a mathematical problem which we solved with the aid of a computer. This problem and its solution form the subject of the present article.

The discussion of the method will be confined to a nickel-chromium alloy consisting of 80% Ni and 20% Cr, which from now on we shall refer to as 80Ni20Cr. The other well-known alloys with high resistivity, such as constantan and manganin, contain the element manganese, which has a relatively high evaporation rate and therefore makes it even more difficult to deposit homogeneous films.

The essential difference between evaporation from the solid phase and evaporation from the liquid phase is that in the former case the rates at which the components evaporate are interdependent. Suppose that all the atoms of the more volatile element, which is Cr, have escaped from the outside atomic layers of the solid; before any more Cr can evaporate the way must be cleared by the evaporation of nickel atoms. This argument only partly represents the actual evaporation process, since the Cr atoms

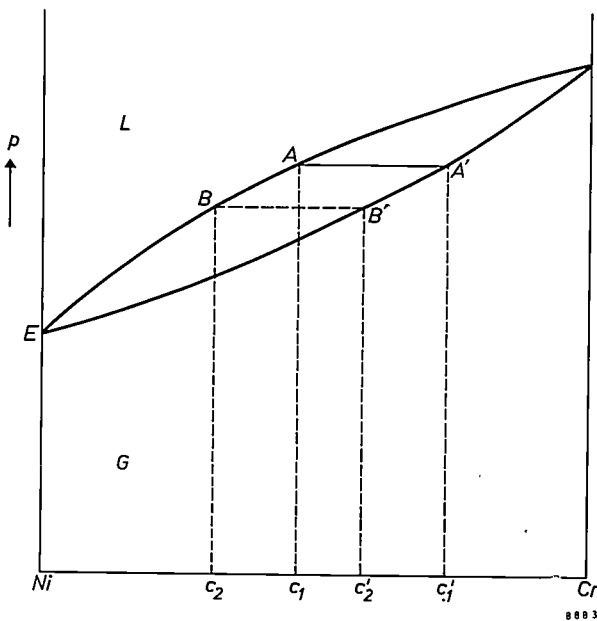


Fig. 2. Phase diagram for the system Ni-Cr, showing pressure p versus Cr concentration c at a constant (high) temperature. The alloy is a vapour in low-pressure region G and a liquid in region L . In the intermediate region the vapour and liquid phases are in equilibrium. If an alloy with a Cr concentration of c_1 is evaporated from the liquid phase, the resulting vapour will have a concentration of c_1' . Because the system is continuously losing Cr, the Cr concentrations in the vapour and in the liquid decrease in the direction of c_2' and c_2 , finally becoming zero ($c = c' = 0$). Evaporation from the liquid phase is unsuitable for vacuum deposition of resistances because the composition of the vapour never attains a steady value.

²⁾ M. Beckerman and R. L. Bullard, Proc. 1962 Electronic Compon. Conf., Washington, p. 53.

³⁾ M. Beckerman and R. E. Thun, Trans. 8th Nat. Vac. Symp., Washington, 1961, Pergamon Press, New York 1962, p. 905. W. J. Ostrander and C. W. Lewis, *ibid.* p. 881.

⁴⁾ M. Schneider, Techn. Mitt. PTT 37, 465, 1959, especially p. 470 ff., or T. K. Lakshmanan, Trans. 8th Nat. Vac. Symp., Washington, 1961, Pergamon Press, New York 1962, p. 868.

can also reach the surface by *diffusion*. But the diffusion process is much slower than convection, and hence, in the solid phase, equilibrium may be established between the rate at which the diffusing chromium atoms are arriving at the surface and the rate at which they are removed by evaporation. Once this steady state has been attained the composition of the vapour will remain constant; and films deposited successively will have the same electrical characteristics.

We have been able to demonstrate that a steady state of this kind can in fact be attained when evaporation takes place from the solid phase. We shall now go into the matter and show how the time taken to reach the steady state is calculated. Little importance need be attached to the conditions actually prevailing in the steady state; if this is associated with a vapour phase of composition other than the desired one, then the remedy will be to modify the starting composition.

The steady state

The sample of 80Ni20Cr to be evaporated from the solid phase — it might take the form of a wire⁴⁾ — has to be raised to a temperature of about 1300 °C (the melting point of the alloy is 1395 °C). At this temperature the vapour pressures of Ni and Cr are 0.25×10^{-3} and 1.8×10^{-3} torr respectively — high enough to ensure reasonably fast evaporation.

If the resulting vapour is condensed in a series of films, the films first deposited prove to have a very high Cr concentration. This is in accordance with Langmuir's formula⁵⁾, which gives the rate W_i in grammes per second at which a component i evaporates from unit area at a temperature T :

$$W_i = 5.85 \times 10^{-2} p_i \sqrt{\frac{M_i}{T}} c_i \text{ g/cm}^2\text{s}, \quad (1)$$

where M_i is the molecular weight of the component, p_i its vapour pressure in torr at T °K, and c_i a fraction indicating its atomic concentration *at the surface*. Ni and Cr have about the same atomic weight (59 and 52 respectively), but the vapour pressure of Cr is a good 7 times that of Ni. Hence the Cr evaporates much faster than the Ni, with the result that the system becomes short of Cr; c_{Cr} de-

creases, as does W_{Cr} and the Cr concentration in the vapour (*fig. 3*). The question now arising is whether the decrease of the Cr concentration in the vapour finally levels off, a steady state thus being attained. Langmuir's formula indicates that it cannot do so unless the concentration c_i *at the surface* attains a constant value; for this to happen, the distribution of concentrations *through the bulk of the wire* would

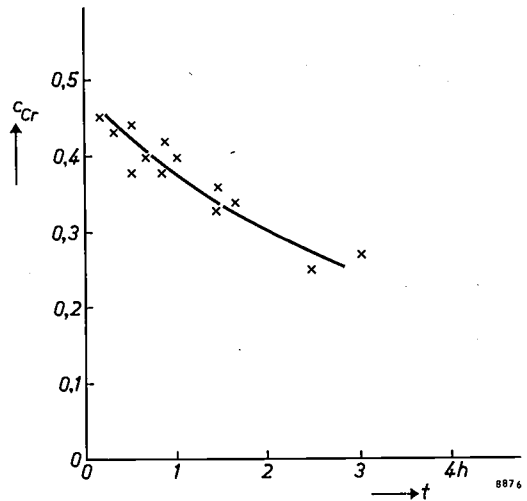


Fig. 3. Vapour composition as a function of time when evaporation takes place from the solid phase. The sample of alloy to be evaporated was in the form of a wire which was in the original state at instant $t = 0$. The vapour composition at various stages of evaporation was determined by X-ray spectroscopic examination of the corresponding condensates (i.e. the deposited metallic films). The spread in the measured concentrations is of the correct order for this method of determination.

have to become stationary. The problem is therefore one of finding whether a concentration gradient can exist in the alloy independent of time.

If the material is in the form of a wire, the surface from which evaporation takes place will be curved, but we shall first consider the simplified case of a wire of infinite radius, which corresponds to a flat interface between a vacuum on the left and a block of alloy on the right which we can imagine extending to infinity, this infinitely large block having a temperature of 1300 °C. An x axis is drawn perpendicular to the interface. For large positive values of x the Cr concentration is constant at 0.2. The evaporation process gives rise to a concentration gradient over the range of small x values (i.e. near the interface), in consequence of which diffusion takes place. The latter process is described by the familiar diffusion equation:

$$\frac{\partial c}{\partial t} = D \frac{\partial^2 c}{\partial x^2}, \quad \dots \dots (2)$$

⁵⁾ This formula is intended to describe the evaporation of single substances. It is permissible to use it for the components of an alloy provided these form homogeneous mixed crystals and provided the alloy is an ideal solid solution. Strictly speaking, 80Ni20Cr satisfies the first condition but not the second; however, our experiments have shown that in regard to evaporation, 80Ni20Cr can safely be regarded as an ideal solid solution.

where t is time and D is the diffusion constant ⁶⁾. Before applying this equation to the present problem we must remember that Ni and Cr atoms are continuously being lost by evaporation, so that the interface is steadily moving to the right. The "stationary distribution of concentrations" that we are looking for cannot strictly speaking be stationary at all, since it will move with the interface, but it must have a form independent of time. The obvious course will therefore be to adopt a moving frame of reference, the origin of which is located in the interface. Let y be the distance of an arbitrary fixed point from this interface, then

$$y = x - vt,$$

where v is the constant rate of displacement of the interface. On substitution of y , (2) becomes

$$\frac{\partial c}{\partial t} = D \frac{\partial^2 c}{\partial y^2} + v \frac{\partial c}{\partial y} \quad \dots \quad (3)$$

In the steady state $\partial c/\partial t = 0$, and that state will be attainable provided the equation

$$D \frac{\partial^2 c}{\partial y^2} + v \frac{\partial c}{\partial y} = 0$$

possesses a real solution. It does in fact prove to have the real solution

$$c = c_\infty - A \exp\left(-\frac{v}{D}y\right) \quad \dots \quad (4)$$

A is a constant depending on the ratio between the evaporation rates of Ni and Cr; the value of c_∞ is 0.2.

Clearly, the shape of the steady-state concentration curve is governed by the ratio v/D , the fall-off being restricted to a narrow surface layer in cases where v is relatively large (i.e. where evaporation is fast) or where D is relatively small. Concentration curves for two values of v/D have been plotted in fig. 4.

The above solution applies to a flat interface. To a first approximation it is also valid for evaporation from a cylindrical wire, provided its diameter is not too small. We were able to determine the distribution of concentrations in a 2 mm wire experimentally, by repeatedly dipping it in a hot mixture of H_3PO_4 and $HClO_4$ and measuring the Cr concentrations in the successive "skins" thus stripped from the wire. A plot of these concentrations has much the same shape as the curves in fig. 4.

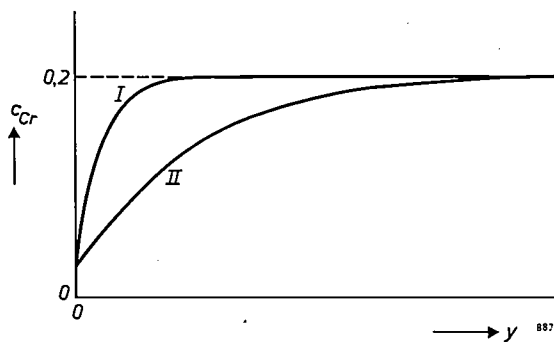


Fig. 4. The steady-state distribution $c(y)$ of Cr concentrations in a block bounded by a flat surface, for (I) large and (II) small values of parameter v/D .

At this juncture we still do not know how long it takes for the steady state to establish itself, but we can at least demonstrate that once it is attained, the composition of the vapour must be the same as that of the bulk of the solid alloy. Over a certain interval of time evaporation will cause the interface to shift through a distance l (fig. 5), the shape of the distribution curve undergoing no change. If evaporation is taking place from a surface of area S , the volume of solid thus removed and converted into vapour is lS . The proportion of Cr in it is given by S times the hatched area in fig. 5; this area is the difference between the areas of two plane figures individually bounded by curves 1 and 2 and extending up to a remote point P on the x axis, and it has a value of $0.2l$ (being a curved parallelogram of height 0.2 on a base l). The Cr concentration in the vapour is therefore 0.2, the same as that in the bulk of the solid alloy.

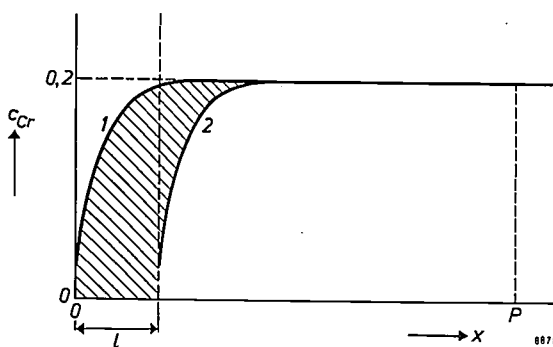


Fig. 5. If, in the steady state, evaporation causes the interface to shift through a distance of l , a volume lS will be removed per area S of interface. The proportion of Cr in this volume is given by S times the hatched area in the above diagram, this area being $0.2l$. Once the steady state has been attained, the Cr concentrations in the solid and vapour phases are both equal to 0.2.

Time required to reach the steady state

We shall now calculate the time that must elapse before equilibrium is sufficiently closely approached to allow the vacuum deposition of the resistive films

⁶⁾ The Ni concentration is greater at the surface than in the bulk of the sample, and consequently Ni diffuses into the interior. However, we feel it is justified to base the argument on the difference between the outward diffusion of Cr and the inward diffusion of Ni. D thus represents the net diffusion.

to be started. For the case of the flat interface this entails solving eq. (3) in its complete form. With the aid of formula (4), in which an estimated D and a measured v value were inserted, it has been established that the layer within which the concentration varies is about 100 μm thick. Unless very thick wires are being used, the curvature of the surface may have a marked influence on the evaporation process right from the start. It will therefore be safer to solve the equation right away for the case of a cylindrical wire. For this purpose we shall adopt the cylindrical coordinates r and φ . Eq. (2) now assumes the form

$$\frac{\partial c}{\partial t} = D \left(\frac{\partial^2 c}{\partial r^2} + \frac{1}{r} \frac{\partial c}{\partial r} \right), \quad \dots \quad (5)$$

where r denotes distance from the axis of the wire.

As an initial condition we shall assume that at the instant $t = 0$ the Cr concentration is 0.2 throughout the sample of alloy. As boundary condition we can take Langmuir's formula, bearing in mind that W_{Cr} is proportional to the slope of the concentration curve at the surface of the wire, at a distance ϱ from the axis. Stated mathematically, the condition is that

$$\left(\frac{\partial c}{\partial r} \right)_{r=\varrho} = kc(\varrho), \quad \dots \quad (6)$$

k being a constant of proportionality. In this case the displacement of the interface is not allowed for by introducing a moving frame of reference, as was done in the case of the flat interface. Accordingly the boundary condition will apply to a *moving* boundary and hence to a variable value of ϱ , viz. $\varrho = r_0 - vt$, where r_0 is the radius of the wire at instant $t = 0$ and v is the rate of displacement of the interface.

Differential equation (5) subject to boundary condition (6) cannot be solved by the usual analytical methods, and we therefore used a computer, arriving at the solution by numerical methods.

Quantities r and t were plotted in rectangular coordinates; fig. 6 shows a quadratic network of points belonging to this system. The distance Δr might represent 1 μm and the interval Δt 40 s. The computer had to work out a Cr concentration satisfying eq. (5) for each point in the network. With this in view (5) was first transformed into a difference equation, in much the same way as is described in one of the other articles in this issue⁷⁾. However, in the present case the value appropriate to a point P was calculated, not from those of the four surrounding points, but from the values at R , S

and T , the three points underneath P , which relate to a stage Δt earlier in the process. The difference equation is

$$c_P = D \left(\frac{\Delta t}{\Delta r^2} - \frac{1}{2r} \frac{\Delta t}{\Delta r} \right) c_T + \left(1 - D \frac{2\Delta t}{\Delta r^2} \right) c_S + D \left(\frac{\Delta t}{\Delta r^2} + \frac{1}{2r} \frac{\Delta t}{\Delta r} \right) c_R \dots \quad (7)$$

If the distribution of concentrations $c(r)$ at a given instant t_1 is known, the machine is able, with the aid of eq. (7), to determine point by point the distribution at instant $t_1 + \Delta t$.

In the t direction the network extends to more than 4 hours. The wire diameter and the diffusion and evaporation rates had values such that the variation in concentration could not, over a period of 4 hours, travel more than halfway along the radius of the wire. At points corresponding to small r values, then, the concentration will not appreciably deviate from 0.2, and there was no need for these points to be computed. The machine worked along each row relating to a given t value in the direction indicated by the arrow in fig. 6. Each time it reached the last point in a row (e.g. P'), boundary condition (6) had to be taken into account. The implication is that (6) had to be satisfied along line a , which shows the decrease in the radius of the wire as a function of time. With the aid of the previously determined concentrations at T' and S' , and of a difference formula derived from (6), the machine first calculated the concentration at Q and then extrapolated to a point R' beyond line a , thus arriving at a fictive

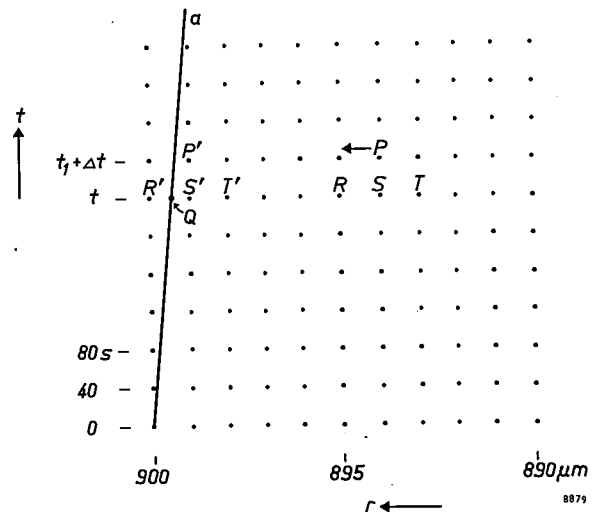


Fig. 6. Grid corresponding to a system of rectangular coordinates in which t is plotted against r . The computer works out the Cr concentration at each point of the grid with the aid of difference equation (7). Curve a indicates the position of the surface as a function of time.

⁷⁾ See C. Weber, Calculation of potential fields and electron trajectories using an electronic computer, pp. 130-143 of this issue.

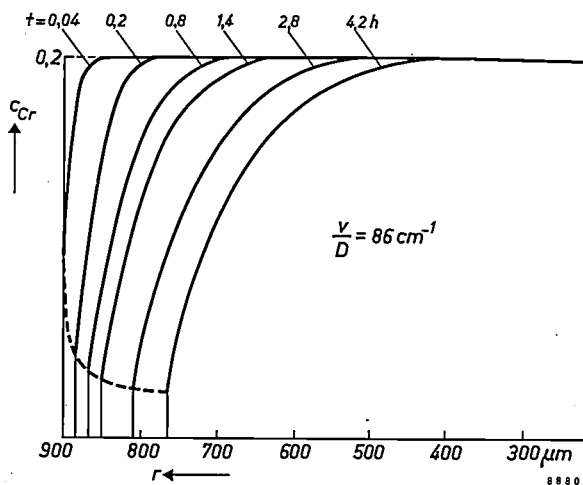


Fig. 7. Computed $c(r)$ curves relating to a constant value of v/D , namely 86 cm^{-1} , showing the distribution of concentrations in the wire after periods of 0.04, 0.2, 0.8, 1.4, 2.8 and 4.2 hours.

concentration for this point which had already evaporated. The machine was then in a position to derive the concentration at P' from formula (7) in the normal way.

Inspection of (7) will make it clear that the machine only had simple mathematical operations to perform. The value of a computer lies first and foremost in the great speed with which it works; in the present case concentrations at about 25 000 points constituting the network had to be calculated. The IBM 650 machine used required only two and a half hours to do this. We have plotted the results graphically for certain values of t (fig. 7). Fig. 8 displays curves derived from fig. 7, showing distributions of Cr concentrations in the vapour and at

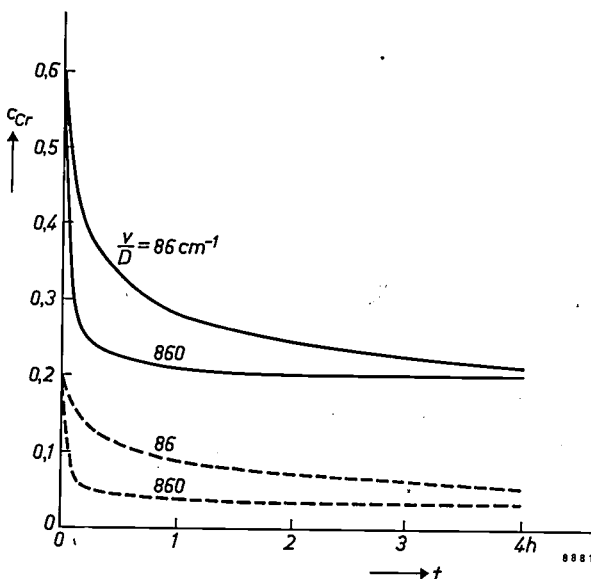


Fig. 8. Computed curves relating to two different values of v/D and showing changes in the Cr concentrations in the vapour (solid lines) and in the surface of the wire (dotted lines) as a function of time.

the surface of the wire as a function of time, and for two different values of parameter v/D .

The complete computation procedure had to be carried out for the two values of v/D indicated in fig. 8 because the diffusion constant of Cr in Ni at 1300°C was not known in advance. Measurements of the radius of the wire showed that this diminished at the rate of $6.1 \mu m$ per hour, thus giving v . The two chosen values of v/D lie on either side of an estimated value based on interim calculations. We were able subsequently to determine the diffusion constant by comparing the computed with the measured distribution of concentrations in the steady state. The two curves appear in the same graph in fig. 9; from these it is possible to deduce a value of 253 cm^{-1} for v/D and one of $0.67 \times 10^{-9} \text{ cm}^2 \text{ s}^{-1}$

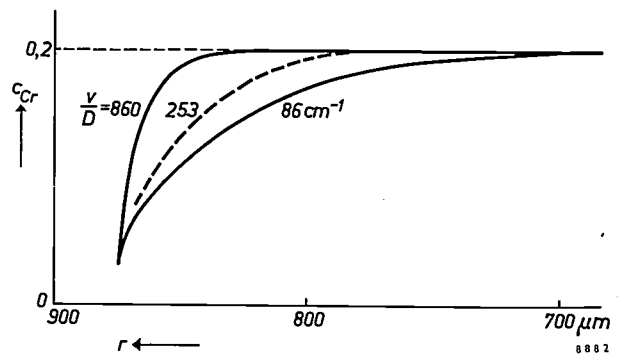


Fig. 9. Computed (solid line) and measured (dotted line) distribution of Cr concentrations in the wire after 4 hours. By comparing these curves it is possible to deduce that the diffusion constant of Cr in Ni at 1300°C is $0.67 \times 10^{-9} \text{ cm}^2/\text{s}$.

for D itself. The diffusion constant now being known, it can be estimated from fig. 8 that in the case dealt with here a time of about 3 hours is necessary for the steady state to be attained. This is in good agreement with experiment.

Fig. 8 indicates that the Cr concentration at the surface attains an equilibrium value of about 0.03. If this value of c_i is inserted in Langmuir's formula (1), we obtain values of W_{Ni} and W_{Cr} which correspond exactly to a vapour composition of 80Ni20Cr. The above ideas about evaporation from the solid phase are thus confirmed.

Summary. One of the techniques used in microminiaturization is the preparation of electronic circuit elements by vacuum deposition of metallic films on substrates such as thin sheets of glass. Resistors cannot be made by depositing the pure metals used for interconnections because the specific resistance of the metals is too low. The deposition of suitable resistance alloys (e.g. nickel-chromium) involves the problem of obtaining a vapour, and hence a condensate, of constant composition. The problem can be solved in various ways. One method, which has been studied in detail, is to evaporate the alloy from the solid instead of from the liquid state. The composition of the vapour will remain constant as soon as the distribution of concentrations in the solid phase has become stationary. With the aid of a computer it is possible to determine the shape of the distribution curve prevailing in the steady state and the time taken for this state to be attained.

DESIGN OF FERROXDURE LOUDSPEAKER MAGNETS

by M. F. REYNST *) and W. T. LANGENDAM *).

538.26:621.395.623.742

Magnets for use in loudspeakers may either be made of metal ("Ticonal" type alloys) or of the ceramic material ferroxdure. No clear dividing line can be drawn between cases where purely technical considerations dictate the use of "Ticonal" and where they dictate the use of ferroxdure. It is usually possible to build flatter constructions with ferroxdure, which may sometimes be advantageous; on the other hand, simpler designs are possible using "Ticonal", in which the external leakage flux is very low.

A constantly recurring problem is that of designing the magnetic circuit. Normally, the problem is to generate a specific magnetic induction in a ring-shaped air gap of given dimensions, in which the loudspeaker coil is to be mounted, and for this purpose the dimensions of the magnet have to be determined. In the design calculations involved it is difficult to take into account fully the influence of the leakage and of various other factors which will presently be discussed ¹⁾. To do this successfully requires experience.

In this article a practical method of calculation is discussed which can be used in designing the magnet system of loudspeakers equipped with a ring magnet of ferroxdure. Systems of this kind are built up from a ring magnet, two soft-iron end plates and a soft-iron core, as illustrated in *fig. 1*. When the dimensions of the system are given, it is possible, using the method to be described, to calculate with sufficient accuracy both the induction in the air gap ($\pm 3\%$) and the induction in the ferroxdure ($\pm 5\%$). The method also makes it

possible to determine the dimensions of the magnet in order to produce a specified induction in the air gap using a magnet of minimum volume. Leakage, etc., is taken into account by means of two coefficients which can be directly derived from the data of the system with the aid of formulae obtained

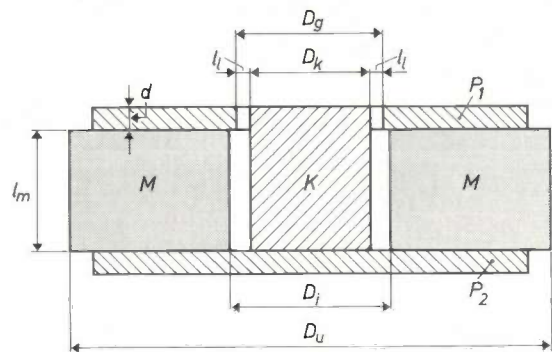


Fig. 1. Magnet system of normal construction for loudspeakers equipped with a ring magnet *M* of anisotropic ferroxdure. A system of this kind is generally flatter than one using a "Ticonal" magnet. The diameters of the soft-iron top plate P_1 and base plate P_2 are somewhat smaller than the outside diameter D_u of the magnet; this gives a gain of a few per cent for the flux in the air gap l_1 between the top plate and the soft-iron core *K*.

empirically from measurements on systems as defined in *fig. 1*. The method should not therefore be expected to yield good results for systems built in different ways. A loudspeaker with a ferroxdure magnet is shown in *fig. 2*.

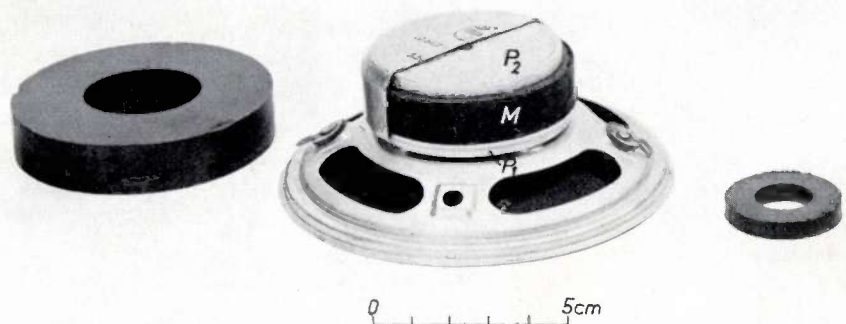


Fig. 2. A loudspeaker with a ferroxdure magnet. Part of the cap has been removed to show the magnet system. Meaning of letters as in *fig. 1*. The ferroxdure rings at either side give an idea of the maximum and minimum dimensions commonly found in loudspeakers. For special purposes larger rings may be used.

*) Industrial Components and Materials Division, Philips, Eindhoven.

¹⁾ For example, G. Hennig, *Dauermagnettechnik*, Franzis-Verlag, Munich 1952, p. 49 ff. Also: A. Th. van Urk, *The use of modern steels for permanent magnets*, Philips tech. Rev. **5**, 29-35, 1940; R. J. Parker and R. J. Studders, *Permanent magnets and their application*, chapter 4, Wiley, New York 1962.

For determining the dimensions of the magnet with minimum volume we have constructed a number of graphs. The calculations for this were carried out in Philips Computing Centre, using a type IBM 650 electronic computer. In the following we shall first present these graphs and comment on their use. We shall then describe how the formulae for designing the loudspeaker system were obtained, and how the graphs for magnets of minimum volume were derived from these formulae.

In connection with the graphs it should be noted that the ring magnet of a loudspeaker will frequently not be designed with dimensions calculated on the basis of minimum volume. Generally speaking, minimum volumes lead to very flat rings with large outside diameters, and with regard to production costs and the necessity for large end plates, they are not the most economic proposition. It is often more important to use ring magnets of standard dimensions than to be extremely sparing with ferroxdure. Even so, it is useful to know the dimensions for minimum volume, since this makes it possible to judge whether the ring magnet designed is unduly large.

Graphs for minimum-volume designs

The graphs from which the dimensions of the magnet with minimum volume can be found are reproduced in fig. 3. In all three graphs the height l_m of the ring is a parameter, while the abscissa is the magnetic resistance R_1 of the air gap multiplied by the permeability of vacuum, μ_0 ($= 4\pi \times 10^{-7}$ H/m):

$$\mu_0 R_1 = \frac{l_1}{\pi D_m d} \text{ m}^{-1} \dots (1)$$

Here $l_1 = \frac{1}{2}(D_g - D_k)$, $D_m = \frac{1}{2}(D_g + D_k)$ and d represent the dimensions of the air gap, expressed in metres (see fig. 1). Plotted vertically in fig. 3a, b and c are respectively the magnetic flux Φ_1 in the air gap, the cross-sectional area S_m of the ring perpendicular to the axis, and the magnetic induction B_m in the magnet. If a particular induction is required in an air gap of given dimensions, then Φ_1 and $\mu_0 R_1$ are given; e.g. $\mu_0 R_1 = 3.5 \text{ m}^{-1}$ and $\Phi_1 = 3.7 \times 10^{-4} \text{ Wb}$ ($3.7 \times 10^4 \text{ maxwell}$). From fig. 3a it follows that $l_m = 12 \text{ mm}$, from fig. 3b we read: $S_m = 41 \text{ cm}^2$ and from fig. 3c: $B_m = 0.217 \text{ Wb/m}^2$ (2170 gauss). For constructional reasons the inside diameter D_i of the ring should be 5 or 6 mm greater than the core diameter D_k . The latter is given as a measure of the annular air gap. In this way the dimensions of the ring are fixed.

If we make $D_i - D_k$ much larger than the 5 to

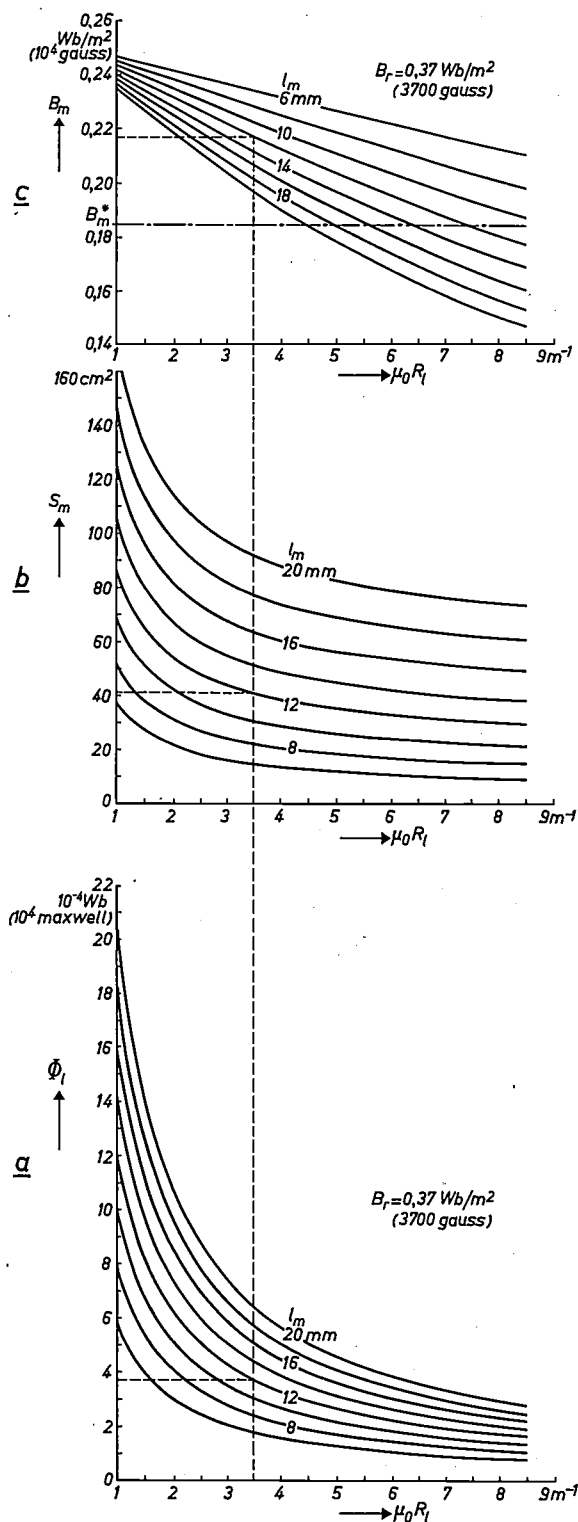


Fig. 3. Graphs for use in designing a ferroxdure ring magnet of minimum volume for a loudspeaker system as in fig. 1. Explanatory comments are given in the text. The graphs relate to ferroxdure 300 R with a remanence B_r of 0.37 Wb/m^2 (3700 gauss). For another value of B_r , e.g. $B_r = 0.38 \text{ Wb/m}^2$, the numbers on the vertical axis in (a) and (c) should be multiplied by 38/37. In (c), B_m^* is the value of the induction B_m in the magnet at which $B_m \times H_m$ is maximum. If the knee in the demagnetization curve (see fig. 5) is at or below $\frac{1}{2} B_r$ (which is assumed in (c)), then $B_m^* = \frac{1}{2} B_r$.

6 mm mentioned, the flux in the air gap will be smaller than required. The construction then no longer belongs to the cases regarded in this article as normal.

Calculation procedure for loudspeaker magnets

To explain the calculation procedure for loudspeaker magnets it is useful to start from the equivalent magnetic circuit shown in fig. 4, which can

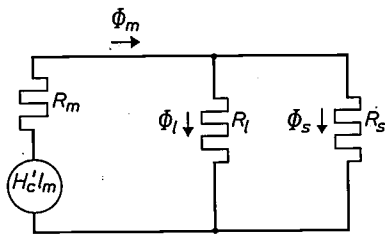


Fig. 4. Equivalent (idealized) magnetic circuit. The magnet is represented by a magnetomotive force of magnitude $H_c'l_m$, in series with an internal reluctance R_m . The external reluctances R_1 and R_s represent respectively the air gap and the path for the leakage flux.

be derived from the conventional theory of magnetic circuits²⁾. The magnet can be regarded as a source of magnetomotive force $H_c'l_m$ in series with an internal resistance (reluctance) R_m , given by:

$$R_m = \frac{l_m}{\mu_0 \mu_r S_m} \dots \dots \dots (2)$$

The significance of H_c' appears from fig. 5, which gives the demagnetization curve of the material used, anisotropic ferroxdure 300 R. The height l_m and the cross-sectional area S_m of the ring magnet are expressed in m and m^2 respectively, and H_c' in A/m. The factor μ_r , representing the relative permeability of the ferroxdure, is equal to the cotangent of the angle which the demagnetization curve makes with the B_m axis when B_m and $\mu_0 H_m$ are plotted on the same scale, as they are in fig. 5. Since this angle is, on average, 43° for ferroxdure 300 R, μ_r is 1.07. For simplicity we put $\mu_r = 1$ in (2), so that

$$R_m = \frac{l_m}{\mu_0 S_m} \dots \dots \dots (3)$$

The fact that the magnetomotive force produced by the magnet must have the value $H_c'l_m$ can be understood by considering two extreme cases, viz the case of a magnetic short-circuit (no air gap) and the case where the reluctances R_1 and R_s are infinitely large. In the first case, according to the equivalent circuit and formula (2), the flux assumes the value

$\Phi_m = \mu_0 \mu_r H_c' S_m = B_r S_m$, which means that the working point corresponds to the remanence. In the second case $\Phi_m = 0$, and hence $B_m = 0$, and moreover the tension has the "open-circuit value" $H_c'l_m$, so that the field strength in the magnet assumes the fictive value H_c' (fig. 5). This corresponds to a working point at H_c' . The equivalent circuit is accordingly based on a straight demagnetization line extrapolated to the point H_c' .

The magnetomotive force causes a magnetic flux through the parallel reluctances R_1 and R_s , which represent respectively the air gap and the path for the total leakage flux. From fig. 4 the flux through the air gap is seen to be given by:

$$\Phi_1 = \frac{H_c'l_m}{pR_1 + R_m} \dots \dots \dots (4)$$

where

$$p = 1 + \frac{R_m}{R_s} \dots \dots \dots (5)$$

and the flux Φ_m through the magnet is given by:

$$\Phi_m = \frac{H_c'l_m}{\frac{R_1}{p} + R_m} \dots \dots \dots (6)$$

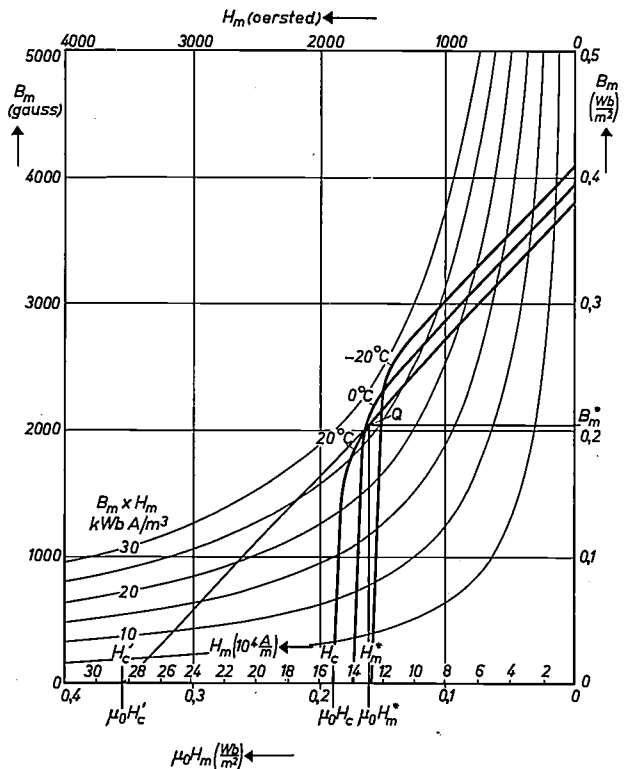


Fig. 5. Demagnetization curve of a sample of ferroxdure 300 R, as used for loudspeaker magnets, at three temperatures. The magnetic state of the material is characterized by a point on the curve, called the working point. This should be above the knee. The inserted hyperbolae, for which the product $B_m \times H_m$ is constant, make it possible to read the value of this product for every point on the demagnetization curve. For the $20^\circ C$ curve $B_m \times H_m$ has a maximum value of about 26.5 kWbA/m^3 (about $3.35 \times 10^6 \text{ gauss-oersted}$) if $B_m = B_m^*$ (about 0.2 Wb/m^2).

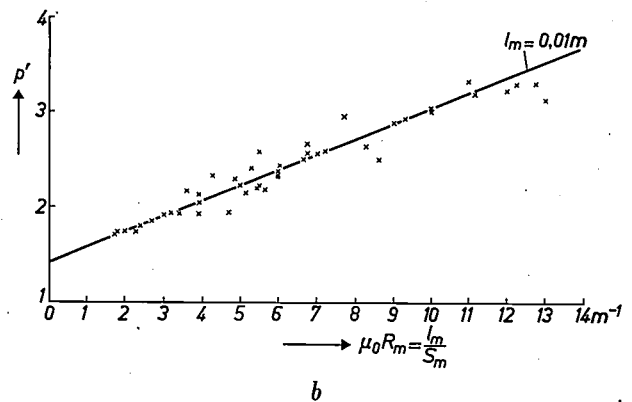
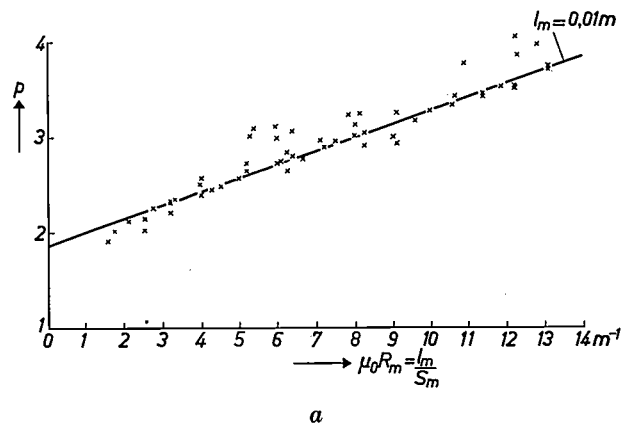
²⁾ See e.g. J. Fischer, Abriss der Dauermagnetkunde, Springer-Verlag, Berlin 1949, p. 103 ff. Also: P. Cornelius, Electrical theory on the Giorgi system, Clever Hume Press, London 1961, p. 145.

where
$$p' = 1 + \frac{R_l}{R_s} \dots \dots \dots (7)$$

For the reluctance R_l of the air gap we use the expression (1), but a simple expression cannot be obtained for the leakage reluctance R_s . If only for that reason, the formulae cannot immediately be used for calculations. Moreover, the use of an equivalent circuit, as in fig. 4, implies a number of simplifying assumptions, which are certainly not correct. Principal among these is the assumption that everywhere in the magnet the field strength H_m and the induction B_m are uniform and parallel to the axis of the magnet. In that case the conventional characterization of the magnetic state of the magnet by a single "working point" on the demagnetization curve (fig. 5) is indeed possible. It is then also assumed that the flux Φ_m through the magnet has the same value in every cross-section. In a certain case, where ferroxdure rings were used in loudspeaker systems, we found by measurements that the flux through the magnet decreased from the foot to the top by 20%; the field in the magnet is thus not as a rule homogeneous. This has also been found by Schwabe, who actually measured the field strength H_m in a system as in fig. 1³⁾. It is therefore not really possible to characterize the whole magnet by a single point on the demagnetization curve⁴⁾. A relation that ought to apply is $H_m l_m = H_l l$ (H_l is the magnetic field strength in the air gap), since both expressions represent the magnetic potential difference between the pole pieces, but this is found to be equally inapplicable. $H_m l_m$ was found by measurement 10 to 35% greater than $H_l l$, which is too great a difference to be explained by the reluctance of the iron pole pieces. It is not surprising that such a simplified theory cannot be a sound basis for calculations that lay claim to some accuracy.

With the object of examining experimentally whether it was possible to modify formulae (4)-(7) so as to make them usable for calculations, we collected data of flux measurements of Φ_l and Φ_m in the air gap of dozens of sizes of loudspeaker systems (such as in fig. 1), and extended these data with measurements specially designed for this purpose. We first reduced the results to a value of

0.37 Wb/m² for the remanence B_r at 20 °C, using the temperature coefficient -0.2% per °C applicable to the remanence of ferroxdure, and assumed Φ_l and Φ_m to be proportional to B_r . Using formulae (4) and (6) we then computed p and p' for all systems. Formulae (5) and (7) for p and p' were disregarded. When we then plotted p as a function of the reluctance R_m of the magnet, we found that points pertaining to the same value of magnet height l_m lay roughly on a straight line. That this was not a coincidence was confirmed by reducing all measurements to the same value of l_m (see fig. 6a; for convenience, $\mu_0 R_m$ is plotted here instead of R_m). We were able to do this by making use of the property that uniform scaling-up of the system does not affect the inductions at corresponding points. It follows, then, that the values calculated for p and p' are independent of the uniform scaling-up of the system, for scaling-up a system linearly by a factor a means multiplying Φ_l and Φ_m by a^2 , l_m by a and R_l and R_m by $1/a$. From (4) and (6) it then follows that p and p' do not change. To transform a point pertaining to a given system to a point pertaining to a similar



³⁾ E. Schwabe, Über die Temperaturabhängigkeit der magnetischen Eigenschaften von Bariumferrit, Z. angew. Physik 9, 183-187, 1957, especially page 186.

⁴⁾ In most points of a magnet B_m and H_m will not be exactly parallel. To establish the magnetic state at a particular point, the components of the vectors B_m and H_m would have to be known. Even for a single point in the magnet the characterization of the magnetic state by a point in the B_m - H_m plane is therefore only an approximation.

Fig. 6. Clusters of points obtained by plotting as a function of $\mu_0 R_m$ the values of p and p' found from measurements on loudspeaker systems, after reducing all measurements to the same height $l_m = 0.01$ m of the magnet. a) Points for p , b) for p' .

system which is linearly larger by a factor a (and where R_m is consequently smaller by the same factor), it is therefore only necessary to divide the abscissa R_m of that point by a (fig. 7). The points for p in fig. 6a were obtained after reduction to $l_m = 10$ mm. The best straight line through these clusters of points is found to give a satisfactory basis for calculations.

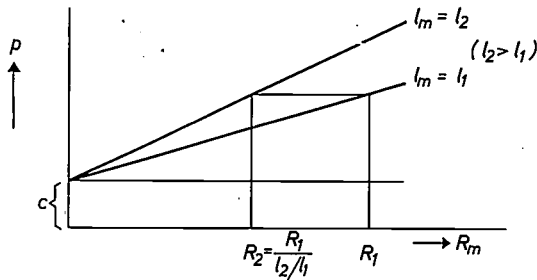


Fig. 7. To derive the straight line for $l_m = l_2$ from the line for $l_m = l_1$ (l_1 and l_2 being constants), all R_m values must be divided by $a = l_2/l_1$. The section c which the lines cut from the p axis is thus independent of l_m , while the slope is proportional to l_m . The equations for the straight lines therefore have the form $p = bl_m R_m + c$, where b and c are constants.

From the above transformation rule it follows that the equation for the straight line must have the form (see fig. 7):

$$p = b l_m R_m + c,$$

where b and c are constants. With due alteration of details the same applies for p' as for p (fig. 6b). The following expressions were found for p and p' :

$$p = 14.2 \mu_0 R_m l_m + 1.86 = 14.2 \frac{l_m^2}{S_m} + 1.86, \quad (8)$$

$$p' = 16.2 \mu_0 R_m l_m + 1.42 = 16.2 \frac{l_m^2}{S_m} + 1.42. \quad (9)$$

p and p' are the coefficients referred to on p. 150. It will be seen that they can indeed be calculated directly from the dimensions of the system. Once this is done, the fluxes Φ_1 and Φ_m in the air gap and the magnet follow from (4) and (6). The inductions B_1 in the air gap and B_m in the magnet are found by dividing Φ_1 by the area of the air gap $\pi d D_m$, and Φ_m by the cross-sectional area S_m of the magnet.

Since, as mentioned on p. 153, the flux through the magnet does not have the same value in every cross-section, it is necessary to agree on the cross-section in which the flux Φ_m is to be measured. In our case this was always at half height of the ring. B_m was calculated by dividing this measured flux by S_m .

Agreement is also necessary on the measurement of Φ_1 in the air gap, since the field in the air gap is not exactly uniform.

We performed this measurement ballistically, using a search coil of height $h = \frac{2}{3}d$ (see fig. 8; d is the thickness of the pole plate — see fig. 1). The coil is introduced symmetrically into the air gap and then withdrawn. Division of the measured flux by the surface area of the coil gives B_1 , and multiplication of this induction by $\pi d D_m$ gives Φ_1 .

For reliable measurements it proved necessary to make core and base plate (fig. 1) from one piece and to make sure that the top face of the core fitted exactly flush with the top plate. It is also necessary to ensure that the induction in the soft iron is nowhere higher than about 12 000 gauss, as otherwise the reluctance of the iron becomes significant.

Caution is called for when comparing results with those of measurements performed elsewhere under somewhat different conditions. If, for example, the flux Φ_1 in the air gap is measured with a coil of exactly the same height as the gap ($h = d$), the value of B_1 then found is a few per cent lower than ours.

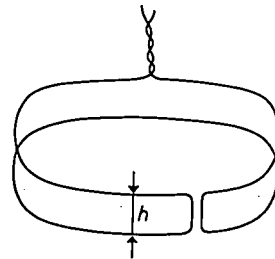


Fig. 8. Form of the search coil used for measuring the flux density in the air gap. The height is $h = \frac{2}{3}d$, where d represents the thickness of the top plate (see fig. 1).

One of the assumptions underlying the equivalent magnetic circuit sketched in fig. 4 is that when H_m is varied the "working point" of the ferroxdure magnet in the H_m - B_m plane (fig. 5) moves along a straight line that intersects the H_m axis at the point H_c' and the B_m axis at the point of remanence²⁾. In fact, however, the working point shifts along the thickly drawn line, which coincides with the straight line only as far as the knee of the curve, at point Q , and which intersects the H_m axis at the point H_c , which represents the coercivity of the material. If the calculation of B_m yields a value corresponding to a point in the portion $H_c'Q$ of this straight line, the values for Φ_1 and Φ_m found in the manner described are too high. It is therefore not sufficient merely to calculate the induction B_1 in the air gap, but B_m must also be determined in order to check whether the working point lies above the knee.

It is even undesirable for the working point to be just above the knee. In the first place, as explained, the working point found is merely an average over the whole magnet. There will certainly be regions where the working point is lower than the average. Secondly, if the magnet has a working point close to the knee of the curve, it is magnetic-

ally not very stable: there is a risk of disturbances from an external field moving the working point into the curve of the knee or below it, resulting in permanent weakening of the magnet ⁵). The risk is greatest, of course, in domains in the magnet where the working point is lower than the average. Thirdly, the demagnetization curve is displaced by a change of temperature (fig. 5). If the working point is close above the knee, it can drop below the knee upon a fall in temperature, again resulting in permanent weakening of the magnet ⁶).

Calculations for minimum volume of the magnet

To find the required minimum volume of the magnet we begin by deriving a general expression for the volume. By inserting $R_m = l_m/\mu_0 S_m$ (3) in the expression (4) for the flux Φ_1 in the air gap we obtain:

$$\Phi_1 = \frac{\mu_0 H_c' l_m}{p \mu_0 R_1 + \frac{l_m}{S_m}}$$

The dimensions of the air gap and the induction required in it are known. $\mu_0 R_1$ (see (1)) and Φ_1 are therefore known, and so is the material constant $\mu_0 H_c'$. After insertion of expression (8) for p in the above formula for Φ_1 we solve S_m and find:

$$S_m = \frac{\Phi_1(14.2 \mu_0 R_1 l_m^2 + l_m)}{\mu_0 H_c' l_m - 1.86 \Phi_1 \mu_0 R_1} \quad (10)$$

In this way the cross-sectional area S_m of the ring magnet is expressed in terms of ring height l_m ; the other quantities are all known. Multiplication by l_m gives the expression required for the volume V_m of the magnet:

$$V_m = l_m S_m = \frac{\Phi_1(14.2 \mu_0 R_1 l_m^3 + l_m^2)}{\mu_0 H_c' l_m - 1.86 \Phi_1 \mu_0 R_1} \quad (11)$$

We must now examine whether V_m shows a minimum as a function of l_m , and at what value of l_m this minimum occurs. For this purpose we differentiate V_m with respect to l_m and equate the result to zero, which gives the following quadratic equation for l_m :

$$28.4 \mu_0 H_c' l_m^2 + \left(\frac{\mu_0 H_c'}{\mu_0 R_1} - 79.2 \Phi_1 \mu_0 R_1 \right) l_m - 3.72 \Phi_1 = 0 \quad (12)$$

Equation (12) always contains only one positive solution l_m , and this is found to correspond to a minimum value of V_m . By solving (12) to find l_m and then calculating S_m from (10) we can find the dimensions of the ring magnet with minimum volume.

The graphs in fig. 3 were calculated by a somewhat different procedure. Eq. (12) was solved not to find l_m but Φ_1 , giving:

$$\Phi_1 = \frac{28.4 \mu_0 H_c' l_m^2 + (\mu_0 H_c' / \mu_0 R_1) l_m}{79.2 \mu_0 R_1 l_m + 3.72} \quad (13)$$

For any set of values for l_m and $\mu_0 R_1$, the value of Φ_1 can be found from (13). Using this result we can also calculate S_m from (10). We can then find p' from (9), followed by Φ_m from (6), after which division by S_m finally yields B_m — all with the assumed values of l_m and R_1 . By carrying out this computation of Φ_1 , S_m and B_m at a fixed value of l_m for a series of values of $\mu_0 R_1$, and repeating this whole procedure for a number of values of l_m , we obtain the families of curves shown in fig. 3a, b and c.

The calculations were performed with a computer of the type IBM 650 for values of l_m increasing in steps of 10^{-3} m from 6×10^{-3} to 24×10^{-3} m; for each of these values of l_m the quantity $\mu_0 R_1$ increased in steps of 0.5 m^{-1} from 1 to 9 m^{-1} . For B_r we chose 0.37 Wb/m^2 , so that $0.37 \tan 43^\circ \text{ Wb/m}^2$ had to be inserted for $\mu_0 H_c'$ (see fig. 5).

The calculation consists of a succession of elementary operations: multiplication, addition, division and subtraction, and could therefore be directly programmed for the computer. The value of 0.37 Wb/m^2 for B_r was chosen because it is the minimum value for ferroxdure 300 R. The real value, however, will nearly always be higher, e.g. 0.38 Wb/m^2 . In that case $\mu_0 H_c'$ will be proportionately higher. If we examine the described procedure for calculating Φ_1 , S_m and B_m at assumed values of l_m and $\mu_0 R_1$, it is easily seen that S_m is independent of B_r , whereas Φ_1 and B_m are directly proportional to it. If, then, B_r is 0.38 instead of 0.37 Wb/m^2 , the graphs can easily be adjusted to this case by multiplying the figures on the vertical axis in fig. 3a and c by $38/37$.

It appears from fig. 3c that the value of B_m applicable to the magnet designed for minimum volume does not generally coincide with the value B_m^* at which the product $B_m \times H_m$ is maximum. This value is usually greater than B_m^* .

The graphs in fig. 3 are based on a straight demagnetization curve extrapolated to the H_m axis. In that case $B_m^* = \frac{1}{2} B_r$. In fig. 3c, B_m^* is therefore taken as 0.185 Wb/m^2 . Frequently the point Q where the knee begins will be higher than $\frac{1}{2} B_r$

⁵) See e.g. p. 19 ff. of the book by G. Hennig, mentioned in footnote ¹).

⁶) F. Tomholt and E. Haes, Temperature dependence of the magnetic properties of ferroxdure 2, Philips Matronics No. 13, 225-228, December 1957; see also the article by E. Schwabe, mentioned in footnote ³).

(as it is, for example, in fig. 5). B_m^* then coincides with the value of B_m which corresponds to Q and is thus somewhat higher than $\frac{1}{2}B_r$.

From the results collected in fig. 3 we can now see that, in view of the simple form of the curves obtained, it would have been sufficient to calculate a much smaller number of points. The shape of the

curves was not, however, known beforehand, nor was it estimated. Since the machine does the donkey work, we are not concerned about one point more or less. The fact that only a small portion is used of the enormous quantity of data delivered by an electronic computer is nothing unusual in physical and technical problems.

Summary. Discussion of a calculation procedure for designing ferroxdure loudspeaker magnets. The leakage and other factors difficult to account for theoretically, such as non-uniformity of the induction in the magnet, are expressed by two coefficients that can be determined from the dimensions of the magnet using experimentally derived formulae. For this purpose data of flux measurements were collected and special experiments designed. With the aid of these coefficients the induction in the air gap of a normal system with given dimen-

sions can be computed with an accuracy of about 3%. Formulae are also derived for calculating the dimensions of a ring magnet of the minimum volume which is required to give a specified induction in a given air gap. These formulae are used as the basis for three graphs from which the dimensions and also the pertaining induction B_m in the magnet can be found. In most cases B_m is seen to be greater than the value corresponding to $(BH)_{\max}$. The calculations for the graphs were carried out in Philips Computing Centre with an IBM 650 electronic computer.

SOLVING A CHESSBOARD PUZZLE WITH THE PASCAL

by A. J. DEKKERS *) and A. J. W. DUIJVESTIJN *).

681.14-523.8:685.854

Introduction

The preceding articles in this issue illustrate the use of the electronic computer for various mathematical problems, such as working out arithmetical formulae, solving equations, integrating differential equations, etc. The fact that the machine has to *compute* in all these problems seems self-evident, but this is not so: it can also be used for non-arithmetical tasks, or at least for tasks which in the first instance have nothing to do with computation. Examples are to be found in the various fields of mathematics itself, such as in topology, abstract algebra (group theory), etc. Then there are "logical" problems and routines such as sorting and collating, which are frequently encountered outside the realm of mathematics, especially in the accounting department of a business enterprise. This is reflected in the installation of two nearly identical computers in Philips Computing Centre — the PASCAL and the STEVIN — the PASCAL being employed for mathematical work proper, and the STEVIN almost exclusively for the various administration departments of Philips¹⁾.

The distinction between arithmetical and non-arithmetical problems is perhaps not strictly tenable. This was the subject of learned debate long before calculating machines existed: some defended the postulate that logical reasoning is nothing but a kind of computation, while others argued that computation is simply a kind of logical reasoning. If one considers the operations which the computer is made to perform in solving all kinds of problems, one may again conclude that the above distinction is indeed superficial. In fact, with all problems which are commonly called non-arithmetical, one is sure to find in the pertaining machine programme the operation of *counting*, probably even more than once. Now, it is evident that this operation is also at the base of the operation of addition and hence of the other mathematical operations proper.

With a view to exploring all the computer's possibilities, it can be useful to attempt to programme it for widely diverse problems, including strictly non-arithmetical ones — again in the commonly under-

stood sense of the word. For example, at the "Studiecentrum voor Administratieve Automatisering" in Amsterdam a programme, commissioned by Euratom, is being designed under the direction of Dr. M. Euwe which will enable the computer to play a game of chess. We have chosen a simpler problem as an object for study — a chessboard puzzle. A procedure for solving this puzzle has been worked out and programmed for the PASCAL. This procedure is discussed below.

Description of the puzzle

The chosen puzzle will now be described²⁾. A chessboard is divided up into 12 pieces as shown in *fig. 1*. Pieces 5a and 5b happen to be identical; all other pieces differ from these and from each other. The problem is to fit these 12 pieces together to form a complete chessboard with the proper alternation of black and white squares. Is there more than one solution? If so, give all possible solutions.

The solution can only be found by *trial and error*, but the number of possibilities to be examined is enormous. Therefore first of all we must systematize the trial procedure so that no possibility is overlooked. Secondly, there is little hope of finding the

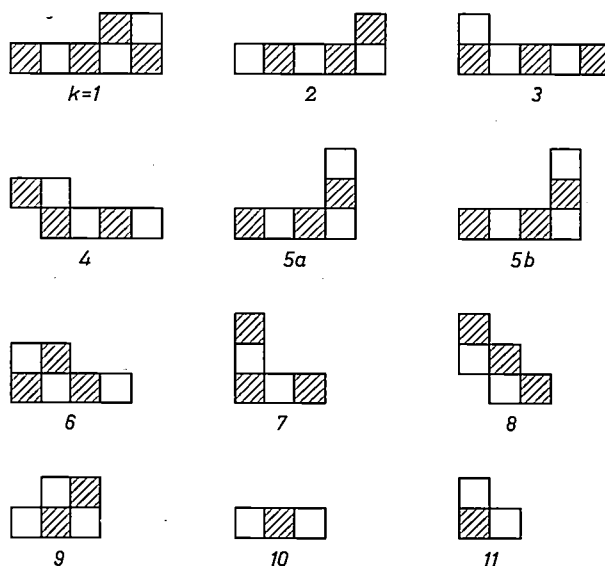


Fig. 1. The 12 pieces into which the chessboard is divided. There are 11 different shapes; one shape (No. 5) occurs twice.

*) Philips Computing Centre, Eindhoven.

¹⁾ See W. Nijenhuis, *The PASCAL*, a fast digital electronic computer for the Philips Computing Centre, Philips tech. Rev. 23, 1-18, 1961/62 (No. 1).

²⁾ The puzzle is one of the Peter Pan series marketed by Cecil Coleman Ltd.

complete solution of the puzzle without the help of an electronic computer, for, as will presently be seen, it is necessary to examine 3 million possibilities, yielding *eleven* different solutions.

Method of solving the puzzle

Systematization of the trial procedure mentioned reduces to trying to cover the 64 squares of an empty chessboard in a fixed order by laying down, again in a fixed order, the given 12 pieces, each in their different possible orientations.

It is essential in this procedure to keep a complete *record* of the squares already filled and of the pieces and orientations already used and tried. For, when each piece is laid (i.e. upon every "move") it is necessary to see that we are not in conflict with the requirements, e.g. that the piece is not covering a square that has already been occupied. If we come to the conclusion that none of the remaining pieces can cover the open square whose turn it now is, without somewhere coming into conflict, then the piece laid in the previous move must be withdrawn and the same piece tried in its next orientation (in the prescribed order), or the next piece tried. When all possibilities have been exhausted with the still available pieces, we must revise the move before last, and so on. All data on every move must therefore be preserved.

For this purpose we can of course make use of the large *memory* of the PASCAL.

We shall now first discuss the systematization of the trial procedure adopted, give an example of how it works, describe how the machine "tries" a piece and how the data are recorded, and finally examine the solving process as a whole.

The systematization

The systematization consists in arranging the possibilities in order of the series of whole numbers, so that a "call" for each possibility can be made by simply increasing an address number in a memory by 1.

To this end the first step is to give the squares on the chessboard serial numbers n . We call the bottom right square $n = 1$, and in our approach to the problem we decide, quite arbitrarily, that this shall be a black square (with due apologies to the chess players among our readers). The numbers n on the bottom row increase from right to left, then on the next row from right to left again, and so on.

The next step has already been seen in fig. 1, where the eleven *dissimilar* pieces are numbered $k = 1, 2, \dots, 11$. (The complication that piece No. 5 is

duplicated will presently be taken into account in a simple manner.)

In principle each piece can be laid on the board in four distinct orientations, each turned through 90° . All the possibilities thus obtained we place in a row and give them serial numbers j . Each of these possibilities we call a "lay" — a piece that is laid experimentally in a particular orientation on the board.

The row of available lays is shown in fig. 2. It can be seen that we have in reality made *two* rows of lays, the "black" and the "white", depending on the colour of the "master square", i.e. the square at the right end of the bottom row of the piece. A lay is used such that this square is laid on the board on the empty field (to avoid confusion a square of the board will be called a field) whose turn it is to be tried (the "trial field"). If that field is white the machine must try a "white lay"; if the trial field is black then the lay has to be black. Further, it will be seen that only two of the four possible orientations of piece No. 1 (fig. 1) are included in the row (black lay 1 and white lay 1). This ensures that once a solution has been found there is no search made later in the solving process for the trivial reverse form of that solution (the form turned 180°). Piece No. 10 (fig. 1) also occurs only twice in the row of lays, but that is simply because only two distinct orientations of this piece are possible (white lays 17 and 18). Altogether there are now 19 black and 21 white lays that have to be tried in their numerical order, j_b and j_w respectively.

Example of attempt at a solution

To illustrate the trial procedure a situation is presented in fig. 3 which arises after a few moves at the very beginning of the solving process. Lay $j_b = 1$ (piece $k = 1$) has been placed on the first trial field $n = 1$, then lay $j_w = 3$ (piece $k = 2$) on the next available trial field $n = 6$, which is white. The next available lay $j_b = 4$ does not go on the next (black) trial field, $n = 7$, for it would project beyond the board; nor do lays $j_b = 5, 6, 7, 8$ and 9, but $j_b = 10$ (piece $k = 6$) does go. The next trial field is on the second row and is white. Lay $j_w = 4$ is not a fit, nor is $j_w = 5$, but $j_w = 6$ is (piece $k = 4$). With lay $j_b = 5$ ($k = 3$) on the next trial field, followed by $j_w = 18$ ($k = 10$), the second row on the board is filled. On the first trial field on the third row we place lay $j_b = 19$ (piece $k = 11$), but now we can go no farther: for the next trial field, which is the last on the third row, there is no suitable lay available.

The action to be taken has already been indicated. The last piece, No. 11, must be withdrawn from the

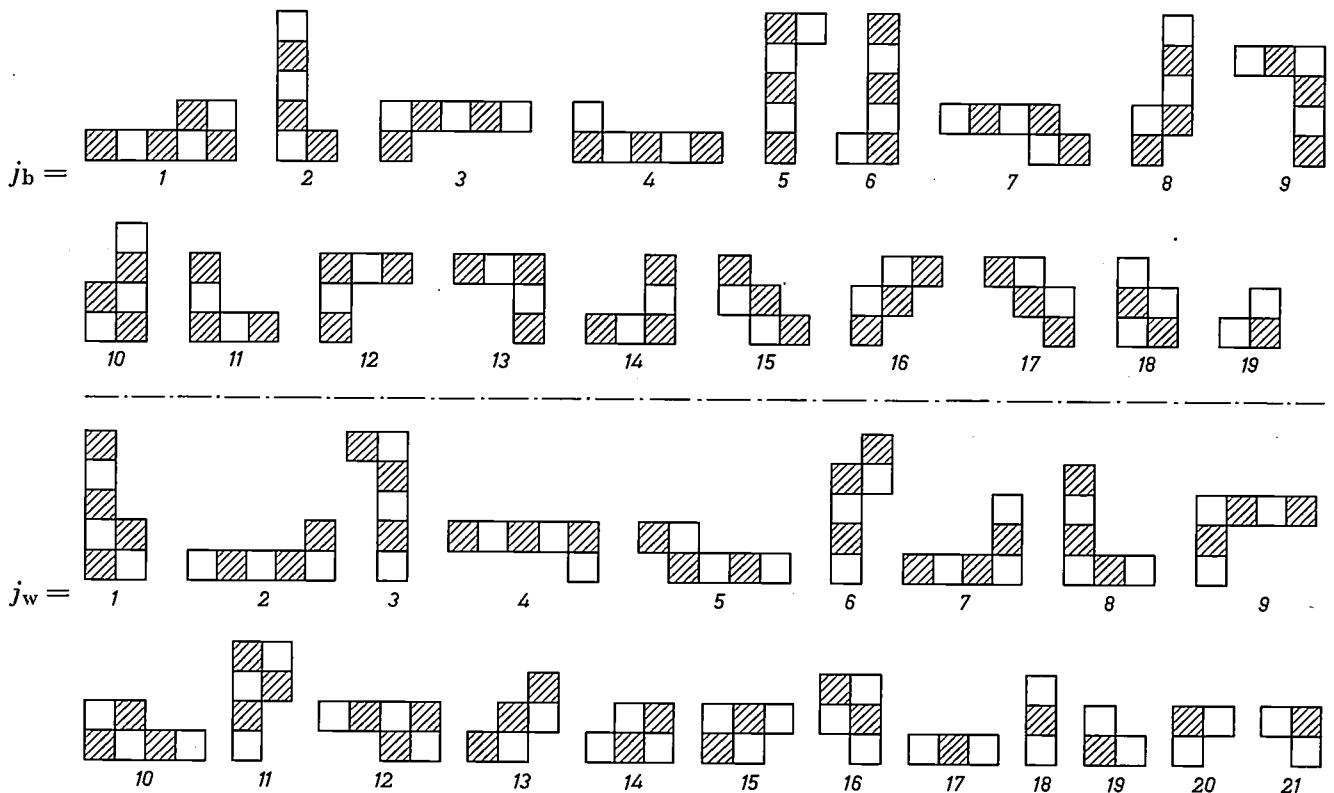


Fig. 2. The complete row of "black" and "white" "lays" (pieces laid in a particular orientation) that must be tried successively in a "move".

board (and the machine must amend all the appertaining steps in the records, which we shall presently discuss). We are now back in the situation of having to fill the preceding trial field, and to do so we must try the next lay j_b after $j_b = 19$. But no such lay exists, and therefore we must also cancel the step preceding the one just cancelled, and so on.

Realization of the trial procedure by the computer; keeping the records

In the example just discussed we have assumed a human operator who can check whether a lay, "called upon" in accordance with the systematic procedure, can be placed or not. The criteria are (1) that no field already occupied can be occupied again, and (2) that no square of a lay may project outside the board.

How can the computer make this check? The checking is linked with the keeping of records, and we therefore cannot avoid describing also part of the latter process.

Each field on the board is assigned a memory location with an address number n . All empty fields receive in their address an arbitrary negative number, viz -1 . When a field is occupied this number is replaced by a positive one, viz the serial number (k)

of the piece in occupation. The changes in the "contents" of the fields, $V(n)$, are recorded by the machine, and this part of the memory forms, as it were, an "operating list".

A field can only be occupied if its address contains -1 . The machine can immediately verify this and thus see that the first criterion is satisfied.

			3	3	4	
	2	2		3	4	4
	6	2	10	3	4	11
	6	2	10	3	4	11
	6	6	2	10	3	4
	6	6	2	1	1	1
	6	6	2	1	1	1

Fig. 3. Situation after the first eight moves have been made in the systematic procedure described. The fields (squares of the board) covered are marked with the number k of the piece covering them. This situation is a dead end, and the last lay must be withdrawn, and even the lay before the last, in order to try further lays.

To enable the computer to check whether the second criterion is satisfied, i.e. that a lay must not project outside the board, we add an extra row of fields along three edges of the board (no extra row is needed along the bottom edge because of our system of laying the pieces). These fields are also assigned memory locations, in which is placed an arbitrary positive number, e.g. +20. The computer can now automatically reject any lay that would cover such a field.

Along the left-hand edge of the board we have added not one but two extra rows of fields. This makes it a very simple matter for the machine to ascertain whether a particular trial field is white or black — it needs this information to decide whether it should try a white or a black lay. With these additions the complete board on which the machine operates is as shown in fig. 4; it has 96 fields and the numbering is consecutive along each row, that is $n = 1, 2, \dots, 96$. It is easily verified that, as a result of this device, all odd fields must be black. The difference between odd and even can immediately be established by a binary machine like the PASCAL by simply checking one bit.

			96	95	94	93	92	91	90	89	88
87	86	85	84						
						24	23	22	
21	20	19	18	17	16	15	14	13	12	11	
10	9	8	7	6	5	4	3	2	1		

Fig. 4. The board on which operations are made. Extra (fictive) rows of fields are added to the chessboard along three edges so that the machine can test whether a lay projects outside the chessboard and ascertain quickly whether a field should be covered with a black or with a white square of a piece.

We have outlined how the machine handles the two criteria. To “try the lays” the machine must be able to determine which fields are occupied by a lay when it is laid. For this purpose we supply the machine with a description of each lay in the form of a row of numbers, which are as it were the co-ordinates of the squares on the lay. This can be done in various ways. In our case we consecutively number the squares of a lay row by row (serial number r) in

the same way as the fields on the board. If for example we take lay $j_b = 1$ (see fig. 2) and place the “master square” of the lay on field 1, then the numbers of the squares will be 1, 2, 3, 4, 5, 12 and 13, and the differences between these numbers give the co-ordinates of the lay, which are: 1, 1, 1, 1, 7, 1, 0. The “master square” of the lay of course needs no co-ordinate. The 0 at the end of the row is added for each lay to indicate that there is no other square to come. It is obvious that no matter where the lay is placed on the board, it will always have the same co-ordinates. The 19 + 21 rows of numbers $U_b(j_b, r)$ and $U_w(j_w, r)$ which describe our black and white lays are stored in the form of two lists in the machine’s memory. Unlike the earlier mentioned list $V(n)$, these are not “operating lists” but “material lists” with permanent contents.

We can now describe how the machine tries the lay whose turn it is. It starts from the number n_0 of the trial field whose turn it is and it adds one after the other all numbers $U(j, r)$ of the lay. In this way it finds the numbers of the series of fields that will be occupied by the lay. The stored content V at all these addresses is examined. If one positive content is found, the lay is wrong. If all addresses contain negative contents, the latter are all replaced by the number k of the lay. This completes the move.

The solving process as a whole; more record-keeping

The final action just mentioned no longer belongs to the trial procedure proper but to the machine’s record-keeping, the writing into list V , and we see at once that the lists mentioned so far are not yet sufficient: another list is evidently needed which will indicate the piece k appertaining to every lay j_b and j_w . We shall call this list $S_b(j_b)$ and $S_w(j_w)$ respectively. Moreover, to keep a complete tally we also need a move counter C (move number = m) and three further lists $M(k)$, $J(m)$ and $N(m)$ (see fig. 5), whose functions may be understood as follows.

When the machine has ascertained whether a given trial field n_0 is odd or even, it must proceed to try on it, one after the other, all black or white lays, as the case may be. Before trying a lay j_b or j_w , however, it must see whether the relevant piece is still available. This is the purpose served by list M , which has eleven memory locations corresponding to the piece numbers $k = 1$ to 11. In each location the machine keeps a record of how many pieces with that number are left; at the start of the solving process, location No. 5 contains a 2 and all other locations a 1. Thus, when the machine wants to try a lay, it first makes sure that the relevant location in M does not yet contain a 0.

S_b		U_b								V		C		M	
j_b	k	j_b	$r =$							n		m	k		
1	1	1	1	2	3	4	5	6	7	1	(-1)		1	(1)	
2	2	2	1	11	11	11	11	0		2	(-1)		2	(1)	
3	2	3	7	1	1	1	1	0		3	(-1)		3	(1)	
.		4	(-1)		4	(1)	
.		5	(-1)		5	(2)	
18	9	18	1	10	1	11	0			6	(-1)		6	(1)	
19	11	19	1	10	0					7	(-1)		.	.	
S_w		U_w								8	(-1)	J		N	
j_w	k	j_w	$r =$							9	+20	m	j	m	n_0
1	1	1	1	10	1	11	11	11	0	10	+20	1	1	2	1
2	2	2	1	1	1	1	7	0		11	+20	2	2	3	2
3	2	3	11	11	11	11	1	0		.	.	3	.	.	3
.
.
20	11	20	10	1	0					95	+20	12	12	12	12
21	11	21	11	1	0					96	+20				

Fig. 5. The complete set of "lists" used by the computer for solving the puzzle. Each list takes up part of the machine's memory. The $S(j)$ lists give for each black or white lay j_b and j_w the number k of the piece. Lists $U(j,r)$ contain the description of all lays, i.e. the co-ordinates of the squares r of each lay. The S and U lists are "material lists" with permanent contents. The others are "operating lists". In V is noted the content of the 96 fields of the board (serial number n); an empty field is given $V = -1$, the fictive fields have a permanent content $+20$, and each field covered by a piece k is given the content $+k$. List M contains information on how many pieces with the number k are still available (0, 1 or 2 — the latter only in the case of piece No. 5). At every move m the machine writes in list N the number n_0 of the trial field whose turn it is, and in list J the number j of the lay whose turn it is (positive for odd n_0 , i.e. for black lays, and negative for even n_0 , i.e. for white lays). The move number m is in the move counter C .

At every move m , the machine writes in list J , which has 12 locations, the number j of the lay being tried — positive for a black and negative for a white lay. If the lay turns out to be wrong, the machine looks up the relevant number noted in J , raises it by 1 and thereby obtains the address $j + 1$ of the next lay to be tried (of the same colour, of course).

The machine is not able to look up the number noted in J so easily as a human operator: it does not "know" offhand which location in J was occupied by the number j . This is where the move counter C comes in. After every successful move the number m in this counter (i.e. in the permanent memory location where the count is made) is raised by 1, and the new number serves as an address for writing and searching in $J(m)$.

It also serves as an address for using the last list, $N(m)$. In this at every move, the machine writes the number n_0 of the trial field on which it is about

to try lays. After a successful move the machine finds this number by checking consecutively through all fields in list $V(n)$ and taking the first in which it finds -1 .

When the machine is in the situation of having tried *all* lays in a move without result, it can now, by looking up the records, go back on its steps as was mentioned at the beginning. First, it withdraws the lay j tried in the last move and continues that move with the lay $j + 1$. In order to do this it decreases the number in the move counter by 1, that is from m to $m - 1$. Using the address $m - 1$ it looks up in list J the number j of the unsuccessful lay, and in list N the relevant trial field n_0 . At the address n_0 in list V it finds the number k of the piece employed. In this location it again writes -1 , and likewise in all locations in V which were covered by the unsuccessful lay j (and which it finds from lists S and U). At the same time it turns to list M and raises

by 1 again the wrongly lowered number at the relevant address k . All it now has to do is to raise by 1 the number j of the lay at the address $m-1$ in list J . The move is then continued with $j+1$, as described.

— the numbers k in that list indicating which piece is on which fields. The result is the first chessboard shown in *fig. 6*.

The machine then proceeds as if move 13 had been a failure. It withdraws the last lay, that is the piece

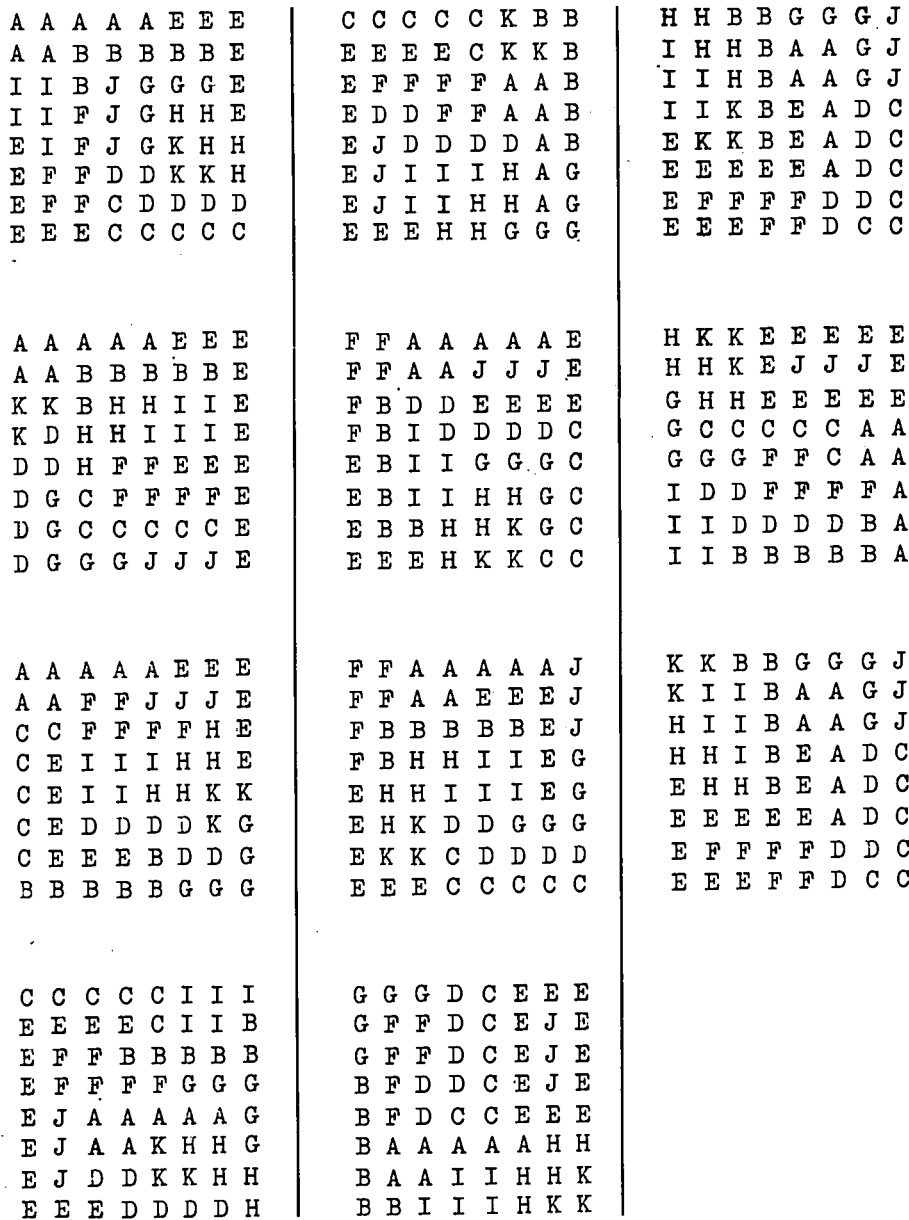


Fig. 6. The eleven solutions of the puzzle typed out by the PASCAL. Instead of the numbers $k = 1$ to 11 we have had the machine print the letters A to K to give a clearer picture. As the typewriter goes from left to right and from top to bottom, the solutions appear here turned 180° with respect to the description given in the text (cf. fig. 3).

Once all twelve successive moves have been made successfully—that is, when a solution is found—the machine is made aware of the fact by the appearance of the number $m = 13$ in the move counter. This is the signal for the machine to type out in chessboard array the complete contents then present in list V

laid in move 12, and tries to replace it by another lay, which of course fails. It therefore cancels move 12 and tries move 11 in another way, and so on.

In this situation the computer's actions seem to be rather "unintelligent", for it is perfectly plain that the withdrawal of one piece or even of two pieces will

not lead to another solution. But the machine can do no more than rigidly keep to the system, which guarantees completeness. We must just accept the few microseconds wasted on moves which to us seem pointless.

After retracing its steps far enough, the machine will have success with a new series of moves. In this way, as mentioned at the beginning, it successively finds *eleven* solutions. These are all shown in fig. 6.

Even after it has typed out the last possible solution, the computer goes on withdrawing moves and trying new ones. It is evident that it has then to go back farther and farther, and upon every move withdrawn it must reduce the number m in the move counter by 1. When all moves have been withdrawn, the number $m = 0$ appears in the counter, indicat-

ing to the machine that all possibilities have been tried and that it can stop.

The PASCAL takes eight minutes to complete the whole solving process, including the return to $m = 0$.

Summary. Electronic computers can be used for solving many problems that — at least in the first instance — have nothing to do with computing in the strictly arithmetical sense. As an example of such a problem a chessboard puzzle is described which served as an object of study for programming the PASCAL. The trial procedure worked out for solving the puzzle involves systematization and record-keeping, for which comprehensive tallies (“material lists”, “operating lists”, etc.) are stored in the computer’s memory. It takes the PASCAL eight minutes to complete the whole solving process, in which it tries some 3 million possibilities and finds eleven essentially distinct solutions.

LISTENING IN TO THE PASCAL

by W. NIJENHUIS *).

681.14

The work of an electronic computer in solving arithmetical problems consists in many cases of running through a succession of repeated cycles, and of cycles within cycles. This will perhaps have become clear from the foregoing articles in this number.

The idea has occurred to many builders of electronic computers to *listen in* to these cycles by making audible through a loudspeaker the passing of numbers through the registers. To this end the loudspeaker is connected to one of the flip-flop circuits in a register of the arithmetical unit. The listener will hear the voltage variations in that flip-flop as the numbers pass. Sometimes the pattern of the voltage variations is so arbitrary and changes so fast that only a hissing noise is heard. Often, however, the cycles mentioned produce a recognizable regularity in the sound, which may even result in a musical tone.

In this way every programme, or part of a programme, produces a characteristic sound by which it can be recognized. The sounds thus offer a means of checking the operation of a computer: the programmer who has become familiar with these characteristic sounds while testing his programme, can later often tell by ear whether the computation is proceeding normally. Use is made of the same facility with the PASCAL, for which purpose the loudspeaker is connected to the last digit (the least significant one) of the S register ¹⁾.

By way of illustration we have brought together on the attached gramophone record some fragments of the sounds which the PASCAL produces in performing the calculations discussed in four of the articles in this issue. The reader will find these fragments in four tracks on side 1 of the record, separated by visible margins; acoustically they are identified by an introductory morse signal. The first fragment consists of three sections, each announced by a morse sign. The fragments are taken from the following calculations:

	} Fourier analysis: — . Smoothing: — . . . Calibration: — . . .
Clover-leaf cyclotron ²⁾	
Corrugated cardboard trim-losses ³⁾ :	
Potential fields and electron trajectories ⁴⁾ :	— — — —
Chessboard puzzle ⁵⁾ :	— — — — —

*) Philips Research Laboratory, Eindhoven.

In the following we shall examine each fragment in turn and comment on the sound pertaining to the various computations. Side 2 of the record contains other sounds produced by the PASCAL, making it possible to go deeper into their relation with the computing operations.

Side 1, first sound fragment (—)

The sound of the first section of the first fragment, which relates to the Fourier analysis of measurement data of the clover-leaf cyclotron ²⁾, can be roughly represented phonetically as:

groom tik t toe-doe-de-doe-da-de-dee-doo-da

followed by a few times "tik-tik", and the whole thing is repeated almost identically a number of times.

During the first "groom" 60 numbers are fed in from the punched tape; these give the results of measurements of the azimuthal variation of the magnetic field at a given radius in the air gap of the cyclotron, for a particular choice of parameters (correction currents, etc.). In the subsequent interval one or two lines are printed, concerning the choice of parameters. During the short "tik" noise the measurements are reduced with the aid of a 7th degree curve (see — . . .) to magnetic induction values, and the deviations from the average field are computed. The composite sound now following is produced by the Fourier analysis proper. It comprises the following operations: the Fourier coefficients of the order 0 (constant term), 3, 6 etc., are successively computed; the sum of the Fourier series obtained up to a certain order, calculated at each of the 60 field points, is subtracted from the field value measured at each point; an autocorrelation is computed for the 60 differences; if this turns out to be too large, the procedure is repeated with the Fourier series to the next higher order. As a rule, working to 6 orders (i.e. using the constant term and 12 Fourier terms)

¹⁾ See Philips tech. Rev. 23, 1-18, 1961/62 (No. 1), especially p. 4.
²⁾ N. F. Verster and H. L. Hagedoorn, Philips tech. Rev. 24, 106-120, 1962/63 (No. 4/5).
³⁾ H. W. van den Meerendonk and J. H. Schouten, Philips tech. Rev. 24, 121-129, 1962/63 (No. 4/5).
⁴⁾ C. Weber, Philips tech. Rev. 24, 130-143, 1962/63 (No. 4/5).
⁵⁾ A. J. Dekkers and A. J. W. Duijvestijn, Philips tech. Rev. 24, 157-163, 1962/63 (No. 4/5).

the autocorrelation of the residue is small enough; the answers are then prepared for printing and are printed in a few lines (the series of sounds "tik tik tik . . .").

In the following sequence of sounds beginning with "groom" the same computation is carried out for the series of measurements relating to the next radius value, and so on.

— . . .

The "smoothing" process produces a sound roughly like:

te-te-te-te-te-feet,

and this is repeated a number of times. The process consists each time in deriving a "smooth" table \bar{y}_i from the originally computed table y_i of one of the Fourier coefficients, for 32 circles i of increasing diameter, such that there is a "smooth" transition between the values of that Fourier coefficient for successive circles. The machine computes \bar{y}_i so that

$$\sum_i [\delta^4(\bar{y}_i)]^2 + K(\bar{y}_i - y_i)^2$$

is minimized with a certain selected value of the factor K ; $\delta^4(\bar{y}_i)$ represents here the fourth difference of \bar{y}_i . First of all K is assigned a large value; \bar{y}_i is then still close to y_i , and the autocorrelation of the differences $\bar{y} - y$ is small. A smaller K makes \bar{y} "smoother" as the 4th differences now have more influence. Each "te" in the sound corresponds to a step in K . The process is stopped when the autocorrelation exceeds a predetermined value, because this indicates that "information" as well as "noise" is also being smoothed away. The process is then repeated for the next Fourier coefficient, and so on.

— . . .

The third section of the first fragment is very short, lasting less than one second. This part is the sound of the complete computation — by the least squares method — of the 8 coefficients of the 7th-degree polynomial, the graph of which is the line of best fit drawn through 22 calibrated points.

The three sections of this sound fragment strikingly illustrate the computing speed of the PASCAL.

Side 1, second sound fragment (— —)

The second fragment, which relates to the problem of cutting-losses, is in marked contrast to the first. Since the process of "linear programming" used for solving this problem was discussed only in broad lines in the relevant article³), we shall not give an explanation of the sounds produced. It can only be noted that each cycle in the sound corresponds to one complete iteration cycle.

Side 1, third sound fragment (— — —)

As mentioned in the relevant article⁴) the calculations of potential fields and electron trajectories were carried out with an IBM 650. Subsequently, however, the calculations have been programmed for the PASCAL, and the sound reproduced in the third fragment relates to a potential field calculation using this machine, for given boundary conditions (given electrode configuration and voltages on the electrodes). The sound somewhat resembles the clucking of a hen; phonetically it can be represented as a frequent repetition of the group (u being pronounced as in the French "la lune")

kree-lu-dlu-dlu-dlu- . . .

Each of these groups corresponds to an iterative cycle over the entire potential field. For some boundary conditions as many as a hundred such cycles are necessary for convergence.

During each "lu", the machine computes the potential in a row of network points on a line parallel to the axis of the electrode system. The pitch of "lu" is determined by the nature of the computing cycle per network point, and its duration by the length of the relevant row of points. The successive "lu" sounds correspond to the successive rows of points at increasing distance from the axis; the individual notes differ in duration since the length of the outermost rows of points, depending on the electrode configuration, differs from that of the rows close to the axis. The first "kree" corresponds to the calculation for all network points on the axis itself; since a different formula underlies this calculation⁵) the pitch of "kree" differs from that of "lu".

The values of all network-point potentials (e.g. 3000 points) are stored in the drum memory of the PASCAL. When a network-point calculation is in progress, use must be made of the (provisional) potential values at points of three successive rows. These "working data" are transferred for this purpose to the ferrite core store. When the computation is completed for all points on a row, the working data of the row nearest to the axis are returned to the drum and the data for a further row of points, more distant from the axis, are extracted from the drum. This transport process causes the "d" sounds in the "lu-dlu-dlu- . . ." series. At the end of an iteration cycle over the entire network, all the working data last used are stored and the data of the row of points on the axis and of the next two rows of points are extracted from the drum. This longer transport accounts for the interval between each group of "kree-lu-dlu- . . ." sounds.

Side 1, fourth sound fragment (— — — —)

The fourth fragment, relating to the solving of a chessboard puzzle⁵), is again quite different in character. First, we hear the sound of the programme being fed in (on punched tape). After about 12 seconds the sound of the actual computing begins. This has entirely the character of noise and continues until the first solution of the puzzle has been found. The PASCAL then types the solution on the typewriter, during which time nothing is heard as nothing is being done in the S register. This period of silence, which in reality lasts about 16 seconds, has been shortened on the record to about 5 seconds. After this the machine is again heard searching for the next solution, which takes about 1 second. When this has been typed out (interval again shortened to 5 seconds) the work of computing the third solutions begins, which takes more than 20 seconds, after which the fragment is broken off. The puzzle has eleven different solutions, which takes the machine altogether about 8 minutes to find and type out. During this time the machine has tried out several million possibilities.

We now turn to side 2 of the gramophone record. The first fragment on this side relates to the search for prime numbers, and for this case we shall give a more detailed explanation of the sounds of computation. The second fragment is musical in character, and the last fragment is a logical complement to it.

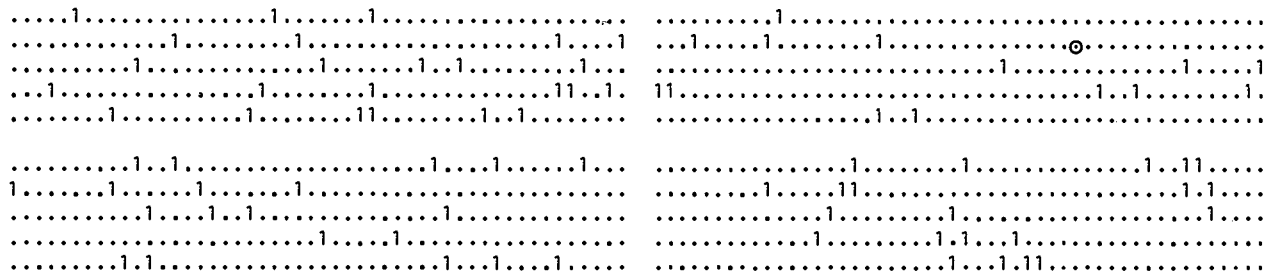
Side 2, first sound fragment

The programme for finding prime numbers served for a time as an example to demonstrate the PASCAL in operation. The odd numbers G, starting from a predetermined arbitrary number, are tried one by

one in the machine to see whether they can be divided by the odd numbers $p \geq 3$. By using slightly different programmes, the results can be presented by the machine in different ways. The most straightforward method is to let the machine print out each prime number found but not the other numbers. This was the procedure adopted for the present recording. Another, very useful method is that of the "prime number pattern": for each number G tried the machine prints a point if G is found to be divisible, and another character, e.g. the figure 1, if G is found to be indivisible, i.e. is a prime number. In this way, beginning for example with the number 34 359 738 000, the pattern in fig. 1 was produced for the thousand succeeding odd numbers; the pattern shows very clearly the arbitrary distribution of the prime numbers. A third method of presenting results may be mentioned, as it will prove useful in the following explanations: in this method the machine resolves all consecutive numbers (which may include the even numbers) completely into their factors and prints out the results. A table is then obtained as shown in fig. 2 for 62 numbers from the pattern in fig. 1, beginning at the place marked with a circle (the number $2^{35} + 1$). This is precisely the series of numbers covered by the machine during the recording of the present fragment. The prime numbers in this series can immediately be recognized (marked by a star in fig. 2).

Let us now consider the computing sound. In the search for large prime numbers the sound begins as a hissing noise, from which a siren-like wailing develops, ending in a short interval of silence while the prime number found is being printed. Sometimes the wail ends in a hissing noise again, indicating that, after some computation, the number tried has proved

34359738000



34359740000

Fig. 1. Prime-number pattern for the 1000 odd numbers between 34 359 738 000 and 34 359 740 000. The number $2^{35} + 1$, with which the machine begins computing in the recorded fragment, is marked with a circle. — Note the "pairs" of prime numbers occurring here and there in the pattern. Experience has shown (although there is as yet no proof) that such pairs continue to occur no matter how far the search is pursued.

	34359738369=3.11.43.281.86171
	34359738371=7.4908534053
	34359738373=59.582368447
	34359738375=3.5.5.5.811.112979
	34359738377=17.19.106376899
	34359738379=97.103.149.23081
	34359738381=3.3.3.3.1427.297263
	34359738383=163.883.238727
	34359738385=5.7.29.33851959
	34359738387=3.13.23.38305171
a	34359738389=20249.1696861
	34359738391=11.3123612581
	34359738393=3.347.1699.19427
	34359738395=5.6871947679
	34359738397=673.51054589
	34359738399=3.3.7.7.77913239
b	34359738401=37511.915991
	34359738403=53.3319.195329
	34359738405=3.5.2290649227
c	34359738407=132949.258443
	34359738409=157.3109.70393
	34359738411=3.17.2221.303341
	34359738413=7.11.13.34325413
	34359738415=5.19.361681457
	34359738417=3.3.3817748713
	34359738419=467.73575457
d *	34359738421=34359738421
	34359738423=3.37.309547193
	34359738425=5.5.1374389537
	34359738427=7.42461.115601
	34359738429=3.31.181.277.7369
	34359738431=419.1583.51803
	34359738433=23.1493901671
	34359738435=3.3.3.5.11.23137871
e	34359738437=29077.1181681
	34359738439=13.67.39448609
	34359738441=3.7.41.39906781
	34359738443=29.1184818567
	34359738445=5.17.3343.120919
	34359738447=3.10529.1087781
	34359738449=1171.29342219
f *	34359738451=34359738451
	34359738453=3.3.19.89.2257687
	34359738455=5.7.43.47.485753
	34359738457=11.107.139.210019
	34359738459=3.11453246153
	34359738461=61.563274401
	34359738463=1571.21871253
	34359738465=3.5.13.173.937.1087
g *	34359738467=34359738467
	34359738469=7.193.25432819
	34359738471=3.3.3817748719
h *	34359738473=34359738473
	34359738475=5.5.15107.90977
	34359738477=3.73.1543.101681
	34359738479=11.17.23.167.47837
	34359738481=617.5003.11131
	34359738483=3.7.439.3727057
	34359738485=5.27541.249517
	34359738487=151.1531.148627
	34359738489=3.3.3.241.5280427
	34359738491=13.19.19.31.59.4003

*

Fig. 2. Table of the printed-out prime factors found for the 62 odd numbers from $2^{25} + 1$ to $2^{25} + 123$.

to be divisible and that the machine has started to examine a following number (and possible further following numbers). This happens in our fragment before the first prime number (*d*) has been found: the machine has then already examined the "difficult" numbers *a*, *b* and *c* (see fig. 2), its efforts each time only being successful with large divisors *p*; the work on the number *c* is particularly audible. Also heard in the fragment is the computation of the prime numbers *f*, *g*, *h*, between which the only rather difficult case is the number *e*.

To explain the strange wail, rising and falling in pitch, we must consider what happens during the programme in the *S* register of the PASCAL, especially as far as it concerns the last digit, to which the loudspeaker is connected. The *S* register is used during the division operations: at the beginning of each division the machine sets the dividend in *S*, and at the end of the division *S* contains the quotient. The programme comprises the steps shown in fig. 3⁶⁾.

The PASCAL takes almost exactly 180 μ s to complete the thickly drawn cycle. Plainly, then, the investigation of most of the numbers *G* (which prove to be divisible by 3, 5, 7 or some other relatively

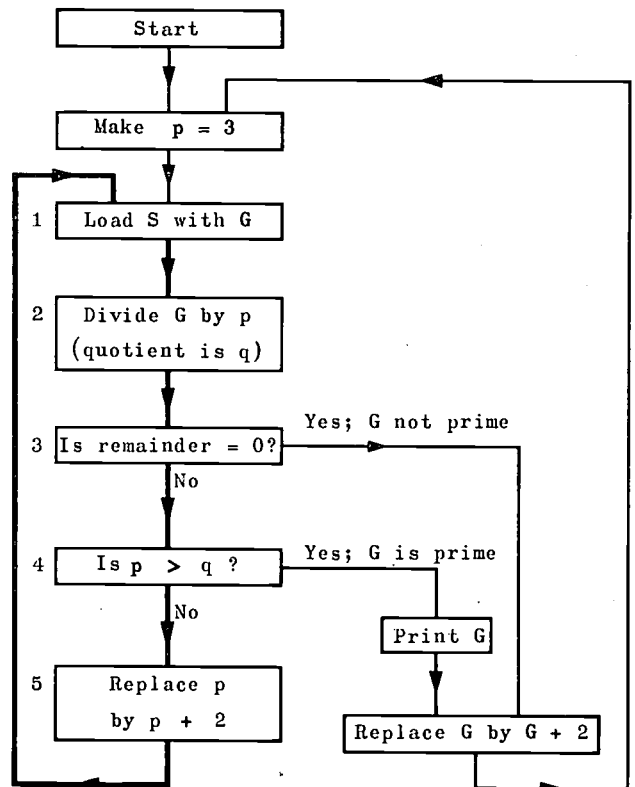


Fig. 3. Computing programme used in the search for prime numbers.

⁶⁾ This programme is certainly not the most economical for the purpose, since *G* is divided by *all* odd numbers *p*, whereas of course the result of a division by, say, 9 or 15 is already established if it is found that the number does not divide by 3 and 5.

small number) takes only a few milliseconds. Only if G is indivisible or contains only very large prime factors is it necessary to continue the divisions up to large divisors p (up to $p \approx \sqrt{G}$ if G is indivisible); the thickly drawn cycle is then repeated numerous times. This is the stage in which the siren wail is heard. All other numbers together contribute only to a hissing sound at the beginning of each such stage.

Let us now consider what happens with the last digit of S during the above-mentioned cycle in the steps 1-5 of the diagram.

- 1) Since G is always odd, the last digit of S becomes a 1. This remains for about $10 \mu s$.
- 2) The quotient q is built up in S . During this process numerous ones and noughts pass the last location of S in a fairly irregular pattern. This takes about $75 \mu s$.
- 3) 4) 5) The quotient q in S remains unchanged, and so too therefore does the content of the last digit of S . This lasts $95 \mu s$.

The voltage variations in the loudspeaker during a single thick cycle thus appear as shown in *fig. 4a* or *b*, depending on whether the quotient q is even or odd. Everything depends now on the alternation between even and odd q .

On a first glance one would expect this alternation to be entirely irregular. The only recognizable periodicity in the loudspeaker voltage is then the fundamental period of the cycle, of duration $180 \mu s$, and all that is strictly repeated in this is the presence of the voltage 1 during step (1), which lasts $10 \mu s$. The next following group of pulses in step (2) is chaotic, and so is the alternation in the succeeding state, lasting during steps (3), (4), (5). All that can therefore be heard is a note of $10^6/180 \approx 5500$ c/s. rich in overtones and not particularly striking.

If the machine is working on a large number G , however, a certain regularity gradually enters into the alternation of odd and even quotients q . To see how this comes about, we consider the graph of the

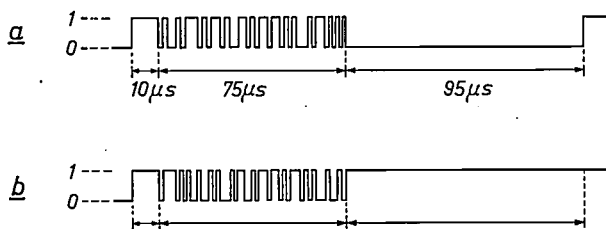


Fig. 4. Voltage variations on the loudspeaker during the thickly drawn cycle in *fig. 3*.
a) Quotient q even, b) q odd.

relation $G = pq + \text{remainder}$: the graph is a stepped curve (*fig. 5*) which follows the hyperbola $PQ = G$ (P and Q are continuous variables instead of the discrete variables p and q) and in which each step arises from an increase of $+2$ in p . We first consider

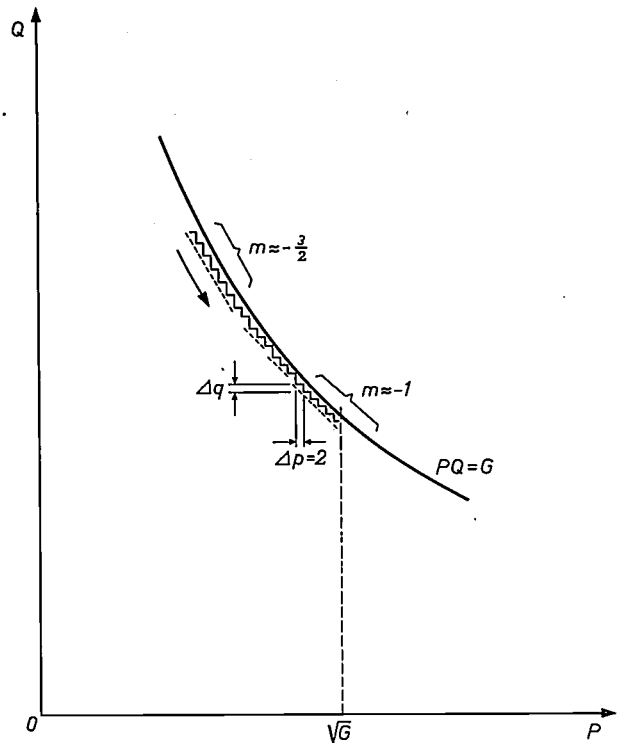


Fig. 5. Hyperbola $PQ = G$ and the stepped curve $G = pq + \text{remainder}$.

the portion of the graph where the hyperbola has the slope $m = -\frac{3}{2}$. If G is large, then each step with $\Delta p = 2$ in this portion has a height $\Delta q = -3$. Here, then, q alternates in every step between odd and even, and these alternations give the loudspeaker voltage a period of $2 \times 180 \mu s$ (the shortest possible in these alternations) and must therefore be audible as a tone of about 2800 c/s. The same apparently applies to the portions of the curve where the slope is $m = -\frac{5}{2}, m = -\frac{7}{2}$ etc., and each of the "summits" in the siren sound means that the machine is then working on such a portion of the hyperbola.

The fact that the sound in between falls and rises again in pitch can be explained in an analogous way. In the portions where the slope m of the hyperbola is equal to a negative integer the quotient for each increase $\Delta p = 2$ takes an even step $\Delta q = 2|m|$, and thus does not alternate between odd and even. A rough estimate shows that the uninterrupted number of repetitions of this even step Δq is of the order of magnitude

$$t_m = k \sqrt[4]{G/|m|^3}, \dots \dots \dots (1)$$

where k is a factor roughly of the magnitude of $\frac{1}{4}$ or $\frac{1}{5}$. After this series of even steps there will be *one* odd step Δq — that is to say a voltage discontinuity in the loudspeaker — and then again a whole series of even steps, the *series* being roughly t_m in length; this will be repeated over and over again (provided G is very large). The voltage discontinuities near the portion of the hyperbola considered, where m is an integer, thus cause the fairly long period $2t_m \times 180 \mu s$, and the sound here has a much lower pitch than before. Gradually the series with even step height Δq become shorter the closer we approach the next portion of the hyperbola where m is half of an odd integer, and longer again the nearer the next portion approaches where m is an integer. Accordingly, the period of the voltage discontinuities grows shorter, and longer again. It is this that produces the “wailing” effect.

For $m = \dots -3, -2, -1$ and $G \approx 2^{35}$ (our recording was made in the region of this G) we find from (1) the respective periods $40 \times 360, 50 \times 360$ and $90 \times 360 \mu s$, i.e. roughly the frequencies 70, 55 and 30 c/s. These should be successively the pitches of the sound in the last three “troughs” before the prime number is found, because from $m = -1$ onwards we have $p > q$ ($p > \sqrt{G}$) and the machine need compute no further (see step (4) in fig. 3) This is difficult to test quantitatively as it is not easy to measure the pitch of the computing sounds reliably, presumably owing to the overtones, whose amount is large and continuously changing and also owing to the continuously sliding pitch. Qualitatively, however, the phenomena are sufficiently explained.

Side 2, second sound fragment

From the foregoing it will be clear that the computer can be made to produce a melody by giving it a suitably designed programme for “computation”. This is in fact a favourite way of letting visitors know that the machine is entering into the spirit of official opening ceremonies at computing establishments and on similar occasions. The programme by which the PASCAL “sings” is illustrated in fig. 6. Every time the machine receives the instruction “Invert S ”, all ones in the (arbitrary) contents of the S register are replaced by noughts, and vice versa. The result is a voltage alternation on the digit to which the loudspeaker is connected. This is periodically repeated at intervals of t microseconds, depending on the “wait” instruction. In the PASCAL, is a combination of a “wait” and a repeat instruction — see p. 14 of the article ¹⁾ quoted above). Thus, a tone is produced having a frequency of $10^6/2t$. The cycle is repeated n times, and therefore the tone lasts

$n \times t$ microseconds. A list of the successive waiting times t to be observed and the number of repetitions n represents the melody to be “sung” and is stored in the machine’s memory.

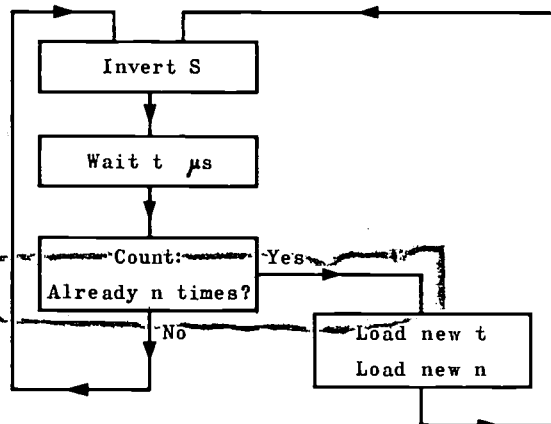
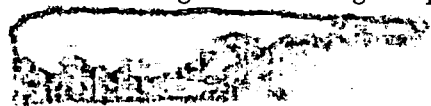


Fig. 6. Programme for “singing” a melody. The melody is stored in the computer’s memory in the form of a list of values of t (pitch) and n (duration of tone).

In this way we have programmed a minuet of Mozart, which is heard as the second fragment on side 2 of the record. The “phrasing” of the melody is produced by introducing suitable pauses. During each tone the voltage variations on the loudspeaker are purely periodic, but because they have a square wave form the tones from the PASCAL have a nasal timbre.

Side 2, last sound fragment

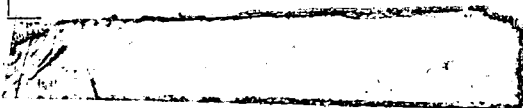
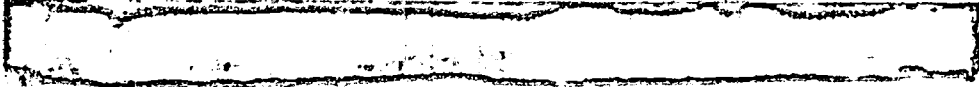
The programmes supplied to the PASCAL always have the form of a series of characters punched into a tape, having the meaning of numbers, which according to a certain code represent operations, or of addresses in the memory, or arithmetical numbers proper. Further to what we have just said about the “singing” programme, when supplying melody information to the machine, we can now take the curious step of not using the list of values of t and n that correspond to the notes of the Mozart minuet (or any other piece of music) but of using an arbitrary piece of programme tape, one character on which being always interpreted by the machine as t and the next as n , and so on right along the tape. We may then expect the machine to produce musical tones that together form a kind of “stochastic music” (tonal combinations that are purely random and thus not predictable by any laws of music), a subject which has been a talking point among composers and music theoreticians in recent years. We achieved this by using as “music information” our test programme for the ferrite core store. The result is the last fragment on our gramophone



record. Although the music is found to be not entirely "stochastic" — there was apparently already too much regularity in the test programme — the effect is perhaps strange enough to prompt reflection on the nature of what we call a "melody".

Summary. By connecting a loudspeaker to one of the flip-flops in the arithmetical unit of an electronic computer, the passage of numbers through the register concerned can be made audible.

In many computing programmes the machine repeats one cycle over and over again; this will often produce in the loudspeaker distinguishable sounds, which in some cases can even be used to check the operation of the computer. A gramophone record attached to the article presents on side 1 some fragments of the sounds produced by the PASCAL in performing computations discussed in four of the articles published in the same issue. On side 2 the computer is heard searching for very large prime numbers, and for this case a closer analysis of the sound is given in the article. The record finally demonstrates the computer "singing" a minuet of Mozart and interpreting an arbitrary programme as a melody; explanatory comments to these sounds are also provided in the article.



Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES

SOME CONSIDERATIONS ON THE NUMERICAL CONTROL OF MACHINE TOOLS

by T. J. VIERSMA *).

621.9-523.8

The article below reproduces more or less verbatim the text of the public lecture given by Dr. T. J. Viersma on 4th October 1962 upon his inauguration as reader at the Delft Technische Hogeschool. With the kind assistance of the author the lecture is supplemented here with a few figures and bibliographical references. The subject dealt with forms a useful introduction to a series of three articles which will shortly be published in this journal, describing a numerically controlled milling machine developed in the Philips Research Laboratories, Eindhoven.

Introduction

Since 1950 the automation of machine tools (milling machines, lathes, etc.) has given rise to a technique in which the information for operating and controlling the machine is supplied in numerical form. This information, frequently fed in on a punched tape, comprises the whole programme of movements and machining operations, cutting speeds and number of revolutions of the machine. Each punched tape thus produces a particular workpiece.

The fundamental problems involved in this new technique can be divided into two categories. The first concerns *data processing*, by which the data of the workpiece (e.g. taken from the drawings) combined with data of the machining method and of the machine itself, are changed into commands to the machine. In the more complicated cases electronic computers are used for this purpose. The second problem concerns *measuring and control technique*, which in this case means the ways of carrying out the commands.

To give some idea of the demands imposed on these techniques we shall first examine some points of general interest in connection with numerical control.

We are concerned here with machine tools in which a workpiece is given a required shape by a cutting tool in machining operations such as milling,

drilling or turning. In these operations the tool moves relative to the workpiece. The motions can be rectilinear, with the workpiece or tool mounted on a linearly moving slide, or they can be rotational, i.e. either the workpiece turns (turning lathe, turning table), or the tool turns (drilling machine, milling machine). Usually there are one or more rectilinear movements combined with rotary movements. In order to control this fairly complicated aggregate of linear and rotational motions, numerous operations have to be carried out which require an enormous amount of information. In numerical control this information is, in principle, fed to the machine without human agency and transformed into the required operations.

An essential part of the information relates to the *dimensions* of the workpiece. In the classical form of machine control the workpiece is given the appropriate dimensions by an operator who consults the drawings and translates the data into corresponding displacements of the slides. In numerical control the cutting data are supplied directly to the machine without the intermediary of drawing and operator. These data and the manner in which they are processed are, of course, directly bound up with the dimensional *accuracy* of the workpiece. The accuracy obtainable is therefore one of the most important aspects of the numerical control of machine tools. In some cases it is a factor of 10 better than in the classical methods.

Another very important part of the data relates

*) Philips Research Laboratories, Eindhoven.

**) The lecture has been published by Uitgeverij Waltman, Delft.

to the *speed* at which the rectilinear and rotary movements are carried out — that is, feed rates and cutting speeds. The latter, to a considerable extent, govern the *productivity* of the machine. In the classical techniques the operator acquires this information during his training, so that the knowledge of metalworking methods and machine tools is, as it were, realized in the operator. In numerical

tool has been evolved, using dozens of cutting tools which are successively brought into position and into operation, sometimes several at a time (*fig. 1*). Numerical control in this case clearly influenced the conception of the machine. On the other hand there are many numerically controlled machine tools where tool selection and changing is left entirely to the operator.

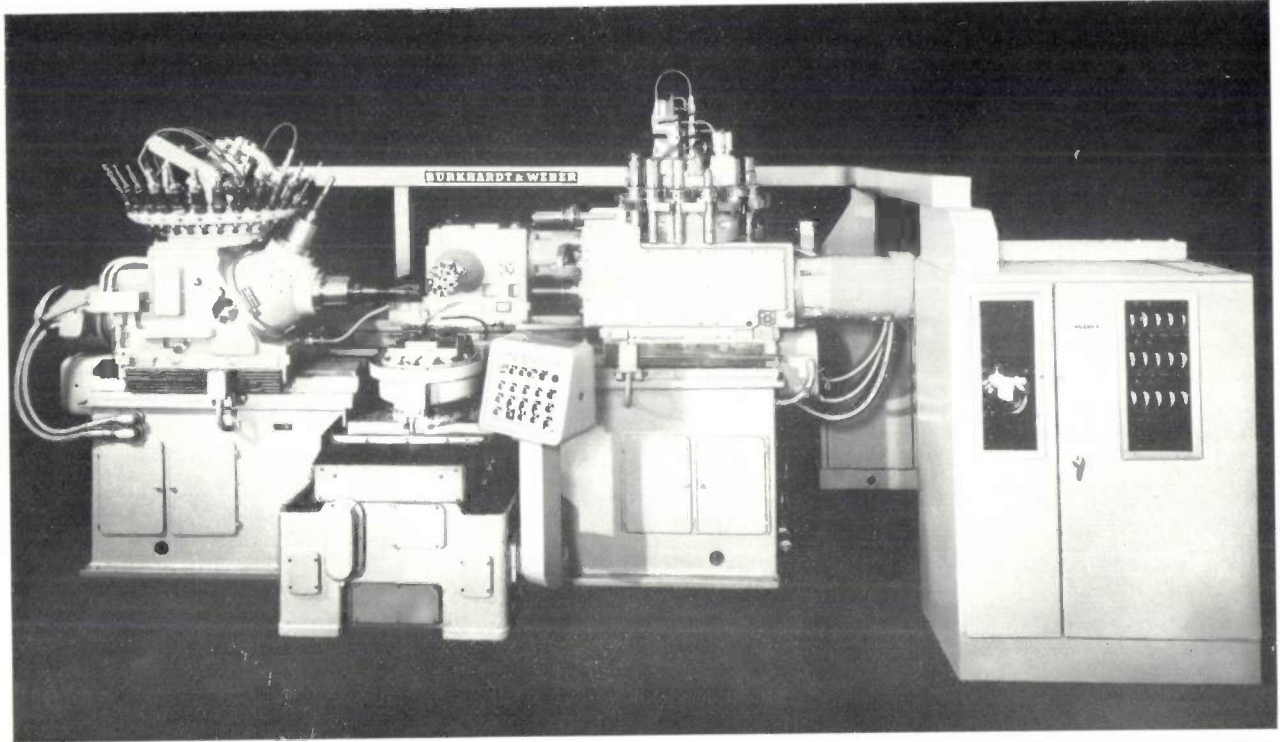


Photo Burkhardt & Weber, Reutlingen/Württs

Fig. 1. Numerically controlled machine tool which can automatically set different cutting tools into operation one after the other. The tools are contained in the magazine on the left of the machine. Bottom right of this magazine are two spindles mounted on a rotary support. Whilst an operation is being carried out with the aid of the horizontal spindle, the cutting tool for the next operation is already being placed in the other spindle, so that little time is lost in tool changing.

control all this information is normally supplied on a punched tape. This means that a large proportion of the operational skill is transferred to the work-planning stage, and thus becomes the responsibility, for example, of the programmer who draws up the programme by which the computer produces the punched tape. This involves many problems, and further advances in the technique of the numerical control of machine tools largely depend on the manner in which these problems are solved.

Other data may relate to the choice and changing of tools. The successive positioning of tools, as in the case of a turret lathe, can also be entirely numerically controlled. There are some spectacular instances where a completely new type of machine

Applications of numerical control

The applications of numerical control fall into two parts¹⁾:

- 1) *Contouring*, in which different controlled movements, which are very accurately coordinated, take place simultaneously during the machining process.
- 2) *Point-to-point positioning*, in which only one controlled movement takes place at a time.

Since contouring and point-to-point positioning differ fairly considerably in regard to data processing and to measurement and control techniques, we

¹⁾ See John D. Cooney, Automatic metalworking: a status report, *Control Engng.* 7, No. 9, 158-173, 1960.
C. A. Sparkes, Automatically controlled machine tools, *Chartered mech. Engr.* 9, 298-305, 1962 (No. 6).

shall deal with these two categories for the most part separately. We shall first consider some general aspects of contouring; some characteristic examples, where it is difficult to achieve the required accuracy with the classical methods and where the productivity of the latter is lower, are aircraft wings, turbine and ship-propellor blades, cams, and dies.

The prototype for such workpieces used to be made by hand. First of all the contour was marked out with a large number of points, for example on a jig borer. Next, a smooth curve was produced through these points by filing, grinding, surfacing, etc.: this being actually a method of interpolation, based on the operator's observation. In spite of ingenious techniques, which can sometimes considerably simplify the operations, nearly all contouring processes depended on these tedious, time-consuming and relatively inaccurate manual methods. The advantages of numerical control are obvious in this connection: mutually coordinated movements are carried out quickly and accurately without interrupting the cutting process. Compared with the classical methods a ten-fold improvement can easily be achieved in both dimensional accuracy and productivity. Planning and data processing are the only obstacles preventing the complete break-through of numerically controlled contouring lathes, for the electronic and mechanical problems no longer present any fundamental difficulties.

To avoid misunderstanding it should be remarked that numerical contouring control is especially suited for single-run and small series production. Starting from an accurately dimensioned prototype, larger series can be produced much more efficiently by one or another copying technique.

Apart from contour machining, point to point positioning has come to the fore in recent years as a subject for numerical control on a variety of machine tools (contour cutting is done almost exclusively on milling machines). The remarkable thing is that this simple form of numerical control was the last to be developed. Examples of machines used for positioning are jig borers, milling machines, drilling machines, capstan and turret lathes.

A characteristic of such numerically controlled machines is that a fairly large number of different setting and machining operations take place in very rapid succession. The great advantage here is that the time formerly spent on setting, measuring, checking readings and thinking is reduced to a minimum, thus also practically eliminating the fatigue factor and the risk of errors. From the metal-working point of view, however, numerical control here offers in principle no new possibilities. Dimen-

sional accuracy, for example, is not essentially improved. The main justification for introducing numerical control here is the improved productivity. The advantages to be gained in this respect will be evident to anyone who has ever made a time and motion study in a workshop to determine what percentage of the working hours is really spent effectively machining. Depending on the kind of workshop, the figure varies from 5 to 25%. With numerical control a very much higher figure can easily be achieved.

Data processing and the technique of measurement and control present far fewer problems in the case of positioning than in contour machining. Positioning in fact involves making three choices:

- 1) Choice of the various possible movements.
- 2) Choice of the beginning and end of the movement.
- 3) Choice of the speed of the movement.

In principle, then, a complete machining pass can be defined by fixing only three numbers. There are not usually more than a few dozen passes, so that the total information remains relatively limited and simple.

Measurement and control aspects

Since important aspects of data processing in numerical control are partly governed by the elements of the measurement and control system, we shall discuss these first.

As we have seen, the most important measuring and control engineering problem concerns the realization of accurately programmed movements. For producing a movement we need a drive system, and an associated measuring system. Furthermore, a comparator is required to see how far the desired movement (input signal f_i) agrees with the movement executed (output signal f_o); see fig. 2. The requirements for the drive and the measuring system depend among other things on:

- a) the variation of the input f_i as a function of time t ;
- b) the permissible deviation $\varepsilon = f_i - f_o$, for example as a function of time;
- c) the possibility of adaptation to the rest of the numerical control system, in particular to the

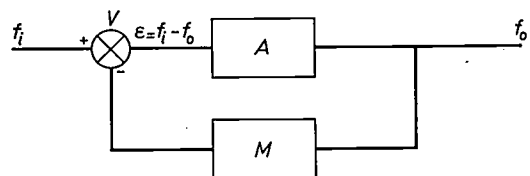


Fig. 2. Block diagram of a feedback control system with drive A , measuring system M and comparator V . f_i is the input signal indicating the required motion, f_o is the motion actually carried out, $\varepsilon = f_i - f_o$ is the deviation.

last phase of the data processing, which governs the form in which the command signal f_i becomes available.

The drive

As regards the drive we first consider the situation with a simple *positioning system*: a jig boring machine. After a hole is drilled and the drill withdrawn, the coordinates of the next hole must be set quickly in order to lose as little time as possible. The input signals f_i are all discontinuous functions of time (fig. 3). The usual requirement regarding the deviation $\varepsilon = f_i - f_o$ is that ε should drop as quickly as possible below a specified limiting value δ (\approx accuracy of measurement). The error ε may temporarily be large provided the value δ is reached quickly. The shortest possible travel time can be achieved if a hydraulic servomotor²⁾ is used in conjunction with a measuring system which supplies information continuously, and if a certain amount of overshoot is permitted. If overshoot is not permissible, the system must be set for a slower response, i.e. a longer travel time.

If a hydraulic drive is thought too expensive, and for example an eddy-current coupling is used instead (having a relatively very large time constant), the travel time again increases considerably. This simple example shows, then, that a compromise must be found between the price of the drive system and the quality obtainable, in this case the travel time.

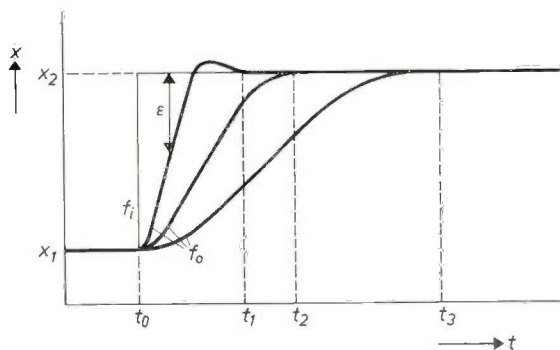


Fig. 3. Input signal f_i and output signal f_o of a positioning control system at various control settings. In the time between two successive operations the tool must be moved with respect to the workpiece from position x_1 to position x_2 . For this purpose an input signal of the form shown is supplied to the control system. The deviation ε must now drop rapidly below a specified limiting value; as a rule overshoot is not permissible. The travel time from t_0 to t_1 , t_2 or t_3 is governed, among other things, by the choice of drive and the setting of the control.

²⁾ T. J. Viersma, Investigations into the accuracy of hydraulic servomotors, thesis Delft, 1961. This thesis has also been published in Philips Res. Repts 16, 507-597, 1961 (No. 6) and 17, 20-78, 1962 (No. 1).

This applies with even greater emphasis to numerically controlled drilling machines, turning lathes and milling machines that are used for the positioning and feeding. Here, as opposed to the situation just discussed, the workpiece continues to be machined during the movement (fig. 4). In this case, then, not only the beginning and end point of the movement are important but also the speed

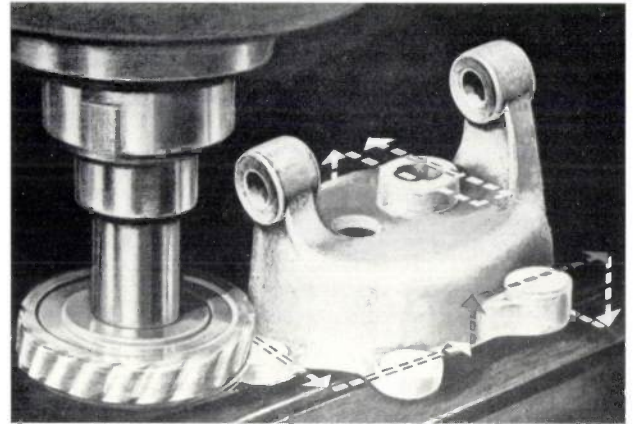


Fig. 4. Example of combined point-to-point positioning and feeding on a milling machine. The milling cutter describes consecutive linear movements as indicated by the arrows. During these movements the workpiece is machined. (Reproduced by courtesy of The Chartered Mechanical Engineer, London.)

of the movement (feed rate). The maximum permissible feed rate is usually governed by the machining technique: life of the tool, surface roughness, etc. It is therefore important that the movement should be carried out at the right feed rate within certain — not very close — limits. The deviation $\varepsilon = f_i - f_o$ may thus still be fairly considerable during the movement, provided that the beginning and end points are accurate (fig. 5). In most cases the specifications of the workpiece do not allow overshoot. These considerations mean that substantially higher demands are made on the drive system than in the previous example. Here too, the hydraulic servomotor is superior but also expensive.

In recent years some European firms have introduced a drive system with preselection of speeds and feed rates, which is closely linked with modern developments in machine tools. The movements for the feed are derived, via a gear box, from an electric motor. This is done by selecting the appropriate combinations of gears by means of electromagnetic couplings. The command signal for this purpose can be obtained in a simple way from the data carrier. A system of this kind has the advantage of being in line with present-day techniques, and can therefore easily prove successful if it is developed

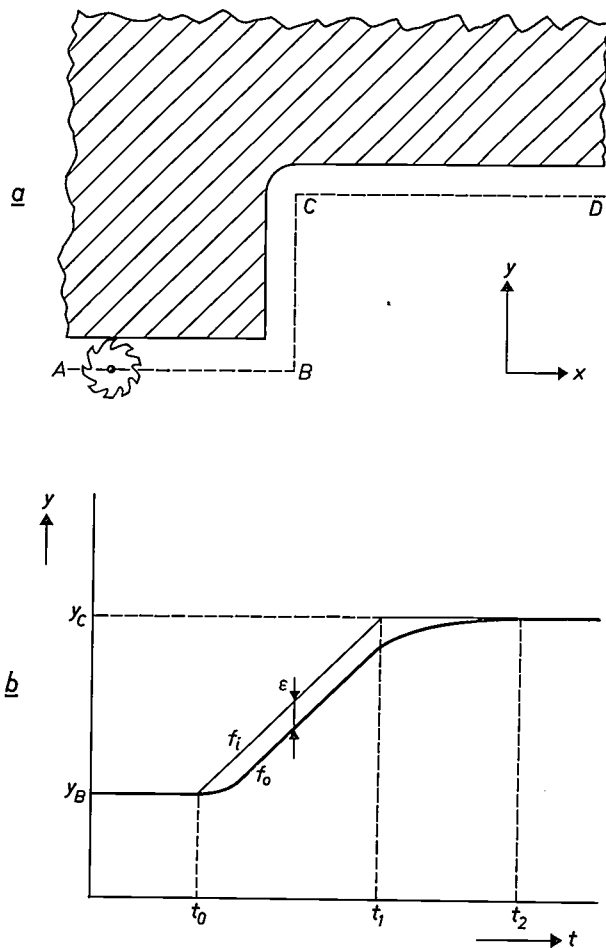


Fig. 5. Combined point-to-point positioning and feeding.
 a) When milling a workpiece of the indicated shape, the milling cutter has to follow the path $ABCD$ with a certain feed rate. The cutter thus describes linear movements in the x - y plane whilst cutting during the movement.
 b) Input signal f_i and output signal f_o of the control system during the displacement of the cutter in the y direction from B to C . The time is plotted on the abscissa. The permissible speed of the movement, which governs the slope of f_i between t_0 and t_1 , is given by the feed rate. During the whole movement, ϵ must remain below a specified limiting value. Overshoot is in general not permissible, so that the last part of the movement (from t_1 to t_2) must be done with a slow feed rate.

in close cooperation with machine-tool manufacturers. From the point of view of control engineering, however, this system is far from ideal. The end point of a given pass is governed by the overshoot of the machine after switching off, by means of an electromagnetic coupling. This overshoot depends among other things on the switching time of the coupling, on the speed at the moment of switching off, and on the load (cutting forces, mass and friction). Without radical measures it is virtually impossible to keep the variations in the overshoot within reasonable limits. Among the most obvious measures are to slow down the feed rate in the last part of the pass, or to introduce a brake, at the cost, however, of a great deal of the

relative simplicity. Another drawback of this method is the fact that power circuits can be a source of interference, which may endanger the reliability of the data processing, particularly with a digital system.

In positioning and feeding operations the hydraulic motor, controlled by an electro-hydraulic valve, is steadily coming to the fore. These motors are continuously variable, and their sensitivity to load variations is extremely low³⁾, mainly as a result of the special valve designs (fig. 6). It looks as if the hydraulic drive will continue to gain ground especially if it can approximate better than at present to the optimum solution from the point of view of control technique. The main disadvantage of the hydraulic drive is that it calls for a machine design which — for many machine tool manufacturers — is entirely new (no gear-box), and to which there is opposition on both practical and psychological grounds.

Finally, if we consider the demands made on the drive system for numerically controlled contouring machines, we see here too that the hydraulic motor with electro-hydraulic servovalves is far superior. What is more, the particularly favourable proper-

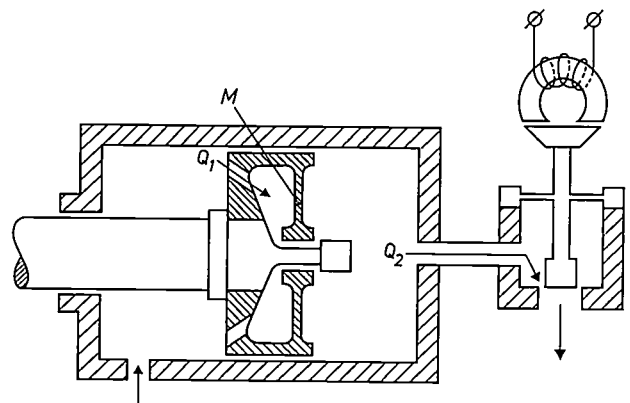


Fig. 6. Example of a hydraulic servomotor with load compensation and an electromagnetic servovalve. In the left-hand chamber of the cylinder, oil is pumped under constant pressure and flows through a port in the piston to the right-hand chamber of the cylinder, which it leaves via the servovalve. The movement of the piston is governed by the difference between the oil quantities Q_1 and Q_2 flowing into and out of the right-hand chamber. Q_1 is constant to a first approximation and Q_2 is proportional to the size of the outlet opening in the servovalve. The speed of the piston therefore changes in proportion to the displacement of the valve. Hydraulic servomotors are by nature relatively insensitive to external loads. The motor represented here is provided with extra compensation. If, for example, a force from the left is exerted on the piston, the pressure in the right-hand chamber rises, causing Q_2 to increase. Since the pressure increase deflects the membrane M , the opening in the piston widens so that Q_1 also rises. If the membrane is properly dimensioned, Q_1 and Q_2 increase by equal amounts, so that the load causes no additional movement of the piston.

³⁾ See chapter IV of the thesis referred to under ²⁾.

ties of the hydraulic drive meet a need in this case in a way that practically rules out other drive systems. In contouring machines stringent demands are made at every instant on the deviation $\varepsilon = f_i - f_o$, irrespective of speed, acceleration and load. The integration time-constant of hydraulic servomotors is of the order of 3 to 10 milliseconds, whereas that of electric servomotors is a multiple of this figure. Consequently a hydraulic servomotor can respond very much faster to changing input signals than an electric servomotor. Moreover the hydraulic servomotor is much less sensitive to load variations, and is capable of supplying a substantially higher power. It is therefore not surprising that the hydraulic servomotor is rapidly ousting its rivals in this field and indeed has already gained the upper hand.

Measuring systems

The situation of measuring systems in relation to numerical control is so chaotic that it is hardly possible here to deal with it completely. We shall therefore confine ourselves to a few main points⁴⁾. The difficulty of measurements in numerical control can be seen from the demands made on an ideal measuring system, which are to measure displacements in the order of metres with an accuracy of microns and to present the results immediately in electrical form. In this respect enormous advances have been made: in 1950 no system could even approximate to these requirements, whereas now, in 1962, there are probably more than a hundred such systems on the market or in development. Given this trend, it is always important to analyse the situation dispassionately and critically.

The ideal method is to make the measurement directly on the workpiece itself, thus excluding such influences as backlash and wear when determining the dimensions. At present this is a practical possibility only on numerically controlled grinding machines. There are some instances where the diameter of the workpiece is measured during the grinding operation itself by means of a capacitive or inductive pick-up. This method is called measurement control. Turning lathes have also been constructed having this form of control by automatic measurement of the workpiece, but in most cases it is far too complicated a concept for most practical applications.

In the great majority of numerically controlled machine tools the invariable practice is to deter-

mine the dimensions of the workpiece from the displacements of tool slides, turntables, etc. This method has the drawback that elastic deflections and play in bearings, in guideways and in the drive system, as well as the wear of the tool, adversely affect the accuracy of the measurement. For want of a better alternative, however, this method is still generally employed.

There are two methods of measuring the displacements of tool slides, etc. — the *direct* and the *indirect* method. The most usual is the indirect method, where a rotation is measured which is mechanically coupled to the linear movement of the slide. For example, it is general practice to take the angular displacement of the leadscrew as a measure of the movement of the slide. This method is simple and inexpensive, but not highly accurate, for the energy transmission inevitably gives rise to wear, which causes backlash, etc. The measuring accuracy that can be achieved is between 5 and 25 μm , depending on the type of machine.

An indirect method which has recently come into use depends on the conversion of linear motion into rotation by means of a rack and pinion. Being distinct from the slide drive, this transmission is subjected to scarcely any mechanical load, and therefore suffers hardly any wear. Although this indirect method is in principle preferable to the one just mentioned, its accuracy is not particularly high. The technique of making the rack and pinion does not in practice allow an accuracy much better than 10 μm .

The best indirect method is undoubtedly that which uses a leadscrew (not the one for the drive) provided with a recirculating ball nut (*fig. 7*). This leadscrew is not self-braking and can thus

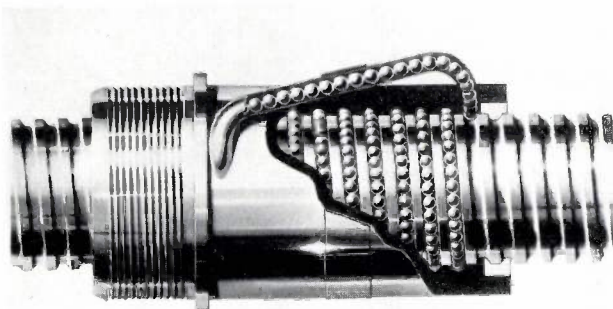


Fig. 7. Recirculating ball nut. The nut and the spindle have helical splines containing balls. The nut moves on these balls along the spindle in the same way as a normal nut along a screwthread. The balls leaving the nut at one end are forced along a channel to the other end and there brought back into circulation. The nut moves along the spindle with a minimum of play and wear. (Reproduced by courtesy of Bristol Siddeley Engines Ltd., Coventry, England.)

⁴⁾ John D. Cooney and Byron K. Ledgerwood, 31 numerically controlled point-to-point positioning systems, *Control Engng.* 5, No. 1, 67-98; No. 2, 99-122; No. 3, 99-114; 1958.

follow the movements of the slide, which it does with great precision. Here again, load and wear are negligible. The backlash may be a fraction of a micron, and the accuracy — depending on length and quality — is between 2 and 10 μm . In those cases where other methods are too inaccurate, this method is to be recommended in spite of its higher cost. It has previously been displayed by only one large European firm at international machine tool exhibitions.

We shall now consider the direct methods of measurement, in which the displacement of the slide is measured without mechanical transmission and converted into electrical signals. The oldest method was developed after 1952 by the National Physical Laboratory in cooperation with Ferranti Ltd. and uses optical diffraction gratings⁵⁾, which move relative to one another. The pitch of these gratings is very small (4 μm). At one side of the gratings a light source is mounted, and at the other side a few photocells. Making use of the moiré fringe effect (fig. 8) it is possible with this system to achieve an accuracy of 1 μm , the gratings being distributed over considerable lengths with a well-nigh unbelievable precision. A similar system is now being used in Germany, but it is very much simpler in design and less accurate (error 10 μm).

Another direct method has been developed by Farrand Optical Co. Inc., and has become widely known under the name "Inductosyn". This system is comparable with a synchro. "Two-dimensional coils" are mounted on a long stationary plate and a short movable plate, the turns on the coils following a special periodic pattern, so that relative movement between the plates can be measured as a change in the mutual inductance (fig. 9). A measuring accuracy of 2.5 μm seems to be obtainable over considerable lengths. Compared with optical and other

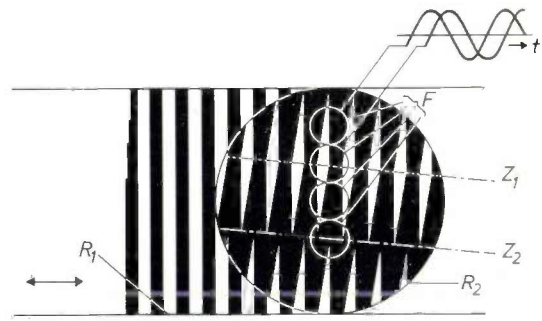


Fig. 8. Optical measuring system using diffraction gratings. The long grating R_1 is stationary; the short grating R_2 , which is turned through a small angle with respect to R_1 , is fixed to the slide whose movement is to be measured. At one side of the gratings a light source is mounted, and at the other side four photoelectric cells F . When the slide moves, the intensity of the light thrown onto the photoelectric cells varies. In the interference pattern light and dark moiré fringes appear, Z_1 and Z_2 , which move roughly perpendicular to the movement of the gratings. The photo-currents of the four photoelectric cells thus differ mutually in phase by 90° . From these currents both the direction and the magnitude of the displacement can be derived with an accuracy of $1/4$ of the pitch of the gratings.

systems, this method has the advantage of being very simple mechanically.

In addition to these high-precision measuring systems, mention should also be made of the systems of "electrical stops" or pins. Although one hesitates to give this the name of measuring system, it cer-

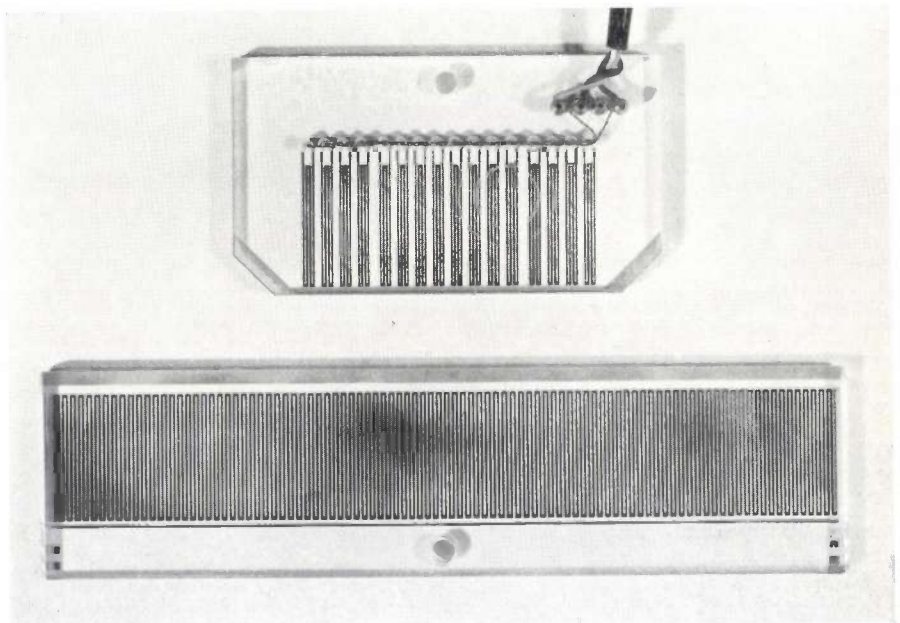


Fig. 9. "Inductosyn" measuring system. Two-dimensional coils are etched on two glass plates. The short plate is fixed to the slide whose movement is to be measured. The long plate, which for clarity is shown here separately, is stationary and mounted very close to the short plate. The system works in the same way as a synchro: the long coil can be regarded as the stator, and the short coil — consisting of two parts mutually displaced by $1/4$ period — is the rotor. One period corresponds to a linear displacement of 2.5 mm. If the long coil is fed with a sinusoidal voltage, two voltages differing by 90° in phase are induced in the short coil. From these voltages the position can be determined with an accuracy of 2.5 μm . (Reproduced by courtesy of E.M.I. Electronics Ltd., Hayes, England.)

⁵⁾ J. Guild, Diffraction gratings as measuring scales, Oxford University Press, London 1960.

tainly deserves to be mentioned as it has come into increasing use in recent years for many simple kinds of positioning machines. Parallel to the slide, various micro-switches are mounted on a shaft, the switches being operated by pins on the slide. The switches control the stopping and starting of the slide. The programme, derived from the information input, determines which movement is started or stopped when a particular switch is activated. In fact, then, there are only a very limited number of preset measuring points within the whole measuring range. Although it may seem primitive and not highly accurate, this system often has advantages because of its simplicity. The accuracy can be increased by the use of a sensing contact, which initiates a slow travel before the last part of the pass. The end of the slow travel is then marked by a cut-out contact or by a mechanical stop.

Apart from our classification into direct and indirect methods of measurement, a distinction can be made between digital and analogue systems, and between incremental and absolute systems. An analogue measuring system delivers a signal which varies continuously with the value to be measured, and this signal — if the system is absolute — also represents the absolute position of the slide. The digital system, on the other hand, measures in discrete steps, and the accuracy can therefore never be better than the "unit of measurement", although this can in theory be unlimitedly small. With incremental data processing, which is often used in combination with this system, only the change in position (increment) and not the position itself is expressed in multiples of the unit of measurement ⁶⁾.

In order to assess the merits of a given measuring system for a given machine, the actual transducer on the machine and the electrical processing of the measuring signal should be considered in relation to the information processing of the whole numerical control system. At the moment it can be said that in general analogue systems are often simpler, though less accurate, and difficult to adapt to the information processing part. Digital measuring systems, on the other hand, which are in themselves usually more complicated, are admirably adaptable to digital data processing, which is being increasingly used. It is frequently the reliability of the digital processing of the measuring signal and the information that decides the choice between absolute and incremental systems. It is to be expected that

incremental measurement — which saves a great deal of expensive electronic components — will win the day, now that the ideal of "foolproofness" seems to be within reach as regards the reliability of digital signal and data processing.

Data processing

The data processing for machine control (*fig. 10*) passes in principle through three stages:

- 1) From workpiece and machine specifications to the data carrier (e.g. punched tape).
- 2) Transfer of data to the machine tool.
- 3) Processing of the data by a computing unit in the machine tool into quantities that are intelligible to the machine parts (the measuring and control system).

In the case of contouring machines a general-purpose computer is indispensable for the first stage. Such computers are usually only an economic proposition in a mathematical computing centre. As mentioned in the foregoing, it is essential to ensure that the work planner and the programmer can form, in a relatively simple manner, the link between the drawing office and the workshop on the one hand, and the computing centre on the other. The further processing of the data in the computing centre presents no fundamental difficulties, since use can be made of well-known computer techniques. Of particular importance here are the phase and form in which the information is transmitted to the machine tool. This transmission should be regarded as a necessary interruption in the complete computing process, part of which takes place in the computing centre and part in the small computing unit in the machine. If the data carrier is a punched tape, it is important to keep the length of the tape and the size of the computing unit within reasonable limits. This simply formulated requirement is more paradoxical than trivial, for a short punched tape implies a voluminous computing unit, and vice versa. Further analysis shows that incremental position information offers considerable advantages here. While the length of the punched tape — compared with absolute position data — is virtually decimated, the size of the computing device is roughly halved. The resultant advantages are so considerable that it is surprising that the incremental system is used on so few contouring machines. Consequently recent exhibitions have contained many examples of machine tools using either very unwieldy lengths of punched tape or large and very expensive computing units.

In point-to-point positioning the situation is entirely different. The work planning is particularly

⁶⁾ This system is employed in a numerically controlled contour milling machine which was developed in the Philips Research Laboratories, Eindhoven. The machine will be described in a forthcoming article in this journal.

simple in this case, consisting essentially of drawing up a table of a few dozen numbers: for each machining phase one number for the position and one for the feed rate. A simple form of data supply is the plug-board, consisting of rows and columns of holes

servomotors. The major point of difference remains, of course, the data processing, which in positioning systems is much simpler. In practice, however, this simplicity appears only in the planning stage. At the machine itself the data processing shows much

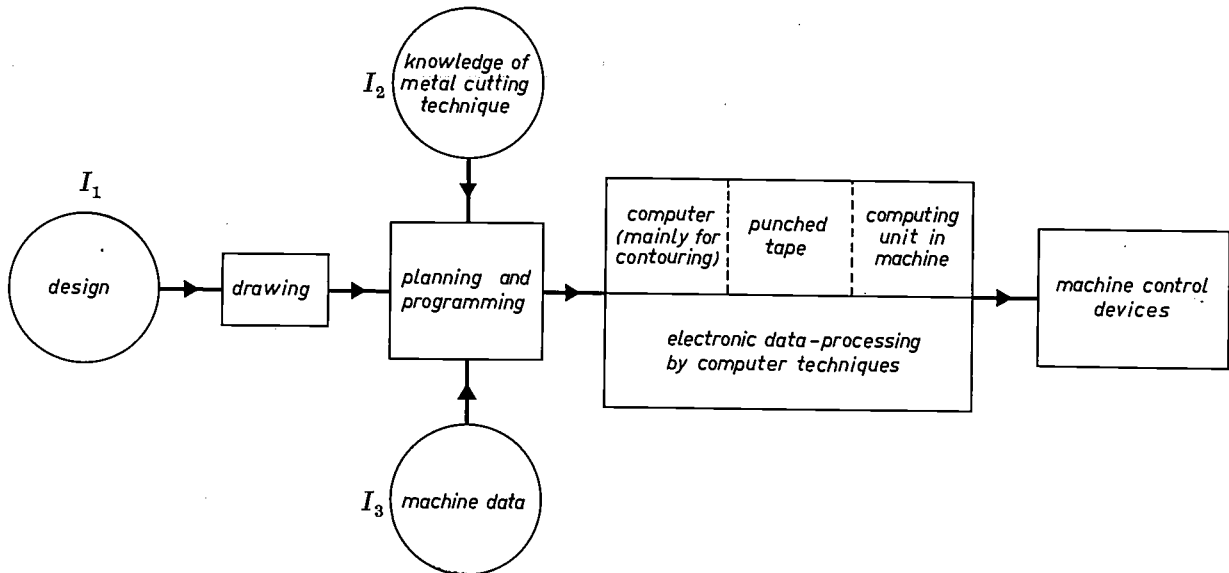


Fig. 10. Diagram of the data-processing system for numerically controlled machine tools. I_1 , I_2 and I_3 are the data sources. The data are conveyed to the machine by punched tape, which constitutes an interruption in the computing process.

into which plugs can be inserted. Each column may then, for example, represent one machining phase, and the plugs represent the kind of operation, the position and the speed. There are numerous versions of such boards and of the appertaining data processing. Some are more intricate than a computing unit belonging to a contour cutting machine, whilst others are so simplified that they can scarcely lay claim to the name of numerical control. A plug-board is almost invariably used in conjunction with a measuring system with a shaft with contacts and "electrical stops". The reason is obvious, for the limited possibility of supplying data using a plug-board ties up excellently with a measuring system of this kind.

In recent years the punched tape has come into increasing use in point-to-point positioning, together with a continuous measuring system. As far as larger and more expensive machine tools are concerned it appears that this combination will soon entirely supersede all other systems. Oddly enough, this brings us back to the same situation as with contouring, both as regards data feed and measuring system. Both systems are thus, from a technical point of view, growing towards each other, a trend which we also noticed in connection with hydraulic

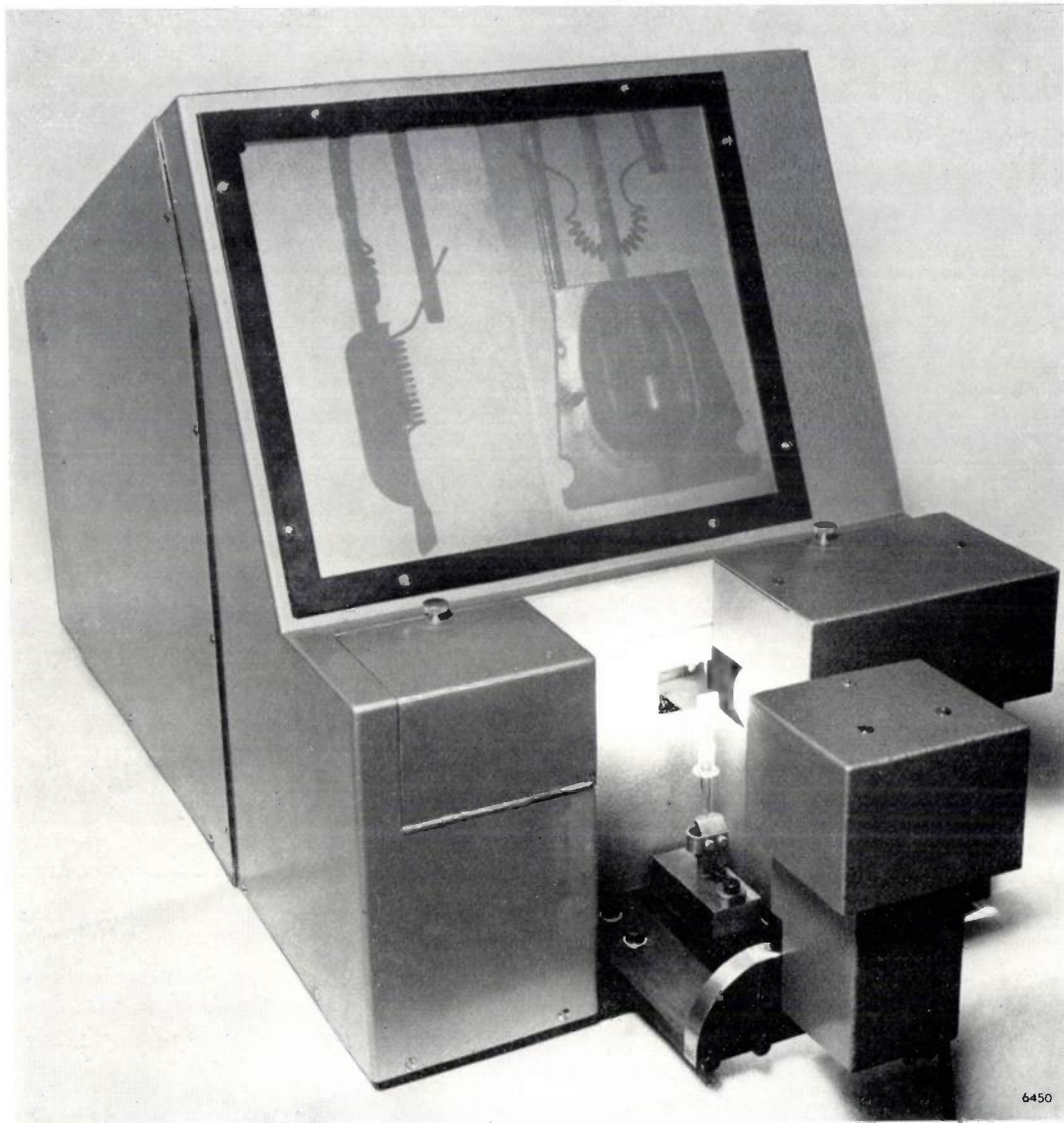
smaller differences in elementary computer technique.

The emphasis is increasingly being placed on (relative) simplicity and reliability, both of which have become ideals capable of realization as a result of transistorization and compact building methods.

Summary. Principal contents of the public lecture presented by the author upon his inauguration as reader at the Delft Technische Hogeschool. In automatic machine-tool control there are considerable advantages to be gained by feeding the data to the machine (e.g. milling machine or lathe) in numerical form. This can considerably improve both dimensional accuracy and productivity compared with traditional techniques. The problems which arise fall into two categories, one concerning data processing and the other the techniques of measurement and control. A distinction is also made between "contouring" and "positioning". As regards data handling the information about the workpiece must be supplied in a form capable of operating the control devices of the machine. In the case of large machines this is done by a general-purpose digital computer. This delivers a punched tape which is fed to a computing unit in the machine itself. This device passes the commands to the control parts of the machine.

A discussion of the engineering aspects of measurement and control shows that in most cases a hydraulic servomotor for the drive gives by far the best results. The numerous measuring systems in use for various machines are briefly reviewed, and it is predicted that digital incremental systems will, in the future, be most widely used for contouring machines, as they are better adapted to the data handling.

INSPECTION OF THE FILAMENT SYSTEM IN CAR-BULB MANUFACTURE



In the manufacture of "Duplo" bulbs for car headlamps, inspection is necessary to check whether the filament for the dipped beam (passing light) is correctly aligned in relation to that for the main beam (driving light). This is difficult to ascertain with the naked eye, particularly as the dipped-beam filament is partly enclosed in a shield. The inspection is carried out with the aid of a specially designed projector which produces side by side on a ground-glass screen a front view of the filament system and a silhouette of the system in side view, both ten

times magnified. The front view is obtained by simultaneous episcopic and diascope projection, using mercury light for the episcopic (incident) illumination and tungsten light for the diascope (transmitted) illumination. This method of projection shows up the contours of the filament system distinctly and with strong contrast on the screen, and also all surface details, as for example the state of the welded joints. The images are bright enough for inspection by daylight.

CRYSTAL GROWTH OF SILICON CARBIDE

548.52:546.281'261

In an investigation into the growth of silicon-carbide crystals ¹⁾ indications were found that the initial stage of this growth consists in the formation of "whiskers" — sometimes of very considerable length — built up from SiC with a hexagonal crystal structure (not the cubic structure, which can appear in later stages of the growth). The whisker formation seems to be independent of the chemical reaction that gives rise to the SiC, and also independent of temperature ²⁾.

¹⁾ W. F. Knippenberg, Philips Res. Repts, to be published shortly.

²⁾ See the article in reference ¹⁾. This also gives a list of references and discusses similar effects previously discovered, including the hexagonal SiC whiskers found by K. M. Merz.

The consequences of this discovery as regards the theory of the stability of various crystalline forms of SiC will be discussed elsewhere ¹⁾. Here we shall only mention a few results of an electron-microscopic investigation undertaken to follow the growth at the earliest possible stage. Another reason why this investigation seemed to us important was that SiC, being composed of two relatively light atoms, is fairly transparent to electrons at the voltages used by us on the electron microscope (about 80 kV). Because of this fact it was possible to obtain fairly detailed information on the internal structure of the whiskers.

Fig. 1 is a micrograph of relatively low magnification which illustrates a number of different forms



Fig. 1. Electron-photomicrograph of whiskers of various types produced in the reaction between carbon and quartz sand when heated for several hours at 1500 °C in an argon atmosphere. The whiskers are formed principally on the carbon particles. They were removed by placing against the particles a piece of copper-wire gauze, made sticky by immersion in a very dilute solution of polyester resin in amyl acetate. If the heating takes place at higher temperatures the crystals quickly grow long and thick enough to be seen with the naked eye.

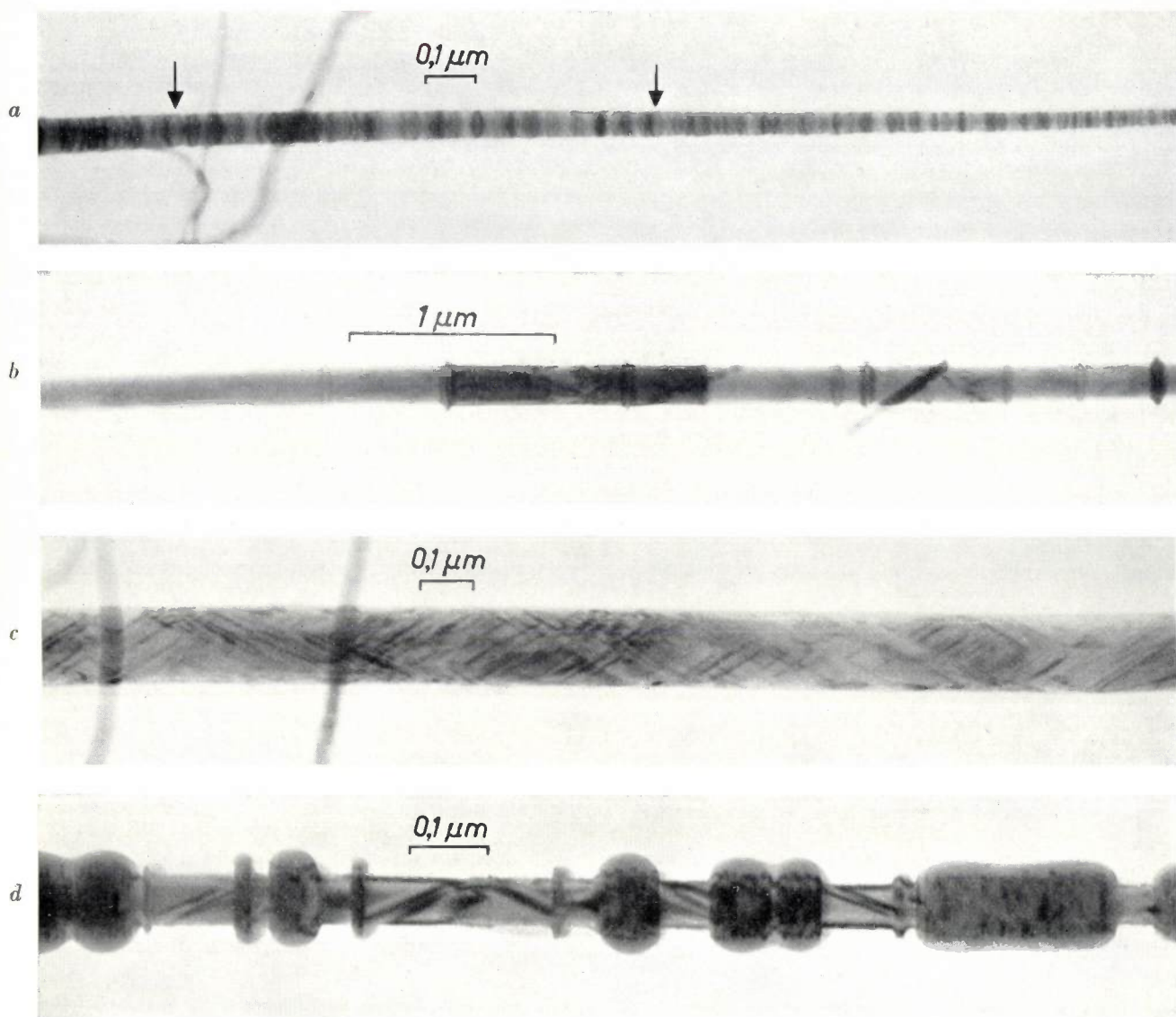


Fig. 2a-d. Electron-photomicrographs of various types of whiskers found when investigating the crystal growth of silicon carbide.

of whiskers. It shows the crystal growth found when carbon and quartz sand, intimately mixed, are heated for several hours at 1500 °C in an argon atmosphere. At higher temperatures the crystals formed grow rapidly to macroscopic size. Conspicuous in the micrograph are some crystals in the form of a string of beads, which occur in different sizes. Some of the other shapes are even more interesting when they are examined under higher magnification. *Figures 2a-d* show four whiskers in shapes that are regularly found and in which curious details are visible.

The tapered, needle-shaped crystal in fig. 2a appears from stereo micrographs to have a hexagonal cross-section. Numerous hexagonal bands can be seen perpendicular to the long axis (the hexagon is particularly clear at the position of the arrows);

when inspected very closely the bands give the impression of being the side faces of truncated hexagonal pyramids, so that the whisker, if one could feel it along its length, would not seem smooth but denticulated. The disturbances in the crystal structure, which cause locally altered interaction with the beam electrons and thus make the bands in the photograph visible, might possibly be attributed to repeated twinning. Although the hexagonal shape mentioned is no proof of a hexagonal crystal structure, it can well be informal evidence of it.

The crystal form in fig. 2b — the specimen in the photograph is about 5 times thicker than the one in fig. 2a — also shows contrasting bands, but much farther apart. The long, smooth sections between successive bands seem to have a perfectly

regular structure. This may be inferred in particular from the fact that the relative brightness of the various sections changes as the crystal is rotated in the electron beam. For, a relatively dark appearance of a section can be understood in the sense that a set of planes in one section of the crystal happens to be oriented in such a way that the incident electrons are reflected according to Bragg's law. The reflection angles can be so wide that the electrons are intercepted by the small aperture objective diaphragm and thus do not contribute to the image on the screen. When the crystal is rotated other sections will successively be oriented with respect to the electron beam in such a way as to give rise to the same effect, and thus these sections go dark — provided the crystal structure of each section is completely regular.

The crystal in fig. 2c is not "whiskery" but has the form of a greatly elongated platelet. This is seen from stereo micrographs. On crystals of this type two groups of parallel bands can be seen which cross each other at an angle of 60° . Even under higher magnifications, the denticulated contours found in the type in fig. 2a could not be seen. To explain the contrast of the bands we at first thought of a Bragg reflection as in fig. 2b. A certain group of lattice planes in a band might be oriented in the reflection position, and the observed picture with the short bands could then be explained by an elastic torsion of the whole crystal. This explanation cannot be correct, however, for when the whisker is rotated in the electron beam the bands are observed to turn with it (although this does alter the contrast). The bands are thus due to a local property of the crystal, as yet unexplained.

In the whisker in fig. 2b a slight thickening can

be seen in the bands which separate the smooth sections from each other. Such a thickening is very clearly apparent in the whisker shown in fig. 2d. The outward form of this type of whisker is, incidentally, the same as found in whiskers which have grown (at higher temperatures) to lengths of several millimetres, and which can therefore be examined using a light microscope. In the latter case it is easy to determine whether the large whiskers of SiC are cubic or hexagonal, because of the double refraction of the light caused by the hexagonal (non-isotropic) crystal. We have found that such whiskers contain alternate cubic and hexagonal parts. The question obviously arises whether this is also the case in the sub-microscopic whiskers as in fig. 2d. It is possible in principle to answer this from electron-diffraction patterns of the various parts of a whisker. We have photographed such patterns, but they are difficult to interpret.

Other investigators have found that in the case of some whiskers, observable under the light microscope, a central dislocation is present in the long axis of the whisker³⁾. It has been thought that this could be related to the mechanism of crystal growth in the form of whiskers. In view of the good transparency of the whiskers investigated here we ought to have been able to see such a central dislocation on our electron-micrographs, but on no single whisker was this detected.

W. F. KNIPPENBERG *),
H. B. HAANSTRA *),
J. R. M. DEKKERS *).

³⁾ For example D. R. Hamilton, J. appl. Phys. 31, 112, 1960.

*) Philips Research Laboratories, Eindhoven.

A MULTI-REFLEX KLYSTRON FOR USE IN MICROWAVE BEACONS

by B. B. van IPEREN *) and J. L. van LIDTH de JEUDE **).

621.385.623.5

At present there are no fewer than some 700 types of 3 cm microwave tubes on the market. Oddly enough, this does not mean that every designer of 3 cm equipment can find a suitable tube for his purposes; the number of types delivering a power of about 10 W is, for instance, very limited. The article below describes a new 10 W multi-reflex klystron for 3 cm waves, which can be frequency-modulated at 50 c/s over a range of more than 200 Mc/s. This makes it especially suited for use as a transmitting tube in marine beacons of the ramark type. Partly because the tube can be mechanically tuned over 8%, there are good prospects also of other specific applications.

In recent years the possibility has been investigated of using suitably situated microwave beacons of the ramark type as an aid to shipping. These beacons would help navigation rather in the way that lighthouses do at night, given good visibility. The beacons are designed so as to produce on the radar screen of all ships in the neighbourhood a very narrow sector filled with bright straight lines (the apex angle being equal to the beam angle of the antenna of the ship's radar receiver, e.g. about 2°), pointing in the direction of the site of the beacon. Just as the beam from a lighthouse is not continuous but has a certain "character" — e.g. every ten seconds two flashes with an interval of two seconds — each microwave beacon is periodically switched on and off in accordance with a particular programme¹⁾, by which it is identified.

In order that the radar receiver on every ship can pick up the beacon (in this case the ship's radar transmitter serves no purpose), the beacon must transmit all frequencies within the civil radar band (9320 Mc/s to 9500 Mc/s). In practice this can be achieved most simply by tuning the beacon transmitter to about 9400 Mc/s (wavelength about 3.2 cm) and frequency-modulating the beam with a swing of about 100 Mc/s²⁾. The modulation frequency used is roughly 50 c/s: thus in the time which the ship's rotating radar antenna

takes to pass the beacon once, the whole frequency band is traversed at least once. Upon each revolution of its radar antenna, every ship can therefore pick up a signal from the beacon. (This is again assuming a beam angle of 2° and an antenna that completes a revolution in about 2 s.) For various reasons it is necessary to use additional frequency modulation with a swing of 10 to 15 Mc/s and a modulation frequency of about 10⁴ c/s.

The above-mentioned requirement that it should be possible to frequency-modulate the beacon within a range of about 200 Mc/s cannot be met by existing types of microwave transmitting tubes at the power of about 10 W required here. The tube described in this article, a multi-reflex klystron, has therefore been specially designed for use in ramark beacons of the kind described. The principle of its operation differs considerably from that of the 12 cm and 8½ cm multi-reflex klystrons previously described in this journal³⁾. In spite of the much shorter wavelength, the new tube has nearly the same high efficiency as the earlier types, i.e. 15 to 20%, at 9 to 12 W. In this range of wavelengths and output powers a normal klystron gives an efficiency of only 5 to 7%.

The new tube is simple in construction, partly owing to the use of a glass bulb, which keeps the production costs relatively low.

It also proved possible, without complicating the construction, to make the tube mechanically tunable over a range of about 8%; as a result the tube is also suitable for other applications.

*) Philips Research Laboratories, Eindhoven.

***) Formerly with Philips Research Laboratories, Eindhoven.

¹⁾ The principle of ramark beacons is described in: The use of radar at sea (editor F. J. Wylie), Hollis and Carter, London 1952. The name "ramark" is derived from: "MARK obtained by RADio means".

²⁾ For a further description of this type of ramark beacon, see J. M. F. A. van Dijk, N. Schimmel and E. Goldbohm, A ramark beacon for use with marine radars, Int. conf. on lighthouses and other aids to navigation, Scheveningen 1955.

³⁾ The klystrons for 12 and 8½ cm wavelengths are described by F. Coeterier in Philips tech. Rev. 8, 257, 1946 and 17, 328, 1955/56. The development of the new tube was also started under the direction of F. Coeterier, who was with Philips Research Laboratories, Eindhoven, until 1959. See: F. Coeterier, Tubes à réflexions multiples, L'onde électrique 36, 917-919, 1956.

The frequency modulation of the carrier with a frequency of about 10^4 c/s and a swing of 10 to 15 Mc/s is necessary for reasons connected with the design of the ship's radar receiver. The echo signal, consisting of pulses of e.g. $1 \mu\text{s}$ width at intervals of $999 \mu\text{s}$ (giving a pulse repetition frequency of 1000 c/s) is detected in the radar receiver, i.e. converted from a series of wave trains into a corresponding series of square-wave pulses. The resultant signals are amplified in the video amplifier and applied to the cathode-ray tube. The video amplifier therefore need not pass frequencies lower than the pulse repetition frequency.

The signal from a remark beacon without the frequency modulation of 10^4 c/s would result, after detection in the receiver, in a square-wave signal having a repetition frequency of 100 c/s, or owing to the 50 c/s modulation the frequency of the transmitter traverses the narrow frequency band to which the receiver is tuned 100 times per second. As we have seen, however, a signal of such a low frequency would not pass the video amplifier. Given a modulation of 10^4 c/s, and a frequency swing such that the beacon frequency upon passing the reception band does in fact pass "in and out" 10^4 times per second, then instead of the 100 c/s signals we in fact obtain corresponding groups of $10\ 000$ c/s pulses. These pulses are passed by the video amplifier and produce on the radarscreen a line of dashes running from the centre to the edge of the screen in the direction of the beacon. The number of dashes in the line obviously increases with the modulation frequency. The value for the latter should therefore be much higher than the pulse repetition frequency.

On the other hand, the modulation frequency should preferably not be higher than 1 or 2×10^6 c/s, as otherwise a property of the remark beacon of considerable practical importance would be forfeited. We refer here to the possibility on the ship of eliminating, or at least strongly attenuating, the beacon signal with the aid of an F.T.C. (Fast Time Constant) circuit, used in many radar installations for suppressing troublesome reflection from raindrops. This can be particularly desirable when the ship is close by the beacon. The effect of switching-in the F.T.C. circuit is to raise the lower cut-off frequency of the video amplifier to a value of say 10^6 c/s.

Operation of the tube

Principle

To make it clear in what respects the new multi-reflex klystron for 3 cm waves differs from those for operation at longer wavelengths, we shall briefly recapitulate the operation of reflex and multi-reflex klystrons. In any klystron an electron beam passes two closely spaced grids M (see fig. 1a) between which an RF field is maintained. In this field the electrons are modulated in velocity, and further on this gives rise to modulation in intensity: later electrons passing M are accelerated and catch up with the earlier ones that have been retarded. In this way an alternating current can be induced in a circuit connected to the two grids I . The frequency of this alternating current is identical with that of the voltage on the electrodes M .

In a reflex klystron (fig. 1b), both M and I are formed by the same pair of electrodes. After passing

M , the electrons enter a retarding field where they are returned. On the way back the "bunches" of electrons meanwhile formed induce an alternating current in M , just as they did in I in fig. 1a. Since M in this case is not formed by a pair of electrodes but usually by a gap in the wall of a resonant cavity, we shall from now on refer to M as the interaction space or gap. The frequency here is nearly the same as the resonant frequency of the cavity.

In a multi-reflex klystron (fig. 1c) a portion of the returning electrons is again reflected in the field between M and the cathode K ; in the figure these again travel to the right and again pass the interaction gap. If all parameters are properly chosen, an electron can be made to pass to and fro several times and to give up a large part of its energy to M . Obviously it is necessary to ensure that each bunch passes the interaction gap at a moment when the electric field retards the electrons. In multi-reflex klystrons the electron beam is held together by means of a longitudinal magnetic field, generally produced with the aid of a permanent magnet.

In looking for a suitable form for the electrostatic retarding field in the tube, conflicting requirements are encountered. On the one hand the electrons that pass M at a given moment as a bunch should also be bunched when they return. (What they do

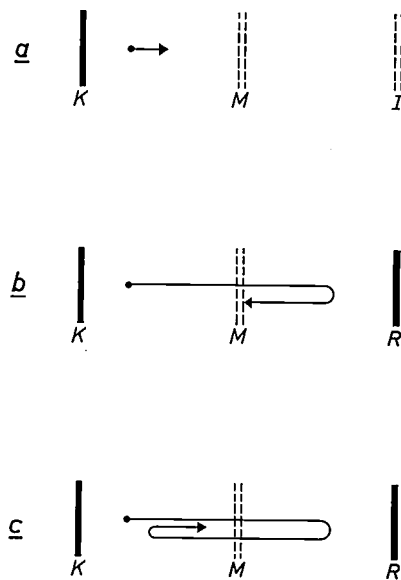


Fig. 1. a) Principle of the operation of a klystron. K cathode, M modulation (buncher) grids, between which the electrons are modulated in velocity. In the field-free (drift) space between M and I the velocity modulation gives rise to intensity modulation. The electrons pass the grids I in "bunches" and induce in I an alternating current.

b) Principle of reflex klystron. The functions of M and I are assumed by one pair of electrodes. A retarding field to the right of M , produced by a reflector electrode R , causes the electrons to return.

c) Principle of multi-reflex klystron. Here the electrons are reflected also in the cathode space and thus travel several times to and fro before disappearing.

outside M is of no consequence here.) This means that the time the electrons take to travel there and back — the transit time — should be the same for all electrons: the transit time, then, must not be dependent on the speed at which the electrons passed M . This requirement is met if the form of the potential is parabolic and provided that M is at the position of the maximum.

On the other hand the velocity modulation produced by M must be changed in the reflected beam to intensity modulation at the position of M . This of course requires a variable transit time.

In the tubes for 12 and $8\frac{1}{2}$ cm wavelengths this problem is solved in the manner sketched in *fig. 2a*. On the right of the interaction gap M , the retarding space is larger than on the left. In the extra part, right of R , the potential curve deviates from the parabolic shape. The whole is so dimensioned that electrons retarded upon their first transit through M do not reach point R at all; those which are accelerated do reach R , where they remain in the space right of R , the "waiting space", long enough to be able to return together with electrons that were retarded in M in the next half period. In returning, these are also of course retarded in M , and never again reach the waiting space: the bunch once formed remains intact.

In practice the parabolic potential curve is approximated by flanking a "drift space", i.e. a field-free space (space of constant potential), by two spaces having a linearly varying potential (constant

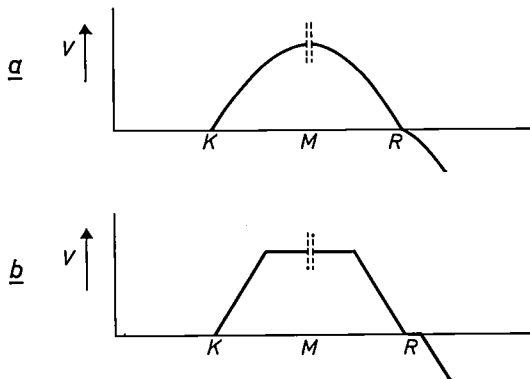


Fig. 2. a) The conflicting requirements of a constant transit time and of bunching can be resolved, in multi-reflex klystrons for wavelengths of about 10 cm, by giving the field the form illustrated in the sketch. Between K and R the potential V varies parabolically — which exactly fulfils the first requirement — but to the right of R the field deviates from the parabolic form. The point R is only reached by the electrons that were accelerated upon their first transit through M . These remain in the space to the right of R — the "waiting space" — for exactly half a period, and return together with the electrons that passed M half a period later and were there retarded. The bunch now formed always passes M in the retarding half period and remains intact because it can no longer reach R . b) In practice the potential form in a) is approximated by combining drift spaces with spaces in which the potential varies linearly.

field strength). The "waiting space" too consists of a field-free part and a part with a linearly varying potential; see *fig. 2b*. To obtain this potential variation at least two repeller (reflector) electrodes are needed: one at R with a hole through it, and the other, which can be solid, more to the right in the figure.

Since it is not possible to make a tube for 3 cm waves on the waiting-space principle by simply scaling-down the dimensions of a tube for $8\frac{1}{2}$ cm waves, a different solution had to be found for the new tube. The principle adopted is illustrated in *fig. 3*. Here, too, the retarding space on the right is larger so as to accommodate the accelerated electrons, but the potential curve remains parabolic. The required intensity modulation is obtained in this case by using three interaction gaps; one of them is again on the axis of the parabola, and the other two are symmetrically disposed on either side of it. Because of their position the latter two gaps do not meet the above-mentioned requirement for obtaining a constant transit time: the change in velocity which an electron undergoes in one of the outer gaps gives rise to a change in the transit time for the relevant half cycle, and thus produces bunching⁴⁾.

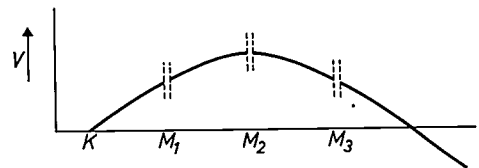


Fig. 3. In the new multi-reflex klystron for 3 cm waves, bunching is not due to a discontinuity in the parabolic form of the potential, but to the use of three interaction gaps M_1 , M_2 and M_3 , two of which are outside the symmetry plane of V .

In the next section we shall consider all this quantitatively. Here it may be mentioned that the three interaction gaps are part of three tightly coupled resonant cavities — one gap in each — which resonate in such a way that the voltages on M_1 and M_3 are mutually in phase and are in anti-phase with the voltage on M_2 . All three gaps are situated in the drift space which here too approximates to the central section of the required field (*fig. 4*).

The reason why a tube operating on the "waiting-space principle" cannot be radically reduced in size is that too high demands would be made upon the cathode. If the earlier described $8\frac{1}{2}$ cm tube were to be proportionately scaled down for operation at 3.2 cm, this would involve making the current density at the cathode about $7 \times$ greater, i.e. about 14 A/cm^2 .

⁴⁾ See also the last of the articles quoted in footnote ³⁾.

The situation is much more serious than appears from this figure, since the cathode in the multi-reflex tube is subjected to heavy ion-bombardment: in travelling repeatedly backwards and forwards, the electrons — just as in a Penning gauge — cause relatively strong ionization.

No cathode can deliver this current density under these conditions, at least not if the tube is to have a long life. In scaling down the tube, then, we are forced to make the cathode relatively larger, which necessitates doing the same to the openings in the electrodes and the resonator. As a result, however, it is no longer readily possible to produce the required field in the "waiting space". The sharp discontinuity in the field, which is required for this space, is found only at the edge of the hole in the electrode. In the middle of the hole the required field form is better approximated the smaller is the hole. In a relatively large hole the field in the greater part of the opening follows a smooth curve.

Nevertheless the electrode configuration of the new tube outwardly resembles that of the old rather more than might be supposed from fig. 4b. Again there are two repeller electrodes. The inner one now serves, however, to eliminate "sagging" of the field due to space charge effects, i.e. to more or less maintain the linearity of the retarding field in spite of the space charge. Without this measure there would be considerable distortion of the retarding field, which would seriously affect the operation of the tube. The difference in function compared with the inner repeller electrode of the 8½ cm tube is immediately apparent from the applied voltage. In the 8½ cm tube this electrode is maintained at cathode

potential, whereas in the new tube its potential is substantially higher.

Analysis of the motion of the electrons

We shall now show how the behaviour of the electron beam can be analysed and the alternating-current component calculated for the new tube at the position of M_2 in the beam. We start from the assumption that the alternating voltage across the gaps is small compared with the potential difference V_0 through which the electrons pass before entering the drift space where the gaps are situated. The change in their velocity inside the gaps is thus relatively small.

In an electron beam which is velocity-modulated at M (fig. 1) by a pair of grids across which the voltage is $Ue^{j\omega t}$, the alternating-current component i_{\sim} , at I , is given by:

$$i_{\sim} = \frac{1}{2} j G_0 \varphi_D U e^{-j\varphi_D + j\omega t} \quad (1)$$

Here G_0 is equal to the direct-current component I_0 divided by V_0 , and φ_D is ω times the length of time which a non-accelerated electron takes to cover the distance MI).

If the electrons after passing M are made to move not in a field-free drift space but in a retarding field having a linearly varying potential, so that they can return to M , the alternating current component at M in the returning beam is then given by:

$$i_{\sim} = -\frac{1}{2} j G_0 \varphi_R U e^{-j\omega_R + j\omega t} \quad (2)$$

Except for the minus sign and the replacement of φ_D by φ_R — the corresponding quantity for the reflecting space — this formula is identical with (1). (The sign is reversed because in a retarding space with a linearly varying potential the slow electrons have a shorter transit time than the fast ones, whereas in the drift space the reverse is the case.) In our calculation, therefore, we can imagine the retarding space to be replaced by a drift space with the same average transit time, provided we reverse the sign of the velocity modulation (rule 1).

If we apply this to a tube in which the electron successively passes through a number of drift spaces, separated by spaces with a linearly varying potential, we can write:

$$i_{\sim} = \frac{1}{2} j G_0 U \left(\sum_n \varphi_{D,n} - \sum_m \varphi_{R,m} \right) e^{-\left(\sum_n \varphi_{D,n} + \sum_m \varphi_{R,m} \right) + j\omega t} \quad (3)$$

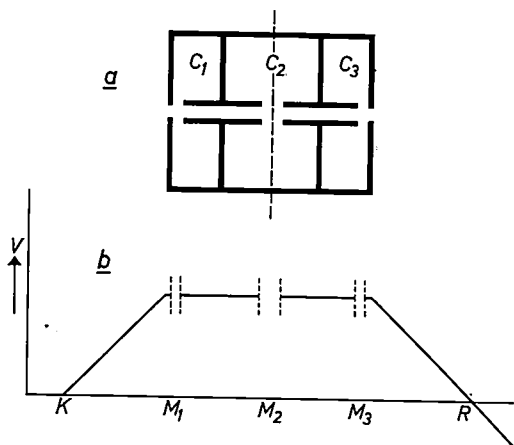


Fig. 4. a) The resonator consists of three tightly coupled resonant cavities C_1, C_2 and C_3 , each of which contains one of the three interaction gaps M . b) Form of the potential. The three interaction gaps M_1, M_2 and M_3 are situated in a central drift space. This is flanked by retarding spaces in which the potential varies linearly. This form of potential variation is again an approximation to the parabolic form.

5) This formula — or at least a formula of which the present one is a simplified version — can be found in numerous books about klystrons, e.g. A. E. Harrison, Klystron tubes, McGraw-Hill, New York 1947 (Ch. 3) or D. R. Hamilton, J. K. Knipp and J. B. Horner Kuper, Klystrons and microwave triodes, McGraw-Hill, New York 1948 (Ch. 9).

Here $\varphi_{D,n}$ holds for the drift space with rank number n and $\varphi_{R,m}$ for the retarding space with rank number m . The exponent contains between parentheses the sum of $\Sigma\varphi_D$ and $\Sigma\varphi_R$, that is ω times the total time taken to traverse the whole system; written before the exponential function is the difference of $\Sigma\varphi_D$ and $\Sigma\varphi_R$, due to the minus sign in (2).

In order to investigate the behaviour of the electron beam in the new tube, we must also take into account the case where there is more than one set of modulation (buncher) grids. The frequency of the voltage on these grids is identical, but the amplitude and phase usually differ. The amplitudes are small compared with V_0 , so that when calculating the resultant alternating-current component we may, as can be demonstrated, simply take the sum of the alternating currents that would be produced by each pair of grids — i.e. each gap — separately. (This is rule 2, the superposition theorem.)

Using the two above-mentioned rules, we can easily calculate for the middle gap the AC component of an electron current in an electron beam which, coming from the cathode in the new tube, has passed through the three gaps and, after reflection in the repeller space, has again passed the gaps on its way back. In doing so we must remember

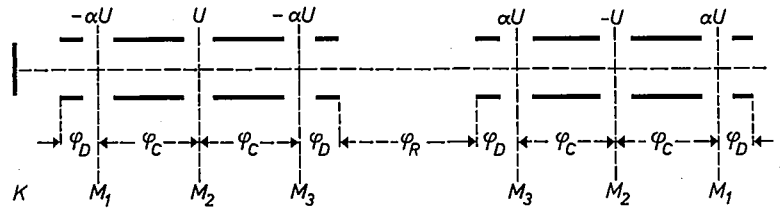


Fig. 5. Diagram for calculating the alternating-current component of the electron beam at M_2 in the new tube. The figure relates to a beam which, coming from the cathode, has passed the three interaction gaps once and has then returned to the gaps after a single reflection. The retarding space is represented by an equivalent drift space with transit time φ_R . The quantities φ_C and φ_D are the transit times in the (actual) drift spaces between and adjacent to the gaps. The voltages on the gaps M_1 and M_3 are opposite in phase to that on M_2 and have a different amplitude.

— and this is the third rule — that the sign of the gap voltage has to be reversed for the returning electrons (fig. 5). This calculation is presented in small print at the end of this part of the article.

An illustration of the electron movement can be seen in the diagram in fig. 6. The abscissa of this figure is equivalent to that of fig. 5. From left to right we again see the cathode space, and four drift spaces separated by the three gaps M_1 , M_2 and M_3 , etc. Since the retarding space (cf. fig. 5) is represented by an equivalent drift space, the whole abscissa can be regarded as the positional axis of the electron: as a measure of all distances the time is chosen which a hypothetical undisturbed electron would need to travel those distances. The unit of time adopted is the oscillation period T of the gap voltage ($T=2\pi/\omega$).

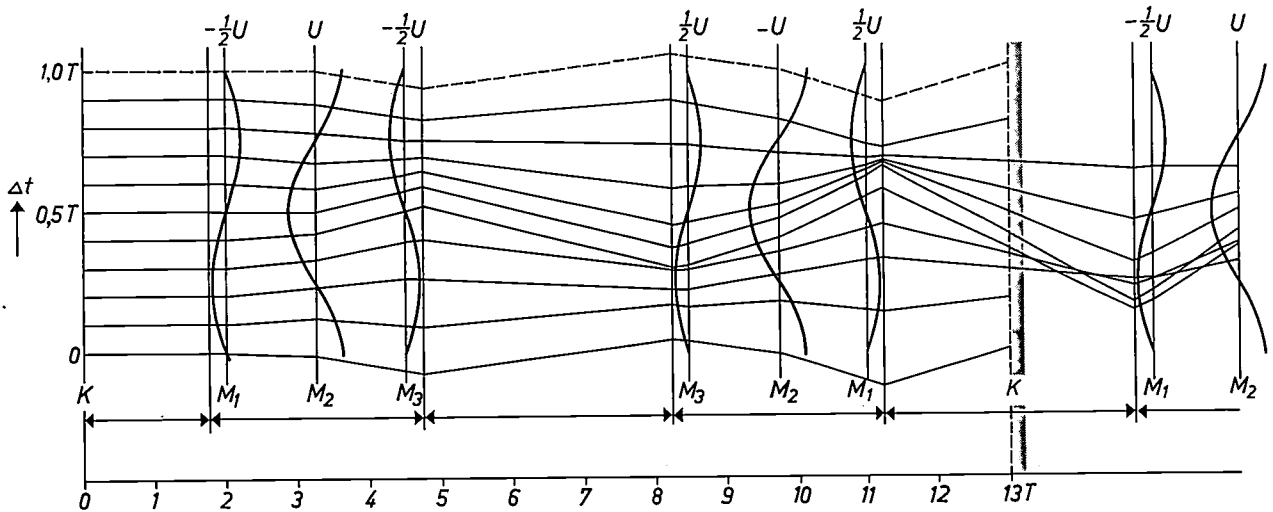


Fig. 6. Transit-time diagram of 10 electrons that leave the cathode K within one oscillation period of the gap voltage. The abscissa of this figure is equivalent to that of fig. 5. The "distances" between the electrodes and gaps are again the times φ which a hypothetical, undisturbed electron would need to traverse them. On the vertical axis is plotted the difference between the time at which a certain electron reaches a given point and that at which the hypothetical, undisturbed electron would reach the same point. Along the verticals representing the

three gaps M_1 , M_2 and M_3 , the potential is drawn which the electrons encounter at the moment of passing these gaps.

Of the ten electrons, seven are retarded after one complete oscillation and form a bunch which, upon the next oscillation, passes M_2 in the right phase of the field; see the part of the figure to the right of the shaded line. Only three are accelerated and return to the cathode. The transit time between the gaps is $1.25 T$, that in the drift space outside the gaps is $0.25 T$ and that in the reflector space $3.5 T$.

On the ordinate is plotted the difference Δt between the time at which a given electron reaches a given point and that at which an undisturbed electron starting at $t = 0$ would reach the same point. For an undisturbed electron the ordinate is thus constant and equal to the time of departure; the behaviour of such an electron is represented in the diagram by a straight horizontal line. For an electron with a greater (but constant) velocity, a line with negative slope is found in the drift spaces, and for a slower electron a line with positive slope. The opposite is the case in the cathode and repeller spaces, for these spaces can only be regarded as equivalent drift spaces if the sign of the modulation has been reversed (rule 1).

Since the ordinate is a time axis, the gap voltages encountered by electrons leaving the cathode at different times can easily be derived in the figure from the sinusoidal line around the vertical lines marking the position of each gap. Since these are values encountered by one and the same electron during its transit, the voltages at the outer gaps are drawn in antiphase (the "distance" being $2\frac{1}{2} T$), although in fact they oscillate in phase.

Fig. 6 shows the behaviour of ten electrons leaving the cathode at equal intervals within a single period, for the case where the amplitude U of the voltage at the centre gap is $0.1 V_0$ and that on the two others is half this value. Measurements of U carried out by us have shown that this choice roughly corresponds to the actual situation in the tube; U is about 380 V — which, owing to the finite gap width, amounts effectively to about 270 V — whilst V_0 is 3000 V. Upon passing each gap the electron's velocity varies, and so too therefore does the slope of the relevant line. As the electron enters and leaves the repeller and cathode spaces, the slope of each line changes its sign in accordance with rule 1.

It can be deduced from the diagram that the behaviour of the electron beam is favourable for the operation of the tube in three respects. In the first place we see that most electrons pass the individual gaps during the retarding phase of the field, so that the beam gives up energy at every gap. It can be seen further that after one complete oscillation (i.e. at K on the right of the figure), the lines have almost returned to their initial height, in other words that the transit time of all electrons is roughly the same for one complete oscillation, and equal to that of an undisturbed electron. Finally, it is seen that after one complete oscillation no fewer than seven of the ten electrons have formed a bunch which passes a given point within half a period, and that all electrons in this bunch

have acquired a velocity lower than that of an undisturbed electron. These electrons therefore can no longer reach the cathode, and remain travelling to and fro between cathode and repeller, thus repeatedly passing the centre gap in the correct phase of the field. The beginning of this process is represented in the part to the right of the shaded edge in fig. 6.

In the initial oscillations relatively little energy is given up to the outer gaps, since the distance from the latter to the centre gap is traversed in $1\frac{1}{4}$ periods. The more energy the electrons lose, that is the slower they become, the longer becomes the transit time, of course, so that a bunch of electrons already retarded gives up energy to each gap every time it passes through one.

The three electrons that remain outside the electron bunch in the first oscillation are, as can be seen from the figure, precisely those which acquired excessive velocity and thus vanish again in the cathode. The bunch able to carry out a number of oscillations evidently consists exclusively of "favourable" electrons.

The accelerated electrons then bombard the cathode, thus raising the temperature of the cathode when the tube is oscillating.

The alternating current i_{\sim} in the electron beam at the centre gap M_2 , after the beam has passed once through the reflector space, is calculated as follows. We omit the recurring factor $e^{j\omega t}$ and so write U for the voltage across the centre gap. The voltage across the other two gaps is then $-aU$, where a is provisionally arbitrarily positive. As described, we imagine the retarding space to be replaced by an equivalent drift space (fig. 5). For the second transit the sign of the gap voltages must be changed, and for calculating the amplitude of i_{\sim} the transit time φ_R must be taken as negative. Using fig. 5 and formula (3) to calculate the contribution of each gap, we immediately find the required current:

$$i_{\sim} = \frac{1}{2} j G_0 U [-(3\varphi_C + 2\varphi_D - \varphi_R) a e^{-j(3\varphi_C + 2\varphi_D + \varphi_R)} + (2\varphi_C + 2\varphi_D - \varphi_R) e^{-j(2\varphi_C + 2\varphi_D + \varphi_R)} - (\varphi_C + 2\varphi_D - \varphi_R) a e^{-j(\varphi_C + 2\varphi_D + \varphi_R)} + \varphi_C a e^{-j\varphi_C}] \dots \dots \dots (4)$$

We shall now try to choose the various transit times as favourably as possible. First of all we note that an electron bunch, once formed, must continuously pass the centre gap in the favourable phase (i.e. in a retarding field). For this purpose the transit time from the centre gap via the retarding space and back again must be $(n + \frac{1}{2})$ periods, where n is an integer. The first condition thus reads:

$$2\varphi_C + 2\varphi_D + \varphi_R = (n + \frac{1}{2}) 2\pi \dots \dots \dots (5)$$

Substitution of this in (4) gives:

$$i_{\sim} = + \frac{1}{2} j G_0 U [(4\varphi_C + 2\varphi_D - \varphi_R) a e^{-j\varphi_C} - (2\varphi_C + 2\varphi_D - \varphi_R) + (\varphi_C + 2\varphi_D - \varphi_R) a e^{j\varphi_C}] \dots \dots \dots (6)$$

The energy given up by the beam to the resonator is governed by the real part of $i\tilde{\omega}$:

$$\text{Re}(i\tilde{\omega}) = \frac{3}{2} G_0 U \alpha \varphi_C \sin \varphi_C.$$

This quantity is roughly maximum when $\sin \varphi_C$ is maximum, which gives the second condition:

$$\varphi_C = (m + \frac{1}{2}) 2\pi. \dots \dots \dots (7)$$

With this choice of φ_C we find moreover that at all gaps the real part of the alternating current in the beam is in anti-phase with the gap voltage, so that the beam in every transit loses energy to the resonator.

Finally, there is a third condition, which is that the oscillation time should be independent of the electron velocity. From fig. 5 it might be deduced that this implies $\varphi_R = 2(\varphi_C + \varphi_D)$. However, formulae (1) and (2), from which we concluded that the reflector space can be replaced by an equivalent drift space, are valid only for relatively small velocity variations. The equality derived from fig. 5 therefore holds only for the first reflection, but not for the last reflections which an electron undergoes before falling on the resonator. Its energy has then decreased on average to about 70% of its initial value, and its velocity v to roughly 85% of the initial value v_0 . We must therefore look for the conditions under which the oscillation time is as constant as possible for electron velocities between v_0 and $0.85v_0$. Now the exact expression for ω times the transit time of an electron from the centre gap to the retarding field and back is:

$$\omega\tau = \frac{v}{v_0} \varphi_R + 2 \frac{v_0}{v} (\varphi_C + \varphi_D).$$

Within a certain range of v values, the transit time is constant if

$$\frac{d\omega\tau}{dv} = 0, \text{ that is if } \left(\frac{v}{v_0}\right)^2 = \frac{2(\varphi_C + \varphi_D)}{\varphi_R}. \dots \dots (8)$$

In our case we can usefully choose $(v/v_0)^2$ halfway between 1 and 0.7, so that in (8) we must put $(v/v_0)^2 = 0.85$. The third condition is thus found to be:

$$2(\varphi_C + \varphi_D) = 0.85\varphi_R. \dots \dots \dots (9)$$

The three conditions so derived still leave us some freedom in the choice of the various transit times. In making that choice we were guided by the consideration that the length of the beam between cathode and reflector should be as small as practicable in connection with the required magnetic field. This led to the choice $m = 1$. The gap distance φ_C , according to (7), is then 1.25. At the smallest possible value of m , i.e. zero, we would then have $\varphi_C = 0.25$ and the alternating-current component would be too small; cf. formula (1). At $m = 2$ or $\varphi_C = 2.25$ the advantage of the somewhat higher current no longer offsets the disadvantage of a heavier magnet.

For φ_D and φ_R we chose the smallest values compatible with the other two conditions. From the conditions (5) and (9) we derive:

$$\varphi_C + \varphi_D = 0.23 (n + \frac{1}{2}) 2\pi.$$

The minimum value of n at which, in combination with $\varphi_C = 1.25$, this leads to a positive value for φ_D , is $n = 6$. It then follows that

$$\begin{aligned} \varphi_C + \varphi_D &= 1.5, \\ \text{yielding } \varphi_D &= 0.25 \text{ and } \varphi_R = 3.5. \end{aligned}$$

For φ_D this is at the same time roughly the minimum possible value that can be obtained in view of the construction.

The resonator

As mentioned, the resonator consists of a combination of three coupled resonant cavities. The

centre cavity is tightly coupled with a fourth one to enable the tube to be tuned. The fourth cavity is partly outside the vacuum, so that the resonant frequency can readily be altered. A highly simplified sketch of the resonator system will be found in fig. 7.

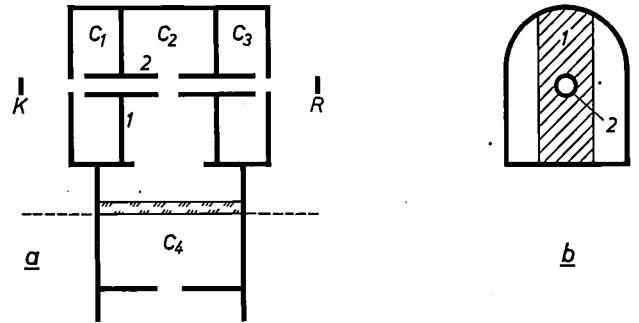


Fig. 7. The resonator of the new multi-reflex klystron. The tube is tunable because the three cavities C_1, C_2 and C_3 (cf. fig. 4) are coupled via C_2 with a fourth cavity (C_4) whose resonant frequency can be varied. The latter cavity is partly formed by a space inside the tube and partly by a section of waveguide to be connected to it (in (a) below the dashed line). The tight coupling between C_1, C_2 and C_3 is obtained by arranging the walls as sketched in (b). The walls are reduced to strips 1 which carry the cylinders 2 through which the electrons travel.

We shall first consider the three-cavity system inside the vacuum. The system may be regarded as a combination of three inductively coupled LC circuits — we shall refer to the L and C as the *equivalent* inductance and capacitance respectively — the first and third of which are identical (fig. 8a). Such a system has three resonant frequencies. It can immediately be seen that at one of these, which we shall call ω_2 , the outer two circuits are oscillating in anti-phase and the middle one is “dead”. For this frequency we can write: $\omega_2 = (L_1 C_1)^{-\frac{1}{2}}$. This symmetrical form of oscillation is ruled out in our tube, because in this case there is no voltage on the centre gap and therefore no coupling with the fourth cavity.

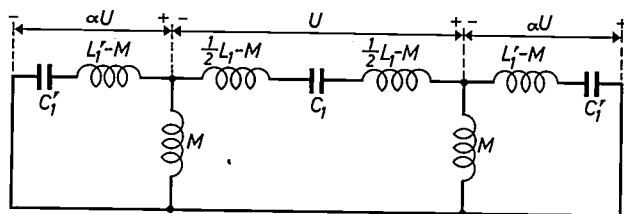
The other resonant frequencies are more easily found by first re-drawing the diagram as in fig. 8b. Since, for reasons of symmetry, there can never be a voltage between points P and Q , we can cut the diagram in half for our purposes (fig. 8c). We then obtain a system of two inductively coupled circuits, thus giving two resonant frequencies. The lower one we call ω_1 and the other ω_3 . For the oscillation with frequency ω_1 , the voltage U on the centre gap is opposite in phase to the voltage αU on the two others. For the oscillation with frequency ω_3 , all three gap voltages are in phase. Our tube is designed so as to produce the mode of oscillation having the frequency ω_1 .

The shape of the cavities is such that the outer ones are equal to half the centre one. This has already been shown schematically in fig. 4. As a result, $a \approx 0.5$, which experience has shown to be a suitable value.

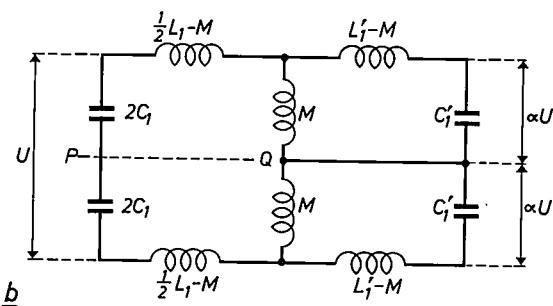
Tight coupling is obtained by making the holes in the walls very large. In fact, as can be seen from fig. 7b, the walls consist of strips. The strips are as narrow as adequate heat conduction allows. A large part of the power is namely dissipated in the cylinders 2.

Because of the fact that the three cavities are tightly coupled, so that we may regard them more or less as a single resonant cavity, the coupling of the centre cavity to the fourth one does not appreciably alter for the factor a . Moreover the frequency ω_3 is sufficiently far outside the tuning range of the tube.

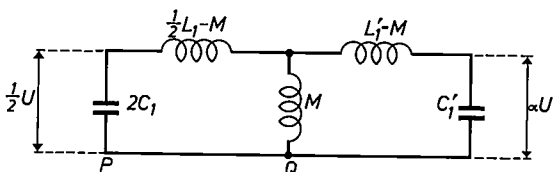
The compromise between the requirements of tight coupling and adequate thermal conduction in the strips led to values of 8900 and 11 400 Mc/s for the respective resonant frequencies ω_1 and ω_3 , and to 10 000 Mc/s for the frequency ω_2 .



a



b



c

Fig. 8. a) Equivalent circuit of the internal part of the resonator. There are three inductively coupled LC circuits, the first and third being identical, corresponding to the cavities C_1 , C_2 and C_3 . The quantities relating to the outer circuits are marked with an accent. M is the mutual inductance.

b) Different form of diagram (a). For reasons of symmetry, points P and Q always have the same potential and may therefore be joined; the circuit can thus be halved, producing the diagram in (c).

Since in the frequency bands around ω_1 and ω_3 we can treat the three-cavity system by approximation as a single resonant cavity, we can quickly find the characteristics of the whole resonator by treating it as a combination of this hypothetical single cavity with the fourth cavity. This combination again has the properties of a system consisting of two inductively coupled oscillatory circuits. If the frequency ω_0 of the fourth cavity is in the region of ω_1 , then instead of this one resonant frequency, there are two, which we shall call ω_{11} and ω_{12} . Likewise, when ω_0 is close to ω_3 , instead of ω_3 there are two resonant frequencies ω_{31} and ω_{32} . (The frequency ω_2 is not disturbed or split, since the appertaining oscillation cannot of course excite the fourth cavity.) In connection with the dimensions of the internal resonator system, it seemed to us to be best to use the higher of the two frequencies ω_{11} and ω_{12} , which we will from now on refer to as ω_{12} .

When the coupling coefficient k between the outside cavity and the inner one is small in relation to unity, for the frequencies of the whole resonator hold:

$$\frac{\omega_{11} - \omega_1}{\omega_1} = \frac{1}{2}[\gamma - \sqrt{\gamma^2 + k^2}] \quad (10a)$$

and

$$\frac{\omega_{12} - \omega_1}{\omega_1} = \frac{1}{2}[\gamma + \sqrt{\gamma^2 + k^2}] \quad (10b)$$

where γ is the relative difference $(\omega_0 - \omega_1)/\omega_1$ between the resonant frequencies of the two circuits⁶⁾. Dividing these expressions by k , we see that

$$(\omega_{11} - \omega_1)/k\omega_1 \quad \text{and} \quad (\omega_{12} - \omega_1)/k\omega_1$$

are functions of γ/k only. A plot of these functions is given in fig. 9. This figure illustrates the effect, well known from the theory of oscillations, that of the resonant frequencies of the coupled resonant cavities (resonant circuits), one is higher than the highest of the frequencies of the non-coupled cavities, and the other lower than the lowest. The lowest value that ω_{12} can theoretically have is thus the value ω_1 , in other words the latter frequency determines the lower limit of the tuning range of the tube.

The upper limit of the tuning range is determined by the losses in the fourth cavity. The higher is ω_0 the smaller is the difference between ω_{12} and ω_0 —

⁶⁾ Expressions (10a) and (10b) can be derived directly from the formula for the two resonant frequencies of a system consisting of two coupled oscillators. This formula can be found in most books dealing with the theory of electrical or mechanical vibrations.

in fig. 9 curve 1 approaches closer to the dashed line as γ/k increases — and the fourth cavity oscillates more strongly. The losses P_i in this cavity thus increase with rising ω_0 , and finally constitute a large part of the power P_i delivered by the inner cavities.

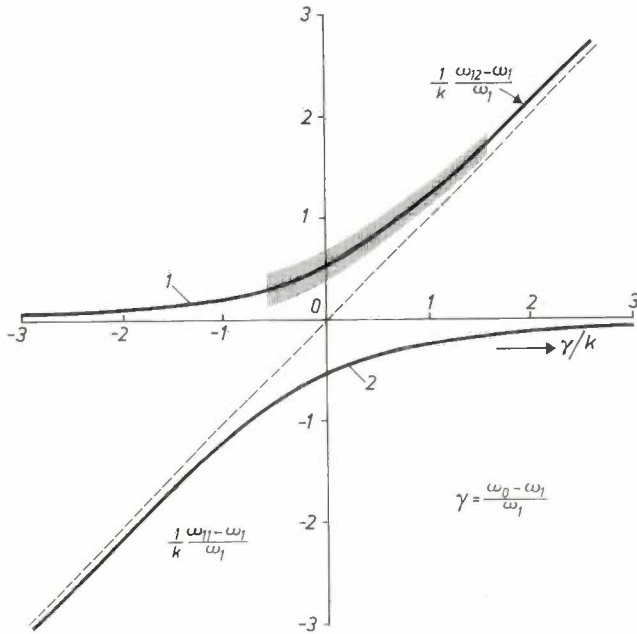


Fig. 9. Explanation of the range within which the tuning can be varied. The tuned frequency ω_{12} is always higher than the resonant frequency ω_1 of the three inner cavities (regarded as a whole) and higher than the resonant frequency ω_0 of the fourth cavity. Where γ is very strongly negative, ω_{12} approaches ω_1 , and where γ is highly positive — in which case $\omega_0 > \omega_1$ — it approaches ω_0 . In the latter case considerable losses occur in the fourth cavity. These govern the upper limit of the practical tuning range (shaded area). The lower limit is theoretically ω_1 . In practice this limit is determined by the lowest value which ω_0 can have.

The output power $P_i - P_l$ is then small, and is smaller the higher is ω_0 . In the new tube it proved possible, given a k value of 0.06, to vary γ/k between -0.6 and $+1.6$ (see fig. 9). This corresponds to a tuning range of 9100 to 9800 Mc/s, a frequency range easily containing the civil radar band (9320-9500 Mc/s). The principal relevant data of the tube are listed in Table I.

Table I. Principal data relating to the 3 cm multi-reflex klystron for use in remark beacons. All voltages are given relative to the cathode.

Voltage on resonator	+ 3000 V
Current to resonator	20 mA
Voltage on open reflector electrode	+ 650 V
Voltage on solid reflector electrode	- 750 V
Mechanical tuning range	9100-9800 Mc/s
Output power at 9500 Mc/s	9-12 W
Output power at the limits of the tuning range	> 7 W

Construction

In describing the construction of the new tube we shall refer to fig. 10, which gives a schematic cross-section. In a way that will be shown presently, the tube is built up from four prefabricated units: the resonator, the electron gun, the reflector and the bulb.

The three resonant cavities C_1, C_2 and C_3 , which form the internal part of the resonator, are made

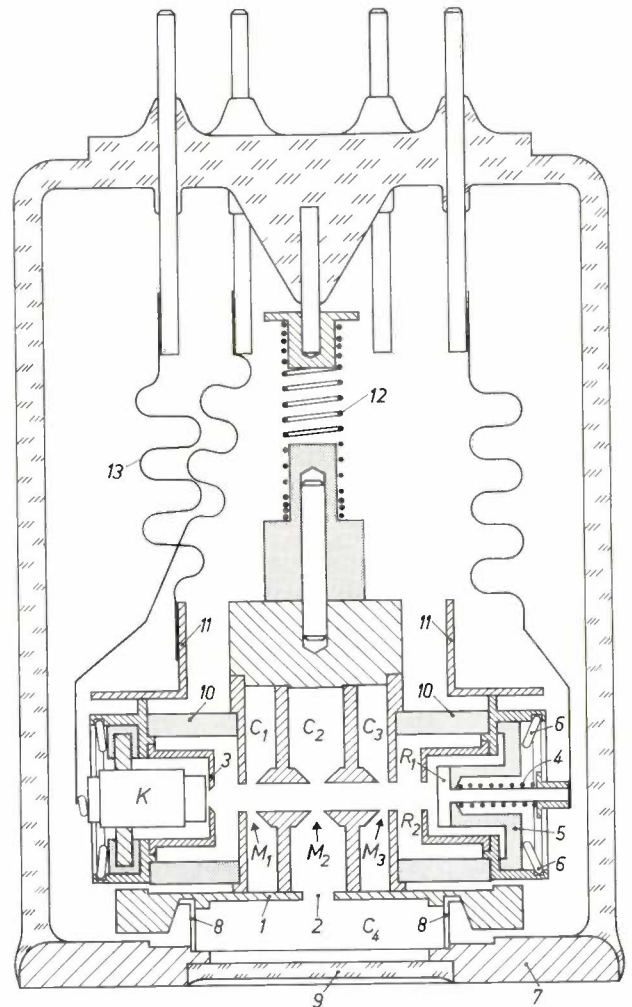


Fig. 10. Schematic cross-section of the tube (simplified). In the centre of the lower part can be seen the three cavities C_1, C_2 and C_3 of the resonator with the three interaction gaps M_1, M_2 and M_3 . The base 1 of the resonator contains a slot 2 which couples C_2 to the fourth cavity C_4 . Left of the resonator are the cathode K and the focusing electrode 3. On the right, the reflector with the solid electrode R_1 and the open electrode R_2 . 4 spring which forces R_1 against the insulator 5. (All shaded parts represent cross-sections of insulators.) 6 spring which retains 5 in position (there is a similar spring in the cathode unit). 7 molybdenum end plate with copper collar 8 and glass window 9 which has low losses in the operating frequency range. 10 spacer rings. 11 flanges fixed to cathode and reflector and between which, apart from the resonator (in the figure above and below the plane of the drawing; see fig. 11), tie rods pass pressing these units against the spacers 10. 12 spring which presses the base 1 of the resonator against the collar 8 (after 8 and 1 have been soldered together, 12 ceases to have any further use). 13 molybdenum strips connecting the electrodes to the base pins (the strips for R_2 and the resonator are not shown).

entirely of oxygen-free copper. The base plate 1, which contains the coupling slot 2, forms at the same time the bottom of the fourth cavity C_4 . As mentioned, the latter is mainly outside the tube. The vacuum-tight seal 9 of the bulb passes through C_4 , but has no essential influence on the electromagnetic field.

The electron gun (to the left of the resonator, in fig. 10) consists of a cathode K and an open electrode 3 which entirely encloses K and serves for adjusting the beam current.

The reflector unit (right of the resonator), consisting of the solid electrode R_1 and the open electrode R_2 , shows a close geometric resemblance to the electron gun. R_1 is pressed by a spring 4 against the insulator 5, which is clamped by a ring 6 in the bushing supporting R_1 .

The bulb consists of a base and a cylindrical wall of glass and an end plate 7 of molybdenum. This metal was chosen because it is non-magnetic and can be sealed to the glass. The end plate contains a round central hole, around which, on the inside of the bulb, a thin copper collar 8 is soldered that fits into a recess in the base plate 1 of the resonator. The hole is made vacuum-tight by the above-mentioned window 9, which is made of a type of glass that has low losses in the tube's frequency range.

When assembling the tube, the first step is to fix the electron gun and the reflector unit to the resonator. They are kept in alignment by insulating spacer rings 10. The whole assembly is held together by two tie rods running parallel with the axis of the electrodes, each rod, at the height of this centre line, passing through a hole in the resonator wall and through holes in the flanges 11. The ends of the rods are provided with springs which press the flanges towards the resonator (cf. fig. 11).

After this the electrodes are connected to the base pins by slack molybdenum strips 13 and the whole assembly is introduced into the bulb. Spring 12 presses the base plate of the resonator against the copper collar 8. Between the plate 1 and the collar 8, a ring of low-melting-point solder has previously been laid, and this melts during the heating of the tube when degassing. This produces between the resonant cavity and the molybdenum disc a connection which is a good conductor both of electric current and heat.

Consequently not only does the fourth cavity have a high quality factor but the power which the beam dissipates on the resonator and the RF power dissipated in the resonator are conducted to the molybdenum plate via a low-thermal-resistance path. This plate is cooled by conduction towards the

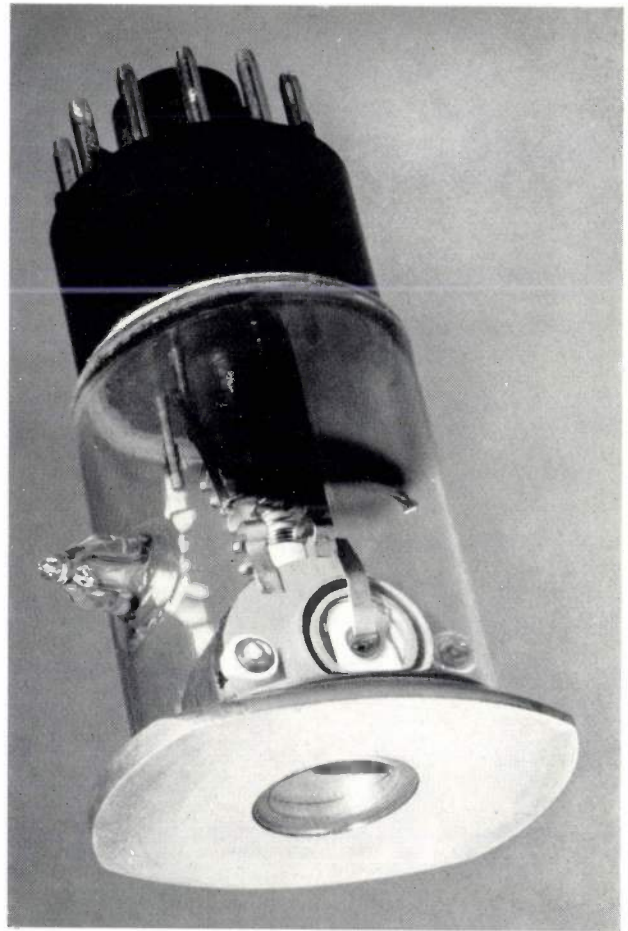


Fig. 11. The 3 cm multi-reflex klystron for microwave beacons. In the lower part, the molybdenum end plate with window for connection to waveguide. Above this, in the centre, can be seen the molybdenum strip connecting the cathode to one of the base pins, and one of the flange plates 11 of fig. 10, showing the ends of the tie rods which press the cathode and reflector units to the resonator.

outer part of the fourth cavity, which is pressed against this plate. Thus the temperature of the resonator remains low, which benefits the life of the tube as well as the frequency stability and efficiency. (At higher temperatures the resistance of the copper is greater and consequently the losses higher.)

Output system

For the purpose of tuning the tube (cavity C_4) and matching it to an external load, the tube must be connected to a small section of rectangular waveguide which is divided into two spaces by a cross-partition containing an opening. One of the spaces constitutes the major part of C_4 . Fig. 12 shows schematically two mutually perpendicular cross-sections of this output system, and also indicates how the tube is connected.

The cavity C_4 is made tunable by introducing into the waveguide a copper tuning disc T which, by

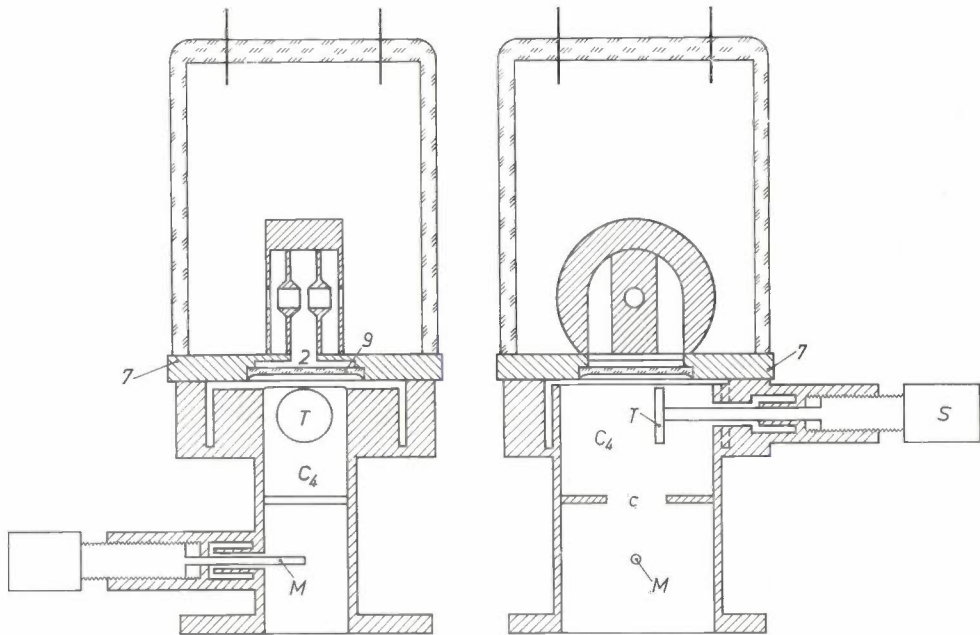


Fig. 12. Schematic cross-section of the tube in two mutually perpendicular directions, with a section of waveguide connected for tuning and matching. Space C_4 forms the external part of the fourth cavity (see fig. 7). This space is coupled by means of the gap c to the rest of the output system. T tuning disc, adjustable with micrometer screw S . For 50 c/s modulation, T is set in vibration by a kind of loudspeaker system (not shown). M matching stub, also on micrometer screw. The figures have the same meaning as in fig. 10.

means of a screw S , can be moved from the narrow wall to the middle of the cavity. When the tuning disc is close to the wall it has little influence on the field in the cavity; when it is moved towards the middle the equivalent capacitance of the cavity is increased and its resonant frequency thus decreased.

For the purpose of 50 c/s frequency modulation, as required for the transmitter of a ramark beacon,

the rod to which disc T is attached is connected with a device that more or less corresponds to the drive system of a loudspeaker. This device causes the tuning disc in the resonant cavity to vibrate axially at the required modulation frequency. To obtain the frequency swing of 100 Mc/s necessary for the beacon, the amplitude of the vibration should be about 0.6 mm. The "loudspeaker system" with the

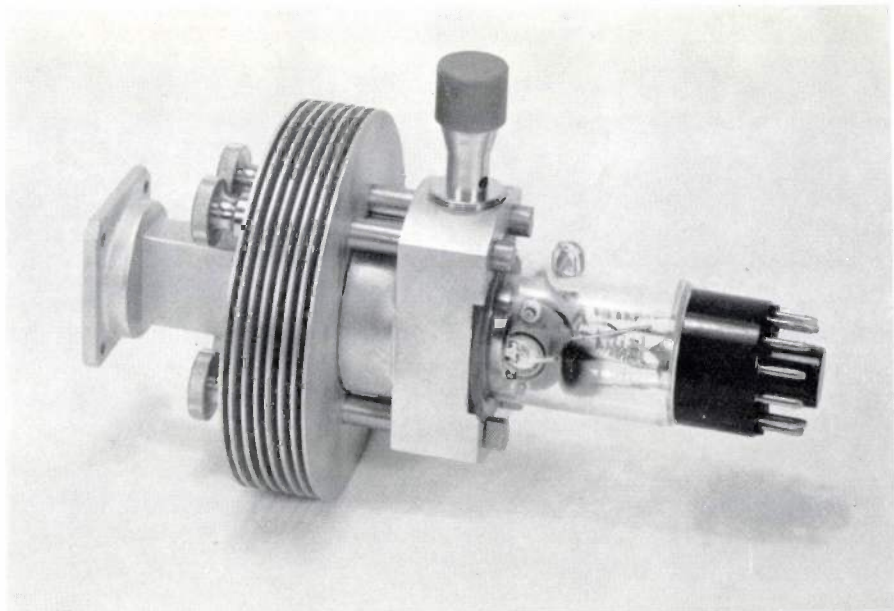


Fig. 13. The klystron coupled to the external tuning waveguide; the cross-section of the waveguide is shown in fig. 12.

attached tuning disc can in turn be moved axially by a tuning screw for adjusting the central frequency.

The other space of the output waveguide contains a matching stub *M*. By varying the depth of this stub in the output waveguide, the impedance of the load can be adjusted for maximum output power.

The output waveguide is mechanically coupled to the tube by means of a normal choke flange, against which presses the molybdenum end plate of the tube. As we have seen, this enables the heat generated in the tube to be removed via the waveguide. *Fig. 13* shows a photograph of the tube, coupled to the output waveguide.

Summary. In remark microwave beacons for sea-going shipping, there is a need for transmitting tubes capable of delivering, at 3 cm wavelength, a power of about 10 W while allowing a frequency swing through the whole civil radar band (9320 to 9500 Mc/s) with a modulation frequency of 50 c/s, and in addition a 10^4 e/s modulation with a swing of 10 to 15 Mc/s. The article describes a new multi-reflex klystron of high efficiency (15 to 20%) which meets these requirements and which is mechanically tunable (over 8%). Bunching is effected by means of a resonator consisting of three coupled resonant cavities, each with its own interaction gap. The required field with a parabolic-

ally varying potential is approximated by a drift space, containing the resonator, flanked by two retarding spaces in which the potential varies linearly. The alternating-current component of the electron beam is calculated, and using this the optimum dimensions of the tube are obtained. The tuning and the 50 c/s modulation are effected by means of a fourth cavity which is coupled to the three others and is largely formed by part of the output system. The tube has a glass envelope with an end plate of molybdenum containing a low-loss glass window. The tube is cooled by conduction via the molybdenum plate to the output waveguide.

ABSTRACTS OF RECENT SCIENTIFIC PUBLICATIONS BY THE STAFF OF N.V. PHILIPS' GLOEILAMPENFABRIEKEN

Reprints of those papers not marked with an asterisk * can be obtained free of charge upon application to the Philips Research Laboratories, Eindhoven, Netherlands, where a limited number of reprints are available for distribution.

3007: H. Draaijer and H. J. Pel: Der Einfluss von Dolomit auf den Säuretaupunkt und auf die Niedertemperaturkorrosion in Ölkesseln (Brennstoff-Wärme-Kraft **13**, 266-269, 1961, No. 6). (The influence of dolomite on the acid condensation point and on the low-temperature corrosion in oil-fired boilers; in German.)

Oil used for heating purposes always contains sulphur. When this oil is burnt in e.g. a boiler, part of the sulphur is oxidized to sulphuric acid, which under certain conditions can condense on the pipes of the economizer and cause serious corrosion of these pipes. Dolomite is sometimes used to neutralize the acid, but it only does so partially, and moreover it forms an undesired deposit on the economizer. The authors have found that the corrosion is highest when the wall temperature of the pipes is 30-35 °C below the condensation point of the acid. The corrosion can be reduced by a factor of 10 by raising the temperature of the walls to above the condensation point.

3008: M. J. Sparnaay: Note on the additivity of retarded London-Van der Waals forces (Original Lectures IIIrd int. congress of surface activity, Cologne 1960, Vol. I, pp. 141-144, publ. Universitätsdruckerei Mainz).

The theory of Van der Waals forces includes a formula for estimating the attractive force between two neutral atoms. If one wishes to use this formula to calculate the total Van der Waals force exerted on one atom by a number of other atoms, one may not assume these forces to be additive. The author calculates the force between two pairs of atoms, using Hamiltonians and a coordinate transformation suggested by London. The result of this calculation can be used to calculate the attractive force between two solid bodies (see also **3009**).

3009: M. J. Sparnaay and P. W. J. Jochems: Measurements of attractive forces between solid bodies (as **3008**, Vol. II, pp. 375-377).

One of the authors has previously published a number of papers on the theory and measurement of Van der Waals forces between solid bodies (see e.g. **2633** and **2900**, where the attractive force between flat and slightly spherical plates is measured, and **3008**). In order to test the theory in different ways, the Van der Waals force has been remeasured by the method described in this publication. A ball of radius about 40 μm is formed on the end of a wire of diameter 20 μm by melting the latter. A plate is slowly moved towards this ball. At a certain distance the Van der Waals force overcomes the elastic forces in the wire, and the ball is pulled against

the plate. This distance, which is of the order of a micron, can be measured with a microscope and is a measure of the Van der Waals force. Because the method is not very sensitive, the measurement must be repeated many times with a large number of balls of various radii.

3010: H. P. Rooksby and C. J. M. Rooymans: The formation and structure of delta alumina (Clay Minerals Bull. 4, 234-238, 1961, No. 25).

Al_2O_3 can exist in several metastable structures as well as in the stable (corundum) structure. The structure of the δ modification, which is very difficult to obtain pure, is investigated here. The authors mention a number of methods of making $\delta\text{-Al}_2\text{O}_3$ in crystalline form. By X-ray investigation it has been found that the structure can be described by a tetragonal unit cell with parameters $a_0 = 7.96 \text{ \AA}$ and $c_0 = 11.70 \text{ \AA}$. By analogy with the occupation of the lattice sites in the spinel lattice, it appears that this unit cell should contain 32 Al vacancies. This is in disagreement with an investigation by Saalfeld, who suggested a unit cell with $21\frac{1}{3}$ vacancies.

3011: H. Risken and H. J. G. Meyer: Contribution of lattice scattering between nonequivalent valleys to free-carrier infrared absorption in semiconductors (Phys. Rev. 123, 416-418, 1961, No. 2).

The conduction band of germanium has a multi-valley structure. The scattering of electrons between non-equivalent valleys plays a role in infrared absorption and hot-electron effects. Examples of such scattering are transitions of electrons from the (111) valley to the (100) valley or the [000] valley. This theoretical article gives a quantum-mechanical calculation of the infrared absorption due to the above-mentioned scattering.

3012: G. Meijer and R. van der Veen: Dual effect of red light on the photoperiodic response of *Salvia occidentalis* (Progress in photobiology, Proc. 3rd int. congress on photobiology, Copenhagen 1960, editors B. C. Christensen and B. Buchmann, pp. 387-388, Elsevier, Amsterdam 1961).

The short-day plant *Salvia occidentalis* (see 2711) can be made to behave as if growing under long-day conditions (16 hour light period) if the night period of 14 hours is interrupted by a brief irradiation with red light (night break). This however only happens if enough blue or infrared light was administered during the day period. The long-day

effect of a night break can be reduced or eliminated by giving red light earlier in the night period.

3013: R. van der Veen and G. Meijer: Critical day-length of the short-day plant *Salvia occidentalis* in red and far-red radiation (as 3012, pp. 389-390).

The short-day plant *Salvia occidentalis* only flowers if the daily period during which it receives light is less than 14 hours (the critical daylength). This period can be reduced to 11 hours by irradiation with red and infrared light (see also 3012, 2855, 2711 and 2678).

3014: C. H. Weijnsfeld, A. Hoogendoorn and M. Koedam: Sputtering of polycrystalline metals by inert gas ions of low energy (100-1000 eV) (Physica 27, 763-764, 1961, No. 8).

Continuation of the investigation of cathode sputtering caused by rare-gas ions; previous results have been published in 2629, 2767, 2840, R 416 and R 425 (thesis, M. Koedam).

3015: C. Haas and M. M. G. Corbey: Measurement and analysis of the infrared reflection spectrum of semiconducting SnS (Phys. Chem. Solids 20, 197-203, 1961, No. 3/4).

The infrared reflection spectrum of *P*-type SnS single crystals has been measured as a function of the wavelength between 2 and 25 μm . The experimental data can be explained in terms of the interaction of light quanta with lattice vibrations and with charge carriers (holes). The refractive index ($n_0 = 3.6 \pm 0.1$), the dielectric constant ($\epsilon = 19.5 \pm 2$) and the effective atomic charge ($e^* = 0.7 e_0$, where e_0 is the charge on the electron) are calculated from the experimental results. It is found that the effective mass of the holes is $m^*_\perp = 0.2 m_0$ ($m_0 =$ mass of electron) for motion perpendicular to the *c* axis, but much greater ($m^*_\parallel \approx m_0$) for motion parallel to the *c* axis.

3016: W. J. Oosterkamp and Th. G. Schut: Image intensification of X-rays in angio-cardiography (Proc. 3rd int. Conf. on medical electronics, London 1960, pp. 487-489, Instn. Electr. Engrs., 1961).

The dose in diagnostic radiology can be reduced if an image intensifier is used, while the increased brightness of the image has made cine-radiography a routine procedure. This publication describes a method having the above-mentioned advantages, used in connection with angio-cardiography and catheterization of the heart. Use of an image inten-

sifier also opens the possibility of X-ray television, if necessary with recording of the televised images on a magnetic image memorizer. (See also Philips tech. Rev. **22**, 1-10, 1960/61.)

3017: O. Bosgra and J. H. G. Roerink: Praktijkproef met een levend avirulent pseudo-vogelpest- (stam B₁) en infectieuze-bronchitis-drinkwatervaccin (T. Diergeneesk. **86**, 1198-1209, 1961, No. 18). (Field test of a living non-virulent drinking-water vaccine against Newcastle disease and infectious bronchitis; in Dutch.)

Before a vaccine for poultry can be used in the Netherlands, it must be given a field test in accordance with the regulations of the Veterinary Service. The vaccines developed by Philips-Duphar against infectious bronchitis and Newcastle disease, which can be administered via the drinking water, have been tested in the prescribed manner on a large number of chickens. 97.5% of the chickens which had been vaccinated with the bronchitis vaccine proved to be immune, while 88% of the controls contracted the disease after infection. The corresponding figures for the Newcastle-disease vaccine are 93 and 95%. The vaccines thus more than comply with the demands of the Veterinary Service.

3018: A. H. Boerdijk: Zero-, first-, and second-order theories of a general thermocouple (J. appl. Phys. **32**, 1584-1589, 1961, No. 8).

The thermocouple considered in this publication consists of two bars of arbitrary form. Each of the properties of the materials (the thermal resistivity α , the electrical resistivity ρ and the Seebeck coefficient S) is represented by a finite number of terms of a Taylor series in T (the temperature) and u (a positional coordinate). A method for deriving a "theory of arbitrary order t " is given, based on a function $T = f(u)$ which satisfies the fundamental non-linear differential equation (obtained by use of the thermodynamics of irreversible processes) and the boundary conditions, neglecting all terms of order greater than t . The order of a term is equal to the sum of the orders of all partial differential quotients of α , ρ and S with respect to T and u that occur in the term.

The method is used to express the electrical output power and the thermal output powers as functions of the current and the temperatures of the junctions, for theories of the zero, first and second orders. The zero-order theory is the same as the normal theory of thermocouples with constant α ,

ρ and S . In the first-order theory, an expression is derived for the efficiency for the production of cold. This efficiency can be increased by making S vary in a suitable way with temperature and position. Finally, the accuracy of the approximations made is discussed.

3019: W. Hondius Boldingh: Quality and choice of Potter Bucky grids, VI. Exposure data for various grids (Acta radiol. **56**, 202-208, 1961, No. 3).

Continued from 2793 and 2990. This publication discusses the relationship between the exposure time and the contrast of medical X-ray photographs when various types of Potter-Bucky grids are used.

3020: G. Diemer: Field effects on photoconductivity quenching (Physica **27**, 979-981, 1961, No. 10).

In this letter to the editor, it is pointed out that the observations made by Kitamura, Kubo and Yamashita on the combined influence of the electrical field and irradiation with infrared light on the photo-sensitivity of CdS crystals can be explained with the aid of generally accepted mechanisms, applied to the model of Klasens and Schön.

3021: O. W. Memelink: De werking en eigenschappen van het vastestofthyatron (T. Ned. Radiogenootschap **26**, 119-126, 1961, No. 3). (The operation and properties of the solid-state thyatron; in Dutch.)

A solid-state thyatron (also known as a pylistor, silicon controlled rectifier or PNP switch) consists of four layers of silicon, alternately P -type and N -type. The voltage-current characteristic shows a close resemblance to that of a normal gas-filled thyatron. The solid-state thyatron described here can block 500 V in both directions. After a description of its operation, some applications as an electronic switch are mentioned. See also Philips tech. Rev. **23**, 272-278, 1961/62 (No. 8/9).

3022: H. Groendijk: Three interpretations of space-charge waves in electron beams (T. Ned. Radiogenootschap **26**, 51-64, 1961, No. 2).

In many types of amplifier tubes for microwaves, such as the klystron and the travelling-wave tube, use is made of an electron beam in which the velocity is modulated by the signal to be amplified. By use of Maxwell's equations and the equations of motion of the electrons, it can be shown that space-charge waves occur in the beam. This article, which discusses a klystron as an example of this class of tubes,

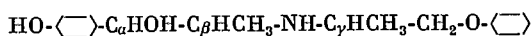
shows how it is possible to acquire an insight into the production and mechanism of the above-mentioned waves by a qualitative consideration of the motion of the individual electrons and the forces which they exert on each other. It appears that the motion of the electrons in a klystron can be described in three different ways. Use is made of these three interpretations in turn to explain the power amplification of a klystron, to show how the noise can be reduced in such tubes, and to show that similar considerations can also be applied to electron-beam tubes in which the signal is transmitted and amplified by other means than velocity modulation.

3023: H. A. Sloom: Luchtsterilisatie-methoden (Pharmaceut. Weekblad **96**, 703-712, 1961, No. 18). (Methods of sterilizing air; in Dutch.)

The author describes and compares the four most important methods of sterilizing air: filtration, electrostatic precipitation, irradiation with ultraviolet light and the evaporation of disinfectants.

3024: J. van Dijk and H. D. Moed: Synthesis of β -phenylethylamine derivatives, VII. The enantiomers of erythro-1-(4'-hydroxyphenyl)-2-(1''-methyl-2''-phenoxyethylamino)-propanol-1 (Rec. Trav. chim. Pays-Bas **80**, 573-587, 1961, No. 7).

The vasodilator (substance which causes blood vessels to expand) "Duvadilan" (Caa 40), whose formula is



is a racemic compound consisting of two enantiomers (optical antipodes, denoted by (+) and (-)). The authors have succeeded in separating this compound with the aid of the enantiomers (+)- and (-)-mandelic acid, and have also elucidated the absolute configuration of the two components of the racemate: (+) = $\alpha\text{R} : \beta\text{S} : \gamma\text{S}$ and (-) = $\alpha\text{S} : \beta\text{R} : \gamma\text{R}$. See also Philips tech. Rev. **24**, 69-79, 1962/63 (No. 3), and 2709.

3025: W. J. A. Goossens and H. J. G. Meyer: Enkele basisbegrippen uit de fysica van halfgeleiders (Ned. T. Natuurk. **27**, 324-344, 1961, No. 9). (Some basic concepts of semiconductor physics; in Dutch.)

The authors give a comprehensible fundamental treatment of the basic concepts of semiconductor physics, which are too often skimmed over. Starting

from Schrödinger's equation, they derive an expression for the energy bands of a perfect crystal lattice. Also mentioned are the difference between conductors, insulators and semiconductors, the band structure of semiconductors, and the velocity, acceleration and effective mass of the electrons. Finally, the influence of lattice imperfections on the band structure is described.

3026: H. G. van Bueren and J. Hornstra: Grain boundaries and the sintering mechanism (Reactivity of solids, Proc. 4th int. Symp. on the reactivity of solids, Amsterdam 1960, editors J. H. de Boer *et al.*, pp. 112-121, Elsevier, Amsterdam 1961).

When sintering metals, it is often desirable to eliminate pores from the material as far as possible. Study of the sintering process has shown that grain boundaries absorb the vacancies which diffuse away from the pores, so that pores will disappear fastest in a fine-grain material. Further consideration shows that the absorption of vacancies by the grain boundaries is coupled with plastic flow (motion of the grains relative to each other). This serves to remove apparent contradictions in the literature on this point. See also 2966.

3027: J. Hornstra: Dislocations, stacking faults and twins in the spinel structure (as 3026, pp. 563-570).

Although dislocations have long played a considerable part in the theory of the crystal lattice, the structure of the dislocations themselves has not been studied much, for lack of experimental methods for so doing. There are however several reasons why such a detailed knowledge of dislocations would be highly desirable. Starting from a model of the spinel lattice and considering only the electrostatic forces, one can predict certain properties of a dislocation, e.g. the glide plane chosen. This choice is determined by the number of partial dislocations into which one dislocation may be split up and the condition of electrical neutrality, among other things. See also 2923.

3028: F. K. Lotgering: Conservation of crystal orientation of hexagonal iron-oxide compounds during solid-state reactions (topotaxy) (as 3026, pp. 584-586).

A contribution to a symposium, dealing with the same subject as already treated in 2938.

Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS

RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES

SHOCK TESTING OF INCANDESCENT LAMPS

by E. W. van HEUVEN *).

621.326.001.45

For use in trams, trains, ships, factories, etc., incandescent lamps are made that are specially designed to withstand shocks and vibrations. The present article is not primarily concerned with these lamps themselves, but with the problem of testing in a sufficiently reproducible way their ability to stand up to shocks.

Effects of shocks and vibrations on an incandescent lamp

Tungsten, the metal from which the filaments of all conventional incandescent lamps are made, is primarily used for this purpose because of its very high melting point. It has two properties, however, that are not so favourable: its density is very high, and its yield point — at the working temperature — is extremely low. These two properties combined make it necessary, when designing a filament, to make careful allowance for the risk of sagging. This applies especially to coiled filaments, which are used either to minimize heat-conduction losses (an important consideration in all gas-filled lamps) or simply to facilitate mechanization of the production process (in the case of vacuum lamps).

Sagging is counteracted on the one hand by "doping" the tungsten wire with special additives to give it a suitable crystalline structure ¹⁾, and on the other by providing the filament with supports at appropriate points. The local cooling caused by these supports is obviously detrimental to the luminous efficiency of the lamp. For this reason, standard incandescent lamps contain only as many support hooks as are necessary to keep the internal stresses, produced in the wire by the force of gravity, just

below the permissible limit. The number of support hooks used also depends on parameters such as the thickness of the wire and the diameter and pitch of the spiral, and all these parameters are chosen with a view to obtaining the highest possible luminous efficiency. That the mechanical stress limit is approached fairly closely in practice can be seen from the slight festoon-like sagging of the filament in lamps that have burned throughout their life in the same position.

When a tungsten lamp is exposed to shocks or vibrations, forces of inertia act on the filament that can be many tens of times greater than the weight of the wire. The said mechanical stress limit may then be exceeded and the filament breaks. But even if it does not break and shows only local deformation (localized stretching of the spiral), the lamp still suffers damage. For a part of the filament that is stretched (*fig. 1*) falls in temperature; in vacuum lamps because that part receives less radiation from other parts of the filament and in gas-filled lamps because per unit length of wire it loses more heat to the gas ²⁾. The local cooling causes a drop in the resistance of the deformed part, resulting in a somewhat higher current. This increases the temperature of the rest of the spiral, so that the tungsten there evaporates faster and shortens the life of the filament.

*) Philips Lighting Division, Eindhoven.

¹⁾ See, for example, J. L. Meijering and G. D. Rieck, The function of additives in tungsten for filaments, Philips tech. Rev. 19, 109-117, 1957/58.

²⁾ See W. Geiss, Improvements in the efficiency of electric incandescent lamps, Philips tech. Rev. 6, 334-342, 1941.

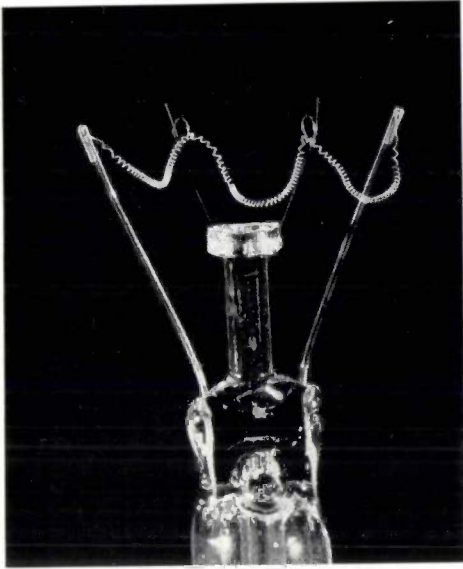


Fig. 1. Filament that has been subjected to shocks, as a result of which parts of the coil are deformed. At the places where the filament has been stretched, the temperature is lower, which causes a drop in resistance and hence a higher current, shortening the life of the filament.

In the case of vibrations, lamp life is adversely affected not only by the varying inertial forces — the maximum amplitude of which determines the damage caused — but also by resonance phenomena and the chattering of the filament in the support hooks.

For use in places where there are constant shocks and vibrations, e.g. in trams, trains, ships and factories, incandescent lamps have therefore been designed with a reinforced construction. The most obvious means of strengthening the filament assembly is to fit more support hooks. Since, however, this inevitably entails some loss of luminous efficiency, one can better revise the whole design of the tungsten-filament lamp and try to give it the required extra strength in some other way. In the first place this implies doing without the gas filling, which was introduced for standard incandescent lamps, at least for wattages of 25 W and above, because it allows the filament to operate at higher temperatures, thus improving the luminous efficiency. Returning to the

vacuum lamp makes it possible to choose a thicker filament wire, having increased mechanical strength. Next, the diameter of the spiral can be reduced and/or the pitch increased. For wattages of 75 W and above, however, the vacuum lamp is no longer a reasonable proposition. A gas filling thus again becomes necessary, but the *coiled-coil* filament, used in standard gas-filled incandescent lamps to improve light output, can now as a rule better be dispensed with as it weakens the filament. For comparison, *fig. 2* shows the filament body of a standard incandescent lamp side by side with that of a lamp of the same wattage with reinforced construction.

Whether the measures adopted to strengthen the construction of a given type of lamp (or rather of its filament assembly) produce the desired result can only be established with certainty after the lamp has proved itself in actual operation. There are times, however — for example where the manufacturer or the consumer has to decide whether a certain type of lamp can be used or a new type has to be developed — when information on the mechanical strength of the type of lamp is required at short notice. To this end numerous test apparatuses have been devised, all of which are based on the same principle of measuring the life of the lamps while they are burning in a frame which is continuously subjected to shocks and vibrations. In view of the widely differing conditions in which lamps are used in practice, mechanical testing of this kind, using an intricate (often rather poorly defined) pattern of forces, is representative only of special cases. Another disadvantage of

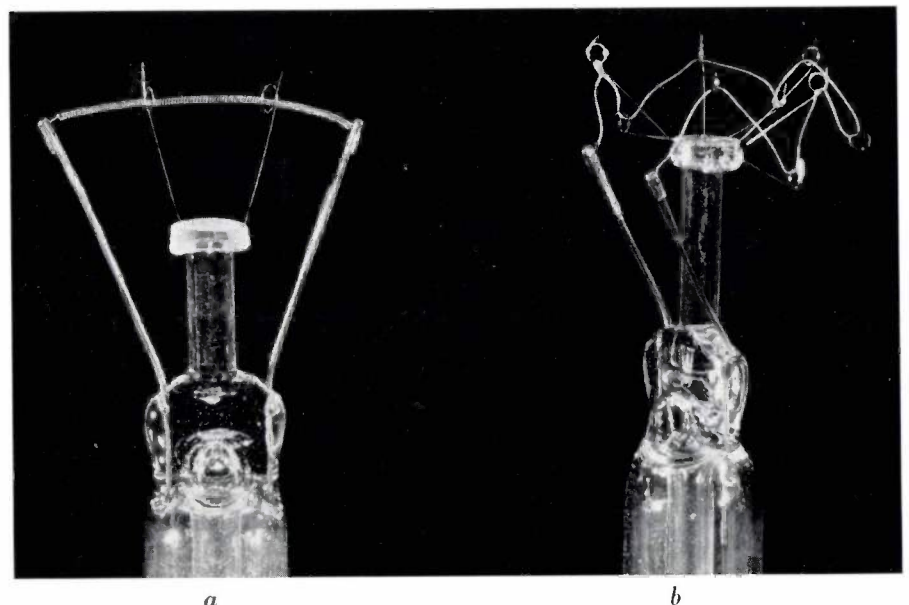


Fig. 2. The filament body of a standard incandescent lamp (a) and of an incandescent lamp with strengthened construction (b), both for the same wattage. In (a) a coiled-coil filament is used, in (b) a single coil. Note the larger number of support hooks in the second case.

existing apparatus is that different examples, even when built entirely to the same specifications, by no means subject the lamps to the same pattern of forces. The results obtained on different models of the same test apparatus cannot therefore be directly compared with one another. Finally, working on the principle mentioned, it is difficult to ensure that the apparatus will retain the same properties as time goes by, for in the long run the apparatus it-

type of lamp to shocks and vibrations of various frequencies. We shall confine ourselves here to one aspect only, that of *shock testing*, and describe the apparatus constructed for this purpose.

In developing this apparatus we devised a means of accurately adjusting the violence of the shocks, making it possible to perform the test with well defined forces, and so to compare results obtained from different examples of the apparatus. Provision was

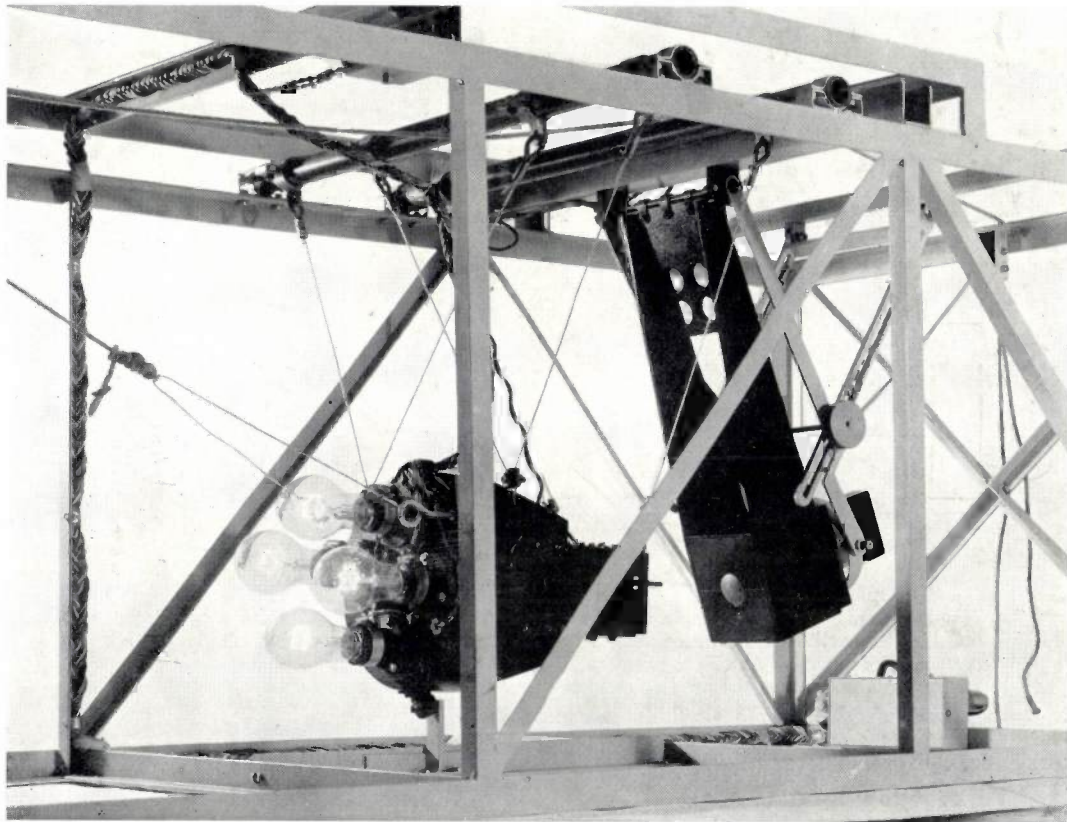


Fig. 3. The apparatus for testing the resistance of tungsten-filament lamps to shocks, built in the strength-testing laboratory of Philips Lighting Division.

self is bound to suffer from the magnitude of the forces that have to be exerted on strengthened lamps to produce a noticeable effect on their life. This means that results obtained with the same apparatus but at different times are not directly comparable either.

A new test apparatus

With the aim of avoiding the drawbacks described, we embarked on the development of a series of new test apparatuses. The forces occurring in practice are nearly always produced by a composition of shocks and sinusoidal vibrations, but we started from the assumption that the mechanical strength of a tungsten filament can be tested by determining its ability to withstand each of the components *separately*. A complete test, then, consists of subjecting a given

also made for varying the violence of the shocks within wide limits, so that the same apparatus could be used for testing many types of lamps.

Fig. 3 shows a photograph of one of the first models built, and fig. 4 gives the set-up schematically. It consists principally of an anvil *A* suspended by four steel wires *D*, and a hammer *H* that can swing about an axis *O*. When both hammer and anvil are suspended in equilibrium they are only just in contact with each other, at point *P*. The apparatus works as follows. Four lamps are fitted in the anvil and supplied with the appropriate working voltage through flexible leads. The hammer is swung upwards through an adjustable angle φ and held in this position by an electromagnet *E*. When the solenoid current is cut off, the hammer is released with zero initial velocity. The shock of the collision

between hammer and anvil imparts to each of the lamps a momentary acceleration; the resultant inertial force acting on the filaments is the well-defined force required for the test. The velocity of the anvil after the collision gradually decreases as the anvil gains in height. At the end of its swing the anvil is intercepted by a special device, which stops it falling back against the hammer. Every-

on mechanical strength than the methods previously employed, and moreover the result is obtained without having to wait until the end of the lamp's life. An incidental advantage is that we now have a means of studying *changes* in the mechanical strength during the life of the lamp, as we can take, as it were, a "snapshot" of the strength whenever required.

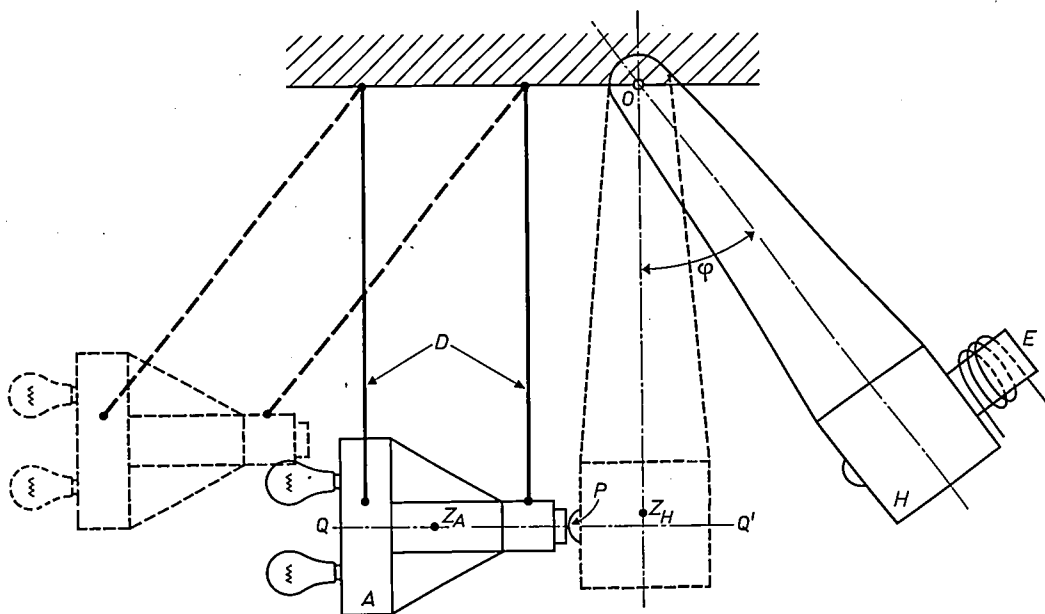


Fig. 4. Schematic representation of the apparatus. Four lamps of the type under test are fitted in the anvil *A*, which is suspended by steel wires *D*. The hammer *H*, which pivots about the axis *O*, is held by an electromagnet *E*. When the magnet current is switched off, the hammer swings against the anvil, producing in it a shock which results in a given acceleration. At the end of its swing the anvil is prevented by a special device from falling back on the hammer. Upon impact the hammer comes completely to rest.

thing is arranged so that the only forces of any consequence acting upon the filament are the inertial forces produced upon impact. The dimensions and weights of hammer and anvil are so chosen that the hammer comes completely to rest upon impact.

The actual test consists in administering a series of shocks, each of equal violence. In a later version of the apparatus this is done automatically, with a frequency of 6 shocks a minute. The effect of the inertial forces is determined by measuring the change in current through each of the four lamps — the current being increased, as mentioned above, if the filament is deformed. The percentage increase of the current through the filament, measured at the end of the series of shocks, thus serves as a quantitative indication of the mechanical strength of the filament. The number and violence of the shocks to be applied are predetermined for each type of lamp and chosen such that the resultant current increase is about 5%.

This test method gives more reliable information

We shall now examine in turn a number of points that received special attention during the development of the apparatus.

Point of impact

The point of impact *P* on the anvil must lie on the line *QQ'* that can be drawn, at the moment of impact, parallel to the velocity of the hammer through the centre of mass Z_A of the anvil (see fig. 4). If the anvil were struck at a different point, the pendular motion described would have superimposed on it a pendular motion around the point Z_A , resulting in unwanted extra forces of inertia on the filaments. The line *QQ'* for the anvil is also a four-fold axis of symmetry. The four lamps under test can therefore be mounted at symmetrical points on the anvil, ensuring that they will all undergo the same acceleration.

The point of impact for the hammer must also satisfy certain conditions. Impact shocks must be prevented from setting up, through the hammer,

reaction shocks in the frame. A reaction shock can be the cause of vibrations in the frames, lengthening of the contact time (see below), wear in the bearings, etc., and can thus jeopardize the reproducibility of the test. Reaction shocks can be eliminated by making the perpendicular distance between the point of suspension O of the hammer and the line QQ' equal to the length of an equivalent simple pendulum with the same period of oscillation as the hammer. The proof is given in an appendix to this article.

Fall angle and contact time

The anvil, which initially has zero velocity, is accelerated upon impact to a velocity v_A , which thereafter gradually decreases to zero. In theory the duration of the impact is infinitely short, but in fact the impact time has a finite value t_c (fig. 5)

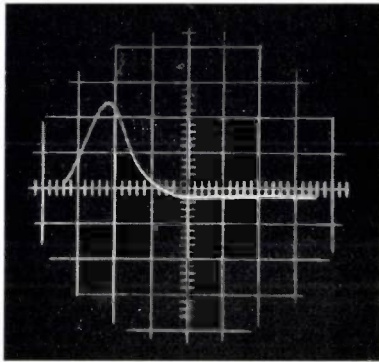


Fig. 5. Oscillogram of the acceleration of the anvil upon impact. The width of a square corresponds to a time of 0.2 millisecond, the height to an acceleration of $100 \times$ the acceleration due to gravity.

due to slight elastic deformation of hammer and anvil at the point of impact. The average acceleration upon impact is consequently v_A/t_c , and it is this that governs the magnitude of the force of inertia which acts on the filaments. For each type of lamp to be tested a suitable average acceleration value is chosen, which must remain constant throughout the test. In this connection the fall angle φ , on which v_A primarily depends, must be accurately adjustable within wide limits, and the contact time t_c should vary as little as possible.

Fall angle

The fall angle φ can be adjusted with an accuracy of about 1° by means of the device illustrated in fig. 6. It consists of a protractor G to which are attached a spirit level W and a needle N .

It can also be seen in fig. 6 that the fork suspension of the magnet hinges around the same axis as the hammer. The magnet itself is mounted in self-

aligning bearings, freely enabling its front face to fit flush against the rear face of the hammer. Mounted in the middle of the front face of the magnet is a microswitch, the contact stud of which projects about 1 mm outside the magnet. The magnet current is switched on only when the hammer touches the contact stud, the magnet thus being able to align itself first without attraction to the back of the hammer. This prevents the hammer being held in the wrong position by the magnet — as far as the backlash in the bearings allows misalignment — which would cause the hammer upon impact to move the anvil laterally. An incidental advantage of this construction is that the current is not switched on until hammer and magnet are practically touching one another, so that there is no risk of their coming together violently, thereby upsetting the adjustment of the fall angle.

The electromagnet must be powerful enough to hold the weight of the hammer with an ample margin at the maximum fall angle. This means that at small angles, where very little attractive force is required, there is a danger that there will be sufficient remanent magnetism to hold the hammer after the current has been switched off. To avoid this, a capac-

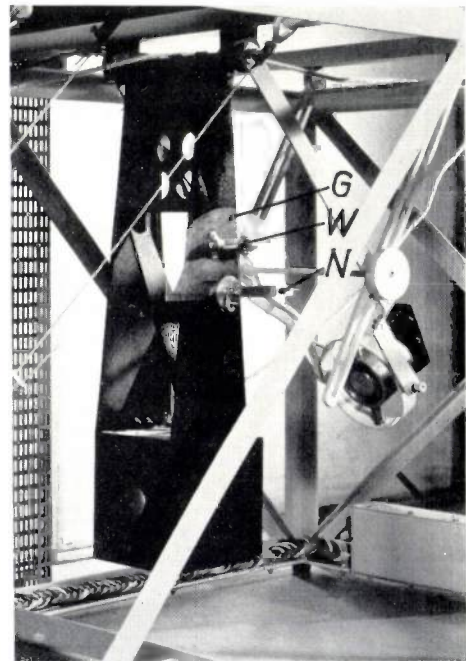


Fig. 6. The hammer with the device mounted on it for adjusting the fall angle φ . Right of the hammer can be seen the electromagnet and its mounting in self-aligning bearings.

The adjustment device is fixed to the hammer when the latter is freely suspended in equilibrium. The protractor G is mounted horizontally with the aid of spirit level W , and the needle N is set to zero on the scale. The protractor is then turned through the required fall angle with respect to the needle, the hammer is drawn up to the magnet, then hammer and magnet are set, by adjusting the magnet suspension mechanism, in the position at which the spirit level is again horizontal.

itor is incorporated in the circuit in parallel with the solenoids, forming with them a resonant circuit. When the direct current is switched off, an exponentially decreasing alternating current flows through the coils, producing a complete demagnetizing effect, so that the magnet is immediately released at small angles as well.

Contact time

The length of the contact time is primarily governed by the hardness of hammer and anvil at the point of impact. To minimize the contact time, the impact faces of hammer and anvil are pieces of hardened steel. Among the other factors affecting it are wear in the bearings and damage to the impact faces. The contact time can be kept sufficiently constant provided the apparatus is built with high precision and is properly maintained.

The contact time is measured with the aid of an electrical contact established between hammer and anvil; this closes a circuit in which a capacitor is charged. The voltage across the capacitor is a measure of t_c . The values found lie between 0.45 and 0.55 millisecond.

Fig. 7 is a nomogram used for giving the average acceleration a predetermined value. After measuring the contact time, one can read from the nomogram the fall angle φ required for the average acceleration chosen. The length of the contact time is kept constantly under observation during the test. If any variation is found, the fall angle is appropriately corrected.

Coefficient of restitution

If we define the magnitude S of the shock as the change of momentum resulting from the impact, then the shock in the hammer is mv_H , where m is the mass of the hammer and v_H its velocity immediately before impact. After impact the hammer has of course zero velocity. During the contact time there is always some slight loss of energy, e.g. due to heat generation, so that the momentum of the hammer is not transferred completely to the anvil. In other words, the shock in the anvil is smaller than that in the hammer by a factor k . This factor is called the coefficient of restitution. Obviously, we want k to be only slightly less than unity, and moreover it should remain constant. By taking special measures — one of them being to make very rigid assemblies of the components of the hammer and the anvil — the value of k , for different apparatuses, can be made to lie between 0.87 and 0.90. The experimentally determined value of k must be taken into account when drawing up the nomogram depicted in fig. 7.

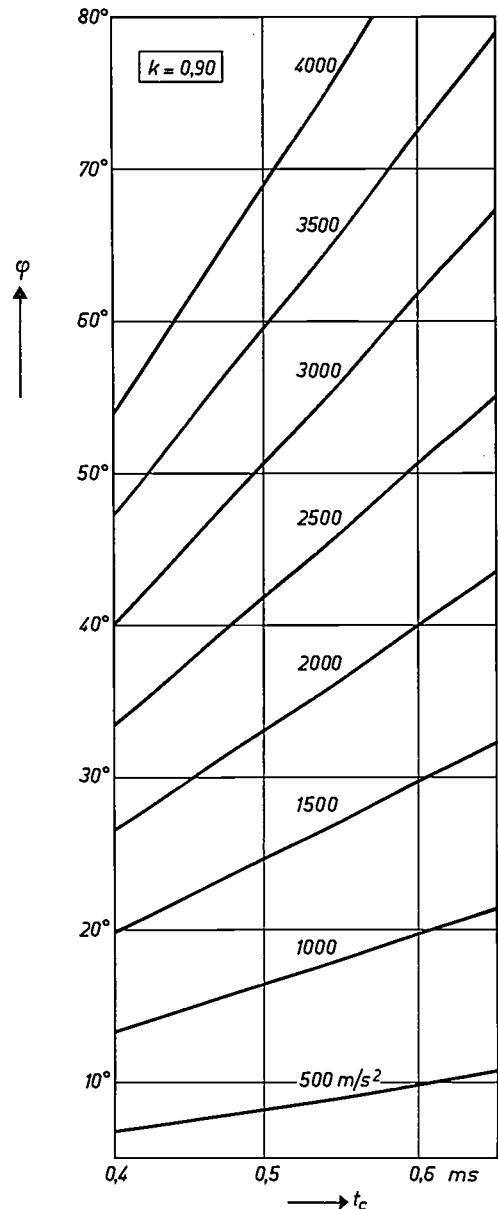


Fig. 7. Nomogram from which the required value of the fall angle φ can be read, as a function of the contact-time length t_c , in order to give the anvil the required average acceleration.

Since k differs from one apparatus to another, a separate nomogram is needed for each apparatus.

Results of measurements

Fig. 8 shows, for a certain type of lamp, a graph of the relative increase of current — averaged over four simultaneously tested examples — as a function of the number of shocks n . The test was carried out at the average acceleration values given in the graph. The sudden upswing of current, corresponding to the bends in each of the curves, is attributable to the formation of loops in the filament, resulting in partial short-circuiting. For testing this lamp type the average accelerations of 3600 m/s^2 and 1350 m/s^2 were unsuitable, the first causing loops to form too

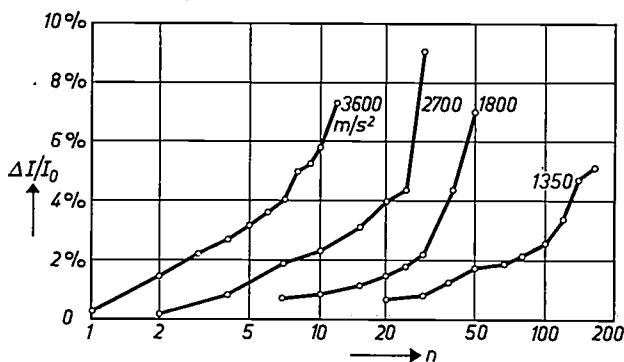


Fig. 8. Relative increase of current through the filament in a given type of lamp, averaged over four simultaneously tested lamps of the same type, as a function of the number of shocks n , for four values of average acceleration.

early, and the second making it necessary to wait a long time before any current increase is observed. The fall angle should thus be chosen so as to produce an average acceleration between these two extremes. An appropriate average acceleration value is found in the same way for every type of lamp.

Appendix: Calculating the position of the impact point on the hammer

Generally speaking the collision between hammer and anvil is accompanied by a reaction shock S_r (fig. 9). We shall demonstrate that, given a suitable perpendicular distance h between the point of suspension O and the line QQ' through the point of impact, this reaction shock can be made zero.

The velocity of the hammer at the moment of impact can be resolved into two components: a linear velocity v_H of the whole body, and an angular velocity ω_H around the centre of mass Z_H (distance z from O). Since the point O is stationary, there exists between these two velocities the relation

$$v_H = z\omega_H \dots \dots \dots (1)$$

Hammer and anvil have the same mass m , and as the head of the hammer has a much greater mass than the shank, the requirement that the hammer should come to rest upon impact is satisfied with sufficient accuracy. This means that the motion of the hammer is entirely converted into a change in the momentum (= shock) S_A of the anvil and into the reaction shock S_r ³⁾. According to the laws of conservation of momentum and angular momentum, we can write

$$mv_H = S_A + S_r \dots \dots \dots (2)$$

and

$$I_Z\omega_H = (h - z)S_A - zS_r, \dots \dots \dots (3)$$

³⁾ Including the coefficient of restitution k in the calculation does not essentially alter the proof.

where I_Z is the moment of inertia of the hammer with respect to the centre of mass. By eliminating v_H and ω_H from (1), (2) and (3) we find for S_r :

$$S_r = S_A \left[\frac{z(h-z)}{I_Z} - \frac{1}{m} \right] : \left[\frac{z^2}{I_Z} + \frac{1}{m} \right].$$

The reaction shock will thus be zero if

$$\frac{z(h-z)}{I_Z} = \frac{1}{m} \dots \dots \dots (4)$$

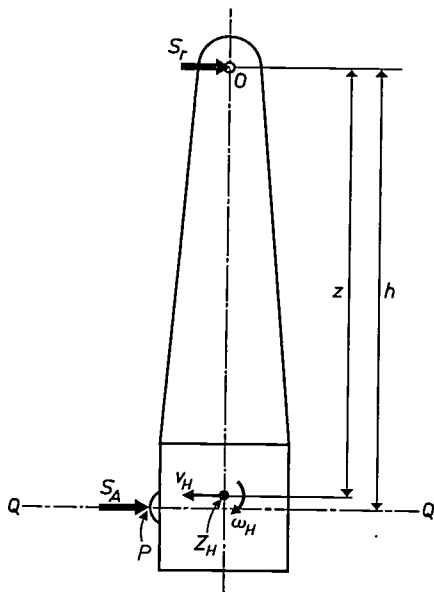


Fig. 9. Diagram used in calculating the most suitable position of the point of impact P .

We can reduce this condition to a simpler form by expressing I_Z as the moment of inertia I_O with respect to the point of suspension, which can be found from the familiar oscillation formula:

$$T = 2\pi \sqrt{\frac{I_O}{mgz}} = 2\pi \sqrt{\frac{l_s}{g}}$$

Here l_s is the length of an equivalent simple pendulum having the same period T as the hammer. Hence

$$I_O = mzl_s.$$

The relation between I_O and I_Z is given by

$$I_O = I_Z + mz^2,$$

so that

$$I_Z = mzl_s - mz^2.$$

Eq. (4) thus reduces to

$$h = l_s,$$

proving that the reaction shock S_r is zero if the distance h is made equal to the equivalent simple pendulum length l_s .

Summary. For use in trams, trains, ships, factories, etc., incandescent lamps are made that are specially designed to withstand shocks and vibrations. This article describes an apparatus used for subjecting such lamps, especially their filament assemblies, to accurately reproducible shock tests. Four lamps at a time are fitted in an anvil suspended by four steel wires. A heavy pendulum (the hammer) is swung through an adjustable fall angle to strike the anvil, the impact producing a shock in the lamps. The percentage increase in the current through the

filament caused by a series of impacts of predetermined force and number serves as a quantitative indication of the filament strength of the lamps. The fall angle of the hammer can be accurately set within wide limits. The duration of contact between hammer and anvil upon impact is measured, and any variation is corrected by adjusting the fall angle. In this way the apparatus can be used for testing many types of lamps under well-defined conditions.

DOUBLE MAGNETIC RESONANCE

by G. E. G. HARDEMAN *).

539.194

The methods of paramagnetic resonance and nuclear magnetic resonance have become firmly established as aids to research into the atomic structure of matter, and even as analysis techniques in industry. Paramagnetic resonance is used for studying, for example, the structure of lattice defects in a solid, and nuclear magnetic resonance yields valuable information on the structure of molecules and crystals. In recent years, experiments have been carried out by which both kinds of resonance are induced in a sample simultaneously. Following earlier articles in this journal on paramagnetic and nuclear magnetic resonance, the present article deals with the theoretical and practical aspects of one of the possible kinds of this "double" resonance in the various forms in which it appears.

Introduction

In addition to the now 15 years old experimental methods of paramagnetic resonance and nuclear magnetic resonance, previously described in this journal^{1) 2)}, in recent years related experimental methods which are based on one form or another of "double" resonance have been developed. One of these depends on the effect that in certain substances nuclear magnetic resonance (n.r.) can be considerably intensified — as we will presently define more exactly — if paramagnetic resonance (p.r.) is induced in the substance at the same time. The converse may also be found, i.e. paramagnetic resonance can be influenced by the simultaneous excitation of nuclear magnetic resonance³⁾.

The intensification of nuclear magnetic resonance, to which we shall confine ourselves in this article, is bound up with the remarkable fact that when paramagnetic resonance is induced in a substance, the spins of the atomic nuclei are aligned to a much more marked extent than would otherwise be the case at the same temperature and in the same external magnetic field. Because of this effect, which is known as "dynamic polarization of nuclear moments", it is possible in some cases to observe nuclear magnetic resonance in samples which, without dynamic polarization, would give an undetectable resonance signal.

In principle, dynamic polarization can be brought about in all substances that contain at the same time paramagnetic centres (e.g. paramagnetic atoms or ions, colour centres or free radicals) and atomic nuclei having a magnetic moment. Both the elec-

trons to which the paramagnetism is due and the nuclei produce a magnetic field in their immediate environment, thus giving rise to magnetic interaction between these particles.

The possibility of dynamic nuclear polarization was first suggested by Overhauser⁴⁾. In 1953 he predicted that if the paramagnetic resonance of the conduction electrons in a metal were induced strongly enough, the nuclear-magnetic-resonance signal of the metal atoms would substantially increase. Shortly afterwards this effect was indeed found by Carver and Slichter⁵⁾ in the metal lithium. The experiment was carried out using a constant external magnetic field of 0.003 Wb/m² (30 gauss), with a frequency of 84 Mc/s for the p.r., and 50 kc/s for the n.r. of the isotope ⁷Li. These alternating fields were generated in two crossed coils in order to minimize mutual induction. The sample was finely powdered to enable the alternating fields to penetrate better into the substance. Without the 84 Mc/s field the n.r. signal remained below the noise level of the detector. When the 84 Mc/s oscillator was switched on, the signal appeared.

After recapitulating the principles of paramagnetic and nuclear magnetic resonance, and at the same time introducing the notation to be used, we shall present in this article a concise account of the mechanisms responsible for dynamic nuclear polarization, illustrated with experiments of our own and of others in which these mechanisms were found in a more or less pure form. Finally we shall describe briefly the instruments required for this kind of investigation.

Paramagnetic resonance and relaxation

Broadly speaking, the possibility of generating paramagnetic resonance is based on the fact that an

*) Philips Research Laboratories, Eindhoven.

¹⁾ J. S. van Wieringen, Philips tech. Rev. 19, 301-313, 1957/58.

²⁾ D. J. Kroon, Philips tech. Rev. 21, 286-299, 1959/60.

³⁾ There are various forms of double resonance. One that has meanwhile become very familiar underlies the operation of the maser. The relation between this and the subject of our article is discussed later (on page 212). Certain forms of double resonance are also used in optical spectroscopy.

⁴⁾ A. W. Overhauser, Phys. Rev. 91, 476, 1953, and 92, 411, 1953.

⁵⁾ T. R. Carver and C. P. Slichter, Phys. Rev. 92, 212, 1953.

electron in a magnetic field (induction B) can be in one of only two quantum states. These two quantum states correspond to different energy levels; the difference in energy between the two levels is ΔE_S . If a system of paramagnetic centres, each containing only one electron that contributes to the magnetic moment, is placed in a magnetic field and exposed to an alternating electromagnetic field whose frequency ν_S is such that $h\nu_S = \Delta E_S$ (where h is Planck's constant), transitions are induced between the two quantum states (fig. 1). If a *relaxation mechanism* is also operative in the paramagnetic substance, i.e. a mechanism that tends in another way to cancel out the change in population of the energy

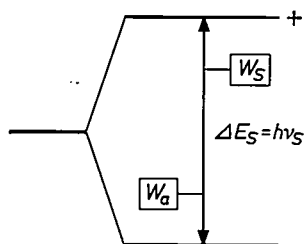


Fig. 1. Energy-level diagram of an electron in a magnetic field. The energy difference ΔE_S is proportional to the induction B . By exposing a sample to an alternating field whose frequency ν_S satisfies the relation $\Delta E_S = h\nu_S$, transitions between the two levels can be induced (paramagnetic resonance). The probability W_a of this happening is proportional to the intensity of the alternating field. The probability of relaxation transitions due to interaction with the crystal lattice is W_S .

levels, the result is that energy can continuously be absorbed from the alternating field. At the frequency ν_S an absorption is then observed which is not present at other frequencies¹). The p.r. spectrum thus consists in the simplest case of a single absorption line. Since ΔE_S depends on the strength of the magnetic field, the frequency ν_S is likewise governed by B . In practice the frequency of the alternating field is generally kept constant (a common value is 10^4 Mc/s, i.e. 3 cm wavelength) while B is slowly varied.

We shall now formulate the principle somewhat more exactly. The energy E_S of a paramagnetic centre (atom, ion) having a magnetic moment μ_S caused by the spins of one or more electrons — we disregard the contribution from the orbital angular momentum of the electrons — is equal to the negative scalar product of the vectors μ_S and B , and has the value $-\mu_S B \cos \vartheta$. Written as a formula:

$$E_S = -\mu_S \cdot B = -\mu_S B \cos \vartheta. \quad (1)$$

Here ϑ is the angle between the directions of μ_S and B . According to quantum theory μ_S is given by:

$$\mu_S = \gamma_S \frac{h}{2\pi} \sqrt{S(S+1)}. \quad (2)$$

Here S is the spin quantum number (the *resultant* spin quantum number for the whole atom), which can assume the values $\frac{1}{2}$, 1, $\frac{3}{2}$, etc., and γ_S is the gyromagnetic ratio. This is connected with the charge e and the mass m of an electron and with the Landé factor g_S according to:

$$\gamma_S = g_S e/2m.$$

For electrons g_S is positive, so that γ_S is negative. (Neglecting the orbital angular momentum, the value of g_S is exactly 2; in reality it is somewhat different.)

If we confine ourselves, as above, to cases where only one electron contributes to the magnetic moment of the paramagnetic centre, we have only one value of μ_S to consider, which is found by substituting in (2) for S the value $\frac{1}{2}$, representing the spin of one electron.

The second condition of quantum theory for (1) is that $\cos \vartheta$ can have only $2S + 1$ values. These are given by the formula:

$$\cos \vartheta = \frac{m_S}{\sqrt{S(S+1)}}, \quad \dots \quad (3)$$

where

$$m_S = S, S-1, \dots, 1-S, -S.$$

If we consider one electron only, then m_S is either $+\frac{1}{2}$ or $-\frac{1}{2}$, and $\cos \vartheta$, as mentioned above, can have only two values. It follows from (2) and (3) that $\mu_S \cos \vartheta$, the component of μ_S in the B direction, is then equal to $\pm \gamma_S h/4\pi$ (i.e. plus or minus one Bohr magneton). This is usually expressed by saying that the spin is either parallel or antiparallel to B . Actually, then, this applies only to the mentioned component of μ_S ; the direction of μ_S itself is at an angle to the B direction. Energy levels and quantities that relate to the case where the spin (the angular momentum) is parallel to B are usually denoted by a plus sign as superscript, the others by a minus sign. It should be added that the angular momentum vector of electrons, owing to γ_S being negative, is oriented in the opposite direction to the magnetic moment.

It follows from the foregoing that:

$$2\pi\nu_S = |\gamma_S| B. \quad \dots \quad (4)$$

The probability $W_a dt$ that an electron will change from one state to the other in the time interval dt by the emission or absorption of a quantum $h\nu_S$ does *not* depend on the direction of the transition. The value of W_a is proportional to the intensity of the alternating field, i.e. proportional to the square of the amplitude of the field. This is due to the fact that the transitions take place solely under the

influence of the alternating field (stimulated transitions). There is no or scarcely any spontaneous emission of quanta at the frequencies involved, and W_a has therefore no connection with such emission. (The same applies to nuclear magnetic resonance, to the amplification of microwave energy in a maser ⁶⁾, etc.)

In order to define more precisely the combination of the absorption of energy from the alternating field and the relaxation effect, and to calculate the magnitude of the power that can continuously be extracted from the alternating field, we must consider the number of electrons per unit volume present in each of the two energy levels (quantum states). For the upper level we shall call this number n^+ (spin "parallel" to B) and for the lower level n^- (see fig. 1). The difference $n^+ - n^-$ is the *spin polarization* Δn ; we use n to denote the total number of electrons $n^+ + n^-$.

If n^+ and n^- are equal, so that $\Delta n = 0$, the number of transitions per unit time is equal in both directions. If $\Delta n \neq 0$, more transitions take place in one direction than in the other. The power P_s extracted per unit volume from the alternating field is then:

$$P_s = -W_a \Delta n E_s. \quad \dots \quad (5)$$

Thus, P_s is positive if Δn is negative, i.e. if n^- is greater than n^+ .

If there were no transitions of another nature, then according to the definition of W_a :

$$\frac{d}{dt} \Delta n = -2W_a \Delta n. \quad \dots \quad (6)$$

The spin polarization Δn would then decrease to zero exponentially with time, and so would P_s . No further absorption would then be observable.

In fact, however, thermal contact exists between the system of electron spins and the lattice of the substance. Consequently, in the absence of an external alternating field, the populations n_{th}^+ and n_{th}^- of the levels will be unequal. We then have a Boltzmann distribution, where

$$n_{th}^+ / n_{th}^- = \exp(-\Delta E_s / kT). \quad \dots \quad (7)$$

This ratio is greater the lower is the temperature. We can then write for the spin polarization:

$$\Delta n_{th} = n_{th}^+ - n_{th}^- = \frac{\exp(-\Delta E_s / kT) - 1}{\exp(-\Delta E_s / kT) + 1} n. \quad (8)$$

If we represent $-\Delta E_s / kT$ by $2\delta_s$, we can reduce eq. (8) for most cases to

$$\Delta n_{th} \approx n \delta_s, \quad \dots \quad (9)$$

⁶⁾ See, for example, J. Volger, Solid-state research at low temperatures III, Philips tech. Rev. 22, 268-277, 1960/61.

for at a resonant frequency of 10^4 Mc/s the quantity $|2\delta_s|$ is very much smaller than 1, even at 4 °K.

Plainly, the probability of a transition being caused by this interaction with the lattice cannot be the same for both directions, for if we denote the probability of an upward transition by W_s^{-+} and the other by W_s^{+-} , it is obvious that in the equilibrium state

$$n_{th}^+ W_s^{+-} = n_{th}^- W_s^{-+}.$$

Writing the product of both probabilities as W_s^2 , we then have:

$$W_s^{+-} = W_s \exp(\frac{1}{2} \Delta E_s / kT) \approx W_s (1 - \delta_s),$$

$$W_s^{-+} = W_s \exp(-\frac{1}{2} \Delta E_s / kT) \approx W_s (1 + \delta_s).$$

After a perturbation, the thermal distribution (7) is restored exponentially, in accordance with a solution of the following equation, which is analogous to (6):

$$\frac{d}{dt} \Delta n = 2W_s (\Delta n_{th} - \Delta n). \quad \dots \quad (10)$$

The characteristic time found in the solution is the "relaxation time" τ_s , which is equal to $1/2W_s$.

With the aid of the foregoing expressions we can now describe the population of the levels in a paramagnetic-resonance experiment, and derive a formula for the absorbed power P_s . In the absence of an equilibrium state, Δn changes in accordance with:

$$\frac{d}{dt} \Delta n = 2W_s (\Delta n_{th} - \Delta n) - 2W_a \Delta n. \quad (11)$$

After equating the right-hand side to zero, we find that in the equilibrium state Δn is given by:

$$\Delta n = \frac{W_s}{W_a + W_s} \Delta n_{th}. \quad \dots \quad (12)$$

The power then continuously absorbed from the alternating field is:

$$P_s = -W_a \Delta n \Delta E_s = -\frac{W_a W_s}{W_a + W_s} \Delta n_{th} \Delta E_s. \quad (13)$$

Since, as stated, W_a is proportional to the square of the amplitude \hat{b} of the external alternating field, its value can vary so considerably in most cases that it is useful to investigate what happens to the right-hand side of (13) in the extreme cases $W_a \ll W_s$ (very weak external field) and $W_a \gg W_s$ (very strong external field). In the first case, according to (12), Δn is approximately equal to Δn_{th} , and we can write:

$$P_s \approx -W_a \Delta n_{th} \Delta E_s.$$

The absorbed power is thus in this case proportional

to \hat{b}^2 . Where the alternating field is very strong, Δn approaches zero, and we have:

$$P_S \approx -W_S \Delta n_{th} \Delta E_S.$$

In a strong alternating field the absorption is therefore virtually independent of the intensity of the field (saturation case). The only factors then governing P_S are W_S — or in other words the relaxation time τ_S — and Δn_{th} , both of which in their turn are temperature-dependent. As can be deduced from (9), Δn_{th} increases with decreasing temperature, whereas W_S usually becomes smaller (the relaxation time increases ⁶⁾). The most suitable value of \hat{b} in p.r. experiments is as a rule relatively small and is determined by a variety of factors.

Nuclear magnetic resonance

What has been said above for electrons applies, with due alteration of details, to atomic nuclei. All we need do is to substitute for the spin quantum number S the nuclear spin I , and the quantities m_p (proton mass) and g_I for m and g_S . Confining ourselves as before to the case where $I = \frac{1}{2}$, we again have two possible spin orientations, i.e. “parallel” and “antiparallel” to B . (Once again, we mean by this that the component of the angular momentum in the B direction is parallel or antiparallel to B .) The level populations are denoted here by N^+ and N^- . The total number of nuclei is N , and the quantity $\Delta N = N^+ - N^-$ is now called *nuclear spin polarization*. For nuclei the gyromagnetic ratio is:

$$\gamma_I = g_I e/2m_p.$$

Here $|g_I|$ has a value that can vary from 0.1 to 10, while γ_I can be either positive or negative. The energy difference in a constant field B is:

$$\Delta E_I = \gamma_I \frac{h}{2\pi} B. \quad \dots \quad (14)$$

This situation is sketched in *fig. 2* for the case where $\gamma_I > 0$; the + state then has the lower energy and the — state the higher energy. The resonant frequency obeys the relation:

$$2\pi\nu_I = \gamma_I B. \quad \dots \quad (15)$$

Since $m_p \approx 2000 m$, the value of $|\gamma_I|$ is several hundred to several thousand times smaller than that of $|\gamma_S|$, so that, given identical B , the resonance effect will appear at a much lower frequency.

Here too there is some form of thermal contact with the lattice, and the level population tends towards the thermal distribution. The governing quantity $\Delta E_I/kT$ we call $-\delta_I$. If $|\delta_I| \ll 1$, we can

write, to a very good approximation, for the case of thermal equilibrium:

$$\Delta N_{th} \approx N\delta_I. \quad \dots \quad (16)$$

Here δ_I has the same sign as γ_I . The discussion of absorption follows the same lines as for paramagnetic resonance. Introducing, by analogy with W_a and W_S , the quantities W_a' and W_I , we can write for the resonance absorption P_I :

$$P_I = W_a' \Delta N \Delta E_I. \quad \dots \quad (17)$$

Given a small W_a' , i.e. a weak alternating field, we can again put ΔN in this expression equal to ΔN_{th} . In a strong alternating field, ΔN again approaches zero, and the absorption becomes:

$$(P_I)_\infty = W_I \Delta N_{th} \Delta E_I. \quad \dots \quad (18)$$

Usually W_I is much smaller than W_S , which means that the nuclear magnetic relaxation time $\tau_I = 1/2W_I$ is as a rule much longer than the paramagnetic relaxation time τ_S . Since, broadly speaking, $\Delta E_I \approx 10^{-3} \Delta E_S$ and $\Delta N_{th} \approx 10^{-3} \Delta n_{th}$, the n.r. absorption is roughly $10^6 \times$ smaller than the p.r. absorption.

Double resonance

From equations (13) and (17) it can be deduced that in both p.r. and n.r. the absorbed power is proportional to the polarization of the relevant particles. This explains why, as mentioned at the beginning, n.r. absorption increases if the nuclear spin polarization can be increased.

The possibility of causing this increase by the simultaneous induction of p.r. transitions — dynamic nuclear polarization — which is the basis of the type of investigation considered in this article, depends on the existence of various kinds of magnetic interaction between the nuclei and the electrons responsible for the paramagnetism. In the following section we shall consider the nature of these interactions in more detail and indicate the effects which they produce. We shall limit our discussion to the two principal forms of the interactions: scalar coupling and dipole coupling.

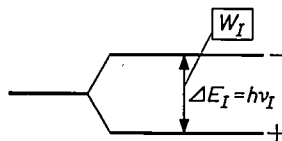


Fig. 2. Diagram of the energy levels of a nucleus in a magnetic field. The energy difference ΔE_I is roughly $1000 \times$ smaller than for an electron at the same B value. Transitions can be induced with an alternating field whose frequency ν_I satisfies $h\nu_I = \Delta E_I$ (nuclear magnetic resonance). W_I is the probability of relaxation transitions.

Dynamic nuclear polarization due to scalar coupling

In many cases the magnetic interaction between electron and nucleus is due to the fact that, in terms of quantum mechanics, the electron has a certain probability of being at the position of the nucleus. If we imagine that the magnetic moment μ_s corresponding to the electron spin is caused by a circular current, then there is equally a certain probability that the nucleus will be surrounded by this current, so that the nucleus experiences a local magnetic field of the electron. This local field is proportional to μ_s . The energy a of this coupling is proportional to the scalar product of the vectors I and S :

$$a = A \mathbf{I} \cdot \mathbf{S},$$

where A is a proportionality constant and (cf. formula 2):

$$|S| = \sqrt{S(S+1)} \quad \text{and} \quad |I| = \sqrt{I(I+1)}.$$

As we shall see, this coupling often causes the p.r. line to be split into a number of components — the hyperfine structure. For this reason A is sometimes called the “hyperfine-structure constant”. For simplicity we again assume that one electron per atom is responsible for the paramagnetism, so that we can say that the nuclei and electrons in the substance occur in pairs, for each of which four states are possible: in both orientations which the nuclear spin can have in relation to the external magnetic field, the electron spin can be parallel or antiparallel to the nuclear spin. We denote these four states by $++$, $+-$, $-+$ and $--$, the first sign relating to the nuclear spin, and $+$ and $-$ meaning parallel or antiparallel to B . If there is any scalar coupling between nucleus and electron, the energy levels of these states cannot be found simply by taking the sum of $\frac{1}{2}h\nu_s$ and $\frac{1}{2}h\nu_I$ (with the appropriate sign); a further quantity $\pm\frac{1}{2}a$ must be added. For a case where $\gamma_I > 0$, this is $+\frac{1}{2}a$ if the spins are mutually parallel and $-\frac{1}{2}a$ if they are antiparallel. Expressed as a formula (cf. fig. 3):

$$\left. \begin{aligned} E_{+-} &= -\frac{1}{2}h\nu_s - \frac{1}{2}h\nu_I - \frac{1}{2}a \\ E_{--} &= -\frac{1}{2}h\nu_s + \frac{1}{2}h\nu_I + \frac{1}{2}a \\ E_{++} &= +\frac{1}{2}h\nu_s - \frac{1}{2}h\nu_I + \frac{1}{2}a \\ E_{-+} &= +\frac{1}{2}h\nu_s + \frac{1}{2}h\nu_I - \frac{1}{2}a \end{aligned} \right\} \quad (19)$$

When an external magnetic field is applied this means that theoretically, in our case, the p.r. line and the n.r. line should each split into two components. In p.r. this is directly observable. In n.r. this is for various reasons not observable, and the

splitting of the levels can only be established by indirect means.

The case represented in fig. 3 is referred to as resolved hyperfine structure. If the coupling extends to a number of nuclei surrounding one paramagnetic centre, numerous different resonant frequencies occur and a single, broadened and structureless line is observed. This is frequently the case in solids.

We can now approach much closer to our objective, which is to explain the intensification of nuclear polarization by the induction of p.r. transitions, if we now consider the following anomalous form of hyperfine coupling. In cases where this anomaly occurs, a single, more or less sharp absorption line is observed instead of a split or broadened line, in spite of the fact that scalar coupling is certainly present. It appears that this effect can only occur in samples in which the nuclei and electrons are capable of moving relative to one another. This is true, for example, in metals, in which the conduction electrons can move in relation to the quiescent nuclei, and also in liquids, in which the paramagnetic centres are capable of movement relative to their environment. In that case the magnitude of A can change considerably in a time that corresponds to the period of the difference frequency $|\nu_{13} - \nu_{24}| = 2a/h$, and the result is that the frequencies ν_{13} and ν_{24} can no longer be separately observed. (The time during which they are emitted is small compared with the time needed to observe one period of the beat frequency.) It then seems as if

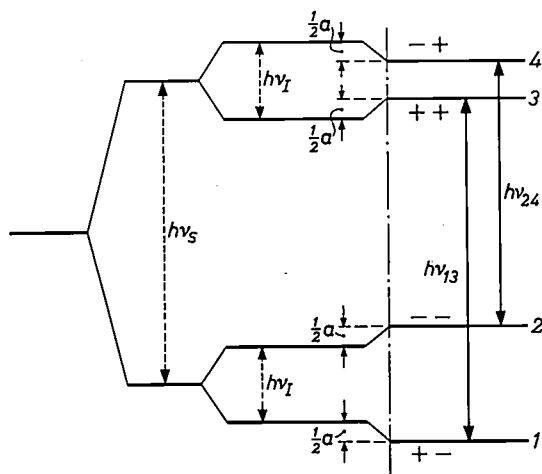


Fig. 3. Right of the dot-dash line: energy levels of the nucleus-electron system in the case of (weak) scalar coupling. Left of the line can be seen how the energy of the four levels is related to $h\nu_s$, $h\nu_I$ and the coupling energy a . As a result of the coupling, the paramagnetic resonance line (original frequency ν_s) is split into two components with frequencies ν_{13} and ν_{24} , which are respectively higher and lower than ν_s by an amount a/h . The signs $+$ and $-$ indicate whether the angular momentum (the spin) is “parallel” or “antiparallel” to B ; the first sign relates to the nucleus, the second to the electron.

A is zero: only the splittings ΔE_s and ΔE_l are left, and $\nu_{13} = \nu_{24} = \nu_s$, resulting in a single narrow resonance line. This effect is referred to as (Brownian) "motional narrowing".

Although there is no perceptible splitting into components in this effect, and it looks on the face of it as if there were no coupling, there is in fact a fundamental difference between both cases. Because the instantaneous values of A are not zero, it is possible here, just as in the case where A is constant (and unequal to zero), for transitions to take place in which the spin orientations of nucleus and electron are both reversed; without coupling only "pure" nuclear and electron transitions are possible. The "combined" transitions can only take place between the states $+-$ and $-+$; the transitions $++ \rightleftharpoons --$ remain "forbidden" (fig. 4). An important difference between the case where A is constant and that where A rapidly fluctuates is that in the first case the transitions $+- \rightleftharpoons -+$ can only be caused by stimulated absorption or emission of quanta under the influence of an external alternating field of suitable frequency — just as in the case of normal paramagnetic or nuclear magnetic resonance (cf. fig. 3) — whereas in the case of motional narrowing this is completely out of the question. The transition here takes place solely without the emission or absorption of radiation.

Thus the allowed "combined" transitions found in samples in which A rapidly varies are an additional kind of relaxation transitions. Like the relaxation transitions already discussed, their transition probability W_0 contains a Boltzmann factor to

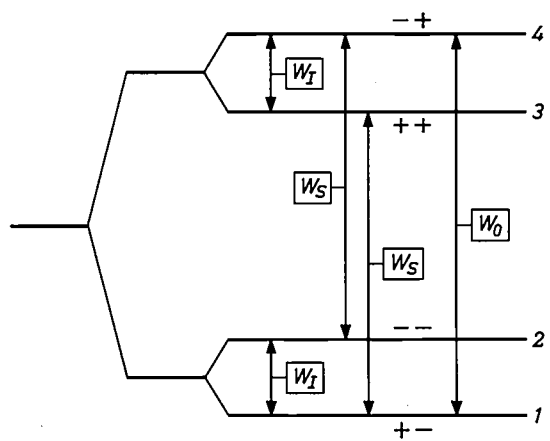


Fig. 4. Energy-level diagram of the nucleus-electron system in the case of rapidly fluctuating scalar coupling. The coupling now causes no shift in the levels, and there is no splitting of the p.r. line (motional narrowing). The transition between the levels 1 and 4 is now, however, no longer forbidden, and forms an extra relaxation transition (probability W_0).

allow for the difference in energy between the relevant levels. The formula is:

$$W_{0(j \rightarrow k)} = W \exp [\frac{1}{2}(E_j - E_k)/kT]. \quad (20)$$

Here W is a constant, and the subscripts j and k denote the energy levels involved.

The Overhauser effect

We shall now consider a sample in which the above-mentioned relaxation transitions $+- \rightleftharpoons -+$ can take place, and have a considerably greater probability of occurring than nuclear relaxation transitions. A qualitative picture of what happens in such a sample is presented in fig. 5.

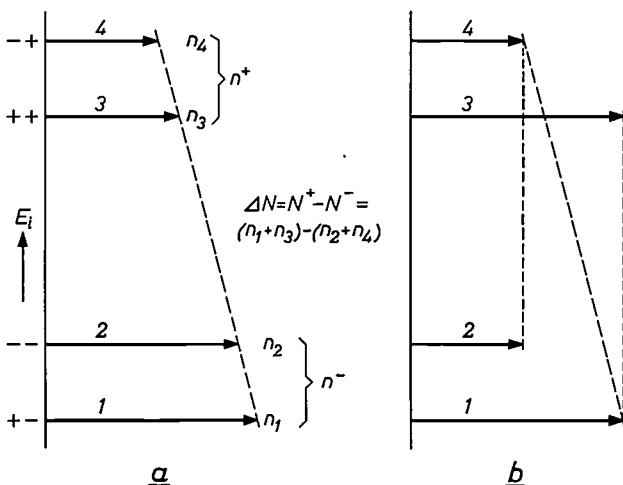


Fig. 5. Explanation of the Overhauser effect. a) Energy-level population in thermal equilibrium. Plotted vertically is the energy E_i of the levels (cf. fig. 4); the length of the arrows indicates the population n_i . (If the energy differences between the levels are small, n_i decreases linearly with increasing E_i ; the decrease is shown here grossly exaggerated.) b) Saturation of the paramagnetic resonance makes $n^+ (= n_3 + n_4)$ equal to $n^- (= n_1 + n_2)$, but owing to the extra relaxation transition between the levels 1 and 4 the ratio between n_1 and n_4 continues to satisfy the Boltzmann distribution. As a result the nuclear polarization, ΔN , which governs the nuclear-magnetic-resonance signal, shows a marked increase.

This figure represents schematically the population n_i of the four possible energy levels. In the absence of an external alternating field (fig. 5a), the population is governed by the value of the Boltzmann factor for each of the four energies (see formula 7). Where the levels have small mutual differences in energy, as in this case, it can be said that the population decreases linearly with increasing energy; see the dashed line. If p.r. transitions are now induced (fig. 5b), sufficiently to produce saturation, the population changes first of all in such a way that n^+ becomes approximately equal to n^- . Owing to the extra relaxation, however, and the high speed at which it takes place, the relation between the populations of levels 1 and 4 remains equal to that required by the Boltzmann distribu-

tion. It can be seen from the figure that the nuclear polarization $(n_1 + n_3) - (n_2 + n_4)$ is now greater. This is called the Overhauser effect.

Using the formulae given above, we can easily calculate the increase of the nuclear polarization. First, according to (20), we have:

$$\frac{n_4}{n_1} \approx 1 + 2\delta_s + 2\delta_I \approx 1 + 2\delta_s. \quad (21)$$

But:

$$\frac{n_4}{n_1} = \frac{N^- n^+}{N^+ n^-}.$$

Consequently, since $n^+ = n^-$ in paramagnetic resonance saturation, in this case holds:

$$\frac{N^+}{N^-} \approx 1 - 2\delta_s,$$

giving for the nuclear polarization:

$$\Delta N \approx -N\delta_s. \quad (22)$$

As we saw earlier, $\Delta N_{th} \approx N\delta_I$ (eq. 16). By saturating the paramagnetic resonance we have thus caused the nuclear polarization in a weak AC field to increase by a factor p , which is given by:

$$p = -\delta_s/\delta_I. \quad (23)$$

If the nuclear relaxation is not negligible compared with the relaxation due to the transitions $+- \rightleftharpoons -+$, then the factor p is smaller or the Overhauser effect may even vanish. The ratio n_4/n_1 cannot then reach the value given by (21). In many cases the condition mentioned is fulfilled. Especially in solids having pure scalar coupling between electrons and nuclei, W_I is usually very small and the nuclear relaxation is effectively the result of the two other relaxation mechanisms (those having the transition probabilities W_s and W_0 ; see fig. 4). It should be remembered, however, that in general among the solids only metals will exhibit an Overhauser effect. In non-metals the mobility is usually too low.

We should add that the Overhauser effect does not always lead to a difference ΔN that has the same sign as ΔN_{th} . When γ_I is negative, the upper level has the larger population, and nuclear magnetic emission takes place instead of increased absorption. The system may then be regarded as a *maser* that amplifies radiation of the frequency ν_I . Where, moreover, there is strong feedback to the coil generating the alternating field required for nuclear magnetic resonance, it is possible in principle for the whole system to start oscillating. One would then have a nuclear magnetic *generator*. Such a system would no longer require an external supply

of energy, and the transmitter used to generate the alternating field could be switched off. The latter has not yet proved feasible; it has, however, proved possible to achieve this with a sample dynamically polarized by the interaction discussed in the next section.

Some experiments on the Overhauser effect

As mentioned at the beginning, the effect was first found by Carver and Slichter⁵⁾, in lithium. In the isotope ${}^7\text{Li}$ with which they worked, the polarization factor p can theoretically have the value 1690. Owing to the limited power of their transmitter (50 W), they were unable to saturate the paramagnetic resonance, but even so the effect was clearly observable. In the metal lithium, then, the time-dependence of the scalar coupling on which the effect depends, is attributable to the high mobility of the electrons. The ${}^7\text{Li}$ nuclei can be regarded here as fixed points in an "electron liquid".

The Overhauser effect in a real liquid was found, again by Carver and Slichter⁷⁾, in a solution of sodium in liquid NH_3 . The dissolved Na atoms are paramagnetic; in this case the electron has a scalar coupling with the protons of the solvent.

In this laboratory the effect was found in a solid, namely the free radical diphenyl-pikrylhydrazyl⁸⁾. The adequate fluctuation of the scalar coupling is presumably attributable here to the fact that the electrons in this substance can easily exchange their mutual spin orientations. As a result a particular state never exists for long, so that even though the electrons themselves are not so highly mobile, the coupling rapidly fluctuates. This also appears from the p.r. line, which is clearly narrower than would be expected if no motional narrowing were assumed. There is not yet any exact theory of this effect.

Combrisson⁹⁾ and co-workers observed the Overhauser effect in silicon doped with phosphor as impurity. The P atoms behave in this case as donors. At a P content of $3 \times 10^{16}/\text{cm}^3$ or less, these donors already give up electrons at temperatures higher than 10 °K. These electrons exhibit scalar interaction with the nuclei of the isotope ${}^{29}\text{Si}$, of which natural Si contains 5% (the bulk is mainly ${}^{28}\text{Si}$). The ${}^{29}\text{Si}$ nuclei have a spin quantum number $I = \frac{1}{2}$, and the ${}^{28}\text{Si}$ nuclei possess no moment. At temperatures above 10 °K the condition for the Overhauser effect is therefore satisfied and an appreciably amplified ${}^{29}\text{Si}$ nuclear magnetic resonance is found. Since

⁷⁾ T. R. Carver and C. P. Slichter, Phys. Rev. **102**, 975, 1956.

⁸⁾ H. G. Beljers, L. van der Kint and J. S. van Wieringen, Phys. Rev. **95**, 1683, 1954.

⁹⁾ J. Combrisson and I. Solomon, J. Phys. Radium **20**, 683, 1959.

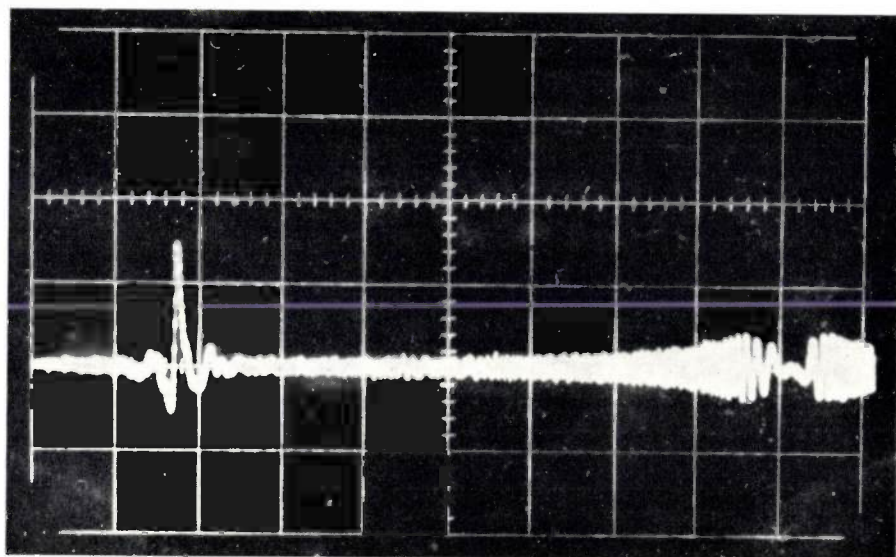


Fig. 6. Overhauser effect observed on silicon containing phosphor as impurity (the positive peak on the left). Nuclear magnetic emission occurs here, caused by the isotope ^{29}Si present in the proportion of 5% in natural silicon. The phosphor content was 10^{18} atoms/cm 3 , the temperature 4.2 °K. Without dynamic nuclear polarization the n.r. peak (a negative absorption peak) would not, in this case, be distinguished from the noise. In the experiment the induction B and the frequency $\nu_{p.r.}$ of the alternating field generating the paramagnetic resonance were kept constant; the frequency $\nu_{n.r.}$ of the alternating field used for generating the nuclear magnetic resonance was varied. On the right can be seen a beat-frequency effect; $\nu_{n.r.}$ here passes the frequency of a calibrating oscillator. (The small peaks beside the n.r. peak are attributable to an instrumental effect.)

the gyromagnetic ratio γ_I is negative for ^{29}Si , the Overhauser effect in this case not only increases ΔN but also reverses its sign in relation to ΔN_{th} ; the nuclear magnetic absorption thus changes to emission, as described above.

At temperatures below 10 °K the electrons remain in local stationary orbits around the P nuclei. The magnetic interaction is then not time-dependent and the Overhauser effect cannot occur. In this temperature range another uncommon effect is found, which will be dealt with below.

Given a 30 times larger P content (about 10^{18} atoms/cm 3) there is still mobility even below 10 °K, since neighbouring P centres are then able to exchange electrons. In this case the Overhauser effect can still be observed at the temperature of liquid helium ($T \leq 4.2$ °K); see fig. 6.

Dynamic nuclear polarization due to dipole coupling

In the foregoing we have considered the interaction between two particles each having a magnetic moment and which, in quantum-mechanical terms, have a certain chance of being at the same position. If they are not, the interaction is not eliminated, but the coupling energy no longer satisfies the simple formula $a = A I.S$; each of the particles is now located in a magnetic field caused by the other particle, and this field is identical with a dipole field. In practice the probability of the particles

being at the same position will seldom be exactly zero and cannot be exactly one, and so a mixture of both interactions is found. Normally, however, one of the two strongly predominates.

The possibility of producing dynamic nuclear polarization by means of dipole coupling between atomic nuclei and "paramagnetic" electrons is again due to additional relaxation transitions¹⁰; see fig. 7. Apart from the transition $+- \rightleftharpoons -+$, which also occurred with scalar coupling, the transition $++ \rightleftharpoons --$ is now "permitted". The probability of this transition we call W_2 ; the component in the B direction of the

sum of nuclear and electron spins changes upon such a transition by two units $h/2\pi$. Furthermore the probability of pure relaxation transitions increases, by an amount W_1^{SI} for the nuclear spin and by W_1^{IS} for the electron spin. Ordinarily W_1^{SI} is much greater than W_1 but W_1^{IS} is much smaller than W_2 ; in other words the nuclear spin-lattice relaxation becomes many times greater but the electron relaxation changes only by a small fraction. If the particles are highly mobile, motional narrowing occurs

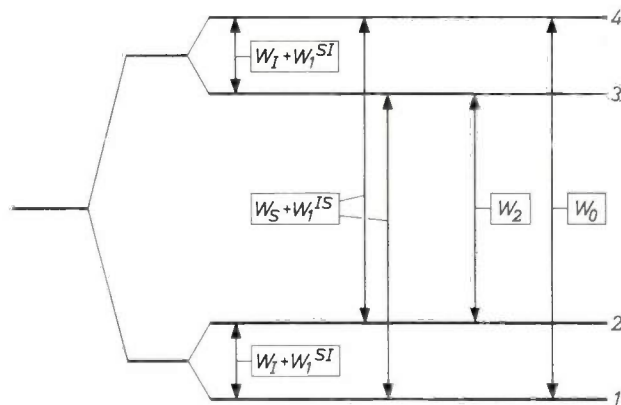


Fig. 7. Energy-level diagram of the nucleus-electron system in the case of dipole coupling with motional narrowing. As opposed to the case in fig. 4, there is also a possible transition between $++$ and $--$ (probability W_2). W_1 and W_2 have both increased, by the amounts W_1^{SI} and W_1^{IS} respectively.

¹⁰ This possibility was first reported by A. Abragam, Phys. Rev. **98**, 1729, 1955.

here too, but this does not prevent the extra transitions. On the contrary, the mobility increases the probability of the transitions, and it is this that makes dynamic nuclear polarization really possible.

The consequence of all this is that, in p.r. saturation, a nuclear polarization occurs of the magnitude:

$$\Delta N = \frac{1}{2} N \delta_S \quad \dots \quad (24)$$

The polarization factor therefore becomes:

$$p = \frac{1}{2} \delta_S / \delta_I \quad \dots \quad (25)$$

As can be seen (eq. 22), the nuclear polarization is half as great as with scalar coupling, and of opposite sign.

Owing to the numerous possible transitions, the derivation of (24) is much more difficult than that of the corresponding formula (22) for scalar coupling. The main lines, however, can readily be indicated. The derivation is based on what is called the "principle of detailed balancing", according to which, in an equilibrium state, as many particles reach an energy level as leave it per unit time. If we apply this principle to the four levels in our case, neglecting W_1 in relation to W_1^{SI} and again assuming equal numbers of nuclei and electrons ($N = n$), we can derive from the four equations obtained:

$$\begin{aligned} (W_0 + 2W_1^{SI} + W_2) (\Delta N - \Delta N_{th}) = \\ = (W_0 - W_2) (\Delta n - \Delta n_{th}) \quad \dots \quad (26) \end{aligned}$$

Quantum-mechanical perturbation theory tells us, however, that where the particles are highly mobile the transition probabilities are given by:

$$\left. \begin{aligned} W_1^{SI} = W_1^{IS}, \text{ call this } W_1, \\ \text{and} \\ W_0 : W_1 : W_2 = 2 : 3 : 12. \end{aligned} \right\} \quad \dots \quad (27)$$

Substituting this in (26) and taking into account that the electron polarization Δn is zero in p.r. saturation, we then find:

$$\Delta N = \Delta N_{th} + \frac{1}{2} \Delta n_{th} \quad \dots \quad (28)$$

The right-hand side of this equation, because $N = n$, is to a good approximation equal to that of (24).

Some experiments on the dipole effect

Experimental proof of the existence of a strong dynamic nuclear polarization, opposed to that found with the Overhauser effect, was provided in 1957 almost simultaneously by a group of European and a group of American research workers. The first group¹¹⁾ found the effect in a solution of potassium-

peroxylaminedisulphonate $K_2NO(SO_3)_2$ in water, and the other group¹²⁾ found the effect in a solution of naphthalene and sodium in 1,2-dimethoxyethane. Both substances in the dissolved state form paramagnetic ions the electrons of which have a dipole coupling with the protons of the solvent. In the naphthalene and sodium solution a p value of -65 was found; the theoretically possible value of -330 was not attained because the power of the transmitter was too low.

After prolonged experiments with the peroxylaminedisulphonate solution, it proved possible to make the nuclear magnetic emission strong enough to compensate the losses in the LC circuit of the transmitter. So in fact the nuclear magnetic generator mentioned in the foregoing was produced¹³⁾. This result was achieved not because a high p value was obtained — it was only 50 compared with a theoretical maximum of 110 — but because a very dilute solution had been used. Consequently the n.r. line was extremely narrow and the intensity of the nuclear magnetic emission (per wavelength interval), which is inversely proportional to the line width, was therefore very high (fig. 8).

As regards the latter experiments, it should be noted that in a powerful external magnetic field, strong scalar coupling with the nucleus of the nitrogen atom of the ion itself occurs, in addition to the dipole coupling of the "paramagnetic" electron with the protons of the solvent. The nitrogen nucleus has a spin quantum number $I = 1$, and can therefore have three orientations in an external magnetic field. The p.r. line is thus split into three components.

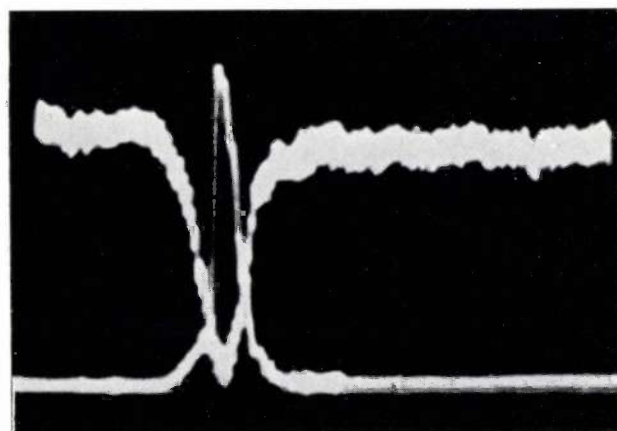


Fig. 8. Nuclear magnetic resonance of potassium-peroxylaminedisulphonate $K_2NO(SO_3)_2$ in water. The negative peak is the absorption peak found for n.r. alone. The induction of p.r. gave rise through dipole coupling (with motional narrowing) to dynamic nuclear polarization, which reversed the peak and made it more than 50 times larger. The latter peak was recorded with less amplification. The line width is only 0.5×10^{-4} Wb/m² (0.5 gauss). (After Allais¹³⁾.)

¹¹⁾ A. Abragam, J. Combrisson and I. Solomon, C. R. Acad. Sci. **245**, 157, 1957.

¹²⁾ L. H. Bennett and H. C. Torrey, Phys. Rev. **108**, 499, 1957.

¹³⁾ E. Allais, C. R. Acad. Sci. **246**, 2123, 1958.

Of course, the transmitter can only be tuned to one of the three frequencies, so that given a sufficiently powerful transmitter the best that can be achieved is: $\Delta n = \frac{2}{3}\Delta n_{th}$, instead of $\Delta n = 0$. The theoretically possible p value, which would be -330 without the line splitting, is consequently reduced to the above-mentioned value of -110 .

If the external field is weak, the scalar coupling with the nitrogen nucleus is in fact advantageous¹⁴⁾. The coupling of the electron with the nitrogen nucleus can then be much stronger than that with the external field. The resonance frequency ν_s of the p.r. is in that case much higher than that which would be found without the coupling to the same (weak) field. The ratio $p = \delta_s/\delta_I$ is likewise increased. In the earth's magnetic field of about 40 A/m (0.5 oersted) a p value of -800 has been reached, again producing a nuclear magnetic generator, at a p.r. frequency ν_s of 55 Mc/s and an n.r. frequency of 2 kc/s¹⁵⁾.

This might be used for accurately measuring weak magnetic fields. The precision is almost entirely governed by the accuracy with which ν_I can be measured. A magnetometer of this kind could be very useful to geologists, for example, for the purpose of tracing slight disturbances in the earth's field.

To avoid misunderstanding we should add that the extremely high p value does not imply that the nuclear polarization itself is very high, resulting in correspondingly strong signals. Because of the weakness of the external field the opposite is the case.

The solid-state effect

In a solid it is not possible, via dipole couplings, to produce dynamic nuclear polarization by inducing paramagnetic resonance. True, the transitions $+- \rightleftharpoons -+$ and $++ \rightleftharpoons --$ are possible here too if there is dipole coupling, but owing to the low mobility W_0 , W_1 and W_2 are much smaller than in a liquid, and are the smaller the greater is the corresponding energy jump. Consequently, equations (27) and (28) are no longer valid for solids, and W_1^{SI} is large compared with W_1^{IS} and W_2 . The excitation of paramagnetic resonance now no longer modifies the nuclear magnetic polarization, because the transitions with probability W_1^{SI} , which restore the thermal equilibrium of the nuclear system, are now predominant.

Dynamic polarization of the nuclear spins in a solid, given dipole interaction, nevertheless proves

to be possible in an entirely different way. To illustrate which phenomena occur, we give a description of the experiment which led to the discovery of this solid-state effect. We shall then give a short explanation of the effect and an account of some recent investigations and applications.

Discovery and nature of solid-state effect

In 1958 Erb, Motchane and Uebersfeld¹⁶⁾ investigated the double resonance of powdered coal upon which a liquid containing hydrogen (e.g. water or benzene) had been adsorbed. Coal contains paramagnetic centres which it was expected would show scalar or dipole interaction with the protons. The p.r. line was excited in the usual way at a fixed microwave frequency by adjusting the magnetic field to the value B_s at which p.r. transitions can take place. As could be expected, dynamic polarization was not observed. Just as in other solids, the magnetic coupling appeared not to fluctuate fast enough to give rise to an Overhauser or dipole effect.

When, however, the magnetic field was altered slightly from the correct setting, a considerably amplified n.r. signal appeared. For $B < B_s$ the amplified signal was negative and for $B > B_s$ positive. Upon further alteration of B the effect disappeared. This effect may be explained as follows¹⁷⁾.

As we have mentioned, given dipole coupling, the transitions $+- \rightleftharpoons -+$ and $++ \rightleftharpoons --$, which are required for nuclear polarization as extra relaxation transitions, are not forbidden in solids, but there is little probability of their occurring through interaction with the lattice. In the same way as normal p.r. and n.r. transitions, however, they can be artificially induced by exposing the sample to an alternating field whose frequency corresponds to the external magnetic field and the relevant energy jump. Although the stimulated transitions occur with less probability than p.r. and n.r. transitions, their probability is greater than that of nuclear spin-lattice relaxation, the latter transition probability being fairly low in solids — considerably lower than in liquids. Consequently the induced transitions bring about a marked shift in the level population. This shift differs depending on which of the two transitions is induced. We shall briefly consider both cases.

The transition $+- \rightleftharpoons -+$, which for $\gamma_I > 0$ corresponds to the transition 1-4 in fig. 7, can al-

¹⁴⁾ J. Freycenon and I. Solomon, *Onde élect.* 40, 590, 1960.

¹⁵⁾ J. Freycenon, *Onde élect.* 40, 596, 1960.

¹⁶⁾ E. Erb, J.-L. Motchane and J. Uebersfeld, *C. R. Acad. Sci.* 246, 2121, 1958.

¹⁷⁾ A. Abragam and W. G. Proctor, *C. R. Acad. Sci.* 246, 2253, 1958.

ready be saturated with a moderately strong alternating field. This makes N^+n^- equal to N^-n^+ . Since the required frequency differs from that at which the p.r. transitions occur, p.r. is not produced. Now provided W_s is not too small, the electrons continue to be governed by the thermal equilibrium condition:

$$\frac{n^+}{n^-} = 1 + 2\delta_s.$$

Since N^+n^- has been made equal to N^-n^+ , it follows that:

$$\frac{N^+}{N^-} = 1 + 2\delta_s,$$

from which it can be seen that the nuclear polarization has increased to:

$$\Delta N = N^+ - N^- = N\delta_s. \quad \dots (29a)$$

The nuclear polarization is thus greater by a factor $p = \delta_s/\delta_I$ and nuclear magnetic emission occurs ($\delta_s < 0$ and $\delta_I > 0$).

The transition $++ \rightleftharpoons --$ corresponds, for $\nu_I > 0$ again, to the transition 2-3 in fig. 7. If we adjust the frequency and intensity of the alternating field so as to saturate the latter transition, we obtain $N^+n^+ = N^-n^-$, giving:

$$\Delta N = -N\delta_s. \quad \dots (29b)$$

The nuclear polarization has thus increased by the factor $p = -\delta_s/\delta_I$, in other words it is just as great as in the previous case, but of opposite sign.

The relation between the foregoing and the effect found by Erb *et al.* is obvious: the result of making the magnetic field weaker was to bring the (constant) transmitter frequency into line with the frequency required to induce the transition $+- \rightleftharpoons -+$; the augmenting of B by the same amount made it possible to induce the transition $++ \rightleftharpoons --$.

The amount ΔB by which B must be increased or decreased depends of course on the difference in magnitude of the energy jumps 1-4 and 2-3 on the one hand and 1-3 (or 2-4) on the other. This difference is precisely ΔE_I (see fig. 2), or $h\nu_I$, so that:

$$\Delta B = \pm 2\pi\nu_I/\gamma_s. \quad \dots (30)$$

At, for example, $\nu_s = 10^4$ Mc/s and $\nu_I = 15$ Mc/s (protons), ΔB is thus roughly 5×10^{-4} Wb/m² (5 gauss).

Since ΔB is fairly small, and since moreover the dipole interaction causes extra broadening of the p.r. line, the two polarization effects cannot always be induced purely and separately. It may well happen that $B + \Delta B$ and $B - \Delta B$ fall within the

region of B values that corresponds to the width of the p.r. line. In such a case both polarization effects will occur simultaneously and oppose one another, so that the p value found is too small. This difficulty can sometimes be minimized by making B very high. In those cases where the line width does not increase with B , ν_I is then sufficiently high to cause $B + \Delta B$ and $B - \Delta B$ to fall outside the p.r. line.

Nuclear spin diffusion

For the dynamic polarization of nuclear spins in general, and for the solid-state effect in particular, it is not necessary that there should be magnetic interaction between every nucleus and an electron. The solid-state effect is already found when there is one paramagnetic centre for every 10^5 to 10^4 nuclei. This means that most nuclei will notice nothing or very little of an electron spin. That the whole system of nuclei can nevertheless be polarized is attributable to the process of *nuclear spin diffusion*. Whilst there is scarcely any diffusion of the atoms themselves in a solid, the nuclei can readily exchange their spin orientation with that of their nearest neighbours. In the solid-state effect, nuclei in the neighbourhood of the paramagnetic centres are dynamically polarized; this polarization is then transmitted to the other nuclei. For this diffusion process to function properly, the characteristic time $\tau_I (= 1/2W_I)$ of the nuclear spin-lattice relaxation must be long, otherwise the extra polarization might vanish during the process, with the absorption or surrender of the corresponding energy from or to the lattice. To make τ_I long enough it is necessary to cool the substance to a low temperature.

Our reason for examining nuclear spin diffusion here is that without this effect it would, as a rule, be impossible to observe the solid-state effect. Where the samples used have such a high concentration of paramagnetic centres that dynamic polarization should be possible even without nuclear spin diffusion, the appearance of the solid-state effect may be prevented by, for example, the p.r. line being too broad, the nuclear spin relaxation too fast, etc. Sometimes too the sample may no longer be paramagnetic at all. In the following description of various experiments we shall therefore place more emphasis on nuclear spin diffusion than on the solid-state effect as such.

Some experiments

The solid-state effect can be observed on many samples, e.g. in polymers in which paramagnetic centres have been created by irradiation with γ rays

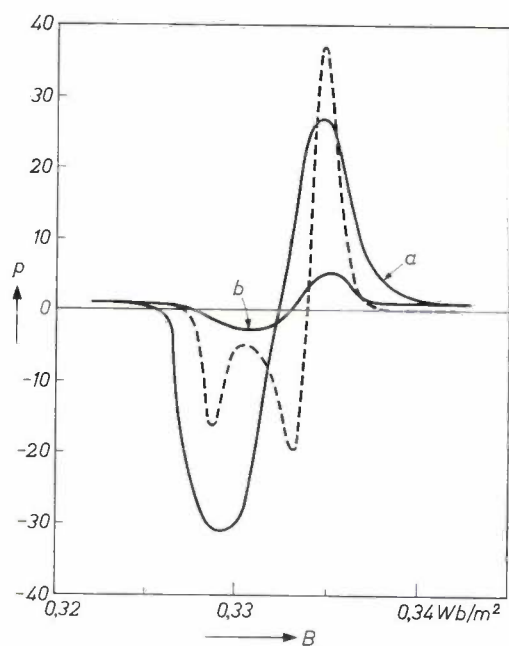


Fig. 9. Dynamic nuclear polarization of ^{19}F in polytetrafluoroethylene after irradiation with fast electrons²⁰). The solid lines indicate the variation of the factor p by which the nuclear polarization has increased, as a function of B for two paramagnetic-centre concentrations (curve a : 1.9×10^{-4} ; curve b : 5.8×10^{-5}). The dashed curve represents the derivative of the p.r. line. As can be seen, instead of one peak at the B value corresponding to ν_S (about 0.333 Wb/m^2), there is a negative peak at a smaller B value and a positive peak at a higher B (solid-state effect). Recorded at 4.2°K .

or with fast electrons, and also on solids containing colour centres or paramagnetic ions as impurities^{18) 19)}. In this laboratory, for example, the solid-state effect has been studied on polytetrafluoroethylene in which paramagnetic centres had been produced by bombardment with fast electrons²⁰⁾. The absorption was studied by the usual practice of varying the external magnetic field while keeping the frequency constant. The measurements were taken at a temperature of 4.2°K . Fig. 9 shows the amplification factor p as a function of B for two different concentrations of centres. The figure relates to the polarization of the ^{19}F nuclei. The dashed curve shows the derivative of the paramagnetic resonance absorption curve, the width of which indicates that the transitions $+-\rightleftharpoons-+$ and $++\rightleftharpoons--$ cannot be separately induced. At the lowest concentration only small p values are found. The nuclear spin diffusion is then not capable of transmitting the polarization throughout the substance.

When the magnetic interaction between the nuclei, on which the nuclear spin diffusion depends,

becomes too weak owing to the fact that the nuclei have a relatively small moment and are far apart, it may take a long time before the spin polarization has spread through the whole sample. This is the case with SiC, which consists largely of atoms that have no nuclear moment. The only isotopes possessing a magnetic moment are ^{29}Si and ^{13}C , which occur in the natural isotope mixture in the proportions of 5% and 1% respectively. If the SiC contains a small quantity of some paramagnetic impurity, paramagnetic resonance occurs²¹⁾ and in that case the nuclei mentioned can be dynamically polarized (fig. 10). At $T = 4.2^\circ\text{K}$, $B \approx 0.3 \text{ Wb/m}^2$, the polarization time — i.e. the characteristic time of the exponential

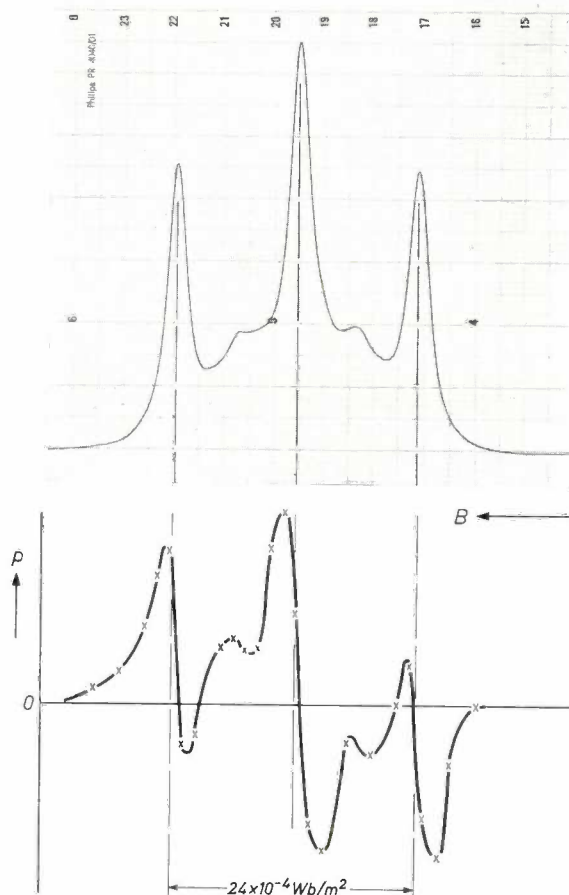


Fig. 10. The solid-state effect in SiC doped with nitrogen. As the spin quantum number of the N nucleus is equal to 1, the p.r. line (upper curve) is split into three. Beside each of the components of the p.r. line in the curve for p (lower curve) can be seen the two opposite peaks which are characteristic of the solid-state effect. The curves relate to the nuclear magnetic resonance of ^{29}Si . The induction B is again along the abscissa. The recording was made at 4.2°K with fixed $\nu_{\text{p.r.}}$.

¹⁸⁾ J. A. Cowen, W. R. Schafer and R. D. Spence, Phys. Rev. Letters **3**, 13, 1959.

¹⁹⁾ M. Abraham, M. A. H. McCausland and F. N. H. Robinson, Phys. Rev. Letters **2**, 449, 1959.

²⁰⁾ G. E. G. Hardeman, Philips Res. Repts **15**, 587, 1960.

²¹⁾ This was found in our laboratories by J. S. van Wieringen: Int. Colloqu. Halbleiter u. Phosphore, Garmisch-Partenkirchen 1956, publ. Vieweg, Brunswick 1958, p. 367-370. For further publications on SiC, see H. H. Woodbury and G. W. Ludwig, Phys. Rev. **124**, 1083, 1961.

function in which the polarization approaches its final value — was found, in a particular sample, to be 80 minutes for the ^{29}Si nuclei and more than two hours for the ^{13}C nuclei (fig. 11). The nuclear relaxation time τ_1 was thus evidently much longer.

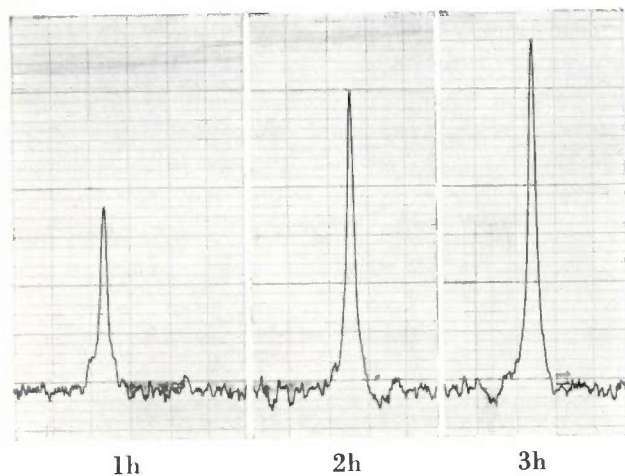


Fig. 11. Nuclear magnetic resonance of ^{13}C in a SiC sample weighing only 80 mg. The nuclear polarization gradually spreads through the whole sample by "nuclear spin diffusion". The curves show the n.r. signal after one, two and three hours from the moment of switching-on the alternating field that induces the p.r. For this recording the frequency was varied, and B kept constant; one scale division corresponds to 2 kc/s. Temperature 4.2°K , n.r. frequency (ν_1) 3.7 Mc/s.

Remarks on the equipment for double-resonance experiments

Apart from the fact that two alternating fields are needed for experiments of the type described in this article, the equipment required is much the same as that used for ordinary nuclear magnetic or paramagnetic experiments. A description of the latter will be found in the articles mentioned in references ¹⁾ and ²⁾. It may be inferred from the theory and experiments described in the foregoing, that here too it is advantageous on the whole to use strong magnetic fields, e.g. $> 0.3 \text{ Wb/m}^2$. The reason is that the energy absorption in paramagnetic and nuclear magnetic resonance increases, according to (5) and (17) respectively, with increasing ΔE_3 and ΔE_1 respectively, i.e. with increasing B . If the aim is to produce the highest possible nuclear polarization, a strong magnetic field is in fact indispensable.

The chief instrumental problem is combining the two fields. The n.r. frequency is relatively low, so that the relevant field can be generated with a coil forming part of an ordinary LC circuit. With fairly weak magnetic fields the p.r. frequency ν_s is also low enough to permit the use of a coil; in that case crossed or mutually screened coils are employed. If the magnetic field is fairly strong, ν is of the order

of 10^4 Mc/s , and the alternating field that induces the p.r. can best be produced with a resonant cavity. Use is then made of a resonant cavity combined with a coil in such a way that both components again influence each other as little as possible.

Motchane, Erb and Uebersfeld ²²⁾ in their above-mentioned experiments used a cylindrical resonant cavity in which an oscillation was excited in the TE 011 mode. The coil was mounted inside the cavity, and so positioned as to fulfil the boundary conditions for producing a standing wave in the cavity ²³⁾. For this purpose the coil had to be wound with almost rectangular turns (fig. 12). The quality factor produced in the cavity was $Q = 4000$. The size of the samples was 0.1 cm^3 , and they were placed centrally in the coil-and-cavity system.

Borghini and Abragam ²⁴⁾, using a frequency ν_s

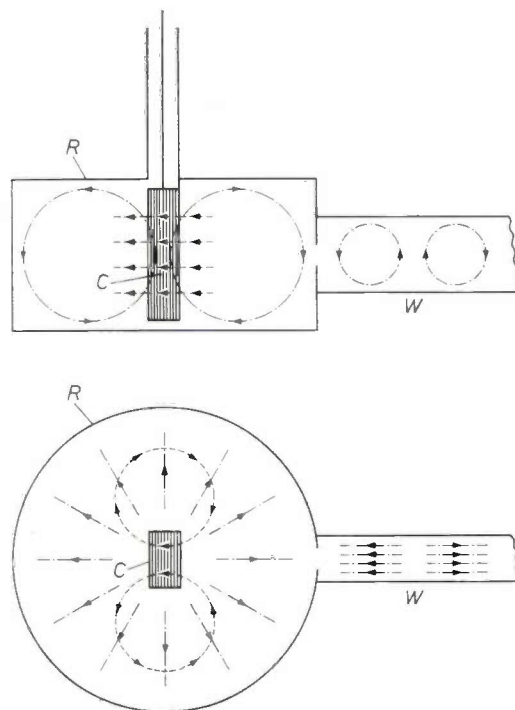


Fig. 12. Combination of resonant cavity and coil for generating the alternating fields for p.r. and n.r., as used by Motchane, Erb and Uebersfeld ²²⁾. In the cylindrical cavity R an oscillation in the TE 011 mode is generated; the dot-dash lines represent the magnetic lines of force. Energy is supplied through the waveguide W . The coil C (magnetic lines of force dashed) is fed via a coaxial cable. The sample is placed in the centre. For the sake of simplicity the lines of force of both fields are drawn as circles.

²²⁾ J.-L. Motchane, E. Erb and J. Uebersfeld, C. R. Acad. Sci. **246**, 1833, 1958.

²³⁾ This implies that the magnetic induction lines of the standing wave must always be parallel to the conducting surfaces, while the electric lines of force must be perpendicular to the surfaces. Nor should any magnetic flux from the cavity field pass through the coil. The coil field and that of the standing wave are accordingly at right angles to one another.

²⁴⁾ M. Borghini and A. Abragam, C. R. Acad. Sci. **248**, 1803, 1959.

of about 3.5×10^4 Mc/s, worked with a rectangular resonant cavity oscillating in the TE 102 mode. The coil was situated in the boundary plane between the two standing half-waves (fig. 13). A small space was left between the coil and the wall of the cavity in order to allow the induction lines of the coil to form a closed loop around the edges of the coil.

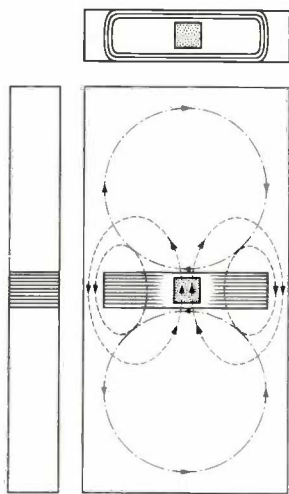


Fig. 13. Combination of resonant cavity and coil as used by Borghini and Abragam²¹). The resonant cavity is made to oscillate in the TE 102 mode (dot-dash line). The coil is situated in the boundary plane between the two standing half-waves. The gap between the coil and the wall of the cavity enables the lines of force of the coil field (dashed lines) to form closed loops. The sample is again in the centre of the system.

In our own experiments we used an oval cavity half a wave in length, at $\nu_s = 10^4$ Mc/s. Fig. 14 shows a sketch of the arrangement. An aperture 2 is cut into the base 3 of the resonant cavity *R*. The coil *C* is wound around the outside of the lower part of the cavity and over the aperture. The magnetic flux of the coil can only pass the aperture in the base and return through the ports *I* in the side walls of the cavity. The sample (volume about 0.1 cm^3) is placed in the centre of the coil above the aperture 2.

In view of the considerable size of the aperture 2, if nothing were done to prevent it the cavity would lose a great deal of its microwave energy, thus greatly reducing its *Q*. For this reason the bottom of the system is enclosed in a brass cap 4. This makes the space behind the aperture totally reflecting, resulting in a *Q* of about 2500.

By enclosing the resonant cavity from underneath, the cap offers the further advantage of enabling the system to be immersed in a liquid, for example liquid helium in low-temperature experiments. The presence of a liquid in the resonant cavity would upset the tuning, and gas bubbles would also interfere with the signal. To ensure good thermal contact

between the samples and the bath, the cavity is filled with helium gas. The cavity is connected by a stainless steel waveguide *W* to the system which detects the paramagnetic resonance. A coaxial line *T*₂ of German silver — a poor heat conductor like stainless steel — containing a thin inner conductor of copper, connects the coil with the system that generates the low frequencies and which detects the nuclear magnetic resonance.

As can be seen, the alternating fields can be produced in various ways. Which method is to be preferred depends entirely on the purpose of the investigation. In general there are three requirements to be met.

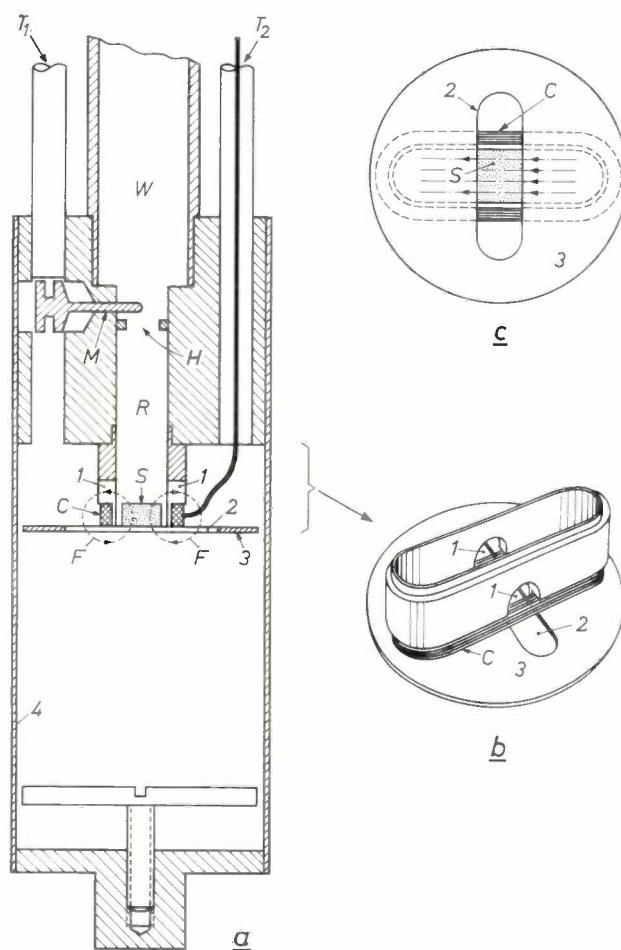


Fig. 14. The system of resonant cavity and coil as used by the author. a) Cross-section of the assembly with a part of the waveguide connected to the resonant cavity. *R* resonant cavity, connected via the coupling slot *H* to the waveguide *W*. *M* matching stub; this can be screwed by means of an eccentric tube *T*₁. *C* coil, fed via a coaxial cable *T*₂. *S* sample. The lines of force *F* of the alternating field generated by the coil, form closed loops through the ports *I* in the side wall of the cavity and through a large oval aperture 2 in the base plate 3. This is shown more clearly in the perspective sketch *b*. The sample is mounted above this aperture and rests on e.g. a quartz plate (not drawn) which is clamped underneath 3. The whole system is enclosed in a cap 4 which totally reflects the microwaves and makes it possible to immerse the system in a bath of liquid helium. c) View of plate 3 from underneath.

1) The resonant cavity must not be unduly large. This is a particularly important point in cryogenic work; with liquid helium, for example; a double cryostat is needed — an inner one for the liquid helium, and surrounding it an outer cryostat containing liquid air for thermal insulation. A large resonant cavity would then mean a very voluminous system, and the magnet would have to have a very wide air gap. An extremely large magnet would thus be needed to generate a sufficiently strong magnetic field.

2) The filling factor should be high, i.e. the sample should occupy as much as possible of the volume containing the alternating field. For observing n.r. it is the coil that should have the high filling factor, and for p.r. the resonant cavity.

3) Both the coil and the resonant cavity should have the highest possible Q . In order to produce dynamic nuclear polarization in a fluid it is necessary to have a strong alternating field with the frequency ν_s . A cavity with a low Q would then make excessive demands on the power output of the transmitter. Moreover, the detection sensitivity of the magnetic resonance (both n.r. and p.r.) is in general proportional to the quality factor.

Finally we remark that the significance of dynamic nuclear polarization as an experimental technique goes beyond its value as a means of investigating magnetic resonance as such. The dynamic polarization of nuclear moments, whether or not with the aid of nuclear spin diffusion, is a technique for the orientation of nuclei that can be useful in other fields of research, and offers advantages over the "brute force" method. The latter requires an extremely strong magnetic field or an extremely low

temperature, or both. Consider, for example, the case of the experiment on polytetrafluorethene (see page 217). Given the same magnetic field without dynamic polarization, the substance would have to be cooled to 0.04 °K to achieve the same spin polarization. There are considerable objections to this, particularly where the nuclei are bombarded with particles such as neutrons, in which case the heat generated soon proves troublesome. Using dynamic polarization the substance need not be cooled below the temperature of liquid helium, whilst the nuclear spin system has the required low effective temperature. Not very long ago the first experimental success was achieved with this application in investigations concerning the scattering of polarized fast protons by dynamically polarized protons in a sample²⁵).

²⁵ A. Abragam, M. Borghini, P. Catillon, J. Coustham, P. Roubeau and J. Thirion, *Phys. Letters (Amsterdam)* 2, 310, 1962 (No. 7).

Summary. After recapitulating the theory of paramagnetic resonance (p.r.) and nuclear magnetic resonance (n.r.), the author deals in detail with the effects observed when both forms of resonance are induced simultaneously. As a consequence of the magnetic interaction (coupling) between the electron spins and the nuclear spins, this can cause dynamic nuclear polarization, resulting in amplification of the n.r. signal. Nuclear polarization can arise through scalar coupling — which occurs when the electron, quantum-mechanically speaking, has a certain probability of being at the same position as the nucleus — and through dipole coupling. The effects differ depending on whether the electrons and nuclei are in relative quiescence or movement. The theory of all these cases is briefly discussed, one or more examples being given, where possible, for each case.

The effects described, apart from aiding the analysis of molecular and crystalline structures, can be turned to use in other fields, e.g. for accurately measuring weak magnetic fields and for nuclear physical research on oriented nuclei. Finally, mention is made of the requirements to be met by the experimental equipment, and some existing equipment is briefly discussed.

HYDROGEN IN IRON AND STEEL

I. SOLUTION AND PRECIPITATION

by J. D. FAST *) and D. J. van OOIJEN *).

546.72:546.11

When new theories are put forward in a particular field of science, they tend for a while to claim the undivided attention of research workers in that field, while experiences that do not at once fit into the new picture are, for the time being, disregarded.

Dislocation theory in metallurgy is a case in point, and led to the neglect of a great deal of earlier knowledge of the effects of hydrogen in iron and steel. On the basis of recent experiments the authors show that it is only by combining dislocation theory with older concepts that a satisfactory picture can be built up, especially in regard to the unexpected ruptures that can be caused by hydrogen in steel structures. This subject has latterly attracted considerable attention, since the hydrogen embrittlement of steel is thought to have been one of the causes of air crashes and other accidents.

Introduction

Of the impurities in iron and steel, hydrogen is one of the greatest sources of danger. In this article, which appears in two parts, we shall review some of the more important aspects of the effects of hydrogen in iron and steel. No attempt will be made to do full justice to the relevant literature, which already comprises over several thousand articles¹⁾. To deal with the subject at any length would thus amount to compiling a thick and almost unreadable book containing numerous contradictory views. Our aim here will rather be to present a coherent picture based on modern insight and on our own experience, and reference will be made to only a fraction of the literature of interest for our purposes. In Part II we shall deal in particular with the role played by hydrogen in the embrittlement of steel.

To begin with we shall discuss the solubility and diffusion of hydrogen and its interaction with dislocations. Attention will be drawn to some effects which appear to be bound up with the formation of gaseous hydrogen at high pressure in lattice defects. It is this high-pressure hydrogen that underlies the deleterious action mentioned.

Solubility and diffusion of hydrogen in iron

The solubility of hydrogen in iron has been shown by many investigations to be proportional to the

square root of the pressure of the gaseous hydrogen with which the metal is in contact²⁾. It follows from this that the solute hydrogen is not present in the metal as molecules H_2 but as atoms H. (The same applies to hydrogen dissolved in other solid or liquid metals.)

Lattice constant and density measurements have shown that the hydrogen is dissolved in metals *interstitially*, that is to say, the solute hydrogen atoms do not replace metal atoms in the crystal lattice but occupy the spaces between the metal atoms, i.e. the interstices.

The solubility of hydrogen in a metal increases with rising temperature, in accordance with well-known thermodynamic laws, if the solution process is endothermic, i.e. one in which, at constant temperature, heat is absorbed from the surroundings. On the other hand the solubility decreases with rising temperature if the solution process is exothermic, i.e. one in which, at constant temperature, heat is given up to the surroundings. As regards the iron-hydrogen system the first holds true: the solubility increases with rising temperature. *Fig. 1* shows the solubility of hydrogen in iron, at 1 atm, as a function of the temperature, according to the literature already cited²⁾.

The fact that the solubility of hydrogen in molten iron is so much greater than in solid iron, as can clearly be seen from *fig. 1*, has consequences of the utmost importance to the manufacture of iron and steel. The sudden drop in solubility upon solidification can cause *porosity* in castings and welds. An

*) Philips Research Laboratories, Eindhoven.

¹⁾ A large proportion of this literature is referred to in "Circular 511" of the National Bureau of Standards (U.S.A.), entitled "Hydrogen embrittlement of steel, review of the literature", 1951, and in the following reviews: Donald P. Smith, Hydrogen in metals, Chicago Univ. Press 1948; P. Cotterill, The hydrogen embrittlement of metals, Progr. Materials Sci. 9, No. 4, 1961; M. Smialowski, Hydrogen in steel, Pergamon Press, London 1962.

²⁾ A. Sieverts, G. Zapf and H. Moritz, Z. phys. Chem. A 183, 19, 1938/39; M. H. Armbruster, J. Amer. Chem. Soc. 65, 1043, 1943; W. Geller and Tak-Ho Sun, Arch. Eisenhüttenw. 21, 423, 1950.

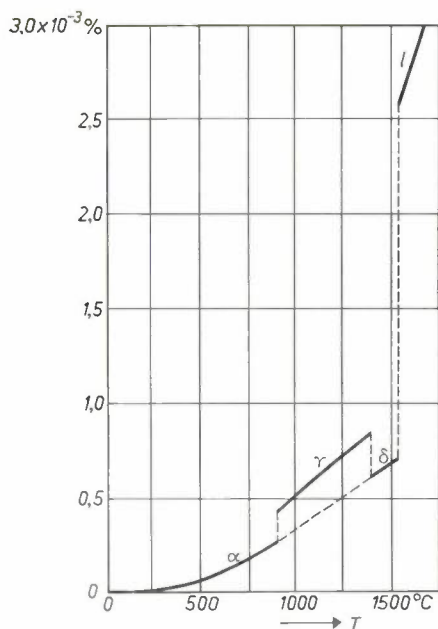


Fig. 1. The solubility of hydrogen in iron at a H_2 pressure of 1 atm in percentages by weight, as a function of temperature in $^{\circ}C$. The symbols α , γ , δ and l denote respectively α , γ and δ iron and molten iron.

extreme case is demonstrated in *fig. 2*, which shows the pores and cavities in an iron bar, formed by melting pure iron in hydrogen at a pressure of 1 atm and pouring it in a water-cooled copper chill.

Among the elements that can dissolve interstitially in metals (H, C, N, O), hydrogen is exceptional in that its diffusion coefficient is very much higher than that of the other elements mentioned; in ferrite (α iron) at $20^{\circ}C$, for example, it is no less than 10^{12} times greater than that of carbon and nitrogen (see *Table I*). This exceptionally rapid diffusion can be explained by assuming that the hydrogen moves as a *proton* from one interstice to another — the diameter of a proton being a mere hundred-thousandth of that of an atom or ion of carbon or nitrogen. This does not exclude the possibility of the hydrogen being present in the interstices as an *atom*; it then only “jumps” as a proton. It is probable that a dissociation equilibrium is involved, of the simple form



where p and e are respectively a proton and an electron.

Table I. Diffusion coefficients D of N, C and H in ferrite.

Temp. $^{\circ}C$	$D(N)$ cm^2/s	$D(C)$ cm^2/s	$D(H)$ cm^2/s
20	8.8×10^{-17}	2.0×10^{-17}	1.5×10^{-5}
100	8.3×10^{-14}	3.3×10^{-14}	4.4×10^{-5}
200	1.7×10^{-11}	1.0×10^{-11}	1.0×10^{-4}
300	5.3×10^{-10}	4.3×10^{-10}	1.7×10^{-4}
400	6.0×10^{-9}	5.9×10^{-9}	2.5×10^{-4}
500	3.6×10^{-8}	4.1×10^{-8}	3.3×10^{-4}
700	4.4×10^{-7}	6.1×10^{-7}	4.9×10^{-4}
900	2.3×10^{-6}	3.6×10^{-6}	6.3×10^{-4}

In agreement with the foregoing it is found that in the metals investigated in this respect — iron³⁾, palladium⁴⁾ and tantalum⁵⁾ — solute hydrogen is electrolytically transported towards the negative pole when an electric current, produced by a DC potential, is passed through the metal.

Traps

We have spoken up to now only of the solubility and diffusion of hydrogen in an *ideal* iron crystal. The iron or steel normally used is polycrystalline, and moreover each crystal contains impurities and lattice defects. The impurities occur both in the form of solute foreign atoms and in the form of separate phases (e.g. as carbide or nitride). At these “imperfections” the hydrogen atoms find sites that are energetically more favourable than the normal interstitial sites, especially in the dislocations (*fig. 3*), and also in the neighbourhood of certain foreign

³⁾ A. Herold, Colloque sur la diffusion à l'état solide (organisé à Saclay, 1958), p. 133, North Holland Publ. Co., Amsterdam 1959.

⁴⁾ A. Coehn *et al.*, *Z. Physik* **62**, 1, 1930; *ibid.* **71**, 179, 1931; *ibid.* **83**, 291, 1933.

⁵⁾ J. Wesolowski, J. Jarmula and B. Rozenfeld, *Bull. Acad. Pol. Sci. chim.* **9**, 651, 1961 (No. 10).

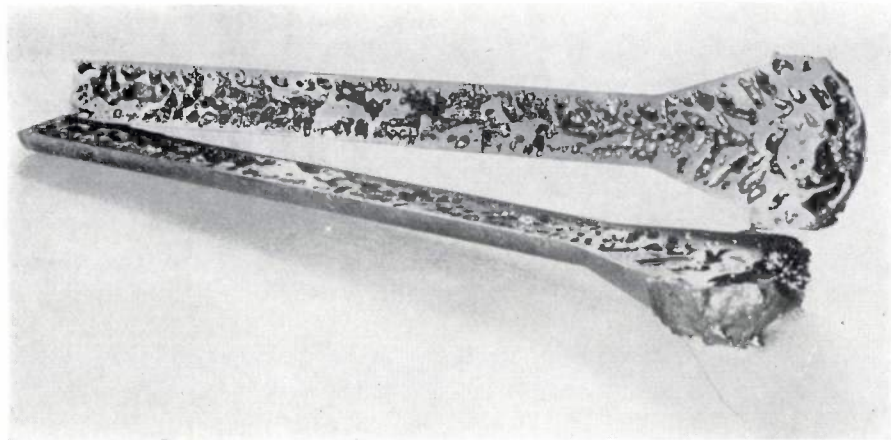
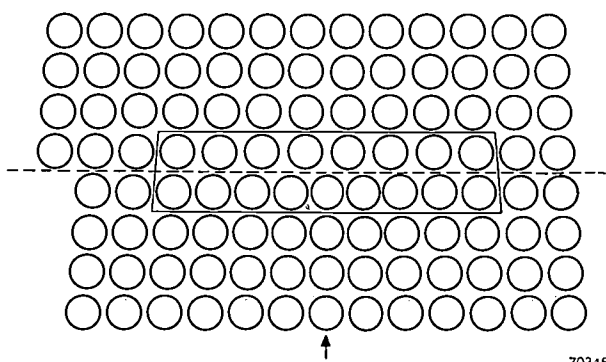


Fig. 2. Sawn-through bar of pure iron, obtained by melting and casting in hydrogen of 1 atm. Owing to the sharp drop in solubility upon solidification, the metal contains numerous gas-filled cavities.



70345

Fig. 3. Schematic representation of an edge dislocation in a simple cubic lattice of metal atoms. The dislocation can be imagined as produced by the forced introduction of an extra plane of atoms (arrow) in the lower half of the crystal. (In reality this situation arises on compression of the lower half of the crystal.) An interstitial atom finds a site in the middle of the outlined area energetically more favourable than a site in an interstice of the unperturbed lattice. The dashed line indicates the slip plane.

atoms, at the grain boundaries and at the ferrite-carbide or ferrite-nitride interfaces. At relatively low temperatures, then, these lattice imperfections will act as "traps" for the H atoms; in other words, the H atoms will spend a very much longer average time at these sites than in the normal interstices. Obviously, this implies a lower diffusion coefficient and a higher solubility. In broad lines, however, the picture is unaltered, since the solubility still decreases with decreasing temperature, and the diffusion rate can still be called exceptionally high.

Interaction of hydrogen with dislocations

The interaction of the interstitial atoms with the dislocations that act as traps can have a marked influence on the mechanical properties of a metal. The interaction of carbon and nitrogen atoms with dislocations in iron has been extensively studied, and we shall briefly discuss this as an introduction to the corresponding interaction with hydrogen, which presents more complications and has not been so thoroughly investigated.

The solute C or N atoms that have diffused to dislocations are unable to leave them at room temperature, so that in the long run, given sufficient atoms, strings of atoms form which extend over the whole length of each dislocation⁶⁾.

The formation of these strings of C or N atoms considerably affects the plastic properties of the metal, since plastic deformation implies the displacement of dislocations. Before they can be dis-

placed they must be detached from the strings of atoms, and this requires an extra stress which, once the dislocations have broken away, is no longer needed. In this way one can understand the occurrence of an upper and a lower yield point in the stress-strain curve of mild steel⁷⁾ (fig. 4a). It also explains why immediately after a slight plastic deformation — that is when the dislocations are still free — no distinct yield point is to be found (fig. 4b). If the metal is then left alone for some considerable time, the C and N atoms again diffuse to the dislocations, so that the upper and lower yield points return and the metal becomes more difficult to deform (harder and more brittle). This spontaneous process is known as strain ageing. It is primarily due to nitrogen atoms, the solubility of nitrogen in iron being much greater than that of carbon.

Further data on the interaction of dislocations with carbon or nitrogen atoms in iron have been derived from *internal friction* measurements. Whereas the C or N atoms in the normal interstices of the iron lattice produce the familiar damping peak, with which the name of Snoek⁸⁾ is associated, the jumps of the atoms located in the stress field of the dislocations give rise to an "abnormal" damping peak at much higher temperatures.

To make this clearer — although we cannot go into details here and must be content with referring to the literature⁸⁾ — it may be mentioned that the

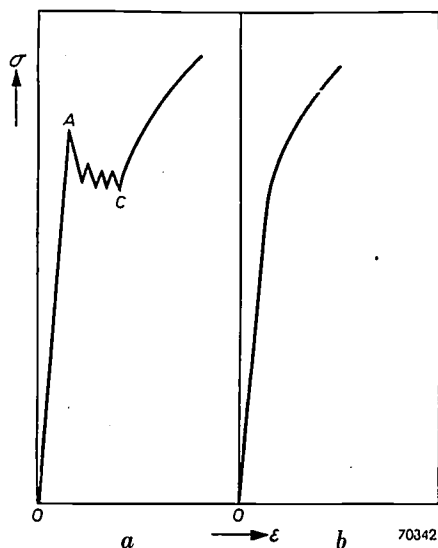


Fig. 4. Stress-strain curves of mild steel. The tensile stress σ is plotted versus the deformation ϵ . a) The steel shows a sharp upper yield point at A. b) After plastic deformation to C and upon renewed loading, the metal shows no distinct yield point.

⁷⁾ A. H. Cottrell, *Dislocations and plastic flow in crystals*, Clarendon Press, Oxford 1953.

⁸⁾ J. L. Snoek, *Physica* 8, 711, 1941 and 9, 862, 1942. See also J. D. Fast and L. J. Dijkstra, *Philips tech. Rev.* 13, 172, 1951/52.

⁶⁾ In body-centred cubic metals, interstitial atoms cause both volume expansion and tetragonal deformation. Consequently they can be taken up in these metals, with an energy gain, in screw as well as in edge dislocations.

position of the maximum of a damping peak of this kind, caused by the jumps of interstitial atoms, corresponds to a specific value of the diffusion coefficient of these atoms. If the maximum occurs at a higher temperature, this means that the value in question is reached at a higher temperature, which in turn implies that the relevant atoms are more strongly bound. That the stronger binding of the atoms is indeed the consequence of their presence in dislocations appears from the fact that the abnormal peak is found only after iron containing nitrogen or carbon has been cold-worked, and thus possesses a relatively high dislocation density. The peak is higher the greater is the degree of plastic deformation, provided the metal contains an excess of C or N. Fig. 5 shows such a peak measured on iron containing nitrogen, and fig. 6 a corresponding peak measured on iron containing carbon.

Conflicting answers have been given to the question whether hydrogen is also one of the elements that can cause internal friction in iron. Gensamer and co-workers⁹⁾ gave an affirmative answer based on damping measurements on a commercial steel charged with hydrogen. Heller¹⁰⁾ also found a damping peak in iron containing hydrogen. He charged wires of fairly pure iron with hydrogen and with deuterium, and at a frequency of 1 c/s found a damping peak at 30 °K in the hydrogen-charged wires, and one at

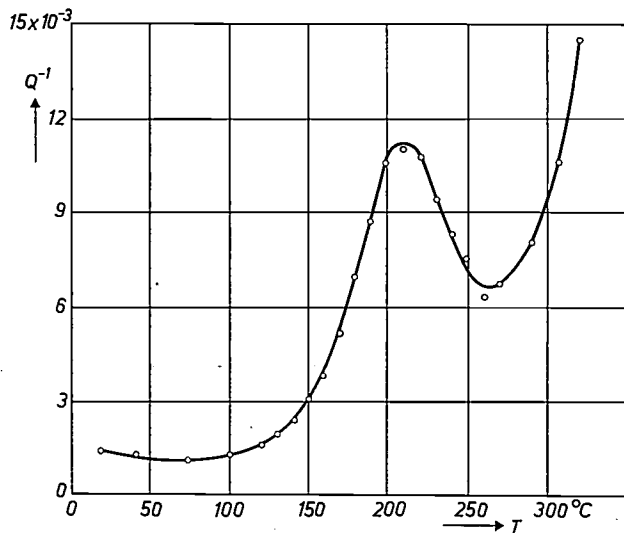


Fig. 5. Abnormal damping peak measured on iron wire after charging with nitrogen in hydrogen containing ammonia at 550 °C, cold-drawing to 60% reduction of cross-sectional area and heating for one hour at 350 °C. The quantity $1/Q$ on the ordinate is the logarithmic decrement of the torsional oscillations divided by π . Vibration frequency 0.13 c/s. (After T. S. Kê, Trans. AIME 176, 448, 1948.)

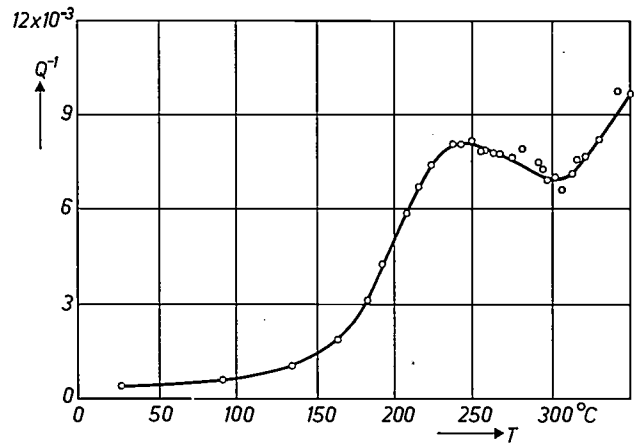


Fig. 6. Abnormal damping peak measured on iron wire charged with carbon in hydrogen containing heptane, cold-drawn to 25% reduction of cross-sectional area and heated for several hours at 250 °C. $1/Q$ is along the ordinate. Vibration frequency 2.2 c/s. (After K. Kamber, D. Keefer and C. Wert, Acta metallurgica 9, 403, 1961, No. 5.)

35 °K in the deuterium-charged wires. He too attributes the peaks to "Snoek jumps" of solute H and D atoms, similar to those of C and N atoms.

After deformation of their hydrogen-containing steel, Weiner and Gensamer⁹⁾ found a damping peak at the much higher temperature of about 105 °K. It seems obvious to assume that the explanation for this must resemble that given for the abnormal peaks caused by carbon and nitrogen; the peak at 105 °K is thus attributed to hydrogen in the dislocations.

The abnormal hydrogen peak, however, is associated with some remarkable effects that are not found in the behaviour of the abnormal nitrogen and carbon peaks. In the first place iron containing hydrogen shows such a peak *also without deliberate plastic deformation*, the peak appearing after several days of ageing at 300 °K. Secondly, the peak disappears spontaneously if the ageing is continued long enough at 300 °K (fig. 7).

Formation of molecular hydrogen in microcavities in iron and steel

The disparate behaviour of iron containing hydrogen as opposed to that containing nitrogen and carbon can be understood as follows. The interstitially dissolved H has the specific possibility of *precipitating in molecular form in lattice imperfections*. Although this property has long been known, its importance until recently was not sufficiently recognized. True, nitrogen and carbon can also precipitate in iron as a separate phase — N as the nitride Fe_3N or Fe_4N , and C as the carbide Fe_3C . However, this nitride or carbide formation does not reduce the free energy of the iron to the same extent as when the N and C atoms bind themselves to dislocations. This

⁹⁾ L. C. Chang and M. Gensamer, Acta metallurgica 1, 483, 1953;
L. C. Weiner and M. Gensamer, Acta metallurgica 5, 692, 1957.
¹⁰⁾ W. R. Heller, Acta metallurgica 9, 600, 1961 (No. 6).

appears, for example, from the experience that strain ageing occurs not only in iron containing solute N or C, but also in iron in which nitrogen or carbon is exclusively present in the form of nitride or carbide. In this case N or C atoms diffuse from the precipitate to the dislocations, implying the entire or partial solution of the precipitate.

As regards hydrogen the opposite is the case: the free energy of the iron falls more when the hydrogen precipitates in the form of H_2 than when it binds itself to dislocations. As a result, H_2 molecules form in all microcavities and lattice imperfections where there is space for them. This can give rise to high local gas pressures, causing plastic deformation of the surrounding metal. The new dislocations produced in this way will absorb part of the hydrogen still in solution, after which the abnormal H peak can appear. Since, however, the free energy decreases even more when H_2 is formed, the abnormal peak finally disappears again, as demonstrated in fig. 7.

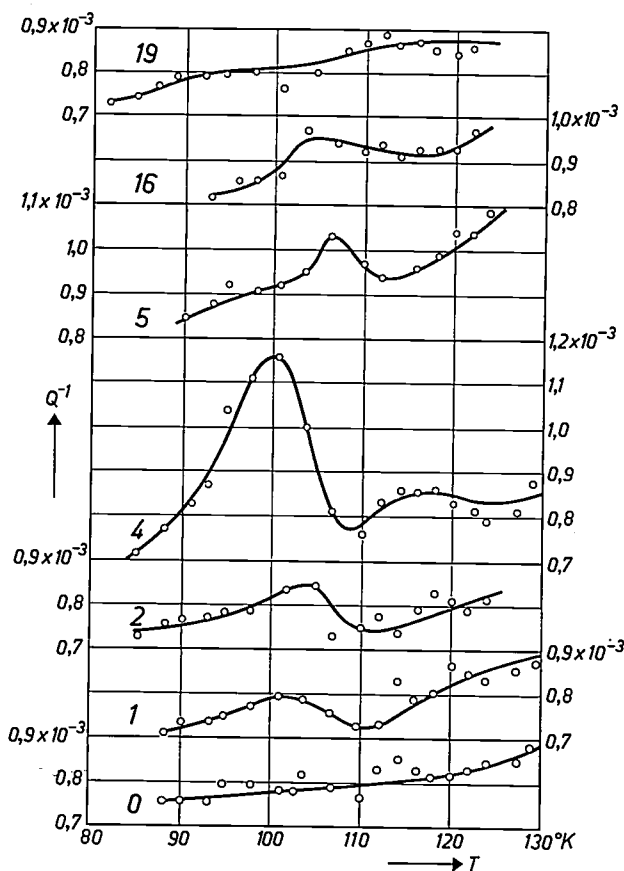


Fig. 7. Internal friction of hydrogen-charged steel, after ageing for various periods of time (vibration frequency 20 c/s). The figures in the graph denote the number of days during which the steel was aged at 300 °K. (After Weiner and Gensamer⁹.) After four days of ageing a marked damping peak is observed at about 100 °K, which vanishes again after further ageing. Effects of this kind, which are not found after nitrogen or carbon charging, can be explained by the formation of hydrogen at high pressure in lattice imperfections.

After charging some commercial steels with hydrogen electrolytically, Rogers¹¹) found that for a certain time they showed no distinct yield point. This remarkable effect can be explained in much the same way as the above-mentioned spontaneous appearance of the abnormal H peak after charging with hydrogen: many new dislocations produced during and after charging, as a result of the plastic deformation upon the formation of H_2 , are not yet anchored by N or C atoms and therefore require no extra stress to set them in motion.

Also connected with the precipitation of hydrogen in microcavities is the experience that, at low temperatures, the (apparent) solubility values found are much greater than might be expected from an extrapolation of the measurements at high temperatures. The difference is too marked to be explained entirely from the binding of H atoms (or protons) in traps at low temperatures. The difference can only be understood from the fact that the equilibrium pressure of the molecular hydrogen steadily increases as the temperature drops (see below), and that the hydrogen, once it has precipitated as H_2 in microcavities, is unable to escape at low temperatures. The hydrogen thus occluded is of course not dissolved.

In view of the high diffusion rate of hydrogen in iron, this possibility of hydrogen occlusion is rather unexpected. At temperatures above about 200 °C the hydrogen will leave the iron rapidly as expected, but at lower temperatures, even though the diffusion rate is still very high, there is no longer any question of the hydrogen actually escaping. (If it were otherwise, hydrogen could not be stored in iron cylinders!) The explanation is that the escape from the iron involves not only diffusion but also surface reactions, in particular the dissociation of molecular hydrogen into atomic hydrogen (at the surface of the cavities) and the converse reaction (on the outside surface). The first reaction in particular is extremely slow at low temperatures¹²). Further particulars will be found in fig. 8.

The precipitation of H_2 in lattice imperfections is also at the root of the numerous detrimental effects caused by hydrogen in steel. It is so important to the understanding of the action of hydrogen in steel that the remainder of this article will be devoted entirely to H_2 formation and its consequences.

If a piece of iron is enveloped at e.g. 1100 °C by H_2 at a pressure of 1 atm, it can be seen from fig. 1 that 0.0006 g of hydrogen (equivalent to 7 cm³ H_2 of 0°C and 1 atm) will be dissolved per hundred grammes of iron. When such a piece of iron is rapidly cooled to 20 °C, the hydrogen content imme-

¹¹) H. C. Rogers, Acta metallurgica 4, 114, 1956 and Trans. AIME 215, 666, 1959.

¹²) If one measures at low temperatures the rate at which hydrogen permeates through an iron wall, or escapes from iron, and the surface reactions are not taken into account, then the diffusion coefficient calculated from the results will be much too small, sometimes even 10 000 times too small. For a detailed discussion of this subject, see J. D. Fast, Philips tech. Rev. 6, 365, 1941, and 7, 74, 1942.

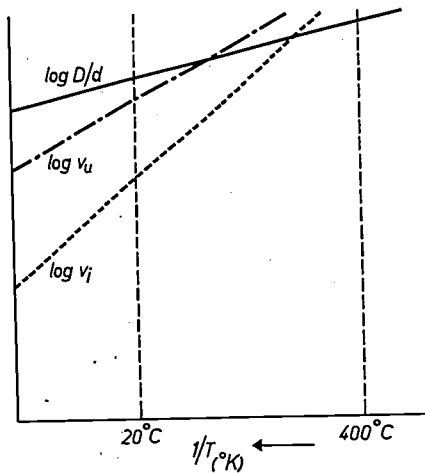


Fig. 8. The figure shows schematically that in the iron-hydrogen system the processes of sorption, desorption and diffusion of the gas have very different temperature coefficients. The diagram shows versus $1/T$ on a logarithmic scale the diffusion rate D/d (solid line), the desorption rate v_u (dot-dash line) and the sorption rate v_i (dashed line). The position of the latter two lines depends to some extent on the surface state, without however altering the essence of the diagram. It can be seen that the desorption and sorption rates decrease faster with falling temperature than the diffusion rate. At room temperature this gives rise to a situation, unexpected at first sight, where on the one hand the permeability of iron to hydrogen is so low that the gas can safely be stored in iron cylinders, whereas on the other hand the diffusion rate can still be called exceptionally high.

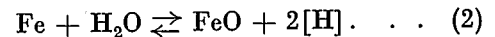
diately after cooling is unchanged, if not near the surface then at least in the interior of the metal. The concentration of dissolved hydrogen is then roughly 10^4 times higher than it would be *in equilibrium* at 20°C and 1 atm hydrogen. The solubility being proportional to the square root of the hydrogen pressure, as we have seen, calculation shows that this high concentration at this temperature can only exist in equilibrium with an H_2 pressure of 10^8 atm. The dissolved hydrogen will therefore attempt to escape from the lattice by diffusion to the outside and to every internal cavity present in the metal.

This calculation still needs a correction: we have not taken into account that at pressures as high as 10^8 atm the hydrogen no longer behaves as a perfect gas, so that the above-mentioned proportionality is no longer valid. If we allow for this in the calculation, we find that the relation between the equilibrium pressure and the H content at 20°C is given by a curve as shown in fig. 9. From this we see that the equilibrium pressure in the case considered is not 10^8 but "only" 10^4 atm.

By exceeding the cohesion of the material, an equilibrium pressure as high as this can easily give rise to ruptures, but only where imperfections exist. In an ideal single crystal of iron even much higher hydrogen concentrations than those mentioned would not lead to rupture. The harmful effects appear only if imperfections are present in the metal (or are

introduced by plastic deformation), in which the precipitation of H_2 up to the dangerous pressure can really take place.

Most commercial steels have hydrogen contents of the order mentioned in the above example, often even higher. The example differs from reality only so far as most of the hydrogen in steel does not originate from hydrogen gas in the atmosphere. One of the principal sources of hydrogen in iron is water vapour, which at high temperatures reacts with both solid and molten iron as follows:



The water vapour may, for example, come from rust on the scrap used in steel-making, from constituents in the slag (lime), from the crucible or furnace wall and from the gas atmosphere present. In electric arc welding with coated electrodes, the coating is the main source of water vapour.

A second and equally important source of hydrogen in iron and steel is the hydrogen produced in atomic form at the surface during galvanizing, pickling or electrolysis processes.

The fact that high H_2 pressures can arise in steel containing hydrogen may be demonstrated by simple experiments ¹³⁾ carried out as long

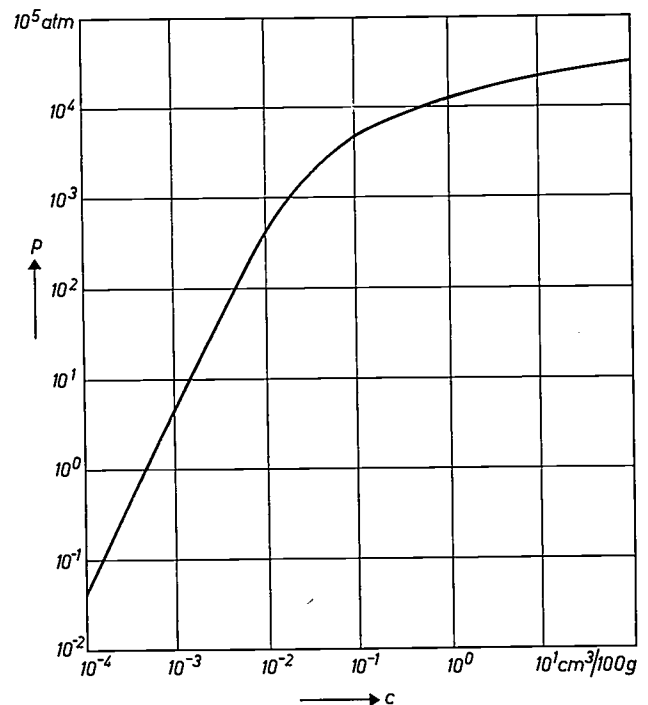


Fig. 9. The graph shows the calculated equilibrium pressure p of molecular hydrogen, taking into account the deviations from ideal behaviour, plotted as a function of the concentration c of atomic hydrogen in iron for a temperature of 20°C . (After G. Vibrans, Arch. Eisenhüttenw. 32, 667, 1961.)

¹³⁾ C. A. Edwards, J. Iron and Steel Inst. 110, 9, 1924; P. Bardenheuer and G. Thanheiser, Mitt. K. W. I. Eisenforsch. 10, 323, 1928.

ago as 1924. The two following experiments are typical. 1) Hydrogen is evolved electrolytically on the surface of a hollow iron cylinder (of small internal volume), enabling the gas in atomic form to penetrate the metal in relatively large quantities. After some time the pressure inside the bore of the cylinder starts to rise, as can be read from a pressure gauge.

In this way the pressure can be observed to build up to several hundreds of atmospheres, after which the experiment is ended

for reasons of safety. 2) By means of an acid or by electrolysis, atomic hydrogen is evolved on the inside of an open iron pan which is enamelled on the outside. After some time the enamel cracks away from the pan — sometimes almost explosively — owing to hydrogen having precipitated at high pressure at the iron-enamel interface.

Difficulties connected with effects of this nature sometimes occur in enamelling processes as a result of the above-mentioned reaction (2) between the steel and the water vapour. The vapour arises from substances used in enamelling, principally from "frit"¹⁴. After cooling, the metal may then be severely supersaturated with hydrogen, which may result in serious damage to the enamel layer.

Closely related to the precipitation of hydrogen under high pressure at the metal-enamel interface is its precipitation at non-metallic inclusions in steel. This may occur, for example, as a consequence of the pickling prior to tin or zinc plating. Part of the atomic hydrogen formed during the pickling process diffuses into the steel and forms H₂ at the inclusions. This can lead to surface blistering if the inclusions are so close to the surface that the hydrogen pressure is able to push the supervening metal outwards by plastic deformation (*fig. 10*).

Formation of high-pressure CH₄

Where iron and steel containing carbon are in external contact with hydrogen gas, methane (CH₄) may be formed internally by the reaction:



¹⁴ D. G. Moore, M. A. Mason and W. N. Harrison, *J. Amer. Ceramic Soc.* **35**, 33, 1952.



Fig. 10. Blisters formed during pickling of steel. In the pickling process atomic hydrogen is produced which diffuses inwards and forms molecular hydrogen at inclusions in the steel. The blisters appear where these inclusions are immediately below the surface. Magnification 10×. The regular pattern of the blisters indicates a regularity of the inclusions, brought about during the rolling of the steel.

At 300 °C the equilibrium constant of this reaction has a value such that a CH₄ pressure of several thousand atmospheres corresponds to an H₂ pressure of 1 atm. At higher temperatures the equilibrium pressure of CH₄ is lower, at lower temperatures higher.

In technology this effect caused many initial difficulties in the large-scale production of various inorganic and organic chemical compounds, e.g. of ammonia, methanol and petrol. Nowadays it is possible to prevent CH₄ forming by using alloying additives in the steel that form highly stable carbides with the carbon.

It is also worth noting that, as reaction (3) is exothermic, the occurrence of this reaction may create the impression that hydrogen in a certain temperature range may dissolve exothermally in the metal instead of endothermally.

In Part II of this article we shall show that the formation of molecular hydrogen under high pressure as discussed here, not only takes place in *large* imperfections (microcavities) but also in defects of *atomic* dimensions. It will be seen that knowledge of this effect makes it possible to understand many of the other harmful effects of hydrogen in steel. Among the more important of these are the reduction of ductility and the appearance of brittle fractures.

Summary. Whereas hydrogen dissolves interstitially in iron in the form of *atoms* and presumably diffuses in the form of *protons, molecules* (H₂) can form in imperfections of the crystal lattice. The latter explains such unexpected effects as the temporary appearance of an "abnormal" damping peak during ageing, and the temporary absence of a distinct yield point after electrolytic charging with hydrogen. The formation of molecular hydrogen in microcavities in iron and steel can give rise to extremely high pressures; this is demonstrated by calculations, by experiments and in practice. In Part II of this article an explanation will be given, on this basis, of the more serious harmful effects of hydrogen in steel, in particular reduced ductility and brittle fracture.

THERMIONIC-CATHODE TESTING

Photo Maurice Broomfield

The emission properties of thermionic cathodes are usually investigated in simple test-valves, containing only an anode besides the cathode. The photograph shows three such valves being evacuated on a mercury-diffusion pump. Behind them is the cold trap, which is kept at about -190°C by liquid nitrogen and condenses the mercury vapour from the diffusion pump. Above the test valves can be seen an ionization gauge¹⁾. In situations like this it is necessary to allow for the fact that the pressure in the ion gauge is not the same as in the valves.

These preliminary tests under greatly simplified conditions are followed at a later stage by tests with the cathodes mounted in the valves for which they are intended.

¹⁾ Philips tech. Rev. 20, 153, 1958/59.

ABSTRACTS OF RECENT SCIENTIFIC PUBLICATIONS BY THE STAFF OF N.V. PHILIPS' GLOEILAMPENFABRIEKEN

Reprints of those papers not marked with an asterisk * can be obtained free of charge upon application to the Philips Research Laboratories, Eindhoven, Netherlands, where a limited number of reprints are available for distribution.

- 3029:** W. L. Wanmaker and C. Bakker: The determination of luminescent properties as a method for the study of diffusion processes in the solid state (Reactivity of solids, Proc. 4th int. Symp. on the reactivity of solids, Amsterdam 1960, editors J. H. de Boer *et al.*, 709-717, Elsevier, Amsterdam 1961).

Luminescence is only produced by a foreign ion (activator ion) in a crystal lattice. If a crystal is first made without activator and then heated in contact with a suitable compound containing the activator ion, the diffusion of the ion in question through the host lattice can be followed by means of the advance of the luminescence. This method has a very wide field of application. A start has been made with the study of the diffusion of Sb and Mn ions in calcium halophosphate, $3\text{Ca}_3(\text{PO}_4)_2 \cdot \text{Ca}(\text{F}, \text{Cl})_2$, which has the apatite structure.

This method can also be used to study the diffusion of ions of the host crystal. In order to do this, the luminescent substance is heated in contact with a suitable compound of the ion in question, e.g. CaCO_3 . Since the intensity of the luminescence is lower when there is an excess of Ca ions, the diffusion can be followed by observing the decrease in luminescence.

- 3030:** H. G. Grimmeiss and H. Koelmans: Analysis of *P-N* luminescence in Zn-doped GaP (Phys. Rev. **123**, pp. 1939-1947, 1961, No. 6).

A *P-N* junction in a semiconductor can emit light when a voltage is applied across it. This article describes one phase of an investigation aimed at increasing the yield of light from this effect. See also **R 398** and Philips tech. Rev. **22**, 360-361, 1960/61.

- 3031:** N. V. Franssen: Eigenschaften des natürlichen Richtungshörens und ihre Anwendung auf die Stereophonie (Proc. 3rd int. congress on acoustics, Stuttgart, 1959, editor L. Cremer, Vol. II, pp. 788-790, Elsevier, Amsterdam 1961). (The properties of natural directional hearing and their application to stereophony; in German.)

A short survey of material which has been fully dealt with in **2889**.

- 3032:** C. M. van der Burgt: Neue keramische Ultraschallwandler und deren Kopplung an die Flüssigkeit (as **3031**, pp. 1219-1221). (New ceramic ultrasonic transducers and their coupling to liquids; in German.)

Ni-Cu-Co ferrites can be used for the efficient production of ultrasonic vibrations. The author discusses the demands these ferrites must meet when high-power vibrations of a frequency between 20 and 50 kc/s are concerned. See also **2902**.

- 3033:** A. Recourt and G. H. F. de Vries: An experimental apparatus for contact microradiography at 200-500 V (Nature **191**, 1185-1186, 1961, No. 4794).

*In contact microradiography (see Philips tech. Rev. **19**, 221-233, 1957/58), the X-ray beam is considerably attenuated on passage through the window. This effect is so strong at low voltages (less than 500 V), whose use is sometimes indicated, that normal contact microradiography equipment simply cannot be used for this purpose. The authors have built an experimental apparatus in which the X-ray source, the sample and the film are all in the same evacuated space, so that the window is completely dispensed with.

- 3034:** J. D. Fast, J. L. Meijering and M. B. Verrijp: Frottement interne dans les alliages ferritiques Fe-Mn-N (Métaux, Corr., Ind. **36**, 112-114, 1961, No. 427). (Internal friction in iron alloys Fe-Mn-N; in French.)

Considerations concerning the relative positions of the three peaks mentioned in the following publication (**3035**). It is concluded from experiments on alloys with varying Mn contents that the main mechanism is based on the presence of Mn-Mn pairs. An N atom is more strongly bound to such a pair than to a single Mn atom, and yet is more mobile in the former case.

3035: J. L. Meijering: Considérations sur l'effet Snoek dans le cas de sites non-équivalents pour les atomes en insertion (Métaux, Corr., Ind. **36**, 107-111, 1961, No. 427). (Considerations on the Snoek effect for the case of non-equivalent sites for the foreign atoms; in French.)

The maximum Snoek damping (see Philips tech. Rev. **13**, 172-179, 1951/52) due to interstitial N atoms in iron is at about 25 °C at a frequency of one c/s. By substituting some Fe atoms by Mn atoms, the peak in the damping curve is broadened and shifted to higher temperatures. Three individual peaks can be distinguished in this broadened peak, the middle one corresponding to damping in the absence of Mn. An N atom in the neighbourhood of an Mn atom will alternately occupy two kinds of octahedral sites, corresponding to two different free energies. This explains the two additional peaks. A somewhat simpler case is also treated, where the two kinds of positions are the octahedral and tetrahedral sites in a lattice containing no Mn.

R 426: J. D. Wasscher: Note on four-point resistivity measurements on anisotropic conductors (Philips Res. Repts **16**, 301-306, 1961, No. 4).

The resistivity of a homogeneous material can be determined with the aid of a slice of the material to which four point contacts are applied. If a current is passed through two of these contacts, a potential difference proportional to this current is produced between the other two contacts. The constant of proportionality, which is a measure of the resistivity, depends on the geometry of the arrangement. The author shows how the results obtained for isotropic material can be applied to anisotropic material, making use of a coordinate transformation proposed by Van der Pauw. The case of four contacts in a straight line is discussed, as is the case of four contacts at the corners of a square; both cases are worked out for two samples, one thin and one thick compared to the distance between the contacts. The square arrangement proves to be the most sensitive to anisotropy. It is also found that the three principal values of the resistivity tensor can be determined from measurements on one single plane at right angles to a direction corresponding to one of these principal values. Corrections for the

finite dimensions of the contacts and the sample are discussed.

R 427: S. Duinker: Short-wavelength response of magnetic reproducing heads with rounded gap edges (Philips Res. Repts **16**, 307-322, 1961, No. 4).

The calculation of the response of magnetic pick-up heads given by Westmijze for a gap with ideal rectangular edges is extended by the author to the general case of edges with a finite radius of curvature. An expression for the gap-width loss factor is found by conformal mappings of head and potential field. The ratio of the radius of curvature of the edge to the gap width occurs as a parameter in this expression; it follows that as the radius of curvature increases the response of the head decreases and the effective gap width increases.

R 428: J. J. Scheer and J. van Laar: Photo-electric emission from cadmium telluride (Philips Res. Repts **16**, 323-328, 1961, No. 4).

The photo-electric effect of CdTe is measured on crystal surfaces and on evaporated layers. Clean crystal surfaces are prepared by cleaving a crystal under high vacuum. Experimental curves of the photo-current as a function of photon energy show a "tail" at long wavelengths; this tail is ascribed to impurities. It follows from the experimental results that the most reliable value of the work function is found from the emission from the surface of a cleaved single crystal. Finally, the possibility of a relationship between the photo-electric properties of CdTe and CdS is discussed.

R 429: W. Albers and K. Schol: The *P-T-X* phase diagram of the system Sn-S (Philips Res. Repts **16**, 329-342, 1961, No. 4).

Two maxima were found in the *T-X_L* diagram during a study of the equilibria between the solid, liquid and gaseous phases in the system Sn-S. The first maximum occurs at the composition of SnS, and at a melting temperature of 881.5 ± 2 °C and a partial pressure of sulphur of 0.033 atm. The second maximum occurs at the composition of SnS₂, and at a temperature of 870 °C and a pressure of 40 atm. Between 10 and 47 at. % S, and perhaps between 70 and 90 at. % S, the different liquid phases are immiscible. In the course of this investigation, a compound has been discovered which probably has the composition Sn₃S₄.

Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES

THE ALLOY-DIFFUSION TECHNIQUE FOR MANUFACTURING HIGH-FREQUENCY TRANSISTORS

by P. J. W. JOCHEMS *).

62L.382.333

The highest frequency at which a transistor can provide a reasonable signal gain depends amongst other things on its dimensions; the higher the frequency, the smaller the transistor has to be. The manufacture of transistors for operation at 100 Mc/s necessitates strict control of dimensions which are of the order of tens of microns and even of microns in the case of base width. Simplicity and reliability are conditions for success in the mass-production of transistors of this kind. The requirements are very largely satisfied by the alloy-diffusion technique, which Philips have now been using for some years in the manufacture of HF transistors.

In connection with microminiaturization, the technique has also been used on an experimental scale for making an almost complete three-stage amplifier out of a single germanium crystal.

The invention of the transistor in 1948 in the Bell Laboratories¹⁾ initiated a tremendous surge of development in the semiconductor field. The first transistors were of the *point-contact* type, and consisted of a thin plate or wafer of crystalline germanium with two wires touching it at points about 50 μm apart (fig. 1). One wire acted as emitter, the other as collector. The name "base" was given

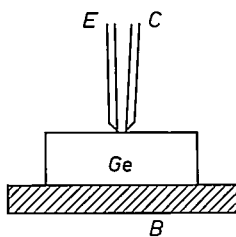


Fig. 1. Schematic representation of a point-contact transistor. Ge germanium crystal. E emitter. C collector. B base contact.

to the crystal itself. Attempts to account for the action of the point-contact transistor (even now, a completely satisfactory explanation has still to be found) brought certain facts to light that formed a

foundation for the theory of *P-N* junctions published in 1948 by Shockley — also of the Bell Laboratories²⁾. In the same paper Shockley described the structure and properties of the *junction* type of transistor. Transistors made in accordance with these theoretical considerations were found to possess the properties predicted for them. It is fortunate that the theory of the junction transistor is well understood³⁾, for this knowledge reveals what possibilities there are of giving a transistor the properties that are desirable in a given application.

Almost all transistors being made at the present time are of the junction type. They consist of a small block, disc or rod of monocrystalline material, generally germanium or silicon, containing the requisite *P* and *N* regions. To perform the various functions allotted to them in radio and television receivers, switching circuits, computers and so on, they have to meet all sorts of requirements as to frequency range, output power, etc. These require-

*) Philips Research Laboratories, Eindhoven.

¹⁾ J. Bardeen and W. H. Brattain, The transistor, a semiconductor triode, Phys. Rev. 74, 230-231, 1948.

²⁾ W. Shockley, The theory of *P-N* junctions in semiconductors and *P-N* junction transistors, Bell Syst. tech. J. 28, 435-489, 1949.

³⁾ See for example F. H. Stieltjes and L. J. Tummers, Simple theory of the junction transistor, Philips tech. Rev. 17, 233-246, 1955/56. The influence of frequency is not dealt with in the article.

ments lead to differences in the dimensions of the transistor and in certain characteristics of the P and N regions of which it is made up, notably their specific resistance, and also in the manner of mounting the crystal. For example, transistors designed to handle high powers are invariably mounted in such a way as to ensure the closest possible thermal contact with the surroundings.

A great deal of effort has been devoted to finding and developing methods which allow P - N junctions to be made at the desired places in monocrystalline material, and which at the same time enable the characteristics and the geometry of the P and N regions to be kept under strict control. In the present article we shall start by discussing some of these methods. They were originally developed for making transistors operating at frequencies below 10 Mc/s. To render them suitable for making transistors that would work satisfactorily at still higher frequencies, these methods have had to undergo all kinds of refinements. In this article we shall deal in some detail with one of the more elaborate techniques. It is a version, developed in Philips Research Laboratories, of the "alloy-diffusion" technique and which for some years now has been employed for mass-producing various types of transistors to operate in a frequency range extending up to 200 Mc/s. These types have become known as "pushed-out base" transistors, for reasons that will shortly become evident.

The alloy-diffusion technique is also suitable for making what are known as "solid-state circuits", a form of microminiaturization. The experimental construction of a three-stage solid-state amplifier will be described by way of example.

Methods of making P - N junctions and junction transistors⁴⁾

The addition of impurities in crystal-pulling

The earlier methods entailed forming the P - N transitions while the monocrystal was being pulled. In one such method the starting material is a germanium or silicon melt in which donor material (Sb, As or P) has been dissolved; naturally, a crystal pulled from this melt will be of type N . However, at a given instant enough acceptor material (Ga, In or B) is added to the melt to make the acceptor concentration greater than the donor concentration. Accordingly, the next section of crystal to grow is of P type.

⁴⁾ For a fuller review, containing many references to the literature, see the article by W. C. Dunlap, *Methods of preparing P - N junctions*, which appears as Chapter 7 in L. P. Hunter, *Handbook of semiconductor electronics*, 2nd impression, McGraw-Hill, New York 1962.

In the making of transistors by this method, the P material is only allowed to grow to the point necessary to give a base of the desired thickness. The melt is then dosed with a sufficient amount of donor to start growing N -type material again. A large number of transistors can be obtained by sawing up the resulting crystal (*fig. 2*). Leads still have

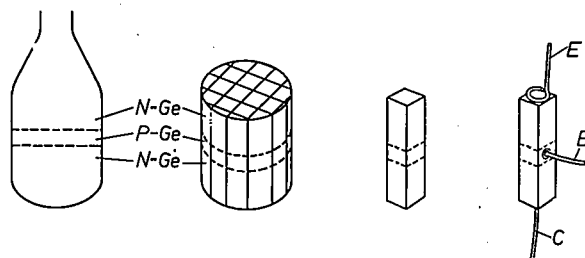


Fig. 2. Diagrams illustrating the manufacture of "grown" transistors. A P - N and an N - P junction have been created in a large monocrystal while it was being pulled from the melt (left). The central portion of the crystal, containing the junctions, can be sawn up into a large number of transistors. To these, electrode leads E , B and C have to be attached. It is no easy matter to attach lead B to the thin base layer (extreme right).

to be fitted to emitter, base and collector of such a "grown" junction transistor, and it is no easy matter to attach a wire to the base layer, which may be only $40\ \mu\text{m}$ thick. For one thing, the neighbouring P - N junctions are easily damaged.

Transistor manufacture took a big step forward when methods were developed for making P - N junctions *after* the crystal had been grown. It now became possible to saw up a large homogeneous block into crystals of suitable dimensions, and to make the required P - N junctions in these smaller units.

The alloying method

The widely used alloying method is one way of producing P - N junctions in a homogeneous germanium crystal. The crystal can either be of P -type or of N -type germanium. If it is of N -type, a pellet of acceptor material — usually indium — is placed upon it, and the whole is heated to 500 or $600\ ^\circ\text{C}$. Having a melting point of about $150\ ^\circ\text{C}$, the indium liquefies while the germanium (melting point about $950\ ^\circ\text{C}$) remains in the solid state. The molten indium takes up germanium into solution and consequently eats into the solid material (see the hatched area in *fig. 3*). The amount of germanium taken up is obviously proportional to the amount of indium and to the solubility of germanium in liquid indium. This solubility, which increases with temperature, can be read from the Ge-In phase diagram. The molten indium stops advancing as soon as it has become saturated with germanium. Various measures can be taken to

ensure that its advance will be even and regular (i.e. over a plane front, see fig. 3), one being to apply the indium to a (111) face of the germanium crystal. The area in contact with the indium pellet can be restricted with the aid of a jig. By a suitable choice of contact area, indium quantity and temperature, one can determine in advance the depth to which the indium will penetrate (h in fig. 3). As the system cools, the solubility of germanium in indium decreases, and the germanium recrystallizes. It is redeposited on the undissolved germanium, refilling the hole in the monocrystal. However, the deposited germanium now contains a little indium, only a tiny residue, it is true, the concentration being of the order of 10^{-4} at. %⁵⁾, but this is enough to make

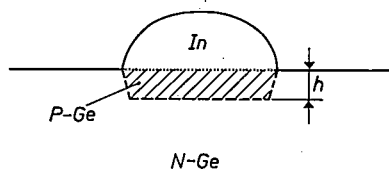


Fig. 3. A P - N junction formed in N germanium by the alloying method using indium (In) as acceptor element. In this and subsequent figures of the same kind, a dashed line represents a P - N junction and a dotted line represents normal electrical contact.

the newly formed recrystallized layer a well conducting P region. A P - N junction has therefore been created.

A transistor can be made by alloying into both sides of the germanium wafer (fig. 4). The base of the transistor is formed by the original N germanium lying between the two P regions. The thickness of the base can be controlled via the indium penetration depths and the thickness of the wafer. It is usual to make the collector region rather more extensive than the emitter region, as otherwise some of the minority charge-carriers injected by the emitter might miss the collector and recombine at

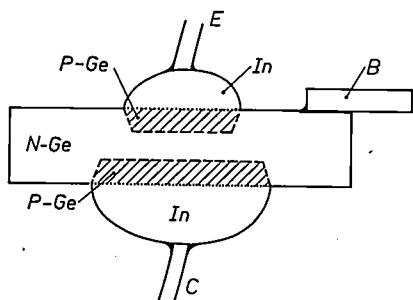


Fig. 4. Section through a transistor made by the alloying method.

the opposite face of the crystal without having made any useful contribution to the action of the transistor.

A great advantage of the alloying method is that attachment of the electrode leads to the indium presents no difficulty. The indium pellets are still present as bulges on the surface of the wafer, offering a convenient means of contact with the emitter and collector regions. It is an easy matter to solder leads to these bulges; the base lead can be attached elsewhere on the crystal. The alloying method is used for many of the transistor types manufactured by Philips.

The diffusion method

A second method of creating P - N junctions in a ready-grown crystal is based on the principle of acceptor diffusion into an N crystal or, alternatively, that of donor diffusion into a P crystal. Generally, donors diffuse through germanium at a much faster rate than acceptors; in silicon it is just the other way about. Accordingly, the starting material usually taken for germanium transistors is P germanium, into which donors are allowed to diffuse. Fig. 5 shows, by way of example, how the donor element antimony diffuses into P germanium held at 780°C , in the presence of antimony vapour in equilibrium with its solid phase at 600°C . The P - N junction occurs at the place where the donor concentration is equal to the acceptor concentration. Under the conditions to which these diagrams relate, an N layer $5.5\ \mu\text{m}$ thick takes a quarter of an hour to form.

Diffusion is a slow process dependent on temperature and time. The accurate control of the diffusion process that is possible, makes this method very suitable for obtaining thin layers reproducibly. Characteristic of such layers is the gradual fall-off in impurity concentration from the crystal surface inwards.

Transistors for high frequencies

The higher the frequency at which a transistor must provide a reasonable gain, the smaller are the values acceptable for the various capacitances in the structure — the diffusion capacitance and the collector- and emitter-barrier capacitances⁶⁾. One way of obtaining small capacitances is to make the transistor smaller. Very small size is accordingly a feature of transistors for high frequencies; for example, a transistor designed to give adequate

⁵⁾ Graphs showing the solubility of various metals in solid germanium and in solid silicon may be found in F. A. Trumbore, Bell Syst. tech. J. 39, 205-233, 1960.

⁶⁾ The factors limiting the frequency range of transistors are dealt with in an article by M. Beun and L. J. Tummers, on transistor behaviour with increasing frequency, which is to appear in this review in the near future.

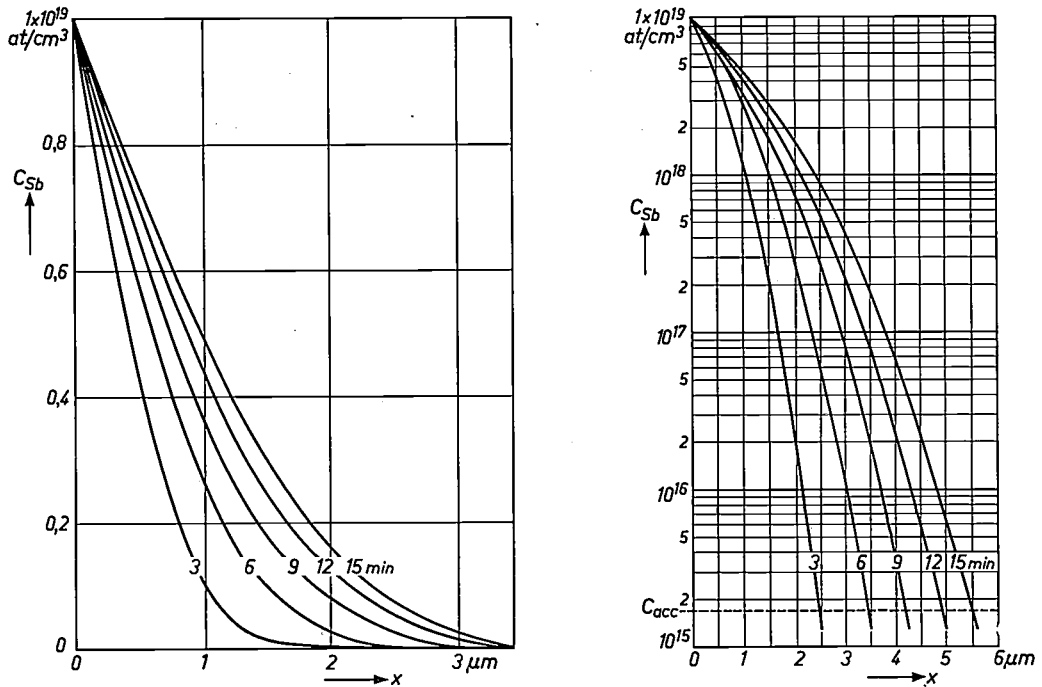


Fig. 5. Diffusion of the donor element Sb into germanium. The germanium is held at 780 °C in Sb vapour that is in equilibrium with solid Sb kept at 600 °C. Under these conditions the Sb concentration at the surface has a value close to the maximum possible at 780 °C, (about 10^{19} Sb atoms per cm^3). The Sb concentration C_{Sb} has been plotted as a function of x , the depth in the germanium crystal, after 3, 6, 9, 12 and 15 minute periods of exposure to the Sb vapour. The concentration scale is linear in the left-hand and logarithmic in the right-hand diagram. The latter reveals that in germanium with an acceptor (e.g. In) concentration of 1.7×10^{16} atoms per cm^3 , corresponding to a resistivity of 2 ohm cm, the N layer attains a thickness of 5.5 μm after 15 minutes.

amplification at frequencies as high as 100 Mc/s must have a base no more than a few microns thick. Since the base current flows *laterally* through the base layer, it will be obvious that a thin base is unfavourable towards the low base resistance that is also a requirement for HF transistors.

Transistors for frequencies above about 10 Mc/s cannot be manufactured satisfactorily by the alloying method described above. It is no easy matter to control the penetration depth of the emitter and collector regions (see fig. 4) so accurately that one can rely on getting a reproducible value of a few microns for the thickness of the N region remaining sandwiched between collector and emitter; deviations from nominal in the wafer thickness and the two penetration depths, will be reflected in the residual base thickness. Another unfavourable aspect of alloying is that the penetration depth would have to be large compared with the base thickness, because wafers thinner than about 50 μm cannot be handled.

But even if these difficulties could be overcome, the designer's difficulties would still not be at an end. In an alloyed transistor the resistivity ρ_B of the base material is always large with respect to that of the collector material. This leads to

conflicting requirements regarding the base resistivity of a transistor to operate at high frequencies. ρ_B has to be large to satisfy certain conditions, and small to satisfy others. Thus a compromise has to be found, but the best possible compromise on ρ_B will still not raise the limit of the usable frequency range beyond about 10 Mc/s.

The way round this is to make the resistivity of the collector material a great deal higher than that of the base. We shall see below that this can be done by using a diffusion technique to produce the base region.

If the resistivity ρ_C of the collector material is much smaller than ρ_B , the collector barrier extends mainly into the material of the base. The thickness of the barrier is proportional to the square root of the voltage across it. Since the base is extremely thin there is a danger of "punch-through" occurring at a relatively low voltage, the collector barrier moving through the base and coming into contact with the emitter barrier. To obviate this danger, and also in order to keep down the base resistance, ρ_B must be made small. But the advantage of so doing is nullified because the smaller ρ_B becomes, the lower will be the voltage at which the barrier breaks down and the higher its capacitance ⁷⁾.

⁷⁾ More light is cast on these questions in the article cited in ⁶⁾. See also the article by B. N. Slade, Device design considerations, which appears as Chapter 10 in L. P. Hunter's book ⁴⁾, and in which a numerical example is worked out.

However, if the material of the collector has a resistivity very much greater than ρ_B , the collector barrier will extend mainly into the collector region. In these circumstances it is possible to choose ρ_C as large as is necessary to raise the collector barrier breakdown voltage, and to lower the collector barrier capacitance to the appropriate safe levels. ρ_B can then be made small enough for an acceptable value of the base resistance to be obtained.

Diffusion and alloying in separate stages

One possible procedure is to take a wafer of *P* germanium whose resistivity makes it suitable collector material, and to allow donors (Sb or As) to diffuse into it, thereby forming a surface layer of *N* germanium (fig. 6a). The *N* layer is etched away at the place where the collector contact is to be made (fig. 6b). The emitter region and the base contact are then produced by the alloying method, on the other side of the crystal (fig. 6c). Two indium pellets can be used for this purpose, one containing gallium in solid solution and the other containing antimony (or arsenic). Alloying is done in the same way as for the alloyed type of transistor. The germanium that recrystallizes on cooling contains impurities additional to indium: gallium is present in the germanium deposited under one pellet, antimony in that deposited under the other. Like indium, gallium is an acceptor element, but the solubility in solid germanium is about one hundred times greater for gallium than for indium. Consequently a strongly doped *P* region, which can serve as emitter, forms under the gallium-containing indium pellet⁸). In the region under the antimony-containing pellet, acceptor atoms (of indium) are overcompensated by donor atoms (of antimony) and this is accordingly an *N* region which will ensure good electrical contact with the *N* layer.

To reduce the capacitance of the collector barrier, as much as possible of the *N* layer is etched away. The result is a crystal of the original *P* material carrying a remnant of the diffused *N* layer and topped by an emitter and a base pellet (figs. 6d and e). Leads are attached to the pellets by soldering.

Regarding the dimensions, it is found that with a base thinner than 3 μm and a collector and an emitter barrier which extend over areas smaller than 0.05 and 0.02 mm^2 respectively, it is possible to obtain a transistor with a good amplifying action at frequencies as high as 100 Mc/s. These requirements as to internal dimensions are satisfied by a transistor whose emitter and base pellets, if equal in size, have the diameter and separation indicated

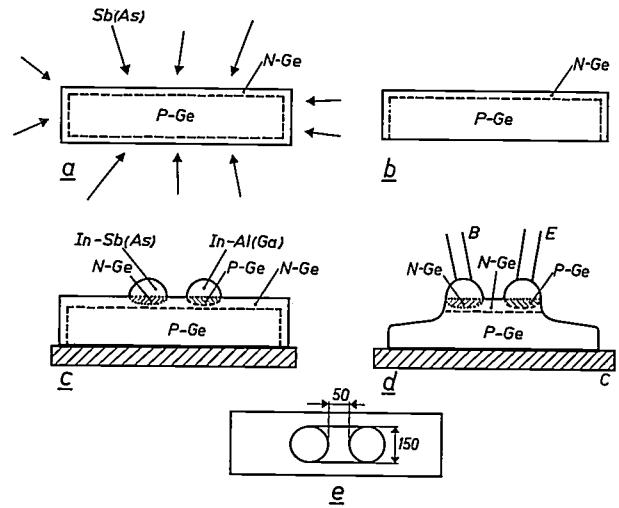


Fig. 6. Manufacture of an HF transistor in separate diffusion and alloying stages. a) Diffusion of a donor (Sb or As) into a *P*-type germanium crystal causes an *N* layer to form on the surface of the crystal. b) The *N* layer is etched away from the underside of the crystal. c) The emitter region (which also serves as emitter contact) and the base contact are made by alloying; a collector contact, which may take the form of a nickel strip, is soldered to the bottom of the crystal. d) Electrode leads *B*, *E* and *C* are then attached, and all the unwanted parts of the diffused *N* layer are removed by etching, leaving an island which constitutes the active zone of the transistor. e) View of the finished transistor from above the dimensions, which are in μm , are appropriate to a type operating at frequencies up to 100 Mc/s.

in fig. 6e. HF transistors manufactured in the manner just described are known as diffused base transistors.

Mention was made of a technological difficulty arising in the manufacture of alloyed transistors, namely the strict control of the alloying process that is necessary if base layers of the required width are to be obtained in a reproducible way. The same difficulty is encountered in the making of diffused base transistors though to a much lesser degree. There is only one reason for uncertainty about the base width, that being the depth to which the alloying front penetrates; the surface *N* layer can be diffused to exactly the prescribed depth. A cross-section through the emitter region, drawn to scale as in fig. 7, reveals how severe are the require-

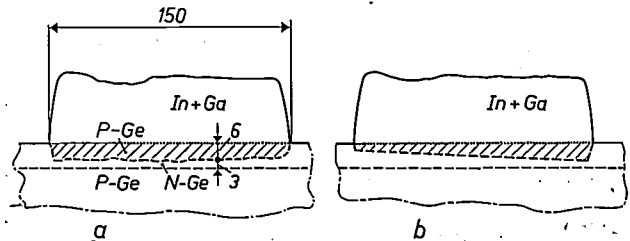


Fig. 7. Cross-section, drawn to scale, of the emitter region of a transistor made by the method illustrated in fig. 6. The drawings reveal the severe requirements to ensure that the alloying front is flat (a), and that this front — a (111) plane — is parallel to the face of the crystal (b).

⁸) F. H. Stieltjes and L. J. Tummers, Behaviour of the transistor at high current densities, Philips tech. Rev. 18, 61-68, 1956/57, in particular p. 67 and 68.

ments that have to be satisfied by the alloying front, not only regarding depth of penetration but also the flatness and parallel alignment with the underlying *P-N* junction. Failure to meet these requirements causes a high reject rate in mass production.

Mesa transistors

The difficulties just mentioned are largely due to the fact that the distance over which the alloying front travels is several times greater than the required residual base width. With normal alloying methods this is inevitable: the alloying must be carried out at a temperature considerably above the melting point of indium in order to ensure adequate wetting of the germanium by the liquid indium. But at this relatively high temperature so much germanium becomes dissolved in the indium that an appreciable depth of penetration is inevitable. There is, however, a technique available in the shape of vacuum deposition of the alloying materials. The emitter can be made by depositing aluminium and the base contact by depositing gold containing a certain amount of antimony. Very small quantities of material can be transferred in this way, so that the depth to which the alloy fronts penetrate are restricted to a fraction of the residual base width. Fluctuations in the alloying front likewise diminish and thus cease to be dangerous. Devices made by this method are called mesa transistors⁹⁾. The emitter and base leads are attached to the alloy films by a special technique known as thermocompression bonding.

The p.o.b. transistor

It has proved possible to get around the above difficulties by adopting what is known as the alloy-diffusion technique¹⁰⁾. In the resulting p.o.b. transistor, the typical transistor layer structure is formed in a single process, contact areas for easy attachment of electrode leads being produced at the same time. The method lends itself extremely

well to mass production, and many millions of p.o.b. transistors for frequencies up to 200 Mc/s are now being produced every year. Here we shall restrict ourselves to describing the principle underlying the manufacturing process.

The starting material, as in the case of the diffused base and mesa transistors that have just been described, is a supply of *P* type germanium wafers whose resistivity of 1 to 2 ohm cm is that required for the collector region. Two lead pellets of diameter 150 μm are placed about 50 μm apart on each germanium wafer. Both pellets contain the donor element antimony, and one of them further contains the acceptor element aluminium. The alloy-diffusion process takes place at 780 °C. The situation is as shown in fig. 8. Both pellets have taken up germanium, antimony was already present in solution, and one pellet additionally contains aluminium. At 780 °C antimony has a fairly fast rate of diffusion in solid germanium whereas aluminium, under the same conditions, has negligible diffusion. Accordingly, the donor element antimony invades the solid *P* germanium and turns it into *N* germanium. The temperature of 780 °C is maintained long enough for the *N* layer under the lead pellets to attain the required base width of 3 μm . But at the same time antimony vapour is escaping from the molten lead and penetrating the exposed part of the germanium crystal (fig. 8). The result is a continuous *N* layer that entirely encloses the crystal.

The holes made by the invading lead fill up again with germanium as the crystal cools. All the deposited germanium contains antimony, and that deposited under the pellet with aluminium contains aluminium as well. This last, an acceptor metal, has a solubility in solid germanium very much greater than that of the donor antimony; this means that the material which recrystallizes from the alumi-

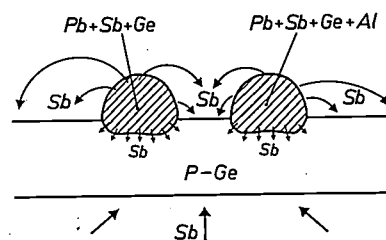


Fig. 8. Situation in a p.o.b. transistor immediately after it has reached the temperature of 780 °C at which the alloy-diffusion process takes place. The lead beads have melted and, in taking germanium into solution, they have eaten into the crystal. The temperature of 780 °C has to be maintained for some time to enable the Sb to diffuse into the crystal. Diffusion takes place in two ways, via the vapour phase and through the contact planes between the liquid lead and the solid germanium.

⁹⁾ The name "mesa" (the Spanish word for "table") is given to this type of transistor because of the shape of the crystal after the etching treatment (see figs. 6d and e). The nomenclature is not altogether apt: transistors of the "diffused base" type, mesa transistors and the p.o.b. transistors discussed below all have a base layer that is formed by diffusion, and all three have an active region of the table shape indicated in figs. 6d and e.

¹⁰⁾ P. J. W. Jochems, O. W. Memelink and L. J. Tummers, Construction and electrical properties of a germanium alloy-diffused transistor, Proc. I.R.E. 46, 1161-1165, 1958. Similar investigations have been carried out by J. R. A. Beale, Proc. Phys. Soc. (London) B 70, 1087-1089, 1957. The alloy-diffusion process has been employed previously by J. J. A. Ploos van Amstel of our laboratories, who by this means produced alloyed transistors with a drift field, and by R. L. Longini, who used the technique to make "hook-collectors" (British patent No. 754404).

niium-containing pellet is *P* type germanium having the low resistivity required for the emitter. The emitter region is thus formed under the Al-containing pellet, whereas under the other pellet an *N* region is formed, ensuring good non-rectifying contact with the base layer produced by diffusion (fig. 9).

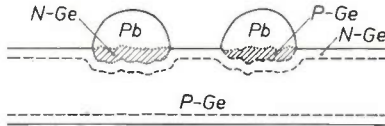


Fig. 9. The situation in a germanium wafer subsequent to alloy-diffusion treatment. The transistor structure has formed automatically. Evaporating from the lead beads, the donor element Sb has diffused into the *P* type germanium crystal, forming a surface *N* layer. This *N* layer is connected to the *N* layers that have formed under the two beads — which have now solidified again — constituting the base layer of the transistor. The hatched area on the right has solidified as *P* germanium, since here Sb atoms outnumber those of the acceptor metal Al; this constitutes the emitter. The hatched region on the left has solidified as *N* germanium and provides good electrical contact with the base. The above sketch is not to scale, the thickness of the diffused *N* layer being greatly exaggerated (cf. fig. 7).

The choice of a good “support metal” — the lead — is all-important. Designedly, the alloy-diffusion process is carried out at a temperature at which the diffusion rate for antimony is fast enough for the base layer to be formed within a few minutes. At this temperature the solubility of germanium in the support metal must be such as to ensure a reasonable penetration depth. Lead satisfies these requirements very well.

In the single-stage alloy-diffusion process, base and emitter contacts are produced as well as the required *P-N* junctions. The effect is as if the base layer were pushed out of the two lead pellets (hence the name “pushed-out base” transistor). Consequently neither the depth of the alloying front nor its flatness or orientation is of much significance and the thickness of the base layer can be accurately controlled. Thus the alloy-diffusion technique offers a neat way round all the difficul-

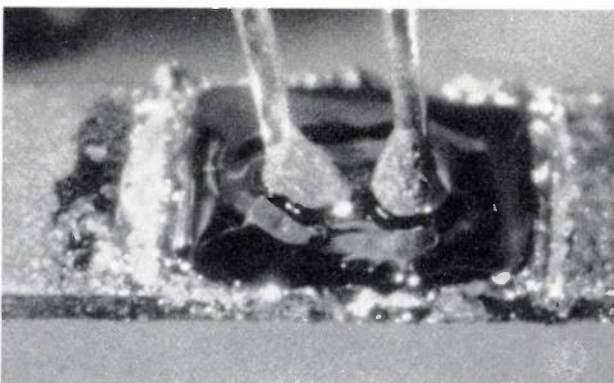


Fig. 10. Microphotograph of a p.o.b. transistor for operation at 100 Mc/s. Magnification 62 \times .

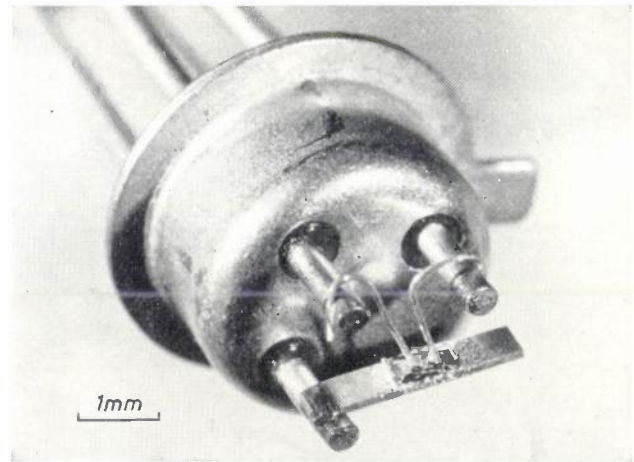


Fig. 11. A p.o.b. transistor mounted on its mechanical base.

ties presented by diffusing and alloying in separate stages. Fig. 10 is a microphotograph of a p.o.b. transistor; in fig. 11 the transistor can be seen mounted on its mechanical base, and fig. 12 shows the encapsulated transistor ready for use.

A refinement to the alloy-diffusion process is to carry out a preliminary diffusion treatment. As in the process by which diffusion and alloying are done in separate stages, the *P* germanium crystal is heated in antimony vapour so that a surface layer of *N* germanium is formed upon it. It then undergoes the alloy-diffusion treatment. The Pb from the pellets passes right through the existing *N* layer, so that the thickness of the base layer formed under the pellets is not affected by the pre-diffusion treatment. The *N* layer formed between the pellets is now thicker and consequently it is a better conductor. The overall result is a reduction in base resistance that improves the properties of the transistor at high frequencies.

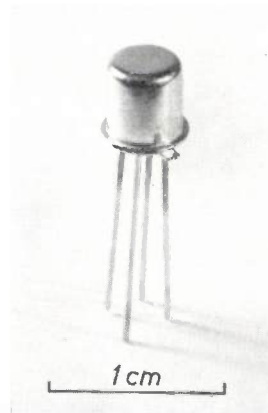


Fig. 12. The same transistor as in fig. 11, ready for use.

A microminiaturized solid-state circuit

In electronics there is a constant endeavour to reduce the size of components and complete circuits, and to increase their reliability by limiting the number of soldered joints. The advent of crystal diodes and transistors has greatly strengthened the trend towards microminiaturization. Not only are the diodes and transistors small in themselves, but

they develop very little heat because they embody no filaments. Heat dissipation is a big obstacle to more compact electronic equipment ¹¹⁾.

What is known as "solid-state circuitry" is one approach to microminiaturization: instead of being made up of discrete elements, the solid-state circuit takes the form of a single semiconducting crystal which has been processed into resistors and capacitors as well as diodes and transistors. Resistors of the desired value can be provided by parts of the crystal with the right shape, size and resistivity. Capacitors consists of *P-N* junctions of suitable dimensions. It is possible in this way to accommodate in one crystal an almost complete circuit having a given function, which may be amplification, pulse inversion, switching, the generation of oscillations or the like. For this purpose use can be made of the alloy-diffusion process of which an account has just been given. By way of example we shall describe here the making, under laboratory conditions, of a solid-state amplifier by means of the alloy-diffusion technique. The good HF properties available with this technique are not exploited in this particular circuit ¹²⁾.

The circuit in question, shown in *fig. 13*, is that of a three-stage amplifier. We shall not discuss the way it functions ¹³⁾, but merely show how the solid-state version is arrived at. The circuit elements outside the chain-dotted line in *fig. 13* are not embodied in the crystal, but are assembled on it as individual components.

¹¹⁾ A systematic review of the various approaches to miniaturization, illustrated with many examples of circuits that have actually been built, may be found in E. F. Horsey and P. J. Franklin, Status of microminiaturization, I.R.E. Transactions on component parts CP-9, 3-19, 1962 (No. 1).

¹²⁾ The necessary experiments were carried out by H. G. Kock.

¹³⁾ D. L. Jones, Directly coupled transistor hearing aid, Mullard tech. Comm. 5, No. 41, 2-9, 1959.

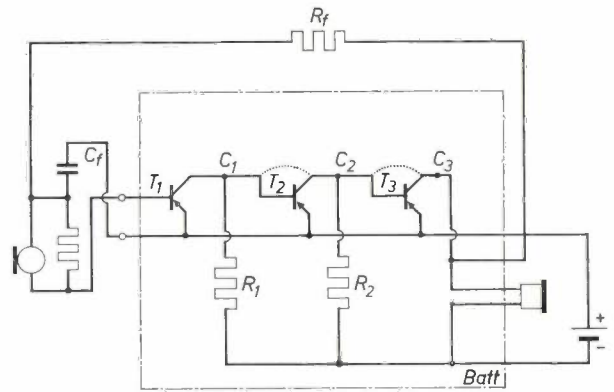


Fig. 13. Circuit diagram of a 3-stage amplifier ¹³⁾ which has been made by a "solid-state" microminiaturization technique. The part of the circuit enclosed by the chain-dotted line is accommodated in a single germanium wafer. The dotted lines between C_1 and C_2 and between C_2 and C_3 indicate shorting paths that have been broken by breaking the *P-N* junction at the places marked *I*, *II* and *III* in *fig. 15*. R_f and C_f provide DC feedback.

The starting material for the solid-state circuit is a germanium monocrystal in which a *P-N* junction has been produced during the pulling process (see p. 232). The section containing the junction is sawn up into slices and these "blanks" are cut to the required shape (*fig. 14*) with an ultrasonic drill ¹⁴⁾. It is important that the grown *P-N* junction should occupy the position indicated in *fig. 14*.

The damaged surface layer caused by the cutting operations is removed by etching, and three transistors are then made by the alloy-diffusion technique at the places indicated in *fig. 15*. In the case of transistors T_2 and T_3 no base pellet is applied since, for reasons that will shortly become evident, no base connection is required. The surface *N* layer formed

¹⁴⁾ E. A. Neppiras and R. D. Foskett, Ultrasonic machining, I. Technique and equipment, II. Operating conditions and performance of ultrasonic drills, Philips tech. Rev. 18, 325-334 and 368-379, 1956/57.

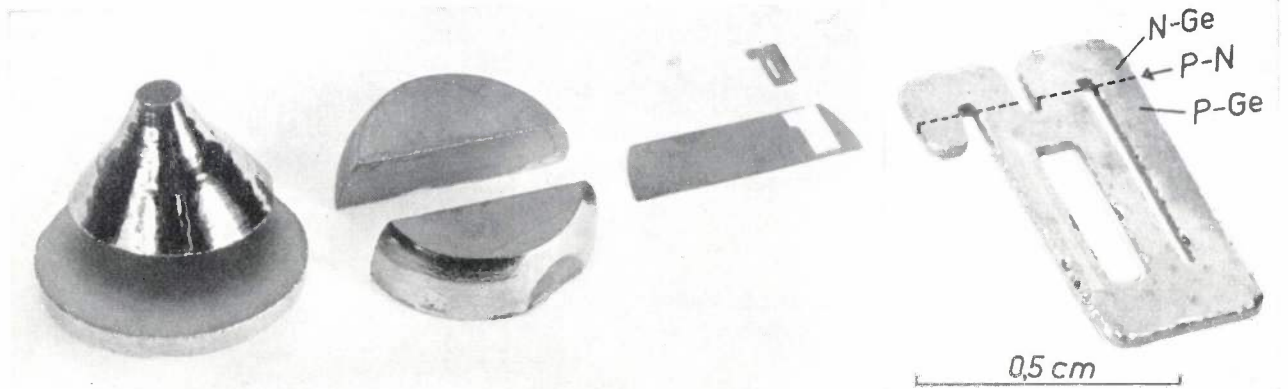


Fig. 14. The wafers are sawn out of a germanium crystal containing a *P-N* junction, and cut to the required shape with an ultrasonic drill. The *P-N* junction must occupy the position indicated in the enlarged photograph on the right.

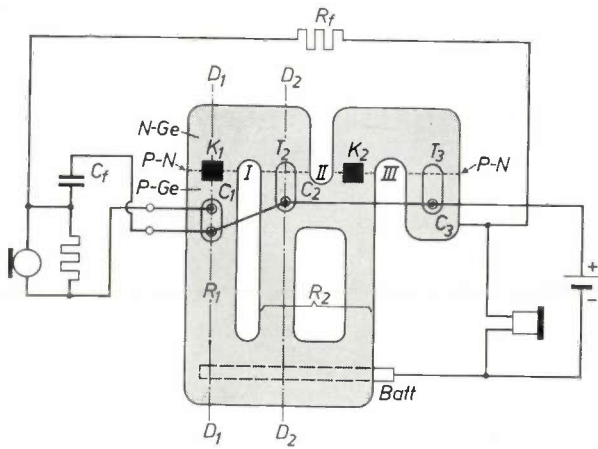


Fig. 15. Diagram to show how the circuit of fig.13 is built into the germanium wafer. Figs 13 and 15 have been drawn in such a way that resistors R_1 and R_2 and transistors T_1 , T_2 and T_3 occupy roughly the same positions relative to one another. K_1 and K_2 represent short-circuits across the grown $P-N$ junction (broken line). (Unfortunately, the letter T_1 has been omitted in the diagram.)

these two transistors (K_2 in fig. 15). Like that of T_2 , the base layer of T_3 extends over the neighbouring grown $P-N$ junction. It will now be clear that the base pellets can be omitted when T_2 and T_3 are made. The grown junction had to be broken at the places marked I, II and III in fig. 15 in order to prevent short-circuits between the collectors (along the dotted lines in fig. 13; cuts I and II break the connection between C_1 and C_2 , cut III breaks that between C_2 and C_3).

The battery contact, marked *Batt*, consists of a strip of metal soldered to the germanium crystal. The connection between the emitters has to be made by attaching wires to the appropriate pellets. Fig. 17 is a photograph of the finished solid-state circuit.

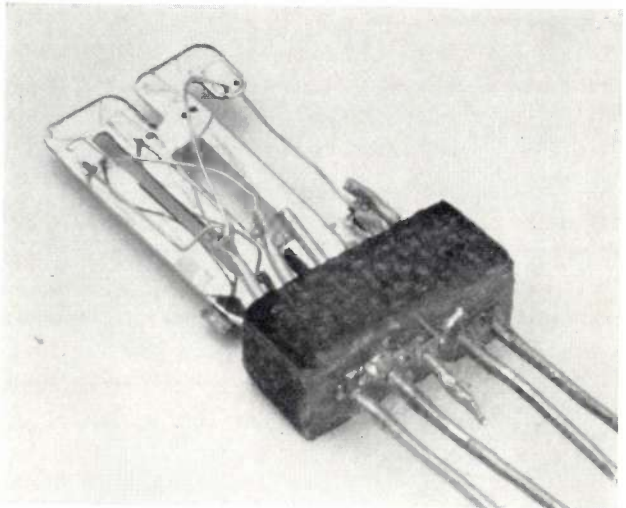


Fig. 17. Experimental version of a three-stage amplifier miniaturized by the solid-state technique. Additional lead beads have been deposited to facilitate measurements. Wire tails, going to a common emitter lead in the circuit mounting, have been attached to each of the emitter beads.

as a result of the alloy-diffusion process is etched away except in small areas around the lead pellets. A cross-section taken along D_1 in fig. 15, cutting through transistor T_1 , is shown in fig. 16a. As is clear from the circuit diagram in fig. 13, the collector of T_1 has to be connected to the base of T_2 without resistance; accordingly, a conducting strip is applied across the $P-N$ junction near T_1 . In practice it has been found that the junction can be short-circuited just as effectively by making a few scratches across it with a diamond. The base layer of transistor T_2 extends over the neighbouring grown $P-N$ junction (as can be seen also in fig. 16b, which is a cross-section taken along D_2 cutting through transistor T_2). By these means a connection without any intervening $P-N$ junction has been provided between the collector C_1 of T_1 and the base of T_2 . The direct connection between the collector of T_2 and the base of T_3 is established by short-circuiting the grown $P-N$ junction between

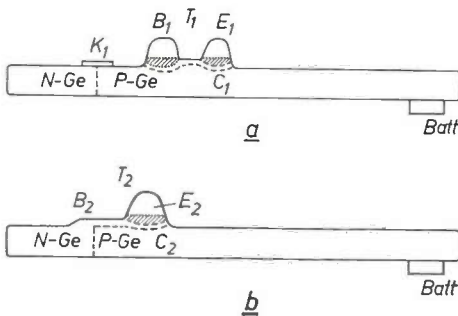


Fig. 16. a) Cross-section through transistor T_1 , taken along line D_1 in fig. 15. K_1 short-circuits the grown $P-N$ junction produced in the original crystal.
b) Cross-section through transistor T_2 , taken along D_2 .

Summary. Almost all transistors nowadays are of the junction type. One way of producing the $P-N$ transitions in junction transistors is to add suitable impurities while the germanium (or silicon) monocrystal is being pulled from the melt. However, the alloying and the diffusion methods are the ones most commonly employed. The higher the frequency at which the transistor is required to provide a reasonable signal gain, the more difficult it becomes to make good $P-N$ junctions efficiently, mainly because extending the frequency range necessitates reducing the dimensions of the transistor. In the alloy-diffusion technique, the $P-N$ junction between emitter and base is formed as a result of an alloying process, and that between base and collector is formed as a result of the diffusion process which takes place at the same time; this, in essentials, is the method developed in Philips Research Laboratories for mass-producing HF transistors. Many millions of transistors capable of operating at frequencies up to 100 and 200 Mc/s are now being manufactured by this method every year.

The alloy-diffusion technique has also been used in experimental work on microminiaturization. As an example the processing of a germanium wafer into an almost complete three-stage amplifier is described.

HEATING A SOLID IN VAPOUR WITH INDEPENDENT PRESSURE AND TEMPERATURE ADJUSTMENT

548.522

In solid-state technology it is sometimes necessary to heat a solid in the vapour of some other substance. This is often the case with semiconductors which are to be given certain electrical or optical properties. For this purpose the vapour temperature and pressure should be independently adjustable within wide limits. A method frequently used to this end is the "double-furnace method." The material to be processed, e.g. cadmium sulphide, is contained at one end of a closed tube, e.g. the left-hand end (fig. 1). The other end contains a substance, e.g. sulphur, which produces the vapour in which the cadmium sulphide is to be heated. The left-hand end is heated in a furnace O_I to a temperature T_1 , and the right-hand end in another furnace O_{II} to a (lower) temperature T_2 . The cadmium sulphide is therefore at the temperature T_1 , and the pressure is equal to the saturation pressure of sulphur at the temperature T_2 . The temperature and pressure of the sulphur vapour are thus in fact independently adjustable.

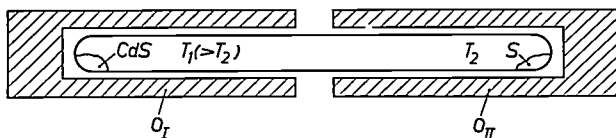


Fig. 1. Double-furnace method of heating a solid, e.g. CdS, in a vapour (e.g. S vapour) with independent temperature and pressure adjustment. The pressure is indirectly controlled by means of the temperature T_2 .

Objections to the double furnace-method are:

- 1) Since the saturation pressure varies strongly with the temperature T_2 , it is necessary to be able to regulate T_2 very exactly in order to adjust the pressure with reasonable accuracy and to keep it constant. Moreover the precise relation between saturation pressure and temperature must be known.
- 2) It may be necessary to "freeze in" a high-temperature equilibrium by very rapidly cooling the treated substance. This is difficult to do with the system shown in fig. 1.

Fig. 2 shows a sketch of a simple and practical system that operates on a different principle, avoiding the above drawbacks ¹⁾. The equipment consists of a vertical tube, the bottom part of which is heated in a furnace O , while the part that projects above

the furnace is cooled by an airstream (thick arrows). The tube is closed by a lid D with rubber gasket, and is connected via an opening in the lid to a source of neutral gas which is kept at an adjustable pressure p . The substance to be vaporized, e.g. sulphur, is initially at the bottom of the tube. When the lower part of the tube is raised to a temperature above the boiling point of sulphur, the latter completely evaporates. In a cooler section in the upper part of the tube the sulphur vapour condenses on the wall (zone B). The liquid sulphur runs down the wall and on reaching the hotter zone A it again evaporates. This circulation process is indicated by thin arrows. The solid under treatment, e.g. CdS, is in a container which is positioned near the bottom of the tube (zone C) by means of a quartz rod attached to the lid. Although the composition of the gas differs at different parts of the tube, the total pressure has everywhere the adjusted value p . In zone C virtually pure sulphur vapour is present, other gases that were initially to be found in C having been gradually removed: gas molecules entering the rising sulphur stream in zone B are carried along by

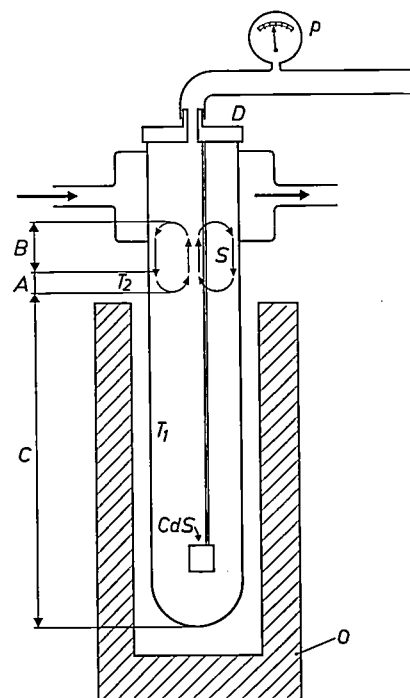


Fig. 2. System in which the temperature T_1 and the pressure p of the vapour are both directly and independently adjustable.

¹⁾ C. Z. van Doorn, Rev. sci. Instr. 32, 755-756, 1961.

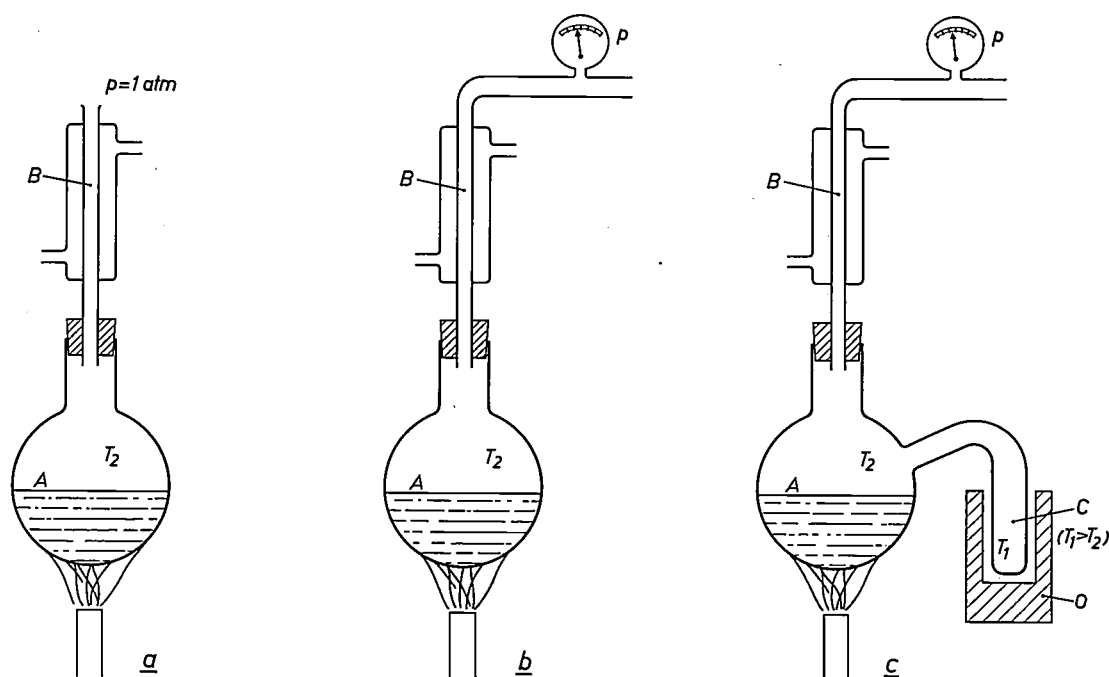


Fig. 3. Illustrating the principle of the system shown in fig. 2. Parts A, B and C correspond to the respective parts in fig. 2.

this stream. The sulphur returns as liquid along the wall, so that finally zone C contains virtually pure sulphur vapour having the adjusted pressure p . The temperature T_1 can be regulated with the furnace independently of p , with the above-mentioned restriction that T_1 must be higher than the boiling point T_2 of sulphur at the pressure p . The boiling point temperature T_2 prevails in zone A; in zone C the vapour is unsaturated.

At any moment the sample can easily be taken out of the tube. All that has to be done is to adjust the pressure in the tube to one atmosphere, after which the lid (which is cool) can be removed together with the sample container. This system thus overcomes the objections to the double-furnace method, and has the additional advantage of automatic purification of the vapour. A drawback of the new method is that it cannot be used for substances that have a high vapour pressure at their melting point. In that case large solid deposits will form at the top of the tube and these can no longer take part in the circulation, which will then finally break down.

We have already used the method successfully for heating solids in Cd and S vapour (in a

quartz tube) and in K vapour (in a nickel tube). Presumably the method will also be useful for vaporized Na, Rb, Cs, Zn, Hg, Se, Te and P.

A further explanation of the principle is given in fig. 3. The familiar set-up in fig. 3a consists of a retort A connected to a vertical condenser B. In A a liquid is brought to boiling point (T_2). The resultant vapour condenses in B into liquid that flows back into A. The vapour in A has the pressure of the outside atmosphere. In fig. 3b the tube to the outside atmosphere is connected to a means of adjusting the pressure p ; the vapour will, of course, also be at pressure p . The temperature T_2 of the vapour is the boiling point of the liquid at the pressure p , and cannot therefore be varied independently of p . This is possible however if, as shown in fig. 3c, the vapour space in A is connected with a side tube C, which is kept at a constant (higher) temperature T_1 by a furnace O. In space C the temperature and pressure are independently adjustable. The spaces A, B and C in fig. 3 correspond to the respective zones in fig. 2.

C. Z. van DOORN *).

*) Philips Research Laboratories, Eindhoven.

THE MAGNETIZATION REVERSAL PROCESS IN SQUARE-LOOP FERRITES

by J. E. KNOWLES *).

538.23:621.318.12

Magnetic components made of a ferrite with an intrinsically rectangular hysteresis loop are widely used in computer stores, and for switching and logic applications. The study of the configuration and motion of domain walls in ferrite grains by means of Bitter patterns has provided valuable insight into the conditions for obtaining "square loop" ferrite giving optimum performance, and into the static and dynamic characteristics of these materials.

Introduction

Magnetic components made of a ferrite with an intrinsically rectangular hysteresis loop now find widespread use in computer stores, and for switching and logic applications. Previous investigations into the properties of such "square loop" ferrites have been mainly macroscopic. It therefore seemed of interest to examine the dependence of the bulk properties on those of the individual crystallites (grains) from which the materials are built up. This is interesting both as an academic problem and because the understanding of the phenomenon might lead to the design of better materials. This article discusses some investigations of the phenomena occurring in the individual grains when the magnetization of the grains reverses under the action of an applied magnetic field.

First of all it will be useful to recapitulate, with reference to *fig. 1*, some terms which were defined in a previous article in this journal ¹⁾ and which will be used in this article too. The solid loop in *fig. 1* gives the magnetization J as a function of the magnetic field H , the latter varying from large negative to large positive values and back. The coercive force H_c , the saturation magnet-

ization J_s and the remanence J_r are well-known quantities. A smaller amplitude H_m of the magnetic field gives rise to the "inner" loop, which is shown as a broken line in *fig. 1*. The "squareness ratio", R_s , is defined by:

$$R_s = \frac{J(-\frac{1}{2}H_m)}{J(H_m)}, \quad \dots \quad (1)$$

where $J(-\frac{1}{2}H_m)$ and $J(H_m)$ are the magnetizations at the fields $-\frac{1}{2}H_m$ and $+H_m$ respectively. The value of R_s is always less than unity and is usually a function of H_m . If the magnetic material is to be used as a storage element, a maximum value of R_s is required which is near to 1. It can be seen from *fig. 1* that a high value of R_s is most likely to occur if the remanence J_r is only a little less than the saturation J_s . In a polycrystalline material, of the kind considered here, the difference between J_s and J_r arises mainly because the magnetization has preferred or easy directions: if the applied field is removed, the magnetization rotates from the direction parallel to H to the nearest easy direction in each crystallite.

It was pointed out in ¹⁾ that in cubic non-oriented materials, the ratio J_r/J_s could be high only if the magnetocrystalline anisotropy predominates over the stress and shape anisotropies. The stress anisotropy is caused by the interaction between the magnetization and stresses in the material and is proportional to the magnetostriction constants; the shape anisotropy is associated here with the porosity of the ferrite. The crystalline anisotropy is a fundamental magnetic property, and consequently preferred directions exist for the magnetization even in ideal single crystals. The anisotropies are expressed in terms of the energy required per unit volume to rotate the magnetization into the energetically most unfavourable direction ²⁾. Both stress and

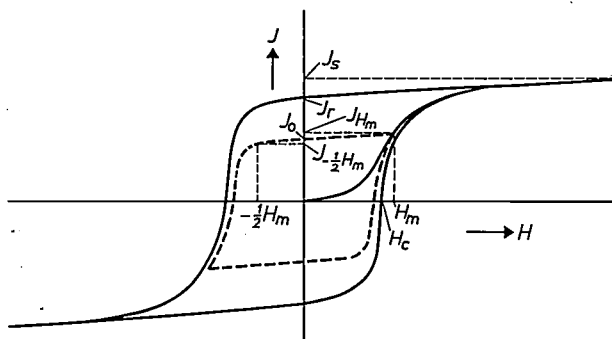


Fig. 1. Hysteresis loops of a ferrite.

*) Mullard Research Laboratories, Salfords, England.

¹⁾ H. P. J. Wijn, E. W. Gorter, C. J. Esveldt and P. Geldermans, Philips tech. Rev. 16, 49-58, 1954/55.

²⁾ See J. J. Went and E. W. Gorter, Philips tech. Rev. 13, 181-193, 1951/52, and J. J. Went *et al.*, Philips tech. Rev. 13, 194-208, 1951/52.

shape anisotropies have in common that at each point in the material there are only two (antiparallel) preferred directions, whereas the crystalline anisotropy in ordinary cubic ferrites gives rise to eight preferred directions in each crystallite, i.e. the four body diagonals of the unit cell (the [111] directions in both senses).

Therefore if only the crystalline anisotropy exists, a cone with a solid angle of $\pi/2$ around any direction of H will certainly contain an easy direction in each grain, whereas other anisotropies require a solid angle of 2π to contain an easy direction. If the magnetic field goes from a high value to zero, the direction of the various magnetization vectors for each grain will spread much less if only crystalline anisotropy is present than if stress or shape anisotropies dominate (fig. 2). This explains why the value of J_r/J_s in the first case will be much higher than in the second case (0.87 and 0.5 respectively).

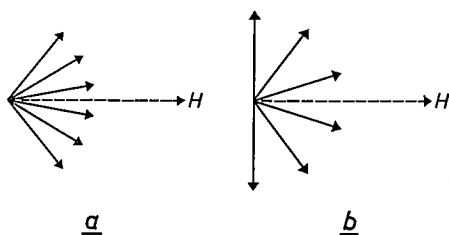


Fig. 2. The distribution of magnetic vectors in a cubic polycrystalline material after removing the field H . a) The crystalline anisotropy is dominant. b) The stress or shape anisotropy dominates.

It was also said in I that a small value of λ_{111} , the magnetostriction constant in the preferred direction, favours a small stress anisotropy, and that the shape anisotropy can be decreased by reducing the porosity of the material. A small value of λ_{111} implies that the magnetostriction produced by low fields will be small.

The above-mentioned definitions and properties serve as a suitable starting point for the investigation described in this article.

The relationship between magnetostriction and squareness ratio

The first experiments were directed to examining the relation between low field magnetostriction, the maximum squareness ratio, and the resistivity, of a range of manganese-magnesium ferrite rings of the same nominal composition, i.e. that of ferroxcube D1, made by Mullard. Six rings (diameter 2 cm) were fired for long periods in different atmospheres to produce specimens in various stages of oxidation or reduction. When the material is fired in a slightly

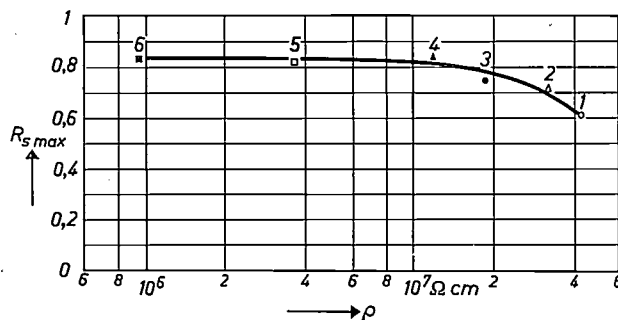


Fig. 3. The variation of maximum squareness ratio $R_{s,max}$ with resistivity ρ of six samples of manganese-magnesium ferrite of the same nominal composition.

reducing atmosphere a small portion of ferric iron (Fe^{3+}) changes to ferrous iron (Fe^{2+}), so that the composition may be thought of as mixed crystals of the original ferrite and ferrous ferrite (Fe_3O_4). The Fe_3O_4 has a lower electrical resistance because the extra electron of the ferrous ion can easily jump from one iron site to another. The electrical resistance is therefore a sensitive indicator of the Fe_3O_4 content.

The rings were numbered from 1 to 6 in order of decreasing resistivity. Measurement of the maximum squareness ratio $R_{s,max}$ showed that its value decreases with increasing resistivity (fig. 3). To measure the magnetostriction in low magnetic fields, the ring under observation was first demagnetized by passing an alternating current of diminishing amplitude through windings wound around the ring. A low magnetic field was then applied to the specimens by passing a direct current through the windings. The minute decrease in diameter (of the order of 50 Å) caused by the magnetostriction was measured by a modified commercial comparator which used a differential transformer with a movable ferrite core as a transducer.

The results of these measurements are summarized in fig. 4. From figs 3 and 4 it is evident that

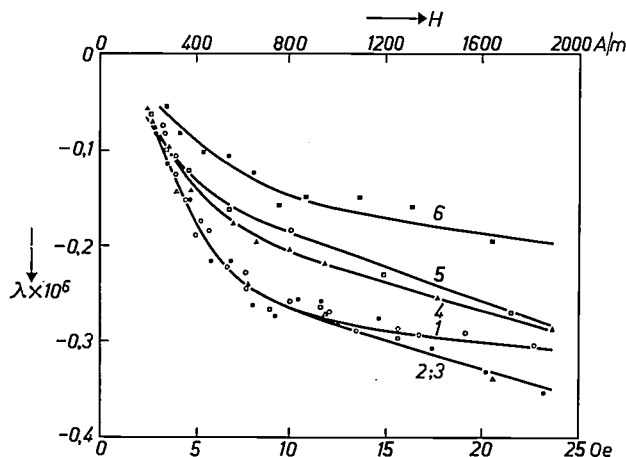


Fig. 4. The magnetostriction λ as a function of the applied field H for the same six samples as in fig. 3 (compare numbers).

rings of low resistivity have a small low-field magnetostriction and a high squareness ratio, whilst rings with a high resistivity show a larger magnetostriction and are less "square". A satisfactory explanation of these results is to assume that the negative magnetostriction of the manganese-magnesium ferrite was offset by the large positive magnetostriction of the small admixture of ferrous ferrite. Thus samples containing a relatively large proportion of ferrous ferrite show a low resistivity, a small magnetostriction and a high squareness ratio. However, the results of these experiments are not entirely conclusive, since the grain structure and crystalline anisotropy of the ferrite may well vary somewhat with the firing conditions.

Domain pattern observations

The most direct way of examining the nature of the magnetization reversal process in a material, as described by its hysteresis loop, is to observe the domain structure, i.e. to study the pattern formed by the walls of the domains into which the material is divided. The domains are all magnetized to saturation, but in different directions. An idealized domain structure is illustrated in *fig. 5a*. Here the magnetization vectors of adjacent domains are oppositely directed, giving two sets of domains. If the total volumes of opposing domains are equal, then the overall magnetization of the body will be zero. When a field is applied, one set of domains grows at the expense of the other (*fig. 5b*) and the body becomes magnetized.

The walls between domains (called domain or Bloch walls) can be observed in the following way. First of all the specimen must be carefully prepared,

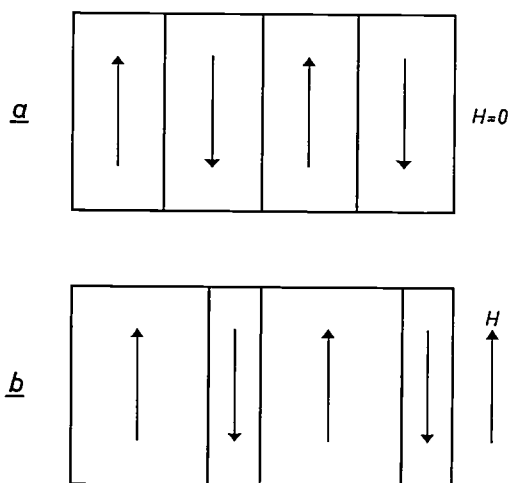


Fig. 5. An idealized domain structure *a*) demagnetized, *b*) after application of a small field H . The domains whose magnetization vector is parallel to H have grown in (*b*) at the expense of the others.

usually by mechanical polishing and then etching away the thin layer of material which has been stressed by the polishing process. This layer must be removed because its stresses introduce a superficial stress anisotropy which spoils the domain pattern that is characteristic of the bulk phenomena. The walls may then be rendered visible by wetting the surface with a colloidal suspension of magnetic particles. The particles are attracted by the high stray fields which lie along the intersection of a domain wall and the surface of the sample, making it possible to see the walls with a microscope. This technique was originated by Bitter³⁾. The visible network of domain walls is referred to as a "Bitter pattern" or as a powder pattern. Other methods of observing domain patterns, such as those depending on rotation of the plane of polarization of reflected or transmitted light⁴⁾, or on the deviation of an obliquely incident electron beam, are not suitable for these relatively thick, insulating ferrite specimens.

The observation of Bitter patterns on ferrites is not as straightforward as on other ferromagnetic substances, since the lack of a suitable etchant makes it difficult to prepare a smooth stress-free surface. After several techniques had been tried, it was found that a good surface could be prepared simply by pressing the prefired powder between highly polished punches in a suitable die and then firing the ring so produced in the usual way⁵⁾.

The progress of the magnetization reversal in a grain of ferrite, in which the magnetization was parallel to the surface, was observed by making a series of photographs of the changing powder pattern during the reversal process. Three of these photographs are shown in *fig. 6*. It may be seen that the grain, which was approximately rectangular in plan, showed two well-defined domain walls which were approximately parallel to the field direction (the black arrow H). Initially the grain was in the "positive" state of remanence. As the field was increased from zero in the negative sense, the two domains which had initially been close together moved apart. They did not move smoothly, however, but in a series of jerks; the Barkhausen jumps. For example, when the field was slowly increased from 98 to 107 A/m (1.27 Oe, to 1.34 Oe), the lower domain wall moved abruptly from the position shown in *fig. 6a* to that shown in *fig. 6b*. On further increasing the field to

³⁾ F. Bitter, *Phys. Rev.* **41**, 507, 1932.

⁴⁾ See C. Kooy, *Direct observation of Weiss domains by means of the Faraday effect*, *Philips tech. Rev.* **19**, 286-289, 1957/58.

⁵⁾ J. E. Knowles, *Proc. Phys. Soc. (London)* **75**, 885, 1960.

110 A/m (1.4 Oe) the right-hand end of the bottom wall moved discontinuously to the position shown in fig. 6c. It will be observed that on the right-hand end of the wall there is a heavy deposit of colloidal particles, and that the wall is markedly kinked. Examination of the complete series of photographs revealed that this region was the site of some form

of obstacle to the free movement of the domain walls, and indeed a number of other obstacle sites were observed over the surface of the specimen. Their origin will be discussed later, but for the moment it may be said that any form of inhomogeneity such as a local change in chemical composition, a region of high stress, or a void, will act as an obstacle to domain wall motion.

It is of interest to obtain the hysteresis loop for a single grain and to compare it with that for the whole toroidal specimen. If it is assumed that the grain is not only rectangular in plan but also rectangular in cross-section, then the magnetization of the grain may be deduced from the relative areas of the domains. This was done for the grains to which the above series of photographs related. The resultant hysteresis loop is shown as the curve *b* in fig. 7, the

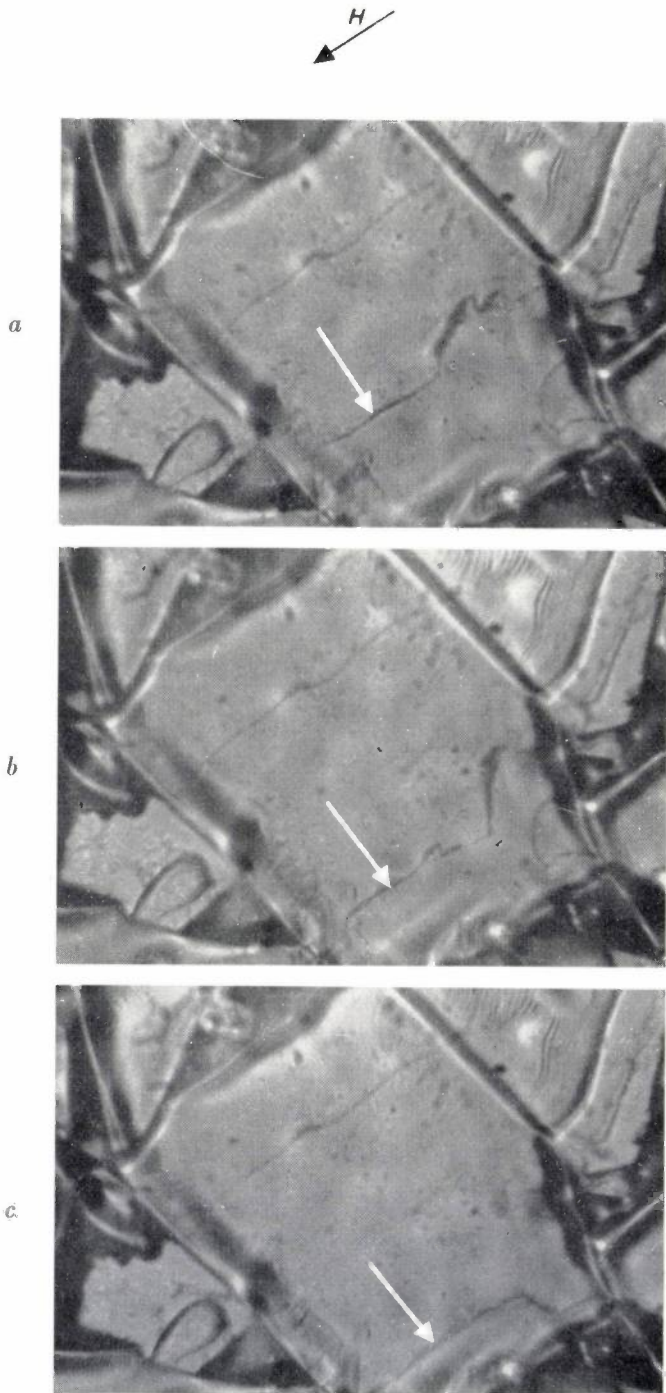


Fig. 6. The movement of a domain wall in a grain of manganese-magnesium ferrite during the magnetization reversal process. The direction of the applied field is indicated by the black arrow *H*; the field increased from *a* to *c*. The relevant wall is indicated by a white arrow.

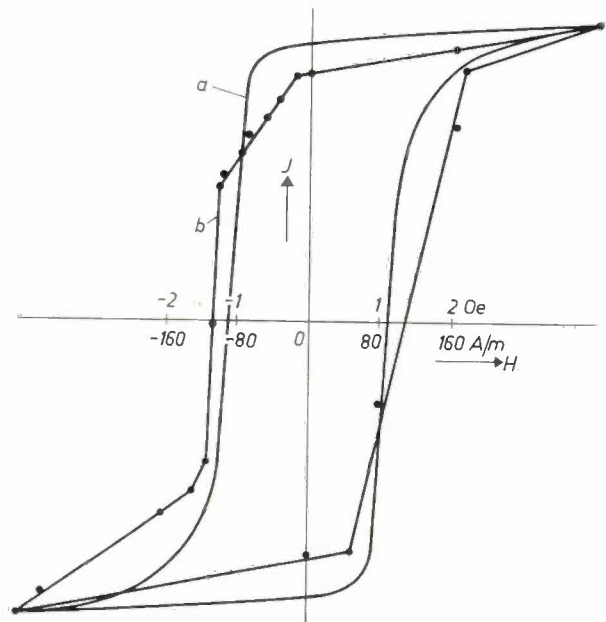


Fig. 7. The hysteresis loop of a polycrystalline ring of manganese-magnesium ferrite (curve *a*), together with that of an individual grain of it (curve *b*). The *J* scale is adapted to make the tips of the loops coincide (top right and bottom left).

points obtained from the complete series of photographs being indicated on it. The measured hysteresis loop of the whole ring is drawn as the curve *a*, whereby the vertical *J* scale has been adapted so that the tips of *a* and *b* coincide. It can be seen that there is an essential similarity between the two loops, which supports the hypothesis that the origin of the "squareness" of the polycrystalline material must be sought in the individual grains.

The hysteresis loop of a single crystallite

Let us consider a cubic-shaped grain which at remanence supports a single domain wall adjacent

to one of its faces. If now a field parallel to the wall and of the correct sense is applied, the wall will tend to move across the grain, and so reverse the magnetization of the grain. Suppose that the motion of the wall is hindered by a series of equally spaced obstacles whose magnitudes are scattered in a random (Gaussian) way about some mean value; the magnitude of an obstacle may be expressed in terms of the minimum field necessary to force a domain wall past it. If now a gradually increasing field is applied, the wall will move freely until it meets the first obstacle in its path, where it will remain until the field becomes sufficiently large to force it past the obstacle. If the obstacles were arranged in order of increasing magnitude, the relevant branch of the hysteresis loop would have the form shown in *fig. 8a*. In reality the magnitude of the obstacles will be randomly distributed, so that the wall will not stop at every successive obstacle but only at those larger than any of the preceding ones. The resultant hysteresis loop will then be rather rectangular, as shown in *fig. 8b*. On the basis of this model, the coercive force will be rather greater than the average magnitude of the obstacles in the relevant grain.

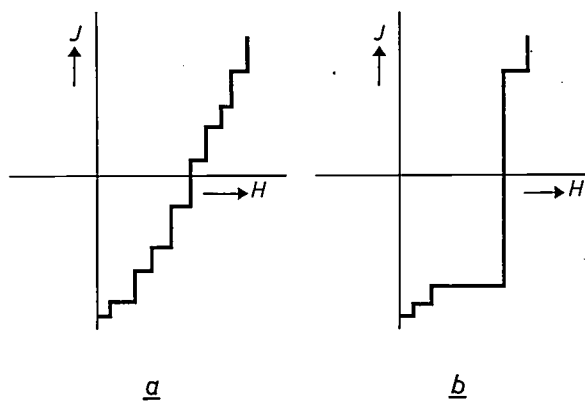


Fig. 8. Possible curves of the magnetization J against the field H for a single grain. When H has grown sufficiently to enable the wall to pass the next larger obstacle, J increases with a jump.

Since the obstacles in the various grains will of course be randomly distributed and their magnitudes may show a wide spread around the mean value, the question arises as to what the hysteresis loop would look like if the average were found from a large number of otherwise identical grains. "Otherwise identical" means that these grains are identically oriented, are traversed by only one wall, and that in all grains the obstacles have the same average magnitude.

In the extreme case where the spread is zero and all obstacles are equal in magnitude, the "average"

hysteresis loop contemplated will be perfectly rectangular. A mathematical treatment based on a theory of Néel ⁶⁾ shows that even for quite large spreads in the magnitude of the obstacles, the average loop shape will still be fairly rectangular provided that the wall traverses a sufficiently large number of them in a grain, say twenty ⁷⁾. This does not of course apply to the individual grains; if the spread in magnitude is large the individual loops may be far from rectangular, and also in general the loops will be unsymmetrical.

Values of the above order of magnitude for the number of obstacles which a wall passes during a magnetization reversal have in fact been found experimentally on specimens of two square loop manganese-magnesium ferrites.

The average number of obstacles that a domain wall crosses during the course of a magnetization reversal can be estimated by measuring the non-linear terms of the magnetization as a function of the amplitude of an applied alternating field. By this method the values 14 and 21 were obtained from two slightly different samples. An independent method ⁸⁾, which compares the time for overcoming an obstacle with the time of complete magnetization reversal, led to a value of 25 for a sample nominally identical with the sample giving the value 21. These results were considered to agree satisfactorily.

The hysteresis loop of the polycrystalline material

Having shown that the individual grains of the type considered above tend to have rectangular loops, the properties of an assembly of randomly oriented grains must now be examined. Let us first consider a simple theoretical model. According to this model the material is supposed to have the following properties:

- 1) The crystalline anisotropy predominates over the stress and shape anisotropies, so that the magnetization of a grain will, for ordinary materials, always lie in a [111] direction. If the easy direction was [100], the results would not be essentially different. This is the criterion for high remanence stated in ¹⁾.
- 2) All grains are of the type considered, i.e. approximately cubic in shape, with one or more walls parallel to a face of the grain. Each grain has a perfectly rectangular hysteresis loop.
- 3) All grains have the same coercive force H_0 .
- 4) There is no interaction between neighbouring grains.

⁶⁾ L. Néel, *Cah. Phys.* No. 12, 1942; No. 13, 1943.

⁷⁾ J. E. Knowles, *Proc. Phys. Soc. (London)* 77, 225, 1961. In this publication a more detailed description of the "averaging" process is given.

⁸⁾ J. E. Knowles, *Proc. Phys. Soc. (London)* 77, 576, 1961.

Let the polycrystalline material be in the remanent state, and suppose that in this condition all the domain walls are adjacent to a face of a grain; some of the domain walls will be parallel to the field; but most will make fairly large angles with it. As described in the introduction, if condition (1) above is satisfied, the directions of the magnetization of the grains are all contained within a solid angle $\pi/2$. For convenience of calculation, this is approximated by a solid cone of semi-vertical angle 55° , the axis of which is parallel to the direction in which the field was applied. Thus the vectors representing the magnetization of the grains can at remanence be sketched as in fig. 9a. Suppose now that a slowly increasing field is applied in the sense tending to reverse the magnetization. Because of the assumed perfect squareness of the loops, no change will occur until the magnitude of the field just exceeds the grain coercive force H_0 , when those few grains where the magnetization lies antiparallel to the field will reverse their magnetization (fig. 9b). In other grains where the magnetization makes an angle ϑ with the field, the field tending to move the domain walls is $H_0 \cos \vartheta$, and these grains will not be switched. As the field increases, grains making progressively larger angles with the field will successively reverse their magnetization (fig. 9c) until, when the applied

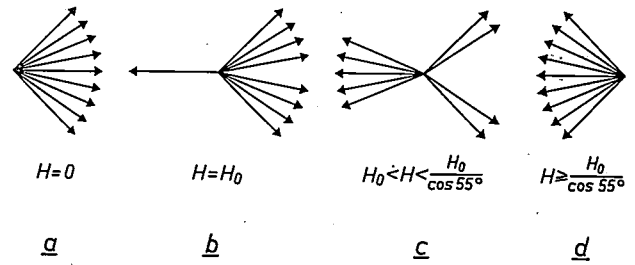


Fig. 9. A two-dimensional representation of the orientation of the magnetization vectors during the reversal process for a specimen in which the crystalline anisotropy is dominant (cf. fig. 2a).

field has the value $H_0/\cos 55^\circ$, all the grains will have reversed their magnetization and the position will be as in fig. 9d. It is easy to calculate the corresponding hysteresis loop; this loop is shown in fig. 10a. The ratio J_r/J_s obtained from this simple model is 0.79; a more exact calculation, taking into account the cubic symmetry of the material, gives the value 0.87.

The effect of reversible movements of the domain walls and rotation of the magnetization in each domain has been neglected. The order of magnitude of these processes can be estimated from the incremental permeability in weak AC fields. In common square-loop ferrites it is not more than a hundred or so. On this basis the loop of

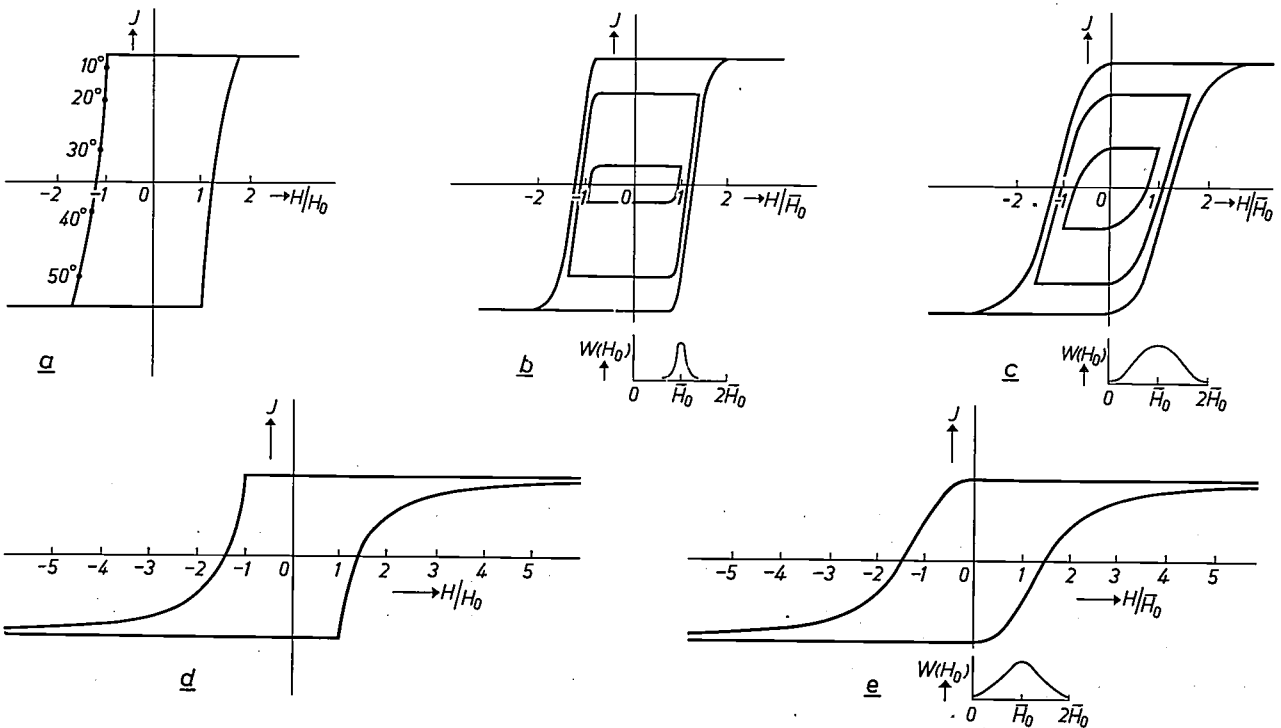


Fig. 10. Theoretical hysteresis loops for polycrystalline samples, neglecting the reversible processes. In a, b and c the magnetocrystalline anisotropy predominates, in d and e the stress or the shape anisotropy is predominant. In a and d each grain has the same coercive force H_0 ; in b, c and e the coercive force of different crystallites is distributed around H_0 according to the shown distribution function $W(H_0)$. The numbers in (a) represent the semi-vertical angle of the cones in which the magnetization reverses (fig. 9).

fig. 10a may be corrected to incorporate these effects by re-drawing in an oblique coordinate system in which the field axis has been rotated about two degrees in the clockwise sense.

It is thus seen that if all the grains had equal coercive forces and rectangular loops, then despite their random orientation the resulting loop for a polycrystalline material is still very "square". We shall now make our model more realistic by dropping the over-simplifying assumption that all grains have the same coercive force, and assuming instead that the grain coercive forces are distributed following a Gaussian function $W(H_0)$ about some mean value, \bar{H}_0 . Figs 10b and 10c show calculated hysteresis loops corresponding to a spread of grain coercive force indicated by the distribution curve under each loop. It may be seen that the effect of a progressively broader distribution is to shear the loop and round off the knee, the remanence remaining unchanged.

Up to now we have considered the case where the crystalline anisotropy predominates. If either the stress or the shape anisotropy is large, then, as described, the solid cone which, at remanence, contains the magnetization vectors, will become a hemisphere. Using the same arguments as before, fig. 10d is obtained for the case where all grains have the same coercive force; the effect of a fairly wide spread in the grain coercive force is shown in fig. 10e. The ratio J_r/J_s is in both cases 0.5. If it is further supposed that there is appreciable rotation of the magnetization in weak fields, then fig. 10e changes to the hysteresis loop of a typical soft magnetic material with $J_r/J_s = 0.5$. Summarizing, it may be said that fig. 10 illustrates how the increasing relative influence of stress/shape anisotropy, together with an increasing spread in the grain coercive force, causes the ideal rectangular loop (fig. 10a) to degenerate progressively into a "commonplace" loop (fig. 10e).

So far the hysteresis loops considered have all been quasi-static characteristics of the material, i.e. curves giving the variation of J when H is slowly varied. In normal use, however, the magnetization of square-loop ferrites is reversed dynamically by a field pulse of very short rise time and the voltage induced on a secondary winding is measured. The *dynamic* behaviour must then be considered. Experiments have shown that this output voltage takes a time to reach its maximum value, which is about half the time taken for the material to reverse its magnetization, that is the "switching time". The model postulated above gives rise to a theoretical output curve which is entirely different in shape, i.e. one which reaches its maximum value instantaneously

at the start of the field pulse, stays constant for some time, and then decays away.

Among the several possible reasons for this disagreement between theory and experiment, the most likely seems to be that the total area of the domain walls, instead of remaining constant throughout the reversal, rises to a maximum near the coercive force and finally decreases to the initial value. This hypothesis is also in agreement with the observation that as the material is taken round a hysteresis loop, the incremental permeability dB/dH passes through a maximum for biasing fields roughly equal to the coercive force.

This permeability, as stated, is caused by two mechanisms, i.e. by reversible wall movement and by rotation of the magnetization in each domain. The contribution of the latter mechanism can be calculated fairly well from the crystalline anisotropy. Crystalline anisotropy measurements on the material in question⁹⁾ indicate that rotations account for about one tenth of the total permeability of 87; the incremental permeability is thus almost entirely due to domain wall motion. Since this permeability is proportional to the total wall area, it seems obvious to assume that a permeability peak corresponds to a maximum in total wall area.

Further domain pattern observations

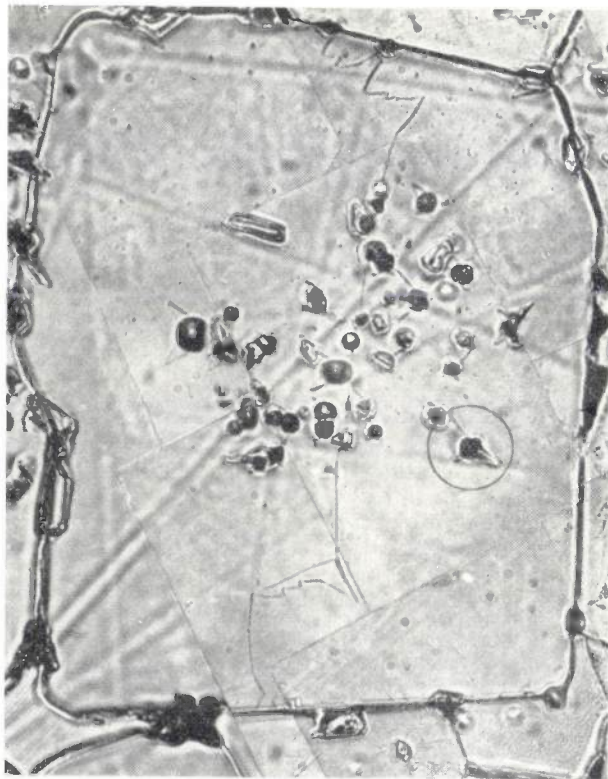
It should be possible to observe directly the above-mentioned wall area. In an attempt to observe this increase in wall area, and also to throw some light on the nature of the "obstacles" mentioned earlier, a further series of domain pattern observations was made¹⁰⁾. A ring of square-loop manganese-magnesium ferrite was fired for a long time to promote the growth of large grains. The surface on this occasion was prepared by polishing with diamond paste and then annealed by re-heating the specimen at the original firing temperature¹¹⁾. This method had the advantage of revealing the pores, which were clustered together near the centre of the grain and also along the grain boundaries.

Fig. 11a shows a photograph of a very large rectangular grain in which the magnetization happened to lie in the surface. It may be seen that the approximately spherical pores are surrounded by Néel spike domains; for pores that are not cut by a wall, this is the wall pattern that corresponds to minimum demagnetization energy. A few of these pores with "spikes" are drawn in the sketch (fig. 11b). The well-defined wall in this grain is composed of a number of lengths of 180° wall — i.e. walls dividing

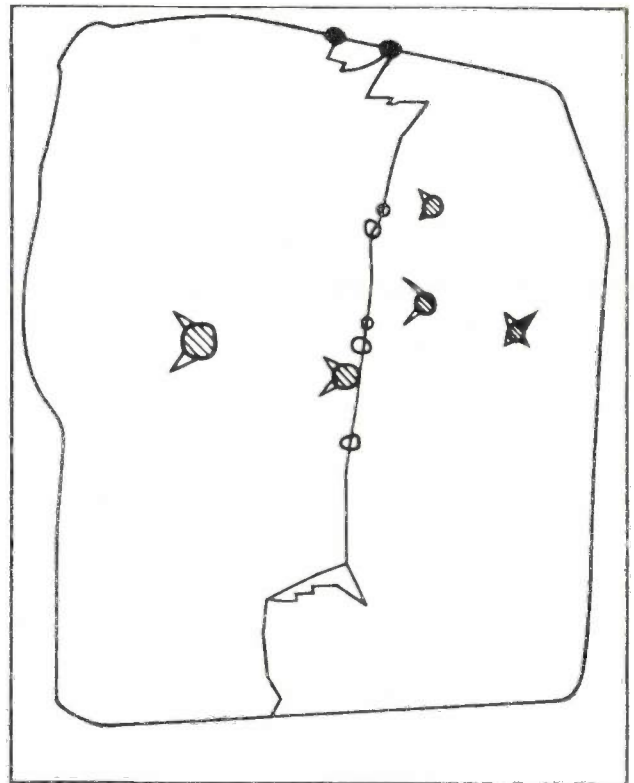
⁹⁾ J. R. Chamberlain, Proc. Phys. Soc. (London) 78, 819, 1961.

¹⁰⁾ J. E. Knowles, Proc. Phys. Soc. (London) 78, 233, 1961.

¹¹⁾ D. J. Craik, Brit. J. appl. Phys. 11, 310, 1960.

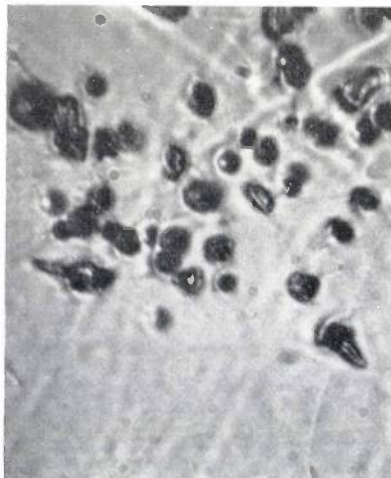
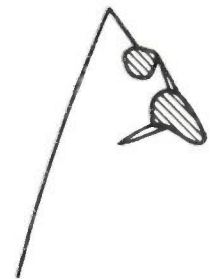
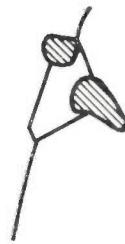
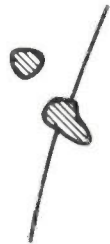


a

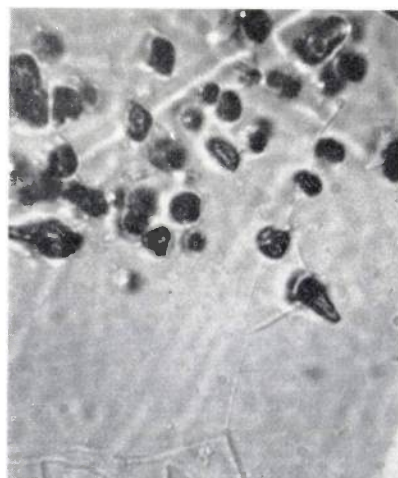


b

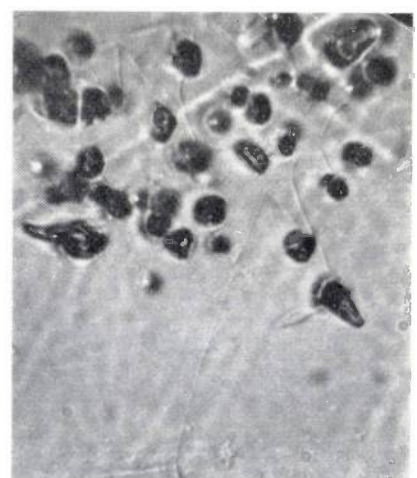
Fig. 11. *a*) A mosaic photograph of a grain of manganese-magnesium ferrite showing a 180° domain wall and Néel spike domains around the pores. For clarity the walls and pores are sketched in *b*.



a



b



c

Fig. 12. The progress of a domain wall through the encircled pore of fig. 11.

domains whose magnetizations make an angle of 180° with each other — joined by short sections of a 71° or 109° wall; this can be deduced from the angles that the parts make with each other. Usually a grain possessed a more complicated domain pattern than that illustrated, but still comprised of 180° walls in conjunction with short closure sections.

When the field was cycled slowly through 240 A/m (± 3 Oe), the wall area was indeed observed to pass through a maximum value when the applied field was more or less equal to the coercive force of the material.

The variation of the domain wall area does not seriously modify the proposed model. Each grain may still be considered to have a rectangular hysteresis loop, and the considerations regarding the orientation of the magnetization vectors still apply. The effect of the increase in wall area upon the theoretical loops for the polycrystalline material will be to make the sides of the loops more nearly vertical.

It was evident that, as suggested in ¹⁾, the pores played a large part in obstructing the free movement of the domain walls. In *fig. 12* a domain wall is shown passing through the encircled pore of *fig. 11*. In zero field (*fig. 12a*) the wall passed straight through the pore, thereby greatly reducing the pores's demagnetizing energy. As an increasing field was applied, the wall first bowed round the pore and then formed the classic configuration, first described by Néel ¹²⁾, (*fig. 12b*; see also the sketch). When the field was further increased to 152 A/m (1.9 Oe) the wall broke free, leaving behind Néel spikes (*fig. 12c*). The photographs also indicate that the wall interacted strongly with the grain boundary, and more particularly with the pores in it.

The determination of domain wall velocity

The time taken for a square-loop ferrite to reverse its magnetization when driven by a pulse field, i.e. its switching time, is determined by the velocity with which the domain walls travel, and by the average distance through which they have to move ¹³⁾. Measurements of wall velocity had earlier been made upon single crystal samples. It is also possible, by using a Bitter pattern technique, to make an estimate of the wall velocity in a single grain of a polycrystalline sample ¹⁰⁾. A similar investigation was carried out on a grain with a domain pattern and an orientation similar to that of *fig. 6*.

¹²⁾ L. Néel, *Cah. Phys.* No. 25, 1944.

¹³⁾ In large pulse fields the magnetization may reverse by rotation within each grain rather than by domain wall motion. See, for example, F. B. Humphrey and E. M. Gyorgy, *J. appl. Phys.* 30, 935, 1959.

This grain was found in the surface of a very small ring (3 mm diameter) which had been prepared by the pressing and polishing method described above. A single magnetizing pulse of known duration and amplitude was then applied, and the distance which the domain wall moved in consequence of this pulse was observed. The domain wall velocity corresponding to a given pulse height was determined by dividing this distance by the total pulse duration. Obviously, the method is not an accurate one and certain precautions must be taken if gross errors are to be avoided. A series of measurements on several grains yielded wall velocities varying by a factor of 3. The most consistent set of observations showed the following relationship between the velocity v (in cm/sec) and the applied field H (in A/m).

$$v = 12(H - 50).$$

This relation is represented by the straight line in *fig. 13*. The figure 50 A/m (0.62 Oe) in the factor between the brackets may be considered to be the grain coercive force.

These values for the domain wall velocity have been independently confirmed by Barkhausen effect measurements on the same material ¹⁴⁾.

Discussion

As stated earlier, in a material which is to have a high remanence the magnetocrystalline anisotropy must predominate over the stress and shape anisotropies. If in addition the material is to have a rectangular hysteresis loop, the individual grains

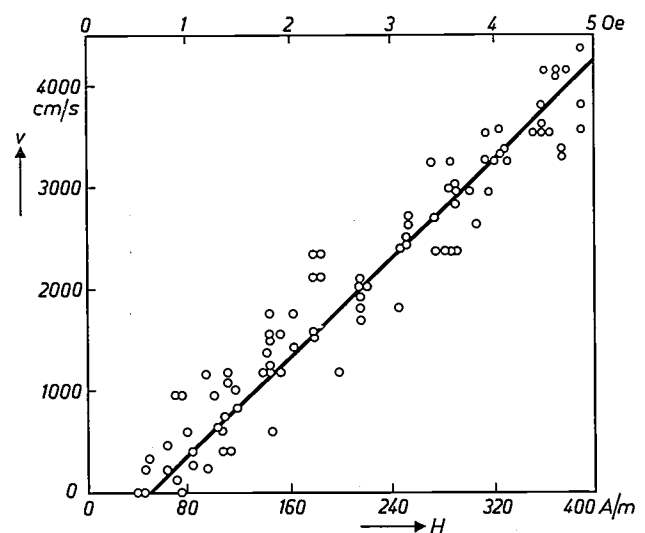


Fig. 13. The wall velocity v , calculated from direct measurements on a single grain, as a function of applied field H for a manganese-magnesium ferrite.

¹⁴⁾ J. Roche (Southampton University), private communication.

should, on average, have rectangular loops with a small spread in grain coercive force. If the first condition is violated the effect is to lower the remanence and produce a loop with long "tails". When on the other hand this condition is satisfied, but there is a large spread in the grain coercive force, loops having a high remanence with round "knees" and more sloping sides are obtained. These remarks, however, apply only to the *intrinsic loop* of the material. This is the hysteresis loop that would be found experimentally with a ring being very thin in the radial direction. The effect of a thick ring is similar to that of increasing the spread in the grain coercive force, thus worsening the squareness ratio; for in practice the field is obtained from a current through a straight conductor placed along the axis and is inversely proportional to the radius of the ring.

It would appear from the work described here that the coercive force of the grains probably arises from the pores. It was estimated (with a large margin of error) that a typical manganese-magnesium ferrite possesses approximately 20 to 200 pores per grain. These figures are not inconsistent with the earlier mentioned figure of twenty obstacles encountered by a domain wall during the course of a magnetization reversal. Obviously, the wall will in some positions intersect perhaps only one pore, i.e. a single small obstacle, and in other positions it will intersect several pores simultaneously, which will behave as one large obstacle. In the above it is of course assumed that a single pore can impede the motion of the whole domain wall. This is very probably the case, since the crystalline anisotropy constant for manganese-magnesium ferrites⁹⁾ is fairly large — about 5×10^4 ergs/cm³ — and thus makes the bending of a domain wall energetically very unfavourable. As the grain coercive force varies only very slowly with the number of obstacles (assuming there are more than about twenty per grain), a spread in the size of the grains will not have a first order effect on the squareness ratio.

Nevertheless, to obtain optimum performance from ferrites in terms of squareness ratio and switching time, uniformity of grain size would seem to be an advantage. On the one side the switching time of the material is required to be short. This will tend to be the case if the distance between the walls is small. An idealized model indicates that this distance will vary only as the square root of the dimensions of the grain¹⁵⁾, so that the grain size

¹⁵⁾ C. Kittel, Rev. mod. Phys. 21, 541, 1949.

would have to be radically reduced to gain significantly in this respect. But if the grain size is greatly decreased, the squareness ratio may deteriorate since the grains will contain so few obstacles that the "average" hysteresis loops, mentioned earlier, cannot be rectangular. There is thus an optimum grain size which corresponds to grains supporting one or two walls. Not surprisingly, commercial square-loop ferrites are found to have grains showing from one to four walls of the type shown in fig. 5.

The other factor governing switching time is the domain wall velocity. This is inversely proportional to a viscous damping factor. The nature of this damping, which is not well understood, has some relation to the damping associated with ferromagnetic resonance. Further, the switching time is proportional to \sqrt{K} , where K is the magnetocrystalline anisotropy¹⁶⁾. The domain wall velocity might thus be increased by decreasing K , but this would have a deleterious effect on the squareness ratio, unless the stress and shape anisotropies could be similarly reduced. Thus at the present time it is difficult to see how a really substantial improvement in the performance of "square loop" materials can be obtained.

Finally, it must be remarked that it would be of great interest if a pore-free polycrystalline ferrite of a standard square-loop composition could be prepared. On the basis of the theories described here, such a material should have a high remanence and a coercive force much less than that of the normal product. The hysteresis loop would not be very rectangular, but the switching time would be extremely short since the motion of the domain walls would no longer be impeded by pores.

¹⁶⁾ J. K. Galt, Phys. Rev. 85, 664, 1952.

Summary It was confirmed experimentally that if a ferrite is to have an intrinsically rectangular hysteresis loop, then the direction of the magnetization in each grain of the ferrite must be determined by the magnetocrystalline anisotropy, and not by stress or shape anisotropies. Domain pattern studies suggested a theoretical model for the magnetization reversal process, which assumed that the hysteresis loop of each grain was rectangular, and that the coercive force varied in a random way from grain to grain. For a rectangular loop it is desirable that all grains should have nearly equal coercive forces. Further studies of domain movements under the action of a magnetic field showed that pores in the material act as obstacles to the free movement of the domain walls: the pores may thus determine the nature of the coercive force. In one grain the domain wall velocity as a function of the applied field was also determined, and was found to correspond to the value deduced from an independent method.

HYDROGEN IN IRON AND STEEL

II. FRACTURING

by J. D. FAST *) and D. J. van OOIJEN *).

539.42:546.11:669.14

The first part of this article made it clear that if hydrogen is present in iron and steel it may cause various harmful effects. Hydrogen taken up during the enamelling of steel may cause cracking of the enamel layer, and hydrogen that penetrates during pickling may give rise to surface blistering. Part II below offers some insight into the origin of the often highly dangerous fracturing of iron and steel under the influence of hydrogen.

On the origin of fractures

During plastic deformation of a metal there arise not only lattice imperfections in the form of dislocations and point defects but also crack nuclei which, under unfavourable conditions, may develop into real cracks. There are various theories on the origin of these nuclei. The hypothesis common to them all is that each crack nucleus is produced by the piling up and coalescence of a number of dislocations under the influence of external shear stresses. In the picture which Zener and Stroh¹⁾ give of this process, the dislocations moving along a slip plane are piled-up against some obstacle or other, e.g. a grain boundary or inclusion. The dislocations at the head of such a piled-up group experience considerable pressure from the dislocations coming along behind them. As a result they can be forced so close together as to merge to form a single dislocation having a large Burgers vector (see *fig. 1*). After exceeding a certain value of the Burgers vector a wedge-shaped void is formed, which can act as a crack nucleus.

For iron Cottrell²⁾ suggested a somewhat different mechanism, in which dislocations move towards each other along two *intersecting* slip planes and coalesce along the junction of the two planes. This gives rise to a wedge-shaped crack nucleus in a cleavage plane (see *fig. 2*). This nucleus will be larger the more dislocations are involved in the coalescence³⁾.

Various theories also exist as to the factors that play an essential part in the catastrophic growth of

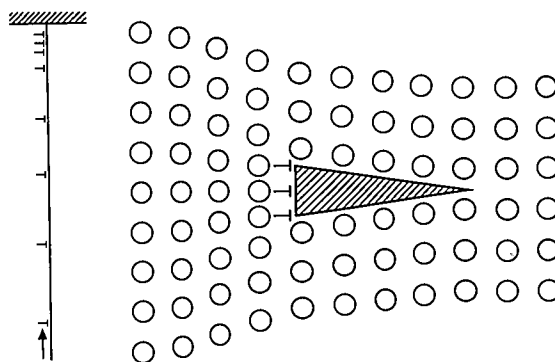


Fig. 1. Illustrating the mechanism by which a crack nucleus forms, as proposed by Zener and Stroh¹⁾. The dislocations move in the direction of the arrow along a slip plane (left) and the first are stopped by an obstacle, such as a grain boundary or inclusion. Under pressure from the following dislocations, three of them coalesce to form a wedge-shaped void — a crack nucleus (right).

a crack nucleus. The formation of a crack nucleus is always accompanied by the building-up of large stresses in its vicinity. The larger these are, the greater is of course the chance that the nucleus will develop into a macrocrack. Among the influential factors is the behaviour of the neighbouring *dislocation sources*. If these sources are easily activated, the plastic deformation will continue and the stresses will diminish before macrocracks have been able to form. If, however, the sources are strongly pinned by foreign atoms, there is a chance that the sources will not become active early enough, with the result that cracks appear. It is known that nitro-

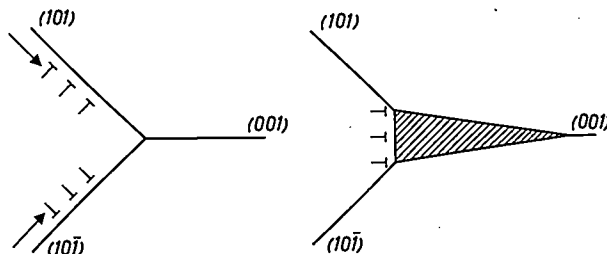


Fig. 2. According to Cottrell²⁾ crack nuclei form in iron on $\{100\}$ planes (the cleavage planes) by the coalescence of dislocations moving towards each other along two intersecting $\{110\}$ planes (the slip planes).

*) Philips Research Laboratories, Eindhoven.

¹⁾ C. Zener, *Fracturing of metals*, Amer. Soc. Metals, Cleveland (Ohio) 1948, pp. 3-31; A. N. Stroh, *Adv. Phys.* 6, 418, 1957.

²⁾ A. H. Cottrell, *Trans. AIME* 212, 192, 1958; see also his article in the congress book "Fracture", Proc. internat. Conf. on the atomic mechanisms of fracture, Swampscott (Mass.) 1959.

³⁾ For a refinement of this theory, see A. W. Sleeswijk, *Twinning and the origin of cleavage nuclei in a iron*, *Acta metalurgica* 10, 803-812, 1962 (No. 9).

gen and carbon may cause embrittlement in this way, and it has been assumed that the same applies to the embrittling effect of hydrogen in steel.

Apart from this hypothesis there are two others worth mentioning. In Part I of this article ⁴) it was stated that high-pressure molecular hydrogen may form in microcavities of iron or steel which is supersaturated with hydrogen. It has been assumed that H₂ might form similarly in any crack nucleus. The resultant internal pressure increases the chance of the nucleus growing catastrophically into a macrocrack.

A third hypothesis concerning the embrittling behaviour of hydrogen envisages an adsorption effect. The adsorption of hydrogen slightly reduces the surface tension of iron, as a result of which the formation of new surfaces — in this case the growth of crack nuclei — will cost less energy in steel containing hydrogen than in hydrogen-free steel.

The question arises to what extent the three above-mentioned effects of hydrogen might contribute to the embrittlement of steel. In an attempt to answer this question we measured the electrical resistance of soft-annealed and also of plastically deformed iron wire before and after charging with hydrogen, on the underlying assumption that the electrical resistance of hydrogen-charged iron wire depends on the form in which the hydrogen is present in the metal, i.e. whether it is dissolved interstitially, whether it is bound to dislocations, and so on. To prevent interference from inclusions and other impurities, our measurements were made on very pure iron ⁵).

The results of this investigation, which will be discussed under the next heading, not only lead to important conclusions with regard to the mechanism of the embrittling behaviour of hydrogen in steel, but also throws new light on the work of other research workers. At the end of this article we shall discuss the consequences of our findings, and consider among other things how the likelihood of fracture is influenced by the rate of deformation and by the temperature. In conclusion we shall examine that dangerous form of steel failure known as delayed failure.

Effects of the electrolytic charging of iron with hydrogen

It is known that the electrical resistivity of iron that contains both dislocations and carbon atoms is higher in the state in which these imperfections

occur separately than in the state in which they are bound to one another ⁶).

Originally our measurements of the electrical resistance of pure iron wires electrolytically charged with hydrogen (in some cases preceded by plastic deformation), seemed to indicate that there is also strong interaction of hydrogen and dislocations. Further investigation, however, revealed the complicating factor that charging with hydrogen leads to macroscopically measurable, permanent changes in the dimensions of the wires ⁷). Examination under a light microscope showed that these dimensional changes are caused by microcavities and microcracks along the grain boundaries (see figs 3 and 4).

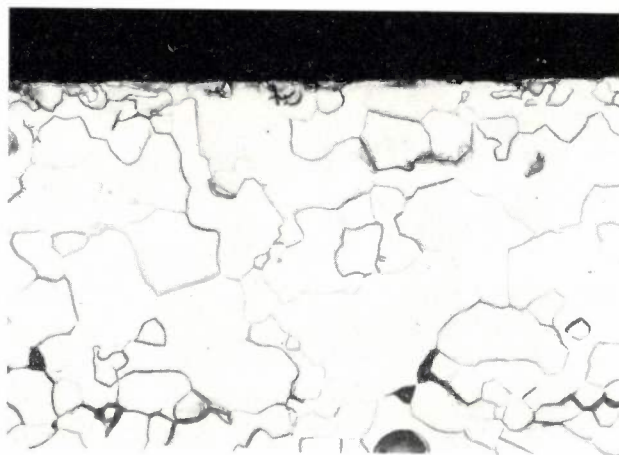


Fig. 3. Soft-annealed iron wire after electrolytic charging with hydrogen. Cracks have formed at grain boundaries owing to the evolution of molecular hydrogen of high pressure. The cracks and grain boundaries were made visible by polishing and etching a longitudinal cross-section. Magnification 100 ×.

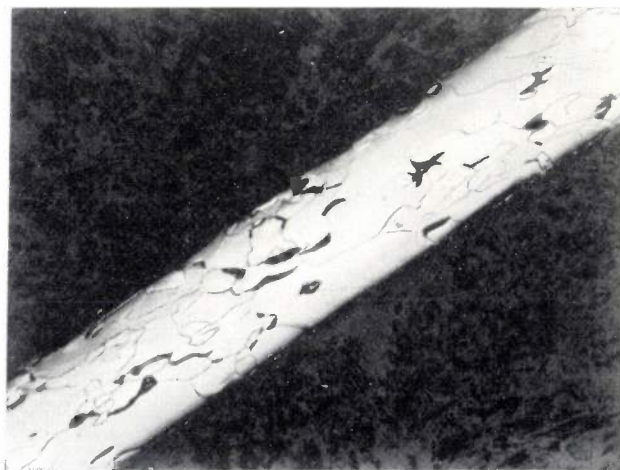


Fig. 4. Cold-worked iron wire after electrolytic charging with hydrogen. The cracks generated, made visible in the same way as in fig. 3, lie here mainly parallel to the long axis of the wire. Magnification 50 ×.

⁴) J. D. Fast and D. J. van Ooijen, Hydrogen in iron and steel, I. Solution and precipitation, Philips tech. Rev. **24**, 221-227, 1962/63 (No. 7).

⁵) Made by the method described by J. D. Fast, A. I. Luteijn and E. Overbosch, Philips tech. Rev. **15**, 114-121, 1953/54.

⁶) A. B. Bhatia, Proc. Phys. Soc. **B 62**, 229, 1949; A. H. Cottrell and A. T. Churchman, J. Iron Steel Inst. **162**, 271, 1949.

⁷) For further details of these experiments, see: D. J. van Ooijen and J. D. Fast, Electrical resistance of hydrogen-charged wires, Acta metallurgica **11**, 211-216, 1963 (No. 3).

The only possible explanation for their presence is that H₂ molecules were formed at these boundaries, giving rise to local gas pressures which exceeded the local cohesive strength of the metal. These H₂ molecules are most likely to form at high-angle grain boundaries, i.e. at the boundaries between crystals that differ considerably in orientation. Fig. 5a shows a model of such a boundary; in fig. 5b can be seen a model of a boundary between crystals differing only slightly in orientation.

Description and discussion of some of our own experiments

Fig. 6 shows for two iron wires (the same as depicted in figs 3 and 4) the relative resistance increments $(\Delta R/R)_{77}$ measured at 77 °K as a function of the charging time (at room temperature). One wire was a soft-annealed wire of 0.5 mm diameter, the

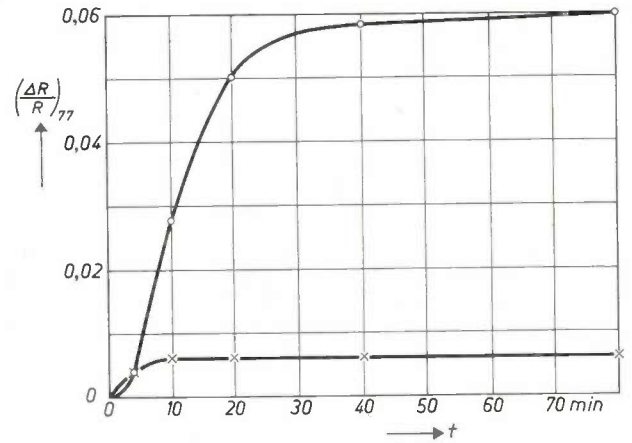


Fig. 6. Relative increase of resistance $(\Delta R/R)_{77}$ at 77 °K as a function of charging time t for soft-annealed and cold-worked iron wire (circles and crosses respectively). The much smaller resistance increase of the cold-worked wire is shown to be due not so much to the interaction of hydrogen and dislocations, as we originally thought, but to the microcracks generated, visible in figs. 3 and 4.

The diameter of the first wire was 0.5 mm; the second wire was drawn from 0.5 to 0.3 mm diameter. The wires were charged at room temperature and at constant current density. The quantity of hydrogen ultimately taken up was about 10⁻³ at. H/at. Fe.

other a wire drawn at room temperature from 0.5 mm to 0.3 mm diameter. The dislocation density in the latter ("cold-worked") wire was considerably greater than in the soft-annealed specimen (the effect of cold-working being to introduce dislocations, whereas soft-annealing removes them). The ultimate percentage of absorbed hydrogen was found by vacuum extraction to be about 0.1 at.% for both wires.

From fig. 6 it can be seen that the relative resistance increase of the soft-annealed wire is many times greater than that of the cold-worked wire. Until we discovered the formation of microcracks along the grain boundaries, it seemed obvious to interpret these observations in the same way as the above-mentioned results found for the presence of carbon instead of hydrogen. Our original conclusion was therefore that the resistivity is increased by interstitially dissolved hydrogen, and that this increase is lower the more dissolved hydrogen is bound to dislocations.

However, after we had discovered that charging with hydrogen causes changes in dimensions, it was plain that this conclusion was no longer valid. Since not only any dissolved hydrogen but also the dimensional changes will affect the resistance of the wires, we must first *eliminate from our results the influence of the changes in dimensions*.

The resistance R of a wire is given by the formula:

$$R = \rho \frac{l}{A}, \dots \dots \dots (1)$$

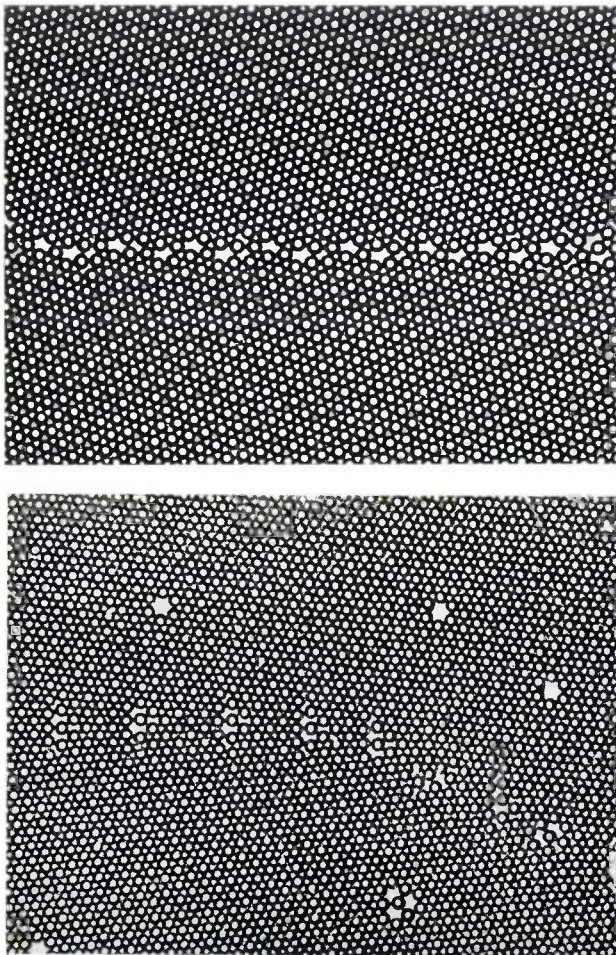


Fig. 5. a) Model of a high-angle grain boundary, obtained using soap bubbles. The boundary is one between two "crystals" that differ considerably in orientation. For iron containing hydrogen it is thought probable that molecular hydrogen forms in the cavities of atomic dimensions existing at such boundaries. b) By the soap bubble method it is also a simple matter to bring together two "crystals" with a small difference in orientation. The difference is then bridged by a series of dislocations. (After W. M. Lomer and J. F. Nye, Proc. Roy. Soc. (London) A 212, 576, 1952.)

where ρ is the resistivity, l is the length and A the cross-sectional area of the wire. The repeatedly occurring quotient l/A will be denoted henceforth by G (from "geometry"). After charging, the increased resistance can be written as:

$$R + \Delta R = (\rho + \Delta\rho)(G + \Delta G) \dots (2)$$

The relative resistance increment is then given by:

$$\Delta R/R = \Delta\rho/\rho + \Delta G/G \dots (3)$$

The two contributions $\Delta\rho/\rho$ and $\Delta G/G$, from which the relative resistance increment is built up, can be separately determined by performing the resistance measurements at two temperatures both before and after charging. The values of $\Delta G/G$ and $(\Delta\rho/\rho)_{77}$ found in this way are presented in Table I for a temperature of 77 °K.

The proof of the above is as follows. Before charging, the difference of the resistances at the temperatures T_1 and T_2 is given by the equation:

$$R(T_1) - R(T_2) = \{\rho(T_1) - \rho(T_2)\} G, \dots (4)$$

where G is regarded as independent of temperature. This is permissible because the change of G , due to thermal expansion, in the transition from T_2 (approx. 77 °K) to T_1 (approx. 300 °K), is very small compared with the change of G resulting from the hydrogen-charging. After charging with hydrogen, we have:

$$R'(T_1) - R'(T_2) = \{(\rho + \Delta\rho)_{T_1} - (\rho + \Delta\rho)_{T_2}\} (G + \Delta G) \dots (5)$$

According to Matthiessen's rule, the increase in the resistivity of a metal, caused by foreign atoms and other lattice imperfections (dislocations, vacancies, etc.) is independent of temperature. Assuming that this rule also applies to the case under consideration (in the article cited under 7) we demonstrate that this is in fact the case, we can write (5) as:

$$R'(T_1) - R'(T_2) = \{\rho(T_1) - \rho(T_2)\} (G + \Delta G) \dots (6)$$

Dividing (6) by (4) yields:

$$\frac{R'(T_1) - R'(T_2)}{R(T_1) - R(T_2)} = 1 + \frac{\Delta G}{G} \dots (7)$$

We can thus calculate $\Delta G/G$ from the results of the measurements. Since $\Delta R/R$ has been measured, $\Delta\rho/\rho$ can then be found with the aid of equation (3).

In reality, $\Delta G/G$ and $\Delta\rho/\rho$ are determined by a somewhat different procedure, since after charging it is difficult to obtain exactly the same temperatures T_1 and T_2 as before charging; for this procedure see 7).

Compared with fig. 6, Table I gives a much less pronounced indication of the existence of an interaction between hydrogen and dislocations. This table shows that the difference between the increases of resistance of the two wires is not primarily due to a difference in the resistivity changes but to the difference in the dimensional changes ($\Delta G/G$). This difference can be explained from the fact that the microcracks in the cold-worked wire do not have

the same orientation as in the soft-annealed wire. As a result of drawing, the grain boundaries of the cold-worked wire are roughly parallel to the long axis of the wire, and so too are the majority of microcracks produced by charging with hydrogen. This can be seen in fig. 4, which presents a longitudinal cross-section of the cold-worked wire after charging. In the soft-annealed wire the cracks are randomly oriented with respect to the axis of the wire (fig. 3), and therefore have a much greater influence on the resistance.

After the effect of the dimensional changes has been eliminated, there remains a relatively small difference in $(\Delta\rho/\rho)_{77}$ between the two wires which might be thought to be attributable to the interaction of hydrogen with dislocations. The following experiment, however, makes even this seem rather improbable.

We kept the soft-annealed wire from the previous experiment for 24 hours at room temperature. During this time the wire lost at least 95% of its hydrogen without showing any change in $(\Delta\rho/\rho)_{77}$. Annealing for two hours at 350 °C was necessary to decrease the value of $(\Delta\rho/\rho)_{77}$ appreciably, viz from 1.1% (Table I) to 0.6%. However, the same decrease is found if an iron wire completely free of hydrogen is given the same heat-treatment after it has been plastically deformed to such a degree as to show the same value of 1.1% for $(\Delta\rho/\rho)_{77}$.

The higher resistivity measured after charging with hydrogen therefore seems to be due not so much to dissolved hydrogen but rather to the plastic deformation of the lattice in the neighbourhood of the microcracks. One can appreciate that this effect would be more pronounced in the soft-annealed wire, which contained far fewer dislocations before charging than the cold-worked wire. This offers a satisfactory explanation for the differences in $(\Delta\rho/\rho)_{77}$ between the two wires, although one cannot entirely exclude the possibility that dissolved hydrogen also plays some part in this connection, though an insignificant one.

Table I. Relative increase of the resistance of pure iron due to hydrogen-charging (one hour in 0.1 n H₂SO₄ doped with 50 mg/litre As₂O₃, to promote hydrogen absorption, at a current density of 0.12 A/cm²) for a soft-annealed and a cold-worked wire. $(\Delta R/R)_{77}$ is the total change at 77 °K; $\Delta G/G$ is the contribution from changes in dimensions; $(\Delta\rho/\rho)_{77}$ is the contribution from changes in resistivity.

	$(\Delta R/R)_{77}$ %	$\Delta G/G$ %	$(\Delta\rho/\rho)_{77}$ %
Soft-annealed	5.40	4.27	1.08
Cold-worked	0.86	0.37	0.49

Let us now turn to the increase of resistance due to changes in dimensions. This increase is a kind of "shadow effect" since it depends on the fact that no current is carried by small regions of the conductor just in front of and just behind a cavity (seen in the direction of the current). The relative increase of resistance of a wire caused by spherical cavities having a relative volume $\Delta V/V$ follows from calculations by Landauer⁸⁾:

$$\frac{\Delta R}{R'} = 1.5 \frac{\Delta V}{V}, \quad \dots \dots \dots (8)$$

where R' is the resistance of the same wire in the absence of the cavities. For the soft-annealed wire $\Delta V/V$ was found to be 2.4% (increase of diameter 1.2%, increase in length negligible). Randomly-orientated lenticular cavities, having a relative volume of 2.4%, will cause $\Delta R/R'$ to be greater than given by equation (8): $(\Delta R/R')_{\text{theor}} > 3.6\%$ with respect to a wire with the same (enlarged) diameter but without cavities. The experimental value of the relative increase of resistance of the soft-annealed wire caused by cavities ($\Delta G/G = 4.27\%$; Table I) relates to the original diameter of the wire. Its value with respect to the enlarged diameter is obtained by correcting for the observed increase in cross-sectional area: $(\Delta R/R')_{\text{exp}} = (4.27 + 2.4)\% = 6.7\%$. The agreement between the theoretical and experimental values is not unsatisfactory.

Experiments of other research workers

From the above experiments and from others carried out by us⁷⁾ it may be concluded that the change in the resistance of an iron wire, measured after charging with hydrogen, is not or to no significant extent, due to the presence of dissolved hydrogen. The cause of the change of resistance is rather the permanent damage which precipitated molecular hydrogen inflicts on the metal in the form of microcracks and plastic deformation. In complete agreement with this conclusion are the results of an investigation into the cause of the broadening of X-ray diffraction lines after electrolytically charging iron with hydrogen⁹⁾. It has been found that this line-broadening agrees entirely with that produced by a few per cent cold working. The "recovery" (elimination of the line-broadening) during heating at 425 °C or 475 °C follows the same pattern in both cases and proceeds with the same activation energy. The latter is very much greater than that of the diffusion of hydrogen in iron, and greater too than that of the surface reactions involved in the escape of hydrogen from iron. This demonstrates that the recovery in question does not depend on the expulsion of hydrogen.

It is interesting to note that the occurrence of permanent damage to iron after charging with

hydrogen was inferred from magnetic measurements as long ago as thirty years by Reber¹⁰⁾. He electrolytically charged flat rings of magnetically soft iron with hydrogen and observed that this caused a considerable increase in magnetic hardness (the maximum permeability and remanent magnetization dropped, and the coercivity rose). Expulsion of the hydrogen left the change in the magnetic properties virtually unaffected.

The cracks found by us along the grain boundaries of iron as a result of electrolytic charging with hydrogen also provide an explanation for certain remarkable phenomena found by other research workers. As an example we mention an extensive investigation undertaken by Simone Besnard¹¹⁾. She charged iron electrolytically with hydrogen from a bath to which Na_2S had been added to promote the take-up of hydrogen. After charging, it was possible to demonstrate the presence of sulphur along the grain boundaries to a considerable depth inside the metal. The phenomenon was carefully studied by partly replacing the sulphur in Na_2S by a radioisotope (sulphur 35). No satisfactory explanation was given, however, for the penetration of sulphur. Our own experiments immediately suggest the explanation that in Besnard's experiments liquid from the electrolytic bath entered the metal via cracks along the grain boundaries.

It was long ago demonstrated by Bardenheuer *et al.*¹²⁾ that hydrogen may cause cracks along grain boundaries in *commercial steel*. They offered an explanation, for example, of the damage produced in steel upon dip-soldering in molten brass, after the steel had been pickled to obtain a clean surface. When the steel is dipped in molten brass, the hydrogen taken up during the pickling process is expelled so rapidly that large cracks form along the grain boundaries, thus enabling the brass to penetrate deep into the steel (*fig. 7*).

What conclusions can we now draw from the foregoing regarding the embrittling behaviour of hydrogen in iron and steel? May we assume, as we did earlier, that this embrittlement is a consequence of interaction between hydrogen and dislocations? From our own experiments, which gave no indications of any strong interaction, this assumption does not seem to be warranted. It seems far more reasonable to suppose that the embrittlement is the conse-

⁸⁾ R. Landauer, *J. appl. Phys.* **23**, 779, 1952.

⁹⁾ A. S. Tetelman, C. N. J. Wagner and W. D. Robertson, *Acta metallurgica* **9**, 205, 1961 (No. 3).

¹⁰⁾ R. K. Reber, *Physics* **5**, 297, 1934.

¹¹⁾ S. Besnard, *Ann. de Chimie* **6**, 245, 1961; S. Besnard and J. Talbot, *Colloque sur la diffusion à l'état solide*, Saclay, 1958, North Holland Publ. Co., Amsterdam 1959, p. 147.

¹²⁾ P. Bardenheuer and H. Ploum, *Mitt. K. W. I. Eisenforsch.* **16**, 129 and 137, 1934; P. Bardenheuer, *Metall* **6**, 351, 1952.

quence of molecular hydrogen forming at high pressure in crack nuclei. This possibility is very plausible now that we know that hydrogen is evolved not only in relatively large internal cavities, but also in the very small cavities existing at grain boundaries. As regards the third possibility referred to, namely that the embrittlement is due to adsorption of hydrogen, we can only say here that this adsorption can be shown on various grounds to have only an insignificant influence¹³).

The conclusion that the formation of molecular hydrogen in crack nuclei is by far the most important of the three effects mentioned, clarifies much that was previously difficult to understand.

standable: if hydrogen is to exercise its harmful action by forming high pressure H_2 in the crack nuclei (produced *during* deformation), the rate of deformation must be low enough and the temperature high enough to give the H atoms an opportunity to reach the crack nuclei by diffusion.

For some part these hydrogen atoms are perhaps only *formed* during plastic deformation, from hydrogen molecules that had previously precipitated in pores or other defects. This view is suggested from experiments of Hofmann *et al.*¹⁴) on the ductility of unalloyed steel (0.22% C) in air and in hydrogen. Tensile test bars of this metal that exhibited high ductility in a normal tensile test in air were found

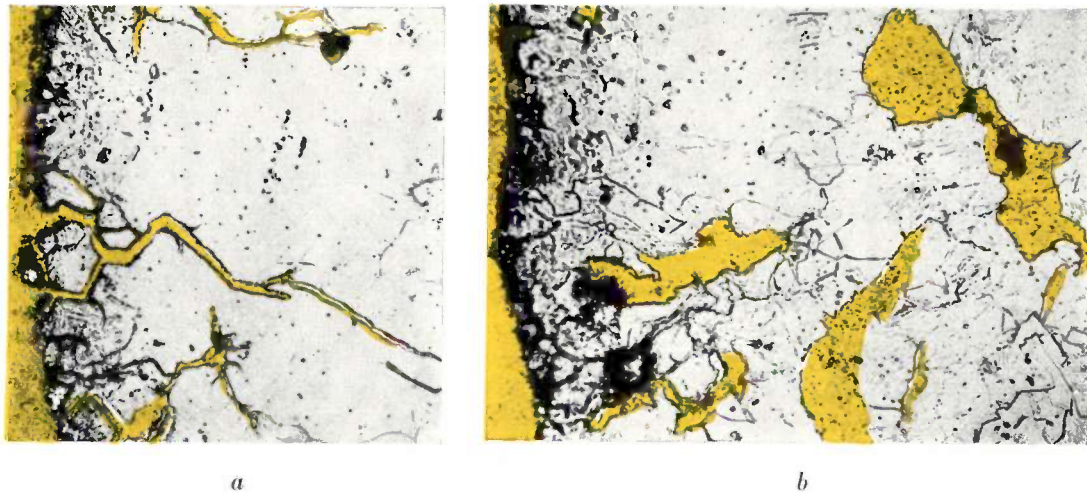


Fig. 7. Hydrogen-charged steel wire after dipping in molten brass (magnification 500 \times). The brass (yellow) penetrates deep into the steel along micro-fissures (a), and fills internal cavities (b). (After Bardenheuer and Ploum¹².)

The influence of hydrogen on the ductility and brittle fracturing of steel

If a steel contains no hydrogen, the chance of brittle fracturing is greater the higher the rate of deformation and the lower the temperature. The explanation is found in the fact that the dislocations in steel are as a rule pinned by foreign atoms (e.g. nitrogen atoms), and that the breaking away of these dislocations is a thermally activated process (see the discussion on page 253).

Where the steel contains hydrogen, however, it may show brittleness which (within limits) is more pronounced the *lower* the rate of deformation and the *higher* the temperature. This is now under-

to have very low ductility in an atmosphere of pure hydrogen at high pressure. The experiments demonstrate that the metal absorbed hydrogen *during plastic deformation* in that gas. This makes it reasonable to assume that also the high-pressure hydrogen contained in the pores may dissolve during plastic deformation and later precipitate in crack nuclei.

In the following we shall consider experiments of various research workers that illustrate the points just discussed. There can be no question, however, of giving anything like a complete survey of the numerous investigations into the brittle fracturing of steel. The literature on this subject has assumed enormous proportions in the last twenty years,

¹³) See e.g. B. A. Bilby and J. Hewitt, *Acta metallurgica* **10**, 587, 1962 (No. 6).

¹⁴) W. Hofmann and W. Rauls, *Arch. Eisenhüttenw.* **32**, 169, 1961; W. Hofmann, W. Rauls and J. Vogt, *Acta metallurgica* **10**, 688, 1962 (No. 7).

largely owing to the occurrence of fractures in many welded-steel ships during and after the second world war. These fractures were frequently of a very serious kind, so serious in fact that several ships broke in two. A case in point, which attracted considerable attention at the time, was the tanker Schenectady which, in 1943, while lying at anchor off Portland quay (Oregon), suddenly snapped in two with a bang that could be heard several kilometres away. It is not known whether hydrogen played any part in this special case, but it is certainly not inconceivable.

Brittleness appearing under test

The following experiments are of special interest concerning the influence of the deformation rate and temperature on the ductility of steel containing hydrogen¹⁵). Test bars of a particular type of commercial steel (spheroidized SAE 1020 steel) were electrolytically charged with hydrogen for one hour in 4% H₂SO₄. After charging, the hydrogen content was about 10 cm³ per 100 grams of metal, which was not, however, distributed uniformly over the whole cross-section of the bar. Fig. 8 shows the

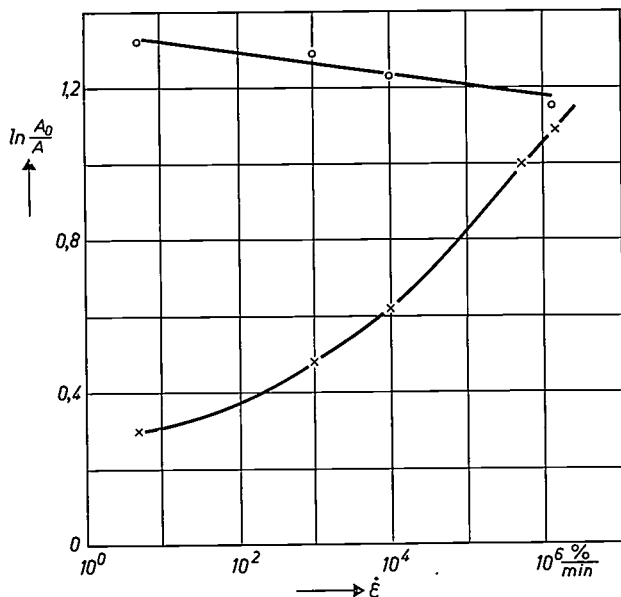


Fig. 8. Influence of deformation rate $\dot{\epsilon}$ on the ductility of uncharged (circles) and of hydrogen-charged steel (crosses) at room temperature. As the deformation rate increases, the ductility of the uncharged steel decreases, but that of the charged steel rises rapidly. The latter can be explained from the fact that at higher deformation rates the hydrogen has less time to diffuse to the crack nuclei.

As a measure of the ductility the true tensile strain $\ln A_0/A$ is plotted as ordinate, A_0 being the initial and A the final cross-section of the test bar at the position of fracture. (After Brown and Baldwin¹⁵.)

¹⁵) J. T. Brown and W. M. Baldwin, *Trans. AIME* **200**, 298, 1954; Taiji Toh and W. M. Baldwin, *Stress corrosion cracking and embrittlement* (editor W. D. Robertson), Wiley, New York 1956, p. 176.

ductility of charged and uncharged steel at room temperature as a function of the deformation rate. The measure of ductility adopted was the *true tensile strain*, given by the natural logarithm of A_0/A , where A_0 is the initial and A the final cross-section of the bar at the position of fracture. At low deformation rates the absorbed hydrogen has a marked embrittling effect. If the deformation rate is high enough, however, the charged steel exhibits the same ductility as the uncharged steel.

Fig. 9 gives a plot of the ductilities of charged and uncharged steel versus temperature for four different deformation rates. Figures 10a and b show how the ductility is affected both by temperature and the deformation rate (in a for uncharged and in b for charged steel). As can be seen, there is a region in fig. 10b, i.e. surface c, in which the ductility of the charged steel — unlike that of the uncharged steel — increases with increasing deformation rate and decreases as the temperature rises. In the light of the foregoing considerations, this behaviour is understandable. To the left of curve *i* the temperatures, and hence the hydrogen diffusion rates, are so low compared with the deformation rates that no appreciable H₂ pressures can form in the crack nuclei during the deformation. To the right of curve *i*, in surface c, the influence of hydrogen becomes more noticeable the higher is the temperature and the lower is the deformation rate.

According to figures 9 and 10, at more elevated temperatures a point is reached where the embrittling action of hydrogen gradually decreases again with increasing temperature. In surface *d* (fig. 10b) the ductility increases both with increasing deformation rate and with rising temperature. This is precisely as expected. In the first place the equilibrium H₂ pressure pertaining to a given hydrogen content drops rapidly as the temperature rises. In the second place the metal loses its hydrogen rapidly at temperatures above 100 °C. A consequence of these facts is that at a given deformation rate noticeable embrittlement occurs only in a *limited* temperature range.

In fig. 9 the limiting temperatures of the brittle zone correspond to the points G and G' (the latter point, at c, has been obtained by extrapolation). Fig. 11 gives these limiting temperatures as a function of the deformation rate for two hydrogen contents. The presence of the hydrogen is noticeable only in the area between the two branches of the curve.

Delayed failure

Particularly notorious are the fractures that may

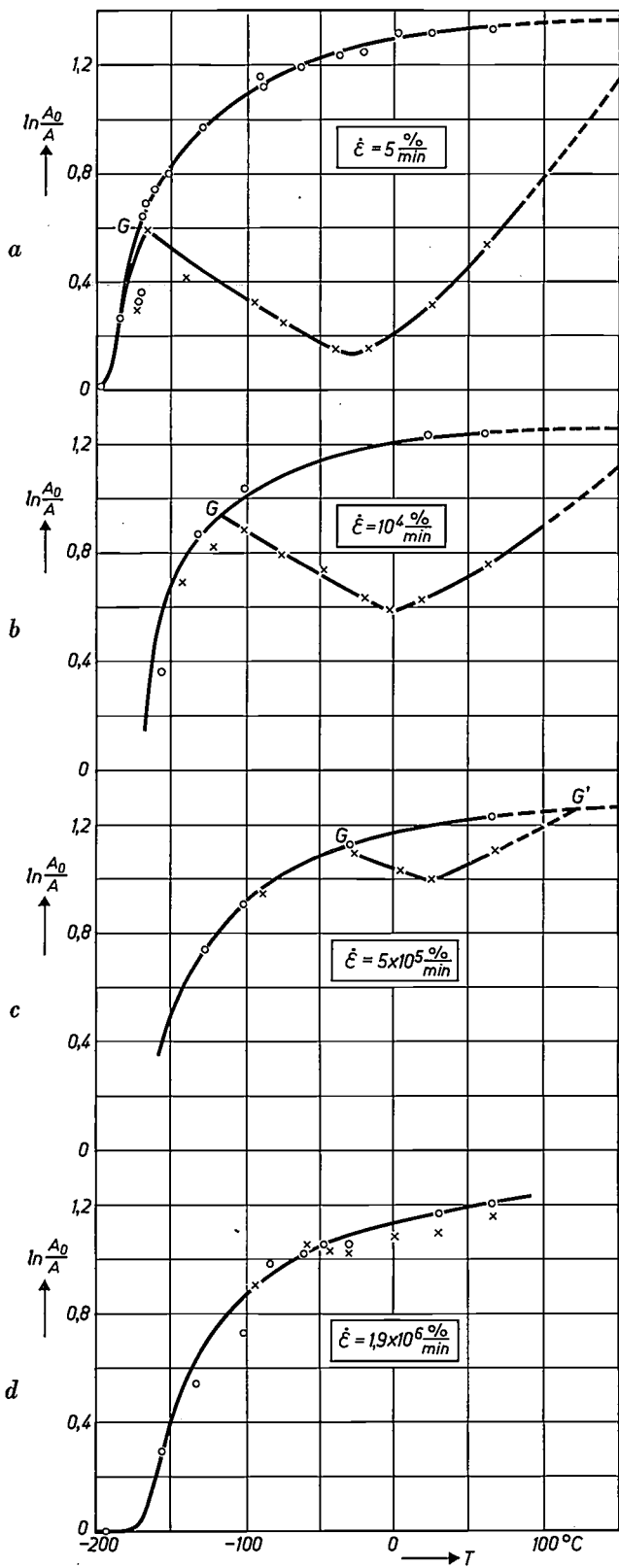


Fig. 9. True tensile strain $\ln A_0/A$ versus temperature at deformation rates of 5%/min (a), $10^4\%$ /min (b), $5 \times 10^5\%$ /min (c) and $1.9 \times 10^6\%$ /min (d). G and G' give the temperatures below which and above which there is no difference in ductility between uncharged (circles) and charged steel (crosses). At higher deformation rates the "brittle zone", due to the action of hydrogen, becomes smaller. (After Brown and Baldwin ¹⁵).

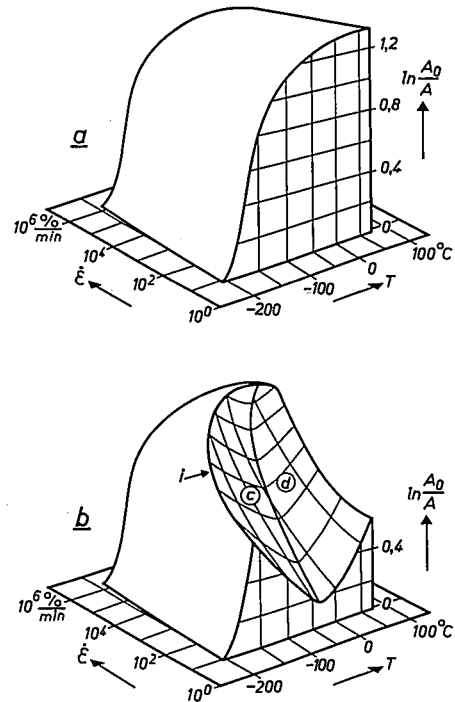


Fig. 10. True tensile strain $\ln A_0/A$ versus temperature T and deformation rate $\dot{\epsilon}$ for uncharged steel (a) and hydrogen-charged steel (b). Left of the curve i the charged steel behaves like uncharged steel, which is understandable since in this region the diffusion rates of hydrogen are low compared with the deformation rates and therefore high H_2 pressures have no time to form in the crack nuclei during deformation. Right of the curve i , in surface c , high pressures are able to form, so that brittleness increases with rising temperature. In surface d the brittleness decreases again with rising temperature. This is due to the fact that the equilibrium pressure pertaining to a given hydrogen content decreases rapidly as the temperature increases. Moreover, at these elevated temperatures the metal rapidly loses its hydrogen. (After Taiji Toh and Baldwin ¹⁵).

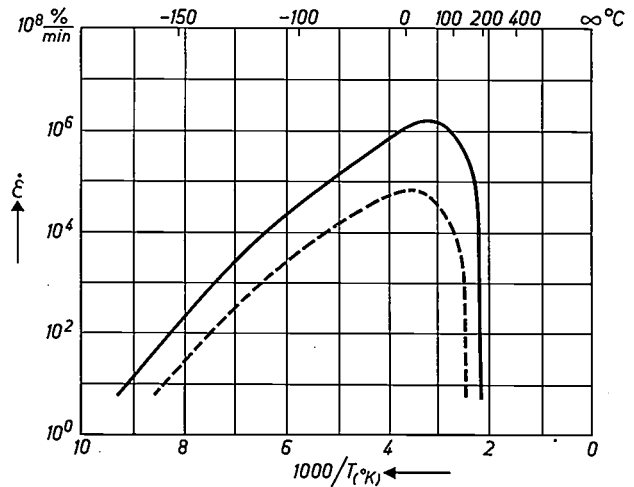


Fig. 11. The limiting temperatures below or above which uncharged and hydrogen-charged steel exhibit the same true tensile strain, as a function of deformation rate $\dot{\epsilon}$. For the solid line the charging time was 1 hour, for the broken line 6 minutes. The hydrogen makes its influence felt only in the areas enclosed by the curves. (After Taiji Toh and Baldwin ¹⁵).

occur in certain high-tensile steels (e.g. AISI 4340 and H11) as used in aircraft construction. To protect these metals against corrosion they were originally given an electrolytic coating of cadmium. Hydrogen taken up by the metal during this plating process gave rise in many cases to fractures when the metal was subjected to external or internal stresses. Characteristic of such fractures is firstly that they may occur under the influence of a static stress far below the yield point, and secondly that the fracturing is often preceded by a long delay period. These fractures can be avoided by ensuring that the metal can absorb no hydrogen, e.g. by applying the cadmium coating by vacuum evaporation.

It seems obvious to assume that the stresses will give rise to slight dislocation movements, leading to the formation of crack nuclei. If hydrogen is present, high-pressure molecular hydrogen can form in these nuclei, which may thus develop into actual cracks.

Of particular interest in this connection are the following experiments in which measurements were made of the electrical resistance of notched steel bars (AISI 4340) that were electrolytically charged with hydrogen and submitted to a constant load¹⁶. Fig. 12 presents the results of the measurements (increase of resistance as a function of time) for

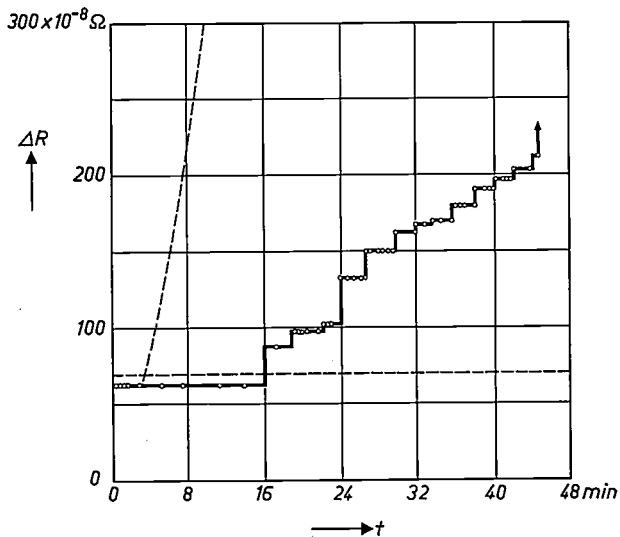


Fig. 12. Resistance increase ΔR of notched steel bars (AISI 4340) after electrolytic charging with hydrogen, versus time t at constant load (127 kg/mm^2) and a temperature of -18°C . After an incubation period of 16 minutes the resistance suddenly increases, presumably owing to the generation of a small crack, after which the resistance stays constant for a while. The phenomenon repeats itself several times before full fracture of the bar (at the arrow). The broken curves represent parts of the curves from figs 13 and 14. (The three figures 12, 13 and 14 are due to Steigerwald, Schaller and Troiano¹⁶.)

¹⁶ E. A. Steigerwald, F. W. Schaller and A.R. Troiano, Trans. AIME 215, 1048, 1959.

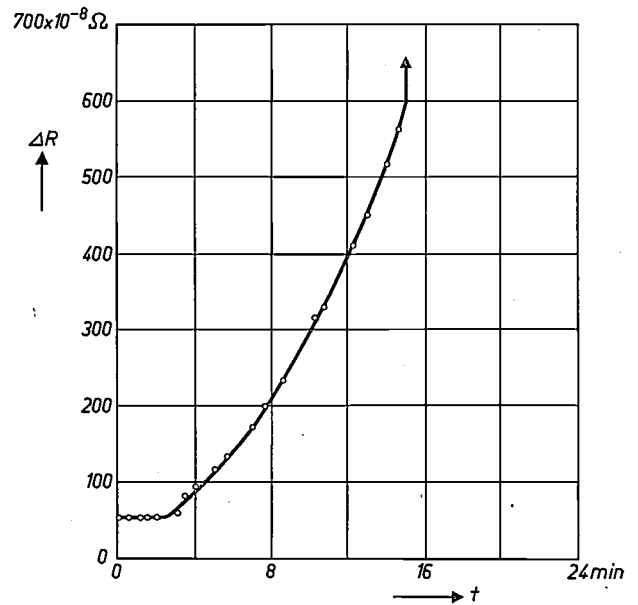


Fig. 13. The same experiment as in fig. 12, at a higher temperature of 27°C . The incubation time is here much shorter (about 3 minutes) after which the crack grows continuously.

a temperature of -18°C and a tensile stress of 127 kg/mm^2 . When the stress is applied the resistance increases by a certain amount and then remains constant for some time. After this incubation period there occurs a sudden rise, followed by a short period of constant resistance. This phenomenon is repeated several times. It may be assumed that each sudden rise in resistance is caused by a crack nucleus growing into a real crack of minute dimensions.

At a higher temperature (27°C) the incubation time is seen to be shorter, and to be followed by a period of continuous crack growth (fig. 13). In this case the discontinuities have apparently become too small to be measured. On the other hand at a lower

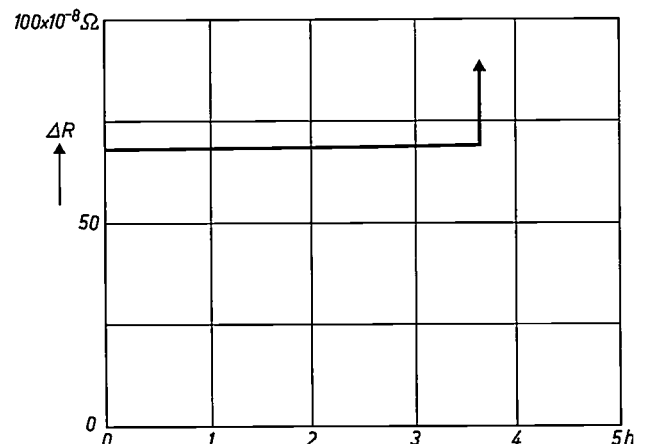


Fig. 14. The same experiment as in fig. 12, at a lower temperature of -46°C . The incubation period is now much longer, about 3.5 hours, followed by a sudden total fracture of the bar.

temperature (-46°C), the incubation period is found to be longer, and to be followed by a sudden complete fracture of the bar (*fig. 14*) — the first crack that forms propagates throughout the whole bar with great velocity.

In conclusion it should be noted that there are many other cases of fracturing which are due to the presence of hydrogen in steel, typical examples being "flakes" or "shatter cracks" and "fish eyes"¹⁷⁾. One of the reasons why these phenomena are not discussed in this article is that our understanding of them is as yet too limited.

¹⁷⁾ An extensive review will be found in: E. Houdremont, *Handbuch der Sonderstahlkunde*, Springer, Berlin 1956, p. 1375-1390. For an investigation in this field, carried out in this laboratory, see P. C. van der Willigen, *Schweissen und Schneiden* 9, 517, 1957.

Summary. In order to get some insight into the causes of the embrittling effect of hydrogen in steel, the authors study the change in electrical resistance of pure iron wires as a result of electrolytic charging with hydrogen. The resistance of soft-annealed wire is shown to increase much more than that of cold-worked wire. This difference is not, or only to a small extent, caused by the interaction of dissolved hydrogen and dislocations but must primarily be ascribed to the different orientation of cracks which form along the grain boundaries during charging. The formation of cracks is shown to be accompanied by plastic deformation of the metal and by changes in dimensions. Based on this investigation and on the existing literature the authors conclude that the deleterious influence of hydrogen on the ductility of iron and steel is mainly due to the formation of molecular hydrogen of high pressure in microvoids, especially in crack nuclei formed by coalescence of dislocations. This gives an explanation for the well-known facts that the chance of brittle fracturing of steel containing hydrogen is greater (within certain limits) the lower the rate of deformation and the higher the temperature, the hydrogen atoms needing sufficient time to diffuse to the crack nuclei formed during plastic deformation. Apart from brittle fractures occurring during testing, attention is also paid to the delayed brittle fractures that can have such serious practical consequences.

ABSTRACTS OF RECENT SCIENTIFIC PUBLICATIONS BY THE STAFF OF N.V. PHILIPS' GLOEILAMPENFABRIEKEN

Reprints of these papers not marked with an asterisk * can be obtained free of charge upon application to the Philips Research Laboratories, Eindhoven, Netherlands, where a limited number of reprints are available for distribution.

3036: W. van Gool and A. P. Cleiren: Bemerkungen zur Natur der Leuchtzentren in aktivatorfreien ZnS-Luminophoren (*Z. Naturf.* 16a, 948-950, 1961, No. 9). (Remarks on the nature of the fluorescence centres in activator-free ZnS; in German.)

ZnS can easily be made fluorescent by calcining it together with a halide. The nature of the lattice imperfection which is responsible for the blue fluorescence which is produced on irradiation with ultraviolet light is not yet known with any certainty. The authors mention five possible imperfections, and discuss why the experimental data are insufficient to allow an unambiguous choice between them. See also 2988.

3037: C. Wansdronk: Miniature condenser microphones (*Proc. 3rd int. congress on acoustics, Stuttgart 1959, Vol. II, pp. 638-640, Elsevier, Amsterdam 1961*).

The author discusses the properties of two small condenser microphones (8 mm in diameter and 55 mm long) which he designed; such microphones are needed for acoustic work and for use in film and television studios. One of the microphones has a preferred direction for the sensitivity, while the other is equally sensitive in all directions.

3038: D. L. A. Tjaden: Some considerations on the recording process on magnetic tape with application of HF bias (as 3037, pp. 758-760).

Measurements of the sensitivity of magnetic recording tape as a function of the amplitude of the signal, and of the "thickness losses" as a function of the recording depth. The results are in reasonable agreement with the theory.

3039: J. Hornstra: Rectificatie van krommen (*Chem. Weekblad* 57, 541-544, 1961, No. 42). (Plotting results as straight-line graphs; in Dutch.)

When plotting experimentally determined relationships between physical quantities, it is often preferable to plot not the experimentally determined quantities but suitable functions of these, so chosen as to make the graph a straight line. Parameters occurring in the relationship can then be determined with greater ease and accuracy. This is illustrated by reference to a number of examples; a discussion of the calculation of the error and of the limits of applicability of the method follows. The author hopes that this publication will fill a gap in the literature which became apparent to him when examining chemical analysts in advanced

mathematics: most of the candidates had no knowledge of this very useful method.

- 3040:** W. J. Oosterkamp and Th. G. Schut: *Magnetische Speicherung von Röntgenbildern* (Elektromedizin 6, 147-152, 1961, No. 3). (Magnetic recording of X-ray images; in German.)

A brief discussion of the use of television techniques in medical X-ray diagnostics. Particular mention is made of the magnetic image memorizer, on which a number of X-ray images can be stored simultaneously. See also **3016** and Philips tech. Rev. **22**, 1-10, 1960/61.

- 3041:** P. A. H. Hart and C. Weber: A transmission-line coupler for a fast wave transverse velocity electron beam amplifier (Nachr.-techn. Fachber. **22**, 358-361, 1961).

Theoretical treatment of a low-noise parametric amplifying tube, in which the signal transmission is effected with the aid of an electron beam. The Cuccia transmission-line coupler and the travelling-wave coupler are compared.

- 3042:** M. T. Vlaardingerbroek: Comparison of noise in microwave triodes and in electron beams (Nachr.-techn. Fachber. **22**, 399-402, 1961).

A brief survey of the subject matter which the author dealt with in detail in his thesis (**R 393**).

- 3043:** G. G. J. Bos: Enkele aspecten van voorraadketens (Statistica neerl. **15**, 489-500, 1961, No. 4). (Some aspects of warehouse chains; in Dutch.)

When goods are not supplied directly from the factory to the customer, but via a number of dispersed warehouses, it is necessary to ensure that each warehouse has sufficient stocks to meet customers' demands. In this paper the author considers the question of when stocks should be ordered for the dispersed warehouses, and in what quantities. A

distinction is made between two systems of organization. In the first there is a chain of warehouses between the factory and the customer; in the second the products are directed to the dispersed warehouses by a central stores management.

- 3044:** J. F. Schouten: *Der Reaktionsablauf beim Menschen* (published in Aufnahme und Verarbeitung von Nachrichten durch Organismen, proceedings of NTG-Fachtagung, Karlsruhe 1961, pp. 49-55, Hirzel, Stuttgart 1961). (Measurement of human reaction times; in German.)

In the Institute for Perception Research at Eindhoven an installation has been developed (called "DONDERS") which makes it possible to carry out in a simple way the numerous observations required for measuring human reaction times. The "DONDERS" is used in combination with a "histometer", an instrument developed at the same time which sorts the observations and presents them in an ordered form. The installation is an illustrative example of the potentialities of electronic engineering in the field of psychophysics.

- 3045:** E. Roeder and G. D. Ricck: *Walz- und Rekristallisationstexturen von dünnem Wolframblech mit und ohne Zusatz* (Z. Metallkunde **52**, 572-576, 1961, No. 9). (Rolling and recrystallization textures of thin doped and undoped tungsten foil; in German.)

Investigation of the structure of rolled tungsten foil, in which the properties of doped and undoped material are compared. The observations were made with an X-ray diffractometer and a metallurgical microscope. Both the same microstructure and the same texture (a rolling texture centred around the [110] direction) are found in the doped and undoped material. The two materials show differences in regard to recrystallization temperature and the shape and dimensions of the grains after secondary recrystallization.

Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES



COLOUR SEPARATION IN COLOUR-TELEVISION CAMERAS

by H. de LANG *) and G. BOUWHUIS *).

535.415:621.397.132

In colour television, light has to be split into a blue, a green and a red component, and this can be done quite satisfactorily by means of interference layers with a colour-selective reflecting action. There are several ways in which the layers can be incorporated in an optical system. In this article the considerations that enter into colour separation design for colour-television cameras are explained and a useful new type of colour separation system is described.

A colour-television camera for normal broadcasting purposes must possess powers of colour discrimination corresponding to the standards of the human eye. This means that essentially the apparatus must consist of three cameras, each with its own spectral sensitivity distribution, the peaks of the three curves occurring at blue, green and red.

In what follows we shall only be concerned with the most widely used form of colour camera, in which the scene to be transmitted is conveyed via three optical "channels" to the photosensitive layers of three exactly similar pick-up tubes, three vidicons

for example. The three channels must have the specified spectral transmission characteristics. Furthermore, the blue, green and red target images of the scene must all be formed as seen from the same point; in optical terms, the three images must have a common entrance pupil. If this requirement were not fulfilled, parallax phenomena would upset the congruence of the three images.

For the actual focusing of the target images, objective lenses are used whose quality in this respect must meet the same requirements as objectives for monochrome television or photography. The separation process, i.e. the splitting into the three optical channels, is effected by colour-selective in-

*) Philips Research Laboratories, Eindhoven.

interference (dichroic) mirrors. The principle of these mirrors, and the way in which they are made, were discussed in an earlier issue of this review ¹⁾.

Possible systems for image formation and colour separation

There is no *a priori* reason for forming the image first and separating the colours afterwards, or the other way about. Some types of optical system that are possible in principle are illustrated in figs. 1, 2 and 3.

In the set-up shown in fig. 1, separation takes place first. Mirror S_R splits off the red component of the light coming from the scene to be transmitted, and mirror S_B splits off the blue from the remaining mixture. Each of the three pick-up tubes has its own objective lens. One drawback of this system (with which experiments were carried out at one time in this laboratory) is that all three objectives have to be changed when it is necessary to switch to a long-distance or wide-angle shot. This makes it very difficult to maintain the optical and mechanical uniformity of the three channels.

Fig. 2 shows a set-up in which the three optical channels share the same objective. Colour separation is effected behind the objective, by dichroic mirrors S_R and S_B . Here no difficulty is involved by objective-switching; the optical uniformity of the three channels is dependent only upon the correct positioning and alignment of the mirrors and pick-up tubes. A disadvantage of this system is that the mirrors necessitate a long working distance between the objective and the tubes. Consequently a reasonable wide-angle effect can only be obtained with the aid of special objectives having a free working distance (back focus) considerably greater than their focal length. The system sketched in fig. 2 is employed in the colour-television camera for medical use that Philips put on the market some years ago. In this camera, the focal length is varied not by changing objectives but by adjusting a Zoom lens (i.e. a multiple objective whose elements can be displaced relative to each other, a continuous range of focal lengths being obtainable in this way). This Zoom lens has a back focus of 140 mm combined with a minimum focal length of only 45 mm.

A final variant is shown in fig. 3. Objective O_1 projects an image of the scene onto the plane occupied by field lens O_2 , which transmits it through a "relay" system consisting of the two objectives O_3 and O_4 onto the green pick-up tube and, via

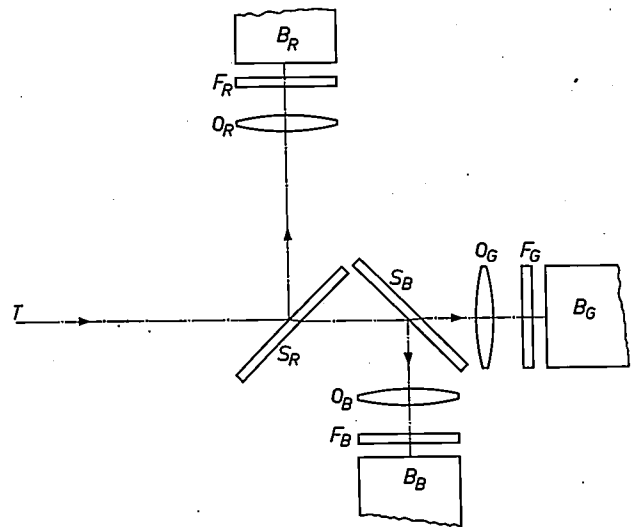


Fig. 1. Colour-television camera set-up in which light from the scene T to be transmitted is first split, by dichroic (or colour separating) mirrors S_R and S_B , into a red, a blue and a green component. Behind the mirrors, objectives O_R , O_B and O_G project images on to three pick-up tubes B_R , B_B , and B_G having exactly similar characteristics. If desired, any unwanted wavelengths can be removed (with some loss of intensity) from the component images by inserting absorption filters F_R , F_B and F_G .

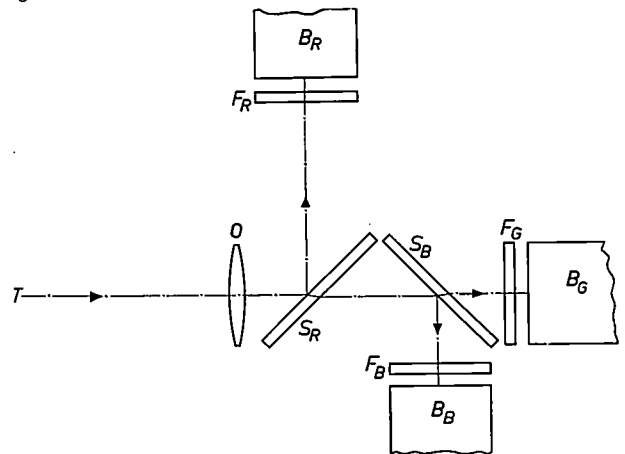


Fig. 2. Colour-television camera set-up in which the dichroic mirrors are placed behind the image-forming element, a single objective O . The other letters have the same meanings as in fig. 1.

mirror S_B , onto the blue pick-up tube. Similarly, the image reaches the red tube by way of objectives O_3 and O_4 and mirror S_R . It will be observed that here the focusing and separating processes have to a certain extent been interleaved. A system of this kind has been built into an experimental camera that is being used in our laboratory. Correction lenses O_{5G} , O_{5B} and O_{5R} have been inserted in order to compensate aberrations produced by the field lens. Considering the somewhat improvised character of its optical system, resolution of this camera is quite satisfactory. In the arrangement shown in fig. 3 no limitations are placed on the back focus of the objectives, and it is possible in consequence to get very short effective focal distances — down

¹⁾ P. M. van Alphen, Applications of the interference of light in thin films, Philips tech. Rev. 19, 59-67, 1957/58.

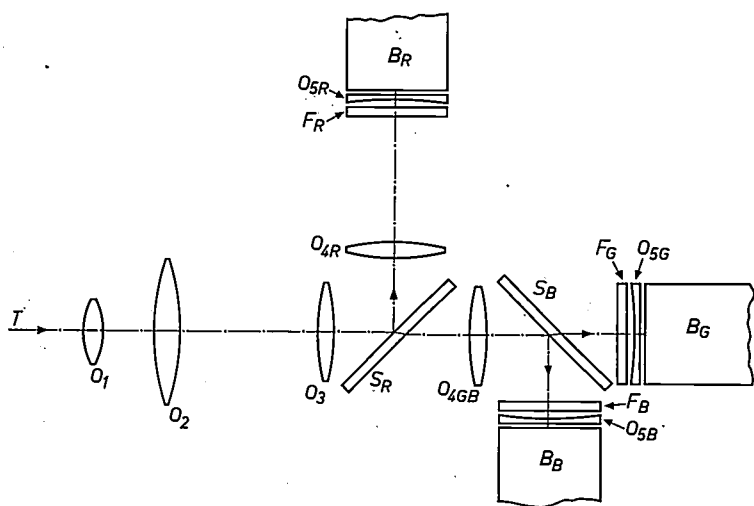


Fig. 3. Colour-television camera set-up in which the image-forming and light-filtering elements are "interleaved". Objective O_1 projects an image of the scene onto field lens O_2 , and this image is relayed by the combinations O_3 - O_{4R} and O_3 - O_{4GB} to the three pick-up tubes. O_{5R} , O_{5B} and O_{5G} are correction lenses. The other letters have the same meanings as in figs. 1 and 2.

to 25 mm at a relative aperture of 1:2.0 in the experimental camera just mentioned, the target image having the dimensions 12×16 mm. A drawback of the system is the large number of glass-air boundaries (about 30), which cause considerable loss of light and loss of contrast, only about 35% of the incident light being transmitted to the tubes.

As already stated, the system shown in fig. 2 is employed in the colour-television camera manufactured by Philips for medical use. The actual layout of the optical elements in this camera may be seen in

fig. 4. In order to limit the dimensions of the camera it was desirable that the three pick-up tubes, with their scanning coils, should lie roughly parallel. Two ordinary plane mirrors, s_1 and s_2 , have been inserted to make this possible. The elements cannot be packed any more compactly in view of the requirement that no part of any beam may be intercepted.

All the methods of colour separation that have been briefly described above involve difficulties inherent in the use of interference layers on flat-plate substrates. The difficulties in question will be discussed in the present article. To a very great extent we have been able to overcome them by developing a new type of colour-separation system, viz a cemented assembly of prisms incorporating colour-selective interference layers.

The article will end with an account of design details of this system.

Drawbacks of conventional dichroic mirrors

The difficulties just referred to are listed below, the remedies also being indicated where these have been found.

- a) As is clear from fig. 4, the dichroic mirrors take up a great deal of space, prejudicing the optical efficiency of the system.

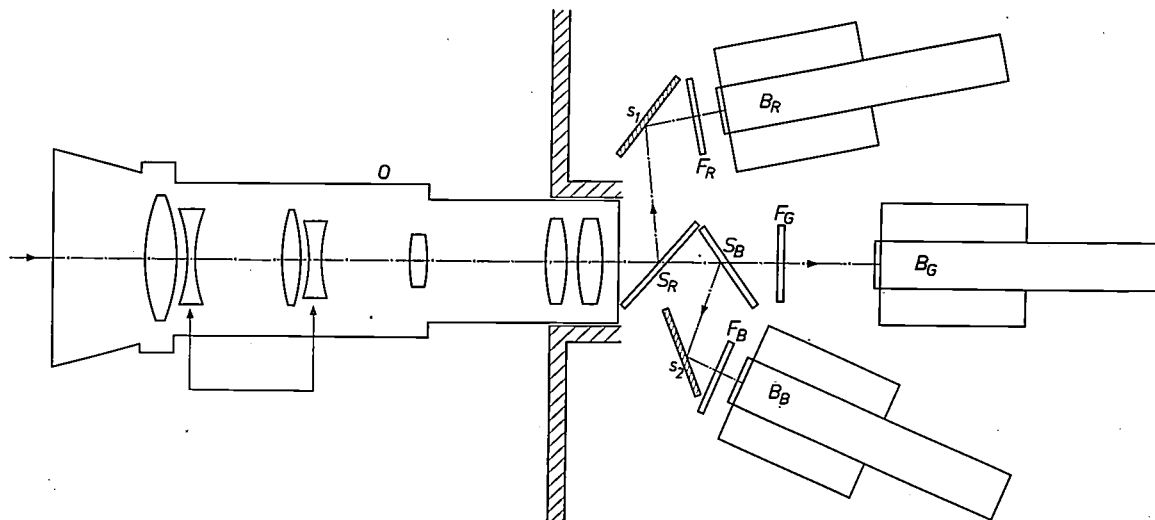


Fig. 4. Layout of components in a colour-television camera which for some years Philips has been manufacturing for medical applications, here shown schematically but approximately to scale. The optical system is based on the set-up in fig. 2, and the letters have the same meaning as in that diagram. The objective O is a Zoom lens. Since the dichroic mirrors take up a good deal of space the lens had to have a large back focus; the lens actually used in the camera has a back focus of 140 mm, a minimum focal length of 45 mm and a relative aperture of 1:4.2. Insertion of ordinary plane mirrors s_1 and s_2 has made it possible to give the three pick-up tubes (whose scanning coils are also outlined in the drawing) a roughly parallel configuration.

- b) The plane-parallel glass plate onto which the interference layers of a mirror are to be evaporated generally introduces aberrations into the transmitted beam. These aberrations — astigmatism and coma — cannot be allowed for when designing the objectives, because the rotational symmetry of the system about the optical axis is upset by the oblique positioning of the glass plate. Aberrations can be diminished by using compensating optical elements²⁾, and an even more effective and convenient remedy is to give the plate a slightly wedged shape; in this way astigmatism and coma can be corrected at the same time³⁾. But at large relative apertures these expedients are inadequate.
- c) Reflection from the rear face of the plate may give rise to a "ghost" image which, though very faint, is liable to be troublesome under certain conditions.
- d) In practice the interference layers are highly susceptible to atmospheric attack and other damage. A big improvement in this respect has been obtained from the use of harder and more resistant layers consisting of materials like TiO_2 and SiO in place of ZnS and Na_3AlF_6 or MgF_2 .
- e) A final, fundamental difficulty arises out of the fact that the properties of an interference layer depend on angle of incidence. This dependence results in all kinds of troublesome effects, particularly at large angles of incidence. In an arrangement like that shown in fig. 4, where the need to avoid partial interception of the beams compels the utilization of fairly large angles of incidence (42° and 35° for the central ray), these effects are something of a nuisance. The whole question is a somewhat involved, but physically interesting one, and a separate section will therefore be devoted to it.

Dependence of reflection properties on angle of incidence

As is explained in the article quoted above¹⁾, a layer of physical thickness d has an effective optical thickness of $nd \cos \varphi$ for a ray striking it at an angle φ to the normal, n being the refractive index of the layer. Because of this, the reflection and transmission curves of the layer will shift towards smaller wavelengths as the incidence angle φ increases. The position of these curves is generally indicated by speci-

fying the wavelength $\lambda_{50/50}$ at which half the incident energy is reflected and half transmitted. The above statement means, then, that $\lambda_{50/50}$ decreases with increasing angle of incidence.

The implications will be clear from a study of fig. 5. The two rays p and q , which originate at different points in the scene to be transmitted, pass through the centre of objective O , and end up at different points P and Q on the photosensitive layer

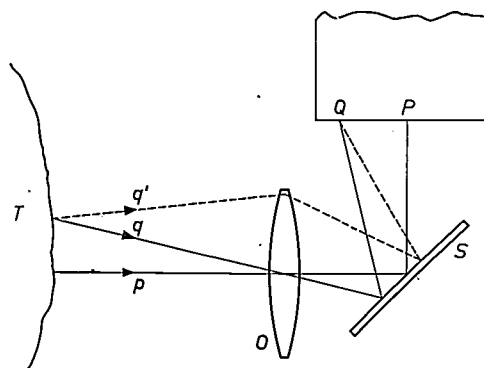


Fig. 5. Rays p and q from two points of scene T impinge on dichroic mirror S at different angles. Since the spectral reflection characteristic depends on angle of incidence, colour rendition is different at points P and Q of the screen image; in other words, there is a spurious gradation of colour across the image.

Rays q and q' , originating at the same point of the scene, likewise strike the mirror at different angles. Consequently their contributions to Q , a single element of the target image, are governed by slightly differing spectral reflection characteristics. The result is a loss of colour discrimination.

of the pick-up tube. The two rays strike the dichroic mirror S at different angles, in consequence of which there is a difference in colour rendition at points P and Q ; in other words, the variation in angle of incidence gives rise to a spurious colour gradation across the target image. Let us now consider rays q and q' , which originate at the same point in the scene but pass through different points in the entrance pupil. These rays likewise strike S at different angles, but they converge to the same point Q on the photosensitive layer. Colour rendition at that point must therefore be in accordance with a reflection characteristic that is the mean of two spectral reflection curves, or, in general, of a whole series, each having a slightly different position in the spectrum; which makes it difficult for the system to discriminate sharply between colours. The effect is particularly troublesome when a high aperture objective is used, but even under ordinary aperture conditions it is quite noticeable. In the camera shown schematically in fig. 4, the relative aperture is 1 : 4.2, and the angles of incidence of rays from the centre of the transmitted scene vary by about 7° on either side of the central ray.

²⁾ L. T. Sachtleben, D. J. Parker, G. L. Allee and E. Kornstein, RCA Rev. 13, 27, 1952.

³⁾ B. Cuny, Revue d'Optique 34, 460, 1955. H. de Lang, Philips Res. Repts. 12, 181, 1957.

Both effects — the variation in colour rendition across the image, and the averaging of a series of reflection characteristics for one image point — can be counteracted by introducing a gradient in the thickness of the layer across the mirror. This is not a very attractive method and, what is more, the required variation in layer thickness in general is not usually the same for the two effects. Indeed, in straightforward cases like the one illustrated in fig. 5, the variation to allow for the one effect would have to be opposite in sign to the variation required by the other effect. If, therefore, the spurious colour gradation is to be compensated, and that is generally the more damaging of the two effects, then an even poorer colour separation will have to be accepted.

In addition to the dependence of effective optical layer thickness on angle of incidence, there is an undesirable effect connected with polarization phenomena to which obliquely incident rays are subject. With increasing angle of incidence the coefficient of reflection becomes greater for light whose electric vector vibrates in a direction perpendicular to the plane of incidence, and smaller for light whose electric vector vibrates parallel with that plane. Accordingly, the reflection coefficient for the "perpendicular component" of natural light arriving at an oblique angle of incidence is greater than that for its "parallel component" ⁴⁾. Fig. 6 is a plot, appropriate to a given case, of the two spectral

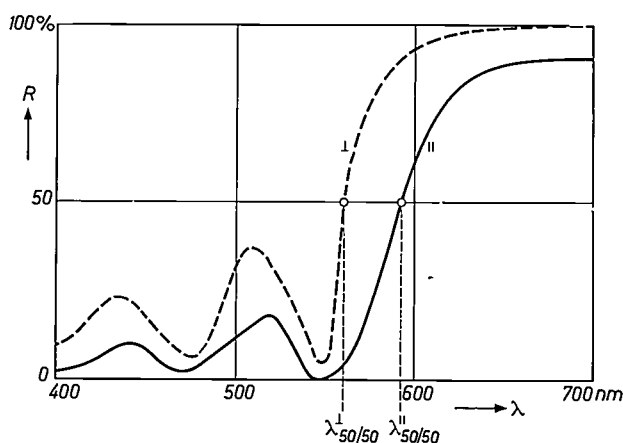


Fig. 6. Reflection coefficients R of a red-reflecting mirror (7 layers of ZnS and Na_3AlF_6 deposited alternately on a glass plate) as a function of wavelength λ for the "perpendicular" and "parallel" polarized components of light arriving at an angle of incidence of 42° . Wavelength $\lambda_{50/50}$ is a good deal shorter for the "perpendicular" component.

⁴⁾ The reflection coefficient is practically zero for the parallel component of light at the Brewster angle; under these conditions, then, transmission for this component is almost 100%, and the interference mirror thus represents a highly effective means of resolving a beam of light into two beams whose vibration directions are at right angles.

reflection characteristics in question (they are known as "mirror characteristics"). In practice, the reflection of unpolarized light entering the camera from the scene to be transmitted will be governed by the mean of the two mirror characteristics.

However, if the light from the transmitted scene is already polarized (as is the case for light reflected from specular surfaces, water surfaces, for example, or perspiring foreheads of actors!), then colour shifts may occur, since the reflection characteristic of the interference layer is dependent on the orientation of the plane of polarization.

As was done in the preceding section, mention will be made of one or two ways of overcoming these undesired effects.

Colour shifts due to polarization effects can be eliminated by setting a $\frac{1}{4}\lambda$ plate, in diagonal position, in front of the colour separation system. It can be demonstrated that when this is done, the perpendicular and parallel components into which the light behind the plate can be resolved are always of equal strength. Hence, whatever the original direction of polarization may be, an effective spectral reflection characteristic is obtained that is the mean of the reflection characteristics for perpendicularly polarized and parallel-polarized light. It is true that this averaging process allows full weight to the poor mirror characteristic for the parallel component at large angles of incidence — but this of course would have happened anyway if the light entering the system was unpolarized. This adverse effect of the parallel-component characteristic can only be countered by keeping the actual angles of incidence as small as possible; this, indeed, is the only course that is correct in principle, since it also gets rid of the difficulties referred to above, arising out of the dependence of the effective optical layer thickness on the angle of incidence. The trouble is, however, as we saw when discussing fig. 4, that the layout of the various camera elements does not allow the angles of incidence to be further reduced to any appreciable extent (i.e. does not allow less oblique positioning of the dichroic mirrors).

Something may be gained in this unfavourable situation by the following elegant method: the high-refractive-index layers of the dichroic mirror can be given a thickness of $\frac{3}{4}\lambda$ (or even $\frac{5}{4}\lambda$, or $\frac{7}{4}\lambda$, etc.), that of the low-refractive-index layers being left at $\frac{1}{4}\lambda$. The result of this is to weaken the dependence of effective layer thickness on angle of incidence. The explanation, in simple terms, is that when light enters a $(\frac{3}{4}\lambda, \frac{1}{4}\lambda)$ mirror of this kind, the greater part of its path lies through a strongly refractive medium; here the rays are closer to the normal than they are

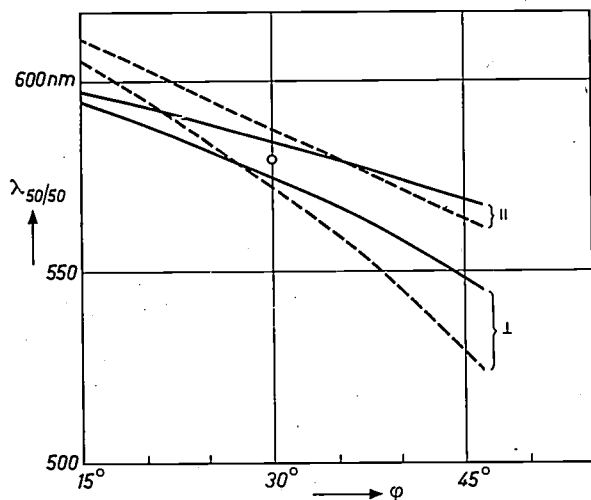


Fig. 7. Colour separation as a function of angle of incidence φ for two mirrors of similar composition to that described under fig. 6, both designed for a $\lambda_{50/50}$ value of 580 nm at an angle of incidence of 30° . In one of the mirrors (dashed-line curves) the strongly refractive and the weakly refractive layers are both $\frac{1}{2}\lambda$ thick; in the other (fully-drawn curves) the strongly refractive layers are $\frac{3}{4}\lambda$ thick. Wavelength $\lambda_{50/50}$ as a function of angle of incidence has been plotted separately for the parallel-polarized and the perpendicularly polarized component. The $\lambda_{50/50}$ characteristics for unpolarized light are given by the mean of this pair of curves. Clearly, the variation in $\lambda_{50/50}$ as a function of wavelength is appreciably smaller for the $(\frac{3}{4}\lambda, \frac{1}{2}\lambda)$ mirror, and there is less difference between this mirror's treatment of the perpendicular and parallel components.

in a weakly refractive medium. The better colour-separation-action properties of $(\frac{3}{4}\lambda, \frac{1}{2}\lambda)$ mirrors is demonstrated in fig. 7.

In employing these $(\frac{3}{4}\lambda, \frac{1}{2}\lambda)$ mirrors it must be taken into account that the red-filtering mirror also has a blue reflection. The layout must therefore be chosen in such a way that the blue component of the light entering the system is split off first by the blue-reflecting mirror; parasitic blue reflection from the following red-reflecting mirror is then of no importance.

Colour separation by means of a prism system

The difficulties referred to above are either completely eliminated or greatly reduced by a new colour-separating system developed in this laboratory, whose principle is illustrated in fig. 8. The colour-selective reflecting layers S_B and S_R have been applied to the faces of prisms which are cemented together. Light entering the system, after crossing a thin air layer, passes through face 1 and falls on to the blue-reflecting mirror S_B . Having been reflected by S_B , the blue component undergoes total reflection from face 1 and leaves the system. The light that has been transmitted by S_B , after crossing a further thin air layer, passes through face 2 and falls on to the red-reflecting mirror S_R . The red light reflected by this mirror leaves the system after under-

going total reflection from face 2. The green component of the light entering the system passes straight through all these prism faces and mirrors.

For all three colours the arrangement of prisms is optically equivalent to a thick plane-parallel plate standing at right angles to the optical axis. Thus the system is free from the aberrations introduced by oblique-lying glass supports, aberrations for which no completely effective remedy is available when the relative aperture is large. Though having an oblique orientation, the plane-parallel air layers within the system are too thin (of the order of 0.1 mm) to cause errors of any consequence. It will further be evident that since the reflecting layers are inside the cemented assembly, they are fully protected from atmospheric attack and other damage and also kept clear of dust. Moreover, if suitable precautions are taken, there is no trouble from "ghost" images arising out of unwanted reflections. A further important feature is the compactness of the camera set-up which is now possible and which is shown in fig. 9. As will be clear from comparison with fig. 4, the mirrors s_1 and s_2 inserted in that system have been superseded here by the totally reflecting surfaces 1 and 2; thus these mirrors are incorporated in the prism without intercepting the incoming light. The compactness of the arrangement thus ob-

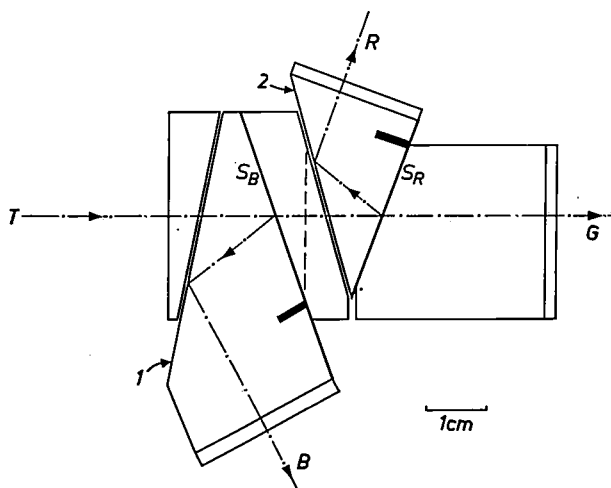
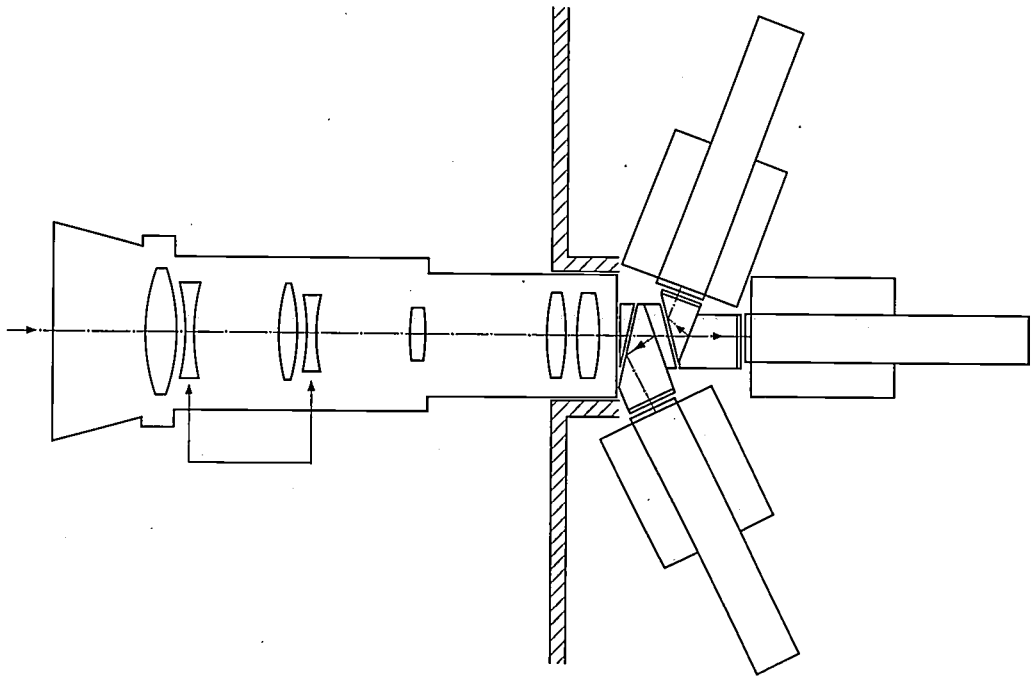


Fig. 8. Schematic representation of a colour-separating system consisting of glass prisms. The colour-selective interference layers S_B and S_R have been applied to the faces of the prisms which are cemented together. The continuity of the glass is broken by the thin air layers interposed in front of faces 1 and 2. Having been reflected by S_B , the blue component B of the incoming light undergoes total reflection from face 1 and leaves the system. The red component R is reflected by S_R and leaves the system on undergoing total reflection from face 2. The green component G passes right through the system. The central ray meets S_B and S_R at the same angle of incidence: 20° . For each of the three components of the incoming light, the whole system behaves like a plane-parallel glass plate lying at right angles to the optical axis.

Before leaving the system, each component passes through an absorption filter which is cemented on to the multiple prism and which performs a similar function to the final filters in figs. 1 etc.

Fig. 9. Layout of components in a colour-television camera similar to that shown in fig. 4 but incorporating a prismatic colour-separating system like that in fig. 8. Use of the multiple prism has resulted in a much more compact layout, an objective having a smaller back focus can be employed, and colour separation is better because the light undergoing reflection arrives at smaller angles of incidence.



tained means that requirements as to the working distance of the objective are much less stringent (see below). Finally, a great virtue of this as compared with earlier designs is that it allows the angles of incidence to be reduced. In fig. 8 the central ray meets S_B and S_R at angles of incidence not greater

than 20° . Admittedly, as it stands this value is not comparable with those of 42° and 35° cited above, where it was a matter of light rays in air incident on interference layers applied to plates: in the present case we are concerned with light rays in glass, and in accordance with Snell's Law, a given angle of incidence in glass is equivalent to that angle of incidence in air whose sine is n times greater, n being the refractive index of the glass. Even so, the gain is appreciable; for $n = 1.52$ we find that a 20° angle of incidence in glass corresponds to a 31° angle of incidence in air.

Fig. 10 demonstrates the improvement in colour separation achieved by reducing the angles of in-

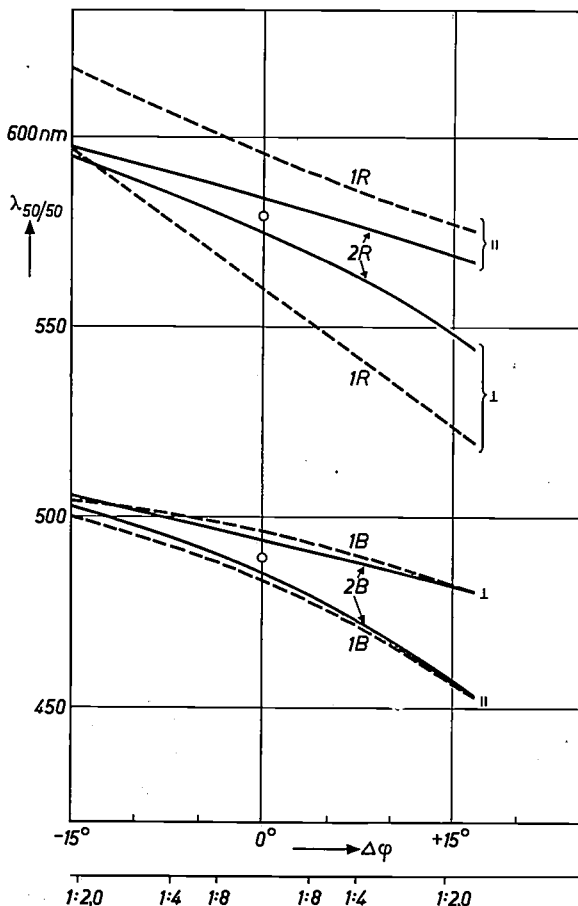


Fig. 10. Colour separation as a function of angle of incidence in two colour-separating systems, plotted in the same way as in fig. 7. Curves $1R$ and $1B$ relate to a system like that in fig. 4, having $(\frac{1}{2}\lambda, \frac{1}{2}\lambda)$ red-reflecting and $(\frac{3}{4}\lambda, \frac{1}{4}\lambda)$ blue-reflecting layers. Curves $2R$ and $2B$ relate to the multiple prism in fig. 8, in which $(\frac{3}{4}\lambda, \frac{1}{4}\lambda)$ interference layers are used for both red reflection and blue reflection.

The quantity along the abscissa, $\Delta\phi$, is the angular deviation from the direction of the central ray (equivalent values in air). This representation is chosen since the central ray itself strikes each of the mirrors under consideration at a different angle of incidence — indeed, the better performance of the mirrors in the multiple prism is attributable to this fact. The values of these angles of incidence are displayed in the table below.

For any angular deviation $\Delta\phi$ can be specified a relative aperture which involves deviations from the central ray up to that value of $\Delta\phi$. A scale of the pertaining relative apertures of the objective has also been drawn along the abscissa.

	Layers on plate		Layers in prism	
	$1R$	$1B$	$2R$	$2B$
Angle of incidence of central ray	42°	35°	20°	20°
	in air		in glass	
Interference layers used	$\frac{1}{2}\lambda, \frac{1}{2}\lambda$	$\frac{3}{4}\lambda, \frac{1}{4}\lambda$	$\frac{3}{4}\lambda, \frac{1}{4}\lambda$	$\frac{3}{4}\lambda, \frac{1}{4}\lambda$

cidence. The graph shows variations in $\lambda_{50/50}$ (explained above) as a function of the angle of incidence; curves 1R and 1B relate to the two mirrors in the system of fig. 4, curves 2R and 2B to those in the multiple prism.

Various practical designs are possible of the principle illustrated in fig. 8. The version drawn in perspective in fig. 11 corresponds quite closely to fig. 8: basically, all that has happened is that the rear portion of fig. 8 (to the right of the dashed line) has been

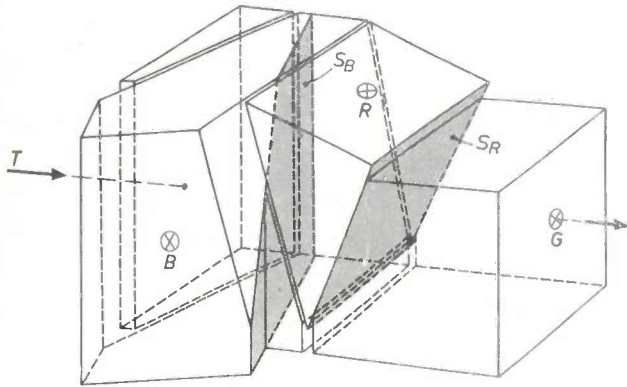


Fig. 11. Practical version of multiple prism based on the principle shown in fig. 8. The shaded faces carry the interference layers. The rear portion of the system in fig. 8, that lying to the right of the dashed line, has here been turned 90° about the axis (central ray G). The points at which the central rays B, R and G of the components leave the system have been indicated for the sake of clarity.

turned through 90° about the optical axis. The resulting arrangement is a convenient one from the viewpoint of camera-tube positioning; the three tube axes no longer have to lie in one plane, as in fig. 9. This prism system is suitable for use with an objective having a relative aperture of 1:2.0 and producing a target image measuring 9×12 mm. The path length of the light through glass is 63 mm, a figure that must be allowed for in the correction of the objective. For the working distance determined by this path length, a Zoom lens is available whose shortest focal length is only 25 mm. The chosen lens type has proved to have an additional advantage in that the exit pupil is practically situated at infinity. Consequently there is very little spurious colour gradation across the plane of the image, of the kind that was explained with the aid of fig. 5.

A prototype of the prism assembly described above has been made in this laboratory, with the collaboration of the Glass Division, for the Electroacoustics Division (which manufactures the cameras). The prototype was used to take the photograph that heads this article.

In camera production on an industrial scale, not much play is available for aligning the pick-up tubes

and it is therefore necessary that the three images projected via the multiple prism should occupy rather closely defined positions. This necessitates a fairly high degree of accuracy in the making of each system component, and a similar degree of accuracy in the cementing operation. The narrow gaps must not of course be filled up with cementing agent, the air layer being necessary for total reflection. To prevent the entry of dust, an airtight strip of the cementing agent is applied around the edges of the prism faces in question.

Another practical version of the multiple prism appears in fig. 12. Underlying the systems in figs. 8 and 11 is the requirement that the central ray should strike the two dichroic mirrors at the same angle. In the system of fig. 12 the principle of equal angles of incidence has been abandoned, on the grounds that the separation of red from green is more important than that of green from blue. This is particularly true of reproducing the colour of the human skin. Accordingly, the angle of incidence at which the central ray strikes the red-reflecting mirror has here been reduced from 20° to 13° (equivalent to an angle of incidence in air φ_a of 20°), and this has meant accepting a larger angle of incidence for the blue-reflecting mirror, namely 25.5° ($\varphi_a = 41^\circ$). The different choice of angles allows the system of prisms shown in fig. 8 to be simplified to some extent: face 1 can become the front face of the multiple prism in fig. 12, and the air layer for face 2 can then be adjacent to S_B . As a result, the number of elements composing the system is reduced to three. However, it is impossible here to rotate part of the

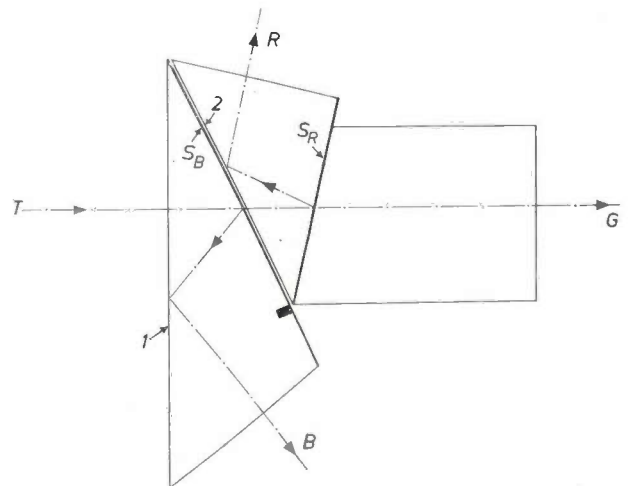


Fig. 12. Another version of the multiple prism. Here the angle at which the central ray is incident on S_R has been made as small as possible, viz 13° , in order that the separation of red from green may be as efficient as possible; this has necessitated increasing to 25.5° the angle at which the central ray is incident on S_B . Consisting of only three parts, this design represents a simplification on that in fig. 8.

system about the optical axis, as was done in fig. 11, so the resulting camera layout will be rather less compact.

Summary In colour television, the light coming from the scene to be transmitted has to be split into a blue, a green and a red component. In most colour-television cameras, including those the Philips Company has been manufacturing in recent years, colour separation is achieved with a system of dichroic mirrors consisting of interference layers deposited on flat glass plates. This form of colour-separation system has a number of inherent drawbacks: the mirrors take up a relatively large amount of space, the objective fitted to the camera has to have a fairly

large back focus, the support plates cause aberrations, the mirrors are exposed to damage, and large angles of incidence have to be reckoned with. The last-named drawback, which is associated with spurious colour gradation across the image, with unsatisfactory colour separation, and with unfaithful colour rendering when the incident light is polarized, is explained at length in this article.

All these difficulties are eliminated or greatly reduced by a newly developed system in which the colour-selective interference layers are enclosed in a cemented assembly of prisms. Several versions are described. Here the effective angles of incidence have been reduced from 42° and 35° to 31° and 31° respectively (or to 20° and 41° in another version of the new system). A graph is given illustrating the resulting improvement in colour-separating efficiency.

FINGERPRINTING VIA TOTAL INTERNAL REFLECTION

by N. J. HARRICK *).

535.394:343.977.33

"The Rays of Light in going out of Glass into a Vacuum, are bent towards the Glass; and if they fall too obliquely on the Vacuum, they are bent backwards into the Glass, and totally reflected; . . . if the farther Surface of the Glass be moisten'd with Water or clear Oil, or liquid and clear Honey, the Rays which would otherwise be reflected will go into the Water, Oil or Honey; and therefore are not reflected before they arrive at the farther Surface of the Glass, and begin to go out of it. If they go out of it into the Water, Oil, or Honey, they go on . . ."

"... more evident by laying together two Prisms of Glass, or two Object-glasses of very long Telescopes, the one plane, the other a little convex, and so compressing them that they do not fully touch, nor are too far asunder. For the Light which falls upon the farther Surface of the first Glass where the Interval between the Glasses is not above the ten hundred thousandth Part of an Inch, will go through that Surface, and through the Air or Vacuum between the Glasses, and enter into the second Glass, . . ."

Newton, Opticks, 2nd (English) ed. 1717, book III, part 1, query 29.

If one looks directly at a finger, the relief of the skin cannot be seen very clearly because there is little contrast between the hills and valleys. There is, therefore, little point in taking a direct photograph of this relief and in fact fingerprinting as applied e.g. in crime detection is still being done by the age-old ink process ¹⁾.

Recently a photographic method has been developed which is capable of producing fingerprints of superior quality and promises to do away with the untidy ink process. The apparatus used is shown in fig. 1. It was based on the phenomenon of frustrated total reflection which was already known by Newton and so aptly described in the lines we have reproduced above. In modern terms, we can describe the phenomenon as follows.

If light travelling in a glass prism strikes a glass-air boundary at an angle θ such that it is totally

reflected, no light will propagate in the adjacent air but there is nevertheless a penetration of the electromagnetic field into it ²⁾. This is indicated by fig. 2; if λ_g denotes the wavelength in the denser medium (glass), the depth of penetration is about $0.1 \lambda_g$ near grazing incidence but reaches indefinitely large values near the critical angle θ_c for total reflection. Owing to this penetration of the field, energy can be extracted from it when suitable matter is brought sufficiently near to the glass surface.

This can be demonstrated by a number of well-known experiments. If water, for example, is placed on a totally reflecting surface of a glass prism, it remains dark, but if a few drops of fluorescein are added to the water, a thin sheet of the liquid adjacent to the prism surface is seen to light up. If a second glass prism is placed at a distance d of a

²⁾ Incidentally, this entails a certain displacement of the totally reflected beam along the reflecting surface. Curiously enough, Newton's ideas also included such a displacement since he suggested that the path of the ray was a parabola with the vertex in the rarer medium (*Principia Phil. Nat.*, book I, prop. 96). Experimental evidence of the displacement was obtained only as recently as fifteen years ago: see F. Goos and H. Hänchen, *Ann. Physik* 1, 333, 1947.

*) Philips Laboratories, Irvington-on-Hudson, N.Y., U.S.A.
¹⁾ For an interesting account of the history and techniques of fingerprinting, the reader is referred to: H. Cummins and C. Midlo, *Fingerprints, palms and soles*, Dover (reprint), New York 1961.

system about the optical axis, as was done in fig. 11, so the resulting camera layout will be rather less compact.

Summary In colour television, the light coming from the scene to be transmitted has to be split into a blue, a green and a red component. In most colour-television cameras, including those the Philips Company has been manufacturing in recent years, colour separation is achieved with a system of dichroic mirrors consisting of interference layers deposited on flat glass plates. This form of colour-separation system has a number of inherent drawbacks: the mirrors take up a relatively large amount of space, the objective fitted to the camera has to have a fairly

large back focus, the support plates cause aberrations, the mirrors are exposed to damage, and large angles of incidence have to be reckoned with. The last-named drawback, which is associated with spurious colour gradation across the image, with unsatisfactory colour separation, and with unfaithful colour rendering when the incident light is polarized, is explained at length in this article.

All these difficulties are eliminated or greatly reduced by a newly developed system in which the colour-selective interference layers are enclosed in a cemented assembly of prisms. Several versions are described. Here the effective angles of incidence have been reduced from 42° and 35° to 31° and 31° respectively (or to 20° and 41° in another version of the new system). A graph is given illustrating the resulting improvement in colour-separating efficiency.

FINGERPRINTING VIA TOTAL INTERNAL REFLECTION

by N. J. HARRICK *).

535.394:343.977.33

"The Rays of Light in going out of Glass into a Vacuum, are bent towards the Glass; and if they fall too obliquely on the Vacuum, they are bent backwards into the Glass, and totally reflected; . . . if the farther Surface of the Glass be moisten'd with Water or clear Oil, or liquid and clear Honey, the Rays which would otherwise be reflected will go into the Water, Oil or Honey; and therefore are not reflected before they arrive at the farther Surface of the Glass, and begin to go out of it. If they go out of it into the Water, Oil, or Honey, they go on . . ."

"... more evident by laying together two Prisms of Glass, or two Object-glasses of very long Telescopes, the one plane, the other a little convex, and so compressing them that they do not fully touch, nor are too far asunder. For the Light which falls upon the farther Surface of the first Glass where the Interval between the Glasses is not above the ten hundred thousandth Part of an Inch, will go through that Surface, and through the Air or Vacuum between the Glasses, and enter into the second Glass, . . ."

Newton, Opticks, 2nd (English) ed. 1717, book III, part 1, query 29.

If one looks directly at a finger, the relief of the skin cannot be seen very clearly because there is little contrast between the hills and valleys. There is, therefore, little point in taking a direct photograph of this relief and in fact fingerprinting as applied e.g. in crime detection is still being done by the age-old ink process ¹⁾.

Recently a photographic method has been developed which is capable of producing fingerprints of superior quality and promises to do away with the untidy ink process. The apparatus used is shown in fig. 1. It was based on the phenomenon of frustrated total reflection which was already known by Newton and so aptly described in the lines we have reproduced above. In modern terms, we can describe the phenomenon as follows.

If light travelling in a glass prism strikes a glass-air boundary at an angle θ such that it is totally

reflected, no light will propagate in the adjacent air but there is nevertheless a penetration of the electromagnetic field into it ²⁾. This is indicated by fig. 2; if λ_g denotes the wavelength in the denser medium (glass), the depth of penetration is about $0.1 \lambda_g$ near grazing incidence but reaches indefinitely large values near the critical angle θ_c for total reflection. Owing to this penetration of the field, energy can be extracted from it when suitable matter is brought sufficiently near to the glass surface.

This can be demonstrated by a number of well-known experiments. If water, for example, is placed on a totally reflecting surface of a glass prism, it remains dark, but if a few drops of fluorescein are added to the water, a thin sheet of the liquid adjacent to the prism surface is seen to light up. If a second glass prism is placed at a distance d of a

²⁾ Incidentally, this entails a certain displacement of the totally reflected beam along the reflecting surface. Curiously enough, Newton's ideas also included such a displacement since he suggested that the path of the ray was a parabola with the vertex in the rarer medium (*Principia Phil. Nat.*, book I, prop. 96). Experimental evidence of the displacement was obtained only as recently as fifteen years ago: see F. Goos and H. Hänchen, *Ann. Physik* 1, 333, 1947.

*) Philips Laboratories, Irvington-on-Hudson, N.Y., U.S.A.
¹⁾ For an interesting account of the history and techniques of fingerprinting, the reader is referred to: H. Cummins and C. Midlo, *Fingerprints, palms and soles*, Dover (reprint), New York 1961.



Fig. 1. Fingerprint recording apparatus. The finger is pressed on a face of a glass prism. On the ground-glass viewing screen an image of the surface relief of the finger at magnification $4\times$ is produced. No ink is used. At the back of the apparatus a plate camera is mounted by which either negatives or instant prints using Polaroid film can be made. (Another apparatus has been built using a 35 mm camera.)

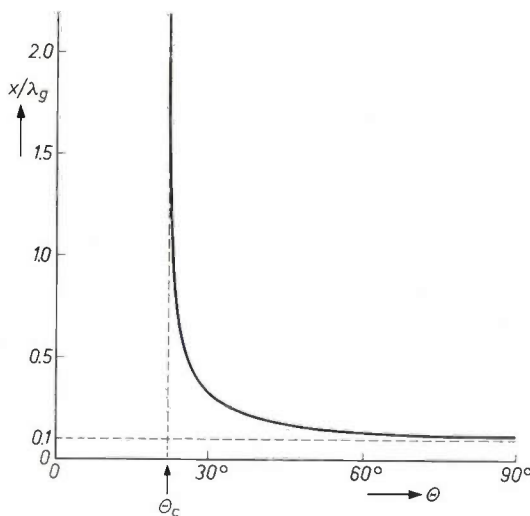


Fig. 2. At total reflection, there is a certain penetration of the electric field into the rarer medium, the field amplitude falling off with distance from the interface in an exponential manner. The penetration depth x , defined as the half-value distance, depends on the angle of incidence θ as shown in the graph. θ_c is the critical angle for total reflection, λ_g is the wavelength of the light in the denser medium.

fraction of a wavelength from the first prism, energy of the penetrating field is used for propagation in the second prism, so that a transmitted beam is obtained. The reflectivity of the prism surface can thus be continuously adjusted between 100% and zero by adjusting the distance d , as shown in *fig. 3*. The latter experiment has given rise to applications such as a *light modulator*, obtained by modulating the distance d and used for telephony on a light beam during the Second World War ³⁾, and a *cold mirror*, based on the fact that the penetration of the field increases with increasing wavelength λ_g , so that by suitable adjustment of d the heat waves are transmitted and only the light is totally reflected.

Our application for fingerprinting is of the same nature: when a finger is lightly pressed on the totally reflecting glass surface the total reflection is frustrated at the ridges of the skin which make contact with the prism (and the now transmitted light

³⁾ Electronics 17, No. 1, 156, 1944.

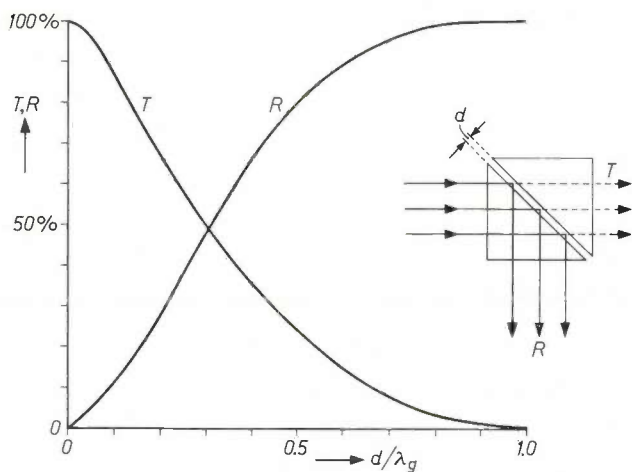


Fig. 3. By placing a second glass prism close to the totally reflecting surface of the first prism (distance d , see insert), energy is extracted from the penetrating field. The fraction of reflected (R) and transmitted light (T) will vary continuously according to this graph with variation of the distance d .

is absorbed in the skin), but not at the valleys or in the pores, where the reflection is still total.

A few points of the apparatus illustrated in fig. 1 should perhaps be explained. A mirror is provided in the optical system, as shown in fig. 4. This has been done in order to eliminate any confusion for those who are used to looking at ink prints, which are reversed images. The image screen (or the photographic plate which can be put in its place) is tilted with respect to the optical axis in order to compensate for the distortion arising from viewing the object (finger in contact with prism surface) obliquely: oblique incidence, of course, is essential for obtaining the effect. The magnification on the image screen is about $4\times$, which is adequate for a detailed examination of fingerprints, since the pores and other details can be seen quite clearly.

In fig. 5 a fingerprint obtained in this way⁴⁾ is compared to one taken by the ink technique. The latter was obtained from the same finger by a police officer, at our request (and was considered a good print). Note the high contrast, improved definition and absence of smudging in fig. 5a. The dots clearly visible in fig. 5a but not in *b* are pores in the ridges of the skin. Fig. 6 is a photo of another fingerprint in which these pores show up even more clearly. A larger portion of the finger can be shown in a print by gently pressing the sides of the finger down towards the surface. It is also possible to increase the visible portion by suitably curving the contact surface of the prism, in which case, however, the required optics will become more complicated because the surface of the prism acts as a curved mirror.

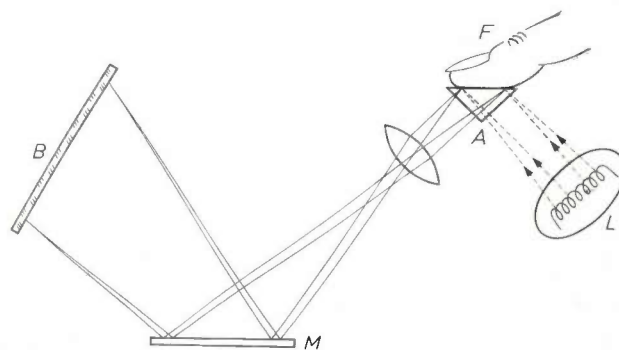


Fig. 4. Schematic diagram of fingerprint recording apparatus. F finger. A prism. L light source. B viewing screen. The mirror M serves a dual purpose: it reverses the image, so that the fingerprints on viewing will be similar to those obtained by the ink technique; and it simplifies the photographing procedure, since arrangements are made so that M can be tilted to deflect the image from B to the focal plane of the camera (not drawn).

⁴⁾ The print of fig. 5a appeared as the cover picture on an issue of the *J. appl. Phys.* where this technique was first described: N. J. Harrick, *J. appl. Phys.* **33**, 2774, 1962 (No. 9).



a



b

Fig. 5. *a*) Fingerprint recorded with the total-internal-reflection apparatus. *b*) Print of the same finger as in (*a*) obtained by the traditional ink method.

Fig. 7 is a print of a portion of an infant's foot, reproduced on the same scale as the fingerprint of fig. 6. Footprints of infants are often used in hospitals for identification. The details on an infant's foot or palm are so fine that they generally do not show up when an ink impression is made; the latter

will reveal only the outline of the foot and creases in the skin.

In addition to the greater clarity and rendering of detail that can be achieved and the cleaner procedure (since no ink is required), it will in some cases be an important advantage over the previous technique that the image can be viewed before printing and the print can be made instantly. Yet with all of these advantages, the apparatus is quite simple.

Using essentially the same arrangement, frustrated total reflection can be put to other important uses, such as microscopic examination of samples in medicine and biology (where in some cases it can replace the technique of dark-field illumination for contrast enhancement) and the study of adhesives. More sophisticated methods still based on the same principle have been used for the study of the infrared spectra of monolayer films and also of the spectra of surface states in solid-state physics (in connection with research on transistors)⁵⁾.

⁵⁾ N. J. Harrick, *Phys. Rev. Letters* 4, 224, 1960 and *Phys. Rev.* 125, 1165, 1962. These and other applications have been reviewed in: N. J. Harrick, *Annals New York Acad. Sci., Conf. on Clean surfaces (suppl. Surface phenomena in semiconductors)*, vol. 101, 928-959, 1963.

Summary. The total internal reflection of light e.g. in a glass prism can be frustrated to an adjustable degree by placing another object close to the reflecting surface. Advantage has been taken of this phenomenon (described already by Newton) for several applications, such as light modulation, cold mirrors, infrared and surface-state spectroscopy, etc. A novel application is the production of high-contrast fingerprints (images of surface reliefs) which is described in this article.



Fig. 6. Another fingerprint, in which the pores in the ridges of the skin are very distinct.

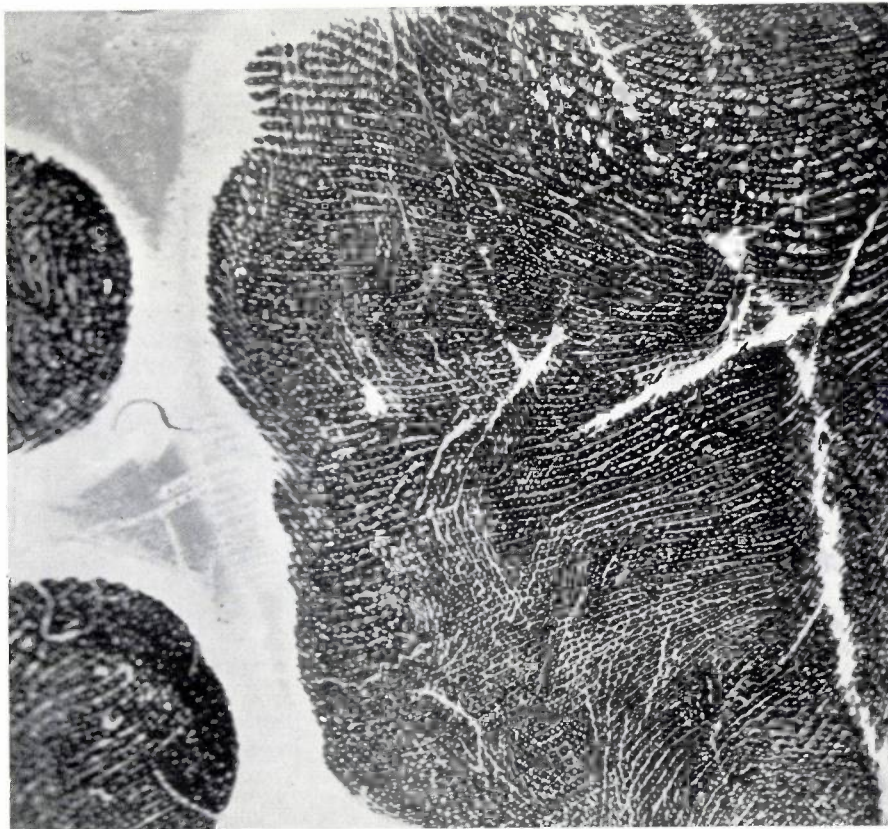


Fig. 7. Print of a portion of an infant's foot recorded with the total-reflection apparatus and reproduced on the same scale as the fingerprint in fig. 6.

DIFFERENCE AMPLIFIERS WITH A REJECTION FACTOR GREATER THAN ONE MILLION

by G. KLEIN *) and J. J. ZAALBERG van ZELST *).

621.317.725.083.6:621.375

In the various Philips laboratories there is a constant demand for special electronic measuring instruments capable of meeting widely diverse requirements. The above authors are members of a research group whose task it is to meet this demand. Among the numerous circuits which they have developed in the course of the years, there are several that give surprising results with simple means. A number of these circuits will be dealt with in a series of articles in this journal, the first of which follows below. It is a sequel to the articles on difference amplifiers that appeared in volumes 22 and 23, and shows the way in which the exceptionally high rejection factor required in a special case was obtained.

The rejection factor of a difference amplifier

In electrical measuring techniques a frequently encountered problem is the measurement of a voltage between two points both of which have potentials with respect to earth which are much larger than the voltage to be measured. This problem is not so easy to solve if the voltage under measurement (the input signal) is so small that it has to be amplified. One of the difficulties is then that the voltage between the output terminals of the amplifier is in general not only a function of the input signal but also a function of the common voltage on the input terminals with respect to earth. This difficulty is overcome by the use of special amplifiers. These *difference amplifiers*, as they are called, have been the subject of several articles in this journal ¹⁾²⁾.

The nomenclature used in those articles will also be used here; see *fig. 1*. The input signal is $2E_{it}$. The voltage $2E_{ut}$ between the output terminals is that to which the meter responds. E_{if} and E_{uf} are the average voltages with respect to earth on the input and output terminals, respectively. E_{if} and E_{uf} will be called the in-phase signals, E_{it} and E_{ut} the anti-phase signals. The letter E may denote either a DC or an AC voltage. The frequencies of E_{it} and E_{if} need not be identical.

In general we can write:

$$E_{ut} = AE_{it} + BE_{if} \quad \dots \quad (1)$$

In this expression A is the amplification which E_{it} undergoes, and the term BE_{if} represents the interfering effect due to E_{if} . In a good difference ampli-

fier B is much smaller than A . The ratio $H = A/B$ is called the *rejection factor* and is one of the most important characteristics of a difference amplifier.

Article I, mentioned under reference ²⁾, gives examples of difference amplifiers where the *guaranteed* ³⁾ value H_{min} of the rejection factor is approxi-

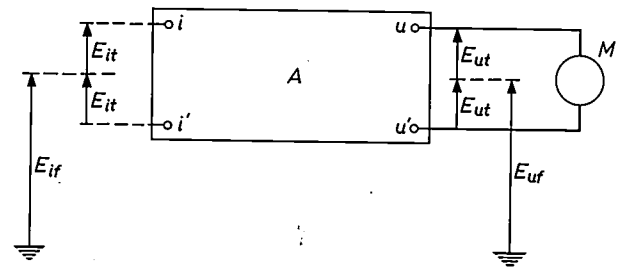


Fig. 1. Difference amplifier A . The signal to be measured, $2E_{it}$, is applied between the input terminals $i-i'$; the amplified voltage $2E_{ut}$ appears between the output terminals $u-u'$, to which a meter M is connected. The voltages on the terminals i, i', u and u' with respect to earth are respectively $E_{if} + E_{it}, E_{if} - E_{it}, E_{uf} + E_{ut}$ and $E_{uf} - E_{ut}$. The voltages E_{it} and E_{ut} are the anti-phase signals; E_{if} and E_{uf} are the (interfering) in-phase signals. E_{if} can be very much larger than E_{it} and have a different frequency.

mately 2×10^4 . This applies to balanced amplifiers in which special measures have been taken to obtain a high degree of symmetry (given perfect symmetry the factor B in eq. (1) is zero and H is therefore infinitely high). The measures consist in using valves having a high amplification factor μ , and incorporating in the common cathode lead an element or circuit having a very high differential resistance R_d

³⁾ A distinction must be made between the guaranteed value and the actual value of the rejection factor. By the guaranteed value H_{min} is understood the calculated value of H for the case where all parameters of the circuit that influence the rejection differ about 10% in the most unfavourable sense from the nominal value. In reality the most unfavourable set of circumstances will seldom be encountered. The actual (measured) value of H will therefore nearly always be appreciably higher than H_{min} ; see, for example, figs 13, 15 and 17 in article I ²⁾.

*) Philips Research Laboratories, Eindhoven.
¹⁾ G. Klein and J. J. Zaalberg van Zelst, General considerations on difference amplifiers, Philips tech. Rev. 22, 345-351, 1960/61.
²⁾ G. Klein and J. J. Zaalberg van Zelst, Circuits for difference amplifiers, I and II, Philips tech. Rev. 23, 142-150 and 173-180, 1961/62 (Nos. 5 and 6).

(this can be achieved in various ways). The mentioned value 2×10^4 is in most cases amply sufficient.

A measurement problem that arose some time ago nevertheless prompted us to study the question whether difference amplifiers could be developed with an even higher rejection factor. The problem concerned a set-up for measuring the Hall effect on test bars of various materials. The set-up is shown schematically in *fig. 2* and explained in the caption. The Hall voltage to be measured, which varied in frequency from 50 to 200 c/s, was extremely small (of the order of a microvolt) and called for a million-fold amplification. This would not have presented any special difficulties if the potential drop in some test bars had not been so large as to give the electrodes a voltage E_{if} with respect to earth that was 10^5 to 10^6 times higher than the voltage to be measured; moreover E_{if} had the same frequency as the Hall voltage. This made it necessary to give the amplifier a guaranteed minimum rejection factor with the exceptionally high value of more than one million. In the most unfavourable case, E_{if} could reach an amplitude of about 200 V.

To be able to guarantee a higher rejection factor than 2×10^4 , it is desirable to increase not only the amplification factor μ but also the differential resistance R_d in the common cathode lead. However, the latter is only meaningful as long as the impedance, formed by the stray capacitance bypassing R_d , is higher than R_d . One can thus obtain better rejection the smaller the stray capacitance and the lower

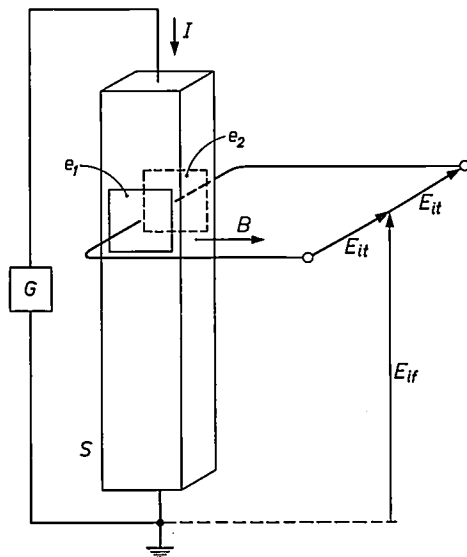


Fig. 2. Measurement of Hall effect. The generator G sends a current I through the test bar S , which is situated in a magnetic field with induction B perpendicular to the direction of I . The Hall voltage to be measured appears between the electrodes e_1 and e_2 , perpendicular to the vectors I and B . Following the nomenclature in *fig. 1*, the Hall voltage is the input anti-phase signal $2E_{if}$, and the potential drop across the test bar gives rise to an input in-phase voltage E_{if} , which is much larger than E_{if} .

the frequency. Along these lines it might perhaps be possible, at low frequencies and given an extremely careful circuit design, to guarantee a rejection factor of 10^5 . In the case just mentioned, however, the rejection factor required was at least 10 times higher. To meet this requirement we therefore had to look for a solution in another direction.

If an AC signal is to be measured — as in this case — we might in principle look for a solution in the use of an isolating transformer between the signal source and the input of the amplifier (which need not then be a balanced amplifier). Two reasons, however, stopped us from adopting this method:

- 1) The input impedance of the transformer should not be much lower than the internal resistance of the signal source. If this resistance is high (as it is in a Hall voltage generator) and if moreover the frequency is low, it follows from this condition that the primary of the transformer would have to have an impractically high inductance.
- 2) It is difficult to screen the windings of the transformer sufficiently from each other to keep the influence of a large in-phase voltage E_{if} small enough.

Combination of a “floating” and an ordinary difference amplifier

In article I ²⁾, an electrically “floating” amplifier was mentioned in passing as a possible solution of the problem (page 142), i.e. a (non-earthed) amplifier which as a whole closely “follows” the interfering in-phase voltage. It was mentioned that this method leads as a rule to complicated and unwieldy constructions, and that in nearly all cases a balanced amplifier with a very large cathode resistance offers a simpler solution.

This most certainly applies if the rejection factor to be guaranteed does not exceed a value of several times 10^4 . If, however, it is to be of the order of 10^6 , then this is an instance where a floating amplifier offers advantages.

The circuit decided upon consists of two difference amplifiers in cascade (*fig. 3*), the first being a floating amplifier. This passes E_{if} as the in-phase signal with virtually no change ($E_{uf} \approx E_{if}$), but considerably amplifies the anti-phase signal $2E_{if}$ which is to be measured. At its output the ratio of the anti-phase to the in-phase voltage is thus very much better, making it possible to use as the second stage (henceforth referred to as the output stage) a normal difference amplifier having a guaranteed rejection factor of the order of 100.

As will be shown presently, it is in fact possible to design the amplifier in such a way that it can follow the in-phase voltage up to a small fraction $1/K$. In such an amplifier the in-phase voltage at the input terminals is effectively reduced to E_{if}/K . Let A_1 be the amplification of the anti-phase signal, and

H_1 be the rejection factor; we then find at the output of the amplifier the following signals:

anti-phase signal . . . $E_{ut} = A_1 E_{it} + \frac{A_1 E_{if}}{H_1 K}$,

in-phase signal . . . $E_{uf} = \left(1 - \frac{1}{K}\right) E_{if} \approx E_{if}$.

Denoting the amplification of the output stage as A_2 and its rejection factor as H_2 , we can write for

obtainable. With these values it follows from the conditions in (3) that where $H_{\min} = 2 \times 10^6$:

$$K > 200 \text{ and } A_1 \geq 2 \times 10^4.$$

In other words, the floating amplifier must amplify the anti-phase signal at least 20 000 times and must follow the in-phase voltage closely to within 0.5%. For the amplifiers presently to be described, $A_1 \approx 50\,000$ and H_2 is greater than 500; consequently $A_1 H_2$

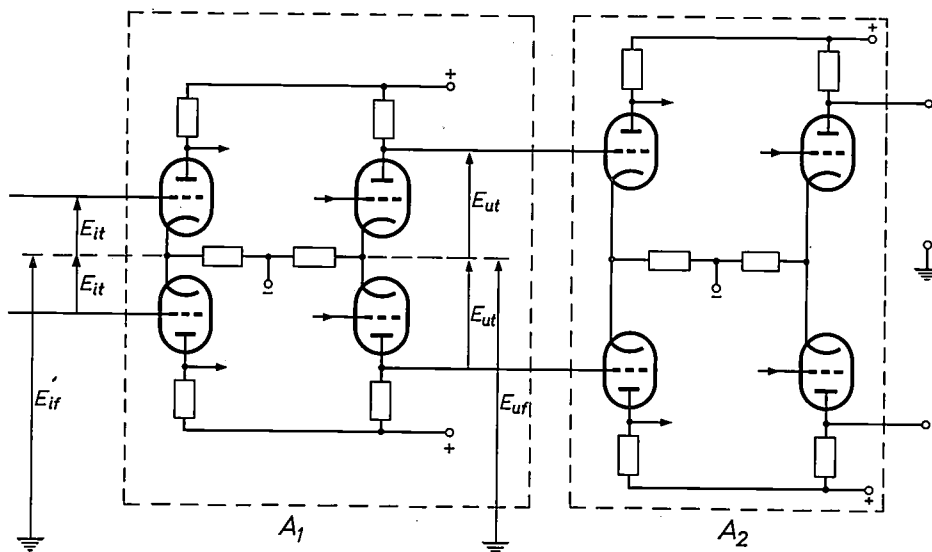


Fig. 3. Cascade circuit of an electrically "floating" difference amplifier A_1 and a normal difference amplifier A_2 ; only the first and last stages are shown in both cases. A_1 has a high gain for the measured signal $2E_{it}$ ($E_{ut} = A_1 E_{it}$) and "follows" as a whole the interfering in-phase voltage E_{if} , so that $E_{uf} \approx E_{if}$. A moderately high rejection factor is therefore sufficient for amplifier A_2 .

the anti-phase signal at the output of this amplifier:

$$A_2 \left(A_1 E_{it} + \frac{A_1 E_{if}}{H_1 K} \right) + \frac{A_2}{H_2} E_{if} = A_1 A_2 \left[E_{it} + \left(\frac{1}{H_1 K} + \frac{1}{A_1 H_2} \right) E_{if} \right].$$

The total rejection factor of the two amplifiers in cascade is thus:

$$H_{\text{tot}} = \frac{1}{\frac{1}{H_1 K} + \frac{1}{A_1 H_2}} \dots (2)$$

From this it appears that both $H_1 K$ and $A_1 H_2$ must be greater than the required rejection factor H_{\min} , even when H_1 and H_2 , as a consequence of parameter values differing by $\pm 10\%$ ³⁾, have their minimum values:

$$H_1 \min K > H_{\min} \text{ and } A_1 H_2 \min > H_{\min} \dots (3)$$

If we suppose the first amplifier to consist of two stages and the second of one stage, then values such as $H_1 \min = 10^4$ and $H_2 \min = 10^2$ are readily

is greater than 25×10^6 . Further, K is indeed roughly 200, while $H_1 \min$ is roughly 2×10^4 . A value of roughly 4 million can therefore be guaranteed for the total rejection factor, which was more than enough for our purpose.

In the following we shall discuss the difference amplifier designed by us for measuring Hall voltages. The problems examined will mainly concern the floating amplifier. At the end of the article we shall touch briefly on the difficulties involved when the signal to be measured is a DC voltage or has a frequency very much higher than 200 c/s.

The floating difference amplifier

The floating difference-amplifier consists of three stages. The first and second stages have identical circuits and together amplify the signals about 50 000 times. Each of the stages consists of two cascodes in balanced configuration and have a triode with cathode resistor in the common cathode lead. A simplified diagram of one stage can be seen in fig. 4. For the most part it is a combination of the

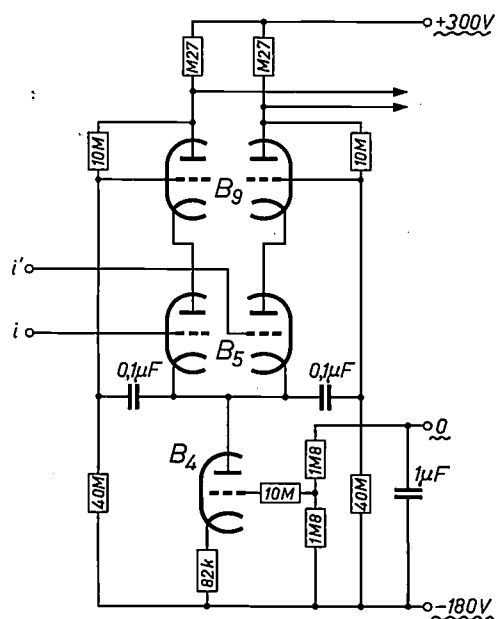


Fig. 4⁴⁾. Simplified circuit of the first stage of the floating difference amplifier. The double triodes B_5 and B_9 (type UCC 85) form a balanced configuration of two cascodes having a high effective amplification factor. Incorporated in the common cathode lead is a triode B_4 ($\frac{1}{3}$ UCC 85) with cathode resistor; this combination offers a high differential resistance. A high amplification factor and high (differential) resistance in the cathode lead are prerequisites for high rejection.

Correct biasing of the "upper" cascode valves is ensured by voltage dividers. The wavy line under +300 V, 0 and -180 V indicates that these DC voltages are "floating", i.e. that an AC voltage with respect to earth (in-phase signal) is superimposed on them.

circuits shown in figures 6 and 8 in article I, where a description will be found of its operation. We shall only recall here that each cascode behaves as a triode having a very high amplification factor μ , and that the triode with cathode resistor provides the high differential resistance in the cathode lead. In this way both conditions for a high rejection factor are fulfilled.

As the successive stages are coupled together capacitively, they do not pass DC signals. At the output it is therefore not directly perceptible whether in the long run the biasing of the balanced cascodes shows disparities. In that case the valves would be operating in the non-linear part of their characteristic, resulting in distortion and hence in an error of measurement. For this reason identical (and constant) biasing of the individual stages is ensured by connecting a voltage divider to the grids of each of the "upper" cascode valves (unlike the situation shown in fig. 6 in article I). As can be seen in fig. 4, these voltage dividers are not connected to the positive terminal of the power supply but to the anode of the relevant valve. This has the effect of producing strong negative feedback for direct voltage, which makes the biasing of the cascodes practically

⁴⁾ In figs 4, 6, 7 and 8 use is made of a space-saving notation for resistance values. For example,

270 means	270 Ω ,
82 k	82 k Ω ,
10 M	10 M Ω ,
1 k2	1.2 k Ω ,
M 27	0.27 M Ω .

equal. Decoupling capacitors between the above-mentioned grids and the common cathode ensure that there is no significant feedback for alternating voltage of the signal frequency. Another result of the decoupling capacitors is that the grids closely follow the voltage on the common cathode, which is a prerequisite for good difference amplifiers.

We shall deal presently with the third stage of the floating amplifier. First we shall consider the measures needed to make this amplifier "float", that is to make it follow the in-phase voltage as closely as possible.

The guiding principle in this connection is that *no point of the circuit should have any perceptible capacitance with respect to earth*. Otherwise, of course, the large in-phase voltage on such a point would give rise via this capacitance to a current that would upset the balance, at the expense of the guaranteed rejection factor. This requirement becomes all the more important the more "sensitive" is the point in question, in other words, electrically speaking, the closer it is to the input terminals and the farther away from the middle of the stages.

To satisfy the above principle the first thing to do is to enclose the amplifier in a metal can (S, fig. 5), which itself follows the in-phase signal with respect to earth. Although the can possesses capacitance (and indeed a fairly high capacitance) with respect to earth, the stray currents flowing via this capacitance are bypassed through an auxiliary cascode (see below) which keeps them out of the amplifier. The can is mounted on elastic strips, which keep it electrically insulated from the earthed chassis and also counteract microphony.

The circuit inside the can has of course various external connections. There is consequently a danger that certain points of the circuit will still "see" earth. These points are: the input terminals, the output terminals, the DC supply terminals, and the cathodes, which possess capacitance with respect to the (externally fed) heaters.

Input terminals. The input terminals are extremely sensitive. They are therefore arranged in such a way (see fig. 5) that their capacitance to earth is only 0.35 pF. They are connected to the object under measurement by cables whose screening is connected to the can, as a result of which the capacitance of 0.35 pF is increased by 0.6 pF per metre length of cable. Using 1 metre of cable, which will usually be long enough, the input capacitance thus remains below 1 pF.

Output terminals. If the (non-floating) output stage were directly connected to the second stage, the rejection of the floating amplifier would be spoilt by the fairly considerable stray capacitance

of the connections between the two amplifiers. Interference voltages might also be induced in these connections. Both harmful effects are smaller the lower are the output resistances in the last stage of

Cathodes. The effect of the cathode heater capacitances can be eliminated by preventing alternating voltage from appearing between the cathodes and the heaters. A circuit for the heater supply has been

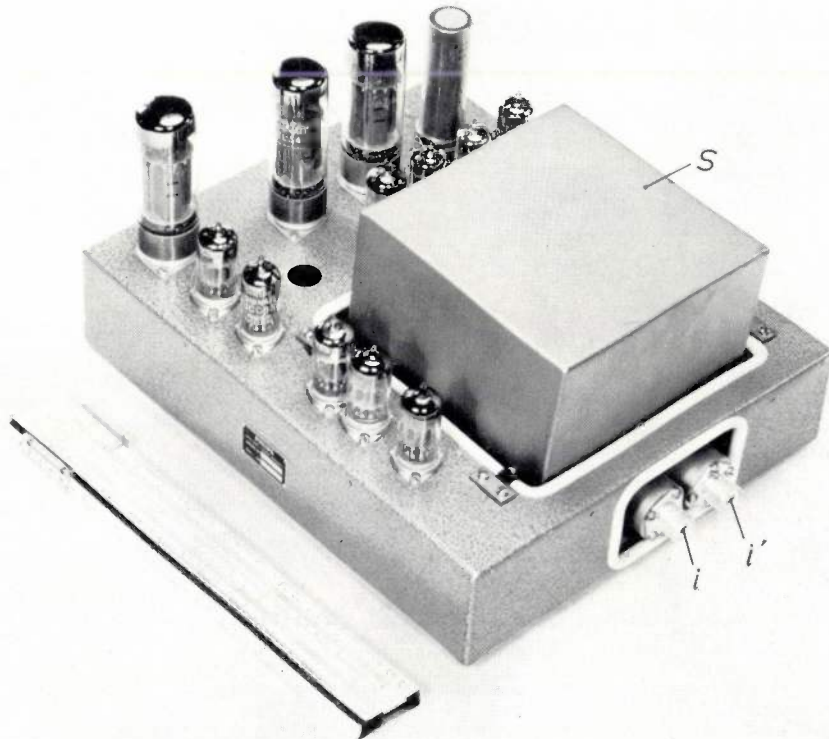


Fig. 5. The difference amplifier with an amplification of about 3.5 million and a guaranteed rejection factor of 4 million. A floating and an earthed difference amplifier are connected in cascade. The floating amplifier is enclosed in a metal screening can *S* and has screened input terminals *i-i'*. The can carries alternating voltage (maximum approx. 140 V r.m.s.) with respect to the earthed chassis, and is mounted on elastic, insulating strips.

the floating amplifier. To this end the amplifier is provided with a balanced cathode follower as a (non-amplifying) third stage with relatively low output resistances ($\approx 1/S$, where *S* is the transconductance of the valves).

DC supply voltages. The floating amplifier requires two DC supply voltages with respect to the cathodes of *B₅* (fig. 4): +300 V and -180 V. The sources of these voltages must follow the in-phase signal *E_{if}*. A simple and adequate, though not particularly practical solution would be to take the voltages from batteries mounted inside the can. In principle the voltages might also be taken from two floating power packs, but in that case the screening between the primary and other windings of the power transformer would have to meet extremely severe demands. It will be shown below how this difficulty was circumvented by using *earthed* power packs in conjunction with a few auxiliary valves.

designed which meets this requirement most satisfactorily. It is simpler than the circuit supplying the floating voltages of +300 V and -180 V, and will therefore be discussed first.

Floating heater-current supply

The circuit shown in fig. 6 neutralizes the cathode heater capacitances by ensuring that the heaters of the relevant valves receive the same alternating voltage (*E_{it}*) with respect to earth as the cathodes. All these valves are double triodes of the type UCC 85, which requires a heater current of 100 mA. This current here is direct current. This is necessary in the first place because it avoids the hum interference which is present to some extent if the indirect cathodes are heated with AC. In the second place, direct current is needed because the floating supply can then benefit from a useful property of pentodes, which is that in their normal operating region their

differential resistance is much higher than their DC resistance.

The heaters concerned are connected in series and are fed from two power-supply units, one of which delivers +550 V and the other -525 V with respect to earth. Connected in series with the heaters there is at one end a pentode B_{15} (fig. 6), which acts as a current source and is biased to an anode current of 100 mA. At the other end is a pentode B_{16} , which works as a cathode follower. The control grid of B_{16} is connected to a point x elsewhere in the circuit (see fig. 8), which is not only at the DC potential that gives the valve its proper operating point but also carries with respect to earth the alternating voltage E_{if} . As a result of the latter the whole chain from the cathode of B_{16} to the anode of B_{15} acquires practically this same alternating voltage: owing to the high differential resistance of B_{15} and the high μ of B_{16} , the latter valve works as a cathode follower, the cathode of which follows the grid voltage very

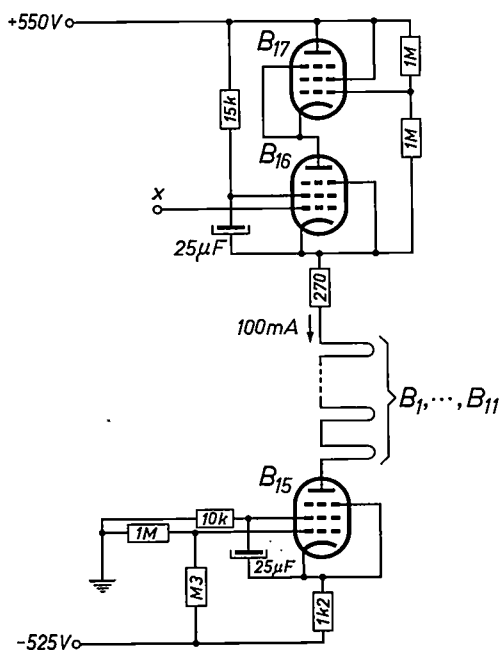


Fig. 6^a). Circuit for supplying in series the heaters of eleven valves type UCC 85 (B_1, \dots, B_{11}) with DC current (100 mA) and for neutralizing cathode-heater capacitances in the amplifier valves.

Pentode B_{15} with cathode resistance functions as current source and is biased for an anode current of 100 mA. Pentode B_{16} works as a cathode follower, owing to its high μ and the high differential resistance of B_{15} . The cathode of B_{16} thus follows the voltage on the control grid; this is connected to a suitable point x elsewhere in the circuit (see fig. 8) which carries the alternating voltage E_{if} with respect to earth. The high differential resistance of B_{15} prevents any appreciable alternating current flowing in the circuit; all eleven heaters therefore receive practically the same AC voltage E_{if} . Since this is also on the cathodes of the valves in the amplifier, no stray currents flow between the heaters and cathodes.

Valve B_{17} is needed to prevent the dissipation of B_{16} from becoming excessive. B_{15} , B_{16} and B_{17} are power pentodes, type EL 34.

closely and thus acquires almost the alternating voltage E_{if} with respect to earth. Moreover, a further result of the high differential resistance of B_{15} is that the AC component in the chain remains limited to a very small value, so that all points of the chain show roughly the same alternating potential with respect to earth. This alternating potential is that which is applied to the cathode of B_{16} , i.e. the signal E_{if} . Since the cathodes of the valves in the floating amplifier also follow E_{if} , the cathode-heater capacitances have no effect.

The DC voltage supply for the circuit is $550 + 525 = 1075$ V, a value which is needed for the amplifiers because of the high amplitude of E_{if} . To prevent the permissible dissipation of B_{16} thereby being exceeded, a second valve is connected in series with it (the pentode B_{17} , circuited as a triode) which takes part of the voltage. The valves B_{15} , B_{16} and B_{17} are type EL 34 power pentodes.

There are altogether eleven UCC 85 valves, whose heaters are series-fed in this way. Seven of them belong to the floating amplifier and two to the output stage; the remaining two are auxiliary valves, which will be discussed presently. The other valves used in the amplifier are E types, the heaters of which are fed with alternating voltage.

Floating DC supply voltages

Fig. 7 shows only the "lower" triodes (B_5) of the first stage of the floating amplifier. The rest of the diagram gives the circuit (simplified) which is needed to make it possible to use earthed power-supply units. These are here the same units from which the heater current is obtained: they deliver +540 V (10 V lower than the point where the heater current is derived, owing to extra smoothing) and -525 V with respect to earth; the latter voltage is highly stabilized.

We shall now show how the appropriate DC voltages are obtained with the circuit in fig. 7, and then discuss the way in which the whole floating amplifier is made to follow the in-phase voltage.

As mentioned, the floating amplifier requires +300 V and -180 V with respect to the cathodes of B_5 , which carry, in addition to a small DC voltage, the in-phase voltage E_{if} with respect to earth. At first sight it might seem obvious to connect that point (the cathodes of B_5) with the screening can. If that were done, however, the capacitive currents flowing via the can would not remain outside the amplifier. For this reason we connect between the points +300 V and -180 V an auxiliary cascode (B_3a - B_3b), the middle of which (point 0) can safely be connected to the can. Since B_3a is a very faithful

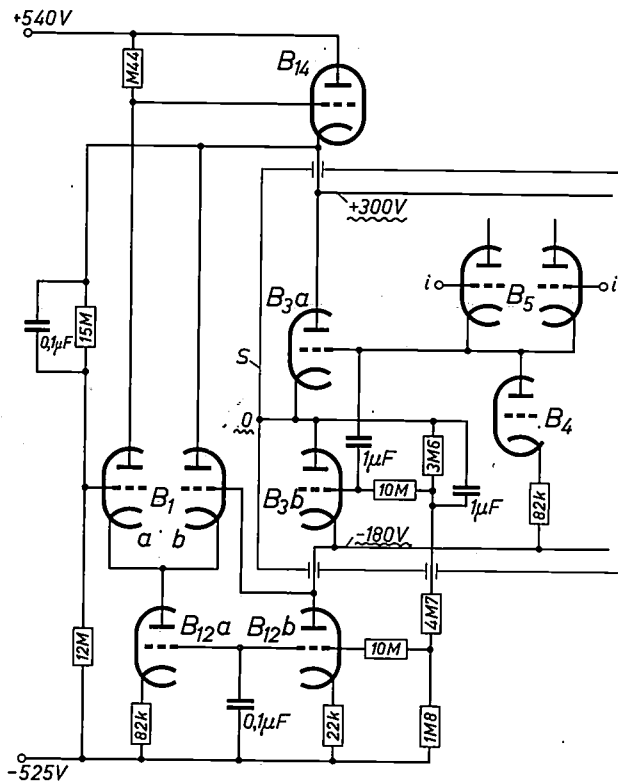


Fig. 7 4). Circuit for delivering the floating supply voltages +300 V and -180 V to the floating difference amplifier. B_5 belongs to the first stage of this amplifier (cf. fig. 4). $i-i'$ input terminals. B_3a - B_3b auxiliary cascode with current source $B_{12}b$ in the cathode lead. Since B_3a is a cathode follower, its cathode (point 0 connected to the can S) has virtually the same DC voltage (≈ 0) and alternating voltage (E_{if}) as the cathodes of B_5 . A voltage divider $3.6 + 4.7 + 1.8 \text{ M}\Omega$ gives a DC voltage of -180 V to the grid of B_3b , which also receives the AC signal E_{if} via a $1 \mu\text{F}$ capacitor. Since B_3b is also a cathode follower, both these voltages appear on the cathode of B_3b (point -180 V).

The series triode B_{14} , controlled by the auxiliary amplifier B_{1a} - B_{1b} - B_{12a} , ensures that the point +300 V acquires both the DC voltage (determined by the voltage divider $15 + 12 \text{ M}\Omega$) and the alternating voltage E_{if} .

cathode follower (its cathode lead contains the very high differential resistance of the triodes B_3b and $B_{12}b$, see fig. 7), the point 0 has virtually the same potential as the cathodes of B_5 .

A voltage divider $3.6 + 4.7 + 1.8 \text{ M}\Omega$ between the points 0 and -525 V ensures that the point -180 V in fact carries this DC potential. This is because B_3b is also a good cathode follower, so that its cathode has practically the same potential as its grid, which is connected to this voltage divider.

The required DC potential at the point +300 V is obtained by using the well-known principle of stabilized power supplies: from the point +540 V the supply is effected through a series triode (B_{14}) which is controlled by an auxiliary amplifier, using the above derived DC voltage of -180 V as reference voltage. This auxiliary amplifier is a difference amplifier, consisting of a balanced arrangement of the double triode B_{1a} - B_{1b} with the triode B_{12a} in the

cathode lead. The grid of B_{1b} is connected to the point -180 V, the grid of B_{1a} to the tap of a voltage divider ($12 + 15 \text{ M}\Omega$) between the points +300 V and -525 V. The voltage division ratio is so chosen that when the point +300 V has the correct potential, the grids of B_{1a} and B_{1b} carry the same voltage. If, for example, the potential of the point +300 V increases, the anode current in B_{1a} increases, the anode voltage of B_{1a} (hence the grid voltage of B_{14}) consequently decreases, and since B_{14} is also a cathode follower, the cathode voltage of B_{14} — i.e. the potential of the point +300 V — likewise drops. In this way, then, the potential of this point is automatically kept constant.

It is now easily seen that the points 0, -180 V and +300 V closely follow the in-phase voltage E_{if} . This is present at the cathodes of B_5 and therefore also at the grid of B_3a and, via a $1 \mu\text{F}$ capacitor, at the grid of B_3b . (A $10 \text{ M}\Omega$ resistor between the grid of B_3b and the tap on the voltage divider of $3.6 + 4.7 + 1.8 \text{ M}\Omega$ prevents undesired loading of the voltage divider via the capacitor mentioned.) Because of the cathode-follower properties of B_3a and B_3b , their cathodes too will follow the in-phase voltage, and these cathodes are connected to point 0 and point -180 V respectively. The reference voltage on the grid of B_3b is therefore the sum of a DC and an AC voltage. Consequently the cathode of B_{14} (the point +300 V) carries with respect to earth not only the correct DC voltage but also the correct AC voltage (the in-phase signal E_{if}).

The complete diagram of the amplifier (without the earthed power-supply units and the heater-supply circuit) is shown in fig. 8. The thin line represents the screening can. The valves have the same numbering as in the previous figures. As can be seen, the amplifier that drives the triode B_{14} in fact consists of a balanced arrangement of two cascodes, giving a greater amplification than the balanced arrangement of two triodes in fig. 7. The control grid of B_{16} (fig. 6) is connected to the point x (fig. 8).

Measurements of the above-mentioned factor $1/K$ at 150 c/s have shown that the point +300 V follows the in-phase voltage to within 0.53%, while the points 0 and -180 V follow this voltage to within 0.36%.

Further details

The output stage

As stated (see fig. 3), the floating difference amplifier (A_1) is followed by a non-floating difference amplifier (A_2) as output stage. The circuit of A_2 is shown on the right in fig. 8. Here too, a balanced configuration of two cascodes is used (B_8 - B_{11} , both

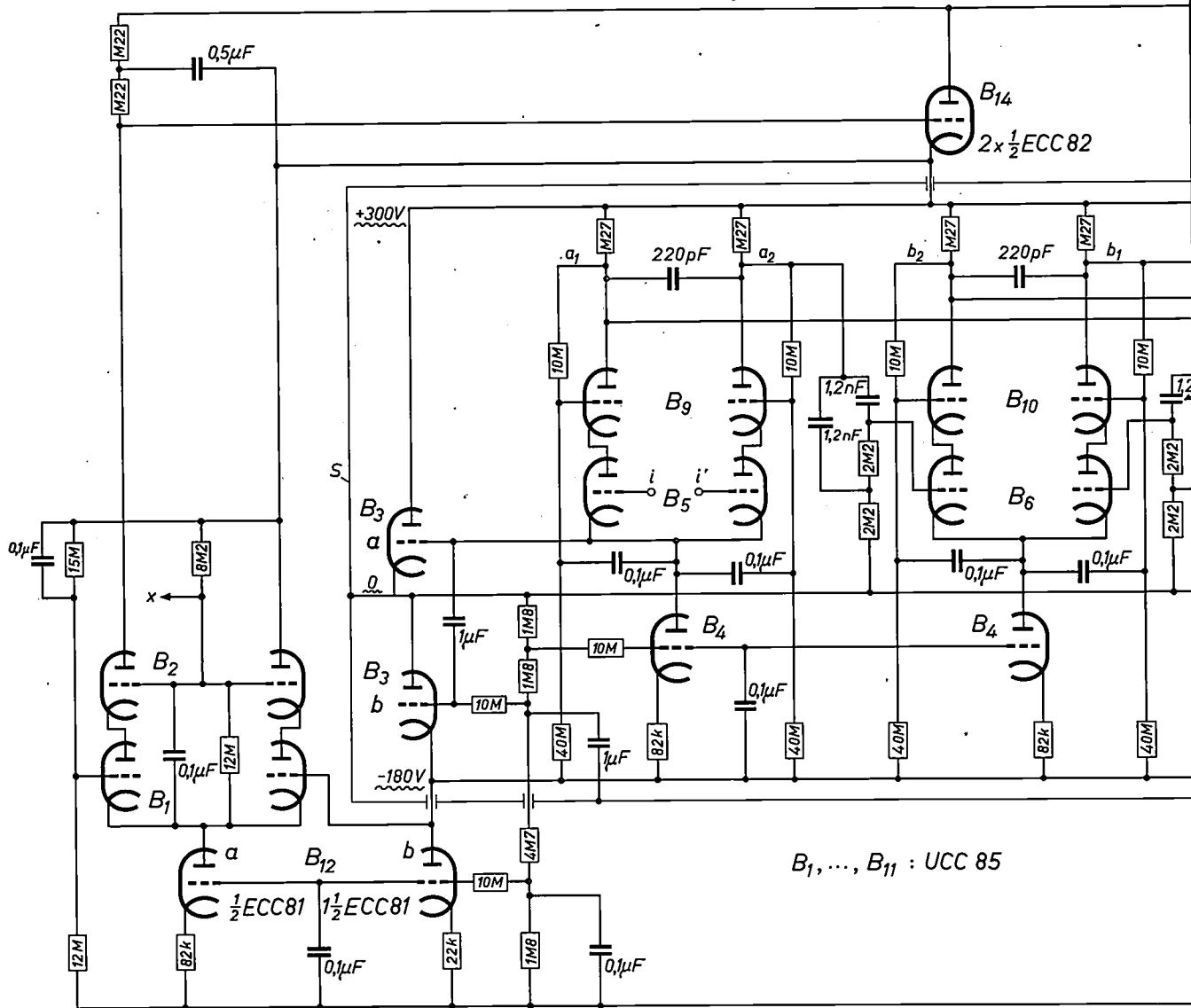


Fig. 8⁴). Complete circuit (but excluding the power supply for the heaters) of the difference amplifier shown in fig. 5. The valves are numbered in the same way ($B_1 \dots B_{14}$) as in the previous figures. $i-i'$ input. $u-u'$ output. S screening can enclosing the floating amplifier. The control grid of valve B_{14} is connected to the point x (fig. 6).

types UCC 85), having in the cathode lead a current source consisting of a triode (one half of an ECC81) with cathode resistor. The circuit was designed to deal with the high amplitude that E_{if} can acquire (approx. 200 V), which called among other things for the high supply voltages. At the output the amplified input signal appears at sufficient amplitude, and sufficiently free from in-phase voltage, that it can be further processed by conventional means.

Stepwise gain control using photoresistors

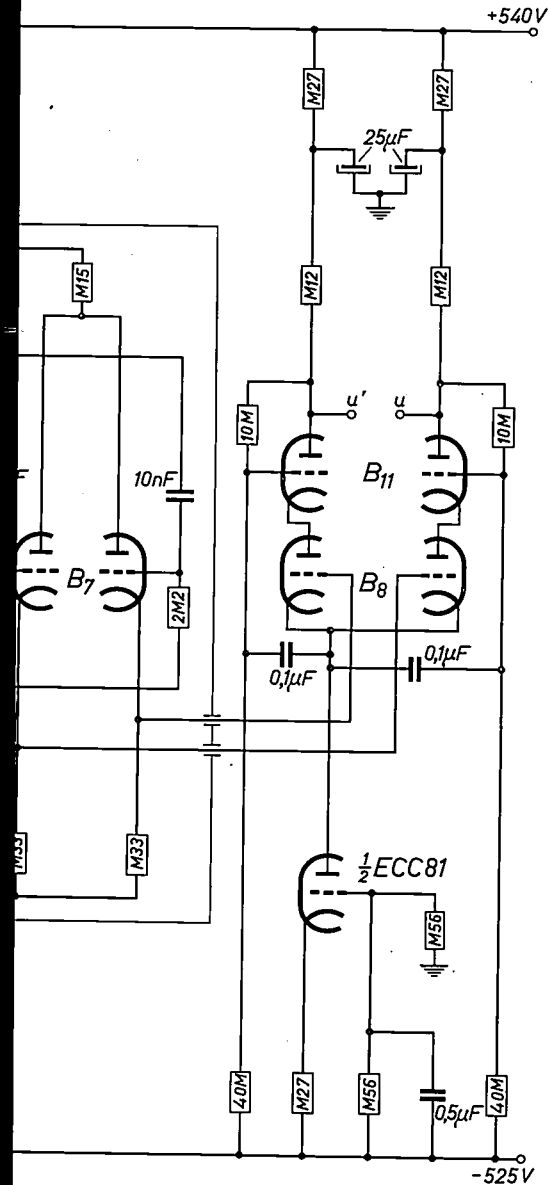
The complete amplification of several times 10^6 is not always needed. Often a gain 10 or 100 times smaller will be adequate, provided the rejection factor is not reduced to a greater extent.

A suitable method of thus reducing the gain is to introduce a resistor between a_1 and a_2 or between b_1 and b_2 (fig. 8⁵): for a 10 times as small a gain a 47 kΩ resistor between a_1 and a_2 , and for a 100 times as small a gain a 4.7 kΩ resistor between b_1 and b_2 .

These resistors are switched in and out of circuit by means of *photoresistors*⁶), which do not have the drawbacks of ordinary switches or relays mounted in a can that electrically and mechanically "floats".

⁵) For methods of gain control, see article II, mentioned in footnote ³), p. 175 *et seq.*

⁶) N. A. de Gier, W. van Gool and J. G. van Santen, Photoresistors made of compressed and sintered cadmium sulphide, Philips tech. Rev. 20, 277-287, 1958/59.



The type of photoresistor used has a dark resistance of more than 100 MΩ, and when illuminated by a small bulb a resistance of only about 150 Ω. Photoresistors of this type can thus serve as switches for resistors with values much higher than 100 Ω and also much lower than 100 MΩ. This condition is fulfilled by the 4.7 and 47 kΩ resistors required in this case. The (earthed) electric bulbs are so positioned that they do not appreciably increase the capacitance with respect to earth of the components inside the can (fig. 9).

Results

The total amplification of the anti-phase signal is about 3.5×10^6 . Despite this very high figure, there is not the least tendency towards oscillation ⁷⁾.

⁷⁾ The explanation will be found in article II, mentioned in footnote ²⁾, p. 177.

During a period of many hours the amplification changes by no more than 0.5%. (If necessary this change can be made still smaller by replacing the carbon resistors, used for the anode load, by metallic resistors.)

As mentioned, a value of 4×10^6 can be guaranteed for the rejection factor. The values measured at 150 c/s on five UCC 85 double triodes at the position B₅ (fig. 8) varied from 8×10^6 to 15×10^6 .

In order to neutralize part of the various interfering voltages (including noise and hum), measures were taken to reduce the amplification outside the signal frequency range (50 to 200 c/s). Below 50 c/s this is done by the RC networks between points a₁-a₂ and the grids of B₆; above 200 c/s by the 220 pF capacitors between a₁ and a₂ and between b₁ and b₂ (fig. 8). The total interference voltage, derived at the input, varied from 1.2 to 2.0 μV on the five UCC '85 valves tested. The frequency pass band can be given a sharper cut-off by conventional methods behind the amplifier (high-Q filters, selective detection).

Owing to the exceptionally low input capacitance (less than 1 pF per metre of cable), even considerable unbalance in the signal source is not able to spoil the rejection factor. If, for example, two signal sources with identical e.m.f.'s are connected, one of which has zero internal resistance and the other an internal resistance of 2000 Ω, it is still possible to guarantee the very high rejection factor of 1.5×10^6 (at 50 c/s).

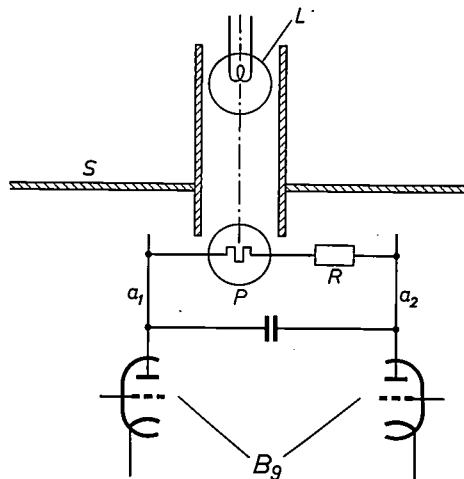


Fig. 9. By connecting a 47 kΩ resistor R between points a₁ and a₂ (fig. 8) the amplification is reduced by a factor of 10. The switch used, which introduces no significant stray capacitance, is a photoresistor P, which has a dark resistance of more than 100 MΩ, and, upon illumination by the electric bulb L, a resistance of only about 150 Ω.

In the same way it is possible to connect between b₁ and b₂ (fig. 8) a 4.7 kΩ resistor, which reduces the gain by a factor of 100.

Other signal frequencies

To conclude we shall touch briefly on the difficulties to be expected if the principles described are to be applied at higher or lower signal frequencies.

The higher the frequencies the greater is the influence of stray capacitances. This means that the floating difference amplifier will not follow the in-phase voltage so closely, and that the rejection factor will therefore be lower. On the other hand, it is pointless to ask for a particularly high rejection factor at high frequencies, for asymmetries in the connections between the amplifier and the signal source make a very high rejection factor in any case impossible⁸⁾.

The lower the signal frequency the better the floating amplifier follows the in-phase voltage, but the more substantial are the ordinary drawbacks encountered at very low frequencies, namely that coupling and decoupling capacitors of very high capacitance then have to be used. Where DC signals are involved the difficulty may arise, especially if

⁸⁾ See figs 18 and 19 of article II, mentioned in footnote ²⁾.

the in-phase voltage (and hence the supply voltage) is high, that the dissipation of some valves, which are required to pass continuously a large current at high voltage, will become excessive. In that case each of the valves involved will have to be replaced by two valves in series, so that the voltage drop — and thus the power to be dissipated — will be divided over the two valves. A similar case arose in the heater power supply (fig. 6): here it was necessary to connect a valve B_{17} in series with B_{16} .

Summary. For measuring Hall voltages (minimum value of the order of $1 \mu V$) in the frequency range from 50 to 200 c/s, an amplifier with a gain of one million was needed. The measuring electrodes could acquire such a high potential with respect to earth (the in-phase voltage) that it was necessary to give the amplifier a guaranteed rejection factor of at least one million. This problem has been solved by using two difference amplifiers in cascade, the first of which is electrically "floating", i.e. closely follows the in-phase voltage (to within about 0.5%). Among the measures discussed are the circuit that neutralizes the cathode-heater capacitance of the valves in the floating amplifier, and the circuit from which the floating DC supply voltages are obtained. The guaranteed rejection factor is 4 million, the amplification about 3.5 million. Resistors for reducing the gain by a factor of 10 or 100 are switched in and out of circuit by means of phototransistors.

THE AMPLISCOPE, AN EXPERIMENTAL APPARATUS FOR "HARMONIZING" X-RAY IMAGES

by E. ZIELER *) and K. WESTERKOWSKY *).

621.397.33

In order to extract the information contained in an X-ray diagnostic photograph as completely or as easily as possible, new methods have been developed for the inspection of such photographs, making use of electronic equipment. In the Ampliscope, an X-ray photograph is inspected with the aid of two television cameras, the two signals being mixed in an appropriate manner before viewing.

The diagnostic use of X-rays depends on the fact that local variations in the thickness and chemical composition of objects in the human body lead to corresponding variations in the attenuation of X-rays passing through the body. When an X-ray photograph is taken this "X-ray contrast" is converted into local differences in photographic density, in a way which depends on the characteristic curve of the photographic emulsion (fig. 1).

Anatomical abnormalities of diagnostic importance are often small, and give rise to a low X-ray contrast. It would seem an obvious idea to use an emulsion with a steep characteristic curve (high gamma) in order to facilitate the observation of such abnormalities. This solution is often impracticable, because the coarse (large area) structure of

the organ under investigation often gives rise to considerable differences in density, on which the details of diagnostic interest are superposed: if the background density lies near the top or bottom of the characteristic curve, where the slope is less steep (see fig. 1), the fine (detail) contrast will in fact be reduced instead of increased. In order to get usable results, one must thus not only increase the fine contrast but also reduce the coarse contrast. This reduction of the coarse contrast is known as "harmonization" of the X-ray image.

Harmonization can also be useful for the observation of details of sufficient contrast in the medium-density range: if adjacent parts of the X-ray image are not very dense, the observer can suffer from glare. In this case too, therefore, it would be desirable to decrease the coarse contrast, while retaining the fine contrast.

Various methods have been suggested for this purpose (and for analogous purposes in fields other than X-ray diagnostics). In this article we will describe a method developed in the X-ray equipment factory of C. H. F. Müller, Hamburg, which was demonstrated already in 1959 at the 9th International Congress of Radiology in Munich. Since then the apparatus involved has been further improved, although it is still in the experimental stage, and the results obtained with it in practice are promising. A further description in this review therefore seems justified. For the sake of simplicity, we shall call this apparatus by the name which has been given to it in the laboratory: the Ampliscope.

The method is based on the use of television techniques. Before describing it in detail, we would like to discuss some previously used harmonization methods, in particular two methods for obtaining a harmonized print of a given X-ray photo. Our method does not give a print, but an image on a television viewing screen. This image can of course be photographed if desired.

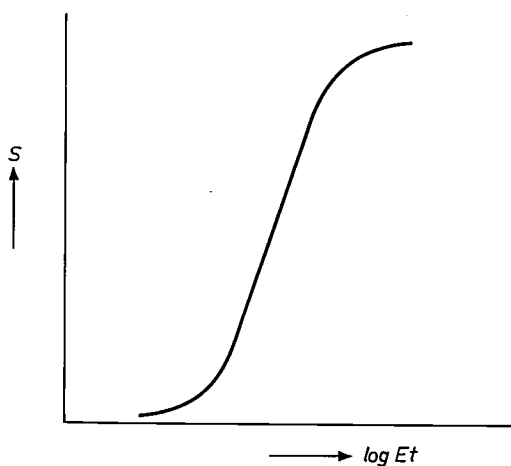


Fig. 1. The characteristic curve of a photographic emulsion. The density $S = \log I_0/I_1$ (where I_0 is the amount of light falling on the negative during observation, and I_1 is the amount transmitted) is plotted as a function of the exposure $E \times t$ (the latter being plotted on a logarithmic scale). In printing, the abscissa is equal to the negative density of the negative (plus some constant), the ordinate then giving the obtained density of the positive print. If the slope of the linear part of the curve is $\alpha = 45^\circ$ (gamma = 1), the brightness values of the original will be faithfully reproduced; if the curve is steeper, the contrast will be increased.

*) C. H. F. Müller GmbH, Hamburg.

Previous harmonization methods

As long ago as 1930 Spiegler¹⁾ suggested a method for making harmonized prints. According to this method, a blurred diapositive is first made from the negative in question, the two are then superimposed and a contact print is made from the combination in the normal way. It may be seen from *fig. 2* that the coarse structure of the object now gives rise to

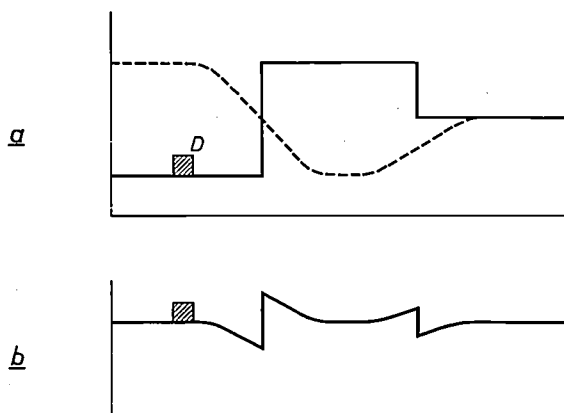


Fig. 2. The principle of Spiegler's method of producing a harmonized photographic print. *a*) The full line represents the local variation of the density of the original negative. *D* is a small detail of low contrast. The other contrasts shown are much coarser. The broken line represents the local variation of the density in the blurred diapositive, which is superimposed as accurately as possible on the negative. *b*) Local variation of the density in the final print. The coarse contrasts are removed (except at the transitions between regions), while the detail contrast of *D* is retained practically entirely.

no differences in the density of the print, except at the boundaries. Details smaller than the blurring of the diapositive are, however, reproduced with practically no decrease in the contrast. If the characteristic curves of the various emulsions used are suitably chosen, this method can give good results; it is however rather complicated and time-consuming, and has therefore never been widely used.

1) G. Spiegler and Kalman Juris, *Fortschr. Röntgenstr.* 42, 509, 1930.

Of recent years, photographic methods of harmonization have been suggested which do not involve making a blurred diapositive first. In these methods, the light source used for making the prints is a fluorescent screen uniformly illuminated by ultraviolet light, the local brightness being altered by means of infrared irradiation (quenching as a result of the Herschel effect) according to the local density of the original²⁾.

A copying method which makes use of electronic equipment (the "LogEtron") has come into use to some extent for radiological purposes, and even more widely for cartography³⁾: the processing of aerial photographs involves problems which are quite similar to those met with in X-ray photography. In this method, a contact print is made of the negative by illuminating it with a flying-spot scanner well known in television⁴⁾; see *fig. 3*. The intensity of the spot of light is controlled during the scanning process by a photocell, which receives the light after it has passed through the negative and the positive material: when a relatively transparent part of the negative is being scanned, the intensity of the light is decreased, and vice versa. This method is much the same as the automatic volume control used in radios, and similarly the lack of inertia in the control unit makes it in principle possible to remove the "modulation" of the image completely. This is of course not what is aimed at: the idea is merely to reduce the modulation. The desired harmonization of the photographic print is now obtained by scanning with a spot of light which is not quite in focus. As in Spiegler's process, the fine contrast in regions smaller than the spot of light

2) United States Radium Corporation (Geneva), Dutch Patent No. 99457. A similar method is described by B. M. Woldringh, *Harmonisation mittels Phosphorographie (Harmonization by means of "phosphorography")*, 9th International Congress of Radiology, Munich 1959, p. 1484.

3) D. R. Craig, *The LogEtron: a fully automatic, servo-controlled scanning light source for printing*, *Phot. Engng.* 5, 219-226, 1954.

4) See e.g. F. H. J. van der Poel and J. J. P. Valetton, *The flying-spot scanner*, *Philips tech. Rev.* 15, 221-232, 1953/54.

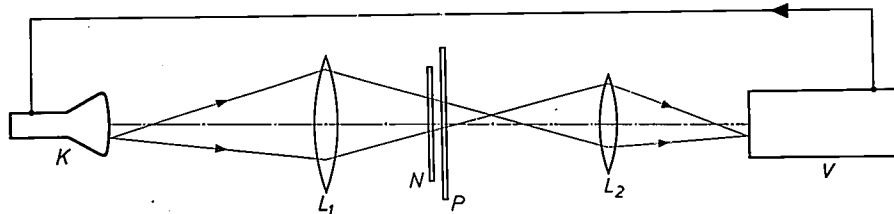


Fig. 3. The principle of the "LogEtron" method³⁾ for making a harmonized photographic print of the negative *N*. *K* tube of flying-spot scanner, whose spot is projected out of focus onto *N* by the lens *L*₁. The light which passes through the negative and the positive material *P* is concentrated by a second lens *L*₂ onto the photomultiplier *V*, whose output signal controls the intensity of the light spot of *K* so that the intensity of illumination of *P*, averaged over the blurred area, remains constant.

is retained while the contrast in larger regions is evened out. Hence, photographic emulsion with a steep characteristic curve can if desired be used for the positive.

Now if we start using television techniques in this way, we may just as well forget about photography altogether for the moment, reproduce the X-ray image on a television screen, and look for ways of harmonizing the image thus produced. This way of producing a harmonized image, if possible, is very attractive since the contrast of a television image can be varied very simply, so that the radiologist can vary the character of the television image continuously while he is looking at it and can thus find the best setting by continually comparing the picture on the television screen with the original. This attractive possibility has been investigated by several other workers besides ourselves.

W. J. Oosterkamp and T. G. Schut, of the Philips Research Laboratories in Eindhoven, have harmonized television images by means of the above-mentioned automatic-volume-control principle⁵⁾: when the image is scanned line by line, a continuous electrical signal (the video signal) is produced, and the average value of this signal over a certain period of time (i.e. over a certain fraction of the length of the line) can be obtained with the aid of a simple integration network. The instantaneous sensitivity can be controlled by suitable feedback of this mean value, so that differences in brightness over distances greater than the above-mentioned fraction are evened out, while differences over shorter distances remain unchanged. It is naturally preferable to take the average of the signal over a segment of the line which lies on either side of the point being reproduced; this can be arranged by passing the video signal through a delay line before adding the harmonizing signal.

In this method, not only do the coarse contrasts disappear, but so do those fine contrasts whose gradient has no component in the direction of scanning. If for example the picture shows a series of black and white stripes, the desired effect is only obtained when the scanning direction is chosen at right angles to the stripes. If this is done, all the stripes become a uniform grey, only their edges being brought sharply into contrast — which is the object of harmonization. The whole structure of stripes, on the contrary, disappears practically completely if scanning is performed parallel to the stripes. This effect can be very clearly observed with the ribs in a picture of the lungs.

Another television procedure for the inspection of X-ray images which we would like to mention in this connection also suffers from this weakness. In this procedure, the video signal from the scanned picture is passed through a differentiating network⁶⁾. This method also emphasizes all *transitions* between regions of differing brightness, the more so as the transition is sharper. Although the principle on which the picture is altered is thus quite different and is not very promising when details are blurred, it is natural to class this method together with that of Oosterkamp and Schut. Here too, the desired effect is produced only so far as those details whose brightness gradient has a component in the direction of scanning.

It will be seen from the following description of our apparatus that we have eliminated this drawback.

The principle of the Ampliscope: double scanning

In the Ampliscope, *two* normal television cameras equipped with vidicons are directed onto the X-ray photo to be inspected (fig. 4). The optical system of one of the cameras produces a sharp image of the photo on the photosensitive layer of the vidicon. The video signal produced by the scanning of this image is used as the basis of the picture to be repro-

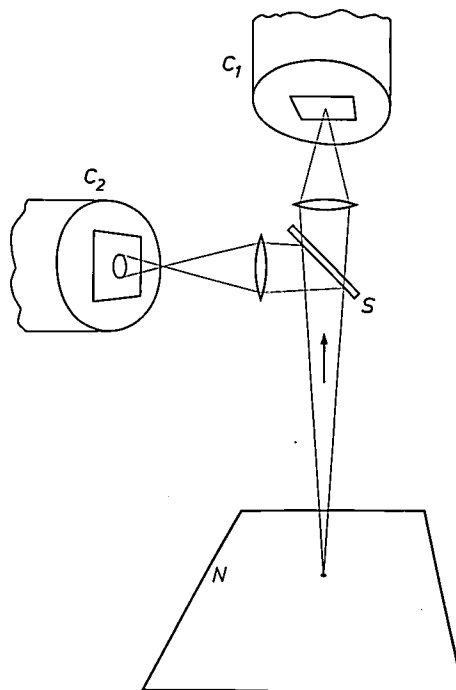


Fig. 4. The principle of the Ampliscope. The X-ray photo *N* to be harmonized is projected in focus onto the photosensitive layer of the vidicon *C*₁, and out of focus, via the half-silvered mirror *S*, onto *C*₂. The two layers are synchronously scanned by their electron beams, which means that the photo is simultaneously scanned by a sharp and by a blurred spot. The video signal proper is provided by *C*₁; the signal from *C*₂ is mixed with this, with reversed sign and variable amplitude, to give the desired harmonization.

⁵⁾ Personal communication.

⁶⁾ K. Bischoff and O. Schott, Fortschr. Röntgenstr. 87, 239, 1957.

duced on the screen of a picture tube. The optical system of the other camera produces a *blurred* image of the photo on the corresponding vidicon; in other words, each point of the photosensitive layer receives light from a certain *region* of the photo, just as if the photo were scanned by a blurred spot of light. The video signal produced by this camera at any given moment is thus a measure of the average density of the film (mean transmission of the X-ray negative) in a region covered by the blurred spot. This signal is added to that from the first camera, with the sign reversed and with a variable amplification. It will be clear from what has been said above that the photo in question can be harmonized in this way. Moreover, the harmonization can very simply be controlled in two ways: in the first place, the extent to which the coarse contrast is removed can be varied by altering the relative amplifications of the two signals. (For example, if they are mixed with equal amplitudes, the whole image produced will have a uniform average brightness.) Secondly, by altering the setting of the optical system of the second camera we can alter the blurring, i.e. alter the area over which the brightness is averaged and thus the size of the details in the original photo in which the contrast is preserved.

The blurred image produced of each point of the original is roughly circular. The direction of scanning (direction of the lines in the picture) is thus no longer of significance, and the character of the harmonized picture is hardly changed if the X-ray photo is turned in its own plane during its inspection. Naturally, this only holds true (and the harmonization will only have the desirable features) if the X-ray photo is projected on the two vidicons in such a way that the images to be scanned are congruent.

The amplification of the contrast

We mentioned in the introduction that the aim of harmonization is to make amplification of the contrast possible. In photographic methods of harmonization, the increased contrast can be obtained directly by making use of the properties of the photographic emulsions (characteristic curve with high gamma). In the television methods, the contrast could be increased in a similar way: in a television image, the relationship between the brightnesses of a given point in the original and the corresponding point on the screen is given by the "transfer function", which corresponds to the characteristic curve of a photographic emulsion. One talks of the "gamma" of a circuit, and this is equal to the product of the gammas of the individual

parts of the circuit. The vidicon tube has roughly $\gamma = 0.5$, the picture tube $\gamma = 2$ to 3 and the video amplifier in the "distortion-free" (linear) region $\gamma = 1$. The gamma of the entire circuit is thus normally approximately equal to 1. It is possible to obtain a transfer function with a greater and moreover a variable value of gamma by introducing an element into the circuit which gives an adjustable distortion⁷⁾. It was found during the construction of the Ampliscope, however, that such an element is not necessary: the desired variable contrast amplification can be obtained much more simply by making use of the "Callier effect".

This effect is based on the fact that the effective density of a film is increased if it is viewed in directional light instead of the usual diffuse light. This is because the attenuation of light in the developed negative is for a large part due to diffusion, and not to absorption; the situation is sketched and further explained in *fig. 5*. Let us consider the two extreme

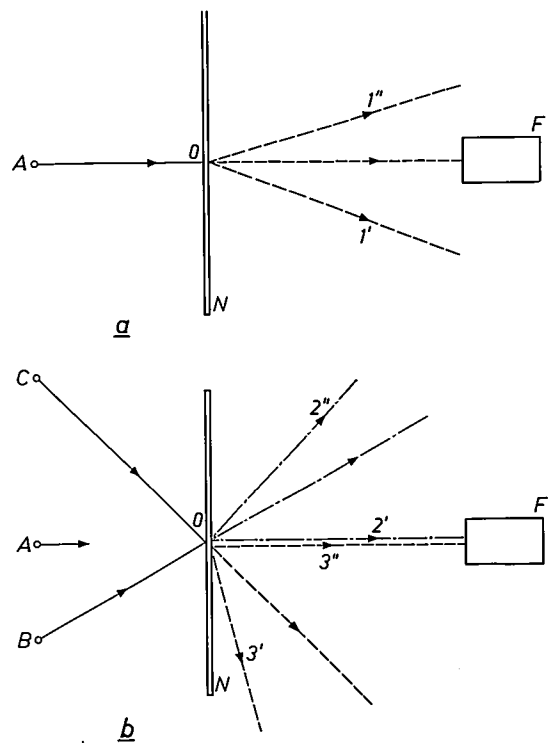


Fig. 5. The Callier effect.

a) If it is desired to measure the density at the point O on the negative N , and if the light I_0 used for this purpose comes from the point A only, then the light I_1 received by the detector F is attenuated not only by absorption but to a large extent by scattering (the portions I' and I'' of the light do not reach the detector).

b) If light also falls on O from the points B and C , then, owing to the fact of the influence of the scattering, a relatively considerable contribution of the light $2'$ and $3''$ also reaches the detector. The total ratio I_0/I_1 is thus smaller than in case a .

⁷⁾ The circuit of such a "gamma corrector" (albeit only for values of gamma less than 1) is described in the article quoted above⁴⁾, pp. 228-230.

cases: completely diffuse illumination, in which the light source subtends a solid angle of 2π at every point of the film, and completely directional illumination, i.e. with a point source of light, in which case the angle in question is zero. If in the first case the difference in density between two parts of the film is $\Delta S = 1.8$, which means that a photocell placed behind one part will get about sixty times as much light as one placed behind the other part, then in the second case we can measure a difference $\Delta S' = 2.9$, which means that one photocell now receives about eight hundred times as much light as the other. The "contrast amplification" $\Delta S'/\Delta S$ (or Callier coefficient, which in this extreme case has the value 1.6) can in principle simply be varied by adjusting the size of the light source. After several preliminary experiments, we adopted an illuminated ground-glass screen with an iris diaphragm for this purpose.

It will be clear that the above argument for a point source of light only holds for the middle of the picture. In order to illuminate *all* parts of the picture with directional light, we placed a condenser lens behind the film, which produces an image of the light source at the entrance pupil of the vidicon camera. In view of the desired large dimensions, it is best to use a Fresnel lens for this purpose.

Some details of the apparatus

Fig. 6 shows the apparatus as a whole. We shall now discuss some details of the construction and circuitry of this apparatus, and finally illustrate the results obtained by a few examples.

Fig. 7 indicates how the various components are arranged. The X-ray photo to be inspected is placed on the (polished) glass plate *H*, which is tilted slightly for the sake of convenience. The above-mentioned Fresnel lens and adjustable light source are behind this glass plate. An image of the photo is produced on the two vidicons of the cameras C_1 and C_2 by their respective optical systems — on C_1 a sharp image and on C_2 a blurred one. To start with we used exchangeable objectives for both cameras, but we later fitted both of them with identical Zoom lenses (which have variable focal lengths) by means of which the scale of the image can be continuously varied from 1 : 1 to 2 : 1. Magnification of the image on the vidicon means that better use is made of the resolution of the television image (which is naturally limited by the line structure⁸⁾).

⁸⁾ In our case, an amplification of 2.5:1 would give the same effective resolution as that of the direct X-ray photo (taken with intensifying foil).



Fig. 6. Experimental model of the Ampliscope. *H* glass plate on which is placed the X-ray negative to be examined. *B* screen of the television tube on which the harmonized picture appears. *1* and *2* controls of the optics of the two television cameras. The degree of harmonization can be varied by means of the knob R_1 . (The picture shown on the screen can be changed from a negative to a positive by simply reversing the signals; see e.g. the article cited in ref. 4.) The knob R_2 controls the contrast amplification, and R_3 the mean brightness of the picture. To the left of *H*, the main switch, a switch for the light source which illuminates the photo to be inspected, and a switch for the lighting of the room are to be seen.

The scanning electron beams in both cameras are controlled by the pertaining control units S_1 and S_2 , whose frequency generators are coupled so as to synchronize the scanning of the two images. The signals obtained from C_1 and C_2 are added in the mixer stage *M*, and the total signal is then applied to the video amplifier of one of the cameras. (In our case, the amplifier of the other camera is not used.) The picture tube *B* on which the harmonized picture appears is placed roughly at eye level, for ease of observation.

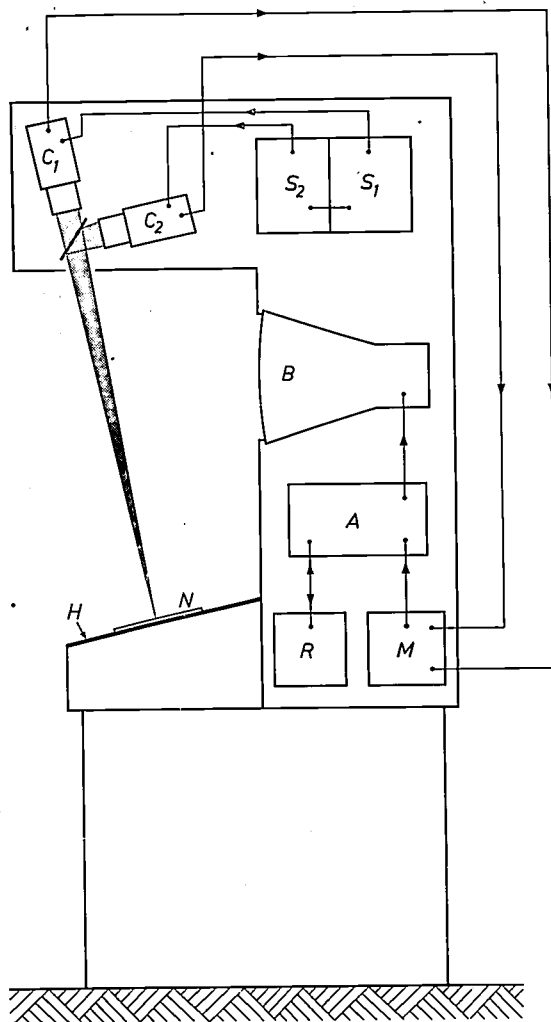


Fig. 7. Sketch showing the construction of the experimental model of the Ampliscope. *H* glass plate on which is placed the X-ray photo *N* to be inspected; the light source and Fresnel lens are situated behind this plate. *C*₁ and *C*₂ two vidicon cameras (type No. GM 4930) with Zoom lenses. *S*₁ and *S*₂ control units for these cameras (GM 4948). *A* video amplifier of one of the cameras. *B* 35-cm picture tube. *M* mixer stage. *R* circuit which automatically keeps the mean brightness of the picture constant at a value which can be adjusted by hand.

It is important that the observer should be able to vary the harmonization continuously and to compare the harmonized picture with the original — as we mentioned above, it was this possibility which induced us to start work on this apparatus. In order to facilitate the control of the harmonization, the control knobs of the mixer stage and of the optical system of the camera producing the blurred image are placed where they can easily be operated by the observer. The contrast amplification is controlled by varying the aperture of the above-mentioned iris diaphragm by means of a Bowden cable; the control knob for this purpose is coupled with a potentiometer which simultaneously alters the luminous output of the lamp so that the total luminous flux which passes through the dia-

phragm is unaltered. Two further measures are taken to enable simultaneous viewing of the original photo (on *H*) and the harmonized picture (on *B*). In the first place, behind the glass plate *H* is placed a second lamp, which can be switched on or off as desired. The light from this lamp cannot fall on the cameras, but can reach the eye of the observer who is able to view the photo from various positions and at various angles (within certain limits). In the second place, the control unit of the picture tube is modified by the addition of an inductive voltage divider so that the 625-line picture does not cover the whole screen of the tube (35 cm diagonal) but only a rectangle of 10×15 cm. It is assumed that this screen will be viewed from a distance of 30-35 cm when the Ampliscope is in use, and by making the image not larger than 10×15 cm we ensure that the line structure does not disturb the viewing of the picture on the tube, and at the same time make the original and the harmonized picture both about the same apparent size. It would of course be even better to replace the 35-cm tube by a smaller one so that the whole screen can be made use of.

We will not discuss the electrical circuitry of this apparatus in detail. As will be apparent from the above, most of the components are standard television equipment, and we think that most of the alterations and additions (mainly to the mixer stage and the control unit) would not be of any particular interest for our readers. The only point we would like to discuss is the regulation of the mean brightness of the harmonized picture.

Every commercially available television receiver has one control knob for varying the "brightness", and another for the "contrast". The first knob alters the bias of the control grid of the picture tube, and thus the brightness of the screen in the absence of a signal (the "black level"). The (objectively) best reproduction of the brightness differences in the original is normally obtained when the black level is made to coincide with the cut-off point of the picture tube's characteristic or to lie slightly above it, and is kept in this position. The "contrast" knob alters the amplification of the video signal applied to the control grid of the picture tube. The contrast is thus not really altered at all, since the transfer function of the circuit is not changed, the gamma remaining approximately equal to 1. This adjustment is however useful in that it allows better use to be made of the available brightness range of the circuit (between the lower limit, which is determined by noise, and the upper, which is determined by overloading of the picture tube) in cases where the brightness range of the original is small.

Neither of these two control methods is very suitable for our purposes. Let us consider the amplification control first: since it is essential to make the best possible use of the brightness range on the screen in order to obtain optimum viewing conditions, it would not be sufficient merely to adjust this by hand now and again. We therefore fitted the Ampliscope with an automatic control of the amplification, which starts with maximum amplification and turns the output signal of the video amplifier down to the permissible value every time the amplitude of the input signal exceeds a certain threshold value. As regards the brightness control, the above-mentioned adjustment of the black level is suitable for normal viewing; but for viewing of a harmonized picture it would seem a better idea to keep the mean density of the picture constant. Therefore we constructed the Ampliscope so that when the video signal of the harmonizing camera is switched on, the manual adjustment of the black level makes way (with a delay of about one second) for automatic control of the mean density. This control circuit contains an RC circuit (time constant 0.13 second) which is used to average the video signal. The difference between this mean value and a constant reference value determines the bias of the control grid of the picture tube. The "brightness" of the picture can still be altered by hand, by changing the above-mentioned reference value.

Another method of double scanning

For the sake of completeness, we would like to devote some space to a description of another method of double scanning, although this was not made use of in the apparatus described here. This method avoids all the problems connected with the above-mentioned stipulation that both cameras should scan precisely the same portion of the original picture; we therefore regard this method as very suitable for further development of the Ampliscope. The set-up is sketched in *fig. 8*. Here the X-ray photo is scanned directly by a flying-spot scanner⁴⁾, whose spot is focused sharply on the photo by a Zoom lens. A selectively reflecting mirror S_1 (which transmits e.g. only green light, and reflects red) is placed in the path of the scanning beam. This reflects the red part of the beam onto the normal double mirror S , which returns it to the original path of the beam via another selectively reflecting mirror S_2 identical with S_1 , as shown in *fig. 8*. The longer path travelled by the red light means that this is no longer in focus when it reaches the photo; the size of the red spot can be easily adjusted by moving the mirror S . The light passes through the X-ray photo (allowing us to profit from the Callier effect) and is concentrated on two photomultipliers by a second lens, a third selective mirror S_3 ensuring that one photomultiplier receives only the red light, and the other only the green. Two video signals are thus produced, just as when two vidicon cameras are used, one obtained by scanning with a sharp spot of light and the other with a blurred spot, and these signals are further processed as described above.

The illumination provided by the scanning light spot will in general be insufficient to allow the observer to view the photo while it is being scanned. A simple solution to this problem is to illuminate the photo with a stroboscope lamp synchronized with the blanking pulses for each scan. Since both photomultipliers are out of action during this pulse, the light from this lamp has no influence on the quality of the harmonized picture.

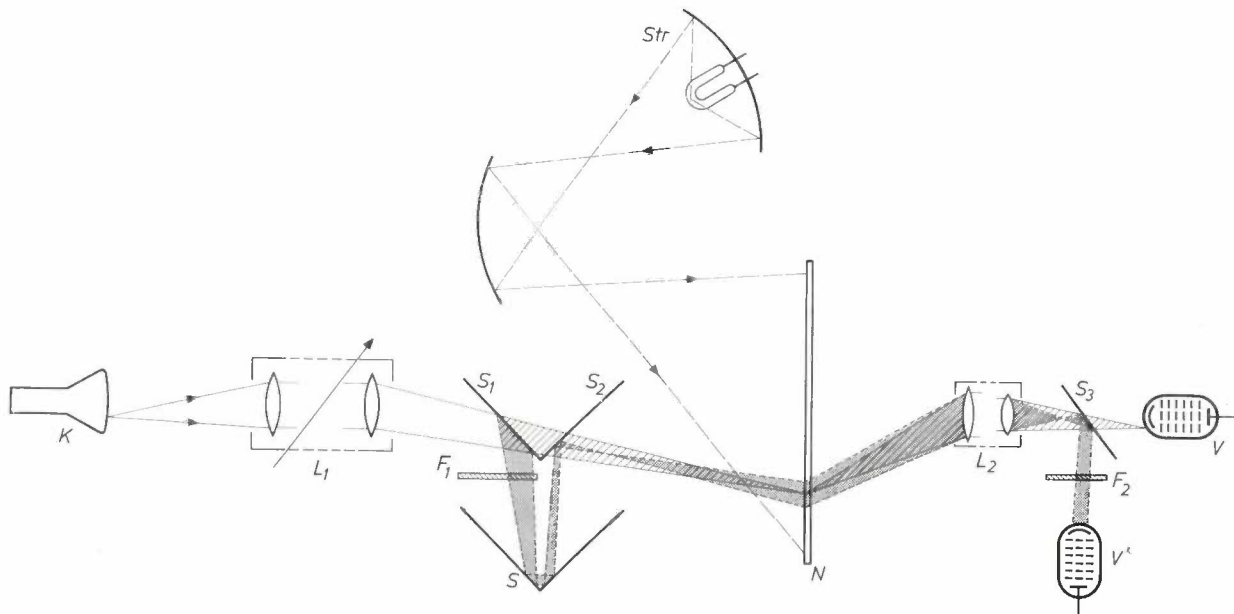


Fig. 8. Double scanning by another method (splitting an optical channel in two). K flying-spot scanning tube. L_1 Zoom lens. N X-ray photo to be inspected. S_1 and S_2 two colour-selective mirrors, which transmit part of the light (e.g. green) and reflect the rest (red). The normal right-angled mirror S returns the red part of the light to the original ray path. (The red light makes a detour.) A sharp green image of the scanning spot and a blurred red image are thus produced on the same point of the X-ray photo N . The transmitted light is concentrated by the lens L_2 onto the two photomultipliers V and V' , the light being again split into two parts by the selective mirror S_3 . If desired, the separation of the colours can be reinforced by the filters F_1 and F_2 . Str lighting system with stroboscopic lamp.

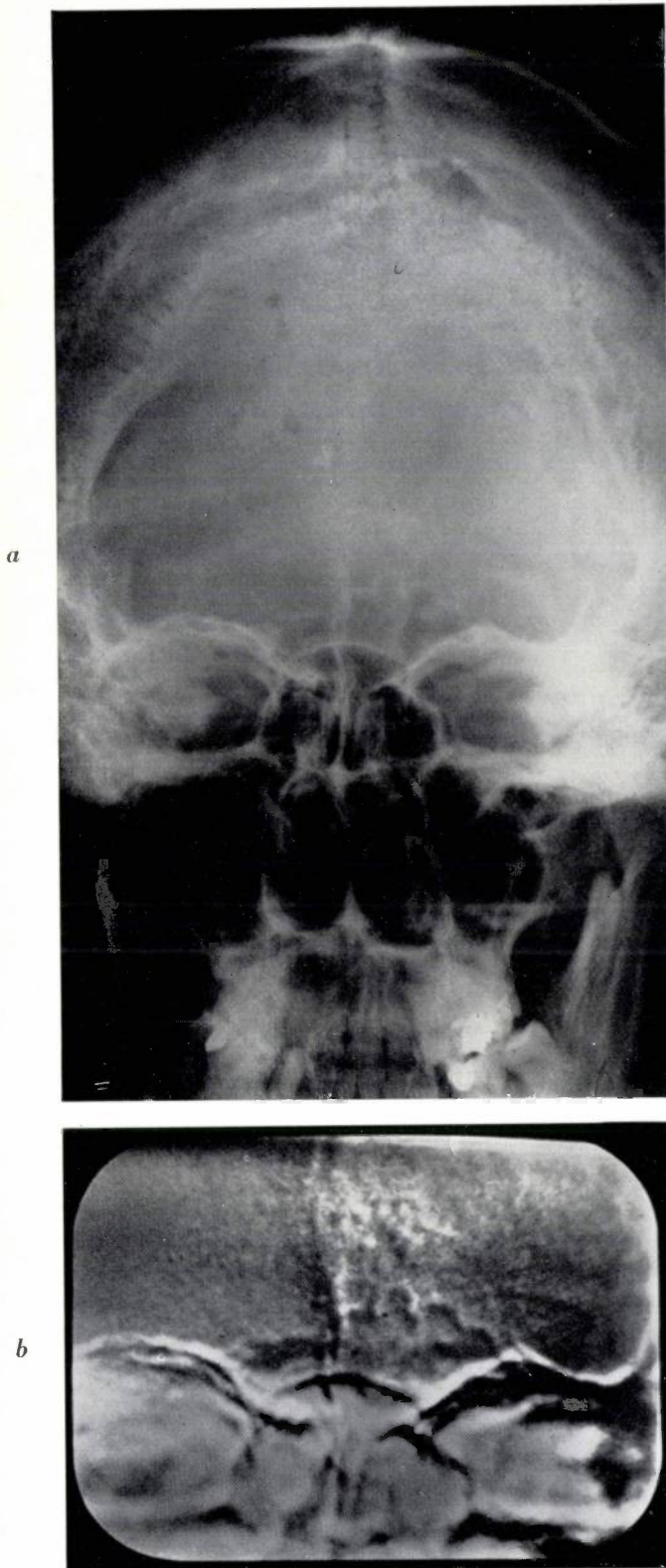


Fig. 9. a) X-ray photo of a skull, observed under normal conditions. b) Photo of the television viewing screen, showing the harmonized version of the middle part of the photo a.

Some results

In order to illustrate the improvement produced by the harmonization, we reproduce here two photos of pictures on the television screen together with

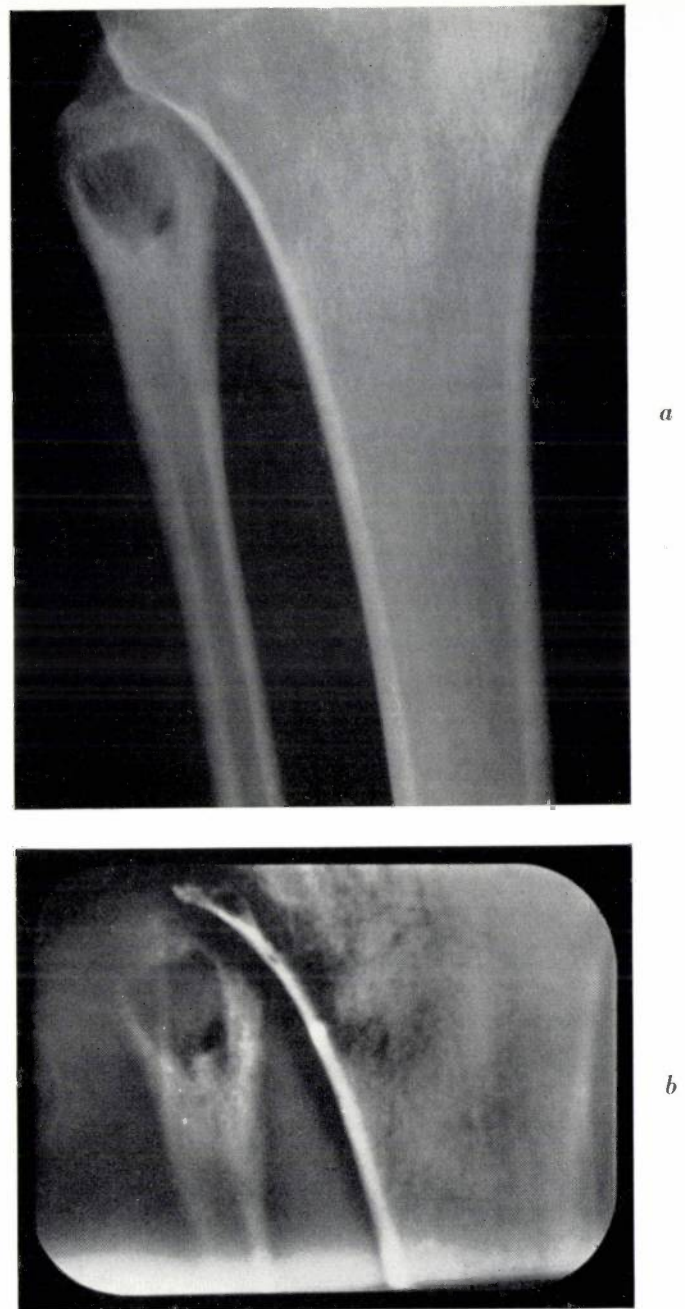


Fig. 10. a) X-ray photo of part of the leg below the knee. b) Photo of the television viewing screen, showing the harmonized version of the upper part of the photo a.

the original X-ray photos. *Fig. 9a* shows a photo of a skull, while *fig. 9b* shows the central part of this photo now completely harmonized (i.e. both signals mixed with equal amplitudes). *Fig. 10a* and *b* are similar pictures of a leg below the knee. In both cases, the effect of the harmonization is very clear.

This does not mean to say that the advantage of this procedure for the radiologist is equally clear. As in many such cases, it must be remembered that the reproductions shown here cannot give an accurate impression of the quality of the pictures. It is known that the range of brightness of a photo reproduced on paper is limited (maximum brightness

ratio about 1 : 10). It is thus impossible to give a faithful impression of the great variations of density in an X-ray negative (density range 1 : 100 or more). The reduction of the coarse contrasts by the harmonization means that the reproduction of a harmonized picture on paper is much closer to its "original" than is the case with a normal X-ray photo. An inexperienced observer looking at these reproductions would therefore tend to over-estimate the improvement produced by harmonization, because full justice is not done to the original X-ray photo.

The judgement of the real diagnostic value of the harmonization must thus be left to the experienced radiologist. So far, the Ampliscope has not been tried out much in practice; but we may refer e.g. to the series of experiments carried out with this apparatus by A. Bonse in the Radiation Institute of the University Dermatology Clinic, Würzburg, where the results were favourable⁹⁾.

As we have just mentioned, harmonized pictures are much more suitable for reproduction on paper, since they bring out the essential information of the original photo, making it no longer necessary to squeeze the whole information of the original into the narrow brightness range of the paper; it may therefore be expected that harmonization will find use for the reproduction of X-ray photos in print¹⁰⁾.

⁹⁾ A. Bonse, not yet published.

¹⁰⁾ This has been pointed out by W. J. Oosterkamp. See *Ärztliche Forschung* 16, I 124, 1962 (No. 3).

The LogEtrón process³⁾ can naturally be used for the same purpose, but the Ampliscope has the advantage, as mentioned above, that the harmonization can be visually checked and adjusted to give the optimum results.

It has already been found that this apparatus can also be of use in other fields than radiology. For example, it proved useful at the Munich Observatory for evaluating photos of the sky taken during an eclipse of the sun¹¹⁾. The pictures, whose contrast was very low because of severe fogging, had to be used for very accurate measurements of the positions of a large number of stars, to allow calculation of the deflection of light in the gravitational field of the sun. In this case, use was not made of the harmonization, but only of the increased contrast produced by the Callier effect.

¹¹⁾ F. Schmeidler, *Naturwiss.* 49, 463, 1962 (No. 20).

Summary. In order to be able to detect low-contrast details on an X-ray photograph, or in some cases to improve the clarity of these details, it may be desirable to reduce or remove the coarse contrasts. In the Ampliscope, this "harmonization" of the picture is produced with the aid of two vidicon cameras. A sharp image of the photo in question is produced on one of these cameras, and a blurred image on the other, and the video signals from the two cameras are mixed in adjustable proportions. By these means, the video modulation in regions which are larger than that covered by the blurred spot is removed or reduced, and this is achieved independently of the direction of scanning the picture. At the same time, the contrast is increased by utilizing the Callier effect. The construction of the apparatus and the way in which the harmonization can be controlled are briefly described in this article, and some results are shown.

ABSTRACTS OF RECENT SCIENTIFIC PUBLICATIONS BY THE STAFF OF N.V. PHILIPS' GLOEILAMPENFABRIEKEN

Reprints of those papers not marked with an asterisk * can be obtained free of charge upon application to the Philips Research Laboratories, Eindhoven, Netherlands, where a limited number of reprints are available for distribution.

R 430: C. Grandjean: Study of germanium by transparency in infrared light; double refraction of silicon and germanium (Philips Res. Repts. 16, 343-355, 1961, No. 4).

The internal elastic stresses in Ge and Si transistors can be traced with the aid of the "photo-elastic" effect. This investigation involves special difficulties in the case of Ge because of the very small value of the photo-elastic constant k (the author measured $k = 0.26 \pm 0.04$ for Ge and $k = 1.97 \pm 0.05$ for Si) and because Ge is transparent only to light of $\lambda > 1.8 \mu\text{m}$. In the set-up described, the detector is a PbS photocell which is sensitive in the wavelength range between 1.8 and 2.8 μm . Since this detector is a point-contact type, the transistor investigated has to be scanned with a light spot; this is done with a Nipkov disc, which has the drawback, however, of a low resolving power. A set-up having a higher resolution is in preparation for measuring very local stresses.

R 431: A. Bril and W. Hoekstra: Efficiencies of phosphors for short-wave ultra-violet excitation (Philips Res. Repts. 16, 356-370, 1961, No. 4).

As a measure of the fluorescence of phosphors one can use either the radiation efficiency (ratio between radiated and absorbed energy) or the quantum efficiency (ratio between the number of radiated and absorbed quanta). These efficiencies have been measured by numerous investigators, but their results differ fairly widely, probably because indirect methods are frequently used. Following earlier investigations, when an absolute and direct method was developed for determining efficiencies with cathode-ray excitation (see Philips tech. Rev. 15, 63-72, 1953/54), the authors now describe a method for short-wave ultra-violet excitation. Use is made of a thermopile, whose sensitivity is independent of wavelength. The article contains graphs and tables showing the measured efficiencies and the spectral distributions of some standard phosphors issued by the National Bureau of Standards, Washington. The efficiencies of other phosphors can be determined by a relative method, also described in this article, in which a comparison is made with a standard phosphor (e.g. MgWO_4).

R 432: W. G. Gelling and J. H. Haanstra: A red electroluminescent ZnSe phosphor (Philips Res. Repts. 16, 371-375, 1961, No. 4).

Description of the preparation of a red electroluminescent phosphor consisting of ZnSe activated with Cu and Al. The luminescence is highly temperature-dependent; with increasing temperature a sudden drop occurs at a point below room temperature. The temperature at which this drop begins can be raised to room temperature by increasing the frequency of the exciting voltage. At 1500 c/s the quantum efficiency of the red emission from ZnSe (Cu, Al) equals that of the green emission obtained under identical circumstances from the standard phosphor ZnS (Cu, Al).

R 433: P. A. H. Hart and C. Weber: Parametric-amplifier electron gun (Philips Res. Repts. 16, 376-388, 1961, No. 4).

Two types of electron gun were built and tested; one a Brillouin type using a magnetically shielded cathode, and the other an immersed-flow type, both of them suitable for use in parametric amplifier tubes having a very low noise figure. For this application the guns are required to have a high perveance ($3 \times 10^{-6} \text{ AV}^{-3/2}$; the perveance is defined as the product of the beam current and the beam voltage to the power of $-3/2$) and also a low beam voltage (6 V). The dimensions of guns built by the conventional methods would be too small for this purpose. The guns described therefore consist of a part with low perveance and high beam voltage, behind which a lens system reduces the velocity of the electrons to that corresponding to the required beam voltage, without impairing the focusing, and at the same time raises the perveance to the value required. For a given beam current, beam voltage and magnetic field strength the diameter of the beam is made as small as possible, resulting in a high space-charge density. This can be confirmed by incorporating the guns in experimental Adler tubes and then measuring the following quantities: the gain at which the pump field is defocused (the saturation gain), and the frequency at which the electrons rotate around an axis parallel to the beam axis, and from which the

beam diameter can be derived. The measured values are in good agreement with those found from the theory.

The minimum noise figure obtained with the Brillouin gun is 1.31 (corresponding to 90 °K) at a maximum gain of 30 dB, and the saturation gain is 50 dB. The corresponding quantities obtained with the immersed-flow gun are 1.4 (116 °K) at 22 dB, and 30 dB, respectively. The Brillouin gun thus gives the best results.

R 434: C. A. A. J. Greebe and W. F. Knippenberg: Grown *P-N* junctions in silicon carbide, II (Philips Res. Repts. **16**, 389-398, 1961, No. 4).

Further investigation of the properties of grown *P-N* junctions in silicon-carbide crystals (see **R 392**). The current density in the junction was found to be inhomogeneous in many cases. For photon energies greater than the band gap the spectral distribution of the *P-N* luminescence, under forward bias, can be brought into agreement with the theory of Van Roosbroeck and Shockley. The article concludes with a description and analysis of a forward current and voltage characteristic, in which several current components can be distinguished.

R 435: O. Reifenschweiler: Sealed-off neutron tube: the underlying research work (Philips Res. Repts. **16**, 401-418, 1961, No. 5).

An account of the investigations that have led to the construction of a simple, compact and transportable neutron source. The definitive version of this sealed-off neutron tube has been dealt with extensively in Philips tech. Rev. **23**, 325-337, 1961/62 (No. 11).

R 436: P. Penning and D. Polder: Anomalous transmission of X-rays in elastically deformed crystals (Philips Res. Repts. **16**, 419-440, 1961, No. 5).

The anomalous transmission of X-rays (very small effective absorption of beams incident on a perfect crystal set to reflect according to Bragg's law) can be explained from the dynamical theory of X-ray diffraction, which is briefly summarized. In *deformed* crystals the diffraction can also be treated dynamically, if the variation of the lattice constant is gradual. This generalized dynamical theory, which forms the bulk of this article is, analogous to the theory of the propagation of light in a medium with a gradually changing refractive index. According to the theory the path of a narrow X-ray beam in a

deformed crystal is curved. For two different deformation patterns (caused respectively by a uniform temperature gradient and by elastic bending) values are calculated of the intensities of the transmitted and the reflected beam that show good qualitative agreement with the experimental data available.

R 437: F. K. Lotgering, U. Enz and J. Smit: Influence of Co^{2+} ions on the magnetic anisotropy of ferrimagnetic oxides having hexagonal crystal structures (Philips Res. Repts. **16**, 441-454, 1961, No. 5).

The Co^{2+} ions in the hexagonal ferrimagnetic oxides $\text{BaCo}_\delta\text{Ti}_\delta\text{Fe}_{12-2\delta}\text{O}_{19}$, $\text{BaCo}_\delta\text{Zn}_{2-\delta}\text{Fe}_{16}\text{O}_{27}$ and $\text{Ba}_3\text{Co}_\delta\text{Zn}_{2-\delta}\text{Fe}_{24}\text{O}_{41}$ may be the cause both of stable and metastable preferred directions of the magnetization vector. These directions can make random angles with the *c* axis. This was found by measuring at low temperature the torque exerted on crystal-oriented samples when rotated in a strong magnetic field. In some cases a special form of hysteresis occurred. The marked influence of Co^{2+} ions is attributed to non-compensated orbital magnetism.

R 438: G. Bosch: On the thermal conductivity of SiC (Philips Res. Repts. **16**, 455-461, 1961, No. 5).

The study of the thermal conductivity of solids has become an important tool in recent years for research into the nature and concentration of lattice imperfections. In this context the author has measured the thermal conductivity *K* at low temperatures *T* of hexagonal (*α*) and cubic (*β*) silicon carbide. The *K* of *β*-SiC is proportional to T^3 , as it is in the majority of solids. The *K* of *α*-SiC, however, is shown to be proportional to T^2 , while the proportionality constant found is so small that it cannot be explained with the existing theories. Similar unexplained phenomena have been reported in the literature.

R 439: G. H. Plantinga: The noise temperature of a plasma (Philips Res. Repts. **16**, 462-468, 1961, No. 5).

Nyquist's theorem cannot be applied to a plasma which is not in thermodynamic equilibrium. It can, however, be used for defining the noise temperature. The author derives from this definition a general expression for the noise temperature by calculating the noise current in the plasma. The expression found holds for an arbitrary isotropic distribution of electron velocities and for a collision frequency which may be a random function of the velocity of the colliding electrons. When the electrons have a

Maxwellian velocity distribution, or when the collision frequency of the electrons is independent of their velocity, the noise temperature is shown to be equal to the temperature of the electrons. Bekefi, Hirschfeld and Brown have derived an expression for the radiation temperature of a plasma, defined with the aid of Kirchhoff's law, which is in principle identical with the expression found here. The latter, however, owing to the method of calculation employed, is applicable under more general assumptions, and can be applied, unlike the other expression, to non-transparent plasmas.

R 440: K. Teer: Investigation of the magnetic recording process with step functions (Philips Res. Repts. 16, 469-491, No. 5).

Most fundamental investigations of magnetic recording have been based on the use of sinusoidal test signals superimposed on a high-frequency bias field. This is an obvious method, well matched to sound recording, which has always been the foremost application of magnetic-recording facilities. With a view to acquiring further knowledge, and partly in connection with new applications in which pulsed signals have to be recorded, an investigation has been made into the unit-step response of magnetic recording systems. For this purpose a step-function was used as a test signal. A theoretical treatment is followed by a discussion of the experiments. The step-function response was measured under various conditions, special attention being paid to the response at very low recording levels and at levels near tape saturation. A comparison of theory and experiment leads in the first case to a conclusion regarding head quality, and in the second case as regards tape quality. Finally results are given of measurements of pulse-correction networks, which may considerably increase the resolving power of the magnetic recording system in practical pulse-recording.

R 441: W. G. Gelling: Switching elements consisting of photoconductors and triggered neon lamps (Philips Res. Repts. 16, 501-506, 1961, No. 6).

A photoconductor combined with a light source can be used as a switch; when the light source is ignited the photoconductor passes current, closing the circuit in which it is incorporated. This article discusses the conditions which a neon lamp must fulfil in order to be used in conjunction with a CdSe photoconductor.

R 442: T. J. Viersma: Investigations into the accuracy of hydraulic servomotors (Philips Res. Repts. 16, 507-597, 1961, No. 6).

This article forms the first part of a thesis (Delft, April 1961; continued in R 444) which deals with the research that has resulted in the design of linear hydraulic servomotors, with very high control accuracy, used for example for controlling automatic machine tools. The first chapter presents a phenomenological discussion of the control error. In quasi steady-state conditions the control error contains three components, proportional respectively to the ram speed, the external load on the ram and the Coulomb or dry friction. The last component gives rise to the highly troublesome dead zone. To minimize the control error, a small velocity-time constant and a large hydraulic rigidity are required. In chapter 2 the oil flow through the ports in the regulating spool is analysed. It is shown that a turbulent flow is needed to obtain a small velocity-time constant. A very large hydraulic rigidity is made possible by applying load compensation. As a result the dead zone can be reduced from the normal value of 15-20 microns to between 0.1 and 0.3 micron.

For analysing the dynamic behaviour, chapter 3 sets up the differential equation for three types of hydraulic servomotor. The difference between the three types can be expressed in terms of various constants occurring in the equation. The derivation of these constants takes into account the compressibility of the oil, the dry friction and the viscous friction. If the dry friction can be disregarded, the differential equation is almost linear in a wide range. After partial linearization of the general differential equation, the servomotor behaviour and the "regulator" behaviour are analysed, using the describing-function method and an electrical analogue. It is demonstrated that the stability conditions are identical for hydraulic servomotors with and without dry friction. Viscous friction is shown to be essential for obtaining the desired damping, and acceleration feedback is also promising in this respect. Apart from the effect of the dynamic dead zone, the remarkable influence of dry friction at low frequencies calls for attention. The results obtained graphically and analytically, using the describing-function method, are confirmed in broad lines by experiments done with the electrical analogue.

The last chapter reports on the use of a hydraulic servomotor in a numerically controlled cam-milling machine, where the favourable properties of these servomotors are used to full advantage. The velocity-time constant amounted to 4 milliseconds,

the resonance-time constant to 1 millisecond. The dead zone was negligible. During the cutting process on the milling machine the control error was 2 microns.

See also the next issue of Philips tech. Rev. (Vol. 24, 320-331, 1962/63, No. 10).

A 40: A. Klopfer: Die Erzeugung von Höchstvakua mit Getter-Ionenpumpen und das Messen von sehr tiefen Drucken (Vakuum-Technik 10, 113-118, 1961, No. 4). (The production of ultra-high vacua with getter-ion pumps, and the measurement of extremely low pressures; in German.)

A description of experiments which show that a getter-ion pump, using a gas discharge with cold cathode, is well suited for producing extremely low pressures in small laboratory vacuum equipments and maintaining them over a long period, without the aid of a diffusion pump. The lowest total pressure hitherto reached was 6×10^{-12} torr. The residual gas, which was analysed with an omegatron, mainly consisted of hydrogen and nitrogen. The set-up used and the pumping method are described. The pressure variation as a function of time and the increase in pressure when the pump is switched off are discussed, and methods of measuring such low pressures are considered.

A 41: B. Lersmacher and S. Scholz: Drucksintern von Hafnium-, Zirkon- und Tantalkarbid ohne Bindephase (Arch. Eisenhüttenw. 32, 421-429, 1961, No. 6). (Pressure-sintering of hafnium, zirconium and tantalum carbides without binding phase; in German.)

Sintered carbides of metals such as hafnium, zirconium or tantalum are widely used in engineering as hard metals. They usually contain about 3 per cent of a metal which serves as a binder. The binder sets a limit of roughly 800 °C to the useful temperature range. By sintering under high pressure it is possible, without using binders, to make hard metals which can be employed up to substantially higher temperatures. A systematic investigation has shown that there is an optimum temperature for the high-pressure sintering process and that the addition of about 1% of certain metals (notably Mn, Co or Ni) accelerates the sintering process. The article examines the influence of pre-heating in vacuum, the duration of grinding in a ball mill using agate balls, and the treatment of the sintering powder with a solution of Mn in hydrochloric or sulphuric acid. The effects observed are discussed in relation to theoretical considerations.

A 42: S. Garbe: Zur Wasserabgabe von Natrium-silikatgläsern (Glastechn. Ber. 34, 413-417, 1961, No. 8). (On the desorption of water from soda-lime glass; in German.)

After the degassing of glass samples in vacuum and measurement with an omegatron of the partial pressures of the released gases, the author determined the quantity of water which soda-lime glasses, containing 7 to 40 mol % Na_2O , give off in the temperature range from 500 to 900 °C. The water dissolved in the glass shows a minimum at 20 mol % Na_2O . The diffusion coefficients of water in the glass samples were derived from the measurements. At 900 °C these coefficients increase markedly with the sodium content, whereas at 550 °C they show minima at 20 mol % Na_2O . The activation energy for the water desorption is shown to be dependent on the sodium content. The results of the investigation are explained on the hypothesis that "water" is dissolved in alkaline glass in two different ways, i.e. (1) in the form of free Si-OH groups, and (2) in the form of hydroxyl groups bound to Si, between which a hydrogen bridge exists.

A 43: P. Gerthsen: Thermoelektrische Anwendung von Halbleitern (Z. angew. Phys. 13, 435-444, 1961, No. 9). (Thermoelectric application of semiconductors; in German.)

The phenomenological theory of thermoelectric effects is treated in so far as it is necessary to understand the subsequent considerations regarding the suitability of materials for technical applications, such as thermo-generators and heat pumps. The author presents the theory of the thermo-e.m.f. in metals and semiconductors in a plausible form. Some special semiconductors are discussed and their application elucidated with examples.

A 44: A. Klopfer: Das Erreichen und Messen von tiefen Drucken (Z. angew. Phys. 13, 480-491, 1961, No. 10). (The production and measurement of low pressures; in German.)

Survey of the methods developed in the last ten years for producing and measuring very low pressures (10^{-12} torr). A limit is set to the minimum pressure not by the pump but by the gas desorption from the walls and components in the vacuum system. After dealing with methods of limiting the gas desorption, the author discusses the various types of vacuum pumps and pressure gauges for extremely low pressures.

A 45: H. G. Reik and H. Risken: Distribution functions for hot electrons in many-valley semiconductors (Phys. Rev. **124**, 777-784, 1961, No. 3).

To give a detailed explanation of the relation between current and voltage in strong electrical fields in semiconductors, e.g. type *N* germanium, a simple parabolic model of the band structure is not sufficient. A more complicated model, called the many-valley structure, has to be used. This article is a contribution to the theory of electron transport, based on the application of Boltzmann's equation to electrons in many-valley semiconductors, in which electrons are scattered in the same valley and between different valleys.

A 46: R. Groth and E. Kauer: Absorption freier Ladungsträger in α -SiC-Kristallen (Physica status solidi **1**, 445-450, 1961, No. 5). (Absorption of free charge carriers in α -SiC crystals; in German.)

The authors have studied the radiation absorption by free charge carriers in SiC crystals (with the common 6 H structure) in the region of wavelengths from 1 to 4 μm at temperatures *T* from 20 to 1500 °C. The logarithm of the absorption coefficient plotted versus $1/T$ gives curves that can be approximated by two straight lines, from which the height of the acceptor levels above the valence band and the distance between the conduction band and the valence band can be derived. In *N*-type crystals the absorption at low temperatures is governed by scattering due to ionized impurities, and at high temperatures by thermal scattering. The manner in which absorption depends on wavelength is in good agreement with the theory. The shift of the absorption edge with temperature was investigated up to 1500 °C.

A 47: R. Groth and R. Memming: Absorption freier Ladungsträger in CdS (Physica status solidi **1**, 650-655, 1961, No. 6). (Absorption of free charge carriers in CdS; in German.)

The absorption of radiation by free charge carriers in *N*-type CdS crystals, in which the concentration *n* of the electrons is several times 10^{17} per cm^3 , was investigated at room temperature and at 90 °K. At 90 °K the absorption coefficient is proportional to λ^3 , as predicted by H. J. G. Meyer's theory for the case of scattering by ionized impurities. At room temperature the scattering is already partly ther-

mal. It was derived from the measurements that the effective mass of the free charge carriers is 0.19 (± 0.02) times the mass of a free electron. It is further demonstrated that when calculating *n* from the formula $R_H = -r/ne$ (R_H = Hall constant, *e* = charge of electron) it is necessary to bear in mind that the factor *r* depends on *n* when the scattering is due to impurities. The relation between *r* and *n* is derived.

A 48: H. G. Grimmeiss, A. Rabenau, H. Hahn and P. Ness: Über elektrische und optische Eigenschaften einiger Chalkogenide von Elementen der IV. Nebengruppe (Z. Elektrochem., Ber. Bunsenges. physik. Chem. **65**, 776-783, 1961, No. 9). (On the electrical and optical properties of some chalcogenides of elements of the IVth sub-group; in German.)

Electrical and optical investigations were made on single crystals of compounds formed from one of the elements Ti, Zr or Hf (these elements all belonging to the sub-group in column IV of the periodic table) and from one of the elements S, Se or Te (elements belonging to the chalcogenides in column VI of the periodic table). The preparation of the crystals is discussed. With the exception of TiS, which proved to be a metallic conductor, all compounds investigated are semiconductors. The compounds are divided into two groups. (1) The trichalcogenides TiS₃, ZrS₃, ZrSe₃ and HfS₃, all of which have the same crystal structure, yielded the most complete information, and in all cases the band gap was determined. (2) The compounds TiS, TiSe, TiTe and Ti₃S₄, having a crystal structure identical with or closely related to that of NiAs, and the compounds TiS₂, TiSe₂ and TiTe₂, having the crystal structure of CdI₂.

A 49: B. Lersmacher, E. Roeder and S. Scholz: Über das Kornwachstum von TaC und NbC unter dem Einfluss geringer Zusätze von Mn, Fe, Co und Ni (Naturwiss. **49**, 35, 1962, No. 2). (On the grain growth of TaC and NbC under the influence of small additions of Mn, Fe, Co and Ni; in German.)

TaC and NbC are interesting because of the properties that make them suitable as a material for filaments in incandescent lamps (among other things they possess a very high melting point). The article describes an investigation into the sintering of these substances with the addition of small quantities of Mn, Fe, Co or Ni.

Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES

A NUMERICALLY CONTROLLED CONTOUR MILLING MACHINE

- I. GENERAL DESCRIPTION
- II. THE COMPUTING UNIT
- III. THE HYDRAULIC SERVOMOTOR

In recent years, increasing use has been made of numerical control in the automation of machine tools. Compared to conventional methods of control this offers many advantages, among which may be mentioned the possibility of obtaining great accuracy and great flexibility. In order to gain experience with such a system, a simple numerically controlled contour milling machine has been built by Philips for internal use. This machine is very accurate: the position of the cutter is controlled to within 2 μm . A general description of the machine is given in Part I of this article, while in Part II an important part of the control system (the computing unit) is discussed. Part III gives some general considerations on hydraulic servomotors, followed by a description of the servomotor used in the present case, which embodies some new principles.

I. GENERAL DESCRIPTION

by J. A. HARINGX *).

621.914.3-522

A few years ago, work started in the Philips Research Laboratories, Eindhoven, on systems for the numerical control of machine tools. These systems are very much in the limelight at present, and particularly in America a whole series of numerically controlled tools have been put on the market. The investigation in Eindhoven has laid the main stress on the attainment of as high an accuracy as possible. One of its results has been the design of an electronic-hydraulic control system into which the data on the shape of the workpiece are fed on punched tape.

In choosing the kind of machine (lathe, milling machine, etc) on which the control system was to be tested in practice, we were guided by an existing need in one of the Philips factories, namely for a *contour milling machine* with an accuracy greater than is normal for such machines. The milling machine constructed for this purpose (*fig. 1*), which was finished about two years ago and is known as

the PRO-PHIL, will be described in this article. We shall start by discussing an important principle used in this machine.

The principle of "derived control"

In a contour milling machine, the cutter must move in such a way with respect to the workpiece that the latter is given the desired shape. In the case of a three-dimensional object, the shape may be given e.g. by an equation of the form

$$F(x,y,z) = 0. \quad \dots \dots (1)$$

With a milling machine where the cutter is moved by a combination of three mutually perpendicular slides, this movement is normally controlled as follows. One of the slides, e.g. that which varies the y coordinate, is periodically shifted a short, fixed distance. While this slide is stationary, the cutter must be moved in the correct manner over the workpiece by the other two slides. A simple way of doing this is by means of what we call

*) Philips Research Laboratories, Eindhoven.

“derived control”. With a constant y coordinate, the cutter must move in the z - x plane along a curve given by the function

$$z = f(x), \dots \dots \dots (2)$$

derived from eq. (1). The slide which moves the cutter in the x -direction is now driven by a motor. During this motion, its displacement is measured,

blades. The contour milling machine to be described here, in which derived control is used, is designed for making surface cams.

The milling machine for the production of surface cams

These surface cams are used to control a winding machine on which grids for electronic valves are

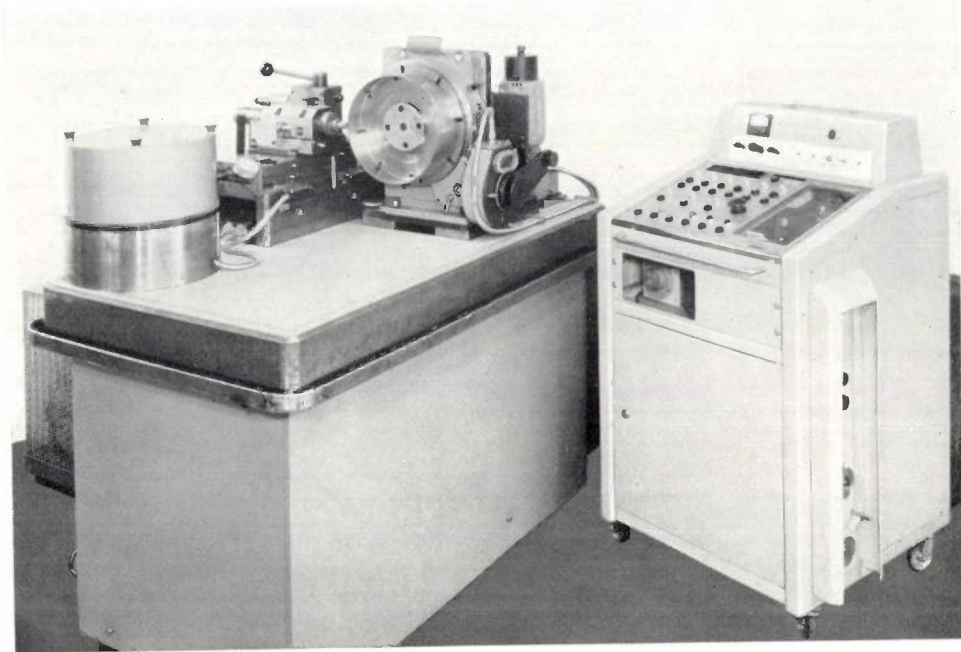


Fig. 1. The numerically controlled contour milling machine.

and a computing unit calculates the necessary displacement of the z slide from these measurements and from data on the form of $f(x)$. The latter are fed separately to the computing unit, e.g. on punched tape. A control system then ensures that the z slide does move in the calculated way. After the y coordinate has been changed, the whole process is repeated with a new function $z = f(x)$ derived from eq. (1).

Derived control thus only needs one control system. The x slide need not move at a constant velocity; it is sufficient to stipulate that it moves continuously in one direction. One disadvantage is that this stipulation imposes a restriction on the shape of the workpiece. Since the x slide may not stop or move backwards, and since the speed at which the cutter can be moved in the z direction is limited, the slope dz/dx of the contour curve can never exceed a certain maximum value. There is however an important group of objects with curved surfaces for which this method is suitable; we may mention flat cams, surface cams, ships' propellers, turbine

made (see figs 2 and 3). The upper edge of this cam must be so formed that the roller which rests on the cam when it is in the winding machine moves in the

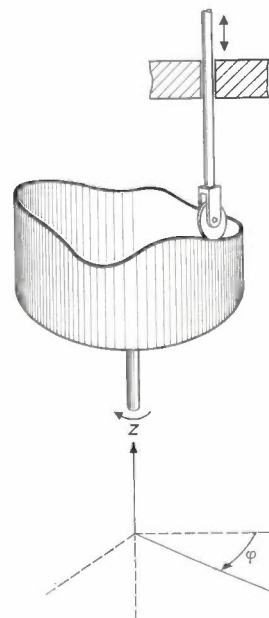


Fig. 2. Sketch of a surface cam, as used in machines for winding the grids of electron tubes. On the raised edge of the surface cam rests a roller which can move in the z direction. The cam rotates, which gives the roller a motion determined by the shape of the cam. This motion in the z direction is used for the control. Each type of grid is given by a certain relationship between the motion of the roller and the angle of rotation φ of the cam, so that a special cam must be made for each type.



Fig. 3. Some surface cams made with the numerically controlled milling machine. The upper one is similar in shape to those used for grid-winding machines, while the two others were made to test the accuracy of the numerical controlled milling machine. The edge which was cut away during milling has been placed above the right-hand one.

desired manner in the z direction as a function of the angle of rotation φ of the cam. In order to give the edge this shape, the motion of the cutter during milling must be given by a function

$$Z = f(\varphi). \dots \dots \dots (3)$$

(In fact, this function is given in cylindrical coordinates; for a given cam, the radius r is constant.)

During the milling, therefore, both a translational and a rotational motion must occur. If derived control is used, only one of these motions need be controlled. We chose to control the translational motion, which is done as follows (fig. 4).

A cutter head is mounted on a horizontal slide S , at right angles to the direction of motion. The cylinder from which the cam is to be milled is fixed on a vertical turntable D whose axis of rotation is parallel to the direction of motion of the slide. This turntable is turned, in one direction only, by an electric motor. The slide is moved by a hydraulic servomotor H which is controlled by a servo-valve whose setting is automatically determined by an electromagnet. During the rotation of the turntable, the current through the electromagnet is varied so that the slide moves in the way needed to give the cam the desired shape. The electronic system which controls this current is a digital system, i.e. it works with quantities which are integral multiples of a given unit; digital systems have been found to be most

accurate. The data on the shape of the cam is fed into the system in numerical form, on punched tape. This punched tape is previously prepared by a digital electronic computer from details of the cam; we used the PASCAL computer ¹⁾ for this purpose. The digital control system makes use of feedback: the displacement of the slide is measured (also in digital form), and compared with the desired displacement.

The motion of the cutter during milling (given by eq. 3) is the same as that of the roller in the winding machine, if the cutter and the roller both have the same diameter. If the diameters are not the same

then the track of the cutter must be calculated; this calculation is also carried out by the PASCAL.

As mentioned above during the discussion of derived control, this method cannot be used to give the cam any desired shape, owing to the one-way rotation of the turntable. A further restriction is added by the finite diameter of the cutter. However, these restrictions did not give any trouble in the production of the above-mentioned surface cams.

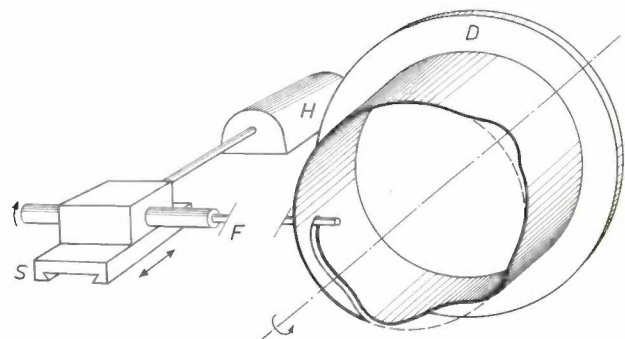


Fig. 4. Perspective sketch of the most important parts of the milling machine. The workpiece is fixed to the vertical turntable D . The cutter head with the cutter F is mounted on the slide S , which is moved horizontally by the hydraulic servomotor H . During the rotation of the turntable, the servomotor is controlled so that the cutter describes the desired path through the workpiece.

¹⁾ W. Nijenhuis, The PASCAL, a fast digital electronic computer for the Philips Computing Centre, Philips tech. Rev. 23, 1-18, 1961/62 (No. 1).

Survey of the most important components

For the purpose of describing the various components of the milling machine, we shall divide the machine into three main parts, viz (fig. 5) the turntable *D* on which the workpiece is fixed, the computing unit *C* which calculates the necessary motion of the slide during milling, and the "step" motor *SM* which brings this motion about.

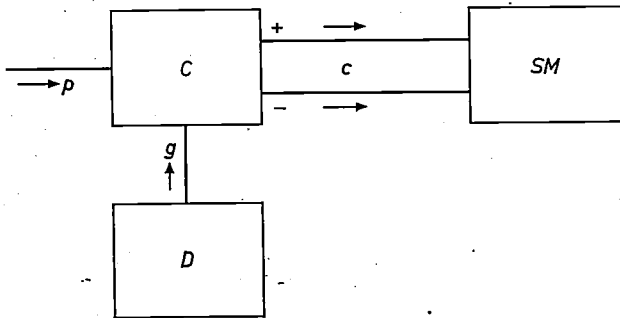


Fig. 5. Block diagram of the milling machine. The computing unit *C* receives the sync. pulses which indicate the rotation of the turntable *D*. The data on the desired shape of the cam are fed into *C* on punched tape *p*. The computing unit delivers command pulses to the step motor *SM*.

The computing unit receives a constant flow of information about the rotation of the turntable. The rotation is measured by an optical system, which delivers a voltage pulse (a "sync." pulse *g*) each time the turntable rotates through an angle of $1/2048$ degree. The displacement that must be given to the slide per sync. pulse (which naturally depends on how far the rotation has proceeded) is calculated by the computing unit from the data on the shape of the cam provided by the punched tape. The calculated displacement is represented by voltage pulses (the "command pulses" *c*), which appear at one of two outputs: a pulse on the + output means that the slide must be moved forward by one length unit (in the present machine $1\frac{1}{4} \mu\text{m}$), while a pulse on the - output means that it should be moved back one unit. From now on, we shall refer to these pulses as "positive" and "negative" command pulses. (This does not refer to their polarity, but to the direction of motion of the slide.) No more than one command pulse can be given per sync. pulse, so that the slide can only be displaced one length unit per unit of rotation (angle unit). The greatest slope which can be cut is thus $2048 \times 1\frac{1}{4} \mu\text{m}$ ($= 2.56 \text{ mm}$) per degree.

The step motor forms a closed loop linear control circuit, composed of a number of electronic and hydraulic units. This control circuit works as a follow-up system which produces the displacements designated by the command pulses. It is called a "step motor" because, as long as the frequency of the incoming command pulses is low, the displacements

are produced step by step. Although at higher frequencies the slide attains a continuous velocity, so that the designation "step motor" no longer applies, this term will be used throughout the rest of the article.

Fig. 6 shows a diagram of the milling machine, in which some components of the step motor are given in greater detail. The most important part is the hydraulic servomotor *H*, which consists of a cylinder in which a piston *Z* can move; the slide is fixed firmly to this piston. Oil is pumped through the cylinder, from the left-hand part to the right-hand part via a small hole in the piston. The amount of oil coming from the right-hand part of the cylinder can be varied by means of the electromagnetically operated servo-valve *R* placed in the oil outlet; the piston is thus given a certain velocity. The servomotor is designed so that the velocity of the piston is proportional to the current through the electromagnet.

The displacement of the slide, like the rotation of the turntable, is measured with an optical system (*M*). This gives voltage pulses at two outputs: one at the + output when the slide moves $1\frac{1}{4} \mu\text{m}$ forwards and one at the - output when the slide moves the same distance backwards. These pulses are used for feedback purposes, and will therefore be called "feedback pulses" (*t*). As with command pulses, we shall distinguish them as "positive" and "negative" feedback pulses.

The step motor also contains a number of electric circuits, viz, a difference register (*DR*), a resistance network (*W*) and a DC amplifier (*A*). The latter provides the current for the electromagnet of the servo-valve. The difference register receives the command pulses and the feedback pulses. The current supplied by the amplifier *A* is determined by the difference between the number of command pulses and feedback pulses received, i.e. by the difference between the desired position of the slide and its actual position. The servo-valve is designed so that the piston is at rest when the difference is zero. A positive or negative difference causes the piston to move until the difference is eliminated. We shall discuss the operation of these circuits in greater detail below.

The computing unit

The operation and construction of the computing unit will be dealt with in part II of this article; we shall therefore give here only a brief statement of its operating principle. The simplest way of producing the command pulses would be to indicate on the punched tape fed to the computing unit how the

z coordinate was to be changed (+1 unit, 0 or -1 unit) for each sync.pulse. This would however use up a lot of punched tape, since one revolution of the turntable corresponds to $2048 \times 360 = 737280$ angle units. However a much smaller number of data are sufficient, since for the objects made on this

pulses t (the measured displacement Δx_0). This register consists of five flip-flops, which form a binary counter. This counter can count both forwards and backwards, by means of voltage pulses. The smallest binary number it can indicate is 00001, and the largest is 11111. The number in the middle

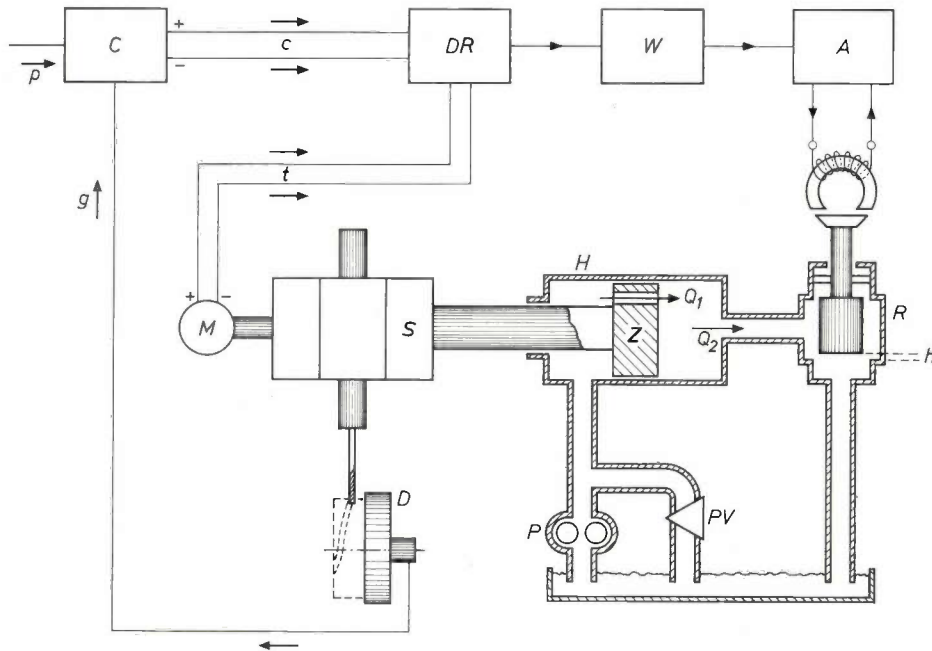


Fig. 6. Diagram of the milling machine, showing some parts of the step motor in greater detail than in fig. 5. C is again the computing unit, which receives both the punched tape p and the sync. pulses from the turntable D . DR is a difference register, which receives the command pulses c from C and also the feedback pulses t from a measuring system M which measures the displacement of the turntable S . The pulses c indicate the desired displacement Δx_i of the slide, and the pulses t the actual displacement Δx_0 . The difference register determines the difference $\Delta x_i - \Delta x_0$. A resistance network W delivers a voltage proportional to this difference. This voltage acts as input signal for a DC amplifier A , which supplies the current for the electromagnet of the servo-valve R which controls the hydraulic servomotor H . P oil pump. PV pressure valve.

milling machine the function $Z = f(\varphi)$ has a simple mathematical form over a certain distance. The displacement is therefore given on the punched tape for intervals of many angle units, and not for every angle unit. The computing unit contains an interpolator, which approximates to the function by first-degree or second-degree polynomials between the values given. The length of the intervals and the method of interpolation are chosen so that the deviation from the proper value is never more than half a length unit. In this way, the length of the interval can be made several hundred angle units.

The step motor

Fig. 7 shows a block diagram of the step motor. The difference register DR forms the difference between the number of incoming command pulses c (the desired slide displacement Δx_i) and feedback

of this range (i.e. 10000) is taken as the zero point of the register, the smallest and largest numbers thus corresponding to the decimal numbers -16 and +15.

The command and feedback pulses are fed to the difference register in the following way. The positive command pulses make the register count forwards, and the negative ones backwards. The feedback pulses work the other way round, i.e. the positive ones make the register count backwards and the negative ones forwards.

The resistance network W is connected to the register, and produces from the voltages of the

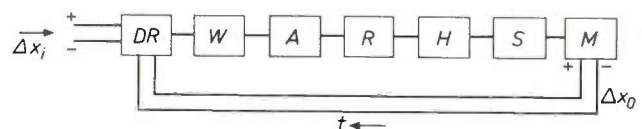


Fig. 7. Block diagram of the step motor. The letters have the same significance as in fig. 6.

flip-flops a voltage V varying from 0 to 2 volt and proportional to $\Delta x_1 - \Delta x_0$. This voltage thus increases discontinuously with the number in the difference register (fig. 8). Since the voltage V acts as the input signal of the amplifier A which gives the activation current for the servo-valve, the ordinate of fig. 8 can also be taken as the velocity v_s of the slide. When the register indicates zero, the slide is at rest, when it indicates $+1$ the slide moves forwards and when it indicates -1 the slide moves backwards.

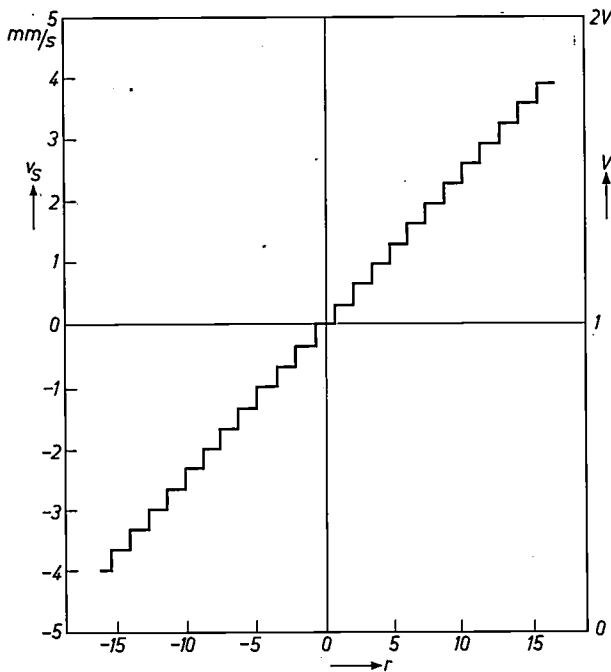


Fig. 8. Graph showing the relationship between the output voltage V of the resistance network and the position r of the difference register (r is proportional to $\Delta x_1 - \Delta x_0$). Since the velocity v_s of the slide is proportional to V , a scale for v_s is also given along the ordinate axis.

We shall now start by considering the state in which the difference register indicates 0 and the slide is at rest. (If the slide should happen to be moving when the difference register indicates zero, it can be stopped by changing the control current; i.e. the curve of fig. 8 can be displaced vertically.) If the difference register now receives a positive command pulse, it indicates $+1$ and the slide moves forwards. As soon as it has moved a distance of $1\frac{1}{4} \mu\text{m}$, the measuring system delivers a positive feedback pulse, which puts the register back one place. It is thus at 0 again, and the slide stops. One positive command pulse thus causes the slide to move forward by $1\frac{1}{4} \mu\text{m}$. Similarly, one negative pulse sets the register at -1 , and the slide moves backwards until, after $1\frac{1}{4} \mu\text{m}$, a negative feedback pulse is produced which sets the register back at 0.

If the register indicates $+5$, e.g. because five command pulses have come in quick succession,

the servo-valve opens further than at $+1$, so that the slide is given a higher velocity. As soon as it has moved one length unit, the measuring system delivers a positive feedback pulse which puts the register back to $+4$, thus reducing the velocity of the slide. The next feedback pulse puts the register back to $+3$, which reduces the velocity of the slide even more, and so on until the register is back at 0. The slide will then have moved five places. If the slide should happen to move too far, e.g. because the system is insufficiently damped, this is corrected by the feedback pulses: as soon as the slide moves $1\frac{1}{4} \mu\text{m}$ too far, another positive feedback pulse is produced, which sets the difference register at -1 . The slide then moves backwards, and the over-displacement is eliminated.

With this system, the position of the slide would differ considerably from the desired position when it had a constant high velocity; however, this deviation is considerably reduced by the introduction of an integrating term in the control. For example, if the register receives 1000 command pulses per second, the slide velocity must be 1.25 mm/s . In order to bring about this velocity, the servo-valve must be out of the equilibrium position, so the difference register must also show a constant deviation from zero during this motion. We see from fig. 8 that at this velocity the register normally indicates $+5$, which means that during this motion the slide will always be $5 \times 1.25 = 6.25 \mu\text{m}$ behind the desired position. This error can be reduced by increasing the steepness of the curve of fig. 8 by a factor of 5, while ensuring that the voltage corresponding to the zero position of the register remains unchanged. The desired velocity is now obtained with the difference register at $+1$, so that the error is only $1.25 \mu\text{m}$. This adjustment is obtained at each velocity by a simple integrating RC circuit placed after the resistance network.

The error due to a change in the velocity is not eliminated in this way. This error could be corrected by introducing a differentiating term as well, but there is no need for this in the present case because the velocity of the slide never changes suddenly during the milling of surface cams.

The integrating term in the control also serves to compensate for drift, which is caused by a number of factors in the control circuit of the step motor. This drift causes the difference register to deviate from its zero position because of the production of extra feedback pulses. This deviation can also be reduced to $1/5$ of the value it would have had in the absence of the RC circuit, in the same way as described above.

After this general description of the operation of the step motor, follows a description of some of its parts in somewhat greater detail.

The hydraulic servomotor

We shall start with the hydraulic servomotor (fig. 6). Let us first consider the state in which the piston is at rest. The amounts of oil Q_1 and Q_2 which flow in and out of the right-hand part of the cylinder are then equal. The servo-valve is then open a distance

h such that the flow-resistance of both ports (the opening in the piston and that in the valve) is the same. Since the same amount of oil flows through both ports, the pressure drop over each port is the same, so that the pressure in the right-hand space is half the pump pressure. The right-hand end of the piston is given twice the area of the left-hand end, so that the forces on the piston now cancel out.

Now let us consider what happens when the piston moves. We shall neglect for the moment the whole load on the piston, consisting of the force which the workpiece exerts on the cutter, and the forces of friction and inertia. If these can be neglected, no force is needed for the motion of the piston, so that the oil pressures on the left-hand and right-hand sides of the piston will again be equal to the pump pressure and half the pump pressure respectively. The pressure drop over each port is thus again half the pump pressure. Since the flow rate of oil through a port depends only on the size of the aperture and the pressure drop, Q_1 remains constant while the flow Q_2 through the servo-valve is proportional to the distance h (fig. 6). As h is increased, Q_2 thus increases, while Q_1 remains constant. The piston now moves to the right, with a velocity proportional to $h-h_0$ (where h_0 is the value of h at which the piston is at rest). Similarly, the piston moves to the left when h is reduced.

In fact, the situation is not as simple as this, because the load on the piston is not negligible. With a dynamic load, pressure variations will occur in the oil that compensate for the load. These variations are however so small that we can neglect them for the purposes of the present consideration (they do however play a role in determining the stability of the system). With a constant load, caused by friction between the piston and the cylinder or by an external force, the piston will come to a new equilibrium position, at slightly different values of the oil pressure and the valve aperture h . It is true that the pressure drops over the ports will now be different, but the same considerations as given above still hold for the motion of the piston in this state. However, the displacement of the equilibrium alters the setting of the control circuit of the step motor, which has an adverse effect on the accuracy.

A further, very troublesome consequence of friction is the dead zone: the servo-valve can move through a certain distance without having any effect on the piston, because of the friction.

In order to meet the highest demands as regards accuracy, certain special constructional details are included in the servomotor to eliminate the above-mentioned effects of the load. These details will be

described in Part III of this article, which gives a general discussion of hydraulic servomotors.

The resistance network

Fig. 9 shows a circuit diagram of the resistance network, which produces from the voltages of the five flip-flops of the difference register the voltage V

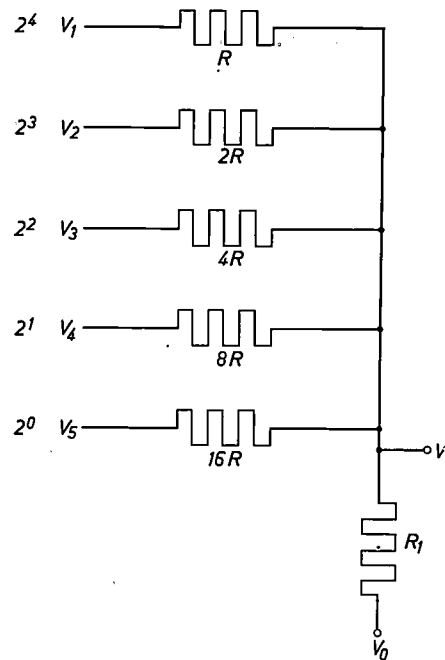


Fig. 9. The resistance network. The voltages $V_1...V_5$ are the output voltages of the five flip-flops of the difference register. V_0 is an adjustable constant voltage. The resistances have been chosen so that the voltage V is a linear function of the number contained in the register.

(see fig. 8) used to feed the DC amplifier for the electromagnet of the servo-valve. By use of Kirchhoff's laws, it is easy to show that the following relationship holds for the circuit shown:

$$V = \frac{16 V_0 + \frac{R_1}{R} (16V_1 + 8V_2 + 4V_3 + 2V_4 + V_5)}{16 + 31 \frac{R_1}{R}}$$

where V_0 is an adjustable constant voltage and $V_1...V_5$ are the output voltages of the flip-flops of the register, V_1 corresponding to the most significant bit (2^4) and V_5 to the least significant (2^0). These latter voltages are either zero or have a certain positive value, depending on whether the bit in question is 0 or 1. We see that the term between brackets in the expression for V gives the number in the difference register. The voltage V is thus a linear function of this number. If we plot V against the contents of the register we get the step-shaped curve

of fig. 8. The whole curve can be raised or lowered by varying V_0 , and its steepness can be altered by changing R_1/R . These two possibilities are made use of to adjust the slide velocity and the stability.

The system for measuring the rotation of the turntable

The turntable is driven by an electric motor via a chain transmission and a worm gear. The rotation of the turntable is measured on the worm gear, a very accurate design being used. Even if there should be a certain amount of play between the worm shaft and the turntable, this would have no influence on the results because the turntable always turns in the same direction. A disc with holes round the edge is fitted on the worm shaft. A light source is placed on one side of the disc, and a photocell on the other, so that the photocell receives a fluctuating luminous flux as the disc rotates. A pulse former produces eight voltage pulses per period of the photocurrent from the photocell. The worm-gear ratio is 1:180 and the disc has 512 holes, so that one pulse corresponds to a rotation of the turntable through $1/2048$ degree. These sync. pulses are fed to the computing unit.

This measuring system is incremental, i.e. it does not measure absolute values but only changes. This could be a disadvantage, since the loss of a pulse cannot be corrected for and would thus give rise to a definite error. However, a sync. pulse is never lost, and it is possible to take adequate measures against the production of parasitic pulses. Moreover, there is a possibility of checking the accuracy of the measurement, since the programmes are also made so that the slide should always be in its original position at the end of a machining cycle, while the turntable should have rotated through a definite angle.

The system for measuring the displacement of the slide

The displacement of the slide is also measured incrementally by an optical method. This measurement must be carried out with exceptional accuracy, i.e. $1:2 \times 10^5$ ($1\frac{1}{4}$ μm in 25 cm). An optical system which was recently described in this journal²⁾ comes into consideration for measurements of this accuracy. It contains two diffraction gratings, one of which is turned a little with respect to the other, and makes use of the moiré fringe pattern produced when these two gratings move relative to one another. H. de Lang of our laboratories has developed a similar system which only needs one diffraction grating.

This has the further advantage that the movement of the grid is less critical than in the two-grid system. A publication on this new system is in preparation.

The construction of the milling machine

Figures 10 and 11 give an idea of the construction of the milling machine; see also fig. 1. The mechanical parts visible in these photos are mounted on a steel plate 10 cm thick. Fig. 10 shows the turntable in the centre of the picture and fixed on it, a cam in the process of being milled. The electric motor which drives the turntable can be seen top right; the worm shaft is mounted horizontally under the turntable. The end of the worm shaft with the chain-wheel on it can be seen bottom right in the photo. The box on the left of the chain contains the system for measuring the rotation of the turntable.

On the left in fig. 10 can be seen the slide with the cutter head mounted on it; the latter has a belt drive. The levers seen behind the cutter head are for moving it towards the cylinder to be milled and fixing it in the desired position. The box on which these levers are mounted also contains an electromagnetically operated safety device which withdraws the cutter head from the object being milled if an error should occur in the control.

Fig. 11 shows the set-up from the rear. The drive and measuring systems for the turntable are now visible on the left. On the right can be seen the hydraulic servomotor with the slide behind it, and in the middle foreground a box containing the servo-valve. The oil reservoir and the pump which pumps the oil to the servomotor are contained under the perforated cover partially visible at the right of the photo. The oil pipes from the servomotor to the servo-valve and from the valve back to the reservoir can also be seen. The electronic control equipment is contained in the cabinet which can be seen in fig. 1 next to the milling machine. This equipment consists of 400 units, in which a total of about 1000 transistors and 1400 diodes are used. The operation of this control equipment will be discussed in Part II of this article.

Results

The surface cams made with this machine fully meet the demands made on them. These cams can now be given a more complicated shape than was possible when they were made by hand, which opens new possibilities for the manufacture of grids for electronic valves. The accuracy obtained comes quite up to expectations. The error in the position of the cutter head is never more than 2 μm ; as far as we know, there are very few other contour milling machines which give such an accuracy. This

²⁾ Philips tech. Rev. 24, 177 (fig. 8), 1962/63 (No. 6).

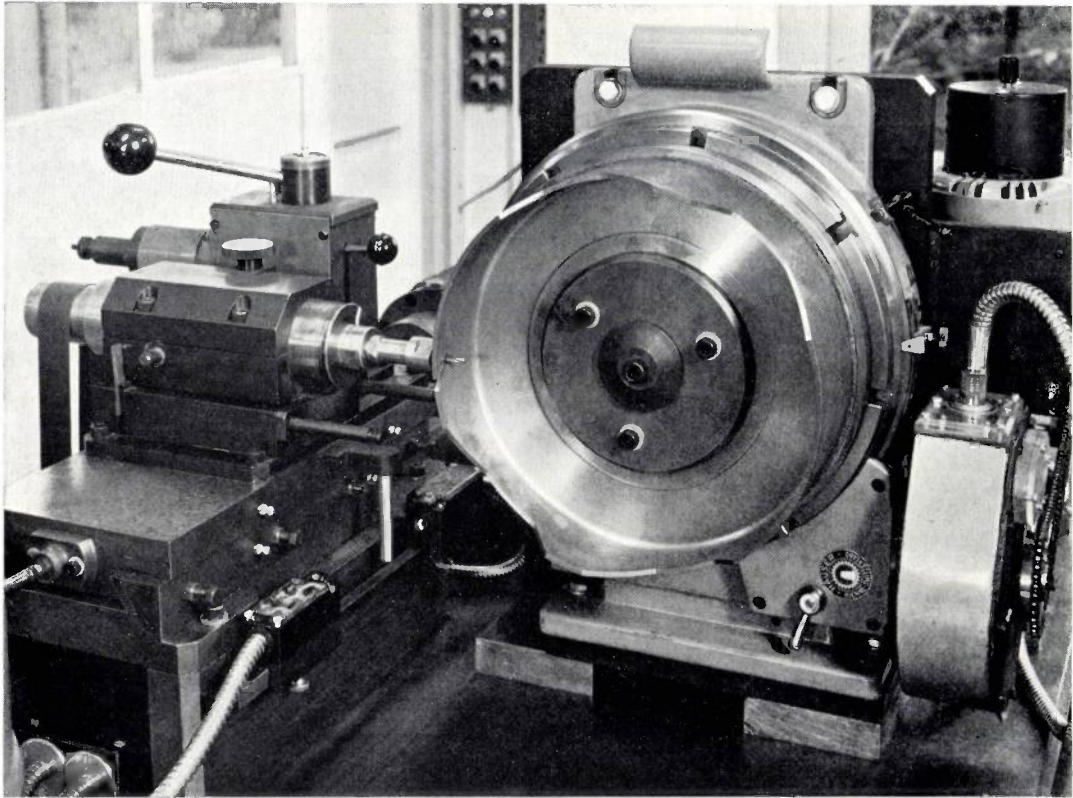


Fig. 10. The arrangement of the milling machine, with a cam in the process of being milled. On the right behind the turntable can be seen the motor with chain transmission for driving the turntable. The system which measures the rotation of the turntable is contained in the box to the right of the turntable. On the left is the cutter head on the slide. Behind the cutter head is an arrangement for moving the cutter to the workpiece and for fixing it in the desired position.

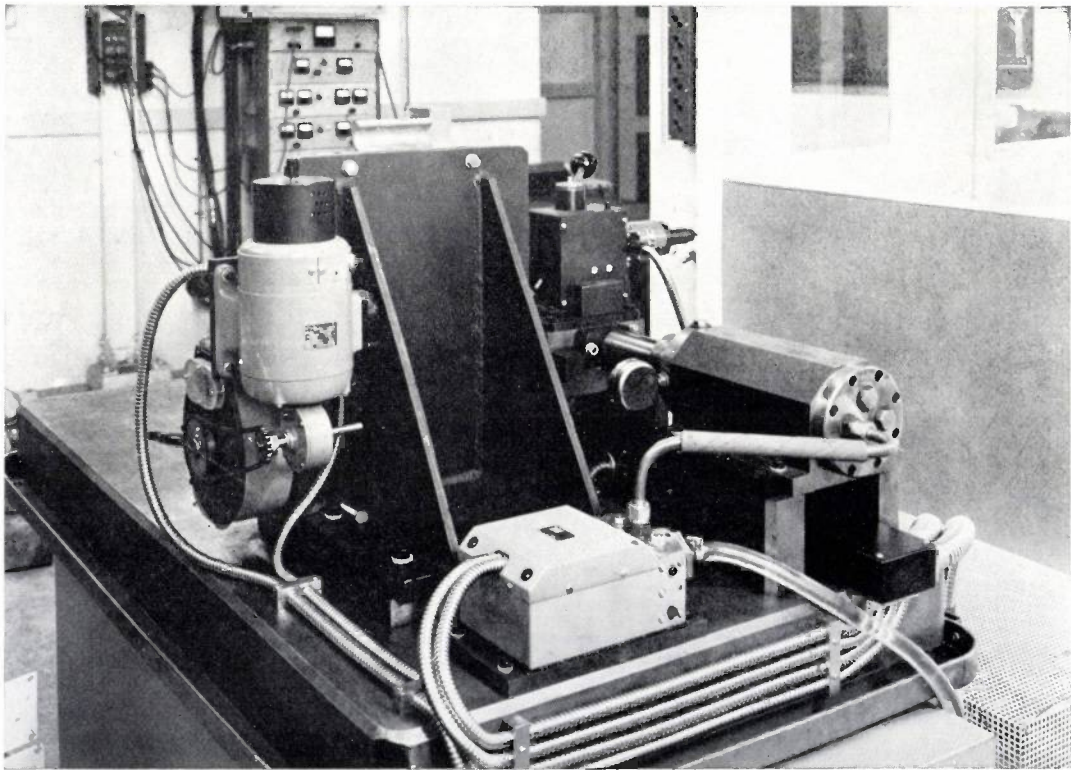


Fig. 11. The set-up of fig. 10 viewed from the rear. On the left the chain drive of the turntable; on the right can be seen part of the slide and in front of it the cylindrical hydraulic servomotor. In the middle foreground the servo-valve which is connected by oil pipes to the servomotor and to the oil reservoir (under the perforated cap together with the pump which gives the oil pressure).

maximum deviation of 2 μm was measured during milling, when cutting forces were low, as is usual in precision work.

The cams were made of cast iron, and each one was milled twice: a preliminary and a final operation. As a consequence of the bending and wear of the cutter during milling, the tolerances of the finished product were greater than the above 2 μm ; the largest error measured was 5 μm .

These sources of error could also be largely eliminated by introducing a second feedback from the cam itself. For example, during machining the dimensions of the section which has just been milled could be measured continuously (we shall not go into how this could be done), so that any error can immediately be compensated for. Alternatively, the shape of the whole cam could be measured before the final cutting, and automatically compared with the

desired shape; the programme of the final cutting could then be corrected in accordance with these findings. These refinements have however not yet been incorporated into our machine

Summary. For internal use Philips have made a numerically controlled contour milling machine for producing surface cams. The electronic-hydraulic control system, which contains about 1000 transistors and 1400 diodes, works on a digital principle. The data on the desired shape of the cam are fed to it on punched tape, which is itself automatically prepared by a digital electronic computer. The milling machine works according to the "derived control" principle. The workpiece is fixed on a vertical turntable, which is given a continuous rotation. The cutter head is mounted on a horizontal slide placed near the turntable. The rotation of the turntable is measured, and at the same time an electronic computing unit calculates the desired motion of the slide from these measurements and the data on the shape of the cam. This motion is produced by a control system which contains a hydraulic servomotor. The accuracy of control of the cutter obtained in this way is 2 μm ; the maximum error found in the finished cam is 5 μm .

II. THE COMPUTING UNIT

by R. C. van OMMERING *) and G. C. M. SCHOENAKER *).

621.914.3-529

In this part of the article we shall give a more detailed description of the operation of the computing unit which was described briefly in Part I. The task of the computing unit is to determine, from data supplied on punched tape, whether the slide must be moved each time the system which measures the rotation of the turntable delivers a sync. pulse. If so, a command pulse must be produced to initiate the desired displacement by means of the step motor.

The computing unit consists of two parts, the interpolator and the internal programme system which will be discussed in turn. After a brief explanation of the numerical notation we shall describe the principle of interpolation used, with reference to a simple linear interpolator.

The numerical notation

The calculations in the interpolator are carried out in the binary notation. It may be assumed that our readers are by now sufficiently acquainted with this notation, which has moreover already been described several times in this journal ¹⁾, so we shall only discuss here how the *negative* numbers are

represented. There are two methods of doing this, the *inverse method* and the *two-complement method*.

In both methods the sign of the number is given by an extra bit (the sign bit) in front of the number, which is a 0 for a positive and a 1 for a negative number. Now in the first method the other bits of a negative number are obtained by inverting the corresponding positive number, i.e. by replacing all ones by zeroes and vice versa. The disadvantage of this method is that the number zero can be represented in two different ways, as may be seen from the following examples:

$$\begin{array}{r} +22 = 0\ 10110 \\ -22 = 1\ 01001 \\ \hline 0 = 1\ 11111 \end{array} + \text{ and } \begin{array}{r} +22 = 0\ 10110 \\ +22 = 0\ 10110 \\ \hline 0 = 0\ 00000 \end{array}$$

This means that it is not possible to add both positive and negative numbers with a single adder: an error would be made each time the zero was passed.

This difficulty does not arise in the second method, where a negative number is represented as the complement of the positive number with respect to a power of two; in the case of a number which consists of n bits (including the sign bit), this power of two is 2^n . The notation for -22 is thus found by subtracting $+22$ from 2^6 .

$$\begin{array}{r} 2^6 = 10\ 00000 \\ +22 = 0\ 10110 \\ -22 = 1\ 01010 \\ \hline 0 = 10\ 00000 \end{array} +$$

*) Philips Research Laboratories, Eindhoven.

¹⁾ See e.g. W. Nijenhuis, The PASCAL, a fast digital electronic computer for the Philips Computing Centre, Philips tech. Rev. 23, 1-18, 1961/62 (No. 1).

maximum deviation of 2 μm was measured during milling, when cutting forces were low, as is usual in precision work.

The cams were made of cast iron, and each one was milled twice: a preliminary and a final operation. As a consequence of the bending and wear of the cutter during milling, the tolerances of the finished product were greater than the above 2 μm ; the largest error measured was 5 μm .

These sources of error could also be largely eliminated by introducing a second feedback from the cam itself. For example, during machining the dimensions of the section which has just been milled could be measured continuously (we shall not go into how this could be done), so that any error can immediately be compensated for. Alternatively, the shape of the whole cam could be measured before the final cutting, and automatically compared with the

desired shape; the programme of the final cutting could then be corrected in accordance with these findings. These refinements have however not yet been incorporated into our machine

Summary. For internal use Philips have made a numerically controlled contour milling machine for producing surface cams. The electronic-hydraulic control system, which contains about 1000 transistors and 1400 diodes, works on a digital principle. The data on the desired shape of the cam are fed to it on punched tape, which is itself automatically prepared by a digital electronic computer. The milling machine works according to the "derived control" principle. The workpiece is fixed on a vertical turntable, which is given a continuous rotation. The cutter head is mounted on a horizontal slide placed near the turntable. The rotation of the turntable is measured, and at the same time an electronic computing unit calculates the desired motion of the slide from these measurements and the data on the shape of the cam. This motion is produced by a control system which contains a hydraulic servomotor. The accuracy of control of the cutter obtained in this way is 2 μm ; the maximum error found in the finished cam is 5 μm .

II. THE COMPUTING UNIT

by R. C. van OMMERING *) and G. C. M. SCHOENAKER *).

621.914.3-529

In this part of the article we shall give a more detailed description of the operation of the computing unit which was described briefly in Part I. The task of the computing unit is to determine, from data supplied on punched tape, whether the slide must be moved each time the system which measures the rotation of the turntable delivers a sync. pulse. If so, a command pulse must be produced to initiate the desired displacement by means of the step motor.

The computing unit consists of two parts, the interpolator and the internal programme system which will be discussed in turn. After a brief explanation of the numerical notation we shall describe the principle of interpolation used, with reference to a simple linear interpolator.

The numerical notation

The calculations in the interpolator are carried out in the binary notation. It may be assumed that our readers are by now sufficiently acquainted with this notation, which has moreover already been described several times in this journal ¹⁾, so we shall only discuss here how the *negative* numbers are

represented. There are two methods of doing this, the *inverse method* and the *two-complement method*.

In both methods the sign of the number is given by an extra bit (the sign bit) in front of the number, which is a 0 for a positive and a 1 for a negative number. Now in the first method the other bits of a negative number are obtained by inverting the corresponding positive number, i.e. by replacing all ones by zeroes and vice versa. The disadvantage of this method is that the number zero can be represented in two different ways, as may be seen from the following examples:

$$\begin{array}{r} +22 = 0\ 10110 \\ -22 = 1\ 01001 \\ \hline 0 = 1\ 11111 \end{array} + \text{ and } \begin{array}{r} +22 = 0\ 10110 \\ +22 = 0\ 10110 \\ \hline 0 = 0\ 00000 \end{array}$$

This means that it is not possible to add both positive and negative numbers with a single adder: an error would be made each time the zero was passed.

This difficulty does not arise in the second method, where a negative number is represented as the complement of the positive number with respect to a power of two; in the case of a number which consists of n bits (including the sign bit), this power of two is 2^n . The notation for -22 is thus found by subtracting $+22$ from 2^6 .

$$\begin{array}{r} 2^6 = 10\ 00000 \\ +22 = 0\ 10110 \\ -22 = 1\ 01010 \\ \hline 0 = 10\ 00000 \end{array} +$$

*) Philips Research Laboratories, Eindhoven.

¹⁾ See e.g. W. Nijenhuis, The PASCAL, a fast digital electronic computer for the Philips Computing Centre, Philips tech. Rev. 23, 1-18, 1961/62 (No. 1).

If we now add +22 and -22, we obtain 10 00000. The first bit of this number has no significance, so that with this method zero is always represented by 0 00000.

A simple rule for determining the notation for a negative number is as follows: invert the positive number, including the sign bit, and add 1 to the least significant bit (the one on the right). The following examples will make this rule clear:

$$\begin{array}{rcl}
 +22 = 0\ 10110 & +45 = 0\ 101101 & \\
 \text{inv. } 1\ 01001 & \text{inv. } 1\ 010010 & \\
 & 1 & \\
 -22 = 1\ 01010 & -45 = 1\ 010011 &
 \end{array}$$

Since, as we shall see, the whole interpolation process as carried out in the computing unit consists of the repeated addition of positive and negative numbers, the two-complement method is used here.

Linear interpolation

Let us suppose that the function $Z = f(\varphi)$, which represents the path followed by the cutter, has the form shown in fig. 1 (φ is the angle of rotation of the turntable, see fig. 2 in Part I of this article). We approximate to this curve in the intervals between φ_0, φ_1 and φ_2 by the straight line segments AB and BC . The values of the function at the points A, B and C (Z_0, Z_1 and Z_2) are for this purpose rounded off to integral values of Z , and the interval width is so chosen that the deviation from the real value of

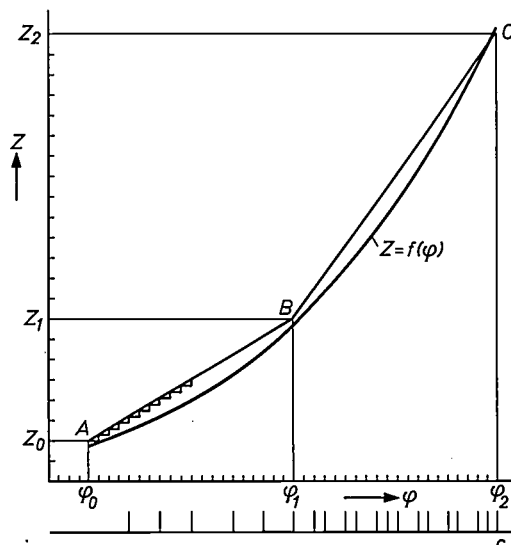


Fig. 1. Approximation to the function $Z = f(\varphi)$ by linear interpolation. The φ axis is graduated in angle units, and the Z axis in length units. The straight line segments AB and BC are so chosen that the deviation from the true value of the function is never greater than one length unit. The contents of the Z register follow a step-shaped curve just under the lines AB and BC . Each time that the value of Z changes by one unit, a command pulse c must be delivered. These pulses are shown under the graph. (In the case shown here, only positive command pulses occur.)

the function is never greater than one unit. The rotation in angle units (an angle unit is the rotation of the turntable corresponding to one sync. pulse) is plotted along the φ axis. The change of Z for each rotation through one angle unit must now be calculated.

The interpolator carries out this calculation with the aid of two binary registers (series of binary elements, e.g. flip-flops, which contain a binary number), the Z register and the ΔZ register, and an adder circuit. At the beginning of the interval AB , the value Z_0 is in the Z register. A number $\Delta Z_0 = (Z_1 - Z_0)/N$, where N is the interval width, expressed in angle units, is placed in the ΔZ register. This quantity is a measure of the slope of the line AB . Now every time a sync. pulse appears, the contents of the ΔZ register are added to those of the Z register. After m additions the contents of the Z register are thus:

$$Z_m = Z_0 + m\Delta Z_0 \dots \dots (1)$$

The value of Z thus increases step-wise just under the line AB as the interval is traversed. After N additions, the value of the Z register has increased by $Z_1 - Z_0$ and thus contains the correct value of Z, Z_1 . For the following interval, the value $(Z_2 - Z_1)/N = \Delta Z_1$ must be placed in the ΔZ register, and the whole process repeated. One number per interval must therefore be supplied to the interpolator on the punched tape.

Each time the contents of the Z register increase or decrease by one unit during the traversing of an interval, the computing unit must deliver a command pulse to the step motor. (As already mentioned in Part I, with the milling machine described here the displacement of the slide per sync. pulse cannot exceed one length unit.) Since only the change of Z is made use of, the total value of the function need not be present in the Z register, which therefore contains only the part of the number after the "binary point", and one bit in front of the point: the "overflow bit". The latter indicates when the Z register "overflows", i.e. when its contents pass from 0.1111... to 1.0000... When this happens, a command pulse must be delivered. The ΔZ register also begins after the point; the number which it contains is always less than unity. This register also contains a sign bit. The contents of the registers are expressed as integral multiples of the smallest unit.

What happens during interpolation will now be shown with the aid of a simple example. Suppose that the interval width is $2^5 = 32$ angle units, that the registers can contain 5 bits and 1 sign bit, and that the ΔZ register contains the number +8. The slope

of the straight line used to approximate to the function is thus $8/32$, so that one command pulse must be supplied every four sync. pulses. If the value of Z at the beginning of the interval is zero, interpolation proceeds as follows:

	ΔZ reg.	Z reg.	overflow bit
	0.01000	0.00000	0
after first addition	„	0.01000	0
„ second „	„	0.10000	0
„ third „	„	0.11000	0
„ fourth „	„	1.00000	1
make overflow bit zero		0.00000	
after fifth addition	„	0.01000	0

After the fourth addition, the overflow bit has become 1, and the value of Z has changed by one unit, so that now a command pulse must be delivered. The sign bit of the ΔZ register indicates whether this change is an increase or a decrease. In the example given here, the command pulse must be produced at the positive output of the computing unit. After this pulse has been delivered, the flip-flop of the overflow bit is reset to 0. After the eighth addition the overflow bit again becomes 1 and the second command pulse is given, and so on. If the ΔZ register contains the number -8 , interpolation proceeds as follows:

	ΔZ reg.	Z reg.	overflow bit
	1.11000	0.00000	0
after first addition	„	1.11000	1
make overflow bit zero		0.11000	
after second addition	„	0.10000	0
„ third „	„	0.01000	0
„ fourth „	„	0.00000	0
„ fifth „	„	1.11000	1

In this example, the Z register overflows after the first and fifth additions, and each time a negative command pulse is delivered. Here too, a pulse is produced every four sync. pulses.

Parabolic interpolation

In parabolic interpolation, the function $Z = f(\varphi)$ is approximated to by parabolas (see fig. 2). The value of Z at the interval boundaries A , B and C is again rounded off to a whole number of Z units, and parabolas are drawn through these points so that the deviation from the true value of Z is never greater than one length unit. It is clear that the interval width can now be chosen much greater than in linear interpolation. The parabolas are also so chosen that the breaks in the interpolated curve at the interval boundaries are as small as possible.

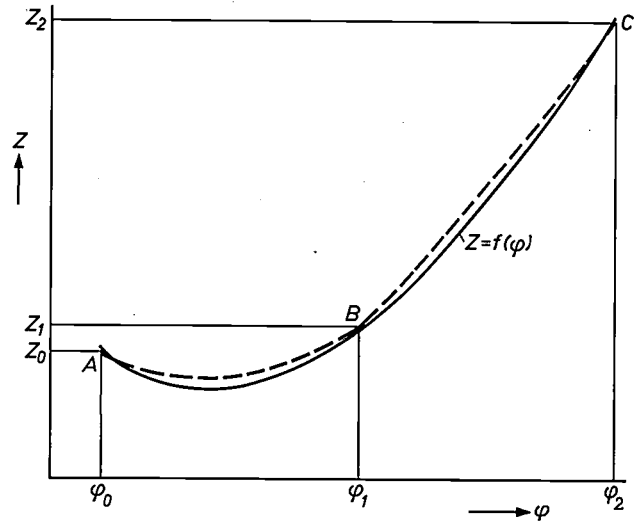


Fig. 2. Approximation to the function $Z = f(\varphi)$ (full line) by parabolic interpolation. The broken lines AB and BC are parabolas, which are chosen so that the deviation from the true value of the function is never greater than one length unit, while at the same time the breaks in the curve between intervals are as small as possible. The interval with parabolic interpolation can be chosen much longer than with linear interpolation.

The parabolic interpolator operates with three registers, the Z , ΔZ and $\Delta^2 Z$ registers (fig. 3). At the beginning of the interval AB (see fig. 2) the Z register contains the number Z_0 . In the ΔZ register is placed a number ΔZ_0 which indicates the slope of the parabola at the beginning of the interval, while the $\Delta^2 Z$ register contains a number $\Delta^2 Z_0$ which is a measure of the curvature of the parabola. The interpolation is now carried out as follows. Each time a sync. pulse appears, the contents of the $\Delta^2 Z$ register are added to those of the ΔZ register, which are then added in their turn to the contents of the Z register. With every sync. pulse, therefore, the slope increases by $\Delta^2 Z_0$. After m sync. pulses, the Z register thus contains the value:

$$Z_m = Z_0 + m\Delta Z_0 + \frac{1}{2}m(m+1)\Delta^2 Z_0. \quad (2)$$

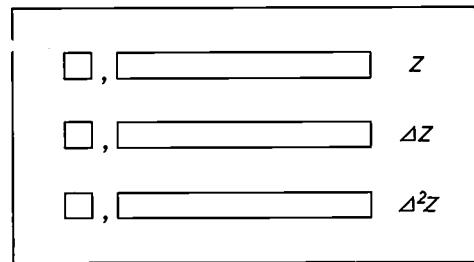


Fig. 3. Sketch of the registers in a parabolic interpolator. The Z register contains an extra bit (the "overflow bit") in front of the point, while the ΔZ and $\Delta^2 Z$ registers contain a "sign bit" in front of the point. At each step of the interpolation, $\Delta^2 Z$ is first added on to ΔZ , and the latter is then added to Z . If this causes the overflow bit of the Z register to become 1, a command pulse must be delivered to the step motor.

If the interval length is N , this value of Z occurs at

$$\varphi_m = \varphi_0 + \frac{m}{N}(\varphi_1 - \varphi_0).$$

It follows from these two equations that Z_m is indeed a parabolic function of φ_m .

At the end of the interval, the Z register contains the value Z_1 , and the ΔZ register contains $\Delta Z_0 + N\Delta^2 Z_0$, i.e. the slope of the parabola of the first interval at the point B . For the second interval, this must be replaced by the slope ΔZ_1 of the new parabola at the point B , while the $\Delta^2 Z$ register must now contain the curvature $\Delta^2 Z_1$ of this parabola. This is done by adding to the contents of these registers the numbers $\Delta Z_c = \Delta Z_1 - (\Delta Z_0 + N\Delta^2 Z_0)$ and $\Delta^2 Z_c = \Delta^2 Z_1 - \Delta^2 Z_0$, which are supplied on the punched tape.

As in linear interpolation, the registers only contain the parts of the numbers after the point. The Z register again contains the overflow bit in front of the point, while the ΔZ and $\Delta^2 Z$ registers have a sign bit (see fig. 3). The command pulses are again produced when the Z register overflows, as described for linear interpolation.

Combined interpolation

The method of interpolation used in the computing unit of our milling machine is a combination of linear and parabolic interpolation, which we shall call "combined" interpolation. In an interval of 256 angle units, 15 points are parabolically interpolated (i.e. they lie on a parabola), while in each of the 16 shorter intervals thus produced 15 points are linearly interpolated. With the path functions met with in the production of surface cams, this method is just as accurate as parabolic interpolation at intervals of 256 angle units (the deviation is again nowhere greater than one unit). The larger number of breaks in the curve causes no trouble, as these breaks are very small. The advantage of combined interpolation is that the register length can be made considerably smaller, which means a saving in the number of flip-flops needed. In the interpolator we used, the registers contained 13 bits after the point.

The interpolator contains the same registers as for parabolic interpolation, but now interpolation is carried out as follows. At the points which are parabolically interpolated, $\Delta^2 Z$ is added to ΔZ , which is then added to Z ; at the intermediate points, which are linearly interpolated, only ΔZ is added to Z . The curvature $\Delta^2 Z$ must thus be 16 times as large as for parabolic interpolation throughout the whole 256 angle-unit interval. At the end of the interval,

two new numbers (ΔZ_c and $\Delta^2 Z_c$) must again be fed to the interpolator to get the new ΔZ and $\Delta^2 Z$. To illustrate the principle of this method, fig. 4 shows the approximation of a function by combined interpolation for a certain interval in which three points are parabolically interpolated and intermediate points are interpolated linearly.

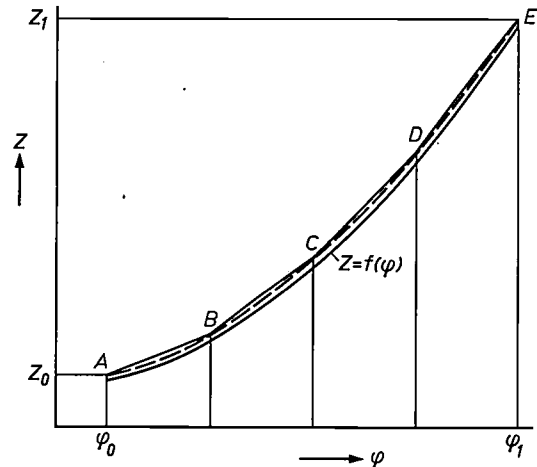


Fig. 4. Approximation to the function $Z = f(\varphi)$ by combined interpolation. In the interval from φ_0 to φ_1 , a number of points (in this example the three points B , C and D) are parabolically interpolated. These points, and the end points A and E , thus lie on a parabola. Linear interpolation is used in the sub-intervals between these points.

If the function has a very slight curvature over a certain range, a lot of punched tape can be saved by using linear interpolation in this range. The interpolator used in our machine can in such cases interpolate linearly with an interval of 8×256 angle units, i.e. 8 times as long as the interval for combined interpolation. We will illustrate the saving thus obtained with some examples.

A complete revolution of the turntable contains 737 280 angle units, each one corresponding to an angle of $1/2048^\circ$; an interval of 256 angle units thus corresponds to $1/8^\circ$. In combined interpolation, two numbers ΔZ_c and $\Delta^2 Z_c$ must be supplied per interval, which takes up two lines of 8-track punched tape. If $\Delta^2 Z_c$ is negative, another line is needed for the sign bit of this number. On the average, therefore, 2.2 lines are needed per interval. The punched tape used contains four lines per cm. If combined interpolation were used for the whole milling operation, the length of punched tape needed would be about 16 metres.

In linear interpolation, only one number need be supplied per interval, so one line of punched tape is sufficient for this purpose. With the linear intervals used here (8×256 angle units, i.e. an angle of 1°), only one line of punched tape is needed per degree

of rotation as opposed to an average of 8×2.2 lines with combined interpolation. If linear interpolation could be used for the whole operation, the length of punched tape needed would be only 90 cm. Actually, combined and linear interpolation are used alternately for making surface cams, and the length of tape needed for one milling operation is about 5 metres.

The punched tape must also indicate when the interpolator must change over from one interpolation method to the other. This is done by means of the "special instructions", which are also used to stop the turntable at the end of the milling. If an extra line is needed for the sign bit of $\Delta^2 Z_c$ in combined interpolation, a special instruction which happens to be necessary at that time can be placed on the same line as the sign bit.

The production of the punched tape

The punched tape is produced with the aid of an electronic computer, in this case the PASCAL¹). This computer must first calculate the track $Z = f(\varphi)$ which must be described by the axis of the cutter during milling. The milling is carried out in two phases, preliminary and final. During the preliminary cutting, the cutter used has a smaller diameter than the roller which rests on the cam in the grid-winding machine. The track of the cutter therefore does not coincide with that of the roller (see page 301) and must also be calculated by the computer. The cutter used for the final machining has the same diameter as the roller, so here this calculation is not necessary.

Once the track of the cutter has been calculated, the computer determines the parabolas used to approximate to this function. As explained above, these parabolas must satisfy the conditions that the error should never be more than one unit, and that the breaks in the approximating curve at the interval boundaries should be as small as possible. The computer carries out this determination by trying successive values of ΔZ and $\Delta^2 Z$. This gives several parabolas per interval which satisfy the first condition. The computer then chooses those parabolas which give the smallest breaks in the curve. The computer then calculates the two numbers ΔZ_c and $\Delta^2 Z_c$ which must be added to the contents of the ΔZ and $\Delta^2 Z$ register at the beginning of each interval.

This method for calculating the correction terms has already become out of date in the four years since this milling machine was designed. A mathematically based method has now been developed for calculating these terms very simply and exactly. This method, which we do not yet make use of in this

machine, is described in an appendix to this article.

Apart from the numbers ΔZ_c and $\Delta^2 Z_c$, the punched tape also contains the special instructions mentioned in the previous section.

The tape puncher, or the tape reader which extracts the information for the interpolator, may punch or read one bit incorrectly, and this can be checked with the aid of "parity bits", as is usual in computers. Each line of the punched tape contains one extra bit which does not convey any information but is merely chosen so that the total number of ones per line is even. Before the information is processed by the interpolator, this parity is checked and if an error is found the machine stops automatically.

The design of the interpolator

The various additions which must be performed in the interpolator are all carried out using the same adder circuit. An internal programme system which will be discussed later makes sure that the additions occur in the correct sequence.

In order to make it possible to use only one adder, the registers have been designed as "shift registers". This concerns the Z , ΔZ and $\Delta^2 Z$ registers and a register for the ΔZ_c correction term which will be discussed below. These registers have the property that when a voltage pulse is applied to each flip-flop of the register via a special input, the information contained in the register is shifted one place to the right (i.e. in the direction of the least significant bit). For example, if a 5-bit register read 10110 to start off with, after one shift pulse it would read 01011 and after two such pulses 00101. During this process, new information may if desired be inserted at the front of the number; for example, after two shift pulses the register can read 11101. The design of such a register will be described below.

A one-bit adder is used, with which the registers are added bit by bit in series. We shall take as an example the addition of the contents of the ΔZ register to those of the Z register. The register length, which is actually 13 bits, is shortened to 5 bits for the purposes of the example (see fig. 5). The Z register also contains the overflow bit, which is always zero before an addition, while the ΔZ register contains a sign bit. We shall indicate the initial state of the Z register by $0Z_5Z_4Z_3Z_2Z_1$, that of the ΔZ register by $\Delta Z_t \Delta Z_5 \Delta Z_4 \Delta Z_3 \Delta Z_2 \Delta Z_1$, and the sum by $S_6 S_5 S_4 S_3 S_2 S_1$. If a given digit is 0, the output voltage of the corresponding flip-flop is zero, while if it is 1 the output voltage is 4 volt. The adder has three inputs, two of which are connected to the least significant bit of the two registers while the third is connected to the output of a flip-flop M ,

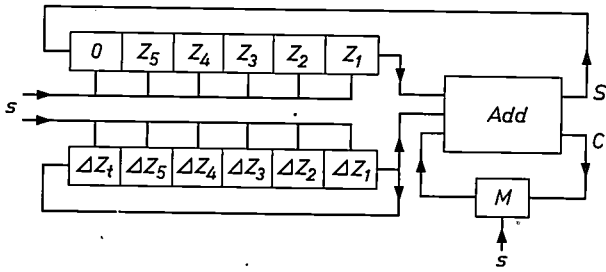


Fig. 5. The circuit for the addition of the contents of the ΔZ and Z registers. The addition is performed bit by bit, with the aid of the 1-bit adder *Add*. Two inputs of this adder are connected to the output of the flip-flop of the least significant bit in these two registers, while the third is connected to the output of a memory flip-flop *M*, which contains the carry from the addition of the previous place (if any). The adder delivers two voltages, one of which represents the sum *S* of the two bits, while the other represents the carry *C* to the following place. The voltage *C* is applied to the input of *M*, and the sum voltage *S* to the left-hand flip-flop of the Z register. If a pulse ("shift pulse") is applied to the inputs marked *s*, the information in the registers is shifted one place to the right; during this process the sum *S* is inserted in the left-hand place of the Z register, and ΔZ_1 in the left-hand place of the ΔZ register. The value of *C* is simultaneously stored in *M*. After six shift pulses, the Z register contains the sum $Z + \Delta Z$, and the ΔZ register the original value of ΔZ .

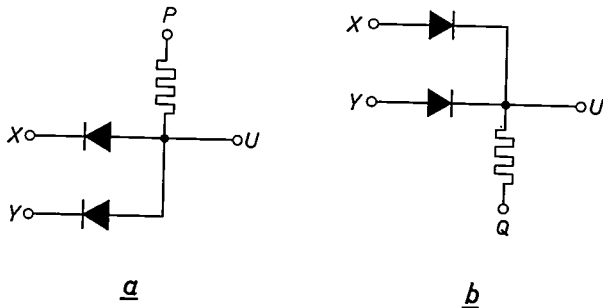
During the addition of the contents of the registers, the "sum" voltage produced by the adder is applied to the left-hand flip-flop of the Z register, while the "carry" voltage is applied to the input of *M*. A shift pulse *s* is now applied to both registers, so that the information which they contain is shifted one place to the right. The sum $S_1 = Z_1 + \Delta Z_1$ is now in the left-hand place of the Z register, which therefore reads: $S_1 0 Z_5 Z_4 Z_3 Z_2$ (see Table I). The number ΔZ_1 is reintroduced at the left-hand side of the ΔZ register, which therefore reads $\Delta Z_1 \Delta Z_5 \Delta Z_4 \Delta Z_3 \Delta Z_2$. A shift pulse is simultaneously applied to *M*, so that the carry *C* is stored here.

Table I. The positions of the Z and the ΔZ register during the adding process in the interpolator after *n* shift pulses.

<i>n</i>	Z -register	ΔZ -register
0	$0 Z_5 Z_4 Z_3 Z_2 Z_1$	$\Delta Z_1 \Delta Z_5 \Delta Z_4 \Delta Z_3 \Delta Z_2 \Delta Z_1$
1	$S_1 0 Z_5 Z_4 Z_3 Z_2$	$\Delta Z_1 \Delta Z_1 \Delta Z_5 \Delta Z_4 \Delta Z_3 \Delta Z_2$
2	$S_2 S_1 0 Z_5 Z_4 Z_3$	$\Delta Z_2 \Delta Z_1 \Delta Z_1 \Delta Z_5 \Delta Z_4 \Delta Z_3$
3	$S_3 S_2 S_1 0 Z_5 Z_4$	$\Delta Z_3 \Delta Z_2 \Delta Z_1 \Delta Z_1 \Delta Z_5 \Delta Z_4$
4	$S_4 S_3 S_2 S_1 0 Z_5$	$\Delta Z_4 \Delta Z_3 \Delta Z_2 \Delta Z_1 \Delta Z_1 \Delta Z_5$
5	$S_5 S_4 S_3 S_2 S_1 0$	$\Delta Z_5 \Delta Z_4 \Delta Z_3 \Delta Z_2 \Delta Z_1 \Delta Z_1$
6	$S_6 S_5 S_4 S_3 S_2 S_1$	$\Delta Z_1 \Delta Z_5 \Delta Z_4 \Delta Z_3 \Delta Z_2 \Delta Z_1$

which serves as a memory for the carry of the addition of the previous two bits. From these voltages the adder forms two output voltages, each of which can be either 0 or 4 volt. One of these indicates the sum (0 or 1), and the other the carry to the next place. These voltages are available with very little delay as soon as the input voltages are applied.

Such an adder is made up of logical circuits. These circuits, which are also used in many other parts of the interpolator and the internal programme system, give a high or a low output voltage depending on whether a certain combination of input voltages is present. Logical circuits are built up of two basic types, the "and" circuit and the "or" circuit (see fig. 6). A more detailed discussion of logical circuits and adders will be found in the article cited in 1).



5134

Fig. 6. The two basic types of logical circuits. a) "and" circuit. The point *P* here has a relatively high positive voltage. The point *U* only has a high voltage if both *X* and *Y* have high voltages. b) "or" circuit. The point *Q* has here a "low" positive voltage. The point *U* then has a high voltage if *X* or *Y* or both have a high voltage. By suitable combination of these two types of circuits, one can make circuits which give a high output voltage only when a certain combination of high and low input voltages is present.

After this first addition, the voltages Z_2 , ΔZ_2 and the output voltage *C* of the memory flip-flop (which was zero for the first addition) are applied to the inputs of the adder. The second shift pulse shifts the information along in the same way. After six shift pulses, the Z register contains the sum $Z + \Delta Z$, and the ΔZ register contains the original value of ΔZ . Each addition requires a series of shift pulses: for the real register length of 13 bits + 1 sign bit, 14 shift pulses are needed.

The interpolator makes use of "clock" pulses for these shift pulses. These are voltage pulses of +3 volt and a frequency of 150 kc/s, and their name refers to another function which they perform, that of synchronizing the various operations carried out in the machine. The clock pulses are applied to the flip-flops of the register via gate circuits, which can be opened by applying a positive voltage to them. These voltages, the "gate" voltages, are supplied by the internal programme system.

The gate voltage also has another function to perform: it brings about the connection between the different registers and the adder. For this purpose, the interpolator contains a selector circuit *Sel* built up of logical circuits, to which the four registers and the adder are connected (see fig. 7). If a gate voltage is applied to the input *I* of this selector circuit, the Z and ΔZ registers are connected to the adder and the

gate voltage is also applied to these registers, so that the addition takes place. Similarly, gate voltages applied to the other inputs bring about the additions $\Delta Z + \Delta^2 Z$, $\Delta Z + \Delta Z_c$ and $\Delta^2 Z + \Delta^2 Z_c$. A buffer register is provided for ΔZ_c , since it proves necessary to read this term from the tape before the addition is performed. The term $\Delta^2 Z_c$ is added to $\Delta^2 Z$ directly from the tape, so that no buffer register is needed here.

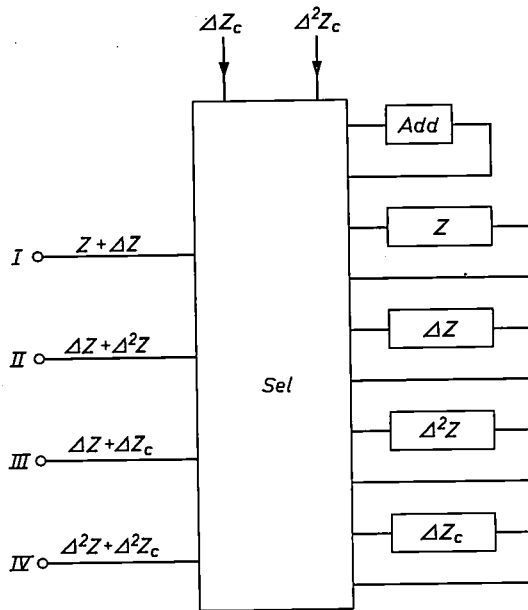


Fig. 7. Sketch of the interpolator with the Z , ΔZ , $\Delta^2 Z$ and ΔZ_c registers, the adder Add and the selector circuit Sel . If a positive voltage ("gate voltage") is applied to one of the inputs shown to the left of the selector circuit, the addition corresponding to this input is carried out.

The internal programme system

When the combined interpolation which has been described above is carried out, various additions must be carried in the right order by the interpolator. The sum $Z + \Delta Z$ must be formed for each of the 256 sync. pulses delivered by the rotation-measuring system in one interval. Further, $\Delta^2 Z$ must be added to ΔZ every 16 pulses, and at the end of the interval ΔZ and $\Delta^2 Z$ must be increased by the correction terms ΔZ_c and $\Delta^2 Z_c$. Moreover, during the interval ΔZ_c and the special instruction (if any) must be read from the tape and placed in a buffer register until needed. In order to bring about the additions, the correct sequence of gate voltages must be applied to certain inputs of the selector circuit of the interpolator each time a sync. pulse appears.

The internal programme system performs all these operations with the aid of two counting circuits, a 6-bit counter called the microprogramme counter P and an 8-bit counter called the interval counter Q . The flip-flops of these counting circuits are connected to logical circuits which, at certain

positions of the counters, deliver voltages which can be used as gate voltages or for other purposes. We shall now discuss the operation of these two counters with reference to fig. 8.

The sync. pulses and the clock pulses are fed to the internal programme system. Each sync. pulse opens a gate circuit G_1 , which lets the clock pulses through to the input of the microprogramme counter P . Now this counter counts from 0 to 63. The logical circuits connected with its six flip-flops all give a voltage which is positive during 14 successive positions of the counter. These voltages are used as gate voltages for the additions. Since the counter P counts the clock pulses which are also used in the adding process in the interpolator, the gate voltages are positive just long enough for one complete addition. The positions of P at which gate voltages are produced are:

- P_1 2 to 15 for the addition $\Delta Z + \Delta Z_c$,
- P_2 18 to 31 for the addition $\Delta Z + \Delta^2 Z$,
- P_3 34 to 47 for the addition $\Delta Z + Z$,
- P_4 53 to 59 for the addition $\Delta^2 Z + \Delta^2 Z_c$.

The terms $\Delta^2 Z$ and $\Delta^2 Z_c$ are at the most 7 bits long, and the registers which contain them are made up of only 7 flip-flops. The fourth gate voltage is therefore shorter than the others: it only lets 7 clock pulses through. After this microprogramme has been carried out, the counter returns to its zero position, whereupon the gate G_1 is closed and the counter is ready to repeat the programme when the next sync. pulse arrives.

We have already seen that not all these additions need to be carried out for every sync. pulse. The interval counter Q determines which additions must in fact take place each time. This counter receives the sync. pulses and counts them from 0 to 255, thus bringing about the division into intervals. The voltages delivered by the logical circuits of this counter at certain positions of the counter operate the gate circuits G_2 , G_3 and G_4 (see fig. 8), which only let through the voltages delivered by the counter P if the addition in question has to be performed. This counter also takes care of the reading of data from the punched tape. The following list shows when the interval counter delivers voltages:

- Q_1 0 for $\Delta Z + \Delta^2 Z$,
- Q_2 0-16-32-...-240 for $\Delta Z + \Delta Z_c$,
- Q_3 80 for reading special instruction and writing it in buffer register,
- Q_4 160 for reading ΔZ_c and writing it in buffer register,
- Q_5 240 for $\Delta^2 Z + \Delta^2 Z_c$ direct from tape.

We see that in this way the addition $Z + \Delta Z$ is performed each sync. pulse, and the addition $\Delta Z + \Delta^2 Z$ every 16th sync. pulse. The additions $\Delta Z + \Delta Z_c$ and $\Delta^2 Z + \Delta^2 Z_c$ must occur in quick succession at the beginning of a new interval. The addition of $\Delta^2 Z_c$ may take place directly after the $\Delta^2 Z$ for the old interval has been added to

should be given. For this purpose, gate G_5 is opened for a short time; if the addition $Z + \Delta Z$ has caused the overflow bit in the Z register to become 1, a command pulse will appear at output c . (The sign of the ΔZ register determines whether this is to be a positive or negative one; we shall not give any details of the circuit which takes care of this.)

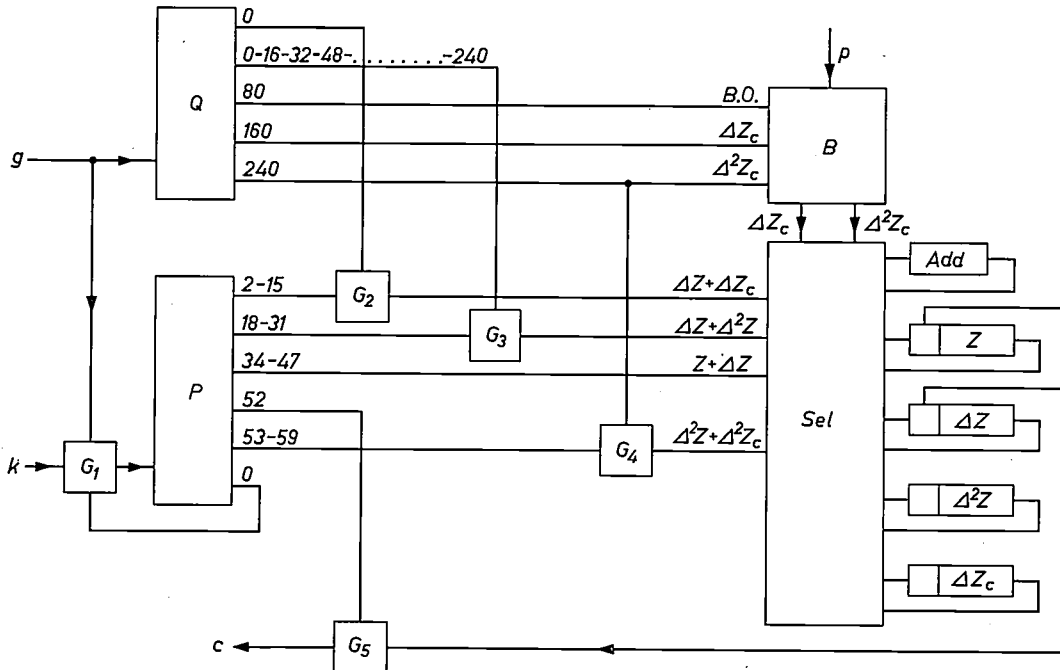


Fig. 8. Block diagram of the control system of the interpolator. Each time a sync. pulse g is received, the gate G_1 is opened and the microprogramme counter P counts 64 clock pulses k . During this period, this counter delivers four voltages which can bring about the various additions involved, with the aid of the selector circuit Sel of the interpolator, as long as the gates G_2, G_3 or G_4 are open. The interval counter Q counts the clock pulses from 0 to 255 in each interval, and opens one of these gates at certain positions of the counter. The counter Q also controls the reading of information from the tape, and the supply of this information p to the interpolator. When the microprogramme counter is at position 52, the gate G_5 is opened. If the overflow bit of the Z register is 1, in this way a command pulse is delivered to the step motor.

ΔZ for the last time, when the interval counter reads 240 (i.e. with the counters at P_4 and Q_5). The addition of ΔZ_c must occur after this, not later than position 0 of the interval counter in the new interval (counters at P_1 and Q_1). The tape reader cannot read two numbers in such quick succession, so only $\Delta^2 Z_c$ is added at the same time as the reading of the tape: ΔZ_c is read at a suitable moment earlier on in the interval, and stored in a buffer register. The same holds for the special instruction, which usually has to be carried out at position Q_1 . The times at which these three data are read from the tape are spread evenly throughout the interval (positions Q_3, Q_4 and Q_5).

When the microprogramme counter reaches 52, a check is made to see whether a command pulse

Some electronic details

Fig. 9 shows the circuit of the flip-flop used for various purposes in the internal programme system and the interpolator, e.g. for each bit of the shift registers. The transistors T_1 and T_2 form a bistable multivibrator, with feedback via the emitter followers T_3 and T_4 . The addition of these emitter followers gives the circuit a very high speed, and allows enough power to be delivered at the outputs a and \bar{a} to supply the logical circuits which are often connected to the flip-flops. In position 1 of the flip-flop, T_1 is cut off and T_2 is conducting, the voltage a is 4 volt and the voltage \bar{a} is 0. The flip-flop changes over to the zero state when a positive voltage pulse is applied to the base of T_2 . Similarly, it changes back again to the 1 position when a voltage pulse is

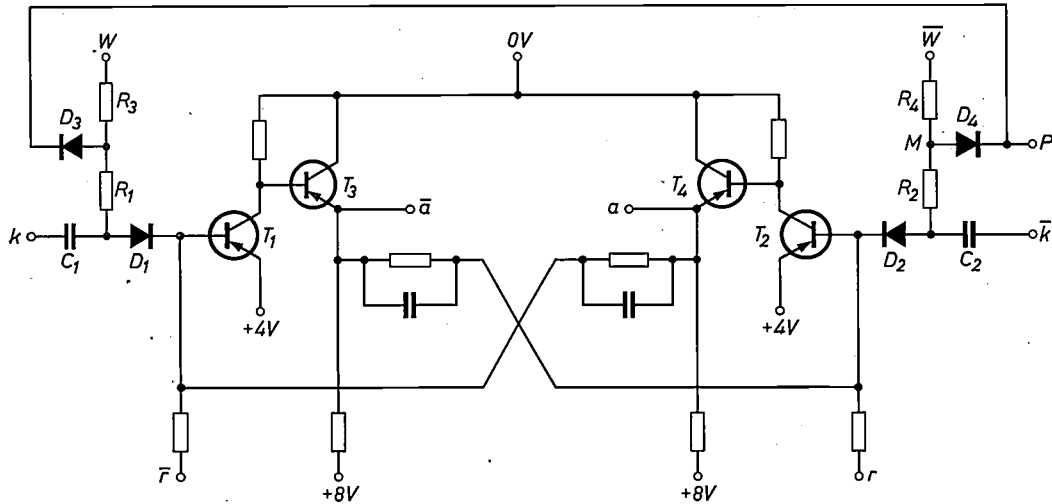


Fig. 9. Circuit diagram of the flip-flop used e.g. in the counters and for each bit of the shift registers. T_1 and T_2 form a bistable multivibrator, with feedback via the emitter followers T_3 and T_4 . Clock pulses are fed to the inputs k and \bar{k} . On receipt of a clock pulse, the flip-flop is set in position 1 if the input P is under a positive voltage and if information is supplied to the input \bar{W} in the form of a positive voltage; if P is positive and the input \bar{W} receives a positive voltage, the flip-flop is set in the 0 position. If a negative voltage appears at the input r or \bar{r} , the flip-flop is set in position 1 or 0 respectively, no matter what its original position and no matter what the voltages at the other inputs.

applied to T_1 . For most applications of the flip-flops, the clock pulses are used for this purpose. The output of the clock-pulse generator is then permanently connected to the inputs k and \bar{k} . Two simple gate circuits, consisting of the diodes and resistors D_1, D_3, R_1, R_3 and D_2, D_4, R_2, R_4 , determine to which side the clock pulse is passed.

Let us suppose that the flip-flop is in position 1 and must be made to change to position 0. A 4-volt positive voltage signal is then applied to the input \bar{W} . If input P also has a positive voltage at this time, diode D_4 is not conducting, so that the voltage at the point M is also +4 volt. The base voltage of T_2 is about 3.8 volt; the voltages on both sides of the diode D_2 are thus about equal, so that the clock pulse is passed through to T_2 and the flip-flop changes over. If however the voltage of P is kept at 0, diode D_4 is conducting; the voltage at M then remains low,

diode D_2 is subjected to an inverse voltage of 4 volt, and the 3-volt clock pulse is not let through. The information about the new position is thus applied to the point \bar{W} or W (for position 0 or 1 respectively), but is only made use of (on arrival of the next clock pulse) if the point P is at a positive voltage (e.g. the gate voltage for a shift register).

As long as the flip-flop must operate in the way we have just described, the points r and \bar{r} are given a voltage of +4 volt, which does not affect the operation of the circuit. If however a negative voltage is applied to r or \bar{r} , the flip-flop is set in position 1 or 0 respectively, no matter what its original position and no matter what the value of the other voltages. This is used e.g. to clear a register.

Of the many applications of the flip-flop, we shall only discuss its use in the shift register in somewhat greater detail (fig. 10). The characteristic feature of

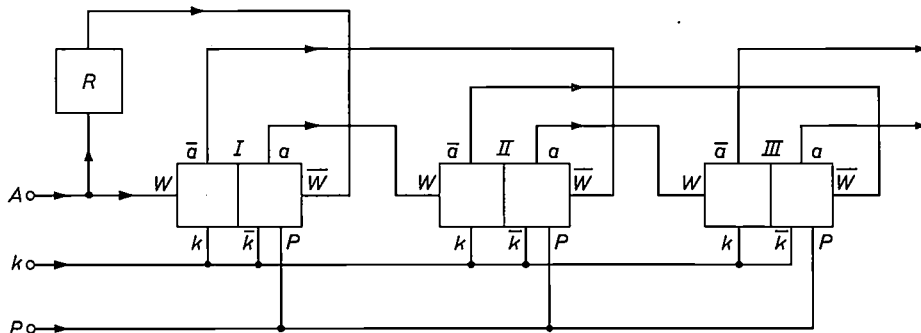


Fig. 10. Sketch showing the connection between the first three flip-flops of a shift register. If a positive voltage is applied to the input P , the information in the register is shifted one place to the right each time a clock pulse appears at the input k .

a shift register is that, as long as a gate voltage is present, the information is shifted one place to the right at every clock pulse. To ensure this, the flip-flops are connected in the following way. The inputs W and \bar{W} of flip-flop *II* are connected with the outputs a and \bar{a} respectively of flip-flop *I*. Now if a gate voltage is present at the point P , flip-flop *II* will take up the same position as flip-flop *I* as soon as a clock pulse comes along. Flip-flop *III* takes over the information from flip-flop *II* in the same way. This happens by the same clock pulse, i.e. at the same moment. There might thus be a certain risk that flip-flop *III* would assume the *new* position of flip-flop *II*. However, each flip-flop has at its input a resistor and a capacitor (R_1-C_1 and R_2-C_2 in fig. 9), which act as a memory thanks to their RC time. The original output voltages of flip-flop *II* are thus available for sufficient time at the bases of flip-flop *III* to avoid this error.

The new information for flip-flop *I* is applied to W via input A (a voltage of 0 or 4 volt, depending on whether flip-flop *I* has to be in the 0 or 1 position). An inverter circuit R forms the inverse of this voltage, which must be applied to input \bar{W} .

The advantage of this circuit is that the clock pulses are applied directly to the flip-flops, and do not need to pass through any other parts of the circuit first. These pulses, the slope of whose leading edges is very important, are thus not deformed much. The gate voltages do have to pass through a large part of the internal programme system, but this does not matter as their form is not so important: it is sufficient that they reach a value of +4 volt on time.

The operation of the milling machine

The electronic circuits of the internal programme system and the interpolator, together with those of

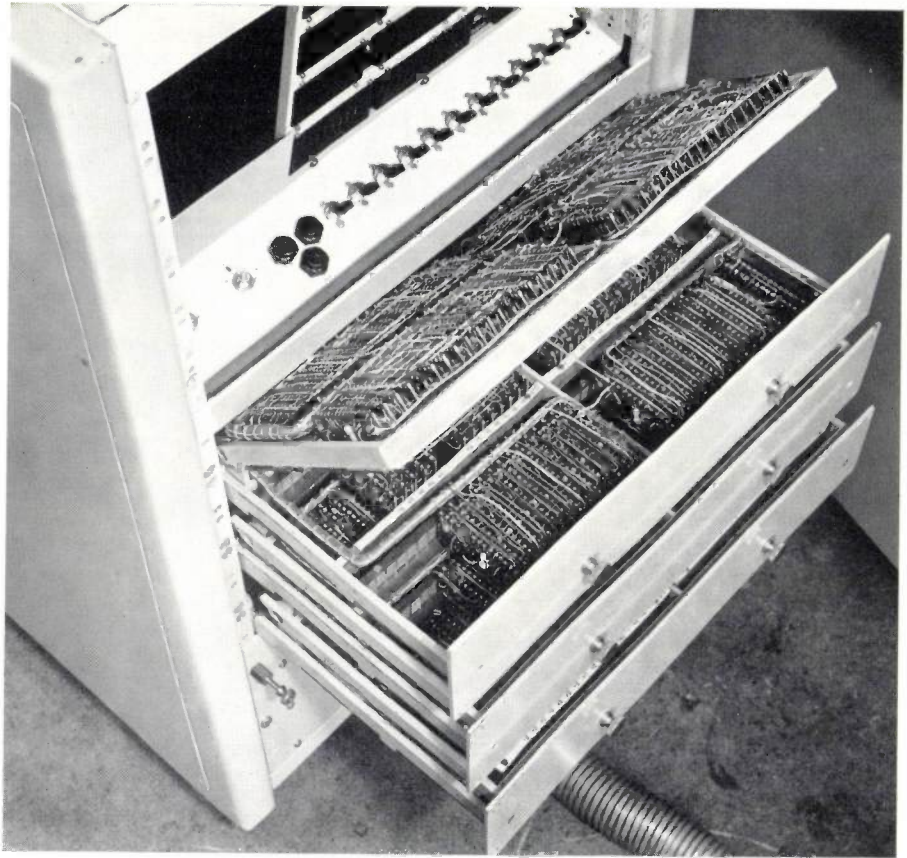


Fig. 11. Back of the cabinet which can be seen next to the milling machine in fig. 1 of Part I, and which contains the various electronic circuits. As may be seen here, the panels on which the components are mounted can be slid out of the cabinet so that the components are easily accessible. The switches at the top of the photo can be used to operate a test programme.

the step motor (difference register, resistance network and amplifier), are mounted in the cabinet shown next to the milling machine in fig. 1 of Part I of this article (page. 300). This cabinet also contains the power-supply equipment for these circuits. The components are conveniently arranged on panels which can easily be reached from the back of the cabinet (fig. 11). Fig. 11 also shows a row of switches to work through a test-programme.

The operating desk for the milling machine is situated on top of this cabinet (fig. 12). On the right, under a glass plate, are two punched-tape readers for the input of information. An eight-channel reader (i.e. eight bits per line) is needed for the input, but at the time this machine was made no good eight-channel reader was available, so temporary use was made of two five-channel readers in parallel (thus leaving two tracks over). Since this photograph was taken, we have replaced these by one single eight-channel reader.

On the left of the panel at the back we see a series of lamps which indicate the position of the difference register of the step motor. To the right of these are

six pilot lights; if the machine stops because something has gone wrong, these lamps indicate the nature of the error. In front of this we see the "zero counter", which indicates the deviation of the slide from its zero position at the end of an operation.

The row of push-buttons in front of this can be used to switch on the various parts of the milling machine separately. The milling machine is set into operation as follows. First of all, only the various

"tion" button. The final corrections to the initial position of the slide can now be made with the switch and push-buttons bottom left on the panel: these can be used to shift the slide in steps of $1.25 \mu\text{m}$ or $20 \mu\text{m}$, corresponding to 1 or 16 command pulses.

After the cutter motor has been switched on, the milling operation can be initiated by pushing the "start" button at the bottom of the panel: the punched-tape reader automatically finds the begin-

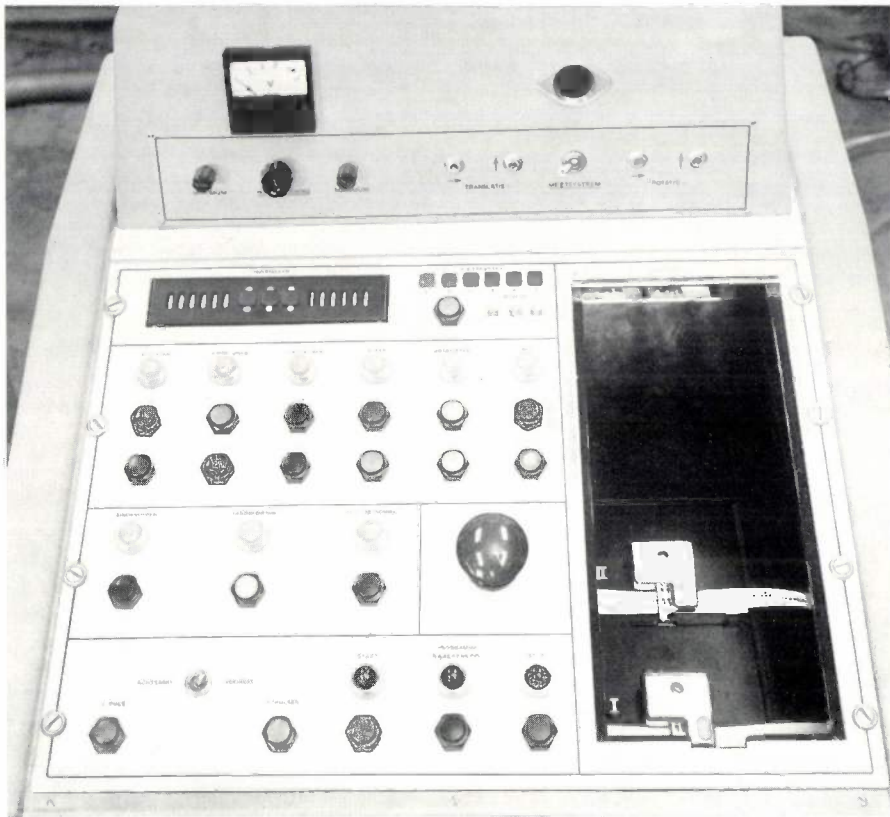


Fig. 12. The operating desk of the milling machine, with the tape readers on the right.

parts of the hydraulic drive are switched on. Then, when the "end-to-end" button is depressed the slide automatically moves to and fro between two fixed end points. This is used to test the hydraulic apparatus, to force any air bubbles out of the oil and to give the slide a chance to run in. The "hand operation" button is then depressed; the slide can now be moved to and fro by hand over a certain range, by means of the three knobs under the meter on the upright part of the panel. This is done to bring the slide as accurately as possible into its initial position. If the slide does not remain stationary in this position, this can be corrected by adjusting the current delivered to the servo-valve with the aid of a potentiometer (not shown in the photo).

The electronic equipment is now switched on, and then the electronic and hydraulic parts of the step motor are connected by pressing the "aut. opera-

ning of the punched tape, the turntable motor is switched on and the computing unit starts delivering command pulses so that milling can begin. As soon as the "start" button is depressed, the switches for the various parts of the milling machine are blocked, so that no damage can be caused by accidentally switching one of them off. When the milling is over, the machine is stopped automatically by a "stop" instruction on the punched tape. It is however also possible to stop the machine during milling with the "stop" button, bottom right on the panel. The cutter is then withdrawn from the material, and the slide returned to its original position. The panel also has a big red "out" button for stopping all parts of the equipment simultaneously in cases of emergency.

The upright part of the panel also contains some instruments used in testing the machine: a meter for the control current through the servo-valve of the

step motor, and to the right of this a small oscilloscope which can be used to inspect the signals from the two measuring systems. If necessary, these can be adjusted with the aid of potentiometers.

Appendix: An improved interpolation system

In this appendix we shall describe the mathematical background of an improved method of interpolation which considerably simplifies the calculation of the data which must be included on the punched tape, and which reduces the length of tape needed.

Interpolation by numerical analysis

The values Z_k of a function $Z = f(\varphi)$ are given at a number of equidistant points φ_k . In order to find the value of the function in the intervals between these points, we construct a table of "central" differences:

φ_k	Z_k	$\delta_{k+1/2}$	δ^2_k	$\delta^3_{k+1/2}$	δ^4_k
φ_0	Z_0				
		$\delta_{1/2}$			
φ_1	Z_1		δ^2_1		
		$\delta_{3/2}$		$\delta^3_{3/2}$	
φ_2	Z_2		δ^2_2		δ^4_2
		$\delta_{5/2}$		$\delta^3_{5/2}$	
φ_3	Z_3		δ^2_3		
		$\delta_{7/2}$			
φ_4	Z_4				

These differences are calculated as follows:

$$\left. \begin{aligned} \text{the 1st difference: } \delta_{k+1/2} &= Z_{k+1} - Z_k, \\ \text{the 2nd difference: } \delta^2_k &= \delta_{k+1/2} - \delta_{k-1/2}, \\ \text{the 3rd difference: } \delta^3_{k+1/2} &= \delta^2_{k+1} - \delta^2_k, \\ \text{the 4th difference: } \delta^4_k &= \delta^3_{k+1/2} - \delta^3_{k-1/2}. \end{aligned} \right\} \dots (3)$$

The values of the function are rounded off to integral values of Z and are expressed in these units; Z_k and the differences are thus whole numbers.

Various formulas are available for interpolating intermediate values with the aid of these differences. We shall use here Bessel's formula²⁾ which we state without proof:

$$Z_{k+p} = Z_k + p\delta_{k+1/2} + \frac{1}{2}p(p-1)(\delta^2_k + \delta^2_{k+1}) + \frac{p(p-\frac{1}{2})(p-1)}{6}\delta^3_{k+1/2} + \frac{p(p+1)(p-1)(p-2)}{48}(\delta^4_k + \delta^4_{k+1}) + \dots (4)$$

where $p = (\varphi - \varphi_k)/(\varphi_{k+1} - \varphi_k)$, and Z_{k+p} is the interpolated value in the interval between φ_k and φ_{k+1} . During the interpolation of a series of values between φ_k and φ_{k+1} , p goes from 0 to 1.

We must now determine how many terms of this series are needed to give the desired accuracy. The values Z_k are normally determined from a function given in analytical form; the maximum rounding-off error is thus $\frac{1}{2}$. We now stipulate that the error caused by truncating the series expansion should also be less than $\frac{1}{2}$, so that the total error does not exceed 1.

In linear interpolation, only the first two terms of Bessel's formula are used:

$$Z_{k+p} = Z_k + p\delta_{k+1/2} \dots (5)$$

It follows from eq. (4) that the remainder will be less than $\frac{1}{2}$ as long as the second differences are less than or equal to 4. The size of the differences is determined by the length of the intervals: a longer interval means bigger differences. Since we want to make the intervals as long as possible, so that the information put on the punched tape can be reduced to the minimum, the interval length for linear interpolation should be chosen so that the second differences are just less than 4.

In parabolic interpolation, Bessel's formula is terminated after the third term:

$$Z_{k+p} = Z_k + p\delta_{k+1/2} + \frac{1}{2}p(p-1)(\delta^2_k + \delta^2_{k+1}) \dots (6)$$

This means that the function is approximated by a parabola in each interval. Here again, the remainder must be less than $\frac{1}{2}$, which in this case implies that the 3rd differences must be less than 62. It will be clear that this makes it possible to use much longer intervals than in linear interpolation.

We shall now apply this theory to the method of interpolation described in this article. The operation of the interpolator remains unchanged, but as we shall see the correction terms can now be calculated exactly in a very simple manner, and less data need be included on the tape for parabolic and combined interpolation than in the method described above.

Linear interpolation

Bessel's formula as used for linear interpolation in the interval between Z_k and Z_{k+1} is (eq. 5):

$$Z_{k+p} = Z_k + p\delta_{k+1/2}.$$

If m points must be interpolated in an interval of N angle units, then $p = m/N$. The equation which we used for linear interpolation in this article (eq. 1) has the form:

$$Z_m = Z_k + m\Delta Z_k$$

in this interval. Equating these two expressions, we find:

$$\Delta Z_k = \frac{1}{N} \delta_{k+1/2}.$$

Similarly, between Z_{k+1} and Z_{k+2} we find:

$$\Delta Z_{k+1} = \frac{1}{N} \delta_{k+3/2}.$$

The correction term which must be added to ΔZ_k between these two intervals is thus

$$\Delta Z_c = \frac{1}{N} (\delta_{k+3/2} - \delta_{k+1/2}) = \frac{1}{N} \delta^2_{k+1} \dots (7)$$

Once the function has been tabulated as shown at the beginning of this appendix, the correction terms can thus simply be read off.

Parabolic interpolation

For parabolic interpolation, Bessel's formula has the form given by eq. (6), again with $p = m/N$. In this article, we used the following equation (eq. 2):

$$Z_m = Z_k + m\Delta Z_k + \frac{1}{2}m(m+1)\Delta^2 Z_k.$$

Equating these two expressions, we find:

$$\Delta Z_k = \frac{1}{N} \delta_{k+1/2} - \frac{N+1}{4N^2} (\delta^2_k + \delta^2_{k+1}), \dots (8)$$

$$\Delta^2 Z_k = \frac{1}{2N^2} (\delta^2_k + \delta^2_{k+1}). \dots (9)$$

The contents of the $\Delta^2 Z$ register remain unchanged throughout the whole interval, while those of the ΔZ register increase during

²⁾ See e.g. D. R. Hartree, Numerical analysis, 2nd Edn., Oxford University Press, London 1958, p. 67.

the interpolation. At the end of the interval, the ΔZ register contains:

$$\Delta Z_k + N\Delta^2 Z_k = \frac{1}{N} \delta_{k+1/2} + \frac{N-1}{4N^2} (\delta^2_k + \delta^2_{k+1}),$$

while the $\Delta^2 Z$ register still contains $\Delta^2 Z_k$ (eq. 9).

At the beginning of the next interval, these values must be replaced by:

$$\Delta Z_{k+1} = \frac{1}{N} \delta_{k+3/2} - \frac{N+1}{4N^2} (\delta^2_{k+1} + \delta^2_{k+2}), \dots (10)$$

$$\Delta^2 Z_{k+1} = \frac{1}{2N^2} (\delta^2_{k+1} + \delta^2_{k+2}). \dots (11)$$

The correction terms are thus:

$$\Delta Z_c = -\frac{1}{2} \Delta^2 Z_c - \frac{1}{4N} \delta^4_{k+1}, \dots (12)$$

$$\Delta^2 Z_c = \frac{1}{2N^2} (\delta^2_{k+2} - \delta^2_k). \dots (13)$$

These two numbers would have to be given by the punched tape. There is however a simpler method: with the aid of eq. (3), we write $\Delta^2 Z_c$ in another form, viz:

$$\Delta^2 Z_c = \frac{1}{2N^2} (2\delta^2_{k+1/2} + \delta^4_{k+1}). \dots (14)$$

It then suffices to add δ^4_{k+1} at the interval boundary, as long as we ensure that the term $\delta^3_{k+1/2}$ is known and is stored in the interpolator. A special memory register is provided for this purpose. The 3rd difference needed for the next interval can be calculated from the 4th difference as follows (eq. 3):

$$\delta^3_{k+3/2} = \delta^3_{k+1/2} + \delta^4_{k+1}, \dots (15)$$

and is then placed in the memory.

A similar argument holds for the combined interpolation described in this article; the method of calculating ΔZ_c and $\Delta^2 Z_c$ is the same as just given.

With this system, the interpolator must carry out three more additions at the interval boundary, and an extra memory register is needed. These disadvantages are however more than outweighed by the advantages, i.e. a considerable reduction in the amount of punched tape needed, and a considerable simplification of the work to be done by the computer which prepares the tape. As mentioned on page 311, the old system

used about 2.2 lines of punched tape per interval. In the method described here, only one number need be added per interval, and this can be put on one line, so that the tape can be shortened to less than half the length. The computer now only has to calculate the differences of the table on page 319, while in the old system it had to calculate a number of parabolas for each interval, and then make a choice between them. It will be clear that this latter process requires a much more complicated programme.

Summary. The computing unit of the milling machine has the task of determining whether the slide on which the cutter is mounted should be moved each time the turntable rotates through one angle unit. This is done by comparing the actual position of the slide with the prescribed position (given on punched tape). If the difference between the two is more than one length unit, the computing unit must deliver a command pulse to the step motor. In order to reduce the amount of punched tape needed, use is made of interpolation. This article starts with a brief exposition of the binary notation used (with negative numbers in the two-complement form), followed by a survey of various methods of interpolation. The combined interpolation method used in this machine is then described: in this method, the path $Z = f(\varphi)$ to be followed by the cutter during the continuous rotation of the turntable is divided up into intervals of 256 angle units. The values of the function at the interval boundaries are rounded off to whole numbers. In each interval, 15 points are interpolated parabolically, and in each of the smaller intervals thus formed 15 points are interpolated linearly. It is then only necessary to provide two numbers per interval, viz ΔZ , a measure of the slope at the beginning of the interval, and $\Delta^2 Z$, a measure of the curvature. (In fact, the method can be made even simpler, as explained in an appendix.) The interpolator used for this purpose contains three registers. At the beginning of an interval, the numbers ΔZ and $\Delta^2 Z$ are placed in the ΔZ and $\Delta^2 Z$ registers, and at the end of the interval the numbers ΔZ_c and $\Delta^2 Z_c$ respectively are added to the contents of these registers, which are now ready for the next interval. These numbers are determined in advance by an electronic computer, which chooses the values so that the error is never more than one Z unit during interpolation. The interpolation is carried out by addition in shift registers, which are actuated by clock pulses with a frequency of 150 kc/s. An internal programme system ensures that the additions are performed in the proper order. This is done with the aid of two counters, a 6-bit microprogramme counter P and an 8-bit interval counter Q . Both counters are connected to logical circuits which at certain positions of the counter deliver voltages and open or close various gate circuits. The article is concluded by a discussion of some electronic details and of the operation of the milling machine.

III. THE HYDRAULIC SERVOMOTOR

by T. J. VIERSMA *).

621-526

Hydraulic servomotors are often used for controlling the motion of machine parts. This is also true of the numerically controlled contour milling machine described in parts I and II of this article. In order to explain the operation of such a motor, we shall consider another example, that of a *copying lathe*, whose control mechanism is simpler. *Fig. 1* shows a cross-section through a copying lathe. The rotating piece of material *1* must be made an exact

replica of the template *2*. The support *3* moves uniformly over the bed *4* of the lathe and carries with it a transverse slide *5* to which the tool *6* is fixed. A follower *7*, which is pressed against the template by a spring *8*, follows the variations in the radius of the template. It is the task of the servomotor to ensure that the tool follows the transverse motion of the probe exactly. The accuracy with which this happens is one of the most important factors determining the copying accuracy.

The hydraulic servomotor consists of a cylinder (an integral part of the transverse slide), a piston *9*

*) Philips Research Laboratories, Eindhoven.

the interpolation. At the end of the interval, the ΔZ register contains:

$$\Delta Z_k + N\Delta^2 Z_k = \frac{1}{N} \delta_{k+1/2} + \frac{N-1}{4N^2} (\delta^2_k + \delta^2_{k+1}),$$

while the $\Delta^2 Z$ register still contains $\Delta^2 Z_k$ (eq. 9).

At the beginning of the next interval, these values must be replaced by:

$$\Delta Z_{k+1} = \frac{1}{N} \delta_{k+3/2} - \frac{N+1}{4N^2} (\delta^2_{k+1} + \delta^2_{k+2}), \dots (10)$$

$$\Delta^2 Z_{k+1} = \frac{1}{2N^2} (\delta^2_{k+1} + \delta^2_{k+2}). \dots (11)$$

The correction terms are thus:

$$\Delta Z_c = -\frac{1}{2} \Delta^2 Z_c - \frac{1}{4N} \delta^4_{k+1}, \dots (12)$$

$$\Delta^2 Z_c = \frac{1}{2N^2} (\delta^2_{k+2} - \delta^2_k). \dots (13)$$

These two numbers would have to be given by the punched tape. There is however a simpler method: with the aid of eq. (3), we write $\Delta^2 Z_c$ in another form, viz:

$$\Delta^2 Z_c = \frac{1}{2N^2} (2\delta^2_{k+1/2} + \delta^4_{k+1}). \dots (14)$$

It then suffices to add δ^4_{k+1} at the interval boundary, as long as we ensure that the term $\delta^3_{k+1/2}$ is known and is stored in the interpolator. A special memory register is provided for this purpose. The 3rd difference needed for the next interval can be calculated from the 4th difference as follows (eq. 3):

$$\delta^3_{k+3/2} = \delta^3_{k+1/2} + \delta^4_{k+1}, \dots (15)$$

and is then placed in the memory.

A similar argument holds for the combined interpolation described in this article; the method of calculating ΔZ_c and $\Delta^2 Z_c$ is the same as just given.

With this system, the interpolator must carry out three more additions at the interval boundary, and an extra memory register is needed. These disadvantages are however more than outweighed by the advantages, i.e. a considerable reduction in the amount of punched tape needed, and a considerable simplification of the work to be done by the computer which prepares the tape. As mentioned on page 311, the old system

used about 2.2 lines of punched tape per interval. In the method described here, only one number need be added per interval, and this can be put on one line, so that the tape can be shortened to less than half the length. The computer now only has to calculate the differences of the table on page 319, while in the old system it had to calculate a number of parabolas for each interval, and then make a choice between them. It will be clear that this latter process requires a much more complicated programme.

Summary. The computing unit of the milling machine has the task of determining whether the slide on which the cutter is mounted should be moved each time the turntable rotates through one angle unit. This is done by comparing the actual position of the slide with the prescribed position (given on punched tape). If the difference between the two is more than one length unit, the computing unit must deliver a command pulse to the step motor. In order to reduce the amount of punched tape needed, use is made of interpolation. This article starts with a brief exposition of the binary notation used (with negative numbers in the two-complement form), followed by a survey of various methods of interpolation. The combined interpolation method used in this machine is then described: in this method, the path $Z = f(\varphi)$ to be followed by the cutter during the continuous rotation of the turntable is divided up into intervals of 256 angle units. The values of the function at the interval boundaries are rounded off to whole numbers. In each interval, 15 points are interpolated parabolically, and in each of the smaller intervals thus formed 15 points are interpolated linearly. It is then only necessary to provide two numbers per interval, viz ΔZ , a measure of the slope at the beginning of the interval, and $\Delta^2 Z$, a measure of the curvature. (In fact, the method can be made even simpler, as explained in an appendix.) The interpolator used for this purpose contains three registers. At the beginning of an interval, the numbers ΔZ and $\Delta^2 Z$ are placed in the ΔZ and $\Delta^2 Z$ registers, and at the end of the interval the numbers ΔZ_c and $\Delta^2 Z_c$ respectively are added to the contents of these registers, which are now ready for the next interval. These numbers are determined in advance by an electronic computer, which chooses the values so that the error is never more than one Z unit during interpolation. The interpolation is carried out by addition in shift registers, which are actuated by clock pulses with a frequency of 150 kc/s. An internal programme system ensures that the additions are performed in the proper order. This is done with the aid of two counters, a 6-bit microprogramme counter P and an 8-bit interval counter Q . Both counters are connected to logical circuits which at certain positions of the counter deliver voltages and open or close various gate circuits. The article is concluded by a discussion of some electronic details and of the operation of the milling machine.

III. THE HYDRAULIC SERVOMOTOR

by T. J. VIERSMA *).

621-526

Hydraulic servomotors are often used for controlling the motion of machine parts. This is also true of the numerically controlled contour milling machine described in parts I and II of this article. In order to explain the operation of such a motor, we shall consider another example, that of a *copying lathe*, whose control mechanism is simpler. *Fig. 1* shows a cross-section through a copying lathe. The rotating piece of material *1* must be made an exact

replica of the template *2*. The support *3* moves uniformly over the bed *4* of the lathe and carries with it a transverse slide *5* to which the tool *6* is fixed. A follower *7*, which is pressed against the template by a spring *8*, follows the variations in the radius of the template. It is the task of the servomotor to ensure that the tool follows the transverse motion of the probe exactly. The accuracy with which this happens is one of the most important factors determining the copying accuracy.

The hydraulic servomotor consists of a cylinder (an integral part of the transverse slide), a piston *9*

*) Philips Research Laboratories, Eindhoven.

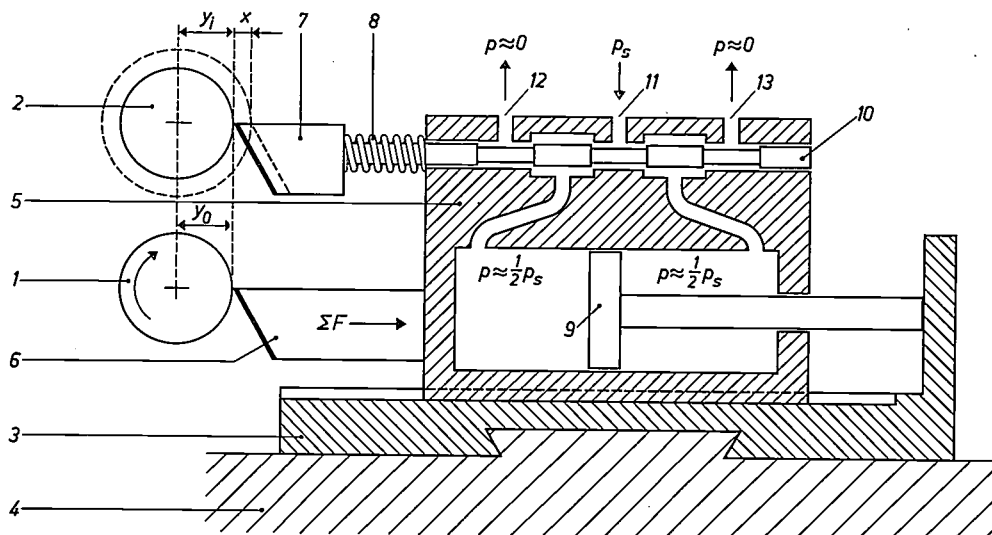


Fig. 1. Cross-section of a copying lathe with hydraulic servomotor. 1 rotating piece of material, which must be made a copy of the template 2. 3 support. 4 bed. 5 transverse slide (and servomotor). 6 tool fixed to the transverse slide. 7 follower, pressed against the template by the spring 8. 9 piston connected to the support by a piston rod. 10 spool of control-valve, which forms a whole with the follower. 11 oil inlet, under constant pressure p_s . 12, 13 oil outlets.

In the unloaded equilibrium state, $y_i = y_0$ as shown in the figure. When the equilibrium is broken (broken lines), y_i is e.g. greater than y_0 , and the follower moves a distance $x = y_i - y_0$ to the right.

and a control-valve 10. The piston is securely fixed to the support. The spool of the control valve is an extension of the follower, and thus moves when the follower moves. The two compartments into which the cylinder is divided by the piston communicate with the oil-supply pipe 11, which is kept at a constant pressure by a pump, and each compartment has its oil outlet (12, 13). When the spool is in its middle position, the oil inlet and outlet (the "ports") of each compartment are equal in size. The volumes of oil flowing through these ports are thus also equal, so that the piston and the cylinder do not move with respect to one another. If the spool moves to the right, the oil inlet of the left-hand compartment is narrowed and the outlet is widened, while just the opposite happens in the right-hand compartment. Less oil is thus supplied to the left-hand compartment, and more is led off, and *vice versa* in the right-hand one. The cylinder (and thus the whole transverse slide) moves to the

right until the (stationary) spool is once more in the middle position with respect to the (moving) cylinder. This explains how the chisel can follow the transverse motion of the probe.

The servomotor drawn in fig. 1 has four controlled ports which regulate the flow of oil. Other types also exist, with two controlled ports, or one. These are sketched in fig. 2a and b. The choice between the various types is largely determined by constructional

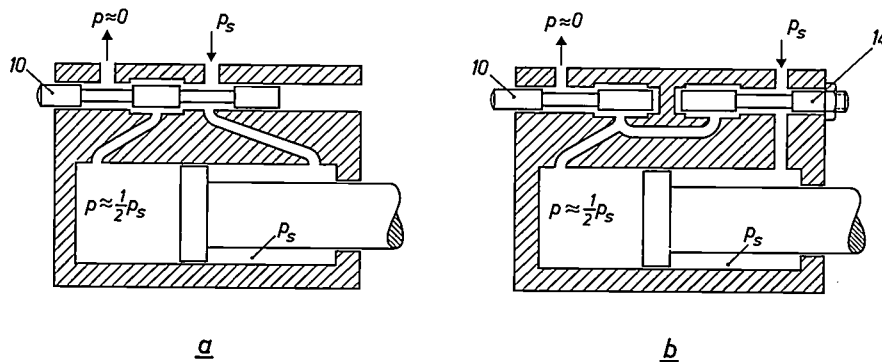


Fig. 2. Apart from the type of hydraulic servomotor shown in fig. 1, which has four controlled ports, there are two other types: a) with two controlled ports, and b) with one controlled and one non-controlled port. In the latter type only the spool 10 moves: spool 14 is fixed in place.

In both these types, the pressure in the cylinder compartment which contains the piston rod is equal to the full pump pressure p_s , while the pressure in the other compartment has half this value, $\frac{1}{2}p_s$, when the load ΣF is zero. Now if the forces on the piston are to be in equilibrium, the piston area on the side of the piston rod must be half that on the other side, i.e. the diameter of the piston rod must be $\frac{1}{2}\sqrt{2}$ times that of the cylinder. This very thick piston rod is thus a characteristic feature of servomotor with either one or two controlled ports. (In servomotors with four controlled ports — see fig. 1 — the pressure on both sides of the piston is $\frac{1}{2}p_s$ at zero load. In this case, if the piston rod is made thin the forces on the two sides of the piston are more or less in equilibrium; the slight deviation from equilibrium can easily be corrected by giving the spool a small constant displacement from the middle position.)

considerations. For our numerically controlled milling machine we used a special version of the type with one controlled port (and one port is not controlled) (fig. 2*b*). For the sake of simplicity, we shall derive the most important properties of hydraulic servomotors¹⁾ for the type with two controlled ports (fig. 2*a*); the three types are so closely related that the considerations given here can be applied practically without change to the other two types.

Hydraulic follow-up system

In the case of the copying lathe, the task of the servomotor is to make the tool (displacement y_0) describe the movement of the follower (y_i) as accurately as possible (i.e. it must act as follow-up system). The difference x between these two displacements:

$$x = y_i - y_0, \quad \dots \dots (1)$$

is equal to the displacement of the spool from its middle position (see fig. 1). Fig. 3 shows how the difference x is formed by means of feedback; y_i can be regarded as the reference signal of the servomotor, x as its input signal and y_0 as its output signal, the whole system forming a closed control loop.

The transverse slide is subjected to various forces, in the first place the reaction of the material on the tool and further the inertial forces (due to acceleration of the slide) and the frictional forces, e.g. between piston and cylinder. All these forces together form the total load ΣF on the transverse slide.

For the moment, we shall completely neglect the load on the slide. In that case, the velocity \dot{y}_0 with which the cylinder moves with respect to the piston is proportional to the control error x :

$$\dot{y}_0 = \frac{1}{\tau_v} x, \quad \dots \dots (2)$$

so that
$$y_0 = \frac{1}{\tau_v} \int x \, dt.$$

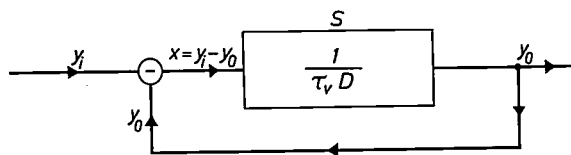


Fig. 3. A servomotor S as an integrating control system. The integrating effect is indicated by the symbol $1/\tau_v D$, where D is the differential operator. y_i is the reference signal, y_0 the output signal, and $x = y_i - y_0$ the input signal (1:1 feedback). In the case of fig. 1, y_i is the dimension of the template to be followed, y_0 the displacement of the tool, and $x = y_i - y_0$ that of the spool from its middle position.

¹⁾ Treated in further detail in: T. J. Viersma, Investigations into the accuracy of hydraulic servomotors, thesis Delft 1961. This thesis also appeared in Philips Res. Repts 16, 507-597, 1961 (No. 6) and 17, 20-78, 1962 (No. 1).

It follows from this last equation that the unloaded servomotor acts as an *integrator*. The quantity τ_v , which has the dimension of time, will be called the integrator constant in this article (in the American literature, it is known as the velocity time constant).

It is known from control theory that in the static case (constant reference signal y_i), an integrating control system can reduce the error x exactly to zero. If however y_i fluctuates, as it does in the copying lathe, x is not zero. The curved line in fig. 4 (both of whose scales are logarithmic) gives the

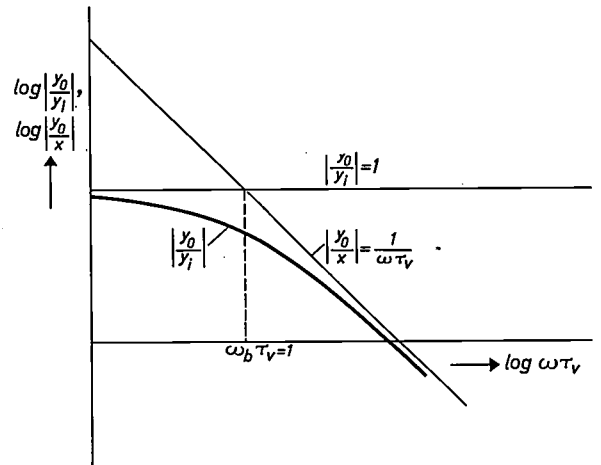


Fig. 4. The curved line gives $\log |y_0/y_i|$ as a function of $\log \omega \tau_v$ for the control system of fig. 3. At low frequencies, this curve tends to the horizontal straight line $|y_0/y_i| = 1$, and at high frequencies to the sloping line $|y_0/y_i| = 1/\omega \tau_v$. These two straight lines intersect at $\omega = \omega_b = 1/\tau_v$; ω_b is called the bandwidth.

variation of the amplitude ratio $|y_0/y_i|$ as a function of $\omega \tau_v$ (where ω is the angular frequency of y_i) for sinusoidally varying y_i . As explained in the caption to fig. 4, at low frequencies this curve tends to the horizontal line $|y_0/y_i| = 1$, while at high frequencies it tends to the sloping line $|y_0/y_i| = 1/\omega \tau_v$. These lines intersect at $\omega \tau_v = 1$. It may be seen from the figure that $|y_0/y_i|$ does not differ much from 1 as long as ω is less than $1/\tau_v$, but decreases rapidly at higher frequencies. $1/\tau_v$ is called the *bandwidth* (ω_b) of the control system. For $\omega \gg \omega_b$, x is approximately equal to $\tau_v \dot{y}_0$, i.e. the control error is proportional to the instantaneous velocity; $\tau_v \dot{y}_0$ is therefore sometimes called the *velocity error*. If ω is greater than ω_b , the control system is no longer able to follow the varying reference signal properly. If the control system has to be able to follow rapid variations well, it must therefore have a large bandwidth, i.e. a small integrator constant τ_v . This conclusion still holds if the loading of the transverse slide is not neglected.

Instead of regarding the performance as a function of the frequency, we can also investigate how the system responds to a discontinuous change in the

reference signal y_i . The values of y_i and y_o in this case are plotted as functions of time t in *fig. 5*. It will be seen that the output signal does not begin to approach anywhere near the new value of the reference signal until a time $t \approx \tau_v$ has passed. This thus happens more quickly as τ_v is smaller.

It follows from the above that it is advisable from all points of view to make the integrator constant τ_v as small as possible: not only is the velocity error $\tau_v \dot{y}_o$ then minimum, but the bandwidth is maximum and the response to a step function is most favourable. We shall shortly discuss the measures which can be taken to reduce the value of τ_v .

The magnitude of the control error x depends not only on the speed with which the cylinder moves, but also on the total load ΣF on the transverse slide (which we have neglected so far). The quantity which determines the influence of ΣF on x is the hydraulic stiffness c , defined as

$$c = \frac{\partial \Sigma F}{\partial y_o} = - \frac{\partial \Sigma F}{\partial x}$$

Since the servosystem is very nearly linear, we may write to a good approximation:

$$c = - \frac{\Sigma F}{x} \dots \dots \dots (3)$$

The total control error is obtained by superposition of the errors according to eqs (2) and (3):

$$x = \tau_v \dot{y}_o - \frac{\Sigma F}{c} \dots \dots \dots (4)$$

As mentioned above, ΣF is the sum of the external force (F_e), the inertial force ($m\dot{y}_o$) and the frictional force (F_f). If we forget about the two latter forces for the moment, we can plot eq. (4) (\dot{y}_o as a function of x , with F_e as a parameter) as shown in *fig. 6a*. However, the situation is more complicated if *dry (Coulomb) friction* (F_f) is present. This type of friction does not depend on the speed with which

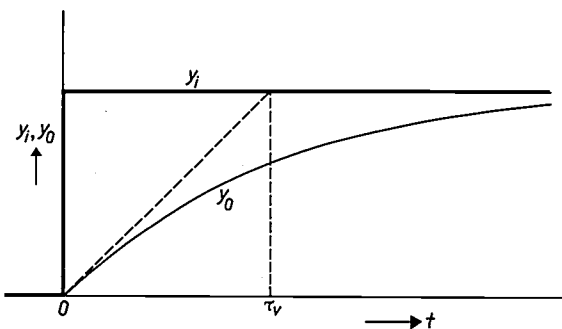


Fig. 5. If the reference signal y_i in the control system of *fig. 3* changes suddenly, the output signal y_o varies as shown here. y_o does not approach within a reasonable distance of the new value of y_i until after a time of the order of τ_v .

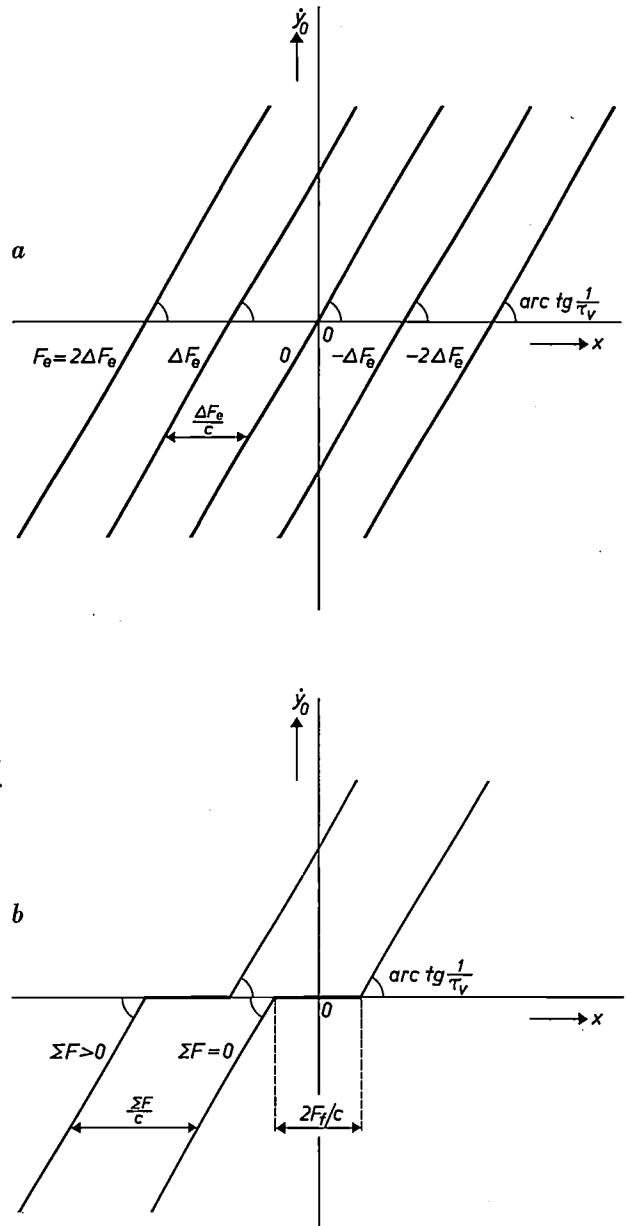


Fig. 6. a) The velocity \dot{y}_o delivered by a servomotor as a function of the input signal x , with the external force F_e as a parameter, if the forces of inertia and friction are neglected. F_e decreases from left to right in equal steps ΔF_e . b) The same, taking the (dry or Coulomb) friction F_f into account. Since F_f changes sign each time the direction of motion is reversed (see *fig. 7*), dead zones of width $2F_f/c$ are produced; these have an unfavourable effect on the operation of the servomotor.

the surfaces in question move over one another, but it changes sign when the velocity v is reversed (*fig. 7*). Thus, whenever \dot{y}_o changes sign, ΣF shows a discontinuous change by a constant amount equal to twice the frictional force F_f ; this means that x also changes suddenly, by $2F_f/c$. If we take the dry friction into consideration, we obtain the curves of *fig. 6b*. This figure shows certain dead zones, where \dot{y}_o is independent of x . Naturally, these dead zones

are extremely undesirable from the control engineer's point of view, since within them displacement of the spool has no effect whatsoever on the motion of the transverse slide.

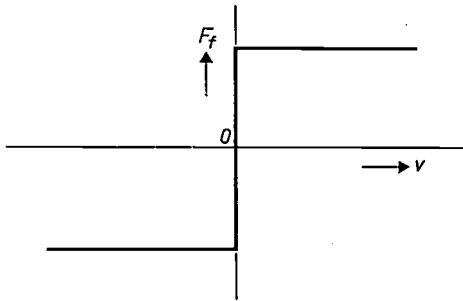


Fig. 7. Variation of the dry (Coulomb) friction. The absolute value of the frictional force F_f is independent of the velocity v with which the surfaces move over one another, but F_f changes sign when the direction of v is reversed.

In order to get an idea of the order of magnitude of these dead zones in modern copying lathes, we performed measurements on twelve such machines which were exhibited at the Techni-Show in Utrecht in 1958. In most cases, the measured values lay between 15 and 20 μm . This means that when these machines are used to copy a cylindrical template, diameter variations of about 35 μm are to be expected on the average! We shall show below how the dead zone in our milling machine has been reduced to about 0.25 μm by means of "load compensation".

Summarizing, we may say that if the control accuracy is to be high, the servomotor must respond rapidly (low integrator constant τ_v) and the dead zone must be small, i.e. the friction must be slight and/or the hydraulic stiffness high. In the following sections we shall discuss how these demands can be satisfied.

High sensitivity through turbulent flow of oil

We shall now derive an expression for the integrator constant τ_v with reference to fig. 8, which shows a hydraulic servomotor with two controlled ports.

If we denote by Q_1 and Q_2 the amounts of oil flowing per second through ports 1 and 2 respectively, then $Q_2 - Q_1$ is the amount of oil which leaves the left-hand compartment per second. If the area of the piston is A , then the velocity \dot{y}_0 is given by:

$$\dot{y}_0 = \frac{Q_2 - Q_1}{A} \dots \dots \dots (5)$$

Now it follows from eq. (2) that $\tau_v = x/\dot{y}_0$, or in the

general (non-linear) case, $\tau_v = \partial x/\partial \dot{y}_0$. Substituting this in (5), we have:

$$\tau_v = \frac{A}{\frac{\partial Q_2}{\partial x} - \frac{\partial Q_1}{\partial x}} \dots \dots \dots (6)$$

Each port consists of an annular slit of circumference πd (see fig. 8) and axial width h_1 . If the oil flows through ports 1 and 2 with velocity v_1 and v_2 respectively, then

$$Q_1 = \pi d h_1 v_1 \text{ and } Q_2 = \pi d h_2 v_2.$$

The flow of oil can be either turbulent or laminar. In the first case, the velocity v of the oil is independent of the width of the port, while in laminar flow v is proportional to the port width.

Let us suppose that the oil flow is turbulent, with velocity v_0 ; then, if the servomotor is unloaded ($\Sigma F = 0$), we have

$$Q_1 = \pi d h_1 v_1 = \pi d (h_0 - x) v_0$$

and

$$Q_2 = \pi d h_2 v_2 = \pi d (h_0 + x) v_0,$$

where h_0 is the port width when the spool is in the middle position and x is the deviation from the middle position. Substituting these expressions in (6) we find:

$$\tau_v = \frac{A}{2\pi d v_0} \dots \dots \dots (7a)$$

If the flow is laminar, v_1 and v_2 are proportional to $h_0 - x$ and $h_0 + x$ respectively. If we put $v_1 = v_2 = v_0$ at $x = 0$, we find, substituting the new expressions for Q_1 and Q_2 in (6), that:

$$\tau_v = \frac{A}{4\pi d v_0} \dots \dots \dots (7b)$$

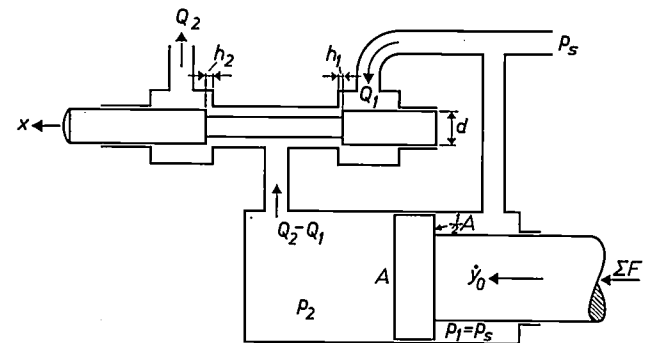


Fig. 8. Hydraulic servomotor with two controlled ports, of width h_1 and h_2 . An amount Q_1 of oil flows per second in through the one, and an amount Q_2 out through the other; the amount of oil in the compartment to the left of the piston thus decreases by $Q_2 - Q_1$ each second. The pressure in the right-hand compartment is equal to the full pump pressure p_s , and that in the left-hand compartment (p_2) is equal to $\frac{1}{2} p_s$ at zero load ($\Sigma F = 0$); this is the reason for the thick piston rod (see caption to fig. 2).

These results also hold for servomotors with *four* controlled ports, as long as the construction is symmetrical. In servomotors with one controlled port and one not controlled, however, the latter has a constant width, with the result that either $\partial Q_1/\partial x$ or $\partial Q_2/\partial x$ is zero; τ_v is therefore twice as big as for the other two types of servomotors.

Summarizing, we may conclude that

$$\tau_v = a \frac{A}{\pi d v_0}, \dots \dots \dots (8)$$

where the factor *a* depends on the nature of the oil flow and on the type of servomotor:

Servomotor with	turbulent flow	laminar flow
4 controlled ports	$a = \frac{1}{2}$	$a = \frac{1}{4}$
2 controlled ports	$a = \frac{1}{2}$	$a = \frac{1}{4}$
1 controlled port	$a = 1$	$a = \frac{1}{2}$

It should not be concluded from this table that one should try to achieve laminar flow in order to make τ_v as small as possible. On the contrary, the fact that *a* is twice as big is more than compensated by the fact that the velocity v_0 can be ten or more times as much with turbulent flow, and as we have seen from eq. (8), τ_v is inversely proportional to v_0 . We must therefore try to obtain turbulent flow through the ports.

Measures to ensure turbulent flow

The geometry of the ports should thus be chosen so that turbulent flow is ensured. There are no data in the literature on flow through annular ports, which are normally used in servomotors, so we started by considering what is known about flow through round and slit-shaped ports in a thin wall (i.e. through a sudden constriction followed immediately by a sudden widening, see fig. 9a and b). The velocity of flow through these ports is given by:

$$v = a \sqrt{\frac{2\Delta p}{\rho}}, \dots \dots \dots (9)$$

where Δp is the pressure difference across the opening, ρ the density of the liquid, and *a* the "discharge coefficient". The value of *a* depends, among other things, on Reynolds' number, *Re*, which plays a role in all phenomena connected with the flow of liquids. (We shall shortly give a formula for *Re*; for the present we shall only remind the reader that *Re* is proportional to the flow velocity *v*.) Wuest has proved that for small values of *Re* the following expression holds ²⁾

$$a = \delta \sqrt{Re} \dots \dots \dots (10)$$

and that flow is then laminar (for which reason the constant δ is known as the laminar flow coefficient). δ depends on the geometry of the aperture, being smaller as the diameter of the liquid column changes

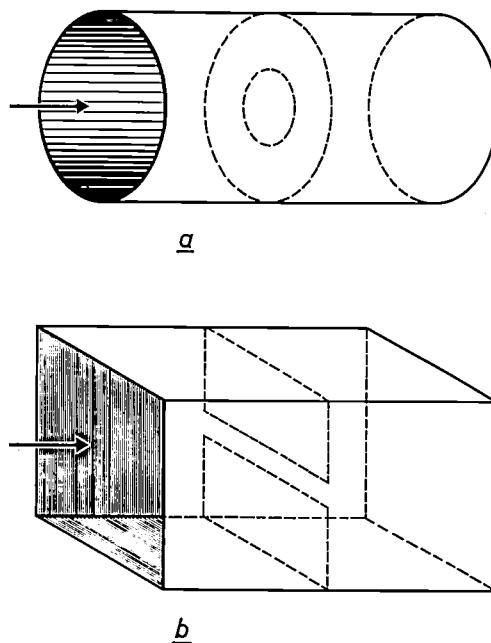


Fig. 9. Two forms of ports for which data on the flow of liquid through them are available in the literature. a) round hole, b) slit, both in an infinitely thin wall at right angles to the direction of flow.

less abruptly at the diaphragm. At large values of *Re*, however, *a* is independent of *Re*. As Von Mises has shown ³⁾, it then has the value:

$$a = \frac{\pi}{\pi + 2} = 0.611 = \alpha_0,$$

and under these conditions flow is turbulent.

The discharge coefficient *a* varies with \sqrt{Re} as shown by the full line of fig. 10. This curve consists of a straight line through the origin given by eq. (10), a horizontal line at $a = \alpha_0$ and a smooth curve joining these two. The two linear parts of the curve intersect at a value of *Re* which we shall denote by \bar{Re} , and which is given by

$$\bar{Re} = \left(\frac{\alpha_0}{\delta}\right)^2 \dots \dots \dots (11)$$

This value of \bar{Re} separates the region of laminar flow (at lower values) from that of turbulent flow

²⁾ W. Wuest, Strömung durch Schlitz- und Lochblenden bei kleinen Reynolds-Zahlen (Flow through slit-shaped and circular diaphragms at low Reynolds numbers; in German), Ing.-Arch. 22, 357-367, 1954.

³⁾ R. von Mises, Berechnung von Ausfluss- und Überfallzahlen (Calculation of efflux and overflow coefficients; in German), Z. Ver. deutsch. Ing. 61, 447-452, 1917.

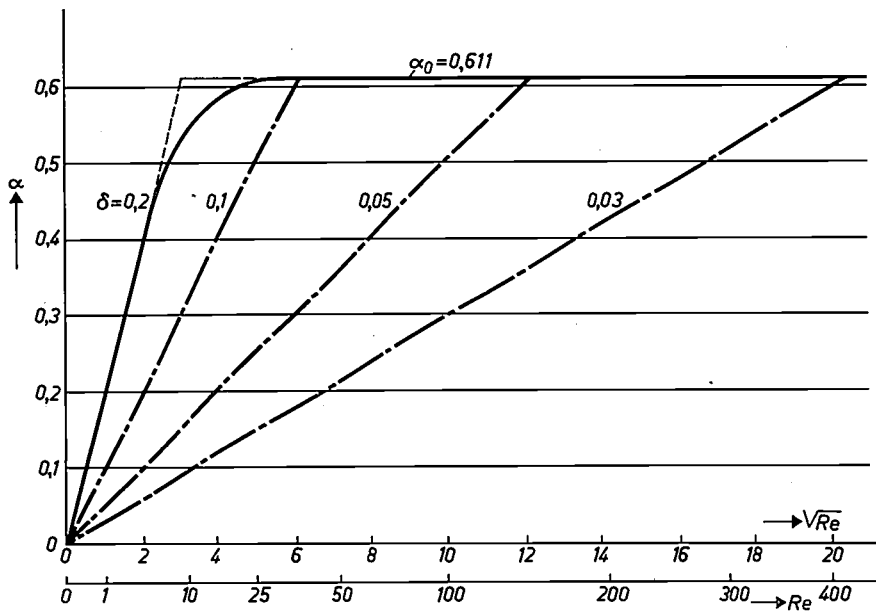


Fig. 10. The discharge coefficient α as a function of \sqrt{Re} ($Re = \text{Reynolds number}$), with the laminar flow coefficient δ as parameter. For laminar flow, α is proportional to \sqrt{Re} , while for turbulent flow $\alpha = \alpha_0 = 0.611$. The transition from laminar to turbulent flow occurs at a lower value of Re as δ increases. For the ports of fig. 9, $\delta = 0.2$.

(at higher values). Using the value $\delta = 0.2$ calculated by Wuest for infinitely thin diaphragms, we find $\bar{Re} = (0.611/0.2)^2 \approx 9$. Johansen⁴) found experimentally: $\delta = 0.17$, $\bar{Re} = 12$. The slight difference between theory and experiment can be explained by the finite thickness of the diaphragms used by Johansen.

The broken lines of fig. 10 show how α varies at values of δ less than 0.2. It will be seen that as δ decreases, laminar flow changes to turbulent at a higher value of \bar{Re} . This is in quantitative agreement with theory, which states that \bar{Re} is about 2300 for an infinitely thick diaphragm.

The experimental confirmation of Wuest's calculations makes it reasonable to expect that eq. (9), (10) and (11) will to a first approximation also apply to flow through an annular slit of the common type shown in fig. 11a, of which a part is shown magnified in fig. 11b. One difference between such an aperture and an infinitely thin diaphragm is that the transition is less abrupt, so that the oil will have a greater tendency to form a boundary layer along the surface, i.e. to laminar flow. This effect is reduced by making the edges past which the oil flows (A and B , fig. 11b) as sharp as possible.

This conclusion is confirmed by experiments which we carried out on a certain hydraulic servomotor, where we determined α indirectly by measuring the integrator constant τ_v . The original value found for

τ_v was 4 millisecond, in good agreement with the value calculated for turbulent flow ($\alpha = \alpha_0 = 0.61$). After some time, during which a large number of experiments were carried out on the motor, it was however found that τ_v had increased to more than 12 ms, corresponding to a decrease of α from 0.61 to about 0.2. It appeared that this change was due to rounding-off of the originally sharp edges by wear. The radius of curvature of the worn edges was about 30 μm while the effective port width was about 50 μm . After a new spool whose edges had a radius of curvature of

less than 10 μm was fitted, the original values of 4 ms and 0.61 were found for τ_v and α . The reason for this is that the wear caused the laminar flow coefficient δ to decrease from about 0.1

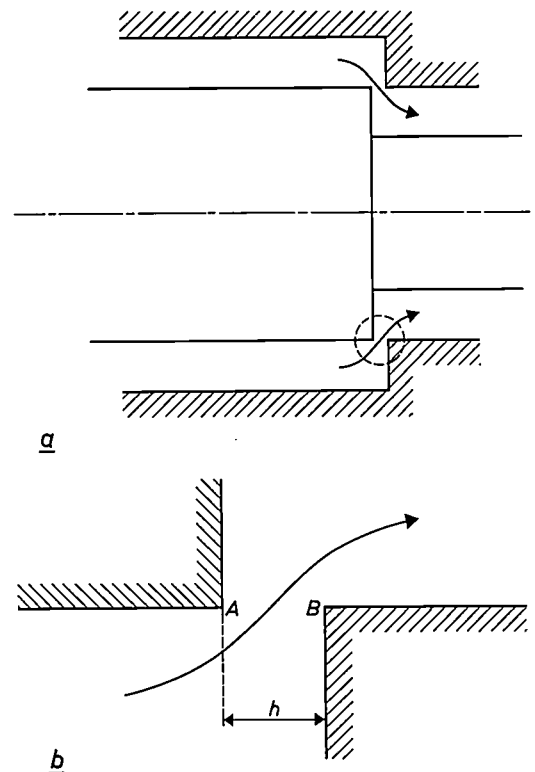


Fig. 11. a) Axial cross-section through the annular type of slit common in hydraulic servomotors. The part inside the dotted circle is shown on a larger scale in (b). The radius of curvature of the sharp edges A and B should be as small as possible, in order to encourage turbulent flow.

⁴) F. C. Johansen, Flow through pipe orifices at low Reynolds numbers, Proc. Roy. Soc. (London) A 126, 231-245, 1930.

to less than 0.03, so that, as may be seen from fig. 10, the value of \overline{Re} at which laminar flow changed over to turbulent flow was increased from 40 to about 400!

Our experiments showed that flow through a port with very sharp edges will certainly be turbulent if $Re > \overline{Re} \approx 20$. This condition can be used to find the minimum value of h , the width of the port: the formula for Reynolds' number is $Re = 2vh/\nu$, where ν is the (kinematic) viscosity of the liquid. Substituting the expression for v of eq. (9) in this equation, we find:

$$Re = \frac{2a_0h}{\nu} \sqrt{\frac{2\Delta p}{\rho}} \geq \overline{Re} = 20, \dots (12)$$

from which the minimum value of h can be calculated.

Temperature sensitivity and oil consumption

We shall close this section with a short remark about a further advantage, and a disadvantage, of turbulent flow.

The oil temperature of a servomotor can increase considerably during operation: an increase of 25 °C is quite normal. This leads to a considerable variation in the kinematic viscosity ν of the oil, and thus of Reynolds' number. As we have seen, α is strongly dependent on Re in laminar flow, but not in turbulent flow. Servomotors with turbulent flow therefore have the important advantage that their operation is not temperature-sensitive.

A disadvantage of a servomotor with turbulent flow is that the flow rate of oil must be high, even when the piston and the cylinder are at rest with respect to each other. The oil pump must therefore have a higher power than in the case of laminar flow. The increased costs brought about by this should be seen as the price which must be paid for increased accuracy; however, these increased costs are only a fraction of the costs of the machine tool for which the servomotor is used.

Increased hydraulic stiffness through load compensation

The load on a hydraulic servomotor has an effect on the oil pressure in the compartments of the cylinder: if a force is applied acting to the left on the piston rod of the servomotor with two controlled ports shown in fig. 8, the pressure p_2 needed to maintain equilibrium will be increased. As we have seen, in the unloaded state this pressure is equal to half the pump pressure p_s ; in the loaded state (total force ΣF), therefore, it will be:

$$p_2 = \frac{1}{2}p_s + \frac{\Sigma F}{A} \dots (13)$$

This imposes limits on the value which ΣF may assume:

$$\Sigma F_{lim} = \pm \frac{1}{2}p_s A,$$

corresponding to $p_2 = p_s$ and $p_2 = 0$. If $p_2 \approx 0$, there is a risk that the air which is always dissolved in the oil at high pressures will be released, or worse still that a vacuum will be produced, leading to cavitation. The system should therefore be designed so that a certain safety margin is always observed, e.g. that the maximum load should never exceed, say, $\frac{2}{3}$ of ΣF_{lim} :

$$\Sigma F_{max} = \frac{2}{3} \Sigma F_{lim} = \frac{1}{3} p_s A \dots (14)$$

Now the amount of oil Q_1 flowing in per second through the one port is given by:

$$Q_1 = \pi d h_1 v_1 = \pi d h_1 a_0 \sqrt{\frac{2\Delta p_1}{\rho}},$$

(see eq. 9), while the flow rate Q_2 out through the other port is:

$$Q_2 = \pi d h_2 v_2 = \pi d h_2 a_0 \sqrt{\frac{2\Delta p_2}{\rho}},$$

where Δp_1 and Δp_2 are the pressure drops across the two ports. In the case shown in fig. 2a, $\Delta p_1 = p_s - p_2$ and $\Delta p_2 = p_2$. It therefore follows from eq. (13) that:

$$\left. \begin{aligned} Q_1 &= \pi d h_1 a_0 \sqrt{\frac{p_s}{\rho}} \sqrt{1 + 2 \frac{\Sigma F}{p_s A}}, \\ Q_2 &= \pi d h_2 a_0 \sqrt{\frac{p_s}{\rho}} \sqrt{1 - 2 \frac{\Sigma F}{p_s A}}. \end{aligned} \right\} (15)$$

These equations show the influence of the load ΣF on the oil flow rates Q_1 and Q_2 , and thus on the relative motion of the piston and the cylinder. In a closed control circuit, this effect is compensated for by a displacement x of the spool. As we have seen above, this displacement gives rise to a control error (for which the load is thus to blame). In order to reduce this error to a minimum, we must thus make $|\partial x / \partial \Sigma F|$ as small as possible, i.e. make $|\partial \Sigma F / \partial x|$ as large as possible. Now $-\partial \Sigma F / \partial x$ is by definition the hydraulic stiffness c , which has already been introduced above. The control error due to a variable load can thus be reduced by making the hydraulic stiffness larger. This also has the effect of making the dead zones narrower (see eq. (4) and fig. 6b). We shall now proceed to calculate c .

Again calling the port width with the spool in the middle position h_0 (so that $h_1 = h_0 - x$ and $h_2 = h_0 + x$), we find from (15), putting $Q_1 = Q_2$:

$$(h_0 - x) \sqrt{1 - 2 \frac{\Sigma F}{p_s A}} = (h_0 + x) \sqrt{1 + 2 \frac{\Sigma F}{p_s A}}$$

or, after a certain amount of manipulation:

$$\frac{\Sigma F}{x} = - \frac{p_s A}{h_0 + \frac{x^2}{h_0}}$$

Now if we design the system (by making h_0 and $p_s A$ sufficiently large) so that x/h_0 is small compared to 1, we find that the hydraulic stiffness is given to a good approximation by:

$$c = - \frac{\partial \Sigma F}{\partial x} \approx \frac{p_s A}{h_0} \dots (16)$$

A similar equation is found for the other types of hydraulic servomotors; but in the case of servomotors with four controlled ports the right-hand side of (16) must be multiplied by 2, while for servomotors with one controlled port it must be multiplied by $\frac{1}{2}$.

As may be seen from (16), the hydraulic stiffness increases as the pump pressure p_s and the cylinder cross-section A are chosen larger, and the port width h_0 smaller. For practical reasons, p_s and A must not exceed certain maximum values, while h_0 may not be too small in connection with the desired turbulent flow (see eq. 12). It is thus not so easy to increase the hydraulic stiffness in this way.

Compensating the loading error through the load itself

A method has been worked out in the Philips Research laboratories for getting round the above-mentioned difficulties, and which in theory allows the hydraulic stiffness to be made infinite. In this method, the displacement of the spool is brought about by the load itself, and not via the feedback.

Fig. 12 shows one way in which this principle can be realized. The servomotor shown in this figure is no longer meant for a copying lathe as in fig. 1, but for a machine tool in which the position of the servo-valve is determined by an electromagnet (14). The current i activating this magnet consists of a DC component I_0 and a variable component $k(y_1 - y_0)$, which is derived in some way from the measured difference between the desired and the actual positions of the slide. The attractive force which the electromagnet exerts on the servo-valve is proportional to i^2 . This force thus consists of a constant and a variable component; if the amplitude $k(y_1 - y_0)$ is small compared to I_0 , the variable component of the force will be very nearly proportional to $y_1 - y_0$. This force is taken up by an elastic membrane (15) to which the spool (10) is fixed; this is the new and essential feature of this system. On the other side of the mem-

brane (to the left in fig. 12) is an oil-filled compartment which is at a pressure p_2 ; this pressure varies linearly with the load ΣF . The load thus has a direct influence on the deflection of the membrane. If the membrane is properly designed this deflection just cancels out the influence of the load.

This can be shown simply as follows. When a force ΣF is applied to the left on the spool, the increase in the pressure p_2 will cause the pressure drop over the outlet to increase and that over the inlet to decrease. One would thus expect the outward flow from the left-hand compartment to exceed the flow into the right-hand one. This effect is however exactly balanced by a slight displacement of the spool to the right under the influence of the increased pressure p_2 . The influence of the load is thus compensated without any change in the "input signal" (the control current i in this case). This means that there is no control error. A change in the load therefore has no influence on $y_1 - y_0$; i.e., the hydraulic stiffness is infinite.

This principle has been embodied in a rather more practical form in the servomotor of the digitally controlled milling machine described in parts I and II of this article. This servomotor is of the type with

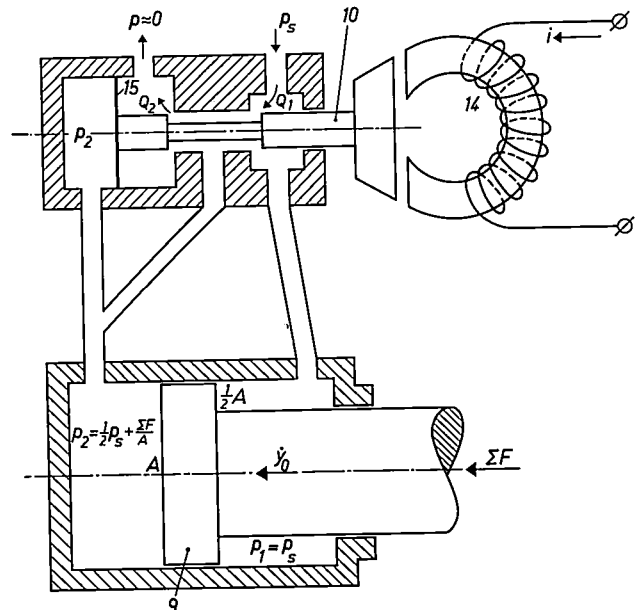


Fig. 12. One method of realizing the principle of load compensation. The spool 10 of the servo-valve is fixed to an elastic membrane 15. The pressure to the left of this membrane is the same as that to the left of the piston 9 (p_2). The pressure p_1 to the right of this piston is equal to the constant pump pressure p_s . When the load ΣF on the piston increases, p_2 also increases; the membrane 15 therefore deflects to the right and displaces the spool to the right. If the membrane is designed properly, this displacement is just enough to compensate for the increase in the load.

In this case, the displacement of the spool proportional to $y_1 - y_0$ is produced by means of an electromagnet 14, the current i through which is also proportional to $y_1 - y_0$.

one controlled port. As can be seen from *fig. 13*, the oil pump is connected directly to one of the compartments of the cylinder. The piston is provided with a chamber 16, which communicates with this compartment via the wide apertures 17 and with the other compartment via the port 18. The width of this port

brane 20 of the servo-valve up, thus opening the port 21 further. The membrane 20 is however designed so that when p_2 increases, the effect of the upwards force on this membrane is compensated by the downwards force on the "shoulder" 22 of the valve 23⁵⁾. The width of port 21 is thus independent of the

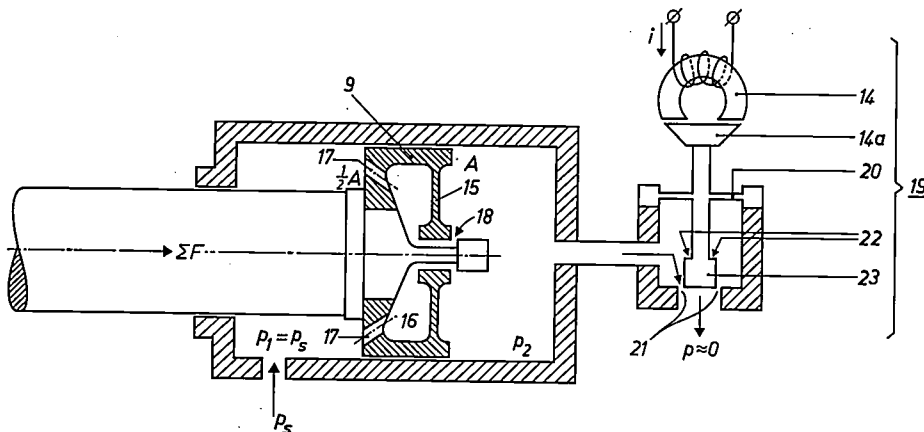


Fig. 13. A more practical way to realize the principle of load compensation, by means of which the dead zone is reduced from 15 to about $0.25 \mu\text{m}$. This type was used in the contour milling machine for surface cams. Membrane 15 (cf. *fig. 12*) is now part of the piston 9, and encloses a chamber 16 inside the piston, which communicates with the left-hand compartment of the cylinder (where the pressure is equal to the pump pressure p_s) via the wide openings 17. Chamber 16 also communicates with the right-hand compartment, via the annular port 18. The width of port 18 depends on the deflection of membrane 15. An electromagnetically controlled servo-valve 19 is placed in the oil outlet; this comprises the electromagnet 14, the armature 14a, the membrane 20 and the valve 23. The width of port 21 depends only on the control current i , which varies linearly with $y_i - y_0$, and not on the oil pressure in the servo-valve 19: the e.g. upwards force on membrane 20 is balanced by a downwards force on the "shoulder" 22 of the valve 23⁵⁾.

depends on the deflection of the membrane 15, which forms part of the piston. A servo-valve 19 is placed in the oil lead back to the sump; the control current i of this valve is again equal to $I_0 + k(y_i - y_0)$.

In the left-hand compartment of the cylinder and in the chamber 16 the oil pressure p_1 is equal to the pump pressure p_s . If we suppose that the piston is subjected to a total load ΣF to the right, then the pressure p_2 in the right-hand compartment will increase linearly with this load. This increased pressure makes more oil flow through the servo-valve, but also causes the membrane 15 to bend a little to the left so that the width of the port 18 increases and more oil flows out of the left-hand compartment despite the decrease in the pressure difference $p_s - p_2$. If the membrane is properly designed, the increased flow of oil through port 18 just balances the increased flow through the servo-valve

It might be thought that the load would also have an influence on the port width of the servo-valve (thus upsetting the compensation), because the increase in the pressure p_2 would press the mem-

pressure p_2 , and therefore also of the load. The only quantity which influences this port width is the control error $y_i - y_0$ (via the current i).

We would remind the reader here of an important constructional characteristic of ports 18 and 21: neither of them contain any metal surfaces which move over one another. There is thus no friction and no wear. *Fig. 14* shows various parts of the servo-valve.

Our method of load compensation can be used in servo-valves of all three types. Complete compensation is only obtained for low loads and small deflections of the membrane. As long as the deflection of the membrane remains small enough for the port 18 to be open, and as long as the total load does not exceed the maximum permissible value ($\Sigma F \leq \Sigma F_{\text{max}}$, see above), this load compensation has a surprisingly good effect. Although the hydraulic stiffness does not actually become infinite, it is made much larger than in the absence of load com-

⁵⁾ Designed by H. L. Günther of Philips' Zentrallaboratorium GmbH, Hamburg Laboratory, at the time at Philips Research Laboratories, Eindhoven.

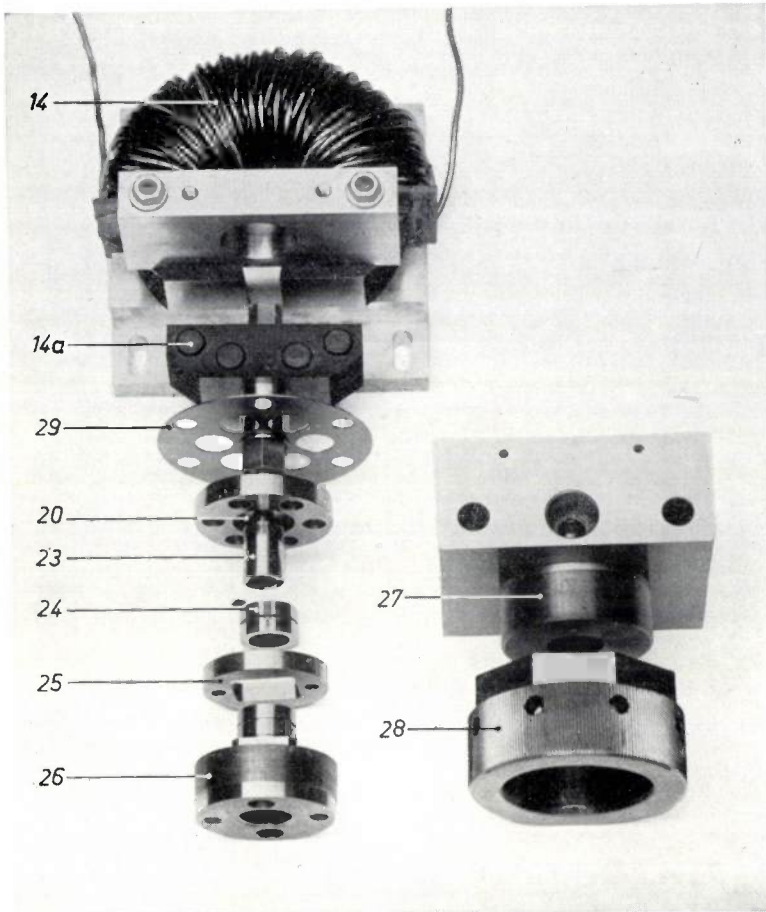


Fig. 14. Exploded view of the servo-valve sketched in fig. 13. The significance of 14, 14a, 20 and 23 is as in fig. 13. The distance between the valve 23 and the seating 24 can be accurately adjusted with the aid of a differential screw (adjustment of the zero point); the screw threads on the support 26 for the seating and the fixed screw 27 differ slightly, so that they move slowly with respect to each other when the grooved nut 28 is turned. (This nut connects 26 with 27). Because 26 may not rotate, it is placed in a square guide 25 which only allows to-and-fro motion. The secondary membrane 29 serves to centre the various parts.

pensation; even under extreme conditions ($\Sigma F = \Sigma F_{\max}$ and $x = h_0$), it is still about ten times larger. It should however be mentioned that no compensation is obtained for rapid variations in the load, because of the compressibility of the oil. This does not matter too much, since the variable component of the load is usually small compared to the constant component.

The results of load compensation (disappearance of the constant control error and of the dead zones) are very important for many applications. For example, the dead zone of machine tools is usually disturbingly large. With a dead zone of e.g. 15 μm , which is not uncommon, it would have been impossible to have limited the control error of the numerically controlled milling machine to 2.5 μm .

Stability through acceleration feedback

As is known, a control system can become unstable under certain circumstances. We must therefore discover when this will happen with the hydraulic servo-system we have been considering.

For this purpose we shall start from eq. (4), which we shall write in a new form:

$$\dot{y}_0 = \left(x + \frac{\Sigma F}{c} \right) \tau_v^{-1} \dots (17)$$

Now ΣF consists of four terms: the external force F_e , the Coulomb friction $-F_f$, the viscous friction $-w\dot{y}_0$ and the inertial force $-m\ddot{y}_0$:

$$\Sigma F = F_e - F_f - w\dot{y}_0 - m\ddot{y}_0.$$

As has been shown elsewhere⁶⁾, two of these terms (the external load and the Coulomb friction) have no influence on the stability; we need not consider them any further, therefore. On the other hand, we must introduce a velocity component into the right-hand side of (17) under dynamic conditions, owing to the compressibility of the oil. This term has the form $\Sigma \dot{F}/c_0$, where c_0 is the stiffness of the column of oil in front of the piston. If we call the length of this column l , its cross-section A and the bulk modulus of the oil E , then c_0 is given by⁷⁾:

$$c_0 = \frac{EA}{l}.$$

Equation (17) must thus be replaced by:

$$\dot{y}_0 = \left(x + \frac{\Sigma F}{c} \right) \tau_v^{-1} + \frac{\Sigma \dot{F}}{c_0}.$$

⁶⁾ T. J. Viersma, Philips Res. Repts 17, 39, 1962 (No. 1).

⁷⁾ This holds for servomotors with one or two controlled ports. Servomotors with four controlled ports have an oil column on both sides of the piston, which makes c_0 roughly twice as great.

Neglecting F_c and F_f , we can to a good approximation rewrite this as:

$$x = \tau_v [\tau_r^2 \ddot{y} + 2\beta\tau_r \dot{y}_0 + \dot{y}_0]. \quad (18)$$

The quantities τ_r and β which appear here are given by

$$\tau_r = \sqrt{\frac{m}{c_0}} \quad (19)$$

and

$$\beta = \frac{\sqrt{mc_0}}{2c\tau_v} + \frac{w}{2\sqrt{mc_0}} \quad (20)$$

τ_r is called the resonance time constant (and is inversely proportional to the frequency of oscillation of the system formed by the moving mass m and the stiffness c_0 of the oil column), and β is called the relative damping. These two quantities, together with the integrator constant τ_v , determine the stability: for the closed control loop of fig. 3 ($x = y_i - y_0$), the condition for stability is:

$$\beta \geq \frac{\tau_r}{2\tau_v} \quad (21)$$

We have shown earlier in this article that for the sake of the accuracy the integrator constant τ_v should be made as small as possible, for which, amongst others, the oil flow should be made turbulent. The stability condition (21) however sets a limit on the reduction of τ_v . Analysis of the frequency characteristics (which we have no space here to describe in any detail⁸⁾) shows that both the stability and the bandwidth ($1/\tau_v$) are satisfactory if $\beta \approx 0.25$ and $\tau_v/\tau_r \approx 4$.

If the damping is too low, there are various means of increasing β . One way would be to decrease the hydraulic stiffness c (see eq. 20). Without going into the possible ways of doing this, we would remark that reducing c is in direct contradiction of the principle derived above, according to which c must be made as large as possible to keep the harmful effects of the load and the dead zones as small as

possible. If load compensation is applied, c becomes so great that the first term on the right-hand side of (20) is negligible compared to the second. One might then try to increase the damping by making the second term bigger, i.e. by increasing the viscous friction.

A solution of a completely different kind is the "acceleration feedback" which has been worked out in the Philips Research Laboratories. In this method, there is feedback of a signal $\gamma\ddot{y}_0$, proportional to the acceleration \ddot{y}_0 (as well as of the displacement y_0 , see fig. 3) to the input of the control system. The term $2\beta\tau_r\dot{y}_0$ in eq. (18) is therefore replaced by $(2\beta\tau_r + \gamma)\dot{y}_0$, so that, if γ is positive, the damping β is apparently increased. It is better to use an electrical signal for this purpose, since there is then no need for any complicated mechanical equipment. For example, a piezo-electric pick-up can be used to provide the signal; this gives a sufficiently strong signal with an acceleration as low as 1% of that due to gravity. The feedback can be arranged by adding a current proportional to this signal to the control current i through the coil of the servo-valve (19 in fig. 12).

This method ensures ample stability in the servo-systems of machine tools, without having to give up the advantages of a high hydraulic stiffness.

The principle of acceleration feedback has been tried out with good success in the servosystem of the contour milling machine described in I and II. However, in this case the viscous friction (see the last term of eq. 20) already provided so much damping that the stability was high enough without this feedback.

Summary. An investigation of hydraulic servomotors has led to the following new conclusions. 1) In order to make the servomotor sensitive, the flow of oil must be made fast (i.e. turbulent). 2) The influence of the load on the control error can be made small by increasing the hydraulic stiffness, which can be done by making the variations in the load itself produce a compensating displacement of the spool of the servo-valve (load compensation). 3) Stability can be guaranteed by bringing the damping up to the desired value by feedback of a signal proportional to the acceleration of the part whose displacement is measured.

Principles 1) and 2) have been made use of in the hydraulic servomotor which forms part of the numerically controlled milling machine developed by Philips.

⁸⁾ T. J. Viersma, Philips Res. Repts 16, 584 (fig. 6.7), 1961 (No. 6).

A SMALL FERROXCUBE AERIAL FOR VHF RECEPTION

by G. SCHIEFER *).

621.396.674.1:621.318.134

In recent years there has been a marked increase in the number of frequency-modulated transmitters broadcasting in the VHF band. As a result, more and more portable radios equipped for VHF are appearing on the market. The relatively long dipole aerial usually needed for good VHF reception with such sets is a drawback which it is now possible to overcome using a new type of ferroxcube.

VHF dipole aeriels

Broadcast, television and innumerable commercial communication transmitters are nowadays making use of the very high frequency band (metre waves). Present-day VHF receivers, compared with those of a few years ago, show a marked advance in electronic engineering, and the sets have become much smaller. Aerials, on the other hand, have not advanced beyond the half-wave dipole¹⁾. The reasons are not difficult to find, for the half-wave dipole has many advantages: in construction it is simple and inexpensive, the real part of the input impedance is of the same order of magnitude as the real part of the input impedance of cable and receiver, the bandwidth is adequate, the losses are low, and the directivity can, if necessary, easily be improved by the addition of passive elements (Yagi aerial). There is one intrinsic drawback, however: in the present trend toward smaller sets the half-wave dipole for metre waves is too long. Nowadays VHF broadcast and radiotelephony receivers of good sensitivity can be made to fit the pocket, but the appertaining dipole aerial must be more than a metre long!

A dipole shorter than a half wave can, it is true, also be used, but in that case the product of efficiency and bandwidth is necessarily smaller, and proper matching calls for the use of coils of exceptionally high quality, which are relatively expensive. Moreover, reception suffers from "proximity effect", for the short dipole reacts mainly to the electrical component of the electromagnetic field and is therefore sensitive to all conductors in its vicinity. Although it is common experience that powerful stations can be received clearly using for an aerial a piece of wire no longer than a finger, such results are more or less fortuitous and bring the technical problem no

nearer to a solution. For these reasons consideration has for some time now been given to the question whether a *small inductive ferroxcube aerial* might give just as good results in the metre wavebands as the familiar *ferroceptor* offers in the medium and long wavebands²⁾.

The following advantages may be expected from a small ferroxcube aerial as compared with a dipole of the same size:

- (1) *Greater product of efficiency and bandwidth*, since the ferromagnetic material concentrates the magnetic field and thus causes a tighter coupling between aerial and field.
- (2) *Smaller proximity effect*, since virtually the only interference with the magnetic component of the field near the aerial can come from ferromagnetic bodies.
- (3) *Absence of "null" directions*. Where the radiation field is horizontally polarized, as it usually is in short-wave broadcasting, an inductive aerial should be set up vertically. In the horizontal plane the aerial then has a circular directivity pattern. Consequently there are no horizontal directions from which the aerial receives nothing (as opposed to electric dipoles for VHF receivers, which, like inductive aeriels for vertically polarized medium and long waves, have directions of zero response (or "nulls"); these can be very troublesome both with fixed and portable receivers).

In spite of these favourable properties, no ferroxcube aeriels for VHF reception have yet been made, for the simple reason that no suitable ferrite was available for frequencies higher than 10 Mc/s. It is only recently that a type of ferroxcube has been developed in the Philips Research Laboratories, Eindhoven, which still has a reasonable permeability

*) Philips' Zentrallaboratorium GmbH, Aachen Laboratory.

¹⁾ See M. Huissoon, The F.M. section of modern broadcast receivers, II. The built-in F.M. aerial, Philips tech. Rev. 17, 348-350, 1955/56.

²⁾ H. Blok and J. J. Rietveld, Inductive aeriels in modern broadcast receivers, Philips tech. Rev. 16, 181-194, 1954/55, in particular page 186 ff.

and sufficiently low losses at 100 Mc/s. This makes it possible to develop small inductive aerials for the broadcast band 11 (87-104 Mc/S). In the following we shall describe an experimental, technically useful ferroxcube aerial which can be mounted in a portable transistor radio of a type now on the market.

The considerations that led to the design of this aerial will be mentioned here only in qualitative terms. The starting point was that small inductive aerials (by "small" we mean considerably shorter than the received wavelengths) have the familiar radiation pattern of the elementary magnetic dipole, the directivity factor of which is, by definition, equal to unity. The directivity factor of a half-wave dipole is then 1.1, in other words scarcely larger. We may conclude, therefore, that an arbitrarily small inductive aerial will give roughly the same performance as the half-wave dipole if it is possible to minimize the losses in the small aerial and in its matching to the receiver. It is precisely this condition, however, that presents the difficulty; for the smaller the aerial, the less tightly is it coupled with the field, and the weaker is the damping of the aerial circuit due to the coupling with the field. In relation to this slight (radiation) damping, the inherent damping of the aerial circuit, due to the losses, is no longer negligible. The result is a low efficiency (see below). The weak damping also results in a small bandwidth and in an extremely high aerial impedance at resonance, which is difficult to match to the receiver input.

An experimental ferroxcube aerial for metre waves

The behaviour of a small inductive aerial, consisting essentially of a conducting loop (resonator), can be studied with reference to the simple equivalent circuit shown in *fig. 1*. Here L is the inductance of the loop, R_r the radiation resistance (which is a measure of the aerial-to-field coupling and can be calculated

from the geometrical dimensions³⁾), R_l is the loss resistance of the circuit, and C the capacitance at which the aerial is tuned to the desired frequency $\omega_0/2\pi$. The aerial impedance is matched to that of the receiver by a loose coupling with a second loop, L' . The relative bandwidth $\Delta\omega/\omega_0$ of the aerial matched to a receiver is then:

$$\frac{\Delta\omega}{\omega_0} = 2 \frac{R_r + R_l}{\omega_0 L}$$

and the efficiency:

$$\eta = \frac{R_r}{R_r + R_l}$$

To obtain a high value of η it is therefore necessary to make R_r large and R_l small. A wider bandwidth, then, can only be achieved by reducing L .

With given loop dimensions, there is only one way of increasing R_r , that is to place inside the loop a ferroxcube rod, which concentrates the magnetic lines of force in the loop. Preferably the permeability μ of the ferroxcube and the length-to-diameter ratio of the rod should be as high as possible⁴⁾. These measures admittedly increase the inductance — which ought to be small for a wide bandwidth — but only in proportion to μ , whereas R_r , as can be demonstrated, increases in proportion to μ^2 . The value of R_l is mainly governed by the loss factor $\tan \delta$ of the ferroxcube; the conduction and dielectric losses in the ferroxcube are negligible, so that $R_l = \omega_0 L \tan \delta$. It is important, then, to choose a ferrite with low losses and to give the loop a small inductance. The latter can be done by suitable design of the loop. Optimum results are obtained if the loop is formed from copper strip which is nearly as broad as the rod is long; the loop thus has the form of a cylinder cut open along its length. Ideally, the tuning capacitance C should be uniformly distributed along the gap, but good results can be achieved with five or six separate capacitors.

Fig. 2 shows a ferroxcube aerial for 100 Mc/s designed on these lines. Note the resonator with the cylindrical loop a and the six capacitors b disposed along the gap. The ferroxcube rod is in two parts, c_1 and c_2 , both of which project by an amount c' outside the resonator; the optimum length of c' depends on the shape of the rod and on the permeability of the ferroxcube. The rod has to be divided into two parts to make room for the coupling loop d , the size of which determines the impedance at the terminals of d .

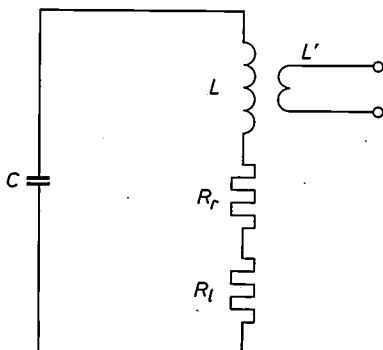


Fig. 1. Equivalent circuit of an inductive aerial. L inductance, R_r radiation resistance, and R_l loss resistance of the loop acting as aerial. C tuning capacitance. L' inductance of coupling loop.

³⁾ See e.g. H. Zuhrt, *Elektromagnetische Strahlungsfelder*, Springer, Berlin 1953, p. 257 ff.

⁴⁾ R. M. Bozorth and D. M. Chapin, Demagnetizing factors of rods, *J. appl. Phys.* 13, 320-326, 1942.

The following are the principal data of this aerial.

Total length of ferroxcube rod	16 cm.
Diameter of ferroxcube rod	2 cm.
Length of resonator	14 cm.
Relative permeability of the ferroxcube . . . approx.	25,
Loss factor $\tan \delta$ of the ferroxcube . . . approx.	0,01,
3dB bandwidth upon matching	1 Mc/s,
Efficiency	} $\begin{matrix} 5\% \\ -13 \text{ dB.} \end{matrix}$

Fig. 3 shows how the aerial is integrated with the radiofrequency block of a receiver for 87-104 Mc/s. The circuit components are mounted directly on the resonator, opposite the gap, thus giving a highly compact construction. Since the bandwidth of the aerial is much smaller than the frequency range mentioned, tuning facilities are necessary. An effective inductive method of tuning is to make one part of the ferroxcube rod capable of axial movement in relation to the other, thus producing a variable air gap between them. This makes the effective permeability of the rod variable, and hence the frequency to which the resonator is tuned. (The alternative, i.e. capacitive tuning, is not so attractive because of the fixed arrangement of the capacitors.) The aerial must be tuned in unison with the RF circuit and (except

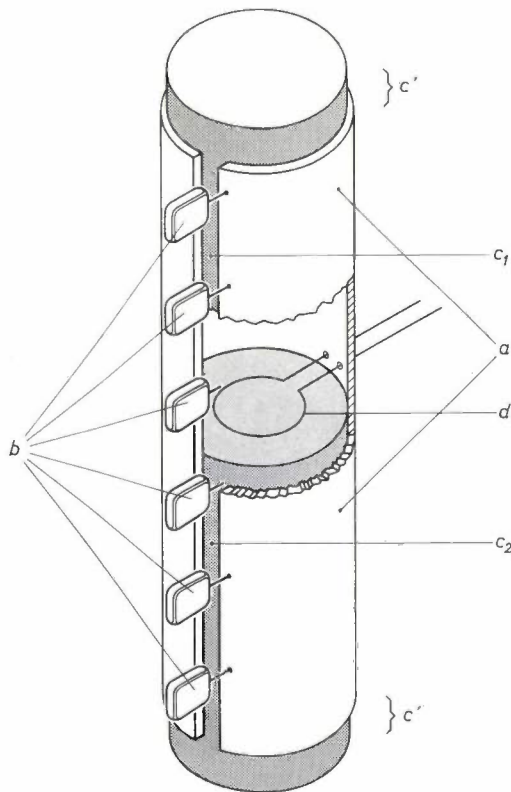


Fig. 2. Experimental ferroxcube aerial for VHF reception. *a* resonator (loop of copper strip); *b* fixed tuning-capacitors disposed along the gap; *c*₁, *c*₂ the two parts into which the ferroxcube rod is divided to leave place for the coupling loop *d*. A length *c'* of the ferroxcube rod projects outside the resonator at both ends.

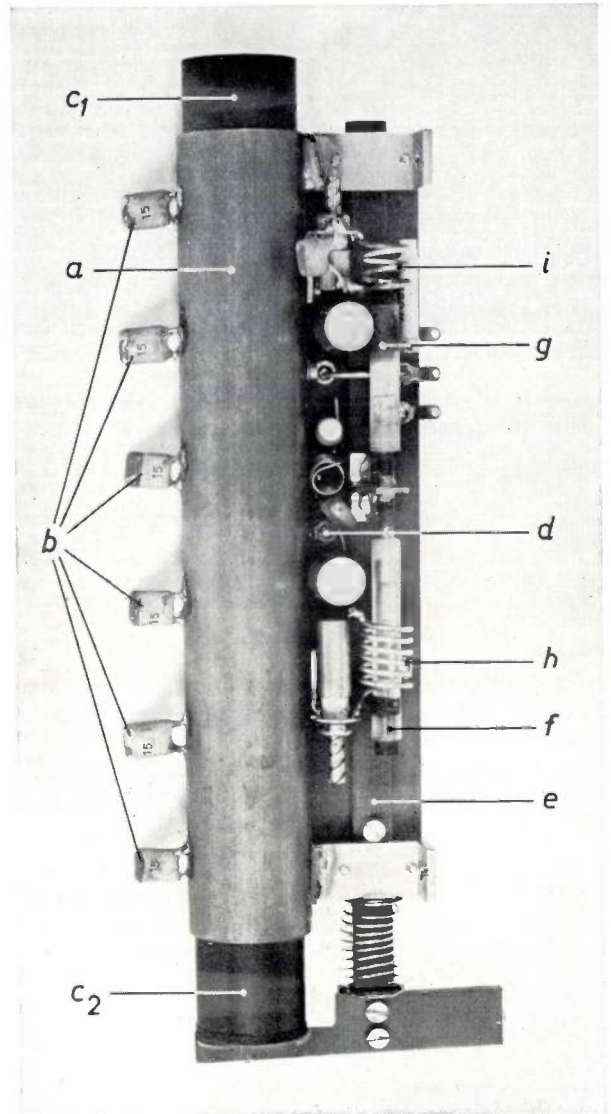
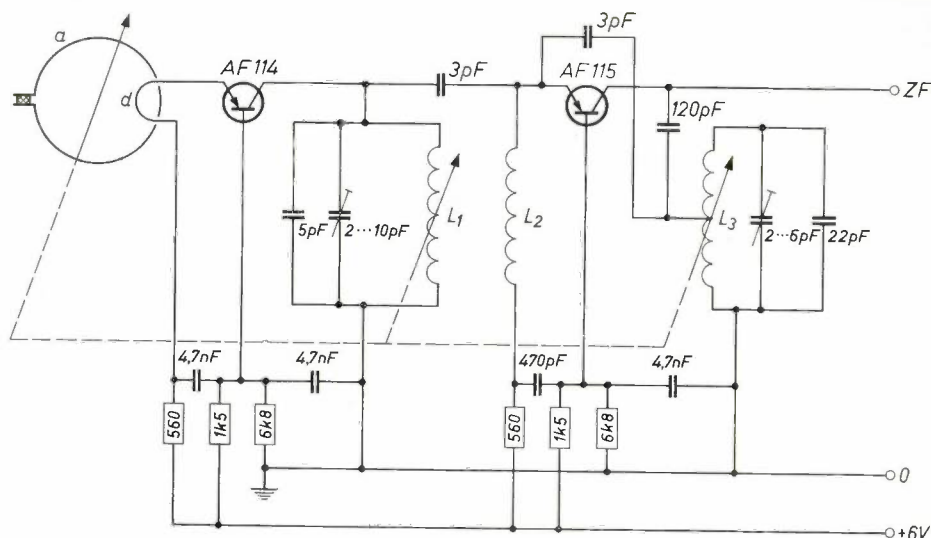


Fig. 3. The ferroxcube aerial integrated with the radio-frequency block of a receiver for 87-104 Mc/s. For the meaning of *a*, *b*, *c*₁, *c*₂ and *d*, see fig. 2. The sliding tuning arm *e* moves part *c*₂ in relation to *c*₁; it also moves the ferroxcube rod *f* relative to the coil *h* of the aerial circuit, and the ferroxcube rod *g* relative to the coil *i* of the local oscillator. Given proper alignment, the tuning of the three circuits is synchronous. See also fig. 4.

for a constant amount) with the local oscillator. This is primarily a mechanical problem. It is solved in this case by causing the arm *e* (fig. 3), which moves one part of the ferroxcube rod, to move two other ferroxcube rods (*f* and *g*) which inductively alter the tuning of the RF circuit *h* and that of the oscillator circuit *i*. With careful adjustment the deviations in synchronism can be kept so small that they remain within the bandwidth of the aerial and of the RF circuit.

The diagram of the radio-frequency block is shown in fig. 4. It consists of an amplifier using a type AF 114 transistor, and a self-oscillating mixer, using a type AF 115 transistor. The diagram otherwise shows no details of interest, but it is worth mentioning that

Fig. 4. Diagram of the radio-frequency block represented in fig. 3. Transistor AF114 operates as RF amplifier; transistor AF 115 is part of the oscillator-mixer stage. The intermediate-frequency amplifier is connected to the terminal ZF. *a* aerial loop, *d* coupling loop. The coils L_1 and L_3 are denoted in fig. 3 by *h* and *i*, respectively.



the circuitry is designed to allow extremely short connections. This makes it possible to raise the gain per stage without endangering the stability.

Fig. 5 shows the interior of a commercially obtainable portable transistor radio, in which the RF block has been replaced by that represented in fig. 4.

The tuning arm, via a cam and lever, is driven from the shaft of the variable capacitor (which is only used for medium- and long-wave reception). It would, of course, easily have been possible to provide a separate drive for the VHF tuning, as is often required.

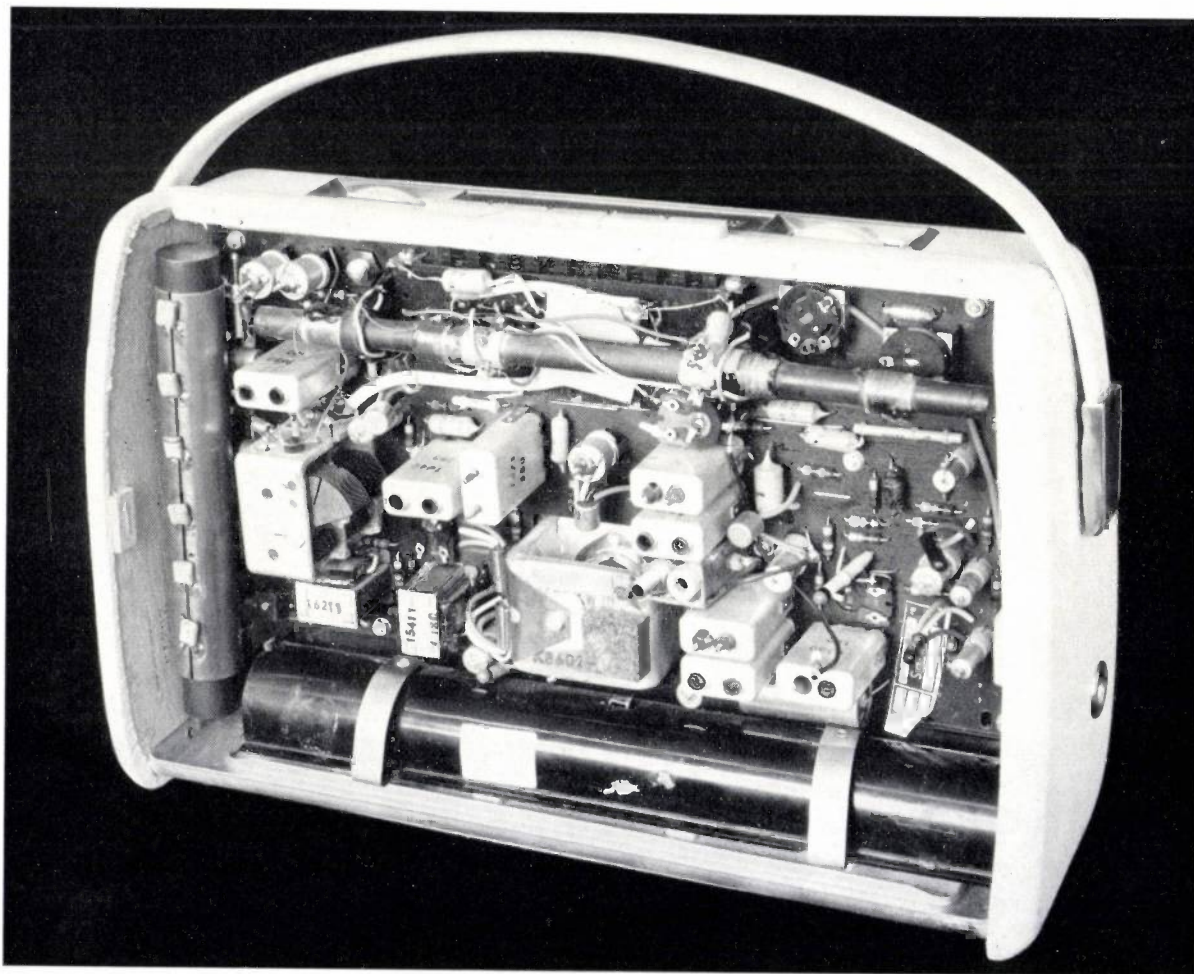


Fig. 5. A commercial transistor-radio (Philips "Annette") fitted experimentally with a ferroxcube aerial for 87-104 Mc/s of the kind shown in fig. 3. The aerial is mounted vertically on the extreme left. The horizontal rod (near the top) is a ferroceptor aerial for medium and long waves.

Results

The reception properties of this set have been compared with those of a similar type of set fitted with the usual half-wave dipole, extendable on two sides to a length of 2×65 cm. A signal generator (with constant modulation) was used to determine the field strength required by each of the receivers to provide the signal-to-noise ratio of 26 dB necessary for good reception. The measurements proved to be exceptionally difficult, since during the day almost the entire band is crowded with radio broadcasting transmissions; measurements at night, when the interference level is much lower, would not have given a true picture. Nevertheless, it was found that the receiver with the ferroxcube aerial needed a stronger field than the receiver with the dipole aerial, but that the difference, depending on the wavelength, was not greater than about 7 to 10 dB. In this connection it should not be forgotten that the (vertical) inductive aerial received equally well from all directions, whereas the dipole was turned in its preferred direction. Given the same signal-to-noise ratio, a dipole less favourable oriented might have needed a stronger field than the ferroxcube aerial.

This result indicates that the small ferroxcube aerial proved more satisfactory in the test than was to be expected from its efficiency of -13 dB. One reason for this is the fact that mass-produced dipole aerials have a not very wide bandwidth, and that their matching is poor, particularly at the limits of the frequency band. A second reason is that, at these frequencies, external noise sources still produce noise contributions which, in relation to the signal, are the same for both aerials. This suggests that an input stage using a low-noise transistor, e.g. type AF 102, would give even better results.

Summary. VHF receivers can nowadays be made to fit the pocket, but they need a dipole aerial more than a meter long, unless a smaller product of efficiency and bandwidth is acceptable, and if the reception is not to suffer from "proximity effect". Compared with a dipole of the same size, a small inductive aerial, using a ferroxcube rod, may be expected to offer the following advantages: a higher product of efficiency and bandwidth, lower proximity effect and reception from all directions (i.e. no horizontal "nulls"). The design of such an aerial has awaited the development of a type of ferroxcube which, at 100 Mc/s, still possesses a reasonable permeability and low losses, ($\mu_r \approx 25$, $\tan \delta \approx 0.01$). A ferroxcube of this kind is now available on a laboratory scale. The article describes an experimental ferroxcube aerial for 87-104 Mc/s, suitable for mounting in a normal, portable transistor radio. Although the sensitivity is lower than that of a half-wave dipole, it is shown to possess the expected advantages.

STUDYING THE INFLUENCE OF LIGHT ON PLANTS



Photo Maurice Broomfield

The effects on plants of such factors as the colour and intensity of the light and the duration of light and dark (daylength effect) are studied on plants growing in special

greenhouses¹⁾ in which an artificial and entirely reproducible climate prevails. The spectral composition of the light — supplied by the tubular fluorescent lamps on both sides — can be adjusted within wide limits; the intensity and duration of irradiation can be fully controlled, as also can the temperature, humidity and CO₂ content of the air in the greenhouse and the composition of the soil. The photo shows a greenhouse in which tomatoes are being grown.

¹⁾ R. van der Veen, A small greenhouse with artificial lighting for studying plant growth under reproducible conditions, Philips tech. Rev. 12, 1-5, 1950.

ABSTRACTS OF RECENT SCIENTIFIC PUBLICATIONS BY THE STAFF OF N.V. PHILIPS' GLOEILAMPENFABRIEKEN

Reprints of those papers not marked with an asterisk * can be obtained free of charge upon application to the Philips Research Laboratories, Eindhoven, Netherlands, where a limited number of reprints are available for distribution.

3046: P. Jongenburger: The influence of plastic deformation on the transverse magnetoresistance of polycrystalline copper, silver and gold (Acta metallurgica **9**, 985-991, 1961, No. 11).

The transverse magnetoresistance $\Delta\rho/\rho_0$, i.e. the relative change of resistivity, of plastically deformed copper, silver and gold wires, caused by the application of a transverse magnetic field, is measured at 20 °K and at 14 °K. When $\Delta\rho/\rho_0$ is plotted versus the magnetic field strength divided by the resistivity in the absence of the magnetic field, the result according to Kohler should be curves which are independent of temperature and purity; this is indeed found. It is now shown that the position of the curves does depend on the degree of plastic deformation of the wires, shifts being found towards both larger and smaller values of $\Delta\rho/\rho_0$. In many cases, annealing below 150-200 °C is found to have little influence on the position of the curves, but annealing at higher temperature causes the curves to return to the position of the curve for non-deformed material. In the case of silver the recovery stage at -120 °C causes an extra displacement which is apparently independent of temperature and likewise disappears upon annealing. These effects can be brought into relation with the dislocations in the plastically deformed material, and may probably be bound up with the anisotropic conduction-electron scattering caused by the dislocations. A satisfactory explanation, however, has not yet been found.

3047: H. G. Grimmeiss, A. Rabenau and H. Koelmans: Some properties of *P-N* junctions in GaP (J. appl. Phys. **32**, 2123-2127, 1961, suppl. to No. 10).

The authors discuss the preparation and some semiconducting properties of single crystals of GaP doped with Zn. This investigation runs parallel with that on AlP doped with Zn, published under No. A 32. See also **3030**.

3048: W. Albers, C. Haas, H. J. Vink and J. D. Wasscher: Investigations on SnS (J. appl. Phys. **32**, 2220-2225, 1961, suppl. to No. 10).

See **R 429**.

3049: H. J. van Daal, C. A. A. J. Greebe, W. F. Knippenberg and H. J. Vink: Investigations on silicon carbide (J. appl. Phys. **32**, 2225-2233, 1961, suppl. to No. 10).

Provisional summary of the investigations on the semiconductor SiC, already dealt with to some extent in **R 392** and **R 434**.

3050: J. C. Balder: Toepassing van tunneldiodes in niet-lineaire schakelingen (T. Ned. Radio-genootschap **26**, 167-179, 1961, No. 4). (Use of tunnel diodes in non-linear networks; in Dutch.)

Paper on the possible uses of tunnel diodes. After deriving the conditions for non-linear behaviour and the quality factor for the switching time, the author discusses various circuits incorporating tunnel diodes: bistable and monostable trigger circuits, the relaxation oscillator and the tunnel-diode pair. In conclusion, a three-phase supply circuit for tunnel diodes is discussed.

Philips Technical Review

DEALING WITH TECHNICAL PROBLEMS
RELATING TO THE PRODUCTS, PROCESSES AND INVESTIGATIONS OF
THE PHILIPS INDUSTRIES



Philips Research Laboratories in the Netherlands are on the move. A new complex of buildings is rising in Waalre, near Eindhoven, to replace the cramped quarters where the Laboratories, growing continuously since their inception about fifty years ago, have been housed until now. In honour of the opening of the first large section of the new complex, a Symposium attended by a large number of distinguished guests, representatives of research institutions in many countries, was held on September 26/27, 1963. At the Symposium, three main lectures and a few general demonstrations

were given and some fifty more-specialized papers read, all by members of our Laboratories.

The present issue of our Review is entirely devoted to these important events in the history of our Laboratories. It contains the introductory address to the Symposium by Prof. Casimir who directs all research activities of the Philips Industries; furthermore, the complete text of the main lectures and summaries of the demonstrations and short papers; finally, an article on the new Laboratory buildings — part of which is shown in the above photograph.

GENERAL INTRODUCTION TO THE SYMPOSIUM AT THE PHILIPS RESEARCH LABORATORIES, EINDHOVEN, SEPT. 26/27, 1963

by H. B. G. CASIMIR *).

The development of science and technology does not follow one single simple pattern but is full of sudden turns and subtle ironical twists. New ideas may be the result of years of painstaking labour or they may arise unexpectedly. They may at once lead to widespread activity or they may lie latent for shorter or longer periods until technology catches up or until they can be incorporated into some other discipline. Sometimes one group, one man even, may pioneer an entirely new line of research; more often similar work is carried out in several centres. Apparently unrelated lines of thought may suddenly become connected and their combination may yield a rich harvest; and what at first appeared to be an insignificant side-line may well become the main issue. We in the Philips laboratories are well aware of all this; perhaps, working as we do in a small country but for a firm with worldwide activities, we are even more keenly aware of the interdependence of our work with that of others than some of our colleagues elsewhere. We have been happy and proud when it has been our privilege to start out alone on a new venture, as in the case of ferrites, or when we were the first to revive almost forgotten ideas, as in the case of the Stirling engine. But we were equally glad when we found that we were no longer alone and that others had adopted our techniques. As a rule our work forms part of research carried out in many different institutions; we have usually been able to make relevant contributions and we have enjoyed sporting competition between laboratories. Even in those cases where we could add little of our own and had mainly to learn and to copy, we could still admire the skill and ingenuity of others. And of course we have always been conscious of our indebtedness to the great scientists in academic institutions. Although in some fields technology may have been ahead of fundamental science, this has rarely been the case in electronics where fundamental ideas were usually formulated several decades before they were applied in industry. Let me quote

a few examples. Fundamental experiments on microwaves were carried out by Hertz in 1888 and resonant cavities were familiar to theorists by the end of the century; the technique of microwaves which incorporated these ideas was developed in the nineteen-forties. The idea of an optical one-way system based on Faraday rotation was proposed by Lord Rayleigh in 1885; the first microwave gyrator which is essentially the realization of this idea for centimetric waves was announced by Hogan in 1952; optical one-way systems based on Faraday rotation are now slowly finding application in connection with optical masers. The relativistic equations of motion of electrons were formulated by Lorentz and Einstein at the beginning of this century; they found hardly any technical application until the advent of betatrons, synchrotrons and linear accelerators. The notion of stimulated emission was introduced by Einstein in 1917; the first maser was made by Townes in 1955. Even in the case of semiconductors two decades elapsed between the introduction of the idea of a positive hole and the discovery of hole survival and the invention of the transistor. There never was *one* Prometheus nor *one* blazing torch but again and again great scientists have wrested from the gods sparks of divine fire which through the care and the toil and labour of many others were kindled into brilliant light.

But on the present occasion we are so glad to be able to receive so many of our good friends, we are so glad to have an opportunity to show our beautiful new laboratories and to explain some of the work we have been doing that, like a child with a new toy, we shall be more than usually egotistical: you will expect us to speak about our own work and our own results, you will perhaps forgive us if we present these results in the context of our own way of thinking. The anthology we offer must needs be incomplete. I hope it will suffice to show that we have done our share of toil and labour to make the lights burn. Perhaps we have occasionally even been able to go out on our own and to capture a few sparks of the original fire. But that is for you to judge.

*) Board of Management, N.V. Philips' Gloeilampenfabrieken.

SOME MAIN LINES OF 50 YEARS OF PHILIPS RESEARCH IN PHYSICS

by H. B. G. CASIMIR.

001.891.53

The field of gas discharges had already a honourable past when it was tackled by Holst and his associates as one of their first major topics ¹⁾ *). It had led to the discovery — or as we Dutchmen prefer to say, to convincing experimental proof — of the electron by J. J. Thomson and to the discovery of X-rays by W. C. Röntgen. *Fig. 1* shows one of the earliest

balance of a sodium lamp by a coating of suitably processed tin oxide which is transparent to visible light but was found to have a high reflectivity in the infrared ⁴⁾. In the late thirties interest in the basic aspects of gas discharges was somewhat on the decline but the rapidly growing field of experimental nuclear physics made good use of ion sources and Geiger-Müller tubes. Even today G.M. tubes have by no means been ousted by photomultipliers-cum-scintillators and solid state devices. As a matter of fact satisfactory stability and reproducibility were only obtained fairly recently ⁵⁾. I show part of the range that is at present manufactured by our factories (*fig. 3*). It was mainly developed by a small branch of our laboratories in the cyclotron institute at Amsterdam. Good fun, building cyclotrons; our latest design of isochronous cyclotron, of which we are rather proud, is one of the items of our Symposium programme ⁶⁾. I also would like to mention scalers, mainly because of the display tube used in many of them, a deflection tube with ten stable positions ⁷⁾ which gives me the opportunity to say a few words to commemorate Dr. J. L. H. Jonker, in whose group this tube was developed, whom many of you have known and to whom the Laboratories are indebted for creating a fine team and a fine tradition. He left us to join the new Technical University here at Eindhoven and was the inspiring leader of its department of electrical engineering

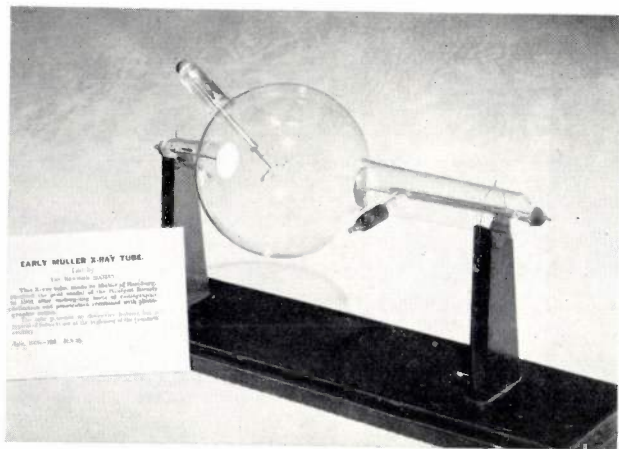
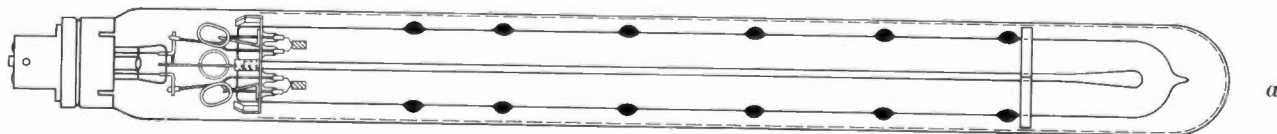
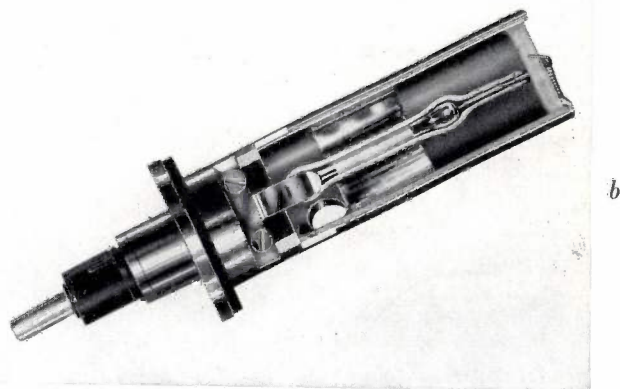


Fig. 1. One of the earliest commercial X-ray tubes, made by Messrs. C. H. F. Müller in 1901.

commercial X-ray tubes, made by Messrs. C. H. F. Müller, a firm that has since become a valuable member of the Philips group. Bohr's theory provided better understanding and hence also better possibilities for technical applications. I have no time to dwell on scientific results of this early work in the Philips Laboratories but should like to mention sodium lamps ²⁾ and super-high-pressure mercury lamps ³⁾ as spectacular technical results (*fig. 2a* and *b*). We have recently been able to improve the energy

*) Literature references are listed at the end of this article.

Fig. 2. a) A modern sodium lamp of high efficiency, obtained i.a. by an infrared reflecting coating of tin oxide (indicated by broken line).
b) A 1000-W water-cooled super-high-pressure mercury lamp (SP 1000 W). This lamp is fitted in a metal reflector enlarging its effective width.



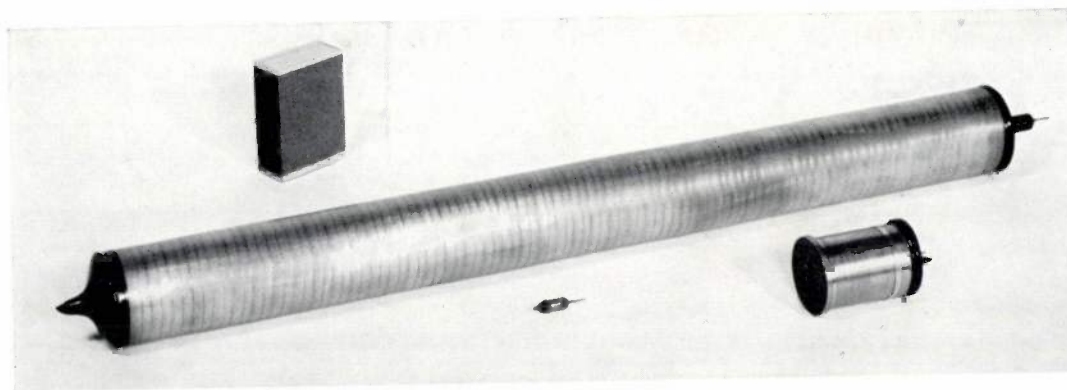


Fig. 3. Three Geiger-Müller counters of different sizes.

until his untimely death this summer. The EIT tube ⁷⁾ is the only surviving member of a family of tubes designed for switching and logical circuits which were made obsolete by the transistor even before they were thoroughly developed.

Back to gas discharges. Radar brought T.R. switches and hydrogen thyratrons and today there is a terrific activity in connection with gaseous lasers, possibilities of controlled fusion, direct conversion of heat into electrical energy, microwave devices and so on, while the older topics are by no means dead. Also, astronomers who in their thinking are so often ahead of us earthbound mortals have started an intensive study of magnetohydrodynamics. Today it is no longer difficult to convince a young man that gas discharges is an exciting subject. Of course you have to call it plasma physics, but as a father of four daughters I have to make worse concessions to fashion. We do not work at present on very high current discharges but we shall show work on several phenomena and devices, including sealed-off neutron generators ⁸⁾, beam plasma interactions ⁹⁾ and optical lasers ¹⁰⁾.

One interesting feature of laser research is that it offers much scope for ingenious optical arrangements; some examples will be shown. Laser modes always remind me of cascade generators. Why? Because they are pseudo-stationary modes in the sense that the wave motion is not only dying out because of energy dissipation but also because waves are leaking away into free space. Such situations are always rather tricky to treat. The most famous example is Gamow's theory of alpha radio-activity (fig. 4). An alpha particle is inside a potential well but can leak out through the barrier. Even in this simple case one had some difficulty to arrive at a satisfactory mathematical treatment. You will remember that the notion of penetration through a barrier encouraged Cockcroft and Walton to proceed with their experiments. Soon afterwards we began

to produce high voltage generators industrially. Fig. 5 shows the first generator we installed at Cambridge and which I think had a satisfactory record ¹¹⁾.

One of our outstanding men in gas discharges was the late Dr. Penning, and one of his contributions made in 1937 was the Penning — or Philips — ionization gauge known to the initiated as „PIG” ¹²⁾. I show you an early model (fig. 6). Later the geometry was modified, but the essential idea remained the same, viz, the use of crossed electric and magnetic fields such that electrons cannot reach the anode from the cathode unless they suffer collisions. Penning emphasized that the device may be used as a manometer — as such it soon found widespread application in our own factories —, as an ion source — he made the first sealed-off neutron tube ¹³⁾ —, and as a pump. I remember clearly the talk he gave in 1936 to the Dutch Physical Society in Amsterdam. For his demonstrations only a rather poor pump was available, but he closed a valve and the manometer did the pumping. I was then a young theoretician reared on quantum mechanics and so the notion that a measurement influences the situation was familiar, but that it should *improve* the situation came rather as a surprise. Today „PIG's” are widely used. One of our competitors has kindly advertized its advant-

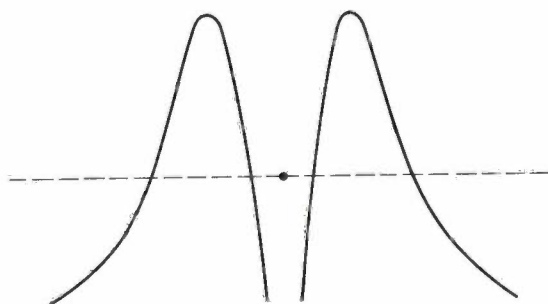


Fig. 4. Symbolic representation of Gamow tunnelling through a potential barrier.

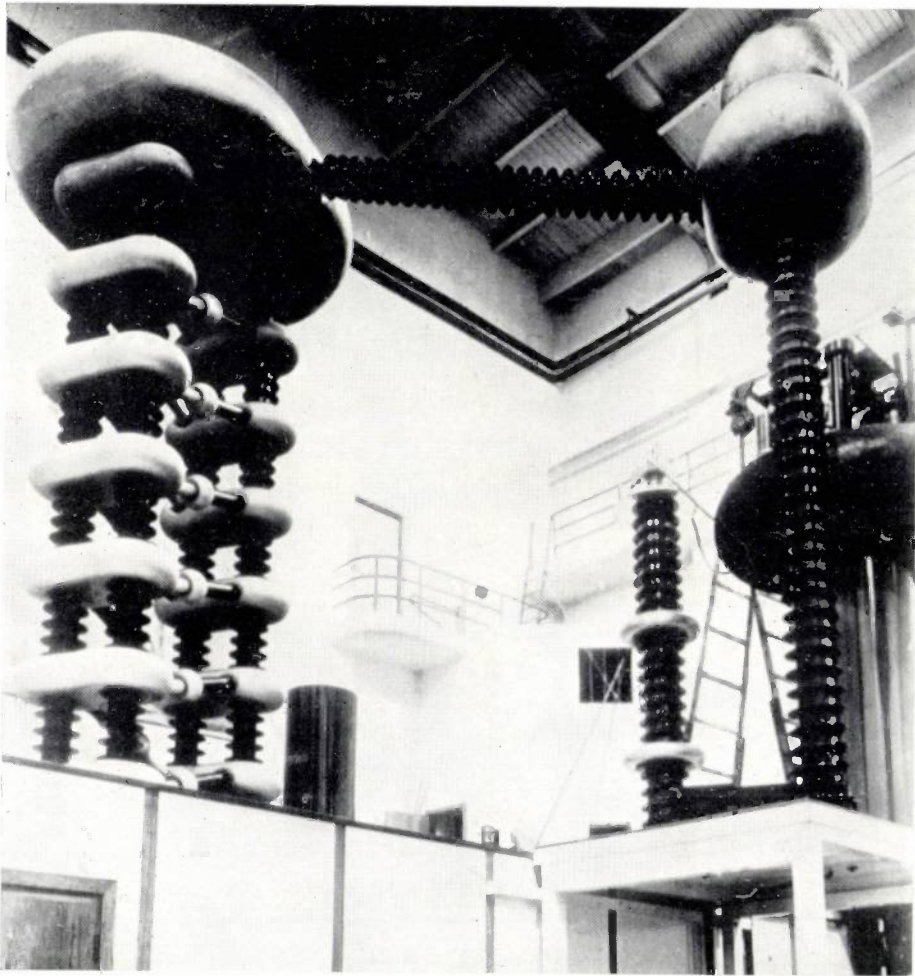


Fig. 5. First Philips cascade generator (1 MeV) installed at Cambridge in 1937.

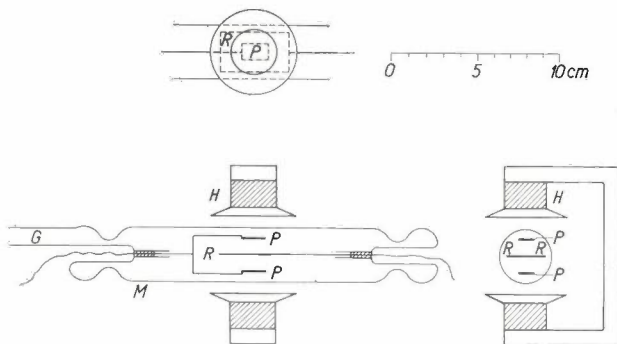


Fig. 6. Geometry of an early model of the Penning ionization gauge; taken from Penning's first publication on this subject.

ages. In our more mercenary moods we might even regret that the patents have run out.

In the thirties our laboratories began to study ferromagnetism, and this has remained one of our major subjects of research ever since. I shall not be able to do justice to the work of Snoek¹⁴) and many others in a quarter of an hour or so but perhaps this does not really matter because many of the results have become staple knowledge. The following dia-

gram (Table I) may help to survey the situation. We can distinguish two classes of materials: metals and alloys on the one hand, semi-insulating oxides on the

Table I. Magnetic materials.

Main properties and use	Materials	
	metals	oxides
permanent	"Ticonal"	hexagonal ferrites (ferroxdure)
semi-permanent memories magnetic recording	thin films thin wire	square-loop ferrites iron-oxide powder (tape)
soft linear transformers, inductances, etc.	silicon iron permalloy etc.	cubic ferrites
soft non-linear magnetic amplifiers voltage regulators	silicon iron permalloy etc.	cubic ferrites
microwave applications		cubic ferrites hexagonal ferrites (ferroxdure, ferroxpiana) garnets
magnetostriction	nickel iron	special ferrites

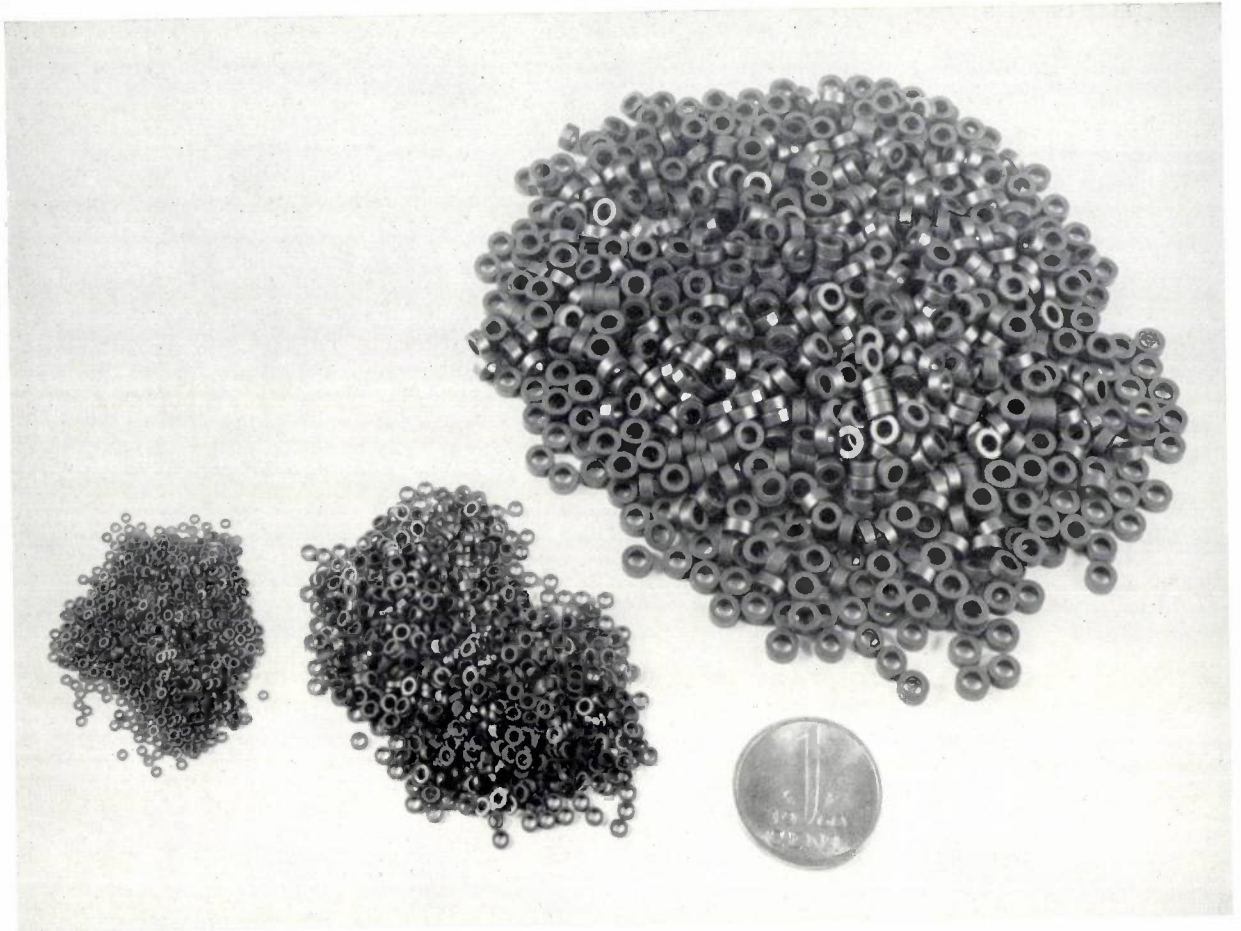


Fig. 7. Rings of square-loop ferrite for magnetic stores. Each pile contains 1000 rings.

other hand. This distinction usually coincides with a more fundamental distinction between ferromagnetic and ferrimagnetic materials. I shall come back to this question a little later. From the point of view of the electrical engineer we can distinguish several fields of application: permanent magnets are used for producing magnetic fields, magnetic materials used for storing information in terms of

magnetization I have called semi-permanent, soft linear materials find application in transformers, filter coils and so on, whereas soft non-linear materials (they may be the same materials used at higher fields) are applied in magnetic amplifiers, voltage regulators, modulators and in oscillators with varying frequency. In the microwave field isolators or unilines are of special importance, and we should also remember

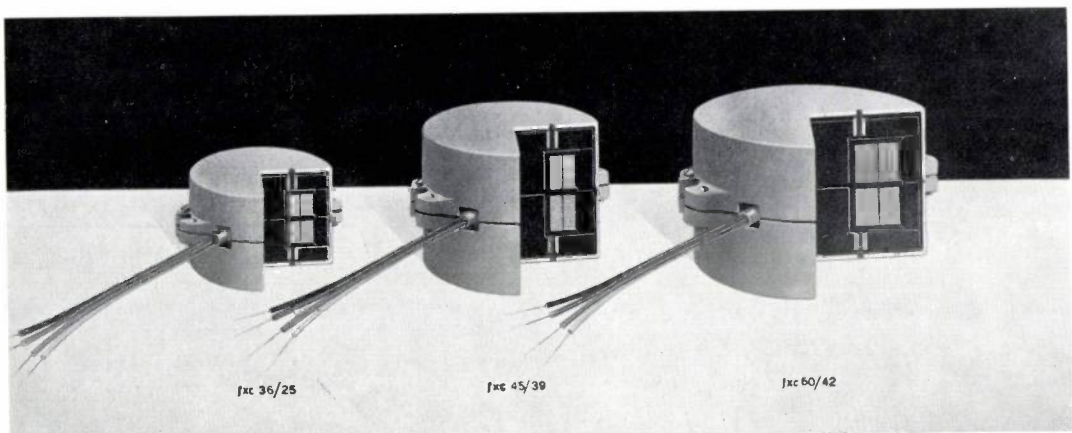


Fig. 8. Three Pupin coils, part of which has been cut away in order to show the windings and the ferrocube core (black).



Fig. 9. Ferrite rings for the 12 accelerator stations of the 33-BeV proton synchrotron at Brookhaven. (Photograph by courtesy of Brookhaven National Laboratory.)

the Suhl type of parametric amplifier. The mentioning of the following examples may illustrate our activities in these various fields. "Ticonal" 5¹⁵): this material was announced in 1937 and its figure of merit $BH_{\max} = 5 \times 10^6$ was obtained by cooling in a magnetic field. "Ticonal" 11¹⁶): this is a similar material but it is now a single crystal and was cooled in a magnetic field parallel to one of the cubic axes. Ferroxdure¹⁷): a substance that is chemically a hexa-ferrite, $\text{BaO} \cdot 6\text{Fe}_2\text{O}_3$, has a hexagonal crystal structure and a very high crystal anisotropy with a preferred direction of magnetization along the hexagonal axis. Its figure of merit is only 3.5×10^6 but a figure of merit does not always tell the whole story. These materials and the physico-chemical basis of their properties will be discussed in more detail by Dr. Vink. Each of these three little piles (fig. 7) contains 1000 tiny rings of a square-loop ferrite

(to be precise it is a Mn-Mg-Zn-ferrite for the largest and a Cu-Mn-ferrite for the two smaller sizes) and therefore the material for storing 1000 bits of information. Ferrite cores are used in many shapes and sizes and one of the earliest examples of their application is found in Pupin coils (fig. 8) which have largely replaced the older coils using highly anisotropic nickel iron ribbon that was magnetically soft in one direction. A rather spectacular application of ferrites is made in large accelerators like the Brookhaven and CERN machines, and I show you one of the rings made for that purpose (cf. fig. 9). Finally my box of exhibits contains a couple of isolators for millimetric waves (fig. 10) and an ultrasonic oscillator of a type that is finding increasing application for ultrasonic cleaning (fig. 11).

Let me now turn to some of the more basic phys-

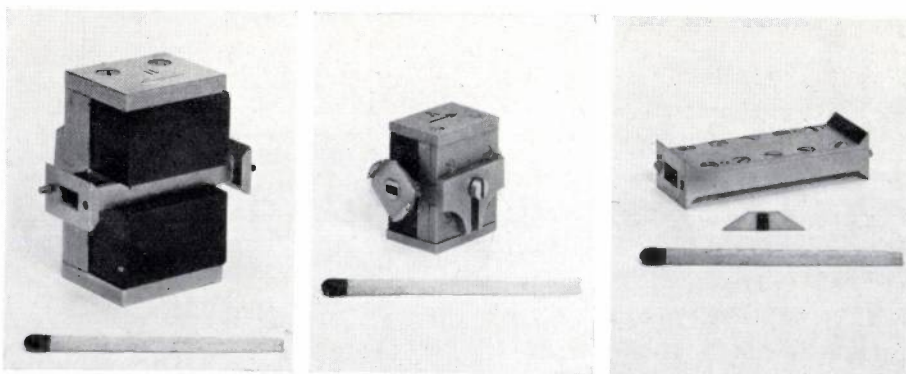


Fig. 10. Isolators for 8-mm waves (left and right) and 4-mm waves (centre).

ical aspects. I mentioned the difference between ferromagnetism and ferrimagnetism. Let me remind you of the basic ideas of ferrimagnetism, ideas that are due to Néel of Grenoble.

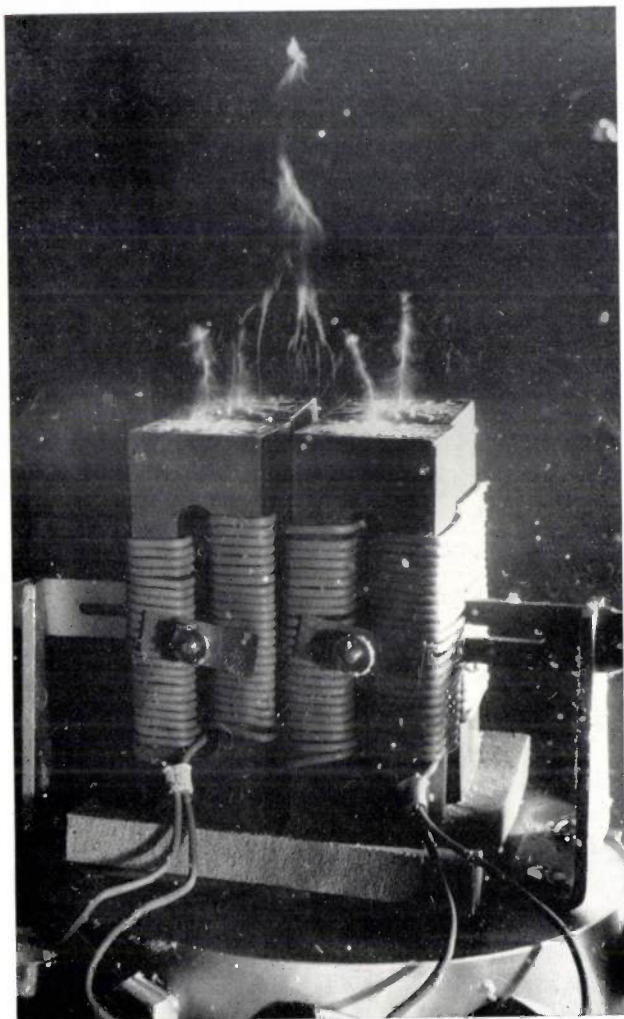


Fig. 11. Ferrite transducer for ultrasonic cleaning, photographed while operating in water.

In ferromagnetism the exchange force between neighbouring spins is such that it tends to orient these spins parallel to one another; in ferrimagnetism the interaction has opposite sign. But if there are various types of sites for the magnetic ions this interaction may yet result in a non-vanishing magnetic moment. In the simplest case there are two types of sites; let us call them A and B sites. Now it may well be that the interaction between A and B spins is larger than either the interaction between B and B and between A and A . Then the state of lowest energy is one in which all A spins are antiparallel to the B spins. For the total magnetic moment we can then write $M = M_A - M_B$. While the relation of these hypothetical A and B sites to actual crystal sites and the distribution of the various

magnetic ions over these sites will be discussed in Dr. Vinks paper, a rather spectacular physical consequence of this theory is illustrated in the following figures (fig. 12a and b). Gorter and Schulkes¹⁸⁾ prepared a lithium-chromiumferrite — the exact composition is $\text{Li}_{0.5}\text{Fe}_{1.25}\text{Cr}_{1.25}\text{O}_4$ — for which the temperature variation of M_A and M_B is such that at a certain temperature the resulting moment is zero. Van Wieringen¹⁹⁾ has measured the g -factor by means of ferromagnetic resonance. Now the magnetic moments are not exactly spin-only values but contain a certain orbital contribution that is different for different ions. Therefore magnetic moment and mechanical moment do not vanish at exactly the same temperature and hence the curious behaviour of the g -factor.

Far more detailed data on the partial moments can be obtained by means of the Mössbauer effect. This method was recently applied by Van Loef²⁰⁾

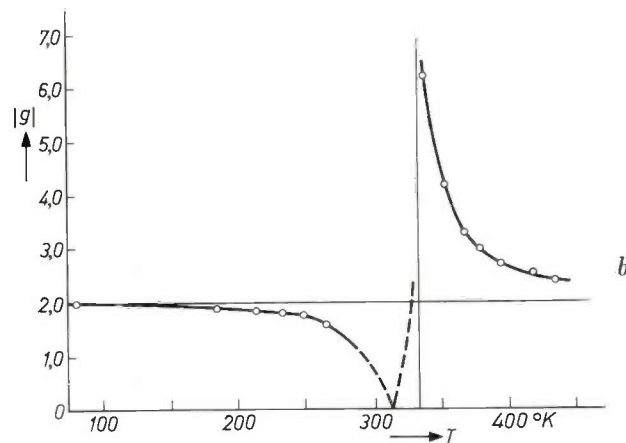
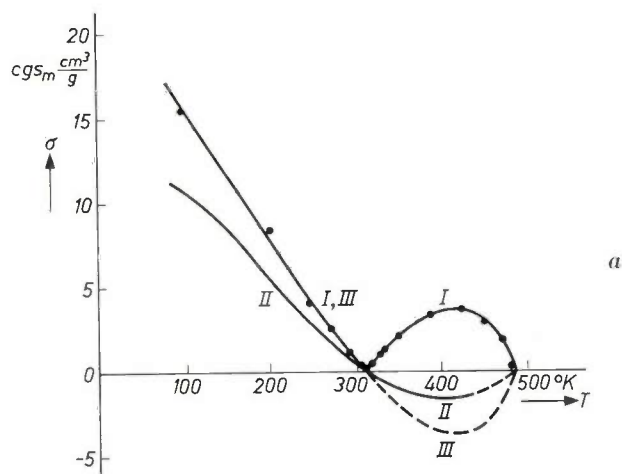
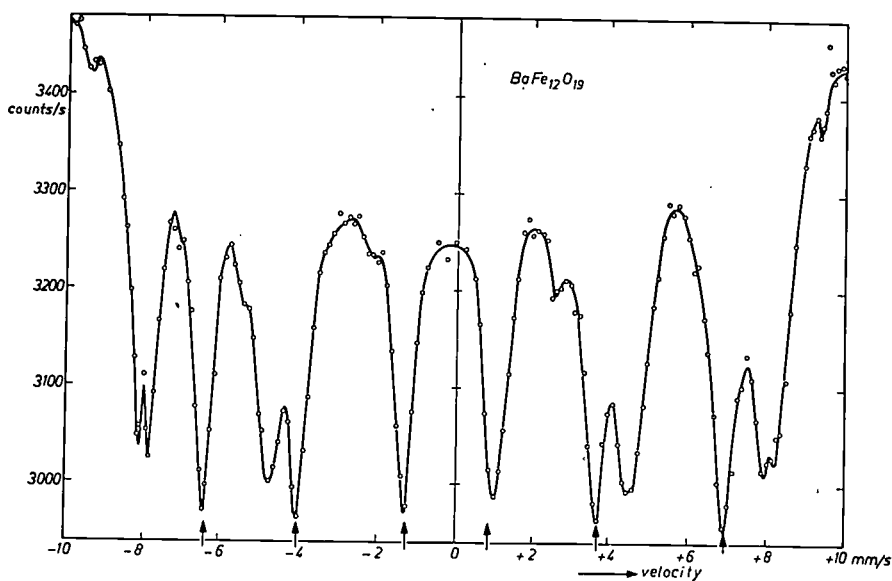


Fig. 12. a) Temperature variation of the saturation magnetization of $\text{Li}_{0.5}\text{Fe}_{1.25}\text{Cr}_{1.25}\text{O}_4$. Note that at a certain temperature the magnetization changes sign. b) Curious behaviour of the g -factor as a function of temperature, in the same substance.



to ferroxdure (fig. 13a). This illustration shows his results, obtained with a cobalt source and a sample of $BaO.6Fe_2O_3$ containing a high percentage of Fe-57. We see not only the normal hyperfine structure but we can distinguish for one hyperfine component four different lines corresponding to four different crystal sites. From crystallographic analysis we know that there are five kinds of sites for the iron ions but one line may well be weak. From measurements at different temperatures Van Loef has calculated the temperature dependence of the partial moments; their resultant moment is in good agreement with experiments (fig. 13b, c). We are convinced that continuation of this work will yield further important information on the partial moments in ferri-magnetic substances.

Before leaving the field of ferromagnetism I should like briefly to recall the ferroxdure-ferroxplana-ferroxcone story. Let us consider a hexagonal crystal and assume that the azimuthal anisotropy of magnetization is slight so that we can write for the contribution of crystalline anisotropy to free energy:

$$F = K_1 \sin^2 \theta + K_2 \sin^4 \theta. \quad \dots (1)$$

Then the direction of preferred orientation depends on the values of K_1 and K_2 . The result is given in the following figure (fig. 14). You will notice that there are regions where $\theta = 0$. That is the ferroxdure case. There are also regions with $\theta = \pi/2$. This is the case of ferroxpiana, of which we still hope that it will be developed into a high-frequency material of outstanding properties. But as was pointed out by Smit²¹⁾ there is also a region where the preferred orientations lie on a cone. Instances of ferroxcone behaviour have really been found. This is very gratifying but so far we have not been able to think

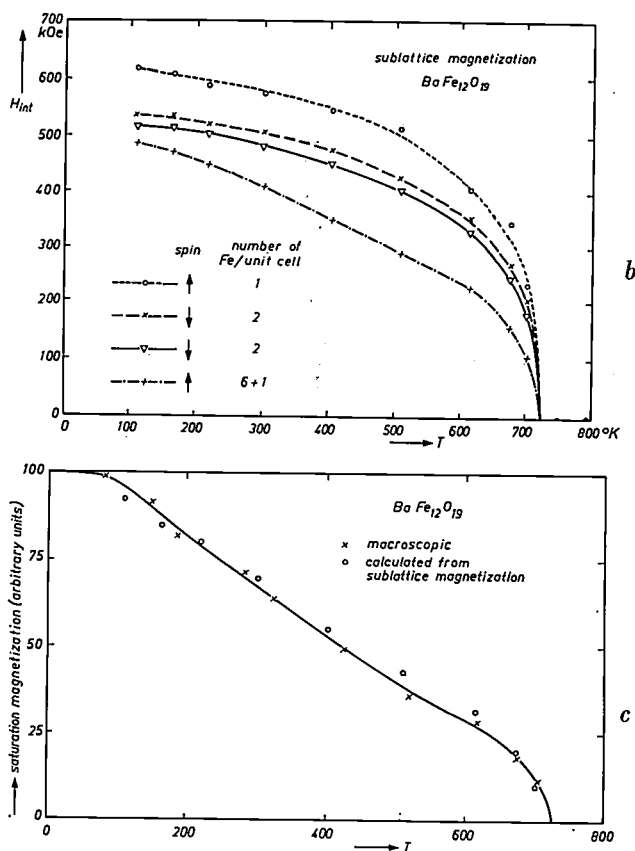


Fig. 13. a) Mössbauer absorption spectrum of Fe-57 in $BaO.6Fe_2O_3$. b) Partial magnetic moments of Fe on four sublattices in $BaO.6Fe_2O_3$ as obtained from the Mössbauer spectrum. c) Comparison between sum of partial moments and measured total saturation magnetization.

of any useful applications. Further material on preferred orientations will be presented by Dr. Enz in the course of this Symposium.

I mentioned optical one-way systems and microwave gyrators at the beginning of this lecture and

showed in passing one or two devices. I have not yet mentioned that the gyrator concept formulated by Tellegen²²) was an important intermediary between Rayleigh's idea and Hogan's realization. Tellegen

the existence of atomic currents. The meaning — and the only meaning — of the magnetization M is that the current density is given by:

$$i = \text{curl } M. \quad \dots \dots \dots (4)$$

Perhaps, when the notion of a spinning electron was introduced one may temporarily have been inclined to think in terms of dipoles, but Dirac's theory showed that this is wrong. In the fundamental state of a hydrogen atom for instance there is a current circulation around the nucleus. This current can to a high degree of accuracy be described by the analogon of the usual macroscopic equation

$$i = \text{curl } \psi^* s \psi, \quad \dots \dots \dots (5)$$

where s is the operator for the spin magnetic moment. This current distribution produces at the nucleus a magnetic field that can be calculated by simple electrodynamics formulae. The magnetic moment of the nucleus orients itself in this field which leads to a hyperfine splitting corresponding to a wavelength of 21 cm. It was the Dutch astronomer Van der Hulst who predicted that hydrogen atoms in empty space should emit this line. It was observed for the first time by Purcell but the Dutch astronomers, who were somewhat delayed because of a fire that destroyed part of their apparatus, were among the first to obtain valuable information about the structure of the galaxy from such observations. *Fig. 15* shows the Dutch radiotelescope at Dwingeloo and I am glad to report that our laboratories have been able to assist in the design of several important pieces of equipment²⁴).

If we regard ferromagnetism as one way of obtaining permanent circulating currents then a comparison with superconductivity suggests itself. There are very fundamental differences. Since the current in magnets is a spin current of the type $\text{curl } \psi^*(\sum s)\psi$ it can never circulate around a hole. That is why permanent magnets are not particularly suitable for producing homogeneous fields along the axis of a long cylinder. A current in a superconductor on the other hand is of the type $\psi^* A \psi$, where A is the vector potential, and can circulate around holes. The difference appears very clearly in the following simple experiment that was carried out by Van Geuns. A thick-walled cylinder of NbSn is cooled below its transition point while the bore is filled by a closely fitting rod of "Ticonal". Below the transition temperature the "Ticonal" is pulled out and a persisting current is started. The following figure gives a schematic drawing of the currents in both situations (*fig. 16*). During our Symposium we shall demonstrate several experiments on superconductivity,

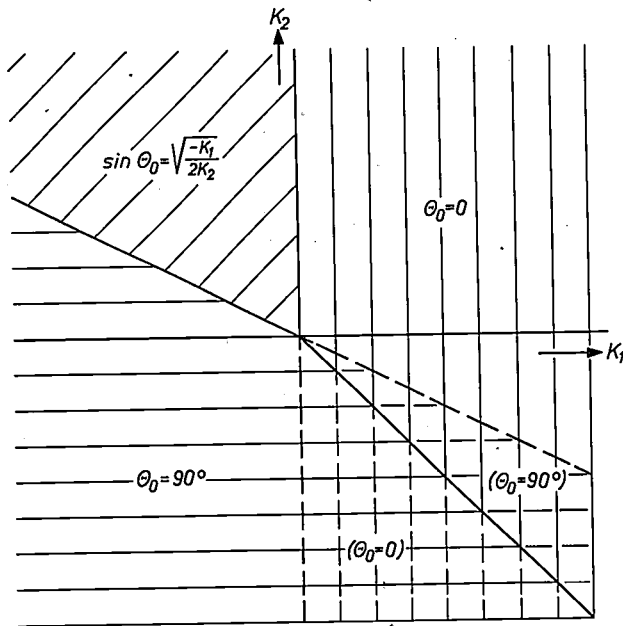


Fig. 14. In different areas of the K_1 - K_2 -plane different types of preferred magnetic orientation are found: preferred direction ($\theta_0 = 0$); preferred plane ($\theta_0 = 90^\circ$); preferred cone ($\sin \theta_0 = \sqrt{-K_1/2K_2}$).

studied from a formal point of view four-pole networks described by:

$$\left. \begin{aligned} V_1 &= A_{11}i_1 + A_{12}i_2, \\ V_2 &= A_{21}i_1 + A_{22}i_2, \end{aligned} \right\} \dots \dots \dots (2)$$

for which the symmetry relation is violated, and especially idealized networks for which $A_{12} = -A_{21}$. He also pointed out that for systems in an external magnetic field possibilities for realizing such gyrators might well exist. I was able to show that the symmetry relations for passive four-pole networks are a special case of Onsager's relations for irreversible processes, which in turn are based on the symmetry between past and future of the equations of motion of the elementary particles and of the equations for the electromagnetic field²³). Therefore the following general formulation holds: if a passive four-pole network incorporates magnetic moments M_1, M_2, \dots and is subjected to external fields H_1, H_2, \dots then

$$\begin{aligned} A_{12}(H_1, H_2, \dots, M_1, M_2) &= \\ &= A_{21}(-H_1, -H_2, \dots, -M_1, -M_2, \dots) \dots \dots (3) \end{aligned}$$

Of course, ferromagnetism makes one also think of radio astronomy. It is well known that the magnetic behaviour of matter must be ascribed to

especially the very elegant unipolar generator conceived by Volger²⁵). The Stirling-cycle engines for producing low temperatures will be demonstrated presently.

Sometimes it seems to me that our successes with ferromagnetic ceramics have in the past temporarily somewhat detracted our attention from the study of single crystals — although Van Arkel²⁶) invented the beautiful method for growing single crystals of refractory metals by vapour phase transport reactions. If there is any truth in this remark we are in any case trying hard to mend our ways. The work of Penning, Polder and Okkerse²⁷) on anomalous transmission of X-rays bears witness of our preoccupation with perfect, dislocation-free crystals. I should also like to mention the work of Hornstra²⁸) who carefully analysed the various possible types of dislocations in complicated crystal structures.



Photograph H. Nienrink, Leiden

Fig. 15. Antenna of the Dutch radiotelescope at Dwingeloo, at its completion in 1956 the largest movable reflector in the world. Several important pieces of equipment were developed in our Eindhoven laboratory.

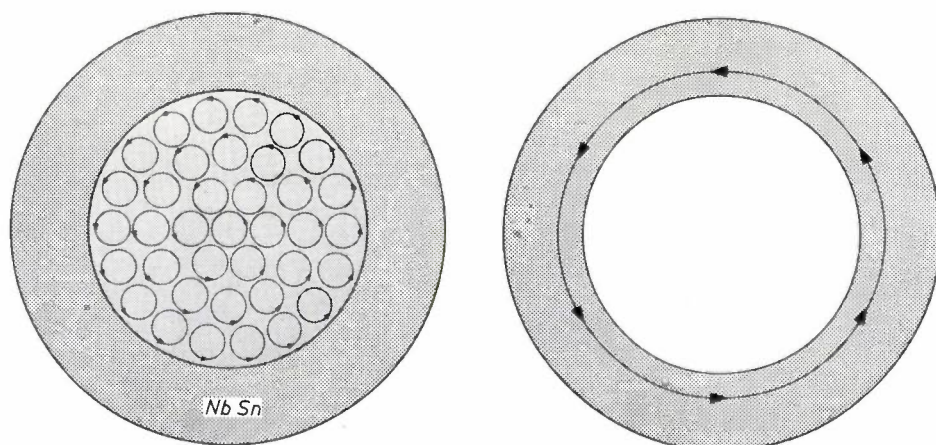


Fig. 16. Generation of a persistent current in a superconducting ring of NbSn by removal of a tightly fitting rod of "Ticonal".

Of course ferromagnetism is just one branch of solid state physics but I will not have time to touch upon our work on semiconducting materials and devices, on luminescence and photoconductivity, on metals and on surface phenomena. I trust these subjects will adequately be dealt with by others.

However, I should like to call your attention to

the importance of solid state work for the design of electron tubes. Our X-ray intensifier²⁹⁾ has become a useful tool in modern radiography (*fig. 17*). Its electron-optical system is simple and efficient but its main distinguishing feature is its combination of a fluorescent screen and a photo-electric layer which required a rather special technique. Our latest television camera-tube, the "Plumbicon", which will be

fundamental research; in fact it was invented by the head of our tube shop, Mr. Lemmens. He was killed in a motoring accident a few months ago and I should like to pay a tribute to his memory. In our work we do not only depend on the achievements of the great scientific pioneers but also on the skill and intuition and perseverance of outstanding craftsmen.



Fig. 17. 9" X-ray image intensifiers during fabrication.

extensively demonstrated, is as far as its principles of operation are concerned similar to the well-known vidicon. Yet it is an entirely different tube; its special performance is the result of a careful study of lead-oxide layers. And, although our centimetric and millimetric tubes require fine workmanship and careful design, they could never have been made without our dispenser cathodes³⁰⁾. We have done quite a bit of work on oxide cathodes, both empirical and theoretical. You may for instance remember the work of Vink and Loosjes³¹⁾, who explained the conduction in porous layers of barium oxide as a conduction in the electron gas that fills the pores. But I should like to emphasize that the first idea of the L-cathode was only indirectly the result of

The interactions we are dealing with in solid state physics are always electromagnetic in origin but electric forces may appear in many different guises. As simple Coulomb forces in heteropolar compounds, as screened Coulomb forces in the more sophisticated versions of the theory of electrons in metals, as exchange forces in homopolar compounds and in ferromagnetism. And they give also rise to attractive forces between neutral atoms and molecules even when these are so far apart that their wave functions do not overlap. Such forces are usually called Van der Waals forces because the forces introduced by Van der Waals to explain the behaviour of non-ideal gases are at least partly of this type. The theory has been given by F. London, who derived the following

formula for the interaction energy:

$$U = -\frac{3}{2} \frac{hc}{R^6} \sum_{n,r} \frac{a_n a_r}{\lambda_n + \lambda_r} \dots (6)$$

Here R is the mutual distance, the summation extends over all excited states of both atoms, λ_n and λ_r are the wavelengths corresponding to the excitation energies and a_n, a_r are partial polarizabilities. These London-Van der Waals forces give in some cases an appreciable contribution to the binding energy in solids and they play an important role in adsorption phenomena.

Hamaker ³²⁾ was the first to point out that Van der Waals attraction between colloidal particles may be significant and later Verwey and Overbeek ³³⁾ developed in detail the theory of stability of lyophobic colloids in which Van der Waals attraction is contending with the repulsion of electric double layers. A further step was made by Overbeek when he found that a suspension of quartz particles was more stable than was to be expected from his formulae. He concluded that at large distances the interaction must decrease more rapidly than the inverse sixth power and he suggested that this might be due to retardation effects. Polder and I ³⁴⁾ were able to treat this problem by means of quantum electrodynamics: we found that London's formula has to be multiplied by a monotonically decreasing correction factor. In the limit $R \gg \lambda_n$ one finds:

$$U = -\frac{23}{8\pi^2} \times \frac{hc}{R^7} a_1(0) a_2(0), \dots (7)$$

where $a_1(0), a_2(0)$ are the static polarizabilities. This limiting case can be treated very simply by considering the change of zero point energy of empty space. Applying similar ideas to resonant cavities with movable walls I derived a formula for the attraction between two conducting plates. The force per cm^2 is given by:

$$K = \frac{\pi hc}{480} \frac{1}{d^4}, \dots (8)$$

when d is the distance between plates. Sparnaay ³⁵⁾ was able to measure this force and found good agreement between theory and experiment. More qualitative evidence for this type of forces was found by another group of our laboratories when they had to approach a metal surface with very thin wires ending in a little sphere.

Personally I think that this connection between colloid chemistry, Van der Waals forces, quantum field theory and zero point energy of empty space and universal attraction between metals is rather amusing. Also I should like to point out that there

exists now a tendency to make smaller and smaller objects. If we ever learn to manipulate dimensions and distances of a few tenths of a micron then forces between metals might have a dominating influence.

In this rambling talk I have already mentioned many subjects and indicated many lines of research. I have spoken from the point of view of a physicist. This I felt I could safely do because the points of view of the chemist and the engineer will be well represented during the rest of this meeting. But I am afraid that one aspect of our work will receive inadequate attention and that is the mathematical aspect which does not lend itself so easily to short lectures or snappy demonstrations. Yet we have ever since the days of Van der Pol a pretty good tradition in applied mathematics. Now I have just time for explaining one little piece of mathematics and my choice is Van der Pauw's theorem ³⁶⁾ because of its outstanding elegance and simplicity. If we have to determine a specific resistance we usually take a prismatic or cylindrical rod (*fig. 18*). Van der Pauw's theorem states that we can take a plate of known and constant thickness but otherwise of arbitrary shape with four small contacts at arbitrary points (1, 2, 3, 4) on the circumference. Then, if we define R_{12} as the ratio of the voltage between 1 and 2 and current between 3 and 4, and similarly R_{23} , we have:

$$e^{-\pi d R_{12}/\rho} + e^{-\pi d R_{23}/\rho} = 1. \dots (9)$$

Let me indicate the proof. First take an infinite half plane. Then all potentials are logarithmic and the proof is elementary. Next map conformally to obtain the desired contour and remember that resistance is a conformal invariant. That is all. Very simple once you come to think of it which apparently no one did before Van der Pauw. And very useful: one measures two resistance values and reads ρ/d from

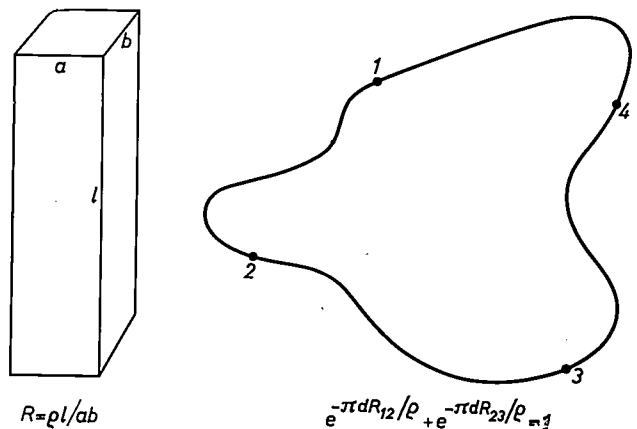


Fig. 18. Illustrating Van der Pauw's theorem.

a set of curves. We do not only apply this in the research laboratory but also in our factories to test specific resistance.

And this is the end of my talk.

BIBLIOGRAPHY

In general, only the earliest publication on each subject has been quoted.

- 1) G. Holst and L. Hamburger, Proc. Kon. Ned. Akad. Wet., Amsterdam, **18**, 872, 1915.
- 2) W. Uytterhoeven and C. Verburg, Bull. techn. Ass. Ing. Ec. polytechn. Brux. **29**, 138, 1933.
- 3) W. Elenbaas, Z. techn. Phys. **17**, 61, 1936.
- 4) M. H. A. van de Weijer, Philips tech. Rev. **23**, 246, 1961/62.
- 5) K. van Duuren, A. J. M. Jaspers and J. Hermsen, Nucleonics **17**, No. 6, 86, 1959.
- 6) N. F. Verster, H. L. Hagedoorn, J. Zwanenburg, A. J. J. Franken and J. Geel, Nucl. Instr. Meth. **18/19**, 88, 1962.
- 7) J. L. H. Jonker, A. J. W. M. van Overbeek and P. H. de Beurs, Philips Res. Repts **7**, 81, 1952.
- 8) O. Reifenschweiler, Nucleonics **18**, No. 12, 69, 1960.
- 9) M. T. Vlaardingerbroek, K. R. U. Weimer and H. J. C. A. Nunnink, Philips Res. Repts **17**, 344, 1962.
- 10) H. G. van Bueren, J. Haisma and H. de Lang, Physics Letters **2**, 340, 1962.
- 11) A. Bouwers and A. Kuntke, Z. techn. Phys. **18**, 209, 1937.
- 12) F. M. Penning, Physica **4**, 71, 1937.
- 13) F. M. Penning and J. H. A. Moubis, Physica **4**, 1190, 1937.
- 14) J. L. Snoek, Physica **3**, 463, 1936; New developments in ferromagnetic materials, Elsevier, Amsterdam 1949.
- 15) B. Jonas and H. J. Meerkamp van Embden, Philips tech. Rev. **6**, 8, 1941.
- 16) A. I. Luteijn and K. J. de Vos, Philips Res. Repts **11**, 489, 1956.
- 17) J. J. Went, G. W. Rathenau, E. W. Gorter and G. W. van Oosterhout, Philips tech. Rev. **13**, 194, 1951/52.
- 18) E. W. Gorter and J. A. Schulkes, Phys. Rev. **90**, 487, 1953.
- 19) J. S. van Wieringen, Phys. Rev. **90**, 488, 1953.
- 20) J. J. van Loef, to be published.
- 21) J. Smit, J. Phys. Radium **20**, 362, 1959.
- 22) B. D. H. Tellegen, Philips Res. Repts **3**, 81, 1948.
- 23) H. B. G. Casimir, Rev. mod. Phys. **17**, 343, 1945.
- 24) C. L. Seeger, F. L. H. M. Stumpers and N. van Hurck, Philips tech. Rev. **21**, 317, 1959/60.
- 25) J. Volger and P. S. Admiraal, Physics Letters **2**, 257, 1962.
- 26) A. E. van Arkel, Physica **3**, 76, 1923.
- 27) P. Penning and D. Polder, Philips Res. Repts **16**, 419, 1961.
- 28) J. Hornstra, J. Phys. Chem. Solids **5**, 129, 1958.
- 29) M. C. Teves and T. Tol, Philips tech. Rev. **14**, 33, 1952/53.
- 30) H. J. Lemmens, M. J. Jansen and R. Loosjes, Philips tech. Rev. **11**, 341, 1949/50.
- 31) R. Loosjes and H. J. Vink, Philips Res. Repts **4**, 449, 1949.
- 32) H. C. Hamaker, Physica **4**, 1058, 1937.
- 33) E. J. W. Verwey and J. Th. G. Overbeek, Theory of the stability of lyophobic colloids, Elsevier, Amsterdam 1948.
- 34) H. B. G. Casimir and D. Polder, Nature **158**, 787, 1946.
- 35) M. J. Sparnaaij, Nature **180**, 334, 1957.
- 36) L. J. van der Pauw, Philips Res. Repts **13**, 1, 1958.

Summary. A lecture in which numerous subjects of research at the Philips Laboratories are briefly touched upon. Gas discharges, ferromagnetism, solid-state physics, network theory are among the main items. The investigations in these variegated fields have resulted on the one hand in industrial products, such as lamps, counter tubes, cyclotrons, ferromagnetic materials and devices of various types, cathodes etc., many of which are illustrated in this lecture, on the other hand they have materialized in contributions to basic physical theory, some of which are discussed.

MODULATION, YESTERDAY AND TOMORROW

by F. de JAGER *).

621.376

Problems of today's communications

Fifty years ago there was no need for a theory of modulation. The problem of transmitting speech or music by means of radio waves had been solved, and most of the physical problems involved were well understood. The people putting these principles into practice found themselves in the happy situation of having available an enormous frequency space, which was very quiet and extensive like the deep blue sea.

Nowadays, however, the situation has changed. The evergrowing need for more channels has filled up this once silent frequency space and we now have a large number of transmitters radiating power in the form of television, speech, radar and control signals. Indeed we are not far from the well-known situation in a cocktail party where everybody is raising his voice in trying to overcome the background noise generated by other people. In telecommunication nowadays this man-made noise presents a far more serious problem than the small atmospheric disturbances, with which we were confronted in the early days of telecommunication. Now we are forced to study the theory of modulation more seriously, aiming at a more effective transmission of information signals.

Three problems arise from this. Firstly, how to communicate in the presence of background noise; secondly, how to use minimum power; and thirdly, how to occupy the minimum bandwidth of the total frequency spectrum available. These are the problems from a theoretical point of view. Connected with this are the practical problems of how to realize the correct physical principles of modulation, preferably without running into very complex apparatus and high costs.

Fortunately, during the last decade, the rapid development of electrical components has supplied us with a great variety of small, cheap and reliable units, such as transistors or magnetic cores, which allow the realization of complex technical processes our ancestors would never have dreamed of. (As an example we need only mention modern high-speed electronic computers.) So, in a way, we are in the situation of a child who has received on his birthday a nice box of building-blocks and now has to determine what should be built with it. In this

respect modulation theory is the art of combining these building-blocks in an efficient way. This, of course, is only possible by keeping in mind the correct physical and mathematical principles, so that, in the end, we are studying physics again.

It may be of interest to look in more detail at some of the principles which have influenced methods of modulation in the past and — as far as we can foresee them — to discuss some future developments.

Early methods of modulation

The original and to some extent most natural way of transmitting information signals is by varying the amplitude of a high-frequency sinusoidal wave. This is called amplitude modulation and it is still applied in broadcast transmitters of long and medium wavelengths. In 1933, however, it was discovered by Armstrong that, generally speaking, a disturbing noise component could less easily affect the frequency than the amplitude of this wave. So Armstrong kept the amplitude at a constant value and transmitted the information signals by varying the frequency. This was called frequency modulation and it proved to be a very efficient way of reducing noise. The method is applied so generally now (for instance in f.m. broadcast transmission, for sound transmission in television, etc.) that it is hard to believe that in the early days of telecommunication these ideas were considered to be revolutionary.

In a physical sense the underlying principle is very simple. Originally the realization required a large amount of apparatus, spread out over a large table in the laboratory. Nowadays it is not difficult to contain the whole system in a matchbox. This illustrates how a new technical principle, which initially demands a rather complex realization, may be realized in a convenient manner after a few years development of the suitable components. In this respect we may be optimistic for future developments of new ideas. Indeed, the continuing progress in solid-state devices permits us to perform more and more functions with a still smaller number of molecules.

Pulse modulation: introducing non-linear techniques

The methods of modulation considered so far all used sine waves for the transmission of infor-

*) Philips Research Laboratories, Eindhoven.

mation, either by varying the amplitude or the frequency. It is, however, equally possible to modulate short pulses, either in amplitude or time. A remarkable fact is that, though this possibility was very well understood theoretically, the practical realization of this only occurred after a rather long time. Probably this is a consequence of the different non-linear techniques which were required for this purpose, and indeed a great deal of this circuitry was only developed during World War II. Besides, the shift from sine-wave techniques to pulse techniques implied more than simply a change of equipment. By tradition, all definitions were based on linear systems and their corresponding, well-known, linear differential equations. Now, with the new pulse circuitry, many circuits contained a non-linear element which could no longer be considered as a small deviation from the linear model, but whose non-linearity was of a very fundamental nature.

In fact, this introduction of a non-linear element in a circuit which for the rest is linear, gave rise to very difficult mathematical problems which have only been partly solved. In this relation may be mentioned the name of Balthasar van der Pol, who was at this laboratory during many years and who investigated many fundamental problems of this kind¹⁾. The fascinating point in these combined linear and non-linear circuits is that various effects have been discovered which have led to important technical applications. For instance, parametric amplification of very high frequencies belongs to this class. Here the variation of a linear network leads to an energy transfer between different frequencies that can be used for amplification.

Regarding the possibility of transmitting information by means of modulated pulses, a close analogy was found to modulated sine waves. Again, the modulation in amplitude was very simple, but did not reduce noise. Modulating the pulses in their position in time, however, could reduce the background noise, just as was found with the method of frequency modulation, mentioned earlier. Analysis of both systems showed that the improvement was proportional to the bandwidth used for transmission.

For some time this led to the belief that any improvement in signal-to-noise ratio required a proportional increase of bandwidth. However, this conclusion proved to be wrong in the light of the development of modern information theory as founded

by Shannon and others²⁾. On a purely mathematical basis, Shannon derived an upper limit for the improvement to be gained if the most efficient coding were used. It showed clearly that this improvement might be far greater than simply proportional and in fact should be of an exponential nature.

Pulse-code modulation: translating speech into numbers

Though we are still far from reaching Shannon's upper limit (which, by the way, would involve very long time delays), there is now one method known in which this exponential improvement with bandwidth can be realized. This is the modulation method known as pulse-code modulation. The underlying idea is that signals like speech or music should not be transmitted in the form in which they are delivered by a microphone, but that they should be transmitted as numbers. This means that at the transmitting end the signals must be analysed in digital form, the digits then being transmitted as pulses which may have either the value 1 or 0, and at the receiver the digits being used for reconstructing the original waveform (*fig. 1*).

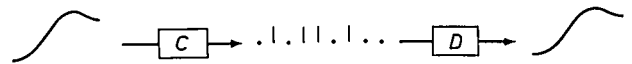


Fig. 1. Functional diagram of a pulse-code-modulation link: *C* coder, for converting speech into a sequence of "1" and "0" pulses; *D* decoder, for reconstructing an output signal which is nearly identical to the original one.

The principle could easily be compared with communicating the information contained in a graphical curve, by means of a table representing the horizontal and vertical coordinates. Indeed in practice the speech signals are measured within small intervals of time and the scale of amplitudes may be divided into, for instance, 128 different levels. Each of these levels can be characterized by a number expressed by 7 digits in the binary scale, because exactly 128 different numbers can be distinguished in this way. With a little bit of sophisticated circuitry these numbers can be converted into voltages again. Using this type of modulation a maximum deviation of 1% or smaller can now be guaranteed in the output signal, whereas the rather large deviation of 50% can be tolerated in the received pulses (*fig. 2*). The efficient use of bandwidth in this case may be demonstrated by the fact that with the addition of one digit to a group of, say, 7 binary digits, the total number of available amplitude levels is doubled,

¹⁾ Balth. van der Pol, Selected scientific papers, North-Holland Publ. Co., Amsterdam 1960.

²⁾ C. E. Shannon, A mathematical theory of communication, Bell Syst. tech. J. 27, 379-423 and 623-656, 1948.

whereas the effect on the resulting pulse frequency is very small, namely in the ratio 8 to 7. In these systems we do not have an exact reproduction of the original signals, but the remaining difference (which is often called "quantization noise") can be made very small by choosing a high enough pulse frequency (50 to 60 kc/s).

Transmitting information in the form of numbers has another advantage in the case of a large number of links put in tandem. The received pulses can be regenerated in each section so that small noise components in the different sections can be eliminated before they have an accumulating effect. In the more conventional modulation systems, where the information is transmitted by varying a physical parameter in proportion to the original signal, there is no possibility of distinguishing between signal and noise and in that case noise entering the different sections will be necessarily accumulated.

By using a special type of computer for coding and decoding speech, a more efficient use of bandwidth can thus be made for obtaining noise reduction than with the older methods of modulation. The question naturally arises whether this is an academic result only or whether it can be put into practice. During the last few years it has become apparent that an important practical problem can be solved in this way. This concerns the overcrowded telephone connections between telephone exchanges in large cities. A large number of telephone connections use the same cable (up to 900 pairs in one cable). The laying of new cable in these cities, however, causes so much trouble that the growing need for more telephone connections should preferably be met by making use of the existing cable network.

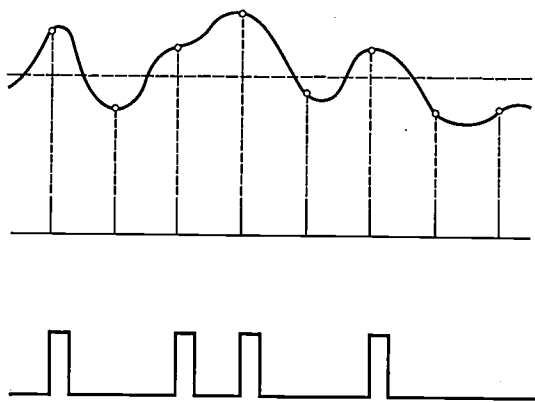


Fig. 2. Regeneration of a pulse series which is affected by noise, in a P.C.M. system. The incoming signal is sampled at the proper instants of time and, in accordance with the level of the detected signal, a unit pulse is, or is not, regenerated. As long as a disturbance does not exceed a given maximum, it can be completely eliminated.

Naturally one could extend the number of connections by using higher frequencies. To this end the mutual interference between different pairs in a cable should be eliminated, as it would otherwise cause serious "crosstalk" between the different channels. Unfortunately, however, this crosstalk is about proportional to the frequency. Using conventional methods of noise reduction the interferences could be suppressed by enlarging the bandwidth proportionally, but if the interferences also increase in this proportion the method is not very effective. So a more rapid improvement is wanted, and this is furnished by pulse-code modulation. In practical systems, as have been installed recently in the U.S.A. ³⁾, pulse frequencies of about $1\frac{1}{2}$ million pulses per second are used in the cables. As a result the number of telephone connections can be extended by a factor 12 and at the same time the quality of the new connections may be equal to or even better than that of the older ones.

Delta modulation: introducing a memory

The realization of a coding procedure using the binary code will always require some form of a computer. The question arises whether the conversion from a continuous signal into a series of "1" and "0" pulses, and the reversed procedure, could not be effected by means of a far simpler device. This problem was investigated some years ago by Schouten and others of this laboratory. It was found that, instead of using a binary code with a large number of digits, the sequence of "1" and "0" pulses could be arranged in such a manner that this pulse series only needs to be applied to the input of a simple electrical network or filter in order to generate the required signal at the output. This method is called delta modulation, and it is based on a simple control mechanism ⁴⁾.

In order to achieve the correct sequence of the pulses three things must be taken into account: the present deviation between the original and the approximation signal, the way in which preceding quantization errors can still be corrected and the changes to be expected in the information signal. So it might be said that each decision of transmitting a "1" or a "0" pulse is based, in fact, on the present, the past and the future. These functions can all be performed by a linear network, functioning as a

³⁾ C. G. Davis, An experimental pulse code modulation system for short-haul trunks, Bell Syst. tech. J. 41, 1-24, 1962 (No. 1).

⁴⁾ J. F. Schouten, F. de Jager and J. A. Greefkes, Delta modulation, a new modulation system for telecommunication, Philips tech. Rev. 13, 237-245, 1951/52.

memory, which is put in a non-linear feedback loop, thus making the circuit self-correcting⁵⁾ (figs. 3 and 4).

In analysing this method of delta modulation the way of transmitting information can, I think, best be illustrated by the following model. If the original information signal corresponds to a winding road, the transmitter may be seen as a car which is driven along this road. Here the pulses

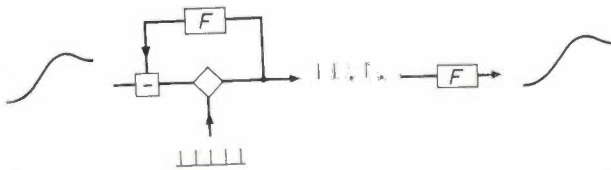


Fig. 3. Functional diagram of a delta modulation link. A "quantized" feedback circuit at the transmitting end produces "1" and "0" pulses. As a consequence of the feedback action, the signal at the output of the filter F resembles the original very closely. The same filter can thus be applied at the receiving end for converting the pulse series into an identical approximation signal.

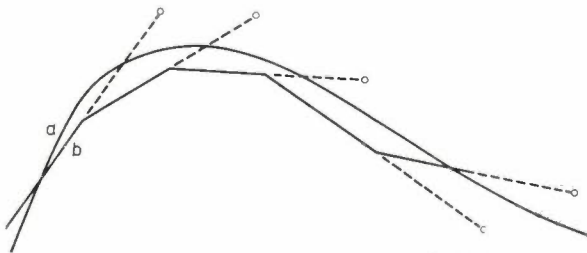


Fig. 4. Approximation of a continuous signal (a) by means of a quantized signal (b), obtained by changing the slope at regular intervals of time. Decisions on changing this slope are not based on the instantaneous deviations between original and approximation signal, but on their predicted values.

given to the steering wheel force the car to follow the road. This indeed implies a feedback loop because the driver watches the road ahead, predicts the movements of the car, and decides from the expected discrepancies whether or not a new pulse has to be applied to the steering wheel. If now identical pulses are given to the steering wheel of a hypothetical car, which is running through a desert, then a fairly good reproduction of the path that is followed by the first car can be obtained.

It was found that speech lends itself very well to this kind of coded transmission. Quite recently, however, the method has also been applied to normal television by Balder and Kramer of this laboratory⁶⁾. The whole computer, which can make 100 million decisions per second, is contained in a small box (fig. 5). This very high speed, by the way,

is only attained by making use of tunnel diodes: very small elements of semiconductor type, which can be switched in one or the other direction in about one thousandth of a micro-second.

The TASI system: using the time distribution of speech

With all modulation methods considered so far, the waveform of the signal that is to be transmitted is arbitrary to a large extent. It should only be limited in amplitude, of course. In modern modulation systems, however, more and more account is taken of the statistical properties of the signal, either in time, frequency or amplitude. For instance, speech is never a continuous stream of information but it consists of a large number of very short signals, often no longer than one tenth of a second. So the time during which signals are really present is very much reduced and one could consider using the silent time intervals for transmitting other information signals. This in fact has been done on circuits using a transatlantic telephone cable. In a large group of telephone conversations the signals are split up into parts and only those parts carrying information are transmitted, after being labelled in such a way that these elementary parts containing information can be directed to the correct receiver. The system was called Time Assignment Speech Interpolation and it was put in service on transatlantic cables by Western Electric in 1960. The result is that the original capacity of 36 channels could be enlarged to 72 channels⁷⁾.

Of course systems like these are very complex, but on the other hand they are very economical too.

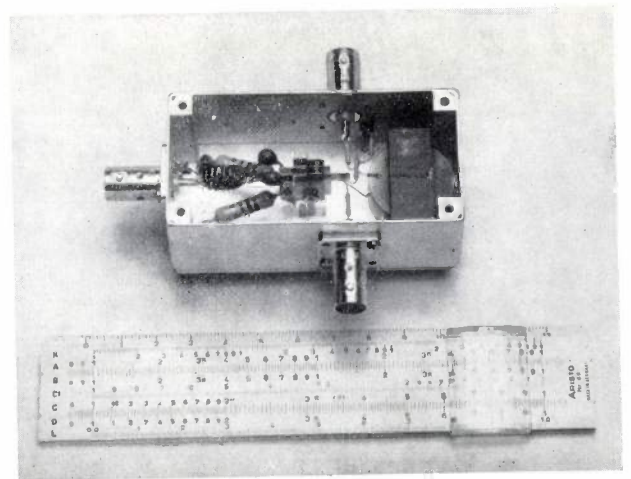


Fig. 5. Tunnel-diode delta modulator for coding television signals (see⁶⁾).

⁵⁾ F. de Jager, Delta modulation, a method of P.C.M. transmission using the 1-unit code, Philips Res. Repts 7, 442-466, 1952.

⁶⁾ J. C. Balder and C. Kramer, Video transmission by delta modulation using tunnel diodes, Proc. IRE 50, 428-431, 1962 (No. 4).

⁷⁾ J. M. Fraser, D. B. Bullock and N. G. Long, Over-all characteristics of a TASI system, Bell Syst. tech. J. 41, 1439-1454, 1962 (No. 4).

Indeed, the installation costs of a transatlantic cable is of the order of 40 million dollar, so that to double the number of channels it is far cheaper to build a somewhat complicated apparatus than to lay a new cable.

Speech bandwidth compression: studying speech parameters

Now we enter the study of the speech signals themselves. In analysing speech it is found that a speech signal can be characterized by a number of parameters like intensity, pitch, spectral distribution, etc. One thing is clear now beyond any doubt. The maximum speed of variation of these parameters is far lower than the speed allowed by the normal telephone channel. This is due simply to the fact that the rate of speech is limited by the velocity of muscular movements. Thus the available bandwidth of a telephone channel is being used in a rather inefficient way. This can even be expressed quantitatively. According to Shannon's coding theorem a normal telephone channel of 3 kc/s bandwidth and 50 dB signal-to-noise ratio has an upper limit of transmitting information of about 50 000 bits per second. Our senses, however, are not able to detect more than about 50 bits per second. Nevertheless, for the exact reproduction of a speech waveform a bandwidth of 3 kc/s is required and this means that a large part of our speech signals does not carry new information but consists of repetitions only. This phenomenon may be observed in any oscillogram of speech (*fig. 6*). Of course, in everyday life this large amount of redundancy is very important, for it enables us to disregard many disturbances.

The technical aspect of this redundancy is that, if we are able to distinguish the proper physical parameters of speech, we could, by transmitting them, carry a speech signal in far less bandwidth than we are accustomed to do now. Theoretically the reduction could be more than a hundred to one, but we should be very glad to achieve a ten-to-one ratio.

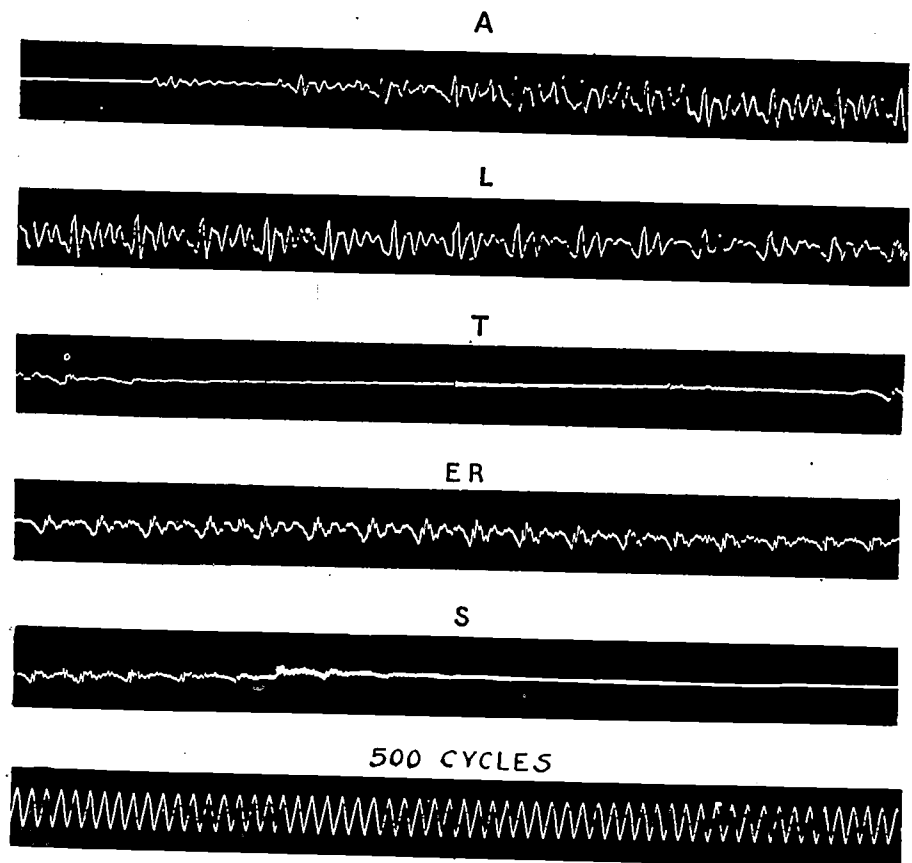


Fig. 6. Waveform of the word "ALTERS", showing the periodicity in speech signals. (Taken from H. Fletcher, *Speech and hearing in communication*, p. 33, Van Nostrand, New York 1953.)

This again would be very economical for use on transatlantic cables, and also in those frequency regions of radio transmission which are now overcrowded.

The study of these speech parameters dates back already to the eighteenth century, when Wolfgang von Kempelen built a famous speaking-machine, using bellows and horns, for reproducing speech sounds⁸⁾. It seems that the reproduction of these sounds was fairly good, as it is reported that his machine could even speak Latin. A more up-to-date speaking machine however was constructed and demonstrated by Dudley⁹⁾ in 1939 (*fig. 7*). This electrical machine could be operated manually, but the necessary parameters could also be extracted automatically from real speech by means of an "analyser" at the transmitting end which was connected to the "speech synthesizer" at the receiving end. This "vocoder", as it was named, was able to

⁸⁾ H. Dudley and T. H. Tarnoczy, The speaking machine of Wolfgang von Kempelen, *J. Acoust. Soc. Amer.* 22, 151-166, 1950. See also a forthcoming issue of *Philips tech. Rev.* devoted to the Institute for Perception Research.

⁹⁾ H. Dudley, Remaking speech, *J. Acoust. Soc. Amer.* 11, 169-177, 1939.

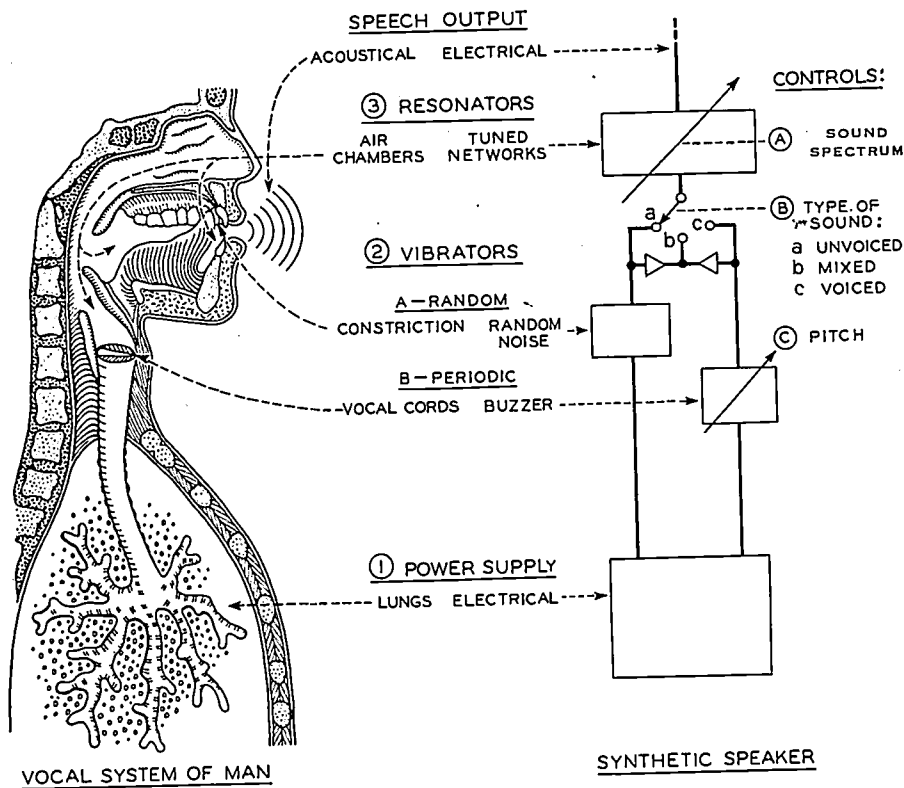


Fig. 7. Functional comparison of synthetic speaker with the human vocal system. (Taken from H. Dudley, R. R. Riesz and S. S. A. Watkins, A synthetic speaker, J. Franklin Inst. 227, 739-764, 1939, fig. 5.)

transmit speech by means of a great number of parameters, each varying so slowly that a bandwidth reduction of a factor 10 could be obtained. The principle on which it was based was: splitting the frequency band of normal speech into a great number of channels, measuring the power in each channel and then regenerating a signal having a corresponding power spectrum at the receiving end. Most of the speech compression systems of today still use the same principle¹⁰); see fig. 8a and b. The results of this first vocoder were very promising, though the remade speech suffered from a lack of naturalness.

Pattern recognition: the missing link

Experiments on speech band-compression systems were continued, but up to now it has been found extremely difficult to achieve natural quality. However, this problem must be solved first, before application to the public telephone system is possible. Indeed, any telephone subscriber expects to be answered by a human voice, and not by a machine reproducing words in a monotonous, impersonal way.

¹⁰) F. H. Slaymaker, Bandwidth compression by means of vocoders, IRE Trans. on Audio AU-8, 20-26, 1960.

Through this research a large amount of information relating to these speech parameters is available now. There are probably two reasons why it has not been possible, so far, to construct a real high-quality vocoder. Firstly, our hearing mechanism makes use of a very great number of nerve fibres which are used for the most part in parallel. Constructing an electrical model would thus require so many circuits that it is technically unattractive. (Possibly the progress in micro-miniaturization may be of help here.) Secondly, however, we have arrived at the far more serious problem that, though we have analysed speech sounds and the mechanism of the ear very carefully,

we really do not know how the received signals are combined and interpreted by the human brain. As long as this is not known, the design of new vocoders will necessarily proceed in a haphazard way (which, eventually, may lead to a solution), but we must not be surprised to find that the analysis of real speech is far more difficult for a machine than for a human observer.

These problems belong, in fact, to the more general problem of pattern recognition: how to build machines which can read and identify informational elements like printed letters, numbers, or acoustic patterns as encountered here. This is a relatively new field of research and it will require the combined effort of mathematical, physical and biological research to solve these problems.

Now suppose one discovers what the "information-bearing elements" of speech really are and that one knows how to develop a system of speech bandwidth compression based on this. There still would arise a curious problem: what will a speech compression system do with the crackling noise of paper that sometimes accompanies the speaker's voice into the microphone? The system is taught to reproduce speech signals only and so it will

convert any noise of this kind into speech. The same effect, by the way, is met in experiments with frequency-band compression systems for television, where noise is not demodulated as noise but as a new type of signal. Practical difficulties of this kind will limit the compression ratio of high-quality vocoders.

As regards practical application of vocoders in communication an intermediate solution may be found, perhaps, in transmitting a small part of the speech spectrum in its natural way and completing this with an additional transmission of parameters representing the larger part of the spectrum. In fact it has been found that especially the lower-frequency region from about 300 to 800 c/s lends itself

reliable operation than trying to obtain security by handling the speech wave in its original form.

Companadors: using the amplitude distribution of speech

In the foregoing examples we have seen how the statistical properties of the time and frequency distributions of speech can be used for increasing the number of channels. From an engineering point of view it is worth-while to ask what can be done with the statistics of the amplitude distribution. It is found then primarily that this knowledge can be used for obtaining a better noise reduction. Generally speaking, the large range of different amplitudes occurring in speech or music can be compressed into

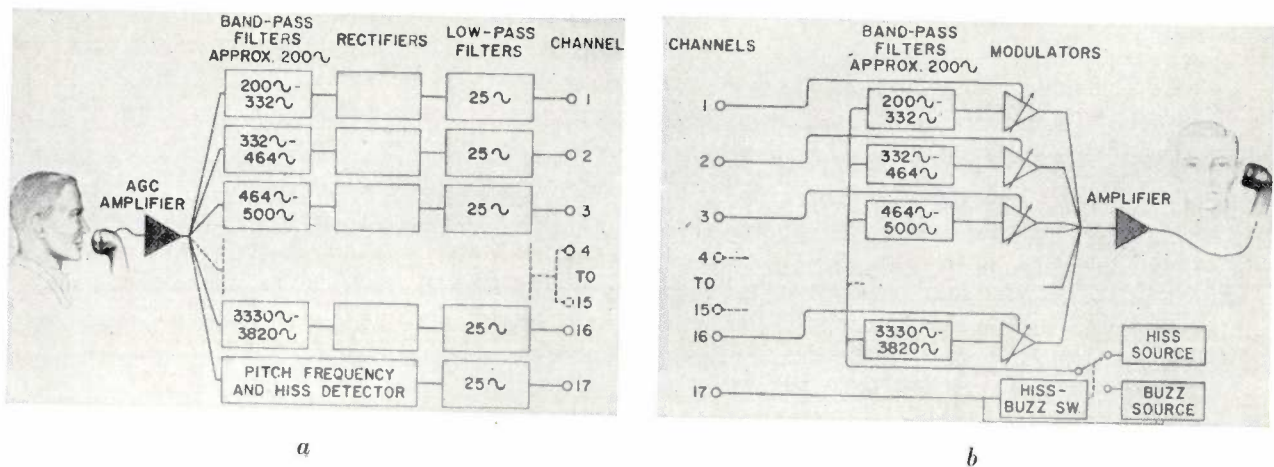


Fig. 8. Sketch of the apparatus in a vocoder system. (Taken from Slaymaker¹⁰.)
 a) Transmitting end.
 b) Receiving end.

conveniently for direct transmission, if it is combined with a parametric transmission of the remaining part of the frequency band. The signals of the first frequency range then constitute a natural base for the reconstructed speech signals and the difficult problem of pitch determination is avoided. The construction of a vocoder is much easier in this way, though the resulting ratio of compression is also much lower.

In spite of the above-mentioned difficulties in vocoder systems it must be admitted that some progress has been made during the last few years in relation to low-quality vocoders. If intelligibility and not naturalness is of primary importance it is now possible to transmit speech signals with only 2000 "1" or "0" pulses per second. In military applications these pulses can be re-coded for the purpose of security and then be transmitted via the normal telephone network. This is a more

a smaller range during transmission, and after reception it can be expanded again. This especially favours the very small information signals which otherwise would have been lost in the noise. Devices of this kind, using a compressor and a corresponding expander, are called companadors¹¹). They are applied more and more in telephony on circuits which would otherwise be useless. In fact a non-linear characteristic is necessary for this purpose, but it must be used in such a way that non-linear distortion of the signal is avoided. This seems a contradiction at first sight, but it is attainable by using different time constants. There are two possibilities: either the signal itself, or a separate pilot signal may be employed in controlling this non-linear mechanism. The first case is applied to normal companadors, where the compression ratio is relatively low,

¹¹) N. Valentini, The dynamics compressor-expander (companion) in telephony, Telettra No. 2, 12-22, Sept. 1954.

mostly two to one in a decibel scale. The second case — that of the pilot compandor — can be used with a very high compression ratio and allows transmission via circuits suffering from a very high noise level. As the operation of pilot compandors can be made independent of the received signal level, these are very suitable at the same time for compensating “fading” effects.

Though it is not our object to go into the details of a compandor, its effect can best be illustrated by means of the following example. If we suppose that the peak value of a speech signal is only two or three times as large as the mean value of the noise, then it will be clear that many parts of the speech signal are drowned by the noise. By using a compandor the noise can now be decreased more or less in proportion to the mean amplitude of the speech, so that also signals of small amplitude — often carrying much information — can rise above the noise. Indeed it has been found that if the accompanying noise is modulated in the same rhythm and in proportion to the amplitude of the speech, then noise and speech voltages may be of the same order before intelligibility is destroyed ¹²⁾.

An important aspect of noise reduction by means of these compandors is that practically no extra bandwidth is necessary here, contrary to the older method of frequency modulation where a proportional increase of bandwidth is required. This means that compandors can also be used in the future for making available more telecommunication channels using radio waves.

Further, the combination of a compandor with other methods of modulation, such as pulse-code modulation, has turned out to be very fruitful. In this case it is only the quantization noise, and not transmission noise, which is of importance and which is reduced in the output signal: this enables us to apply lower pulse frequencies. If we may mention delta modulation again, it was found quite recently by Greefkes of this laboratory that the feedback circuit which is already used in delta modulation can, at the same time, perform the function of a compandor ¹³⁾. In that case two different information signals are transmitted by the generated pulse series, one relating to the detailed structure of the speech wave, and the other to the level of the speech. Referring again to the model of a car following a road, it might be said that the sensitivity of the steering mechanism is adjusted here to the mean

curvature of the road. With this system of “continuous delta modulation”, as it was called, it was found possible to reduce the necessary pulse frequency from about 60 to 30 kc/s and even with the rather low pulse frequency of 16 kc/s a 20 dB change in the level of the input signal could still be tolerated.

In this way compandors are used in a process of “adaptation” that is largely analogous to adaptation in the biological sense. It enables us to handle information signals of widely varying signal strengths, in more or less the same way as the intensity of light entering the human eye is controlled by a process of adaptation. We can still learn much from these biological processes, for instance as regards the time constants involved, and the research on different aspects of perception and adaptation will, therefore, be of great value for the design of optimum communication systems for use in the future.

Compatibility

In the laboratory it has been shown that a great variety of compandors is useful for noise reduction, each having its specific advantages. In “point-to-point” communication one is free to choose any of these. However, for general application, one is faced now with the difficult problem of “compatibility”. It is desirable that most transmitters and receivers can communicate without much alteration of apparatus and this was comparatively easy when only amplitude and frequency modulation were employed, but with the introduction of new ideas the technical developments have diverged so much that achieving compatibility is an almost hopeless task. This is one of the main reasons why many improvements in the method of modulation have been held back for a long time before being put in practice. The need for more and better communication links, to be expected in the near future, will soon force us, however, to apply only the best principles that are known ¹⁴⁾.

Considering the last mentioned aspect of compatibility, an important problem is found in broadcast transmission. Most of the transmitters apply amplitude modulation and all receivers are constructed for this type of modulation. However, in this case half of the bandwidth is wasted by using the symmetrical waveform that is generated in the normal way. It has already been known for 40 years that single-sideband transmission is to be preferred as regards power and frequency demands,

¹²⁾ F. de Jager and J. A. Greefkes, “Frena”, a system of speech transmission at high noise levels, Philips tech. Rev. 19, 73-83, 1957/58.

¹³⁾ Belgian Patent No. 620450.

¹⁴⁾ G. Jacobs, Radio interference — suicide or challenge, IRE Trans. on radio frequency interference RFI-4, No. 2, 21-23, 1962.

and so this method is commonly applied in systems of carrier telephony. However, a somewhat more complicated receiver is necessary in this case and this difficulty has excluded the use of single-sideband transmission for broadcast purposes. Several solutions have been proposed for generating a "compatible" single-sideband signal capable of being received by all normal broadcast receivers, but the distortion involved is usually too large for high-quality reproduction of music. Van Kessel and others from this laboratory have recently found an elegant solution to this problem, based on a careful mathematical analysis. A normal broadcast transmitter need be only slightly changed here in order to produce single-sideband components of such phase and amplitude that any receiver detects the signal wave form without distortion¹⁵). The result is a more effective use of transmitter power and a better reproduction of the higher frequencies in any conventional receiver.

Television

Another field in which it is worth-while to look for the essential information is that of television. Although television pictures are two-dimensional, their information is transmitted from point to point, which requires the rather large bandwidth of about 5 Mc/s. This bandwidth (sufficient for more than 1000 telephone channels) would permit us to change the whole television picture 25 times per second, which would neither be perceived nor enjoyed by any human observer. Considering the statistics of the signal we find, again, that most parts of the signal are repeated many times. In reducing the amount of available information one is faced, however, with the difficulty that the observer's attention may be directed to any part of the picture, and this requires the picture quality to be maintained at a sufficient high degree over the whole area.

For reducing bandwidth one can in this case either use the correlation which is found to exist in time, or in place, of the succeeding signals, or one can make use of the slowness with which brightness variations can be perceived by the human eye. Several systems have already been investigated, for instance by Teer in this laboratory¹⁶). Some results have been found with respect to colour television. For instance, it was found possible to add colour

information without increasing the bandwidth. In that case a number of signals, containing information about the different colours, can be transmitted simultaneously in the same frequency band, and it is possible to suppress mutual interferences by choosing special phase and frequency relations of the subcarriers used¹⁷). Thanks to the integrating effect of the human eye the picture appears natural. As the system is compatible, a picture in black and white can still be obtained with a conventional receiver.

Data transmission

So far we have considered the modern aspects of transmitting continuously varying information signals such as speech, music or television. During the last decade, however, a growing interest is found in the transmission of telegraphic signals, especially in the form of high-speed telegraphy or data transmission. Indeed, with the growing application of computers, equipment for data-processing, etc., a large amount of digital information is available now which often needs to be transported over long distances. By means of the thousands of pulses that can be transmitted per second via a telephone circuit, the speed of transmitting this information can be much higher, of course, than with a normal telephone conversation.

So we find more and more that the telephone network, originally designed for speech, is used now by machines which are talking in digital languages. An air-line reservation system, answering questions within one second, represents a typical example of this kind. Reliability is often of primary importance in these cases and so error-detection and error-correction schemes are usually applied. Just as the redundancy which is present in normal speech was found to be very useful for a reliable transmission, we may apply a redundancy in this digital information by introducing check-bits. It was pointed out by Golay¹⁸), and in more detail by Hamming¹⁹), that by applying special codes, these check-bits could not only be used for detection of errors, but also for automatic error correction. Nowadays a great interest is taken in the scientific aspects of error-correcting codes and it is a remarkable fact that many ideas which have originated in a branch of pure science, namely number theory, find here practical applications²⁰).

¹⁵) Th. J. van Kessel, F. L. H. M. Stumpers and J. M. A. Uyen, A method for obtaining compatible single-sideband modulation, E.B.U. Rev. Part A, No. 71, 12-19, 1962.

¹⁶) K. Teer, Investigations into redundancy and possible bandwidth compression in television transmission, Philips Res. Reports 14, 501-556, 1959, and 15, 30-96, 1960.

¹⁷) K. Teer, Colour-television transmission, Electronic and Radio Engr. 34, 280-286, 326-332, 1957.

¹⁸) M. J. E. Golay, Notes on digital coding, Proc. IRE 37, 657, 1949.

¹⁹) R. W. Hamming, Error detecting and error correcting codes, Bell Syst. tech. J. 29, 147-160, 1950.

²⁰) W. W. Peterson, Error-correcting codes, M.I.T. Press, 1961.

In this way modulation theory is now closely related to the storage of information in memories, the way of handling this information by means of logical circuits and the problem of performing these functions in the most effective manner. Again the study of statistics has entered the field, though it is not the statistics of the signal to be transmitted, but the statistics of the transmission path which is now of primary importance.

The transmission of pulses through the normal telephone network, however, presents its own difficulties. This network has been designed for speech transmission and many peculiarities which do not have the slightest effect on speech (such as the suppression of very low frequencies or the influence of phase distortion) may do much harm in the case of pulse transmission.

If a new telephone network were to be designed, then the requirements for digital transmission would certainly be taken into account. But the large investments in telephone circuits all over the world force us to work the other way round. So we have to look for modulation methods which permit the transmission of digital conversations via speech circuits, in a more effective manner than is mostly used nowadays. One might even say that the possibility of using and connecting data-processing equipment in the future depends on the development of better methods of modulation, giving faster and more reliable operation. At present there is much research in this field. Consider speed, for example; a system was built in this laboratory for transmitting 4000 bits per second, using only part of the available bandwidth in a normal speech circuit and transmitting information at a rate very close to the theoretical maximum²¹⁾; see *fig. 9*.

Many problems, however, remain to be solved before the specific troubles on telephone connections are overcome. Some of them are related to the problems of radio transmission, others are of a quite specific nature. For instance the large and unpredictable phase distortions which are encountered on switched telephone connections present serious difficulties. In this case the disturbing voltages are generated by the signal itself, and so improving the signal-to-noise ratio by increasing the signal power is of no use here, as the interfering voltages are growing proportionally. So we must look for other means which are effective and not too complicated.

Looking into the future

Comparing modulation today with that of fifty years ago, we find that we can do much more; but we are also faced with many more problems. Development in the next fifty years may be expected, I think, to follow two different trends. The first is the more practical application of theoretical results obtained in "information theory". This will probably lead to the design of systems that approach more and more the theoretical limits, thereby increasing reliability and using less power and less frequency space than is found in existing systems. (In satellite communications for instance this is very urgent. By making use of special amplifiers — masers — one has already nearly reached the physical limit of attainable signal-to-noise ratio.) The second aspect is the closer relation between man and machine, such that the possibilities of adaptation and control in electrical systems can be used in relation with the sensory behaviour of man. For the design of these adaptive systems we need many more results of the research in "perception" (as investigated for instance by Schouten and his co-workers at the "Instituut voor Perceptie Onderzoek"²²⁾). The combination of the results in these two fields — information theory and perception — will probably be the most important characteristic of future development. So far the combination of the two has presented some difficulty, which may be due to the fact that problems and results in the first field are usually considered as "exact", whereas the results in the second field are of a more "subjective"

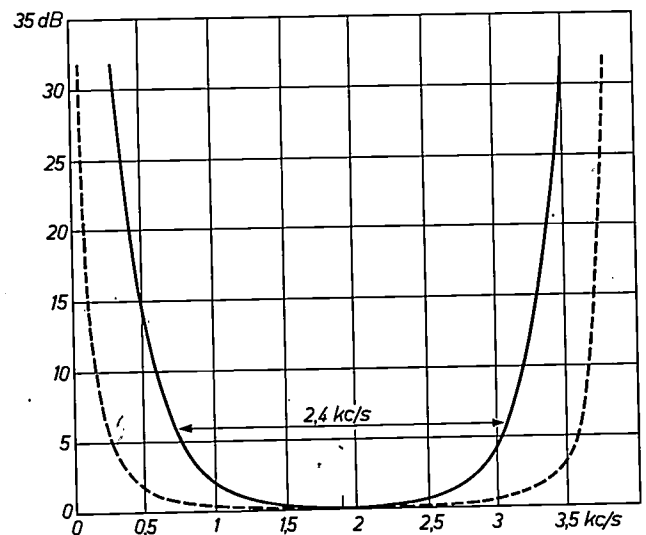


Fig. 9. Bandwidth used for transmission of data signals at a speed of 4000 bits/sec, compared with the bandwidth (dotted line) of a normal telephone channel (see²¹⁾).

²¹⁾ F. de Jager and P. J. van Gerwen, CO-modulation, a new method for high-speed data transmission, IRE Trans. on information theory IT-8, S 285-S 290, 1962 (No. 5).

²²⁾ J. F. Schouten, The Institute for Perception Research, to appear in a forthcoming issue of Philips tech. Rev. devoted to this Institute.

nature. But if enough people can be found who are interested in both of them, we may expect revolutionary developments in modulation systems for the future.

During many years engineers have transmitted electrical signals carrying information, without considering the type of information contained in them. Now research is directed towards finding the "information-bearing elements" of these signals. This is not a question of academic interest only. Indeed, the need for improving modulation systems is so very urgent, in view of the acute expansion in telecommunications, that engineers are forced to study these fundamental aspects of human communication more seriously. The solution of these problems, however, will require more mathematical, physical and biological analysis — together with the art of good guessing.

The large number of communication channels spread over the world (and sometimes even reaching other planets) may be looked upon as an enormous nervous system which mankind has built in only half a century and which is going to be used more and more. By means of modulation theory we try to discover how to use these nerves most efficiently.

Summary. In this article it is shown how the ever increasing demands on communication links have led to the introduction of new methods of modulation. After considering frequency modulation, pulse code modulation and delta modulation special attention is given to those methods of modulation which make use, to a large extent, of the statistical properties of the information signals: bandwidth compression, companders. Modern research is directed here towards finding the information-bearing elements in these signals. Finally some problems related to compatibility, and to the transmission of digital information, are considered.

It is expected that in the future the combination of results found in information theory and in perception research will lead to the design of still more effective methods of modulation.

CONTRIBUTIONS OF THE PHILIPS RESEARCH LABORATORIES TO SOLID-STATE CHEMISTRY

by H. J. VINK *).

54-16

Introduction

In the past four decades scientists in the Philips Research Laboratories have been engaged in various areas of research in chemistry such as colloidal chemistry, surface phenomena, electrochemistry, organic chemistry, photo-chemistry, analytical chemistry, biochemistry, glass technology and solid-state chemistry. Especially the field of solid-state chemistry has been extensively studied. At this Symposium it seemed therefore worthwhile to take a retrospective view of the work in this field; to contemplate the problems that have been posed, the methods that have been used to solve them and the results that have been achieved.

The choice of the various topics of research is naturally determined by the fact that this laboratory serves the electronics industry. Chemistry has always played an important role in the development of electrical engineering and electronics, but it is especially in the last 30 years that the importance of chemistry to electrical engineering and even more to electronics has been made abundantly clear. The control and preparation of materials has always been of prime importance but the way chemistry nowadays is able to develop materials with new combinations of known physical properties or materials possessing physical properties that were hitherto unknown or did not occur in such a high degree makes it indispensable for the electronics industry. It is not going too far to state that in chemistry lie the foundations of the rapid expansion and increasing importance of the electronics industry.

However the technical importance of the various studies undertaken in this laboratory will only lightly be touched upon in the course of this review. A more general point of view will be taken here, viz, the interrelation between physical and chemical properties and exact chemical composition.

The physical and chemical properties of a substance depend on its chemical composition. On the other hand the electrical, optical, magnetic and chemical properties can often be conveniently used to determine the chemical composition. Therefore an interplay of physical and chemical methods can often lead to an understanding of the interrelation between physical and chemical properties and exact

chemical composition. This interplay requires not only a close cooperation between chemists and physicists, but also, for both of them, a certain ability to understand each other's "language", methods and problems.

Research in this field is largely governed by two questions:

- 1) What kind of sites in what kind of crystal structures must be occupied by what kind of atoms (or groups of atoms) in what state of valency in order to provoke a certain (desired) combination of physical and chemical properties?
- 2) Is it possible to develop certain groups of compounds and ways of preparing them such that desired concentrations of specific kind(s) of atoms each in a proper state of valency do indeed occupy the right kind of site, in order to control reproducibly certain desired physical and chemical properties?

In this paper the approach towards these questions will have a chemical bias and will be divided into four main areas, viz, a) crystal chemistry, b) internal charge compensation, c) gases and metals, and d) thermodynamics.

Crystal chemistry

Cation distribution in spinels; experimental determination

An investigation that proved to be the starting point of a whole array of studies was begun by Verwey¹⁾ 2) 3) 4) 5) in the middle thirties^{**}). The then known oxides of iron and other elements of the transition series showed quite a variety of interesting electrical and magnetic properties. Especially the fact that some of the magnetic oxides were also insulators made it appear worthwhile to start an investigation of these related oxidic compounds (Snoek⁶⁾). This investigation soon centred around the spinels, a large proportion of which can be described by the general formula XY_2O_4 . The spinel structure is cubic and can be characterized by a cubic, almost close-packed arrangement of oxygen (O^{2-}) ions in which the metal ions lie on certain sites surrounded tetrahedrally by four O^{2-} ions and on

*) Philips Research Laboratories, Eindhoven.

**) Literature references are listed at the end of this article.

certain octahedral sites surrounded by six O^{2-} ions whose centres form the apices of a slightly distorted octahedron (fig. 1). The packing of the O^{2-} ions can be described by a parameter u . For $u = 0.375$ ($3/8$) the oxygen ions form an ideally cubic close-packed arrangement. Usually the value of u is somewhat larger than 0.375.

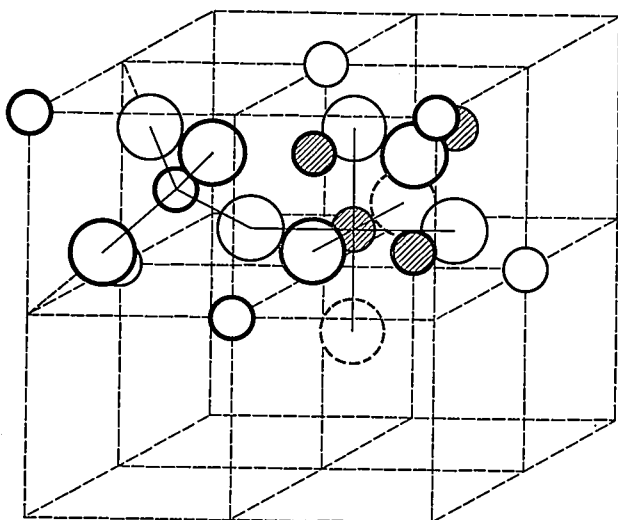


Fig. 1. Unit cell of spinel structure. The position of the ions in only two octants is shown. The dashed circles belong to other octants. The drawn lines indicate the fourfold and sixfold coordination of the respective metal-ion positions. Large circles: oxygen ions; small hatched circles: metal ions on octahedral sites; small unhatched circles: metal ions on tetrahedral sites. The figure is drawn for $u = 0.375$.

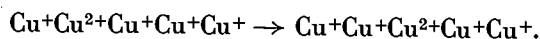
The distribution of the cations among the sites available may be:

- "Normal" with X in the tetrahedral positions and $2Y$ in the octahedral positions. This may be formulated by $X[Y_2]O_4$, using brackets to indicate cations in the octahedral positions.
- "Inverse" $Y[XY]O_4$.
- "Mixed" $X_x Y_{1-x}[X_{1-x} Y_{1+x}]O_4$, with $0 < x < 1$. The metal ions that together with oxygen can form spinel structures must have radii between 0.5 and 1.0 Å, but as there seems to be no restriction as regards valency, a great number of spinels exist.

It soon appeared that the magnetic and electrical properties and their combination depend strongly on the distribution of cations over the two types of sites. It is therefore of interest to know whether there is a relation between the kind (and valency) of an ion and a possible preference for one of the two types of site. The programme, therefore, was first to find out experimentally, for a number of spinels, how the ions were distributed and secondly to give an explanation for these experimental facts so that predictions could be made of the cation distribution and hence the physical properties of other spinels. With

regard to the first part of this programme not only binary spinels but also mixed solutions between them, ternary or quaternary spinels had to be investigated in order to determine the relative preference of the various ions for a specific site.

X-ray diffraction immediately suggested itself for this kind of study. The distribution of the cations in a number of spinels has indeed been determined in this way, especially by comparing ratios of intensities of pairs of reflections that are more sensitive to a change in distribution parameter than to a change in u . However in view of the small difference in scattering power between the metal ions of the first transition series or between different ions (Fe^{2+} , Fe^{3+}) of the same atom, this method can only be used for a limited number of cases. For other spinels the study of the electrical and/or magnetic properties can very often be used to determine the cation distribution. This is an example of the fact, mentioned in the introduction, that physical properties that are to be controlled by the exact chemical composition are helpful in their turn in the determination of the chemical composition. The electrical conductivity for instance was used to establish the distribution of the Fe^{2+} and Fe^{3+} ions in Fe_3O_4 ^{7) 8)}. This spinel has a very high electronic conductivity of the order of $\sigma = 10^2 \Omega^{-1}cm^{-1}$ at room temperature. Now it was known that substances like Cu_2O which, when pure, show a high resistivity, could have their electronic conductivity considerably increased by a slight deviation from stoichiometry, for instance an oxidation. This had been explained by Wagner and co-workers by assuming that in oxidized Cu_2O not only Cu^+ ions are present but also Cu^{2+} ions. As a consequence of this presence of homonymous ions in a different valency state, electronic conductivity was possible because electronic transport could now occur by exchange of valency electrons:



Verwey and De Boer ^{7) 8)} explained the high electronic conductivity of Fe_3O_4 by assuming it to have the inverse structure $Fe^{3+}[Fe^{2+}Fe^{3+}]O_4^{2-}$ rather than the normal structure $Fe^{2+}[Fe^{3+}Fe^{3+}]O_4^{2-}$. As a matter of fact, the Fe ions of different valency occupying one type of crystallographic site (octahedral) will be more similar to the situation in oxidized Cu_2O than the case of the normal spinel structure where Fe^{2+} and Fe^{3+} ions occupy different crystallographic sites, viz, the tetrahedral and octahedral sites respectively.

Another physical property that can often be used is the saturation magnetization at 0 °K. Néel's theory of ferromagnetism postulated a strong

negative exchange interaction (TO) between ions on tetrahedral (T) and those on octahedral sites (O) and weaker negative TT and OO interactions. In simple cases, i.e. for which the TO interaction is dominant, it is possible to find the saturation magnetization (σ) at 0 °K, expressed in Bohr magnetons (μ_B) per formula unit by simply subtracting the sum of the ionic magnetic moments occurring on the T sites from those on the O sites or vice versa. Let us take as examples the compounds $MgFe_2O_4$ and $NiFe_2O_4$ both of which, because of their X-ray diffraction patterns, were thought to be inverse spinels $Fe^{3+}[Mg^{2+}Fe^{3+}]O_4^{2-}$ and $Fe^{3+}[Ni^{2+}Fe^{3+}]O_4^{2-}$. However the accuracy of these determinations was not very high. Therefore the saturation magnetization was determined for these compounds in order to substantiate the rather vague results of the X-ray diffraction investigations. According to these results and applying Néel's rules, the saturation magnetization values one expects to find are $\sigma = 0$ for $Fe^{3+}[Mg^{2+}Fe^{3+}]O_4^{2-}$ and $\sigma = \mu_{Ni^{2+}} = 2\mu_B$ for $Fe^{3+}[Ni^{2+}Fe^{3+}]O_4^{2-}$. However the values found were $1.1 \mu_B$ and $2.3 \mu_B$ respectively.

Now the value $\sigma = 1.1 \mu_B$ for $MgFe_2O_4$ could be changed by different heat treatments. For instance quenching from 1250 °C increased its value to $\sigma = 1.4 \mu_B$. Gorter⁹⁾ explained these results by assuming that Mg-ferrite actually shows a "mixed" distribution, $Fe_{1-x}^{3+}Mg_x^{2+}[Mg_{1-x}^{2+}Fe_{1+x}^{3+}]O_4^{2-}$ with $x = 0.11$ and 0.14 respectively. The same explanation however does not apply to Ni-ferrite. Its saturation magnetization does not change with different heat treatments. Gorter therefore had to assume that $x = 0$ for Ni-ferrite, and further that the magnetic moment of Ni^{2+} also had a contribution from the orbit leading to a total spin moment of $2.3 \mu_B$. This value of the total spin moment of Ni^{2+} has been amply verified in other experiments. These simple and well known cases are given here to demonstrate how careful one must be in interpreting physical phenomena into chemical formulae. The facts that 1) ions can occur in more than one (unknown) state of valency, 2) the distribution in ternary systems has to be described by two parameters x and y and 3) the dilution by non-magnetic ions interferes in an unknown way with the interaction between the magnetic ions, all underline the care necessary in interpretation.

Cation distribution in spinels; theoretical explanation

Romeijn^{10) 11)}, Gorter^{12) 13) 14)} and Blasse¹⁵⁾ prepared a great number of binary, ternary and quaternary spinels and studied the cation distribution using the methods mentioned above. Their com-

bined results proved to be of an enormous complexity. This is not very surprising however as there are a great number of different kinds of energies to be taken into account. In the course of the years the following contributions to the lattice energy have been considered: a) electrostatic Coulomb-forces, b) Born repulsion, c) ordering, d) individual site preference of some cations having a tendency to be surrounded by a certain configuration of oxygen ions, e) crystal-field effects, f) ligand-field effects and g) polarization.

Of the corresponding energies only the Coulomb energy^{16) 17)} and the ordering energy¹⁸⁾, by making certain assumptions, could be calculated with any accuracy (Verwey, Heilmann, De Boer and Van Santen). It appeared however that this Coulomb energy depends on the values of the oxygen parameter u and the length of the unit cell edge a (both themselves depending on the actual distribution) to such an extent that it proved to be virtually impossible to predict any cation distribution from this effect alone. Geometrical arguments replacing calculations of Born repulsion energies¹⁰⁾ also did not lead to any conclusion. Further, the energies calculated for ordering phenomena proved to be a substantial fraction of the Coulomb energy, for cases where one or both of the sublattices are occupied by ions differing in charge. This means that there must exist a tendency to the formation of superstructures in these spinels. Such a superstructure will generally be accompanied by lattice distortions. Superstructures have indeed been found in a few cases: $Fe^{3+}[Fe^{2+}Fe^{3+}]O_4^{2-}$ ¹⁹⁾, $Fe^{3+}[Li_{0.5}^+Fe_{1.5}^{3+}]O_4^{2-}$ ²⁰⁾, $Fe^{3+}[V_{1/3}Fe_{5/3}^{3+}]O_4^{2-}$ ²¹⁾, $Li_{0.5}^+Fe_{0.5}^{3+}[Cr_2^{3+}]O_4^{2-}$ ¹⁴⁾ and $In_{2/3}^{3+}V_{1/3}[In_3^{3+}]S_4^{2-}$ ²²⁾, V being the symbol first for a vacancy on an octahedral and then on a tetrahedral site. However the transition temperatures of these ordered phases to a disordered one are too low compared with the high ordering energies involved. It has been argued by Van Santen²³⁾ that in ionic crystals the energy difference between long-range and short-range order is small, due to the fact that the Coulomb interactions are long-range interactions. This implies that the transition temperature will be low and also that above this temperature a fair degree of short-range disorder will still exist.

Already Verwey and Heilmann pointed out that apart from these electrostatic effects, individual site preferences of some cations must also be taken into consideration, such as the tendency of Zn^{2+} , Cd^{2+} , Ga^{3+} , In^{3+} (and Ge^{4+}) to be surrounded by four anions (sp^3 hybridization).

On the other hand it was suggested by Van Santen and Van Wieringen²⁴⁾ and by Romeijn¹¹⁾

that the behaviour of ions of the transition series with partially filled shells such as Ti^{3+} , V^{3+} , Cr^{3+} , Mn^{4+} , Mn^{3+} , Fe^{3+} , Mn^{2+} , Co^{3+} , Fe^{2+} , Co^{2+} , Ni^{2+} , Cu^{2+} would be also governed by energy gained by orbital splitting due to the crystalline field.

Other workers made quantitative calculations about this effect. For the ions of the transition series they calculated energies of octahedral site preferences. From these values it could be deduced that ions with $3d^5$ and $3d^{10}$ configuration (Mn^{2+} , Zn^{2+}) should show no site preference at all, whereas the ions with $3d^3$ and $3d^8$ configuration (Cr^{3+} and Ni^{2+}) would show a preference for octahedral sites, in good agreement with experimental data. It is very unsatisfactory however that the cation distribution should be controlled by energies which at their highest (Cr^{3+} : 40 kcal/mol) have values less than 1% of the electrostatic Coulomb energies.

Blasse¹⁵⁾, using the method of molecular orbitals, combined the effects of specific tendency to covalent binding and of crystal field stabilization. This approximation — the ligand field theory — leads to the following result, differing somewhat from that of the authors just mentioned: ions with $3d^3$ configuration prefer octahedral sites, ions with $3d^5$, $3d^6$, $3d^7$, $3d^9$ or $3d^{10}$ configuration prefer tetrahedral sites, whereas ions with $3d^0$, $3d^1$, $3d^2$, $3d^4$ and $3d^8$ do not show any pronounced preference. These results are of a qualitative nature only. It is therefore impossible to compare the energies due to this kind of arguments with those of the Coulomb interactions. It is reasonable to assume that they will be larger than the energies calculated for the crystal field stabilization but they will again be small compared to electrostatic energies.

Blasse also brought anion polarization into the picture; and he was able to show that polarization of the oxygen ions as a consequence of these being asymmetrically surrounded by differently charged cations can be decisive for the actual cation distribution.

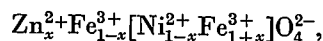
For instance in $Fe^{3+}Ni^{2+}Rh^{3+}O_4^{2-}$, from every other point of view one should expect the configuration $Fe^{3+}[Ni^{2+}Rh^{3+}]O_4^{2-}$. Actually however, part of the Ni^{2+} ions occupy tetrahedral sites, in this way increasing the asymmetry with respect to charge of the immediate surroundings of the oxygen ions and thus increasing the energy gained by polarization.

Cation distribution in spinels; molecular engineering

Although therefore from a theoretical point of view the situation around the cation distribution in spinels as a whole leaves something to be desired, yet

the chemist now has enough experimental data and theoretical considerations at his disposal to be able to perform "molecular engineering" with this structure. In this way a great many products could be developed differing in initial permeability, saturation magnetization, electrical resistivity, electrical losses, magnetic losses etc. It appeared that for all spinels the product of the static initial permeability and the dispersion frequency is substantially constant (Snoek's relation²⁵⁾). We will not go into all this but give just one very simple example of molecular engineering in this field of magnetic materials.

When $ZnFe_2O_4$ ($Zn^{2+}[Fe^{3+}Fe^{3+}]O_4^{2-}$), having a total magnetic moment of zero, is introduced into $NiFe_2O_4$ ($Fe^{3+}[Ni^{2+}Fe^{3+}]O_4^{2-}$), it gives rise to mixed crystals with a higher magnetic moment than $NiFe_2O_4$ ²⁶⁾. The distribution of the different ions in the mixed crystals, according to the rules mentioned above will be



leading to a saturation magnetization per molecule of

$$(1-x)\mu_{Ni^{2+}} + 2x\mu_{Fe^{3+}}.$$

This is also the case when $ZnFe_2O_4$ is introduced into other spinels (fig. 2).

For higher concentrations this effect disappears,

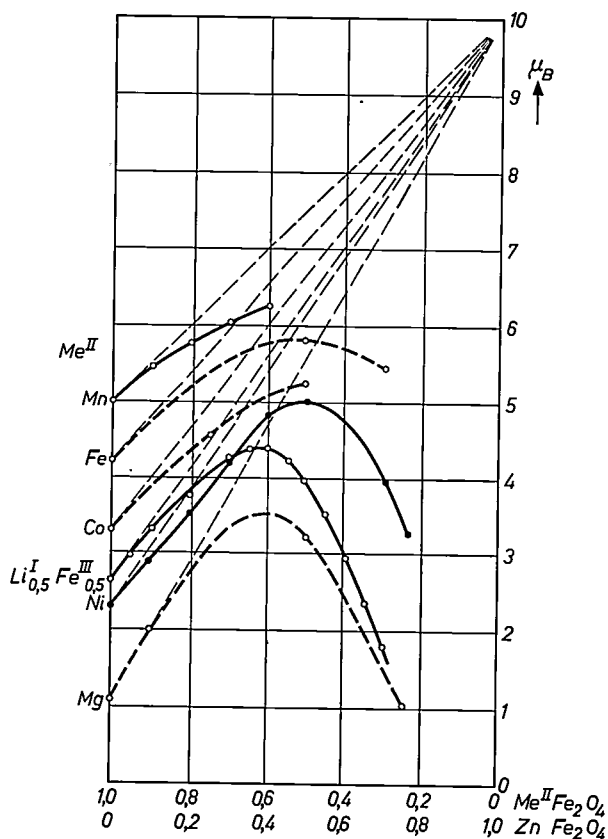


Fig. 2. The saturation magnetization in Bohr magnetons per "molecule" for mixed crystals of Zn ferrite with other ferrites.

the dilution with non-magnetic ions on the tetrahedral sites making the TO exchange small compared with the OO and TT interactions.

Hexagonal and trigonal ferrites; determination of their structure and composition

In extending the study of spinels Braun, Jonker and Wijn found an interesting group of new compounds in the system BaO - MeO - Fe₂O₃ (fig. 3)^{27) 28) 29) 30) 31) 32) 33)}, where Me stands for a divalent metal like Fe, Mn, Co, Zn and Mg.

These ferrites all crystallize with hexagonal crystal structures. Their compositions are given by the points M, W, X, Y, Z and U in fig. 3 and can be represented by the formulae:

M	BaFe ₁₂ ^{III} O ₁₉
W	BaMe ₂ ^{II} Fe ₁₆ ^{III} O ₂₇
X	Ba ₂ Me ₂ ^{II} Fe ₂₈ ^{III} O ₄₆
Y	Ba ₂ Me ₂ ^{II} Fe ₁₂ ^{III} O ₂₂
Z	Ba ₃ Me ₂ ^{II} Fe ₂₄ ^{III} O ₄₁
U	Ba ₄ Me ₂ ^{II} Fe ₃₆ ^{III} O ₆₀

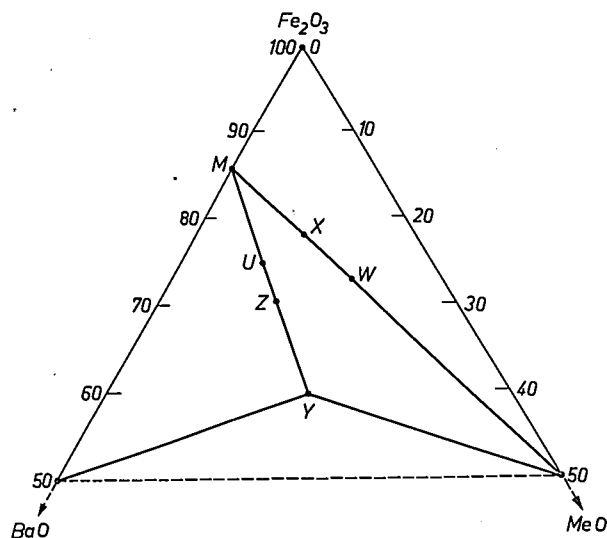


Fig. 3. Part of the phase diagram BaO-MeO-Fe₂O₃, where Me stands for a divalent metal such as Mn, Fe, Co, Ni, Zn or Mg.

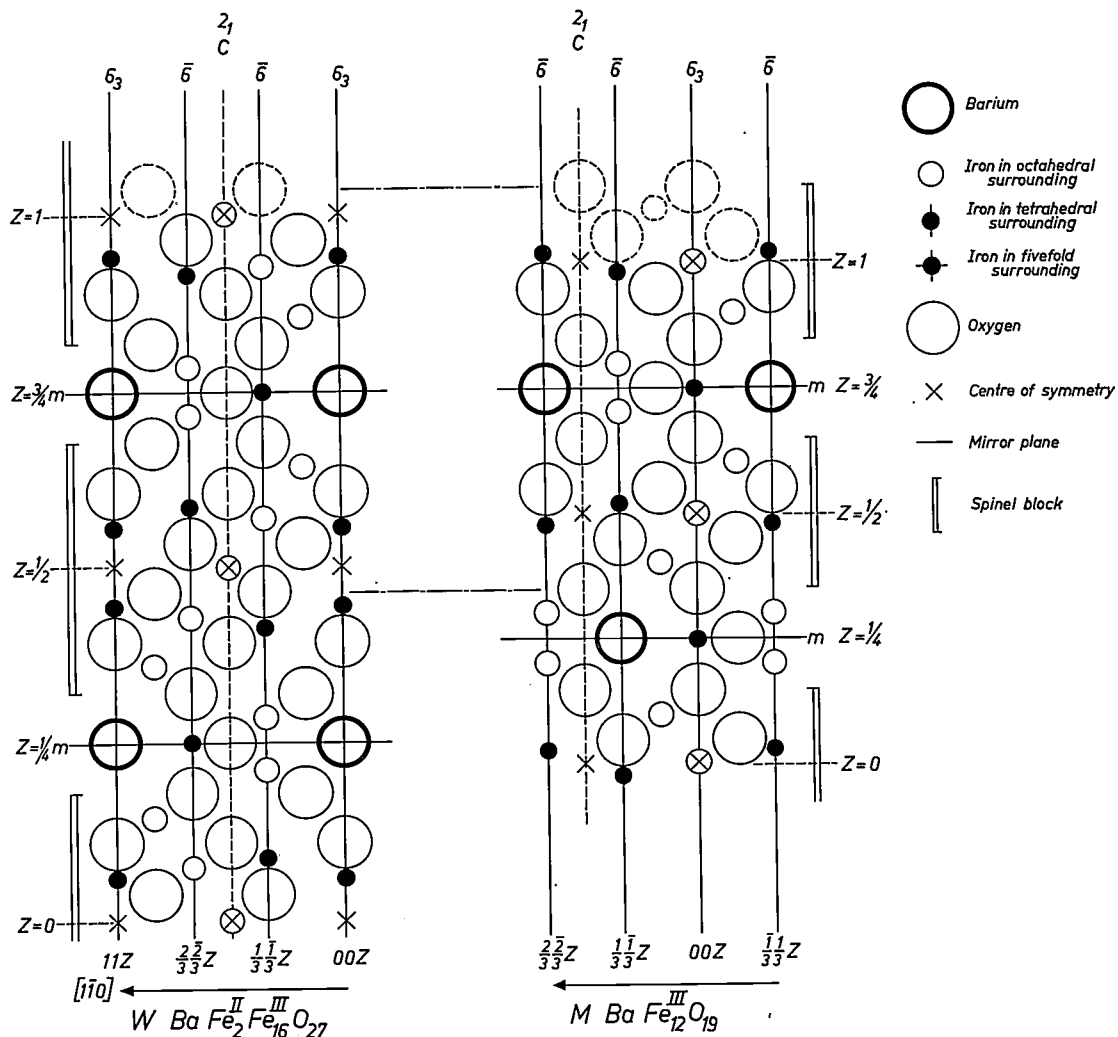


Fig. 4. Unit cell of W compared with that of M. Only atoms and symmetry elements in a mirror plane containing the c-axis are shown. The relative positions of all atoms in nine layers

and those of the interstitial atoms between are identical in the two structures. The spinel blocks include 4 oxygen layers in M and 6 oxygen layers in W.

The preparation of those crystals often leads to a mixture of a number of compounds. Therefore it would have been impossible to determine the very existence and the precise composition of these compounds by the usual chemical methods. This study however was greatly facilitated by the fact that contrary to the magnetic spinels these hexagonal compounds show a strong magnetic anisotropy. By making use of this effect the finely powdered reaction products could be oriented in such a way that some of the compounds would have their hexagonal axis parallel and others perpendicular to an applied magnetic field. An X-ray diffraction pattern of such an oriented powder is greatly simplified. The different reaction products could now be "labelled" so to speak and in this way it was possible to improve the reaction procedure such that single crystals of the various compounds could be obtained. Braun ³⁰⁾ ³¹⁾ was then able to determine the crystal structure of these compounds and from this their chemical composition. The complicated chemical formulae of these compounds have thus been determined not directly

which would have been impossible, but by means of X-ray diffraction. This again exemplifies how the chemist can make use of physical properties to determine not only the composition but the very existence of complicated compounds.

In the determination of these structures use was made of the fact that a strong relationship appeared to exist between the crystal structure of these hexagonal ferrites and the well known spinel structure. These structures can be described in different ways ³⁴⁾. In this paper the notation of Braun ³¹⁾ will be followed. In the cubic close-packing of the O^{2-} ions in spinels, trigonal O^{2-} layers perpendicular to a [111] direction have a mutual distance of 2.33 Å. In the unit cell (32 O^{2-} ions) there are six of these layers.

Now the crystal structures of M, W, and X also contain such close-packed layers. However in some of these layers (the "Ba-layers") one quarter of the O^{2-} ions is substituted in an ordered way by Ba^{2+} ions that are about as large as the O^{2-} ions. The compounds can generally be described by an alternation of blocks with pure spinel structure and consisting of 4 or 6 layers followed by blocks of 1 or 2 "Ba-layers". The "Ba-layers" are connected with the spinel layers in a "hexagonal" way. For instance in the structure of M one can distinguish blocks of 4 layer with cubic close-packing of O^{2-} ions and a spinel-like distribution of Fe^{3+} ions over the tetrahedral and octahedral spinel sites between the O^{2-} ions (S_4) followed by one Ba-layer (B_1) also containing Fe ions. The unit cell consists of two equivalent parts each containing 5 layers with, in total, 19 O^{2-} , 1 Ba^{2+} and 12 Fe^{3+} ions, and with a thickness of $5 \times 2.33 = 11.6$ Å. Similarly the compound W contains blocks of 6 spinel layers (S_6) alternating with a "Ba-layer" (B_1). Here the spinel layers, as well as Fe^{3+} ions, also contain 2 Me^{2+} ions (see fig. 4). X consists of a regular alternation of layers with M and W structures. Therefore one has:

$$\begin{aligned} M &= M_5 = B_1 S_4 B_1 S_4 B_1 \dots \text{thickness } 11.6 = 5 \times 2.33 \text{ \AA} \\ W &= W_7 = B_1 S_6 B_1 S_6 B_1 \dots \text{thickness } 16.4 = 7 \times 2.34 \text{ \AA} \\ X &= X_{12} = B_1 S_4 B_1 S_6 B_1 \dots \text{thickness } 28.0 = 12 \times 2.33 \text{ \AA} \end{aligned}$$

In the compounds Y, Z, U one encounters as a new building block, that of the double Ba-layer (fig. 5):

$$\begin{aligned} Y &= Y_6 = B_2 S_4 B_2 \dots \text{thickness} \\ &14.7 = 6 \times 2.42 \text{ \AA} \\ Z &= Z_{11} = B_2 S_4 B_1 S_4 B_2 \dots \text{thickness} \\ &26.2 = 11 \times 2.38 \text{ \AA} \\ U &= U_{16} = B_2 S_4 B_1 S_4 B_1 S_4 B_2 \dots \text{thickness} \\ &37.8 = 16 \times 2.36 \text{ \AA} \end{aligned}$$

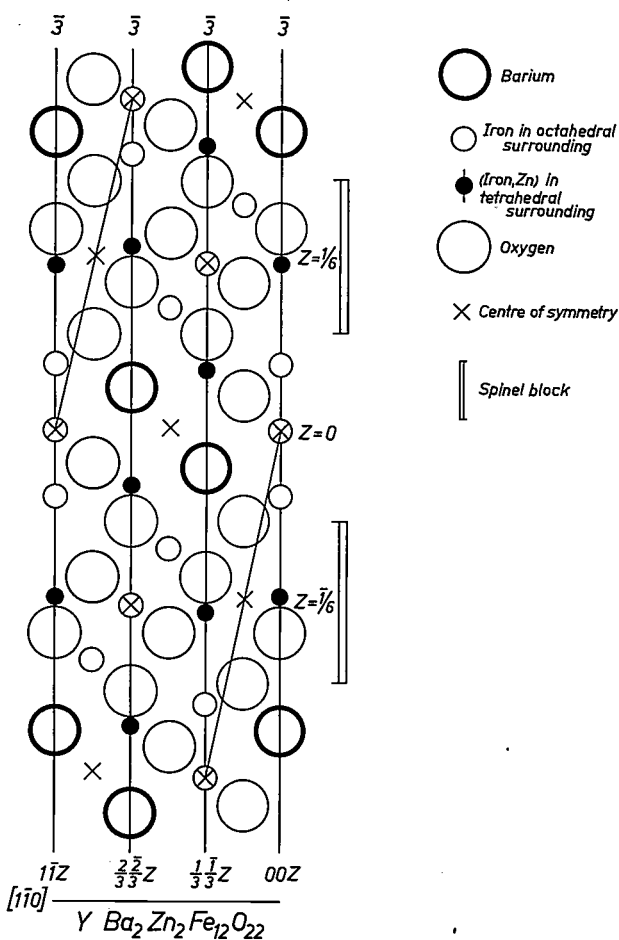


Fig. 5. Two thirds of a unit cell of Y. Only the atoms and symmetry elements in a mirror plane containing the c-axis are shown. The spinel blocks include 4 oxygen layers. The rhombohedral repeat distance is indicated.

Hexagonal and trigonal ferrites; magnetic properties

An elegant confirmation of the structure of these compounds is the fact that Gorter³⁵⁾ was able, by straightforward reasoning, to derive from these structures the ferrimagnetic character of the W, Y and Z compounds and to calculate the values of their saturation magnetization.

It has been said already that the new magnetic compounds show a magnetic anisotropy. There are two extremes of this anisotropy. For some of the compounds the ion-spins are all directed parallel to the *c*-axis, whereas in others they are perpendicular to this axis. In the first case the components are magnetically hard which gives us a valuable new class of materials in addition to that of the magnetically soft spinels. In the second case, components with a preferred plane of magnetization, the magnetization is virtually free to assume any direction in that plane, but energy is required for rotation out of the plane. Such a preferred plane occurs with the Y group and also with those compounds of the Z and W groups in which Co-ions are completely or partly substituted for Me. Table I illustrates this phenomenon more clearly.

Table I. Compounds with a preferred direction (\uparrow) or a preferred plane (\perp).

Group	Divalent ion substituted for Me					
	Mn	Fe	Co	Ni	Zn	Mg
W	\uparrow	\uparrow	\perp	\uparrow	\uparrow	\uparrow
Z	\uparrow	\uparrow	\perp	\uparrow	\uparrow	\uparrow
Y	\perp	\perp	\perp	\perp	\perp	\perp

An important feature of the compounds showing a preferred plane is the fact that here the relation of Snoek mentioned earlier is no longer valid and that the compounds can therefore be used at much higher frequencies³⁴⁾.

It is very intriguing how complicated structures like those mentioned above can be formed out of a melt or by sintering of the oxides and that an irregular alternation of the different but very similar blocks of layers never seems to occur. No crystal was found with an averaged structure, so mistakes in growing occur much less often than one would expect from the close relationship of the layers.

Polytypes in SiC; growth phenomena and stability

This problem of how to explain the "memory" of a growing crystal for such a complicated ordered sequence of layers is encountered in a perhaps more pronounced degree in the well-known polytypes of SiC, although in these structures a slight disorder

can often be found. SiC has been a subject of study in this laboratory for a number of years. In the course of these studies Lely³⁶⁾ developed a method of growing large pure SiC crystals in the laboratory, by recrystallization of polycrystalline material under its equilibrium Si pressure. Knippenberg³⁷⁾ studied especially the growth phenomena of SiC polytype crystals, using an improved Lely technique and also various other methods of preparation. One of his results was that independent of the way of preparation a definite statistical relationship seems to exist between the polytypes formed and the temperature of preparation. A qualitative picture of this relation is sketched in fig. 6. As can be seen, cubic SiC could

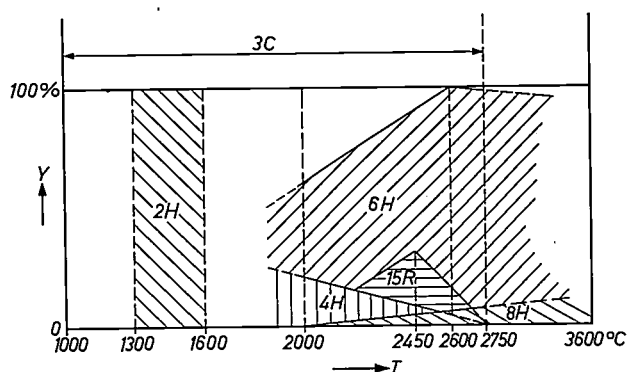


Fig. 6. Relation between structure of SiC and temperature of preparation. The relative amounts *y* of the different non-cubic structures (H and R) formed at the temperature *T* of preparation are shown qualitatively. Cubic (3C) SiC (unstable) is formed at all temperatures. The highest temperature of preparation was 2750 °C.

be formed at all temperatures but is always unstable. The cubic modification grows only under conditions of supersaturation and then in a "skeleton" formed by multiple twinning of anisotropic lamellae.

Polytypes in SiC; variation of the band gap

It is a remarkable fact that for these very closely related structures the variation of the band gap with the polytype is surprisingly large. Table II gives a comparison of the band gaps for the different polytypes. The second column in this table shows the way in which for every polytype the hexagonal (h) and cubic (c) layers follow each other along the *c*-axis.

One sees that the band gap seems to vary in a

Table II. Variation of the band gap in SiC with the polytype.

Polytype	Sequence of layers	Band gap (eV)	α
2 H	hh	> 3	0
4 H	hchc	3.1	0.50
15 R	hchcc	2.9	0.60
6 H	hchccc	2.86	0.67
24 R	hchcccc	2.72	0.75
3 C	cccc...	2.2	1

systematic way. This can be seen more clearly if one plots the band gap against the ratio α (Table II) of the number of cubic layers to the total number of layers per unit cell (fig. 7).

It is known that in mutually comparable structures there exists a relationship between band gap, lattice energy and single-bond energy, the band gap increasing with increasing lattice energy or increasing bond energy. Applying this to the series of polytypes of SiC one may expect an increase in single-bond energy and lattice energy for the respective polytypes going from the bottom to the top in Table II. The fact that cubic SiC is unstable seems to be in line with this correlation.

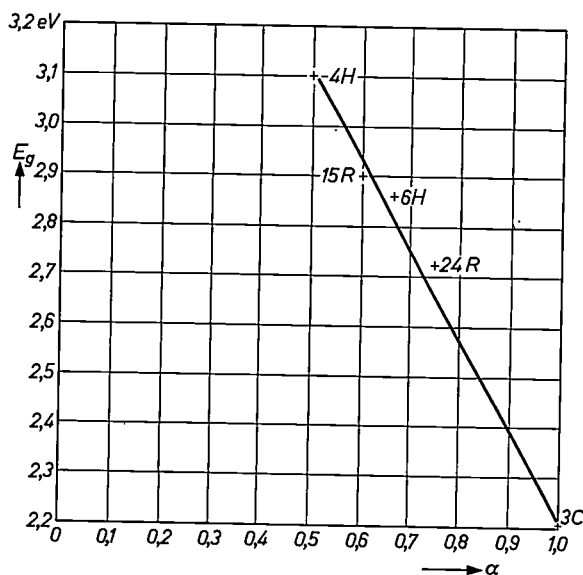
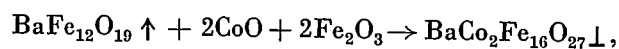


Fig. 7. Plot of the band gap E_g of various polytypes of SiC as a function of the ratio α of the number of cubic layers to the total number of layers per unit cell.

Topotactical reactions

We will now return once again to the system BaO-MeO-Fe₂O₃ (with Me = Zn, Co, etc.). Under proper conditions these compounds can be converted into each other by solid-solid reactions. For instance, when a pellet consisting of 1 equivalent BaFe₁₂O₁₉, 2 equivalents CoO and 2 equivalents Fe₂O₃ is fired for some hours at 1250 °C, the following solid-solid reaction takes place:



the signs \uparrow and \downarrow indicating a preferred direction and a preferred plane of magnetization respectively. In this way a sintered aggregate of hexagonal crystals is obtained. A remarkable discovery, made by Lotgering³⁸⁾, was the fact that when the small grains of BaFe₁₂O₁₉ (M) present in the original pellet had been previously oriented by means of a static magnetic field, the *c*-axes of the new crystals (W) show the same orientation as the

c-axes of the original M crystals, which disappear during the reaction. In the initial mixture only the M crystals are oriented. It is evident therefore that the preferential orientation of the W crystals is caused by chemical reactions that take place on the surface of the M crystals.

Such reactions are called topotactical reactions. A topotactical reaction starts as an epitaxial process; there is however, an essential difference from epitaxy. The final product of an epitaxial process is an intergrowth of the new phase and the substrate, whereas during topotactical reactions (like the one just mentioned) the substrates disappear completely during the reaction. In the example given, the amount of preferentially oriented material increases considerably during the reaction. It was found in further study that not all reactions of this kind lead to well oriented products, diffusion of reactants and intermediate products playing an important role. Of course as the reactions all take place at temperatures (1000-1300 °C) that are far above the Curie temperatures (300-500 °C), the magnetic fields of the ferrimagnetic grains originally present have disappeared and cannot play any part in the orientation.

Internal charge compensation

Controlled valency

Verwey and De Boer⁷⁾, combining Wagner's conception of the occurrence of differently charged ions in non-stoichiometric compounds (giving electronic conductivity) and their own explanation of the conductivity of Fe₃O₄, generalized these ideas by stating that electronic conductivity could be expected for any compound in which equivalent crystallographic sites were occupied by ions of the same element occurring in two states of charge differing by not more than one electron unit.

Starting from these notions Verwey, Haaijman, Romeijn and Van Oosterhout^{39) 40)} developed the conception of controlled valency. According to this principle the introduction on Ni-sites of a monovalent ion like, say, Li⁺ into the lattice of a non-conducting compound such as NiO, containing divalent Ni²⁺ ions, should bring about the formation of an equal amount of Ni³⁺ ions, in order to maintain electro-neutrality. This being the case, the lattice would then contain, on equivalent crystallographic points, Ni ions differing in charge by only one electronic unit. Therefore this Li-containing NiO should show electronic conductivity. This was in fact found. The conductivity of, for instance, pure NiO could be increased from 10⁻¹⁴ Ω⁻¹cm⁻¹ to 10² Ω⁻¹cm⁻¹ by the introduction of 10% Li.

Of course one condition has to be fulfilled, viz, the foreign ion (here Li^+) must show a smaller tendency to change its charge than the ions (Ni^{2+}) of the base crystal. Further, in order to get appreciable effects, the ionic radius of the foreign ion must not differ too much from that of the ion to be substituted, thus facilitating the incorporation.

It will be clear that in this way both *N*- and *P*-type semiconducting compounds can be obtained. The incorporation of Li^+ into NiO leading to the formation of Ni^{3+} ions among the Ni^{2+} ions gives *P*-type conduction, whereas incorporation of Ti^{4+} ions into Fe_2O_3 leading to the formation of Fe^{2+} ions among the Fe^{3+} ions will give rise to *N*-type conductivity.

Table III shows a number of mainly oxidic systems where this principle has been applied. Column 1 contains the formula of the poorly conducting mother substance, column 2 the substance containing the ion which replaces the bold-type ion of column 1. Columns 3 and 4 specify respectively the normal ions and the ions of deviating valency that are presumably present in the final system.

Table III. Systems in which ions of deviating valency have been introduced.

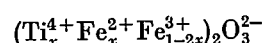
1	2	3	4	Crystal structure
NiO	Li_2O	Ni^{II}	Ni^{III}	rocksalt
CoO	Li_2O	Co^{II}	Co^{III}	
MnS	Li_2S	Mn^{II}	Mn^{III}	
CaTiO_3	La_2O_3	Ti^{IV}	Ti^{III}	perovskite
SrTiO_3	La_2O_3	Ti^{IV}	Ti^{III}	
BaTiO_3	La_2O_3	Ti^{IV}	Ti^{III}	
$\text{CaMn}^{\text{IV}}\text{O}_3$	La_2O_3	Mn^{IV}	Mn^{III}	
$\text{LaMn}^{\text{III}}\text{O}_3$	CaO	Mn^{III}	Mn^{IV}	
$\text{LaMn}^{\text{III}}\text{O}_3$	SrO	Mn^{III}	Mn^{IV}	
$\text{LaFe}^{\text{III}}\text{O}_3$	SrO	Fe^{III}	Fe^{IV}	
ZnFe_2O_4	TiO_2	Fe^{III}	Fe^{II}	spinel
MgFe_2O_4	TiO_2	Fe^{III}	Fe^{II}	
NiFe_2O_4	TiO_2	Fe^{III}	Fe^{II}	
CoFe_2O_4	TiO_2	Fe^{III}	Fe^{II}	
Fe_2O_3	TiO_2	Fe^{III}	Fe^{II}	hematite
Fe_2O_3	SnO_2	Fe^{III}	Fe^{II}	
Fe_2O_3	WO_3	Fe^{III}	Fe^{II}	
SnO_2	Sb_2O_5	Sn^{IV}	Sn^{II}	rutile
TiO_2	Ta_2O_5	Ti^{IV}	Ti^{III}	
MgWO_4	Cr_2O_3	W^{VI}	W^{V}	wolframite

However, the systems represented in Table III are by no means the only ones for which this principle of controlled valency works. Further research made it clear that this principle is quite general and is also valid for instance in compounds like CdS ⁴¹⁾, PbS ⁴²⁾, CdTe ⁴³⁾ ⁴⁴⁾, which on introduction on cation sites of ions like Ag^+ may become *P*-type and of ions like Ga^{3+} , In^{3+} , Sb^{3+} may become *N*-type. Not

only can foreign cations change the valency of the cations of the base crystal, but also foreign anions can bring this about. For instance, the introduction of Cl^- ions on S^{2-} sites in CdS may lead to the formation of an equal amount of Cd^+ ions between the normal Cd^{2+} ions and consequently to an *N*-type semiconductor ⁴¹⁾.

Experimental determination of electron-energy levels of impurities in Fe_2O_3

As has been said the foreign ion must show a smaller tendency to change its charge than the ions of the base lattice. Jonker ⁴⁵⁾ has investigated this requirement more quantitatively for Fe_2O_3 . Introduction of Ti gives rise to the formation of



and *N*-type conduction. However, at first sight a second formula $(\text{Ti}_x^{\text{III}} + \text{Fe}_{1-x}^{\text{III}})_2\text{O}_3^{2-}$ not involving any change of charge of the Fe ions, cannot be excluded. To compare the two possibilities one has to consider the energy balance of the reaction



The energies involved are the following.

- 1) The difference of the ionization potentials: $I_3(\text{Fe}) - I_4(\text{Ti})$.
- 2) The polarization energies $\sum E_p$ caused by the electric field around the ions Fe^{2+} and Ti^{4+} having a charge that differs from that of the Fe^{3+} ions occurring in Fe_2O_3 . Rough calculations show that this polarization energy for a pair of Fe^{2+} and Ti^{4+} ions is about 16 eV.
- 3) Crystal-field stabilization energies $\sum E_{cf}$ being of the order of 1 eV.

This leads to

$$E = I_3(\text{Fe}) - I_4(\text{Ti}) + \sum E_p + \sum E_{cf}.$$

Samples of Fe_2O_3 to which different foreign atoms where added in a concentration of 2 at. % were all fired in air at 1300 °C, and their resistivities measured at room temperature. The results are given in Table IV.

It can be seen from Table IV that two groups of elements can be distinguished, a first group leading to low values of ρ and a second group having no influence on the resistivity. Obviously for the second group of elements the value of E of the reaction is negative.

Neglecting the relatively small values of $\sum E_{cf}$, the value for $I_3(\text{Fe}) - I_4(\text{f.a.})$ must be compensated by $\sum E_p$, the polarization energies, in order to have controlled valency. From Table IV it follows that for Fe_2O_3 , $\sum E_p$ is of the order of 17 eV, in good agree-

Table IV. Influence of 2 at. % substitution of foreign ions on the resistivity of Fe_2O_3 . For each kind of ion the fourth ionization potential I_4 and the difference between this and the third ionization potential of iron $I_3(\text{Fe})$ is given.

Foreign atom (f.a.)	ρ $\Omega \text{ cm}$	$I_4(\text{f.a.})$ eV	$I_3(\text{Fe}) - I_4(\text{f.a.})$ eV
Ta	1.3	33.3	- 2.7
Zr	0.7	34.0	- 3.4
Nb	7.9	38.3	- 7.7
Ti	0.6	43.2	- 12.6
Sn	0.4	46.4	- 15.8
V	1.5×10^6	48.0	- 17.4
Cr	1.0×10^8	50.6	- 20.0
Mn	7.0×10^7	53.7	- 23.1
Fe	1.7×10^4	56.2	- 25.6

ment with rough calculations. Attention must be drawn to the fact that the polarization energy which is difficult to calculate, is in this way determined experimentally by measuring resistivities. The results of this can further be used to construct an electronic energy-level diagram for Fe_2O_3 with several types of impurities. As will be discussed further on, it is generally assumed that for transition-metal oxides this picture consists not of bands but of localized levels. The distance E_g between the occupied Fe^{3+} levels and the empty Fe^{2+} levels is found from

$$-E_g = I_3(\text{Fe}) - I_4(\text{Fe}) + \sum E_p + \sum E_{cf}.$$

Here $\sum E_{cf}$ cannot be neglected and is about 2 eV. Taking $\sum E_p = 17$ eV and $I_4 - I_3 = -25.6$ eV one gets $E_g \approx 6.5$ eV. Fig. 8 gives the corresponding electronic energy-level diagram of Fe_2O_3 .

It can be remarked that the results obtained in

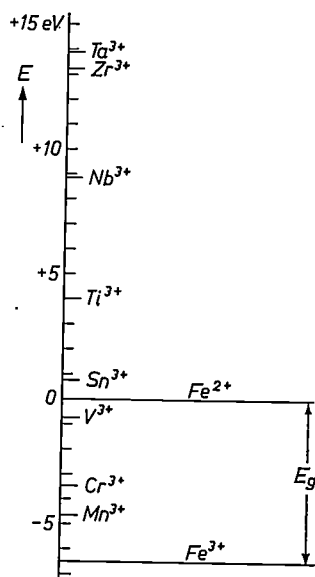


Fig. 8. Electronic energy-level diagram of Fe_2O_3 with the levels of substituted trivalent foreign ions.

this experiment and also the results of quite a number of other experiments in solid-state chemistry can often be qualitatively predicted by the chemist's "chemical feeling". Now this chemical feeling amounts to an almost unconscious knowledge about the relations between ions dissolved in water. That this feeling nevertheless can account for the results obtained can be understood quite easily when one realizes that the hydration energy of the ions (6-11 eV) is about equal to the polarization energy (6-9 eV) in these oxidic and other ionic compounds.

Calculation of electronic energy levels in NiO

Van Houten ⁴⁶⁾ was able to calculate the electron-energy scheme of NiO. He assumed 1) the crystal to be build up from Ni^{2+} and O^{2-} ions and 2) localization of the electrons. By taking into account a) the ionization energies or electron affinities of the ions involved, b) the Madelung energies, c) the polarization energies and d) the crystal field stabilization of the transition metal ions, a distance of 5.4 eV was found between the highest filled level (that of the Ni^{2+} ions) and the lowest empty level (that of the Ni^+ ions). This large distance gives, just as did that of 6.5 eV for Fe_2O_3 , a quantitative explanation of the fact that pure oxides of this type exhibit a very high resistivity.

Mobility of electronic charge carriers in transition-metal oxides

De Boer and Verwey ⁸⁾ drew attention to the fact that pure stoichiometric oxides of the transition metals are nearly all insulators. These oxides, however, contain metal ions with incompletely filled 3d or 4f shells. Especially for those ions having incomplete 3d shells the collective electron treatment of wave functions of Bloch-Wilson would lead one to expect metallic conduction for these oxides. It became clear that instead of the usual band picture of electron energies, these oxides must show sharp localized levels.

As has been said already the conductivity of these oxides could be enhanced substantially by the incorporation of an excess of one of their components or of ions with a deviating valency. Now the semiconductors made in this way mostly show a high positive temperature coefficient of conductivity:

$$\sigma = \sigma_0 e^{-Q/kT}.$$

The activation energy Q can have values lying between 0.1 and 0.5 eV and decreasing with increasing deviation from stoichiometry or concentration of foreign ions ⁴⁶⁾. It has long remained obscure whether this activation energy was caused by the energy

needed to liberate the charge carriers from their centres or by the fact that in contradistinction to semiconductors like Ge, the mobility involves an activation energy.

To investigate this question Jonker⁴⁷⁾ measured the conductivity of mixed crystals of CoFe_2O_4 with Fe_3O_4 or Co_3O_4 , as a function of temperature and concentration. In the mixed crystals $\text{Co}_{3-x}\text{Fe}_x\text{O}_4$ in the vicinity of $x = 2$, one has for $x > 2$ an Fe^{2+} content of $x - 2$ and for $x < 2$ a Co^{3+} content of $2 - x$. Now applying the rules of occupation of tetrahedral and octahedral sites mentioned in the second section it can be concluded that for $x > 2$ ($\text{Fe}^{3+}[\text{Co}_{3-x}^{2+}\text{Fe}_{x-2}^{3+}\text{Fe}_{x-2}^{2+}]\text{O}_4^{2-}$) one has electron conduction among the Fe ions and that for $x < 2$ ($\text{Fe}^{3+}[\text{Co}^{2+}\text{Co}_{2-x}^{3+}\text{Fe}_{x-1}^{3+}]\text{O}_4^{2-}$) one has hole conduction among the Co ions.

With this series of mixed crystals, therefore, it is possible to vary the conductivity and even its nature without introducing centres, the concentration and kind of charge carriers being given directly by the ratio of the Co and Fe contents (fig. 9). By also

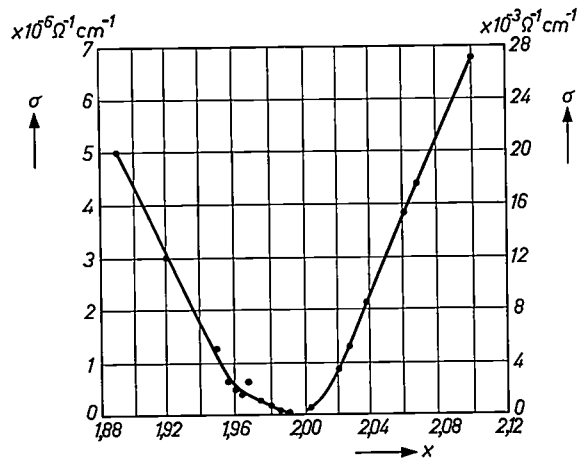


Fig. 9. Conductivity σ at room temperature of a series of mixed crystals of Co_3O_4 and Fe_3O_4 ($\text{Co}_{3-x}\text{Fe}_x\text{O}_4$) in the neighbourhood of the composition CoFe_2O_4 ($x=2$). For hole conductivity (excess cobalt: $x < 2$) the left ordinate has to be used; for electron conductivity (excess iron: $x > 2$) the right ordinate.

measuring the conductivity as a function of temperature it was found that for this compound the mobility showed an exponential dependence on temperature. At room temperature the mobilities turned out to be very small, e.g. $\approx 10^{-4}$ cm^2/Vs for electrons and $\approx 10^{-8}$ cm^2/Vs for holes, their activation energies being about 0.2 and 0.5 eV respectively. By measuring the Seebeck coefficient as a function of charge-carrier concentration (fig. 10) the number of available sites for the charge carriers proved to be about 10^{22} per cm^3 , this being roughly equal to the number of metal ions per cm^3 . This number is quite different from that to be expected from the

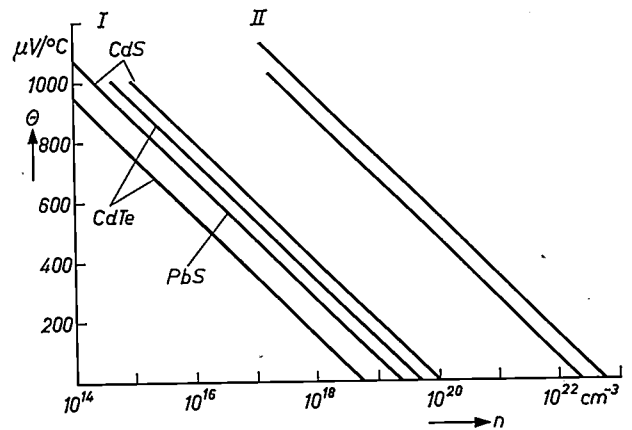


Fig. 10. Schematic survey of the Seebeck effect Θ of a number of compounds as a function of electron or hole concentration n , extrapolated to the intersection points with the abscissa. Group I consists of the Ge-type semiconductors CdS⁴¹⁾, PbS⁴²⁾ and CdTe⁴³⁾. The two lines of group II enclose the measurements on NiO and CoFe_2O_4 ⁴⁶⁾ ⁴⁷⁾.

usual band theory of semiconductors ($\approx 10^{19}$) and in good agreement with the assumption that the charge carriers jump from one metal ion to the next, as should be the case for the localized electron-energy levels. The activation energy Q therefore refers in this case to the mobility and must be considered as being the energy q required to overcome the Landau trapping of the electron, i.e. the polarization of the lattice around the slowly moving charge carrier (see fig. 11).

In oxides that have been made conductive by means of a slight deviation from stoichiometry or by introduction of foreign atoms (centres), at low temperature the charge carriers will be trapped near

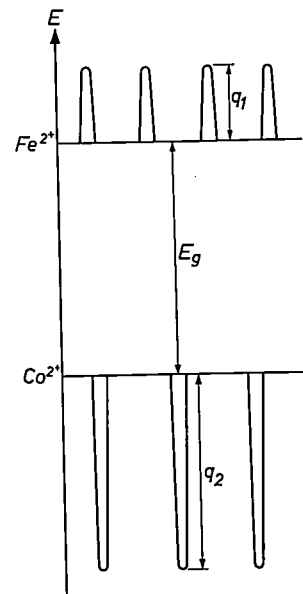


Fig. 11. Energy-level scheme for $\text{Co}_{3-x}\text{Fe}_2\text{O}_4$ ($x \approx 2$). For pure CoFe_2O_4 at low temperature the Fe^{2+} level is empty and only Fe^{3+} is present, whereas the Co^{2+} level is full and only Co^{2+} ions are present. q_1 and q_2 are the activation energies needed for jumps of electrons or holes.

these centres and there will be no conduction. However, in contradistinction to the usual band type semiconductor, here the trapped electron because of its localization, together with the deviating charge of the centre may form a permanent dipole, which can reorient itself by local jumps of the trapped carrier between equivalent positions around the centre. Dielectric loss measurements (Volger) yield the relaxation rate $1/\tau$, which is given by:

$$\frac{1}{\tau} = \frac{1}{\tau_0} e^{-\mu/kT}$$

The value of the activation energy, μ , can be found by carrying out loss measurements as a function of temperature. It need not be equal to the activation energy q governing the motion of the dissociated charge carrier. From low temperature measurements on $\alpha\text{-Fe}_2\text{O}_3$ with a slight excess of Fe, a value for μ of 0.005 eV was found⁴⁸⁾.

Van Houten⁴⁶⁾, measuring Seebeck coefficients in NiO heavily doped with Li as a function of temperature, was able to calculate from these measurements the energy needed to liberate the holes from their $(\text{Li}+\text{Ni}^{2+})$ centres. It was found to be low: $E_a \approx 0.035$ eV. This seems to be the case for high concentrations in a large number of compounds of this kind.

Activators and co-activators in phosphors

We have seen that charge compensation plays an important role in the concept of controlled valency. This same idea of charge compensation was used by Kröger⁴⁹⁾ 50) 51) 52) in elucidating the complicated state of affairs in the ZnS-type phosphors. He was able to demonstrate that the so-called activator atoms Ag, Au and Cu actually were effectively negatively-charged monovalent cations (Ag^+) replacing divalent Zn^{2+} ions (Ag'_{Zn}), their effective negative charge being compensated by Cl^- ions (coactivators) on S^{2-} sites (Cl'_{S}) brought into the crystal from the NaCl flux used as mineralizer. In the symbols dots denote effectively positive centres and the dashes effectively negative centres. If the activators Ag^+ and Cu^+ really occupy Zn sites and not for instance interstitial sites, then the trivalent ions Al^{3+} , Ga^{3+} , In^{3+} , Sc^{3+} substituting Zn^{2+} ions should also be able to act as coactivators, which indeed proved to be true. Kröger, Hoogenstraaten⁵³⁾ and Klasens⁵⁴⁾ and later Van Gool⁵⁵⁾ using these notions were able to explain for this kind of phosphor a surprising number of physical effects related to luminescence, including those of decay, trapping, killing, quenching and thermal glow.

General defect chemistry

So far we have encountered two ways of compensating the deviating charge of a foreign ion incorporated in a crystal structure. The first (controlled valency) takes place by a compensating change of the charge of the ions of the host crystal, the second by the simultaneous incorporation of other foreign ions having a valency such that the deviating charges of the foreign ions compensate each other. There is still another possibility. A well-known example is that of the incorporation of divalent Ca^{2+} ions into NaCl crystals (see for instance Haven⁵⁶⁾), where the effective positive charge of the Ca^{2+} ion replacing an Na^+ ion (Ca'_{Na}) is compensated by the effective negative charge of a missing Na^+ ion, that is by an effectively negatively charged Na^+ -ion vacancy (V'_{Na}). Now Kröger and Vink⁴¹⁾ 57) 58) 59) 60) asked themselves how it comes about that sometimes the deviating charge of the foreign ion is compensated by a corresponding change of valency of the ions of the base lattice, whereas in other cases the host lattice creates effectively charged atomic defects (like vacancies or interstitials) to arrive at the same end.

As a consequence of thermodynamical requirements the pure compound contains all kinds of imperfections: free electrons, holes and effectively charged and uncharged atomic imperfections. Starting from ideas developed by Schottky and Wagner, Kröger and Vink were able to work out a general theory showing how the concentrations of the different electronic and atomic imperfections depend on the concentrations of foreign ions in the lattice and on the composition of the gas phase that is in equilibrium with the crystals at the temperature of preparation. In this theory the effectively charged imperfections are of prime importance. In a simple case, for instance pure PbS, one has four of those effectively charged imperfections: free electrons, holes, effectively negatively charged lead vacancies V'_{Pb} and effectively positively charged sulphur vacancies V_S . The concentrations depend on each other via relations similar to those of the usual law of mass action, e.g.

$$n p = K_i$$

for the electronic imperfections, with n representing the concentration of the electrons and p that of the holes, and

$$[\text{V}'_{\text{Pb}}][\text{V}_\text{S}] = K'_s$$

for the atomic imperfections. The brackets denote concentrations. The influence of the composition of the gas phase can be expressed by

$$p_{Pb} = K_r [V'_S] n \quad (\text{or } p_{S_2}^\dagger = K_{ox} [V'_{Pb}] p),$$

with p_{Pb} the partial pressure of the lead atoms in the gas phase that is in equilibrium with the pure PbS and p_{S_2} that of the S_2 molecules.

Electro-neutrality must always be maintained:

$$n + [V'_{Pb}] = p + [V'_S].$$

This nomenclature and also a more thorough investigation of the deeper-lying aspects of this theory have been discussed in a paper by Kröger, Stieltjes and Vink ⁶¹).

The actual calculation of equilibria involves simple but somewhat tedious and complicated mathematics. Brouwer ⁶²) published a method by means of which the calculations can be substantially simplified. It appears that for large ranges of the partial pressure of the components only two imperfections in the above equation for electro-neutrality have concentrations that need to be taken into account. So for relatively high lead pressures one has

$$n = [V'_S],$$

i.e. Pb excess in the solid and *N*-type conductivity. For relatively high sulphur pressures (low lead pressures) one has however

$$p = [V'_M],$$

which means S excess and *P*-type conductivity. In between there is a region where $n = p$. Now for PbS containing some Bi on Pb sites one has for the condition of electro-neutrality

$$n + [V'_{Pb}] = p + [V'_S] + [Bi'_{Pb}]$$

(trivalent Bi^{3+} ions on Pb^{2+} sites giving rise to effectively positive Bi'_{Pb} centres).

It can now easily be shown that for sufficiently high concentrations of Bi two new simplified forms arise for the condition of electro-neutrality together with their corresponding ranges of lead pressures. For high lead pressures one has

$$n = [Bi'_{Pb}],$$

i.e. the effective charge of the foreign ion is compensated by an electronic defect (change in valency). In other words, for high lead pressure Bi is incorporated according to the principle of controlled valency, with avoidance of lead vacancies. On the other hand, for low lead pressures one has

$$[V'_{Pb}] = [Bi'_{Pb}],$$

i.e. the effective charge of the Bi ion is compensated by an effectively charged atomic imperfection.

Application of defect chemistry in CdS, PbS and SnS

By means of this theory it was possible to explain in a qualitative and quantitative way how in the systems CdS, ZnS, PbS, CdTe and SnS the properties of the compounds depend on the kind and concen-

tration of the impurities introduced into the crystal and on the partial pressures of the components during the preparation. For instance, it was possible to explain why CdS containing Ga on Cd sites and prepared under sulphurizing conditions (low Cd pressures) has a bright red color and is an insulator, whereas with the same amount of Ga but prepared under reducing conditions (high Cd pressures) the yellow colour of pure CdS does not change with the incorporation of Ga, and the substance shows *N*-type conductivity, the charge carriers at room temperature having the same concentration as the Ga atoms ⁴¹).

The explanation amounts to the fact that under "sulphurizing" conditions the Ga^{3+} ions are incorporated into CdS in the same way as the Cd^{2+} ions are incorporated into NaCl, i.e. by the formation of effectively charged atomic imperfections, viz, Cd vacancies. It is those vacancies that are responsible for the red colour. Under reducing conditions, however, the Ga^{3+} ions are incorporated according to the principle of controlled valency, viz, by the formation of an equal amount of Cd^+ ions, which amounts to free electrons in this energy-band type of semiconductor.

Fig. 12 gives the concentration of electrons (n) and

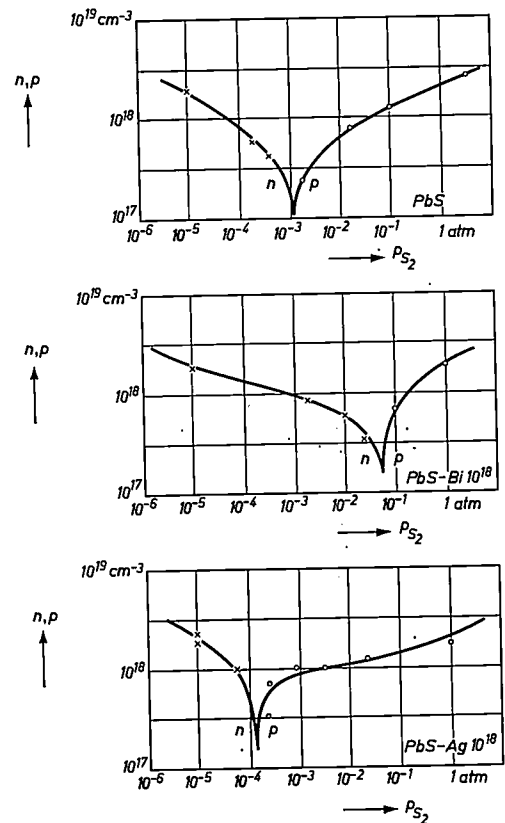


Fig. 12. Concentrations of electrons (n) and holes (p) in undoped PbS crystals (upper graph), in PbS crystals doped with 10^{18} Bi atoms per cm^3 (middle graph) and with 10^{18} Ag atoms per cm^3 (lower graph) as a function of the sulphur pressure p_{S_2} during heating at $1000^\circ K$, as measured after subsequent quenching. Logarithmic scales along both axes.

holes (p) in doped and undoped crystals of PbS, equilibrated under various sulphur pressures at 1000 °K, as measured after quenching ⁴²). From the figure for pure PbS it can be seen that the plot is nearly symmetrical, showing a sharp minimum in the carrier concentration at a certain sulphur pressure; at higher and at lower sulphur pressures the concentrations increase rapidly, whereas still farther from the minimum the curves flatten and reach a constant slope, corresponding to an electron or hole concentration proportional to $(P_{S_2})^{\mp 1/4}$. The type of carriers is different on both sides of the minimum, the material being N -type in the region of low sulphur pressures and P -type in the region of high pressures. When bismuth ions are present (Bi_{Pb}^+) it is seen that there is a certain range of pressures where the crystals contain an electron concentration which is about equal to that of the added foreign ions. Similarly, there is a range in which crystals doped with silver (Ag_{Pb}') contain a concentration of holes about equal to that of the added silver ions. In these ranges the foreign ions are incorporated according to the principle of controlled valency. Bloem ⁴²) was able to explain all these features and others not mentioned here by the application of the theories mentioned above.

For the pure substance these ideas permit the determination of the limits of the usually very narrow region of stability of solid semiconductors with respect to other phases as a function of the composition of the vapour phase and of the temperature along the three-phase equilibrium: solid semiconductor — vapour phase — other liquid or solid phase. Fig. 13 gives as an example the limits of stability of solid SnS as a function of temperature, the pressure of the coexisting vapour phase also being given ⁶³).

Refined molecular engineering in phosphors and related compounds

Van Gool ⁵⁵) pointed out that one of the advantages of the considerations of this theory is the fact that the experimenter becomes aware of the large number of possibilities of explaining an often *limited amount* of experimental data, and is thus prevented from jumping to conclusions. For instance, in CdS containing chlorine on sulphur sites the requirement of electro-neutrality may read

$$n + [V_{Cd}'] + 2[V_{Cd}''] + [(Cl_S V_{Cd})'] = p + [V_S] + 2[V_S''] + [Cl_S],$$

leading to quite a number of possible simplified electro-neutrality conditions with their corresponding ranges of pressures. The knowledge of the simplified

conditions that occur in reality as a function of the partial pressures of Cl_2 and S_2 gives a great deal of information about the system under consideration.

In this connection Van Gool ⁶⁴) gave examples of how to calculate the partial pressures of the relevant components from thermodynamical data; in particular he calculated how, in the system $ZnS-H_2S/H_2/Cl_2$, as a function of temperature, the partial pressures p_{Cl_2} , p_{S_2} and p_{Zn} depend on the initial pressures p_{H_2S} , p_{H_2} and p_{HCl} or p_{Cl_2} of the gas stream to be equilibrated with ZnS .

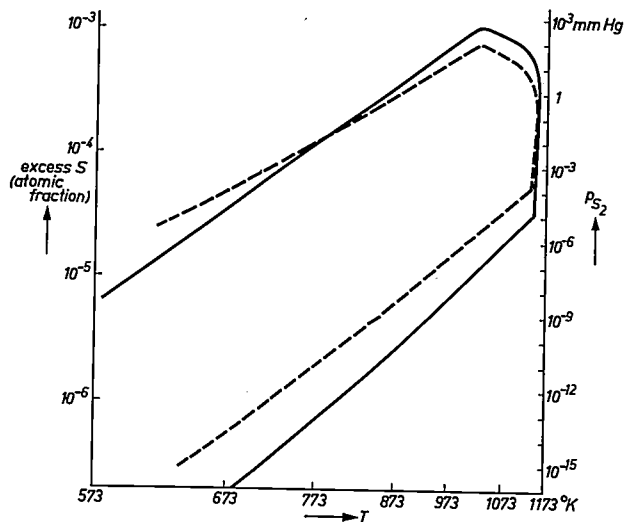


Fig. 13. Limits of stability of solid SnS. The solid curve (left vertical scale) gives the composition of solid SnS at the three-phase equilibrium: solid SnS — vapour phase — other liquid or solid phase, as a function of the temperature T . The composition is given in atom fractions of excess S atoms. The broken curve (right vertical scale) refers to the partial pressure p_{S_2} of S_2 molecules in the coexisting vapour phase. T is plotted reciprocally ($1/T$ scale), the quantities on the vertical scales are plotted logarithmically.

It will be clear from the foregoing that by controlling in this way the composition of the gas phase in equilibrium with the compound at its temperature of preparation, together with a thorough purification and doping of the compound, it is possible, in principle, to perform a very much refined sort of molecular engineering. However, as has been said, the number of possible defects is considerable. It therefore appears likely that although the chemist is nowadays able to prepare all kind of phosphors, photoconductors and other devices of almost any desired property, the explanation of what he is really doing is still lacking in many cases. Only by a much more careful study of the correlation between a *large variety* of physical properties (such as absorption, luminescence and its dependence on temperature and intensity of stimulation, glow curves, photoconductivity, quenching, electron-spin resonance) and the exact chemical composition obtained by this refined molecular engineering, a real insight

into the nature of centres responsible for the different effects might eventually be reached.

Gases and metals

Internal friction; Snoek effect

In 1939 Snoek⁶⁵) discovered one of the causes of internal friction in solids. The solid solution of nitrogen in iron will be taken as an example. At zero stress the nitrogen atoms will be distributed in a completely disordered way over the interstitial sites. On the application of a unidirectional stress some of them will jump into sites in the immediate neighbourhood, if their energy will be lower in these sites. The relaxation time τ for this redistribution is determined by the jump frequency of the nitrogen atoms and, as a consequence, τ obeys the equation

$$\tau = \tau_0 e^{U/RT},$$

where U is the activation energy for the diffusion of nitrogen atoms. Because of this relaxation time, the strain lags behind the stress when a periodically varying stress is applied to the metal. The logarithmic decrement of the mechanical oscillations divided by π is called the loss angle ($1/Q$). The maximum loss angle will be found for $\omega\tau = 1$, ω being the angular frequency of the imposed mechanical oscillation.

The phenomenon of internal friction has been applied in a wide range of investigations⁶⁶), mainly because of two facts. Firstly, as only jumps over atomic distances are needed to give the inelastic effects, information related with movements over small distances can easily be obtained. For instance, Fast and Verrijp⁶⁷) were able to measure the diffusion constant of nitrogen in iron at room temperature down to -35°C , the values lying between 10^{-16} and 10^{-20} cm^2/s . Secondly, as the effect is very sensitive and proportional to the concentration of the jumping atoms, very small concentrations can be detected. For instance, concentrations of 0.001% by weight of N in Fe can be detected with an accuracy of 10%⁶⁸).

Application of internal friction in the study of the system Fe-Mn-N

In this laboratory Fast and Verrijp have been studying the systems Fe-C and Fe-N extensively, mainly by means of this effect. Only one example will be given, viz, the interaction of interstitial N with substitutional Mn⁶⁹). Fig. 14 gives the loss angle as a function of temperature for Fe-N and for Fe-Mn-N with 0.5 and 2% Mn respectively, the N-concentration for the three systems amounting to about 0.03% by weight. One sees that on introduction of Mn new maxima come into existence lying to the left and to the right of the maximum of the pure

Fe-N system. The new maximum at the higher temperature is interpreted as being caused by a nitrogen atom jumping around a Mn atom. The free enthalpy of formation of $\text{Mn}_{2.5}\text{N}$ being larger than that of Fe_4N , the N atom is more tightly bound around a Mn atom, thus causing the corresponding maximum

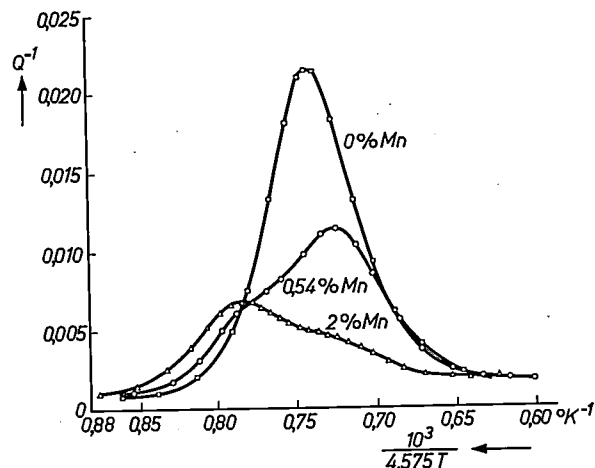


Fig. 14. Internal friction of iron containing 0.03% N with 0%, 0.54% and 2% Mn. The loss angle Q^{-1} is given as a function of the temperature T . The maximum loss angles lie at 22° , 29° and 6°C . Oscillation period 1.2 s.

loss angle to lie at a higher temperature. The new maximum to the left of the original one, occurring at higher Mn contents, is thought to be caused by an N atom jumping around a pair of Mn atoms lying next to each other. It can be made plausible that jumps around a pair of Mn atoms require less activation energy than those around Fe atoms⁶⁶). The bonding of an N atom to such a pair of Mn atoms however is still stronger than that to a single Mn atom. This can be seen in fig. 15, where the loss angle is given as a function of temperature after different annealing times at 100°C . During this annealing the N precipitates out in the form of Fe_4N . One sees that the free N and the N jumping around a single Mn atom disappear more quickly than the N attached to a Mn pair.

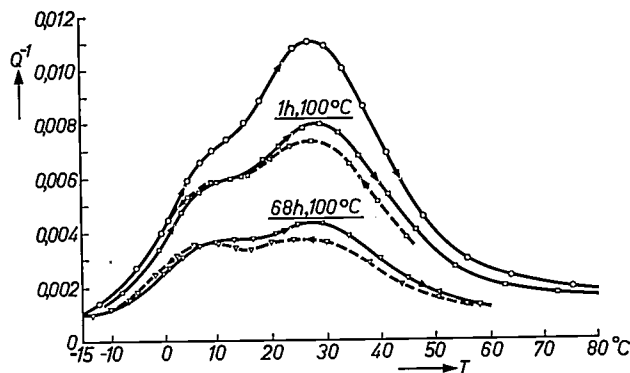


Fig. 15. Internal friction of iron containing 0.03% N and 0.54% Mn. Upper curve: before tempering; other curves: after 1 hour and 68 hours of tempering at 100°C in order to partially precipitate the nitrogen in the form of Fe_4N .

Hardening by precipitation and internal oxidation

The precipitation of N into very small particles of Fe_4N causes a considerable increase of the hardness of iron. This precipitation hardening, as it is called, not only occurs in the system Fe-N (or Fe-C), but also in a number of other cases. For instance, Ag containing some Mg (1.5%) can be hardened considerably (the Vickers hardness increasing from 40 to 170) by internal oxidation of the Mg. The experimental and theoretical aspect of this internal oxidation has been given a great deal of consideration by Meijering and Druyvesteyn^{70) 71) 72)}. The conditions for internal oxidation are evident:

- 1) The oxygen must diffuse into the metal quicker than the foreign metal diffuses out. This means that one must have $C_0D_0 > C_{\text{Me}}D_{\text{Me}}$, C_0 being the surface concentration of oxygen atoms in equilibrium with the O_2 pressure or with an external oxide layer (e.g. MeO), C_{Me} the bulk concentration of the foreign metal and D_0 and D_{Me} the respective diffusion constants.
- 2) The affinity of the dissolved metal for oxygen must be greater than that of the base metal.

With regard to the first point it can be said that in most metals D_0 of the interstitial oxygen atoms will be larger than D_{Me} of the substitutional metal atoms but that in a limited number of metals only, e.g. Ag, Cu, Ni, Fe, Ti and Zr, the solubility of oxygen is large enough to obtain $C_0D_0 > C_{\text{Me}}D_{\text{Me}}$. The second point, however, cannot be fulfilled by Ti and Zr. Therefore only Ag, Cu, Ni and Fe can be subjected to internal oxidation. This process was subjected to a rigorous mathematical treatment by Meijering and Druyvesteyn, one of their results being that to a good approximation the concentration of the oxide molecules for a dissolved divalent metal Me is given by

$$C_{\text{ox}} = C_{\text{Me}}(1 + C_{\text{Me}}D_{\text{Me}}/C_0D_0).$$

The concentration of oxide molecules is therefore somewhat larger than that of the original foreign metal. In a metal that has been subjected to a complete internal oxidation this cannot be the case throughout the volume. It is to be expected therefore that in a strip of metal being oxidized ("inter-

nally") from both sides, somewhere in the middle there must exist a region showing almost no oxidation. This region can be made visible by etching and can be seen in *fig. 16*. This "internal reaction" can be used as an extremely sensitive tool to see whether atoms like O and F dissolve and diffuse at all in certain metals. Use is thereby made of the fact that by etching, the transition between the "oxidized" part of the volume and that part of the metal that has not yet been reached by the diffusant can be made visible, as we have seen already.

This was used by Meijering⁷³⁾ to show that fluorine can be dissolved in and diffuses through silver. By this method it could also be demonstrated that oxygen dissolves in iron⁷⁴⁾, although its concentration is so small that even by means of the method of internal friction its presence could not be detected.

Controlled precipitation in ferromagnetic alloys

Precipitation also plays a predominant role in making ferromagnetic alloys with high $(BH)_{\text{max}}$ values, because by this method, in principle, particles can be prepared small enough to contain only one single domain separated by non-magnetic material. In ferromagnetic alloys of a special composition the particles may be precipitated in the form of needles, thus minimizing the demagnetization. By carrying out this precipitation in a single crystal under the application of a static magnetic field, all needles and their magnetic moments can be oriented. Thus a material with a very high value of $(BH)_{\text{max}}$ can be obtained (Jonas, Meerkamp van Embden, Luteijn, De Vos, Fast and De Jong⁷⁵⁾⁷⁶⁾⁷⁷⁾). With alloys consisting of 35% Fe, 34% Co, 15% Ni, 7% Al, 4% Cu and 5% Ti, from which needles rich in Fe and Co are precipitated in a matrix rich in Ni and Al, values of $(BH)_{\text{max}} = 12 \times 10^6$ gauss oersted have been obtained. The orientation of these precipitates can be clearly seen by using the electron microscope (De Jong, Smeets and Haanstra⁷⁸⁾). *Fig. 17* shows two such electron-microscope photographs.

Preparation of pure metals by means of decomposition of gaseous compounds

We will now enter the field of interaction between gases and metals. This is another area of research that has been investigated extensively by the Philips Laboratories. The research started in 1923 when Van Arkel^{79) 80)} published a method to prepare pure and single-crystal tungsten. In this method the tungsten is deposited on a hot wire at 1600 °C by dissociation of WCl_6 . The Cl_2



Fig. 16. Ag with 1.5 weight % Mn, 48 h on 800 °C in air; etched $\text{NH}_4\text{OH} + \text{H}_2\text{O}$. Magnification 10 \times . Although oxidation is completed, a "middle-line", poor in oxide, persists.

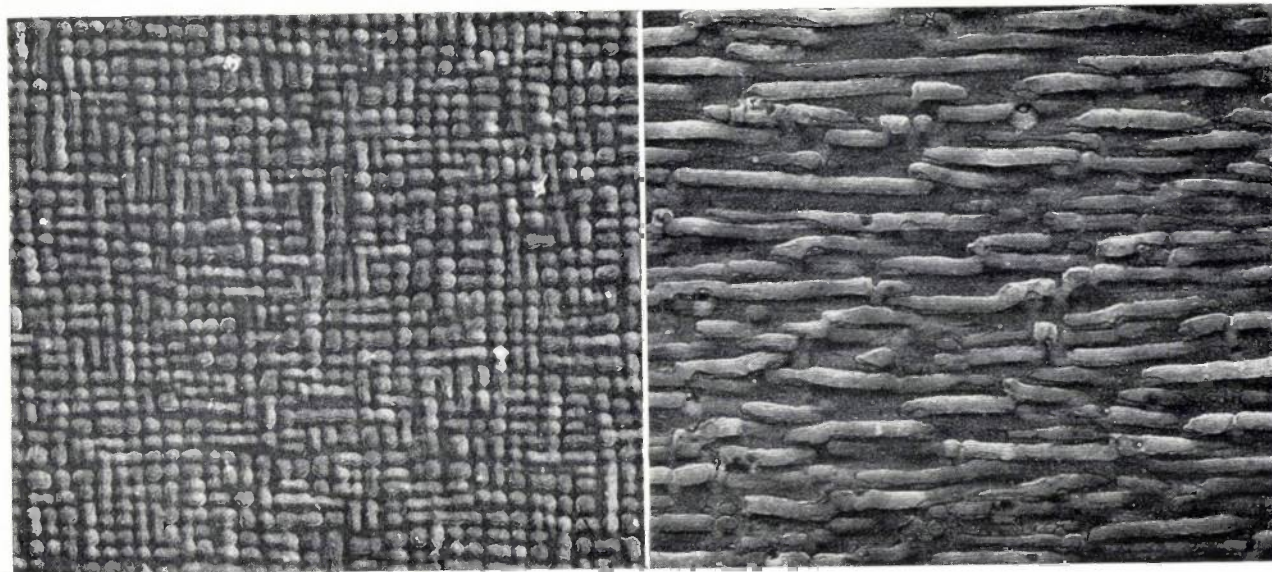


Fig. 17. a) Electron-microscope photograph of "Ticonal" with $(BH)_{\max} = 11 \times 10^6$ gauss. oversted. The regular arrangement of the needles of the precipitates is shown in a plane perpendicular to their length. Magnification $70\,000\times$.
b) The same as (a), but now the arrangement of the needles is shown in a plane parallel to their length.

liberated by this dissociation then reacts with W powder in another part of the reaction tube kept at 400°C . The reaction product WCl_6 again dissociates at the hot wire and so it is understandable that by means of a relatively small amount of chlorine a large quantity of W powder can be converted into a thick wire of single crystals of W of extremely high purity.

In the last few years this method under the name of vapour-phase transport reactions has drawn new attention, as it often proves to be an excellent way of preparing pure single crystals or epitaxial layers of compounds.

Shortly after Van Arkel's experiments De Boer and Fast showed that this method could also be applied to the metals Ti, Zr and Hf^{81) 82) 83) 84) 85)} by using the tetra-iodides instead of the chlorides. Up to then these metals had been considered to be very brittle by nature. It appeared, however, that when prepared by the method just mentioned these metals were all very ductile indeed, so much so that it was first thought that perhaps they were other modifications. Very soon however it appeared that the brittleness of these metals is caused by the presence of small concentrations of oxygen, nitrogen or hydrogen.

In Zr powder for instance every grain of metal is covered by a thin layer of oxide. On pressing and sintering the powder this oxide dissolves in the metal causing it to become brittle. Ti, Zr and Hf can dissolve very high concentrations of these gases. Zr, for instance, may contain up to 30 at. % of oxygen⁸⁶⁾ and 20 at. % of nitrogen⁸⁷⁾, while Ti and

Hf also can contain high concentrations of oxygen and nitrogen.

Permeation of gases through metals

It was further found that the dissolved atoms of these gases are present at the interstitial sites of the metal crystal and have a relatively high mobility. It is, however, impossible, once they have been dissolved in the metals (Ti, Zr or Hf), to drive these gases out by mere heating. Therefore although their solubility is high and their absorption is rapid, the permeation of oxygen and nitrogen through Ti, Zr or Hf is impossible because for that process also a third step, that of extraction, must be easy. This permeation of gases through metals is another aspect of the interaction between gases and metals that was carefully studied in these laboratories⁸⁸⁾.

For instance, the case of hydrogen and iron⁸⁹⁾ is in some way opposite to that of oxygen and zirconium. As De Boer and Fast showed, hydrogen atoms move very easily interstitially through iron and it is also easy for them to leave the metal and recombine to H_2 molecules. On the other hand it is only possible to introduce H atoms into the iron at room temperature when by some other means the H_2 molecules are first dissociated into atoms. This can be done either by dissociation at a hot W wire or in an aqueous solution with a pH value low enough for H atoms to be formed at the Fe interface.

Owing to these and other effects the permeation of gases through metals proved to be a complicated

process consisting of at least five successive stages, viz, absorption, adsorption, diffusion, de-absorption and de-adsorption, the theoretical aspects of which were carefully analysed.

Investigation of getters; the system Th-Al-Ce

As has already been indicated, Zr, that at high temperatures can absorb huge quantities of gases, is protected at low temperatures by a very thin layer of oxide. This property makes it useful as a getter. At low temperatures it can be mounted in the tube to be gettered and later on at any moment it can be heated to absorb gases. It goes without saying that "gettering" in general has been a subject of investigation. For instance, a great deal of research work has been carried out in order to understand the getter action of the Ceto getter. This getter contains mainly Th, Al and Ce. To understand its action Van Vucht made a study of the binary systems Th-Al, Al-Ce and Ce-Th and the ternary system Th-Al-Ce, using X-ray analysis and thermo-analytical and metallographical methods. It was found that pure Ceto has the structure of the intermetallic compound Th_2Al in which about one quarter of the Th atoms is replaced by Ce⁹⁰.

Absorption of hydrogen in Th_2Al

To examine the action of this type of getter in more detail, the absorption of H_2 by the well-defined compound Th_2Al was carefully analysed by means of the determination of hydrogen-absorption isotherms at different temperatures, X-ray diffraction of Th_2Al containing hydrogen, neutron diffraction of deuterium solutions in Th_2Al and proton magnetic resonance⁹¹). It was found that the unit cell of Th_8Al_4 apart from other holes, contains 16 equivalent holes tetrahedrally surrounded by Th atoms arranged in 8 pairs of twin holes, having one tetrahedral plane in common, namely the plane parallel to the base of the cell (fig. 18). This lattice can be filled with H atoms up to the composition of $\text{Th}_8\text{Al}_4\text{H}_{15.4}$. It could be shown that the hydrogen atoms are all situated in the holes mentioned above.

Measurements of the hydrogen equilibrium pressures at various temperatures yielded isotherms which (around 225 °C) showed three points of inflexion (fig. 19). Two of these near the compositions $\text{Th}_8\text{Al}_4\text{H}_4$ and $\text{Th}_8\text{Al}_4\text{H}_{12}$ were considered to be remnants of plateaux, indicating two-phase regions at lower temperatures. This would imply that $\text{Th}_8\text{Al}_4\text{H}_8$ had to be considered as an intermediate hydride. X-ray diffraction showed that at room temperature a phase separation did in fact occur, the compositions of the two phases being about Th_8Al_4

and $\text{Th}_8\text{Al}_4\text{H}_8$. However, even at 83 °K a similar region was not found between $\text{Th}_8\text{Al}_4\text{H}_8$ and $\text{Th}_8\text{Al}_4\text{H}_{15.4}$. The changes in the cell dimensions with increasing hydrogen content seemed to indicate that at $\text{Th}_8\text{Al}_4\text{H}_8$ the hydrogen atoms had a strong tendency to order.

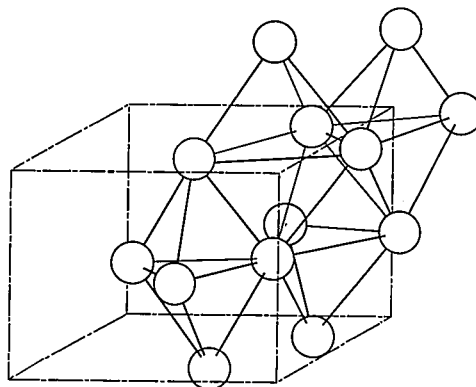


Fig. 18. Structure of Th_2Al (only Th atoms are shown). Arrangement of tetrahedral holes as "double holes".

Thermodynamics

The results obtained so far by Van Vucht could be further corroborated by a more theoretical analysis of the hydrogen-absorption isotherms⁹²). A general expression for the Gibbs free energy could be set up in which the enthalpy contained a term corresponding to the energy of a single hydrogen atom in a "double hole" and a term corresponding to the energy of a pair of hydrogen atoms in a double hole. If ϑ is the fraction of holes occupied by a proton and q the fraction of holes that are occupied and

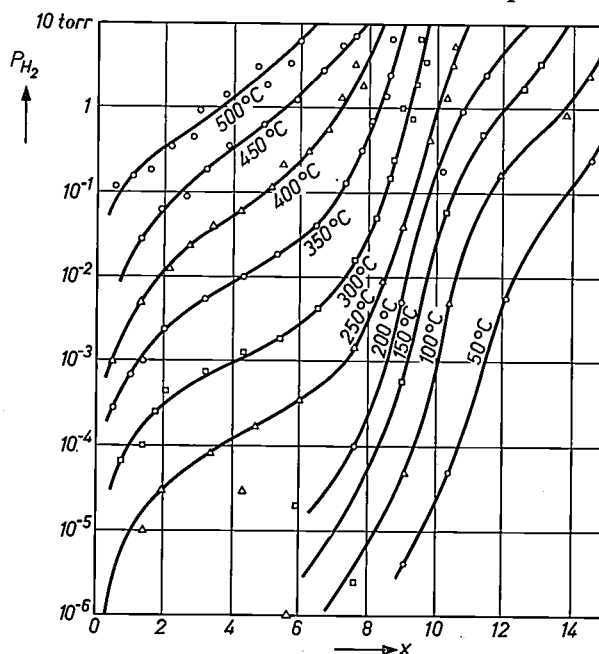


Fig. 19. Isotherms of the hydrogen equilibrium pressure P_{H_2} above $\text{Th}_8\text{Al}_4\text{H}_x$. This pressure is plotted logarithmically against x .

associated with a vacant twin, then at equilibrium

$$\frac{\partial G}{\partial q} = 0 \quad \text{and} \quad \frac{dG}{d\phi} = \frac{1}{2}\mu_{\text{H}_2} = \frac{1}{2}\mu_{\text{H}_2}^0 + \frac{1}{2}RT \ln p_{\text{H}_2}.$$

From these conditions a general expression for the hydrogen-absorption isotherms as a function of hydrogen concentration could be obtained. Comparing these with the experimental isotherms it could be safely concluded that a strong tendency exists for the hydrogen atoms to fill the "double holes" with one atom only, which is in agreement with the indications obtained from the changes of the cell dimensions as a function of hydrogen content.

Retrograde solidus curves

The power of this method of pure thermodynamical reasoning has been often demonstrated by Meijering⁹³). It was possible to derive an expression enabling one to predict (e.g. from eutectic data) whether in a T, x projection of a solidus-liquidus equilibrium of a binary system, the solidus would show a retrograde character or not^{93) 94)}. The principle is very simple indeed. In a binary system two phases (solid and liquid, with compositions x_S and x_L atom fractions) are in equilibrium when at a certain temperature

$$\frac{dG_S}{dx_S} = (G_S - G_L)/(x_S - x_L) = \frac{dG_L}{dx_L}.$$

Now for coexistence of the two phases this must be valid also on changing T . Remembering that $\partial G/\partial T = -S$ and $G = H - TS$, the equation of Van der Waals for coexistent phases is obtained:

$$\frac{dx_S}{dT} = \frac{dS_S/dx_S - (S_S - S_L)/(x_S - x_L)}{d^2G_S/dx_S^2}.$$

The denominator being always positive, retrograde solidus lines $dx_S/dT > 0$ will be found for positive numerators. Taking for the entropy only Gibbs's entropy of mixing one gets for S_S and S_L the curves given in fig. 20. It was thus possible with this to explain many retrograde solidus curves and to predict others. For instance, from fig. 19 it could be predicted that the number of retrograde solidus curves would increase enormously with the application of refined techniques of determination of small solid solubilities, a prediction that in the course of the years for the cases of Ge and Si has been amply verified.

Regular solutions; the system Ni-Cu-Cr

A powerful tool used in those thermodynamical considerations^{95) 96) 97) 98)} was the concept of regular

solutions. For binary regular solutions the Gibbs free energy can be written:

$$G = ax(1-x) + RT \{x \ln x + (1-x) \ln(1-x)\}.$$

i.e. adding a simple extra enthalpy term to the entropy of mixing. For a ternary regular system one gets:

$$G = axy + bxz + cyz + RT(x \ln x + y \ln y + z \ln z),$$

with x, y and z the atom fractions of the three components. The use of this relatively simple function was demonstrated for instance for the system Cu-Ni-Cr. From the known binary systems between body-centred cubic Cr and face-centred cubic Ni and Cu the ternary diagram of fig. 21a could be expected, i.e. a small body-centred region in the Cr corner, a face-centred region along the Cu-Ni side and a two-phase region in between. However, metallographical and X-ray investigations^{99) 100)} revealed the existence of an extended region of a three-

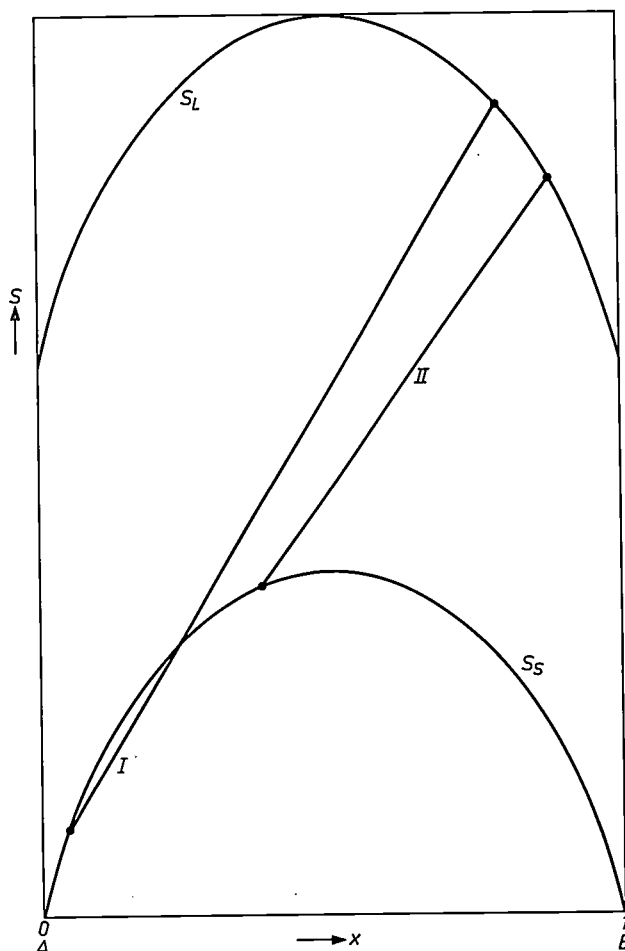


Fig. 20. The lower and upper curves give the entropies S_S and S_L of the solid and the liquid phase respectively as a function of the atom fraction x of B in the binary system A-B. Gibbs's expression for the entropy of mixing is used in both cases, and 1.1R is taken for the entropy of melting per mole. On applying Van der Waals's equation for coexistent phases it is seen that in example I one has $dx/dT > 0$, i.e. a retrograde solidus curve, whereas in example II this not the case.

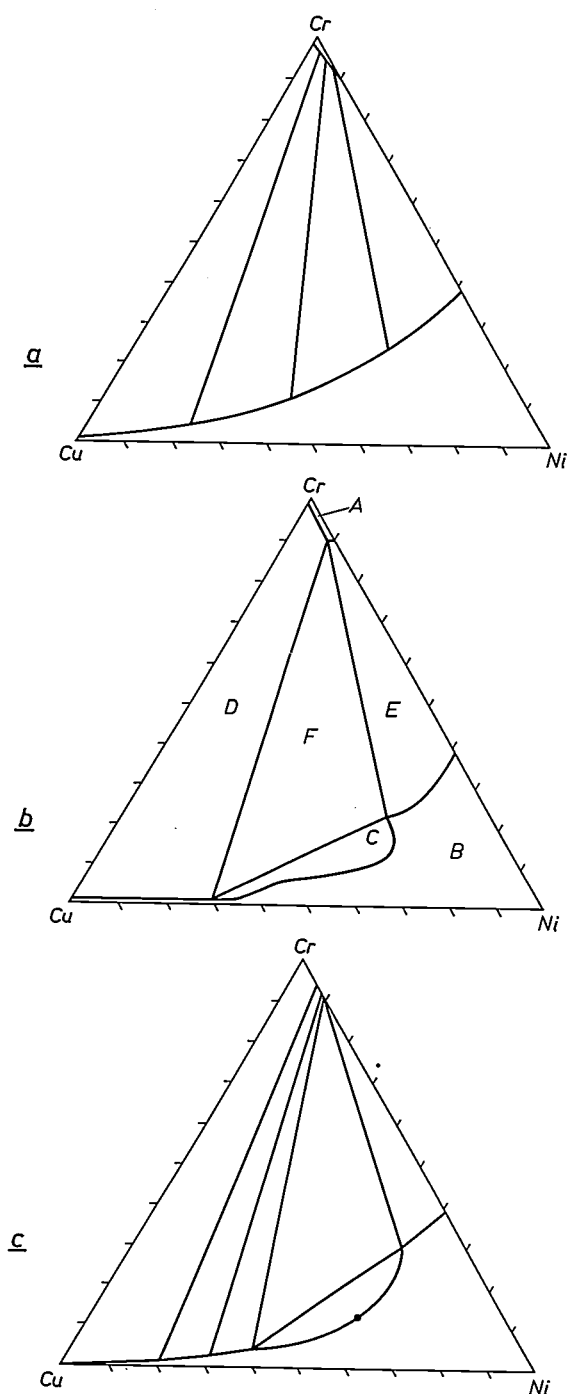


Fig. 21. a) Phase diagram of Cu-Ni-Cr at 930 °C obtained by graphical interpolation of the solubilities in the binary border systems.

b) Real phase diagram of the system Cu-Ni-Cr at 930 °C as determined experimentally.

A: one phase; body-centred.

B: one phase; face-centred.

C: two phases; face-centred.

D: two phases; body-centred and Cu-rich face-centred.

E: two phases; body-centred and Ni-rich face-centred.

F: three phases.

c) Phase diagram of the system Cu-Ni-Cr as calculated for 1200 °K. Between the Cr-rich and Cu-rich phases two connodes have been drawn. The dot denotes the critical point where the Cu-rich and the Ni-rich phases become identical.

phase equilibrium: a Cr-rich body-centred phase and two face-centred phases (fig. 21b). It now proved possible to calculate the values of a , b and c from the binary systems Ni-Cr, Cu-Ni and Cu-Cr respectively. Taking into account also the difference in free enthalpy between body-centred and face-centred Cr, it was possible, using the above expression for G , to calculate a phase diagram (1200 °K) resembling in many aspects the experimental one (fig. 21c).

Only a few items of the application of thermodynamics have been given in this relatively short section. It will be clear that in all areas of chemistry reviewed in this paper the use of thermodynamics is indispensable. For instance the work of Kröger and Vink and that of these authors together with Stieltjes mentioned in the section on internal charge compensation is in essence the application of thermodynamics to the defect chemistry of homogeneous and heterogeneous phase equilibria. The same can be said of the calculations by Van Gool also dealt with in that section.

Conclusion

In this paper a limited number of topics of research related to solid-state chemistry have been briefly dealt with. It is hoped that an impression has been given of the chemical climate in this laboratory. One important point, however, has not yet been stressed and certainly must not be overlooked. Many workers in this laboratory have been mentioned in this paper. It must be emphasized that they could not have done what they did without the continuous help, in the form of papers, lectures and discussions, from scientists in other laboratories. Only a very few of them have been mentioned. Every research worker knows, however, that progress is possible only by means of close cooperation between scientists all over the world.

BIBLIOGRAPHY

- 1) E. J. W. Verwey, *Z. Kristallogr.* **91**, 65, 1935.
- 2) E. J. W. Verwey, *J. chem. Phys.* **3**, 592, 1935.
- 3) E. J. W. Verwey and M. G. van Bruggen, *Z. Kristallogr.* **92**, 136, 1935.
- 4) A. E. van Arkel, E. J. W. Verwey and M. G. van Bruggen, *Rec. Trav. chim. Pays-Bas* **55**, 331, 1936.
- 5) E. J. W. Verwey, A. E. van Arkel and M. G. van Bruggen, *Rec. Trav. chim. Pays-Bas* **55**, 340, 1936.
- 6) J. L. Snoek, *New developments in ferromagnetic materials*, Elsevier, Amsterdam 1947.
- 7) E. J. W. Verwey and J. H. de Boer, *Rec. Trav. chim. Pays-Bas* **55**, 531, 1936.
- 8) J. H. de Boer and E. J. W. Verwey, *Proc. Phys. Soc.* **49**, extra part, 59, 1937.
- 9) E. W. Gorter, *Nature* **165**, 798, 1950.
- 10) F. C. Romeijn, *Philips Res. Repts* **8**, 304, 1953.
- 11) F. C. Romeijn, *Philips Res. Repts* **8**, 321, 1953.
- 12) E. W. Gorter, *Philips Res. Repts* **9**, 295, 1954.
- 13) E. W. Gorter, *Philips Res. Repts* **9**, 321, 1954.

- 14) E. W. Gorter, Philips Res. Repts 9, 403, 1954.
 15) G. Blasse, Philips Res. Repts, to be published.
 16) E. J. W. Verwey and E. L. Heilmann, J. chem. Phys. 15, 174, 1947.
 17) E. J. W. Verwey, F. de Boer and J. H. van Santen, J. chem. Phys. 16, 1091, 1948.
 18) F. de Boer, J. H. van Santen and E. J. W. Verwey, J. chem. Phys. 18, 1032, 1950.
 19) E. J. W. Verwey and P. W. Haaijman, Physica 8, 979, 1941.
 20) P. B. Braun, Nature 170, 1123, 1952.
 21) G. W. van Oosterhout and C. J. M. Rooijmans, Nature 181, 44, 1958.
 22) C. J. M. Rooijmans, J. inorg. nucl. Chem. 11, 78, 1959.
 23) J. H. van Santen, Philips Res. Repts 5, 282, 1950.
 24) J. H. van Santen and J. S. van Wieringen, Rec. Trav. chim. Pays-Bas 71, 420, 1952.
 25) J. L. Snoek, Physica 8, 426, 1941.
 26) E. J. W. Verwey, P. B. Braun, E. W. Gorter, F. C. Romeijn and J. H. van Santen, Z. phys. Chem. 198, 6, 1951.
 27) J. J. Went, G. W. Rathenau, E. W. Gorter and G. W. van Oosterhout, Philips tech. Rev. 13, 194, 1951/52.
 28) H. P. J. Wijn, Nature 170, 707, 1952.
 29) P. B. Braun, Nature 170, 708, 1952.
 30) G. H. Jonker, H. P. J. Wijn and P. B. Braun, Philips tech. Rev. 18, 145, 1956/57.
 31) P. B. Braun, Philips Res. Repts 12, 491, 1957.
 32) G. H. Jonker, H. P. J. Wijn and P. B. Braun, Proc. Instn. Electr. Engrs 104 B, Suppl. No. 5, 249, 1957.
 33) G. H. Jonker, XVIe Congrès int. Chim. pure et appl., Paris 1957, section de chimie minérale, p. 117, Birkhäuser, Basel 1958.
 34) J. Smit and H. P. J. Wijn, Ferrites, Philips' Technical Library, Eindhoven 1959.
 35) E. W. Gorter, Proc. Instn. Electr. Engrs 104 B, Suppl. No. 5, 255, 1957.
 36) J. A. Lely, Ber. Dtsch. Keram. Ges. 32, 229, 1955.
 37) W. F. Knippenberg, Philips Res. Repts, to be published.
 38) F. K. Lotgering, J. inorg. nucl. Chem. 9, 113, 1959; 16, 100, 1960/61.
 39) E. J. W. Verwey, P. W. Haaijman, F. C. Romeijn and G. W. van Oosterhout, Philips Res. Repts 5, 173, 1950.
 40) E. J. W. Verwey, in: Semi-conducting materials, (Proc. Conf. Univ. Reading), p. 151, Butterworths, London 1951.
 41) F. A. Kröger, H. J. Vink and J. van den Boomgaard, Z. phys. Chem. 203, 1, 1954.
 42) J. Bloem, Philips Res. Repts 11, 273, 1956.
 43) D. de Nobel, Philips Res. Repts 14, 361, 1959.
 44) D. de Nobel, Philips Res. Repts 14, 430, 1959.
 45) G. H. Jonker, Proc. int. Conf. on semiconductor physics, Prague 1960, p. 864, Academic Press, New York 1961.
 46) G. H. Jonker and S. van Houten, Halbleiterprobleme 6, 118, 1961 (Vieweg, Braunschweig).
 47) G. H. Jonker, Phys. Chem. Solids 9, 165, 1959.
 48) J. Volger, Disc. Faraday Soc. No. 23, 63, 1957.
 49) F. A. Kröger and J. E. Hellingman, J. Electrochem. Soc. 93, 156, 1948.
 50) F. A. Kröger, J. Opt. Soc. Amer. 39, 670, 1949.
 51) F. A. Kröger and J. E. Hellingman, J. Electrochem. Soc. 95, 68, 1949.
 52) F. A. Kröger and J. Dikhoff, Physica 16, 297, 1950.
 53) W. Hoogenstraaten, J. Electrochem. Soc. 100, 356, 1953.
 54) H. A. Klasens, J. Electrochem. Soc. 100, 72, 1953.
 55) W. van Gool, Philips Res. Repts Suppl. 1961, No. 3.
 56) Y. Haven, Rep. Conf. Defects in crystalline solids, Bristol 1954, p. 261, Phys. Soc. London 1955.
 57) F. A. Kröger and H. J. Vink, Physica 20, 950, 1954.
 58) F. A. Kröger and H. J. Vink, in: Halbleiter und Phosphore (Semiconductors and phosphors), Proc. int. Colloq. Garmisch-Partenkirchen 1956, p. 17, Vieweg, Braunschweig 1958.
 59) F. A. Kröger and H. J. Vink, Solid State Physics 3, 307, 1956.
 60) F. A. Kröger and H. J. Vink, Phys. Chem. Solids 5, 208, 1958.
 61) F. A. Kröger, F. H. Stieltjes and H. J. Vink, Philips Res. Repts 14, 557, 1959.
 62) G. Brouwer, Philips Res. Repts 9, 366, 1954.
 63) W. Albers and H. J. Vink, Chem. Weekblad, to be published.
 64) W. van Gool, to be published.
 65) J. L. Snoek, Physica 6, 591, 1939.
 66) J. D. Fast, Métaux, Corr., Industr. 36, 383, 431, 1961.
 67) J. D. Fast and M. B. Verrijp, J. Iron Steel Inst. 176, 24, 1954.
 68) J. D. Fast and M. B. Verrijp, J. Iron Steel Inst. 180, 337, 1955.
 69) J. D. Fast, J. L. Meijering and M. B. Verrijp, Métaux, Corr., Industr. 36, 112, 1961.
 70) J. L. Meijering and M. J. Druyvesteyn, Philips Res. Repts 2, 81, 1947.
 71) J. L. Meijering and M. J. Druyvesteyn, Philips Res. Repts 2, 260, 1947.
 72) J. L. Meijering, Trans. Metallurg. Soc. AIME 218, 968, 1960.
 73) J. L. Meijering, Rev. Métall. 54, 520, 1957.
 74) J. L. Meijering, Acta metallurgica 3, 157, 1955.
 75) B. Jonas and H. J. Meerkamp van Embden, Philips tech. Rev. 6, 8, 1941.
 76) A. I. Luteijn and K. J. de Vos, Philips Res. Repts 11, 489, 1956.
 77) J. D. Fast and J. J. de Jong, J. Phys. Radium 20, 371, 1959.
 78) J. J. de Jong, J. M. G. Smeets and H. B. Haanstra, J. appl. Phys. 29, 297, 1958.
 79) A. E. van Arkel, Physica 3, 76, 1923.
 80) A. E. van Arkel and J. H. de Boer, Z. anorg. allg. Chem. 148, 345, 1925.
 81) J. H. de Boer and J. D. Fast, Z. anorg. allg. Chem. 153, 1, 1926.
 82) J. H. de Boer and J. D. Fast, Z. anorg. allg. Chem. 187, 177, 1930.
 83) J. H. de Boer and J. D. Fast, Z. anorg. allg. Chem. 187, 193, 1930.
 84) J. D. Fast, Z. anorg. allg. Chem. 239, 145, 1938.
 85) J. D. Fast, Z. anorg. allg. Chem. 241, 42, 1939.
 86) J. H. de Boer and J. D. Fast, Rec. Trav. chim. Pays-Bas 55, 459, 1936.
 87) J. D. Fast, Metallwirtschaft 17, 641, 1938.
 88) J. D. Fast, Chem. Weekblad 38, 2, 19, 1941.
 89) J. H. de Boer and J. D. Fast, Rec. Trav. chim. Pays-Bas 58, 984, 1939.
 90) J. H. N. van Vucht, Philips Res. Repts 16, 1, 1961.
 91) J. H. N. van Vucht, Philips Res. Repts 18, 1, 21, 35, 1963 (No. 1).
 92) J. H. N. van Vucht, Philips Res. Repts 18, 53, 1963 (No. 1).
 93) J. L. Meijering, The physical chemistry of metallic solutions and intermetallic compounds, Proc. Symp. Nat. Phys. Lab. 1958, Vol. 2, paper No. 5A.
 94) J. L. Meijering, Philips Res. Repts 3, 281, 1948.
 95) J. L. Meijering, Philips Res. Repts 5, 333, 1950.
 96) J. L. Meijering, Philips Res. Repts 6, 183, 1951.
 97) J. L. Meijering and H. K. Hardy, Acta metallurgica 4, 249, 1956.
 98) J. L. Meijering, Acta metallurgica 5, 257, 1957.
 99) J. L. Meijering, G. W. Rathenau, M. G. van der Steeg and P. B. Braun, J. Inst. Metals 84, 118, 1955/56.
 100) J. L. Meijering, Chem. Weekblad 51, 438, 1955.

Summary. A review is given of contributions of the Philips Research Laboratories to solid-state chemistry in the past four decades. Four areas are briefly reviewed, viz, crystal chemistry, internal charge compensation, gases and metals, and thermodynamics. Stress is laid on the fact that whereas the physical and chemical properties of a substance depend on its chemical composition, the exact chemical composition can often be determined only by making use of these physical and chemical properties. This interplay of physical and chemical methods, leading eventually to an understanding of the interrelation between physical and chemical properties and exact chemical composition, is illustrated by means of a number of diverse examples.

THE PLANNING OF THE NEW COMPLEX OF BUILDINGS FOR PHILIPS RESEARCH LABORATORIES IN THE NETHERLANDS

by M. J. JANSEN GRATION *).

727.5:658.2:697/699

Immediately after the second world war, Philips Research Laboratories at Eindhoven were faced with a severe shortage of space, and this at a time of an expected increasing tempo in scientific research. By way of a short-term remedy, the first measure was to add a new storey to one of the oldest laboratory buildings (built in 1924); subsequently a new three-storey wing including a lecture theatre was built, and a number of small temporary structures were erected in the laboratory grounds¹⁾. Yet all these rather expensive alterations and extensions were not enough to keep up with the fast growth of Philips research activities and, apart from that, there was at last no more building space left.

The Research Laboratories, built in 1924, lay originally in tranquil rural surroundings, but they have gradually become hemmed in as a result of the rapid expansion of the Philips factories and of the Eindhoven urban agglomeration. In the form of noise, vibration, dust and electrical interference, this has constituted a steadily increasing hindrance to research.

In 1958 a decision was taken to build a new laboratory complex outside the town.

The design of the new complex

From the decision to design a new complex the following questions arose:

1. How big would it have to be? In other words, what numbers of staff would have to be accommodated in the near future?
2. Was the whole staff to be accommodated in a single building, or would it be preferable to build several laboratory blocks?
3. Were the buildings to be low (i.e. the usual two or three storeys) or were multi-storey blocks preferable?

In 1958, when the decision was taken, our staff numbered about 1350. To provide a guide as to the size of the complex, it was estimated that by 1970 the staff would number 3000, each man being allotted a floor space of 25 m² net (his own working space)

and 50 m² gross (including a share of corridors, lifts, toilets etc.). The original idea was to erect low buildings; with a building density of 20%, this meant that a 185-acre site would be required.

We managed to secure a site of about 275 acres in the municipality of Waalre, a few kilometres to the south of Eindhoven. As may be seen in *fig. 1*, the site is bounded by a stream, a disused railway and the future main traffic artery from Antwerp to the Ruhr, now under construction. The former

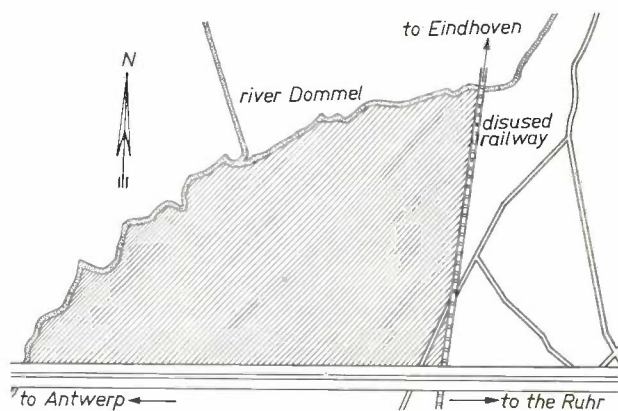


Fig. 1. Ground plan of Waalre site, on which the new complex is being built for Philips Research Laboratories.

railway track is being turned into a tributary road to Eindhoven, so that other Philips sites, factories and offices will be easily accessible from the new laboratory site. There will be no direct road connection with the main highway in order to avoid attracting even more traffic to the vicinity of the laboratories.

The laboratory staff are divided up into a fairly large number of groups covering different research fields. Would their work suffer if the groups were accommodated in separate buildings? A study of this question led to the conclusion that such was in fact to be expected, since the groups are strongly interlinked by personal contacts. Two factors had however to be kept in mind: even if all staff were accommodated in one large block it would still be difficult, with such large numbers, to facilitate all the personal contact that might be desirable. The splitting up of staff into a number of buildings is therefore of less consequence than first considerations would suggest. Moreover, part of a single gigantic laboratory building would inevitably remain un-

*) Philips Research Laboratories, Eindhoven.

¹⁾ A ground plan of the old laboratory complex and a view from the air can be found in Philips tech. Rev. 13, 31 and 2, 1951.

occupied for several years, and this would mean a lower return on outlay. It was finally concluded that a wise decision would be to split the complex into four "sectors".

We use the word "sector" to represent the combination of one large laboratory block, a workshop, and one or more subsidiary buildings for special experiments which, for safety or other reasons, cannot be performed in the large laboratory. A reading room, lecture theatre and canteen also form part of each sector. Each sector is thus intended to be substantially self-supporting.

Elaborating this master plan we envisaged an additional building at the centre of the site for the offices and general services. This building is also intended to include a large restaurant, a central library and a spacious auditorium, suitable for big meetings and congresses. The restaurant will serve hot meals, and each of the sector canteens will supply snacks. The reading room of each sector will only contain the latest books and periodicals relevant to the scientific fields covered by the sector: the central library will hold the complete collection of books and bound volumes of periodicals, including the earlier ones.

A closer study of the matter led to a decision that the four large laboratories should be multi-storey buildings of moderate height. One of the considerations that operated here is that a tall building with sufficient high-speed lifts provides better facilities for personal contacts between staff than a low building. In the latter case the walking distances will discourage mutual consultation. Greatly increasing the height of the buildings, on the other hand, would involve sacrificing a large proportion of the available floor space for lifts etc.

The above paragraphs indicate the general concept we had formed of the new laboratory complex, and the time had now come to get in touch with architects and other experts in order to work out an acceptable design. As was to be expected from a project on this scale — on which, moreover, little guidance was available from past experience or from the literature — this involved a great deal of time and effort. Some impression will now be given of the requirements, often conflicting, which had to be met and reconciliated.

As mentioned in the foregoing, the laboratory management was anxious that there should be the widest opportunity for contact between members of the staff. It was further necessary, within any one sector, that any member of that sector should have easy access to (a) the workshop, (b) the reading room, (c) the canteen and (d) the car parks.

It was also considered desirable to have covered passageways between the various sectors.

Further, a most important requirement was that work in the laboratory should be able to proceed with the minimum of interference from disturbances. This implied measures against noise, vibration, electrical and magnetic interference, etc. In this connection we specified, amongst other things, that the laboratories must be built at a distance of more than 300 m from the main highway. According to measurements by American experts, which have been confirmed by our own investigations, at that distance the noise from a highway is reduced to the average background noise in a quiet residential area. Also, we wanted the laboratory workrooms to face north wherever possible, so as to minimize complications arising from temperature fluctuations due to direct sunlight. Dust might be another nuisance, and the requirement of minimizing this had a considerable bearing on the final design.

Finally it was required that the laboratories should be flexible in operation.

Arrayed against these requirements were others of an architectural and constructional nature, concerning (a) the rational employment of the available ground, (b) flexibility from the standpoint of town planning, (c) aesthetic aspects and (d) clarity of layout. Some of these latter requirements will be discussed in more detail. Although in the first place it had to be ensured that the complex would be able to accommodate a staff of 3000 by 1970, it was also necessary to reserve enough land for expansion in the more remote future. The basic plan had to be flexible enough to allow for such expansion. There were a number of questions to be borne in mind in this connection. Would it be possible to deviate from the basic plan if desired when the complex came to be expanded, without undue difficulties in regard to service supply mains on the site, the smooth handling of traffic and so on? Would it be possible to use other block sizes, materials and architectural forms and methods for future sectors while still retaining the layout of the basic plan? As regards aesthetic aspects the plan entails a number of high buildings which in due course will occupy the skyline as seen from the nearest part of Eindhoven (which is fortunately only a small residential area). In view of this filling-up of the landscape it is important that the laboratory complex should have an open character. The aspect of the buildings from the site itself also requires due consideration. A monolithic appearance has to be avoided. And from the practical standpoint, if the occupants are to feel themselves at

home and to be able to find their way about easily, it was important that the complex should have a simple and clear layout.

No less than 50 designs were proposed, and each was carefully studied to see to what extent it satisfied the requirements. Five of the designs survived this first scrutiny. To select the best out of these the following procedure was adopted. An "appraisal sheet" was drawn up in which the characteristics to be considered were listed. In respect of each characteristic the experts concerned awarded a rating number from 1 to 5:

- 1: The plan is unsatisfactory from this particular viewpoint.
- 2: —
- 3: The plan is reasonably satisfactory.
- 4: —
- 5: The plan is more satisfactory than the others.

These ratings were multiplied by weighting factors in accordance with the relative importance of the viewpoints under consideration. All the weighted ratings thus obtained for a given plan were then added together.

The appraisal sheet for the plan that was finally approved is given in *Table I*. This plan scored 427; the other four plans scored 357, 342, 320 and 294. *Fig. 2* shows a model of the plan finally accepted.

Designing the laboratory buildings

It was now possible to make a start on designing the actual laboratory buildings. Obviously, various studies were carried out first of all, involving visits to laboratories of similar size and discussions with designers and building experts and with scientists working in modern research centres. In what follows we shall confine ourselves mainly to discussing the evolution of the plans for the multi-storey laboratory blocks.

It is impossible to predict the fields in which a laboratory for fundamental research is going to be actively engaged at a point in the somewhat distant future. Indeed, it is difficult to take account of the current programme at the time when building operations start; after all, the building of a laboratory of this size takes about three years, and experience has taught that a research programme can undergo quite a few changes within that space of time. Our guiding principle was therefore that we must build a *universal* and *flexible* laboratory which it would be possible to adapt to any research programme.

By "universal" we mean a block in which each laboratory workroom is suitable for work in chemistry, physics and electronics, these being the three

Table I. Appraisal sheet for the plan finally approved. The plan has been awarded ratings ranging from 1 to 5 (in the first column) in respect of characteristics specified in the sheet; the weighting factor in the second column is an index of the importance of that particular characteristic as compared with others. Ratings and weighting factors are multiplied and the products (third column) added to give a yardstick for comparison of the plans. The design to which this sheet refers scored a total of 427 as against scores of 357, 342, 320 and 294 for four other designs.

	Rating	Weighting factor	Product
1. Accessibility			
<i>From laboratory</i>			
a) Other sectors	2	4	8
b) Sector workshop	5	5	25
c) Other sector workshops	3	1	3
d) Car parks	4	2	8
e) Central office building and central library	5	3	15
f) Canteen and lecture theatre	4	1	4
<i>From sector workshop</i>			
g) Other sector workshops	2	2	4
h) Whether covered passageways are feasible	2	3	6
i) Connection block WA with sector I	4	4	16
2. Freedom from disturbances			
<i>Distance from main highway</i>			
a) Laboratories (noise and vibration)	4	4	16
b) Sector workshops (vibration)	4	4	16
c) Central office building (noise)	5	3	15
<i>Distance from Eindhoven road</i>			
d) Laboratories (noise and vibration)	4	4	16
e) Sector workshops (vibration)	4	4	16
f) Central office building (noise)	4	3	12
g) Influence of prevailing wind	5	1	5
h) Transmission of noise and vibration between buildings of the group	5	1	5
i) Orientation of laboratory workrooms	3	2	6
3. Town-planning considerations			
a) Aspect from nearby residential area (Bennekel)	5	1	5
b) Aspect from Eindhoven road and entrances to complex	5	5	25
c) Aspect from main highway	4	4	16
d) Aspect from site itself	5	3	15
e) Clear basic plan	5	6	30
f) View from buildings	5	6	30
4. Economics			
a) Rational use of site	2	1	2
b) Service supply mains on site	4	2	8
5. Flexibility			
a) Changes possible in basic plan	5	6	30
b) Changes possible in laboratory buildings	5	8	40
6. Security			
a) Siting of central office building from security viewpoint	4	3	12
b) Number of entrances to complex required	3	1	3
7. Clarity of layout			
	5	3	15
Total			427

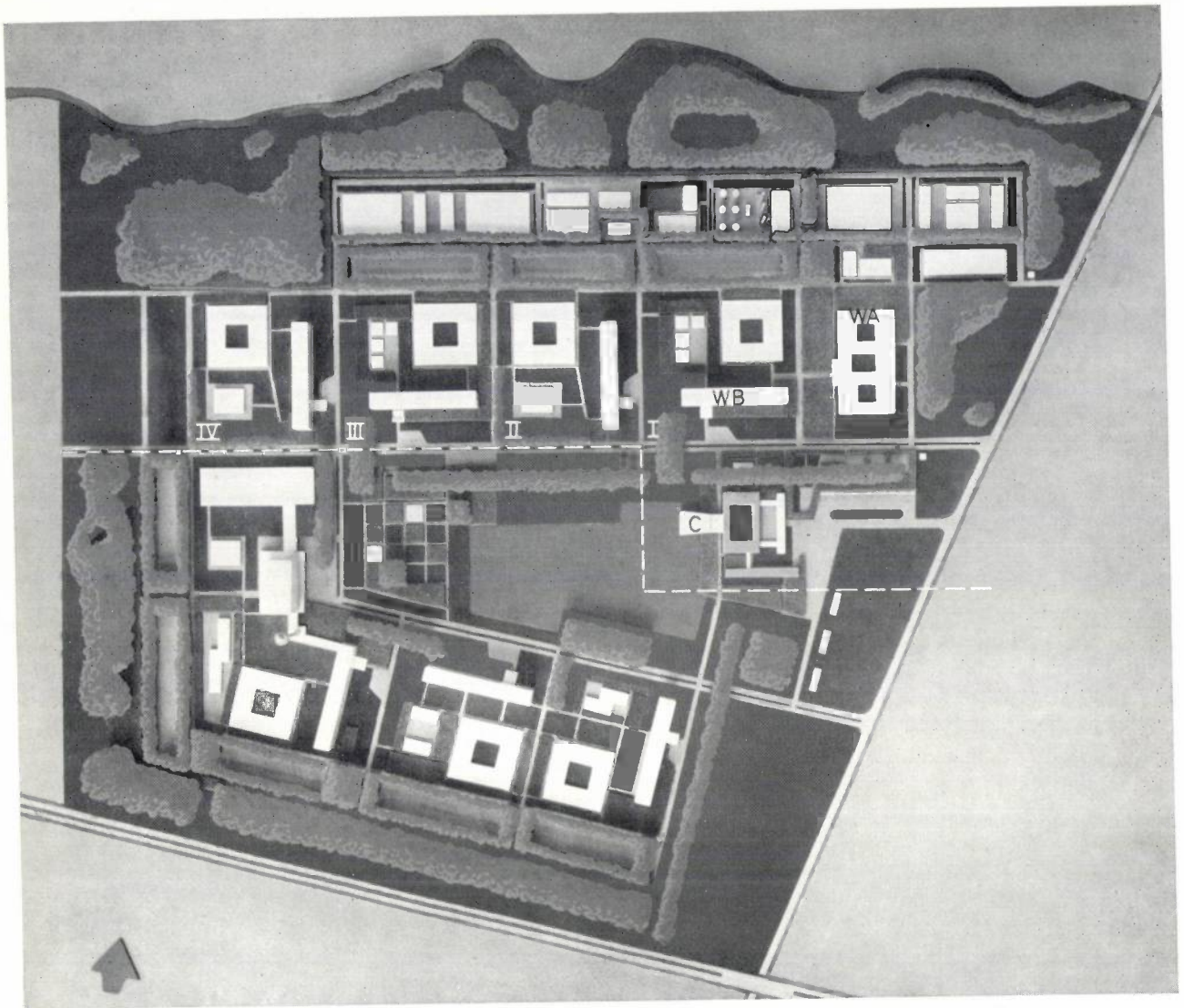


Fig. 2. Model of the new laboratory complex. Above the dotted line are the four sectors I, II, III and IV which, by about 1970, will accommodate a staff numbering 3000. Under the dotted line, buildings that may be erected in the course of a possible extension in the more remote future. Block WA was completed in 1960. Sector I with laboratory WB has now been officially taken into use. C is the central office building.

main classes of research work of interest to Philips. For the sake of flexibility, the aim was to split up the building into rooms with standardized sizes, whose construction would allow subsequent modifications. This would facilitate any regrouping necessitated by a change in programme.

For this purpose the best design is one in which a certain structural unit is repeated throughout the building. The structural unit settled on comprised a standard laboratory workroom and a standard study room, separated by a section of corridor. The repetition interval or module was fixed at 3.85 m, the standard laboratory workrooms having a depth of 7.5 m and the standard study rooms a depth of 5 m (the floor areas are 28 m² and 19 m² respectively); see the plan in fig. 3. The ceiling height was fixed at 3.5 m.

The module was chosen on the basis, amongst other things, of the dimensions of the furniture to be used and the space it was desired to allow between items of furniture (see fig. 4). The furniture is standardized, having been specially designed for

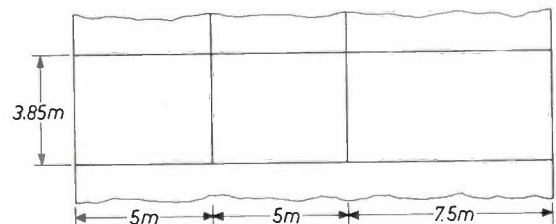


Fig. 3. Ground plan of structural unit of which the multi-storey laboratory blocks will be built up; it is formed by a standard laboratory workroom (right) and a standard study room, on either side of a section of corridor. The repetition interval or "module" is 3.85 m.

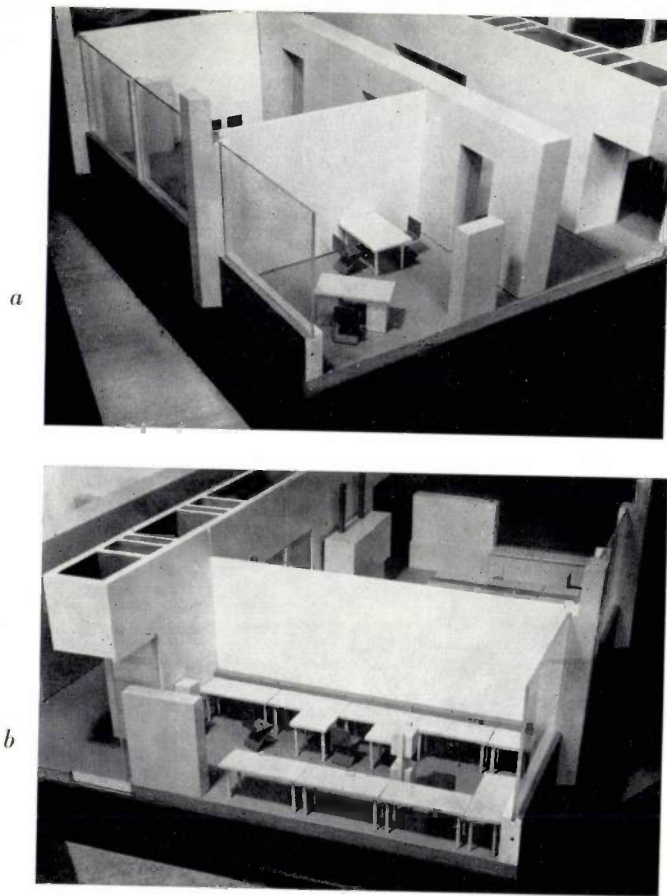


Fig. 4. Model used to arrive at module distance on the basis of furniture dimensions and other data, a) photographed from the study room side; b) from the workroom side.

the new laboratories after a thorough investigation of the use of the laboratory furniture in the existing laboratory.

It was found that in consequence of the modular system just described, the reading room, lecture theatre and canteen and associated kitchen were difficult to fit into the laboratory block. The architects found a very attractive solution to this problem, accommodating these facilities in a low-lying section adjoining the laboratory block — see fig. 5. At this stage it was also decided that the laboratory block would be eight storeys high.

Before a start was made on building, the design was tried out by erecting what amounts to a “slice” of a laboratory block (fig. 6). This small experimental structure stood in one of the quadrangles of the old laboratory, and was used among other things for ascertaining the suitability of the

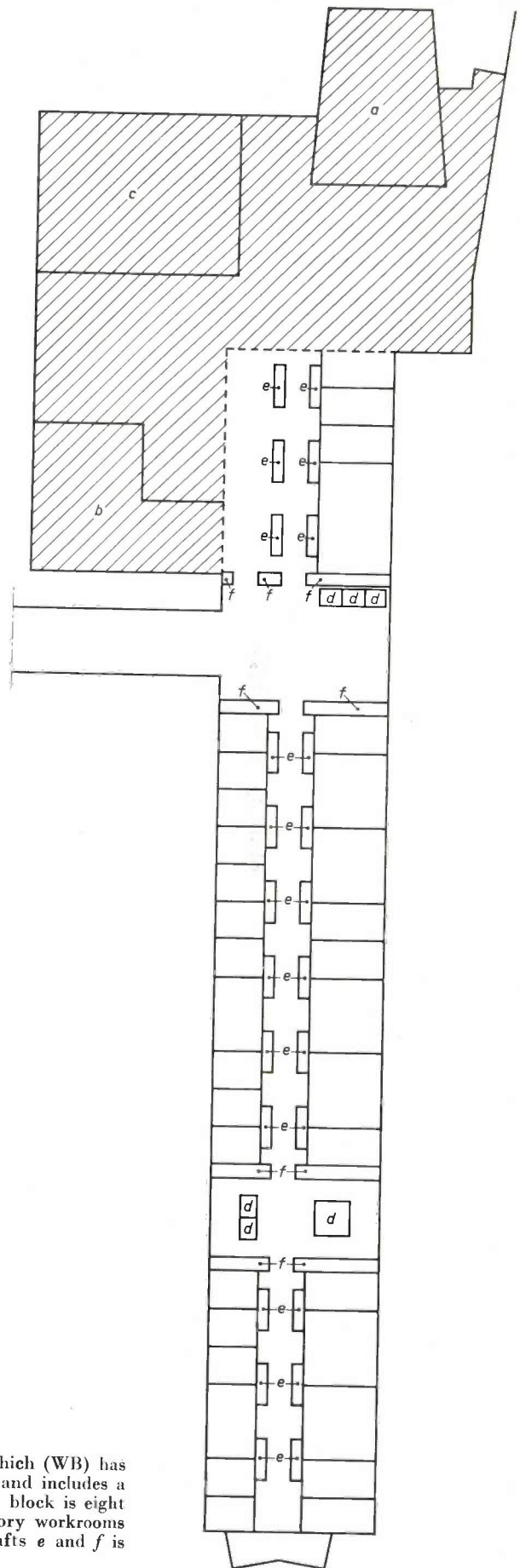


Fig. 5. Ground floor plan of the large laboratory blocks, one of which (WB) has now been completed. The hatched portion is one storey high, and includes a lecture theatre a, a reading room b and a canteen c. The rest of the block is eight storeys high. The study rooms are on the left-hand side, the laboratory workrooms on the right. The spaces marked d are lift shafts. The purpose of shafts e and f is explained below.

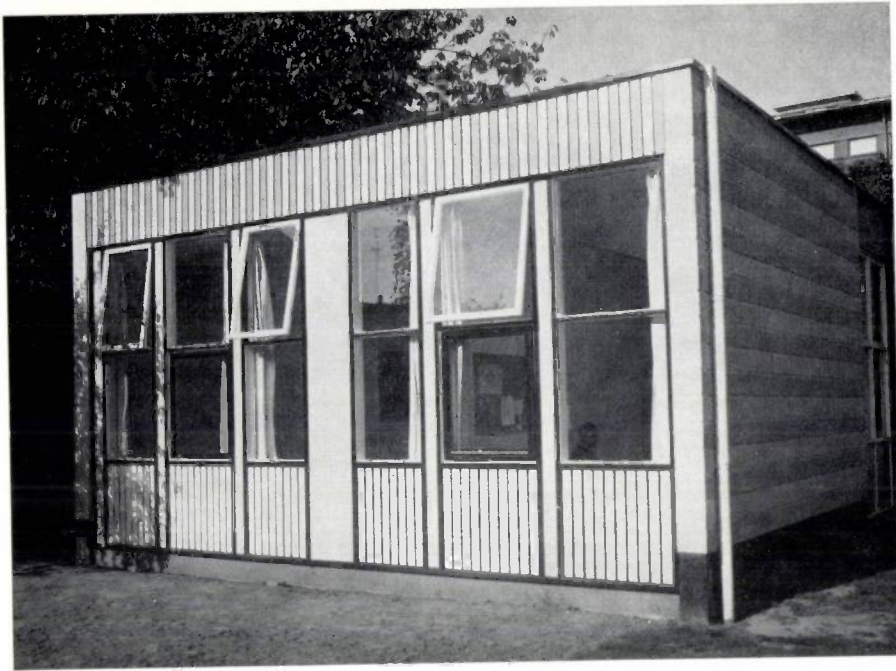


Fig. 6. A small trial structure erected in one of the quadrangles of the old laboratory.

furnishings; it was on the basis of the findings thus obtained that the module was finally settled on.

The erection of a building coded WA on the Waalre site (this building was referred to in the appraisal sheet, and also appears in the ground plan in fig. 2) constituted a trial on a much bigger scale. Block WA is a low laboratory complete with workshop and services. Once the provisional

target has been attained, at which time 3000 staff will be working in the new complex, block WA will be employed as a technological laboratory. This pilot building may be seen in *fig. 7*; to the left in this photograph is block WB, the first of the big laboratory blocks, on the point of completion. Block WB belongs to sector I which has now been officially taken into service.

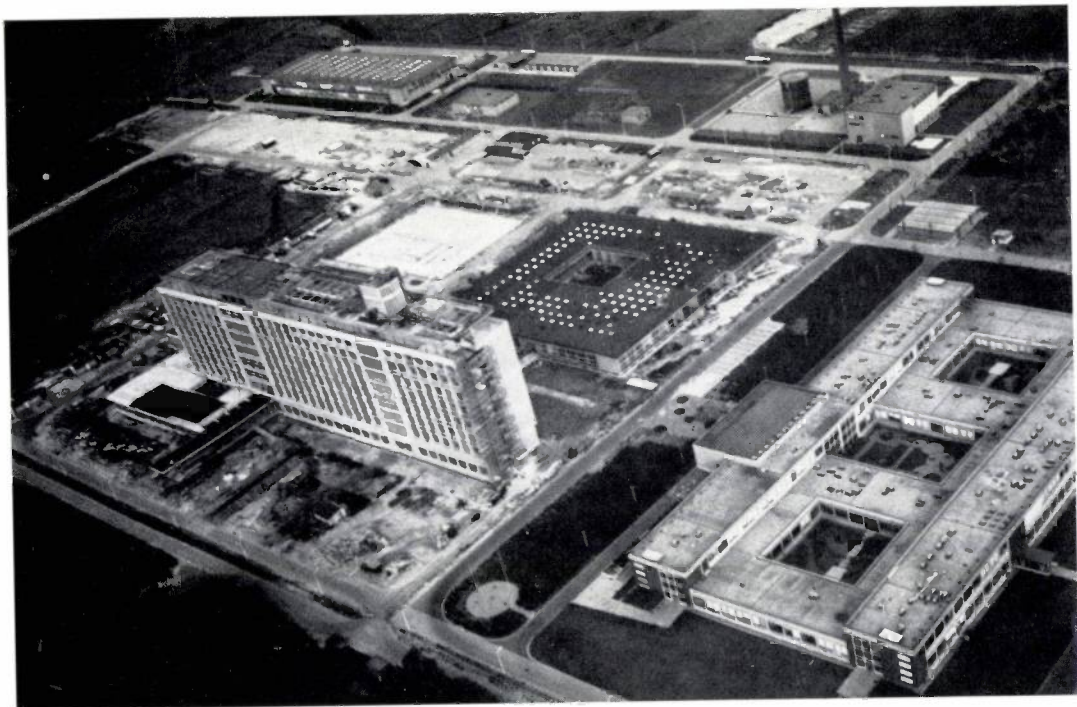


Fig. 7. To the right of this view of the Waalre site is seen the laboratory block WA, built as a pilot project; to the left is sector I with the block WB, the first of the multi-storey buildings.

Electricity, gases, water and drainage

A great deal of attention was devoted to the supply system when the laboratories were being designed. *Vertical shafts* (marked *e* in fig. 5) were planned, through which electric cables and supply pipes for water, coal-gas, nitrogen, oxygen etc. and drainage pipes could be led from the cellar direct to each room. Twelve shafts suffice to supply all the laboratory workrooms. It will be noticed in fig. 5 that 12 shafts have also been built on the study room side. The reason for this is the great shortage of space still prevailing, in view of which it was felt necessary that for the time being a large proportion of the study rooms should be used as laboratory workrooms. When the study rooms finally revert to their intended role it will be possible for the shafts supplying them to be converted into cupboards.

The advantage of a vertical as against a horizontal supply system is that cables and pipes do not need to be laid to a given room until it is definitely certain that they are required; in a building with the horizontal supply system this is not easily done without causing some disturbance to neighbouring rooms. The two systems have been compared in regard to the outlay required: for an eight-storey building like ours, a 12-shaft vertical system repre-

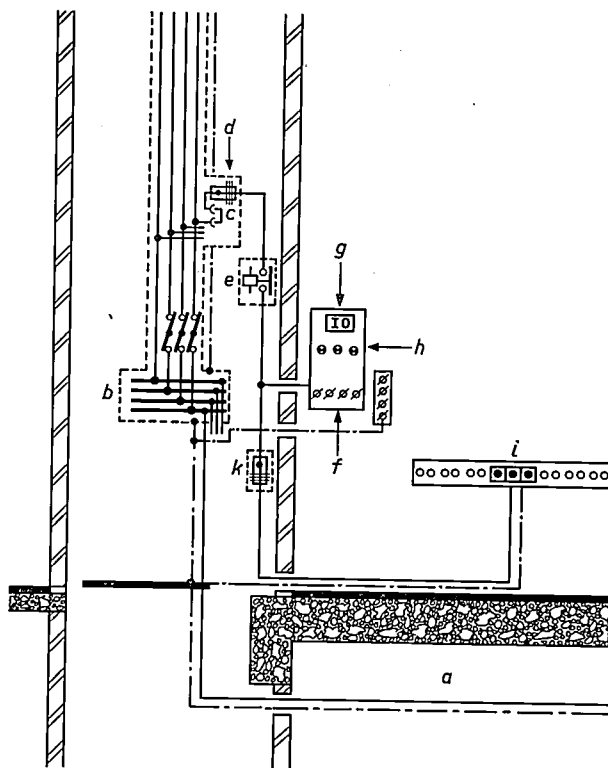


Fig. 9. Vertical cross-section of a supply shaft on the ground floor, and the adjacent corner of a laboratory workroom. Electricity is supplied by a 3-phase system in star connection, providing voltages of 220 V and 380 V. The mains cable runs from the cellar *a* and is tapped beneath each shaft, the branch supply passing through a 4-pole switchbox *b* rated for 200 A. The current from the switchbox is carried by encased copper strips. The copper strips in the shaft are tapped in turn for individual room supplies, which pass through a 4-pole mechanical switchbox *c*, with built-in fuse *d* rated for 60 A and a magnetic switch *e*. In each room is a panel containing terminals *f* for the three phases and the star point, push-buttons *g* controlling the magnetic switch, and three neon lamps *h* indicating which terminals are live. Each room is further equipped with connection boxes *i* containing six 220 V plug-sockets and three 10 A cut-outs; these supplies pass via fuses *k* rated for 25 A, the cables being housed under the wall service strips. A connection to the earth line (chain-dotted line) is also available in each room. If necessary for certain measurements, a special earth line can be installed.

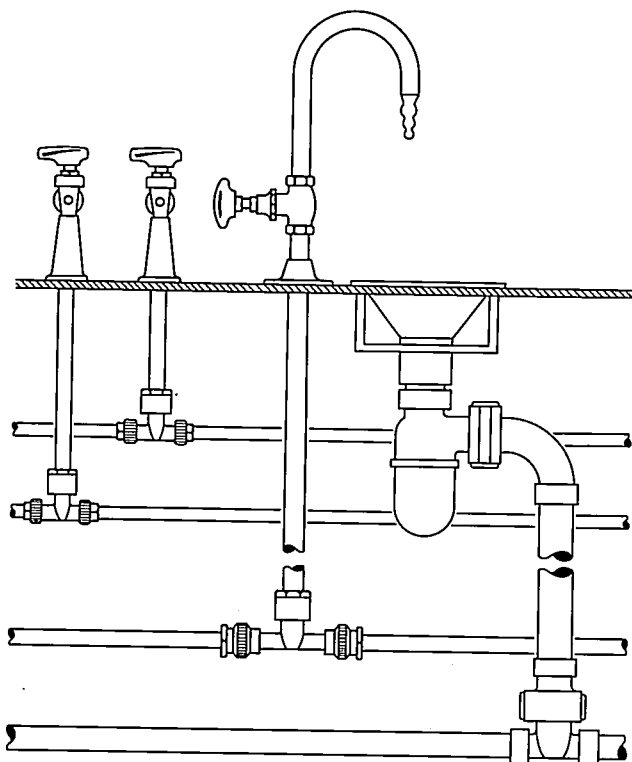


Fig. 8. In the multi-storey laboratory blocks, electricity, water, coal-gas, nitrogen, oxygen and other supplies are piped in via vertical shafts to the rooms of all eight storeys. In each room, individual supply lines are housed under panels ("service strips") 20 cm wide (see fig. 12). Under these panels, the pipes can be tapped at any desired point, in the manner shown schematically in the figure.

sents a saving of a good 50%; the addition of the 12 provisional shafts at the study room side has still allowed a saving of about 25% to be made.

With an eye to minimizing the dust nuisance, the pipes entering the rooms from the shafts are accommodated under wall panels ("service strips"). Special couplings are available which enable the pipes to be tapped at any desired point (fig. 8). The electric cables are also accommodated under these panels. Each room draws its electricity supply from a switchboard which is connected up via the shaft to a 3-phase 380 V/220 V supply in the cellar (for details, see the caption to fig. 9). If direct current is required, it can either be obtained locally from a rectifier or, if a high power is wanted, from a DC generator in the cellar.

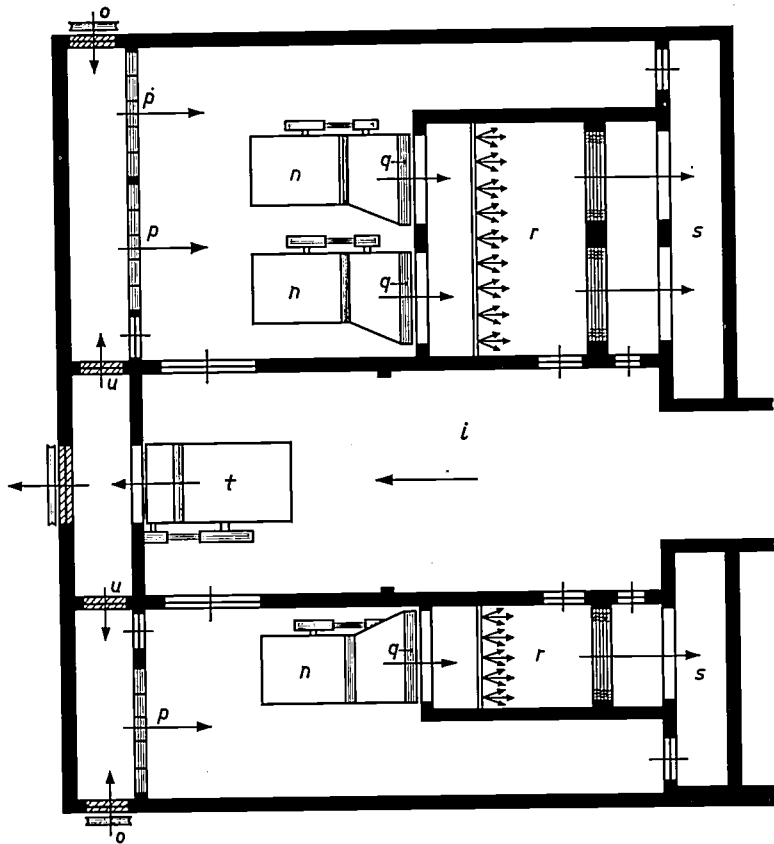


Fig. 10. Plan of ventilation chambers on the roof. The fans *n* draw in outside air through grilles *o* via dust filters *p*, and blow it over heaters *q*, cooling equipment (not yet installed) and humidifier *r* to the openings *s* of the vertical shafts via which the air is circulated throughout the buildings (see fig. 11). The return air collects in reservoir *i* and is blown out by fan *t*. By opening valves *u* to a certain extent it is possible to circulate through the building a mixture of return air and fresh air from outside in the proportion 2 : 1; this conserves heat in wintertime. A complete air change is obtained by closing valves *u*.

Air conditioning

The air in the laboratory has to be regularly changed, heated or cooled as necessary, and given the right degree of humidity. Much of the equipment necessary for this is housed on top of the eighth storey, on the roof of the block — see fig. 10. In the centre is a return air reservoir *i* which is connected with the outside air and with all the rooms in the block. The function of this reservoir will now be explained with reference to fig. 11.

In that diagram a number of laboratory workrooms and study rooms are shown of which *b* and *d* are provided with special exhaust arrangements for fume cupboards, ovens and the like, while those marked *a*, *c*, *e* and *f* lack these facilities. From laboratory workrooms *e* and *f* the air circulates via grilles *g* and shafts *h*, passing from there into the return air reservoir *i*; the air from study rooms *a* and *c* reaches the reservoir by way of the corridors *k*. Having collected in the reservoir, the used air can be freshened by the admixture of outside air and recirculated. The fan *t* expelling the used air, and fans *n* drawing in fresh air, may be seen in fig. 10. When valves *u* are open the latter fans also draw used air from the reservoir and pass a 1 : 2 mixture of fresh and used air into the openings *s* in the vertical shafts via the heating and humidifying equipments lettered (*q*) and (*r*) respectively (and via

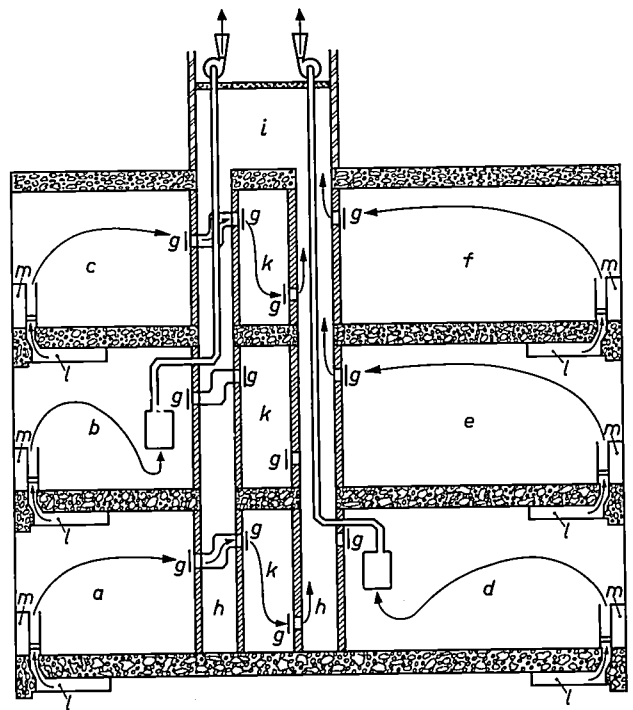


Fig. 11. Vertical cross-section (schematic) through three storeys to show how air circulates within the building. The inflow air grilles are fed via horizontal ducts *l* and vertical shafts (marked *f* in fig. 5) from the ventilation chamber on the roof. Study room *b* and laboratory workroom *d* are equipped with fume cupboards having their own exhaust fans. For the other letters, see text.

cooling equipment if necessary). These special air shafts run from the upper storeys right down to the cellar and are linked to all the rooms by a system of horizontal ducts. The outlets of the ducts in question may be seen in fig. 11 (marked *l*); they supply the inflow air grilles *m* on the window side of the rooms. Air is blown into the rooms at a velocity

in the room) it is necessary to close air grille *g* whenever the central fans come into action. This has the further consequence that the fume cupboard must not be turned off while the central fans are working; otherwise air pressure in the rooms will be above normal, with the risk that fouled air from that room will spread over the whole building,

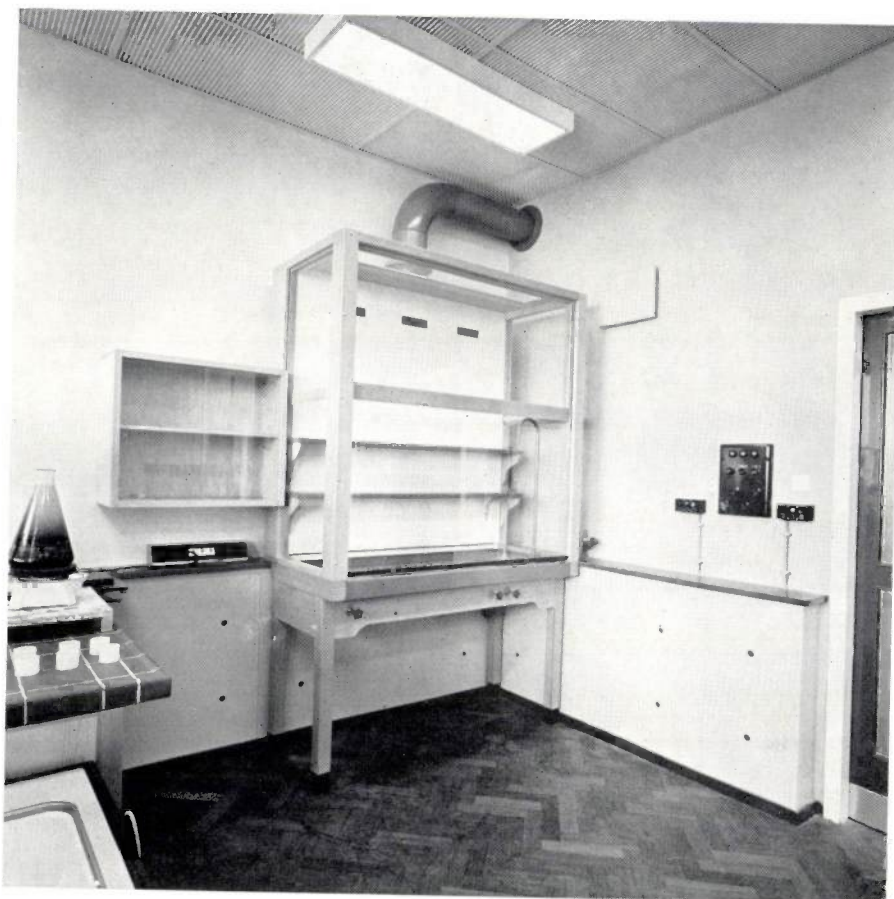


Fig. 12. View of a laboratory workroom equipped with fume cupboard. Left of the latter is a 10 A mains socket box and to the right a 60 A supply box. Above right is the normal ventilation grille which functions when the fume cupboard is not in use. Note also the wall panels under which the supplies are housed.

low enough (about 1.5 m/s) to avoid noise nuisance. In winter the recirculation of used air from the reservoir has the advantage of conserving heat. In summer the air is completely renewed, for which purpose valves *u* in fig. 10 are closed.

The foregoing applies to the ventilation of laboratory workrooms and study rooms not equipped with fume cupboards and the like. The complete air volume of study rooms is changed six times an hour, that of laboratory workrooms ten times an hour.

In rooms equipped with *fume cupboards*, as are *b* and *d* in fig. 10, a special exhaust duct with its own fan conveys the air directly to the exterior (see also fig. 12). To prevent the fume cupboard from "blowing back" (as a result of reduced air pressure

escaping when the door is opened or leaking through cracks when the door is shut. Two courses might be adapted to prevent this happening. Firstly, the air grilles could be permanently closed and the special exhaust arrangements could be turned on and off at the same time as the central fans. Alternatively, the grilles can be automatically closed when the fume cupboard is turned on, and opened when it is stopped. The latter course makes it possible for the fume cupboard to be used as required, and is therefore cheaper; it is the method that will be adopted for the main laboratory blocks.

In most rooms the air removal capacity is not unlimited, being governed by the maximum of six outflow ducts that can be placed in the shaft. This limitation does not apply to rooms on the top

storey, where outflow ducts going direct to the exterior can be provided. Many of these rooms will serve for experimental chemical work. In rooms on other floors destined for chemical work special measures can be taken so that the air can if necessary be replaced at a faster rate than in the ordinary workrooms.

As regards heating, it has to be borne in mind that many laboratory workrooms contain various heat sources in the shape of ovens, burners, high-power transmitting installations and the like. It is therefore necessary that the occupier should himself be able to regulate the temperature of the room. As stated above, fresh air from outside passes through the heating equipment marked *q* in fig. 10. This installation is adjusted to supply heat at a rate such that the temperature of air blown into the rooms will not exceed 16 °C. Each room has separate facilities for further heating in the form of a pair of central-heating pipes running horizontally along the window side of the room — see fig. 13. The inflow air passes over the pipe marked *a*; the other pipe, marked *b*, heats the room by air convection. The overall effect is uniform heating of the rooms and efficient circulation of the air. The highest room

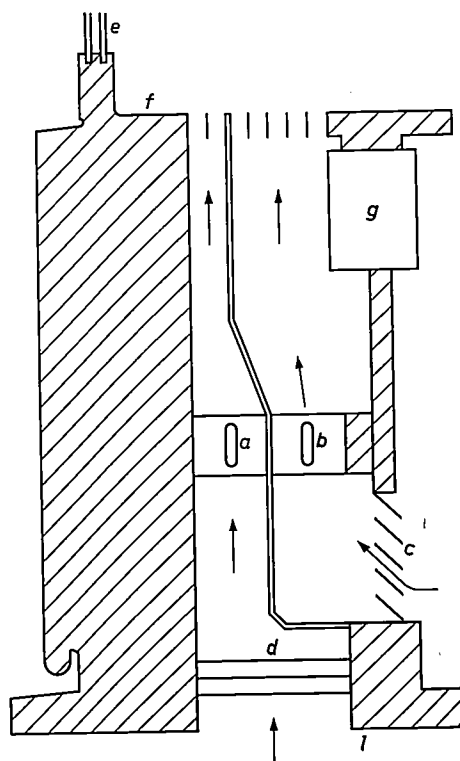


Fig. 13. Cross-section through inflow and convector compartment installed in each room. The air blown into the room passes over the central-heating pipe *a*; at the same time, air from the room, entering via louvres *c*, passes along central-heating pipe *b*. The overall effect is efficient circulation of the air and uniform heating of the room. The valve *d* serves for regulating the rate at which air is blown in; a faster air renewal rate is necessary where fume cupboards are present. *e* double-glazed window. *f* window sill. *g* mounting board e.g. for telephone junction box.

temperature attainable is 22 °C. In summer the air drawn in from outside can be cooled when necessary. Provision has been made for the cooling equipment to be installed either in the ventilation chambers on the roof, or above the ceiling in the rooms themselves. The humidifier marked *r* in fig. 10 is designed to produce a relative humidity of 55% at a room temperature of 20 °C.

After working hours the ventilation system is switched off, and heating of the rooms is by the central-heating pipes only.

Measures to obviate noise, electrical interference and dust

High standards of freedom from noise in laboratory workrooms and, of course, in study and conference rooms, were required. As already pointed out, attention was given to the noise nuisance right from the start. For example, the siting of the laboratories was chosen with an eye on the distance to neighbouring roads. To minimize the nuisance of noise from workshops, these have been housed in separate buildings linked with the associated laboratories by covered passageways. Effective noise suppression was also a consideration in the design of the air circulation system described above and in the choice of building methods and materials — the ceilings, for example, consist of perforated aluminium plates whose surfaces are broken up by a relief pattern, and which are covered with sound-deadening material; the rooms are partitioned off by stretcher walls plastered on both sides.

The transformers for the laboratory electricity supplies are a considerable source of electrical interference. They are installed in the sector workshop, away from the laboratory buildings. Faraday cages are available for interference-free measurements. Where an extremely low level of interference is required the cages can be set up in the cellars, which have a ceiling height of 5 m; below ground level the interference is weaker. The filters employed for suppressing interference in electricity supply lines can then be earthed via short leads giving a low earth resistance.

One dust-prevention measure to which we have devoted much attention is the careful filtering of the ventilating air (by the filters marked *p* in fig. 10). Furthermore, it was to minimize the production and collection of dust that metal ceilings have been used. The various cables and pipes are housed under wall panels for the same reason. Also, the polishes and cleaning powders employed by the janitorial service have been specially selected with an eye to freedom from dust.

Landscape gardening on the site

To conclude, a word may be said about the new landscape to be created in order that the laboratory complex may fit without dissonance into the future agglomeration formed by Eindhoven and its suburbs. On the land adjoining the neighbouring Dommel stream there is a natural woodland in which various species of poplar and willow are represented. The landscape architect intends to plant on this land other kinds of tree that are native to the region, such as elder, hazel, alder, birdcherry, oak, elm and sycamore. As one moves from the edge towards the centre of the site, the wooded area will give way to groupings consisting of more cultivated varieties of plants.

All that has so far been accomplished at the Waalre laboratory complex is the result of close cooperation between the architects and engineers bureau of Philips and a group of members of the Research Laboratory staff. Realizing that what they

had in hand was a big and very costly project, they jointly undertook prolonged studies of the problems associated with the preparation and execution of the first phase of the operation. To the best of our ability we have incorporated their findings into the actual building of the laboratories. We hope that practice will confirm the correctness of these findings, and that the new laboratory will prove to be fertile soil for the research activities of the Philips concern.

Summary. A brief sketch is given of the problems arising in the design of the new complex for the Philips Research Laboratories in Waalre (near Eindhoven). This complex will have to accommodate a staff of about 3000 by 1970. One of the planned four big laboratory blocks has now officially been taken into use. It was necessary to reconcile rather diverse requirements — on the one hand demands from the laboratory management that there should be good facilities for personal contacts between staff members, and freedom from noise, vibration, dust and electrical interference; on the other hand all kinds of architectural and constructional requirements. Certain special aspects of the selected design are discussed, including the vertical shaft system for electricity, gas and water supplies, the air conditioning system, and various measures taken to suppress noise, electrical interference and dust.

SUMMARIES OF DEMONSTRATIONS AND SHORT PAPERS *)

The papers, presented at the Symposium of Sept. 26/27, 1963 at Eindhoven by members of the Philips Research Laboratories, were classed in four sections: A, mainly fundamental; B, devices and materials; C, systems and measuring; D, biochemistry and perception.

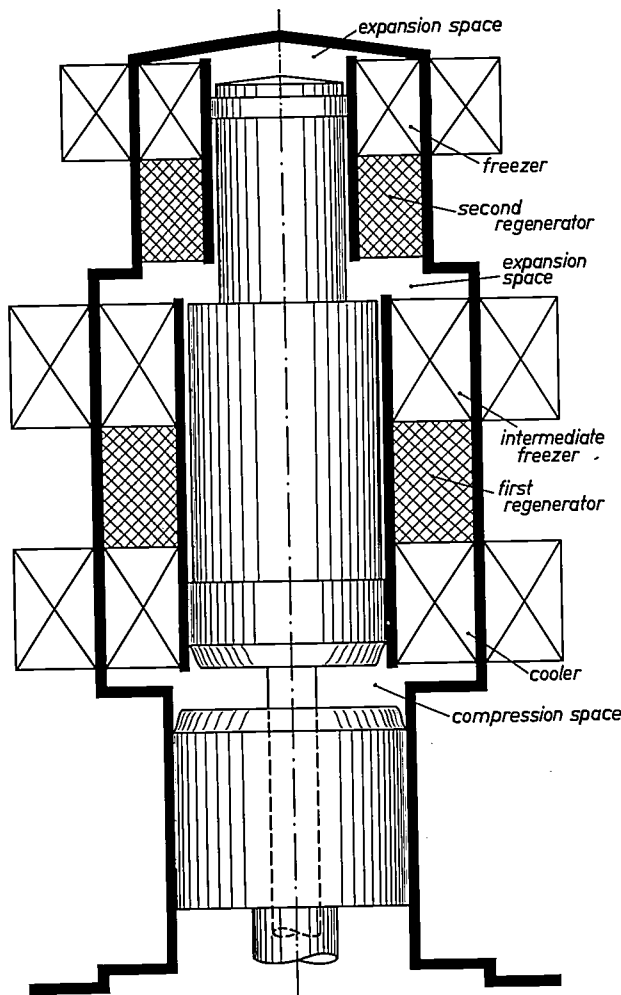
DEMONSTRATIONS

1. A gas-refrigerating machine for the temperature region of liquid hydrogen, by G. Prast.

After a short historical survey of the development of the thermodynamical cycle known as Stirling cycle a description is given of the cycle as it is used in the Philips gas-refrigerating machines. These machines which Philips originally developed for use in the liquid-air temperature region are at the moment well-known in cryogenic technique, primarily due to their compactness and easy handling. Some of their features are: No impurities accumulate in the system as the cycle, consisting of

alternately compressing and expanding a fixed amount of gas, is performed in a closed space. No valves are used in the system. The machines run completely unattended. The efficiency compares favourably with other refrigeration processes used in this temperature range.

A discussion of the properties of the conventional Philips-Stirling system around 20 °K is followed by the description of a modification of the cycle which exhibits all the above advantages in that region too. This modification consists of adding a second expansion space with its heat exchangers to the elements of the conventional type (see illustration). In each of the two expansion spaces part of the gas is totally expanded, the first stage serving as a pre-cooler for the second expansion stage. A prototype of such a "three-space machine" shows quite a high efficiency at 20 °K and reaches 12 °K as its lowest temperature. This machine is demonstrated by cooling a superconducting disc below its critical temperature, the super-conductivity being shown by having a magnet floating over the disc. With this machine we have obtained a general-purpose refrigerator in the region of liquid-hydrogen temperatures. By adding a closed Joule-Thomson circuit the temperature range can easily be extended towards the helium region. Special adaptations include general purpose cryostats for temperatures above 4 °K, liquefaction and recondensation of neon, hydrogen and helium, and cryopumping.



2. The "Plumbicon", by E. F. de Haan and S. L. Tan.

The "Plumbicon" is a pick-up tube operating on the same principles as the vidicon but making use of the photo-conductive properties of PbO. The most significant advantages of the "Plumbicon" are the low dark current (which solved the problems of black-level non-uniformities); the high speed of response, which is independent of light intensity; and the high sensitivity. This means that this type of camera tube, which has the same virtues of simple construction and easy operation as the well known Sb_2S_3 -vidicons can satisfy the demands of broadcast television even at low lightlevels (1 lux on face plate).

The linear light-transfer characteristic up to high signal currents (1 μA) and the satisfactory spectral response ($\lambda_{\text{max}} = 6400 \text{ \AA}$) make the "Plumbicon" especially suitable for a colour television camera. A colour camera using this tube has been constructed which is 4 times as sensitive as a 3" image orthicon colour camera and has a better signal-to-noise ratio and a higher contrast range.

*) The subject-matter of a number of these papers will in the near future be dealt with in articles planned for publication in this Review.

Several varieties of the "Plumbicon" have been developed, or are now being developed. In addition to the standard type to which the above data refer, a type of "Plumbicon" is being developed which has an increased red sensitivity ($\lambda_{\max} = 7500 \text{ \AA}$). By using this type of "Plumbicon" the overall sensitivity of the colour camera is increased about 3 times ($f4.2; 60 \text{ fc.}$) over that using the standard type. By shifting the spectral sensitivity to a

still longer wavelength an infrared "Plumbicon" has been realized without sacrificing the high sensitivity and the high speed of response. The demand for still simpler operation in certain applications has also led to the development of an electrostatically focussed "Plumbicon". In the near future the development of still other variations as e.g. a "Plumbicon" for X-ray detection and a $\frac{1}{2}$ " "Plumbicon" can be expected.

SECTION A, MAINLY FUNDAMENTAL

A1. Study of surface states by photo-electric emission, by J. van Laar and J. J. Scheer.

The spectral distribution curves of photo-electric emission from semiconductors can generally be divided into two parts: a tail of relatively low quantum efficiency at the long wavelength side followed by a steeper rise at higher photon energy. This second part can be ascribed to a bulk process in which valence band electrons are excited into high energy levels in the conduction band and subsequently escape into vacuum. This has been proved by measurements on vacuum-broken silicon crystals of different dope content.

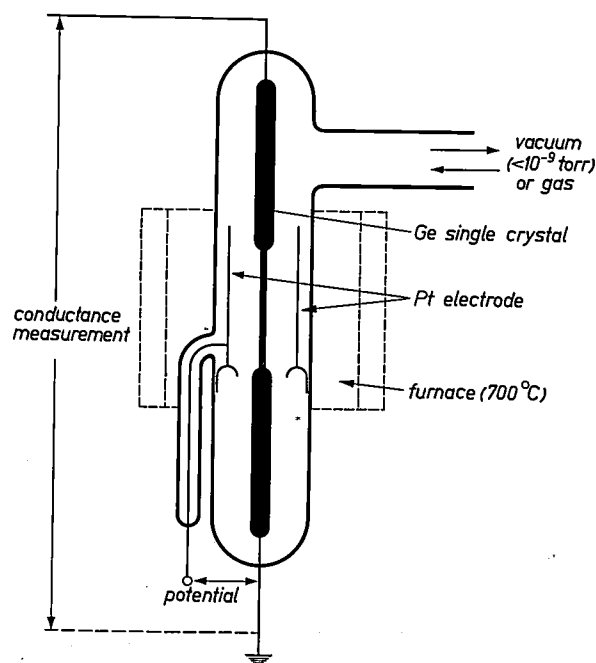
The tail part of the emission curve may be interpreted in several ways. In order to clarify the physical mechanism underlying this part of the emission, several experiments were carried out which strongly suggest that the tail emission originates from filled surface states. It appears that the quantum efficiency of this emission is proportional to the $5/2$ power of the difference between photon energy and threshold energy.

A2. "Clean" germanium surfaces, by M. J. Sparnaay and A. H. Boonstra.

Completely outgassed germanium single crystals were used in high vacuum ($< 10^{-9}$ torr) or in a known gas atmosphere for the determination of surface properties. The germanium crystals were wire shaped: length 25 cm, width 1 mm to 50 micron. The width could be decreased by means of a gas etching procedure: oxygen was introduced at 700°C , giving rise to the formation of GeO which is volatile at that temperature. Complete evacuation, followed by cooling at room temperature (where all the measurements were done), resulted in a surface which we call a "clean" surface. The conductance of the single crystal "wire" was measured as a function of the pressure of various gases and vapours. Also the field effect was measured. To this purpose a platinum cylinder was built around the germanium single crystal in the vacuum system (see illustration) and the conductance of the germanium was measured as a function of the applied potential difference between the platinum and the germanium. This potential difference could be made of the order of 10 kV of either sign. Considerable technical dif-

ficulties had to be met before measurements could be carried out.

The results of the measurements were as follows: there is a weakly *P*-type space-charge region at the cleaned surface of germanium; small amounts of various gases ($< 10^{-3}$ torr O_2 , H_2O , CH_3OH , NH_3 , CO , Hg) lead to an increased *P*-type character. Other gases (H_2 , A , N_2) have no effect. Addition of O_2 , H_2O , CH_3OH , NH_3 to higher pressure (10^{-2} torr) leads to a decrease of the *P*-type character. Occasionally, by suitable combinations of gases, weakly *N*-type space charges can be obtained. Field effect measurements decided whether the space-charge regions were *P*-type or *N*-type. At the cleaned surface the field effect is extremely small, both with applied positive and negative potentials. Therefore in this case most of the induced charges must be stored in (fast) surface states of donor and acceptor character whose density is about 10^{14} cm^{-2} . The existence of donor states must be held responsible for the weakness of the *P*-character of our "clean" surfaces. This adds to the evidence against those interpretations which assume "clean" surfaces to be strongly *P*-type.



A3. On the electrical conduction of semiconducting barium titanate, by G. H. Jonker.

According to Haayman, Dam and Klasens¹⁾ semiconducting BaTiO₃ can be obtained by substitution of small amounts of large trivalent ions like La³⁺ for Ba²⁺, or of small pentavalent ions like Nb⁵⁺ or Sb⁵⁺ for Ti⁴⁺. A maximum of the room temperature conductivity of about 0.1 Ω⁻¹cm⁻¹ generally appears at a concentration of 0.2-0.3%. Above 0.6% the conductivity disappears. This effect can be explained as a transition from compensation of the foreign ions by electrons to compensation by vacancies.

A remarkably high resistivity of 10⁴-10⁶ Ωcm appears above the Curie temperature of BaTiO₃ (~120 °C). This property is not found in single crystals²⁾ but only in polycrystalline samples. From the frequency dependence³⁾ and the field strength dependence⁴⁾ of the resistivity it is now clear that this effect is caused by grain boundary resistances. Moreover, Gerthsen and Hårdtl⁵⁾ have demonstrated the presence of such resistances in a direct way. The origin is mainly surface oxidation.

On cooling the sample this resistance disappears near the Curie point. For some samples this is a gradual effect that can be explained by the theory of Heywang⁴⁾ as an influence of the dielectric constant on the height of the electric barrier. Other samples, however, show a sudden transition at the Curie point, that can be understood by considering the typical ferroelectric properties of BaTiO₃. Most samples show a mixture of these two effects.

1) P. W. Haayman, R. W. Dam and H. A. Klasens, German Patent No. 929 350 (1955), Dutch Patent 15th February 1957.

2) G. Goodman, J. Amer. Ceram. Soc. 46, 48, 1963.

3) O. Saburi, J. Phys. Soc. Japan 14, 1159-1174, 1959.

4) W. Heywang, Solid-state Phys. in Electronics and Telecomm. 4, 877, 1960. W. Heywang, Solid-state electronics 3, 51, 1961.

5) P. Gerthsen and K. H. Hårdtl, Fachausschuss Halbleiter, 1963 (in press).

A4. The critical fields of superconducting lead in relation to its normal state resistivity and its lattice defects, by W. F. Druyvesteyn and D. J. van Ooijen.

The residual resistivity (ρ_r) of a superconductor in the normal state may have a drastic influence on its behaviour in the superconducting state. For an ideal superconducting wire (i.e. one with low ρ_r) the superconducting state is described by zero resistivity ($\rho = 0$) and by a magnetization $M = -H/4\pi$, if a longitudinal magnetic field (H) is applied. If H exceeds a critical value (H_c), then the normal state having $\rho = \rho_r$ and $M = 0$ is restored. The theory as developed by Ginzburg, Landau, Abrikosov and Gorkov (GLAG) predicts a critical value (ρ_{cr}) for ρ_r . If $\rho_r \geq \rho_{cr}$, then the superconducting behaviour is described by: $M = -H/4\pi$ up to a critical field $H_{c1} \leq H_c$; $-H/4\pi$

$\leq M \leq 0$ for $H_{c1} \leq H \leq H_{c2}$ (mixed state); $M = 0$ for $H > H_{c2}$ (normal state).

There are several ways to increase ρ_r of a superconductor, e.g. alloying, cold-working or irradiation. In the latter cases ρ_r is increased by lattice defects: dislocations and point defects. Experiments in several laboratories on superconductors with a high ρ_r have indeed shown the existence of superconductivity at fields well above H_c and a magnetization lying between $-H/4\pi$ and zero, in accordance with theory.

From many experiments it was concluded that the critical current density that destroys superconductivity in the presence of a magnetic field below H_{c2} depends strongly on the degree of cold-working of the superconductor. There exist theories which make use of some of the basic concepts of the GLAG-theory (mixed state) and which can explain the influence of cold work on the critical current. On the other hand alternative explanations of the experiments have been suggested in which the dislocations serve directly as superconducting paths.

Our experiments, which were carried out in order to test the GLAG-theory and to investigate the role of lattice defects, have hitherto yielded the following results:

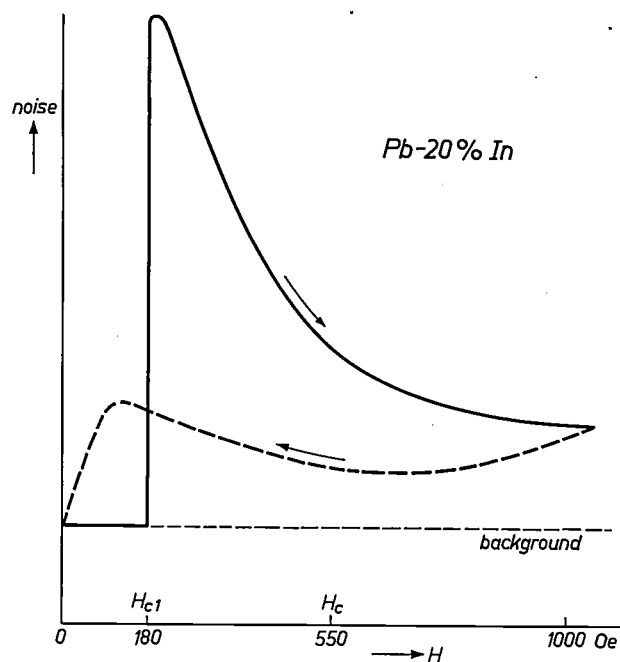
- 1) Magnetization measurements on a series of lead-indium alloys give, in accordance with experiments by Livingston¹⁾, a value of ρ_{cr} and values of H_{c1} and H_{c2} that are in excellent agreement with theory.
- 2) Resistance measurements on the same alloys show superconductivity ($\rho < 0.01 \rho_r$) in fields far above H_{c2} . This cannot be explained by the GLAG-theory.
- 3) Although ρ_r of pure lead, after cold-rolling or neutron-irradiation, both at 78 °K, is smaller than ρ_{cr} , superconductivity above H_c is found using the resistance method. $\Delta H/\Delta \rho_r$ ($\Delta H =$ increase of the critical field, $\Delta \rho_r =$ increase of ρ_r) has the same value after cold-rolling and after irradiation and exceeds many times the values of $\Delta H/\Delta \rho_r$ found for $\rho_r > \rho_{cr}$ on the lead-indium alloys.
- 4) From recovery experiments on the cold-rolled and irradiated wires it was found that the critical field and the critical current density (at fields $<$ critical field) are raised not only by the dislocations but also by the point defects. It seems that the residual resistivity caused by the defects is the crucial parameter, rather than the type of defect.

1) J. D. Livingston, Phys. Rev. 129, 1943, 1963.

A5. Analogon of Barkhausen noise in superconductors, by D. J. van Ooijen and W. F. Druyvesteyn.

In the presence of an external magnetic field a superconductor has a negative magnetization due to partial exclusion of the field (Meissner effect). On increasing the applied field strength large scale penetration of the field, finally leading to destruction of the superconductivity, will occur. This field

penetration is coupled with a change of the magnetization of the superconductor. It was found that in a pick-up coil surrounding the superconductor a noise e.m.f. is generated which accompanies the penetration of the field. From this observation it is inferred that the penetration of the field does not proceed smoothly but by means of small discontinuous jumps. The effect bears resemblance to the well-known Barkhausen noise occurring in a ferromagnetic. With a soft superconductor, e.g. soft-annealed pure lead, noise is observed during the superconducting-to-normal transition (or vice-versa), i.e. at the field strength H_c , the thermodynamic critical field. With a superconductor of the second kind on the other hand, e.g. a lead-indium alloy, the noise is found to start at a field strength $H_{c1} < H_c$ where the superconductor begins to split up into a mixture of superconducting and normal regions. The noise persists, though at decreasing intensity, up to field strengths $H > H_c$. On decreasing the field the noise is observed down to zero field strength, so the sample exhibits hysteresis (see the graph). A recording of the noise, as observed on several samples, is demonstrated.



A6. A superconducting homopolar dynamo for the production of large persistent currents and the problem of energizing a coil, by J. Volger.

In the near future larger persistent currents than have been used hitherto might be required to energize superconducting coils. However, leads bringing very large currents (>300 A) into the helium cryostat cause trouble, in that even with optimum lead design the inevitable Joule-heating and heat-conduction losses consume very large quantities of liquid helium. In this situation a superconducting DC generator in the helium bath would be the solution.

Our laboratory is working on a new type of generator which may be considered as a superconducting variant of the homopolar dynamo. Its essential feature is a moving pattern of a normally conducting region and of magnetic field lines penetrating this region in a superconducting sheet of metal of a special topology. In fact, dynamo and coil form an essentially triply connected body in which the law of conservation of magnetic flux has to be obeyed even when the body (or rather the superconducting part of it) is being deformed. It is shown how the said pattern transfers flux into the coil and how the relevant contour of integration is

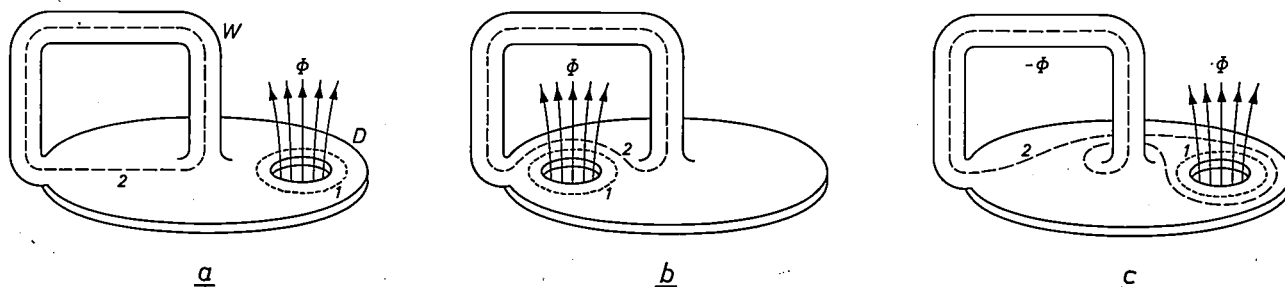
wound around the normal zone (see drawings a-c).

A laboratory model of the dynamo is demonstrated showing the generation of currents fed into a superconducting magnet coil in the persistent current regime. "Slip"-phenomena of the dynamo when running at higher speed are discussed. The recoil torque due to the generated current is demonstrated. Finally, some further prospects of superconducting coils are discussed.

A7. Anomalous transmission of X-rays in perfect crystals, by D. Polder.

Anomalous transmission of X-rays is the phenomenon of the apparent absorption coefficient of the X-rays in perfect crystals being very small if the beam suffers a Bragg-reflection during the passage through the crystal. The phenomenon is explained by the dynamical theory of X-ray diffraction.

In this laboratory Okkerse¹⁾ carried out experiments on dislocation-free germanium crystals. He confirmed among other things that the size and the vibrations of the atoms cause the remaining absorption of the anomalously transmitted beam and he determined the corresponding Debye temperature. He also extensively studied the effect of pur-



posely introduced dislocations and of surface damage²⁾.

Penning and the author³⁾ extended the theory to elastically deformed crystals. They calculated the curved path of the beam in such crystals and the corresponding increased absorption and showed that a radius of curvature of the atomic planes of one kilometer can easily be detected. Relevant experiments by Okkerse⁴⁾ on crystals elastically deformed by bending or by a temperature gradient are in detailed agreement with theory and lead to an independent determination of the structure factor.

Van Bommel⁵⁾ used a divergent X-ray beam from a point source close to a slice of perfect germanium. The intensity pattern produced on a photographic film behind the slice is characteristic for anomalous transmission. Changes of the pattern caused by bending the specimen are precisely predicted by the theory for deformed crystals⁶⁾.

Consecutive anomalous transmission and Bragg reflection on the two legs of one L-shaped perfect germanium crystal have been used by Okkerse⁷⁾ to demonstrate details of the dynamical theory and to evaluate again the structure factor. A change of 10^{-7} radians in the relative orientation of the two legs can be detected. The set-up also makes possible an easy demonstration of the refractive index for X-rays, which is smaller than 1 in materials such as paraffin.

Some of the subjects mentioned in this summary are discussed in the paper.

¹⁾ B. Okkerse, Philips Res. Repts. 17, 464, 1962.

²⁾ H. Goemans and B. Okkerse, unpublished lab. report.

³⁾ P. Penning and D. Polder, Philips Res. Repts. 16, 419, 1961.

⁴⁾ B. Okkerse and P. Penning, Philips Res. Repts. 18, 82, 1963.

⁵⁾ A. J. van Bommel, to be published.

⁶⁾ D. Polder and P. Penning, to be published.

⁷⁾ B. Okkerse, to be published.

A8. Polytypism of SiC: growth, stability and properties of SiC crystals, by W. F. Knippenberg.

The growth of crystalline silicon carbide and the polymorphism of this substance were investigated by subjecting the products of certain chemical reactions to heating procedures. The structure and chemical composition of the crystals formed were analysed and correlated with the growth conditions.

In an isothermal environment in contact with the equilibrium vapour, polycrystalline silicon carbide develops by parallel sheet growth into thin lamellae. Whereas the structure with cubic symmetry behaves as an unstable phase and as such can be formed at all temperatures, the hexagonal and trigonal modifications show a statistical dependence on temperature.

The stability and the occurrence of the polymorphic forms of silicon carbide are finally brought into relation with the sensitivity of the electron configuration of the substance to slight changes in structure and with the kinetics of seed formation and growth. It appears that considerable modifica-

tions of present growth theories have to be introduced in order to explain the observed growth phenomena.

A9. Sensitization of the photoconduction of anthracene by organic dyes, by B. J. Mulder.

The photoconductivity of anthracene as measured with water-electrodes on both sides of a plate-shaped crystal has its origin at one of the crystal surfaces (the one connected with the positive electrode). Absorption of a light quantum in the interior of the crystal gives rise to an excited state (exciton) that can migrate to the surface. There the exciton decays by a process in which an electron from the anthracene is sent into the solution and a positively charged anthracene molecule constituting a mobile hole is left behind. The hole moves through the crystal to the negative electrode on the opposite face where it is discharged. Experiments performed in this laboratory have shown that adsorption of certain organic dyes at the positive crystal-surface increases the observed photocurrent. Under this condition excitons in or near the surface transfer their energy to dye molecules and the resulting excited dye molecules can take up electrons from the anthracene with a very high efficiency. This proposed mechanism is supported by the observation that the increase of the photocurrent is only observed with dyes whose absorption spectra overlap at least to some extent the emission spectrum of anthracene, a necessary — though not sufficient — condition for energy transfer. An even stronger argument in favour of this picture is provided by the finding that a photocurrent in an anthracene crystal can also be obtained by direct optical excitation of the dye molecules by light which is not absorbed by the anthracene.

A10. Preferred magnetic orientations in hexagonal oxides, by U. Enz.

The anisotropy properties of some hexagonal oxides are investigated. It is well known that the anisotropy energy F_K can be described by a power series. Using only the first and second term $K_1 \sin^2 \theta + K_2 \sin^4 \theta$ of this series one can show that there exist three different types of preferred orientation of the magnetization vector with respect to the crystal axis: a preferred direction, a preferred plane and a preferred cone. Several materials are discussed: $BaFe_{12}O_{19}$ (ferroxidure) has a preferred direction and $Ba_2Zn_2Fe_{12}O_{22}$ (Zn_2Y) a preferred plane. $Ba_2Co_2Fe_{12}O_{22}(Co_2Y)$ shows a preferred cone with a semi-vertical angle of about 70° at low temperatures and a plane at higher temperatures. An anisotropy with six-fold symmetry is found in this substance. Another recently discovered example of a material having a preferred cone at low temperatures is $Ba_2Ni_2Fe_{12}O_{22}$ (Ni_2Y). The anisotropy constants of all these compounds have been determined by torque measurements. A demonstration of the anisotropic properties of the material

$\text{Ba}_3\text{Co}_2\text{Fe}_{24}\text{O}_{41}$ (Co_3Z) is given. This material is unique in the sense that it exhibits a preferred cone, a preferred plane and a preferred direction in subsequent temperature regions.

A11. Application of non-conventional classical mechanics in the design of an AVF cyclotron, by H. L. Hagedoorn.

In the design of an azimuthally varying field (AVF) cyclotron many branches of modern science and technology have to be applied. A part of the design which relies essentially on pure classical physics concerns the calculation of the particle orbits in a magnetic field of complex form, which leads to the design of the shape of the magnetic field required for obtaining the desired orbits.

Though in principle the explicit equations of motion of the charged particles in the magnetic field could be used directly, great simplification and much insight is gained if a Hamiltonian formulation of the problem is used as a starting point. From this it can be shown that the particle orbits in a given magnetic field can be simulated by a thin wire through which a current flows. As this is obviously not a method for finding the finally required shape of the magnetic field explicit mathematical expressions for the particle orbits have to be found. Even these are obtained without solving the equations of motion, which are very complex non-linear differential equations. Instead, starting from an expression for the Hamiltonian of the particles in an unconventional form, canonical transformations are applied which eliminate the time dependency and finally result directly in the desired mathematical expressions for the orbits.

A12. Optical studies on gas lasers, by H. de Lang.

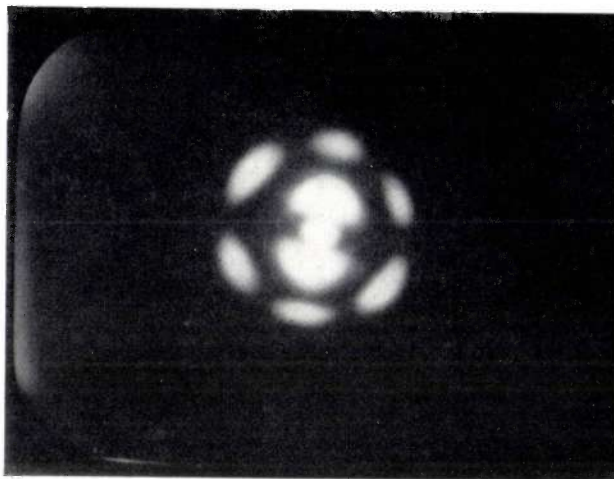
Gas lasers containing totally reflecting prisms are discussed and demonstrated. Instead of the usual 90° -roof prisms, a non-degenerate configuration is shown with roof angles substantially deviating from 90° . This offers the possibility of non-degenerate, cyclic, non-walk-off plane-wave propagation together with greater ease of adjustment. In particular a He-Ne laser ($1.15 \mu\text{m}$) with a tetrahedral configuration of four plane prism boundaries is discussed and demonstrated which operates without requiring any critical adjustment at all.

Another study concerns the problem of investigating the behaviour of single modes of gas-laser oscillations. In order to be able to separate a single mode from the simultaneously occurring ones, a confocal interferometer of one metre length (resolution 6×10^8 at $1.15 \mu\text{m}$) was constructed and is used as a frequency analyser. This interferometer is demonstrated while analysing the output of a small He-Ne laser at $1.15 \mu\text{m}$. Furthermore, the

polarization properties of the oscillations of a gas laser under the influence of a longitudinal magnetic field are demonstrated.

A13. Mode patterns of a laser interferometer, by M. J. Offerhaus.

The beam of the small He-Ne gas laser developed in this laboratory¹⁾ shows mode patterns which run through a cycle as the channel is dilated over $\frac{1}{2}$ wavelength by heating. In order to explain this cycle of mode patterns a theory is presented involving the solution of the scalar wave equation for a Fabry-Pérot cavity between two parallel small flat circular mirrors. The beam is required not to depart much from a flat wave moving along the cylinder axis, with an amplitude distribution narrowly concentrated around this axis.



It is found that each wave function has a cross-section with a Laguerre-Gauss pattern; it moves along the cylinder axis as a travelling wave with an effective wavelength slightly dependent on the symmetry order and on the size of the pattern.

The resonance condition for a standing wave determines the size of each pattern. The selection of a number of modes out of those available is determined by various factors: width of the cylinder, phase errors at the pattern edge, spatial distribution of the gas discharge and competition between overlapping modes. On the other hand, interference between certain slightly overlapping patterns might reduce the phase error and favour superpositions such as that shown in the photograph (rings with radial symmetry numbers $m = 0, 1, 3$).

The various patterns found experimentally and the order in which they occur with varying channel length show good agreement with our theoretical considerations.

¹⁾ H. G. van Bueren, J. Haisma and H. de Lang, *Physics Letters* 2, 340, 1962.

SECTION B, DEVICES AND MATERIALS

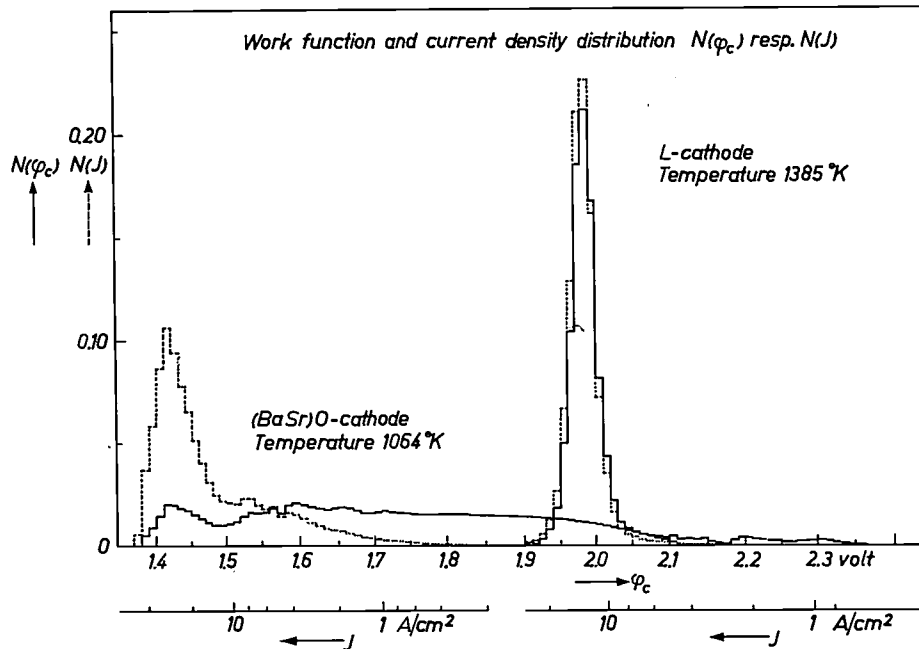
B1. Inhomogeneity and electron energy distribution of thermionic cathodes, by C. G. J. Jansen, A. Venema and Th. H. Weekers.

In the theoretically derived thermionic emission equation it is supposed that the cathode has a homogeneously emitting surface. The best approach to such a cathode, placed in a plane parallel diode, is the single crystal surface. Practical cathodes, however, possess in general a quite inhomogeneously emitting surface, not only because of the differences in work function along the surface, but also as a result of their roughness and their resistance.

The influence of the differences in work function and roughness, leading to a patchiness of the cathode surface, on the I - V characteristic of a diode has been the subject of many theoretical studies. Experimental evidence of the existence of patches

This saturated current distribution can be transformed into an apparent work function distribution. In this way the characteristics of a tungsten single crystal, L-cathodes, impregnated cathodes and (BaSr)O-cathodes have been investigated as a function of their life. The tungsten single crystal has given a very sharp work function distribution. From the other types of cathodes, which are polycrystalline and porous, the L-cathode turns out to be the best approximation to the single crystalline surface (see the graph). The work function distribution of the (BaSr)O-cathode broadens considerably towards high work functions, moreover the main fraction of the saturated current has to be ascribed to a relatively small number of low work function areas.

Additional information can be obtained from the electron energy distribution. The retarding-field



has been obtained by electron-optical means, but this technique supplies only qualitative information about the cathode surface state. Moreover the operating conditions in an electron emission microscope are generally different from those found in electron tubes.

The present authors describe another method which supplies this information in a more quantitative way for different kinds of practical cathodes operating under nearly normal circumstances.

The measuring device consists of a flat cathode and a parallel anode with a small hole (diameter 10 μm). While drawing saturated current under square microsecond pulse conditions the cathode surface is scanned by the anode hole. The current measured in a collector, placed behind this hole, gives the saturated current distribution along the cathode surface with a resolution of about 25 μm .

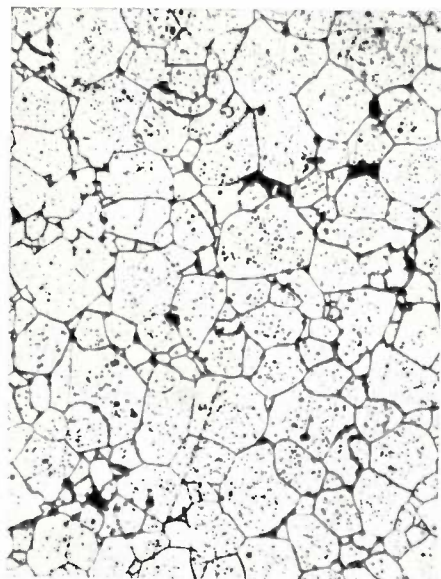
characteristics of the collector current behind the anode hole (in this case with a diameter of 100 μm), while drawing a constant cathode current, supplies knowledge of the resistance and of the surface potential distribution of the cathode. Moreover a knee found in these characteristics of (BaSr)O-cathodes indicates, in accordance with the work function distribution measurements and with enlarged electron-optical pictures, that only a small fraction of the cathode surface consists of areas with low work functions. These areas, however, in the saturated and in a large part of the space-charge-limited region of the I - V characteristics supply the main fraction of the current.

This behaviour of a (BaSr)O-coating affects not only the A -value of the emission formula but also the shape of the I - V characteristics and the noise from the cathode.

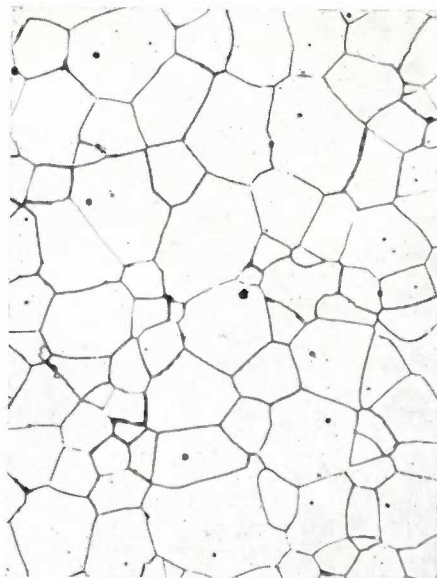
B2. Magnetic deflection coil design using the third-order aberration theory, by J. Kaashoek.

In colour-television reproduction not only the resolution but, dependent on the display device, also the colour registration and purity are very sensitive to deflection errors, these becoming more severe with increasing deflection angle. This has stimulated the research into deflection systems which in combination with the display tube under study give optimum picture quality.

The deflection error theory has already been developed by Glaser and Wendt some twenty years ago. Owing to the fact that practical application of this theory requires very extensive numerical calculations, the only direct result at that time was a better general understanding of the nature of deflection errors. When it was found that there exists a rather simple relationship between the coil shape and the field distribution, some qualitative coil design information could also be derived from the theory. Full use, however, of the theoretical calculations can only be made when the errors are determined quantitatively, which has become possible by the introduction of the electronic computer. As an introduction to this paper the deflection theory is outlined briefly. A discussion then follows of the theoretically calculated third-order errors as regards their dependence on the angle of deflection, the beam parameters and the magnetic field distribution. A general conclusion is that error-free deflection is impossible. In most cases, however, the absence of only some of the errors is sufficient, and in fact, specific deflecting fields do exist which give good overall results. These field distributions can be derived from numerical calculations of the deflection errors and a good knowledge of them has proved to be very helpful in simplifying the coil design procedure. As an example of the practical application of the theory some aspects of the design of a deflector for the shadow-mask colour-television display tube will be dealt with.



a



b

B3. Amplification and generation of microwaves by beam-plasma interaction, by K. R. U. Weimer.

The occurrence of space-charge waves in a system consisting of an electron beam and a plasma is discussed. Under certain circumstances, when the field frequency is smaller than or equal to the plasma frequency a growing density modulation can occur in the system. This growing modulation can be used for the amplification of microwaves. A tube showing the amplifying effect of beam-plasma interaction is demonstrated. This tube resembles a two-cavity klystron operating at 4.2 Gc/s where in the drift space a mercury discharge can be sustained. The modulated beam interacts with the plasma and due to the interaction a notable amplification is observed.

Beam-plasma systems are attractive because no external circuit is necessary to obtain a growing modulation. Therefore, these systems may in the long run lead to interesting sources of mm-wave energy. However, in order to be able to realize such a device an excitation of the modulation without using a premodulated beam is required. A new and simpler method of modulation is in fact applied in another tube that resembles a travelling wave tube. This tube is filled with mercury vapour and the electron beam itself creates a plasma that extends over the whole length of the beam and is sufficiently dense to obtain beam-plasma interaction at 4.2 Gc/s.

B4. Low-porosity ferrites, by A. L. Stuijts.

In the technology of polycrystalline ferrites the heat treatment plays an essential role. During this heat treatment two reactions take place: a solid state reaction of the mixed oxide powders forming the spinel phase, and a sintering reaction which forms a coherent body from the powder.

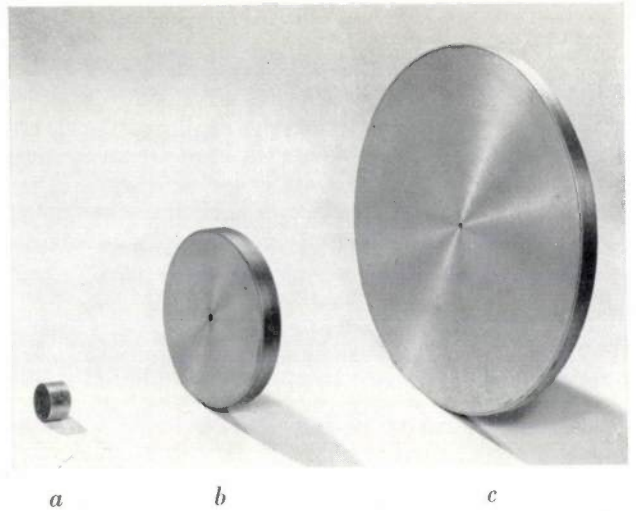
The sintered body has several microstructural characteristics which can vary. As a great number of magnetic, electrical and mechanical properties depend on these microstructural characteristics, control of the sintering process is of primary importance. The mechanism of the material transport during sintering, which results in a densification of the body and thus a decrease in porosity, is understood rather well. However, ceramic technology is still far from being able to provide the tools to change any microstructural characteristic deliberately and independent from the others. The main disturbing factor is the grain growth which occurs during the heat treatment.

As a result of fundamental studies of the sintering and grain-growth behaviour of ferrites prepared from different raw materials and studied in relation to their processing, it has been possible to fabricate practically pore-free ferrites. The choice of the raw materials is a key-point, with special emphasis on purity and particle size, and a complete adaptation of the technology to these powders, mainly in terms of homogeneity. The chemical composition also plays an essential role. The lowest porosities are obtained with compositions based on nickel-zinc ferrites with a small excess of divalent metal oxide with respect to the stoichiometric composition. There need not be a second phase present, however. The photomicrographs on the preceding page (magnification $400\times$) show the microstructure of a) a normal grade and b) a low-porosity ferrite.

The low-porosity ferrites obtained by us show very interesting properties. They are for example excellent for use in magnetic recording heads of intricate construction. The piezo-magnetic coupling coefficient is exceptionally high. The high-frequency magnetic properties are discussed in paper No. B5.

B5. Low-loss high-permeability ferrites in the 1000 Mc/s range, by J. Verweel.

By means of the ceramic techniques described in paper No. B4, very low-porosity NiZn-ferrites can be prepared. When magnetized in a suitable static field H , these ferrites have at the same time a reasonably high value (10-20) of the transverse R.F. permeability and losses which are remarkably low, both being rather independent of frequency up to 2000 Mc/s. These properties are demonstrated by displaying the variation of the resonance frequency

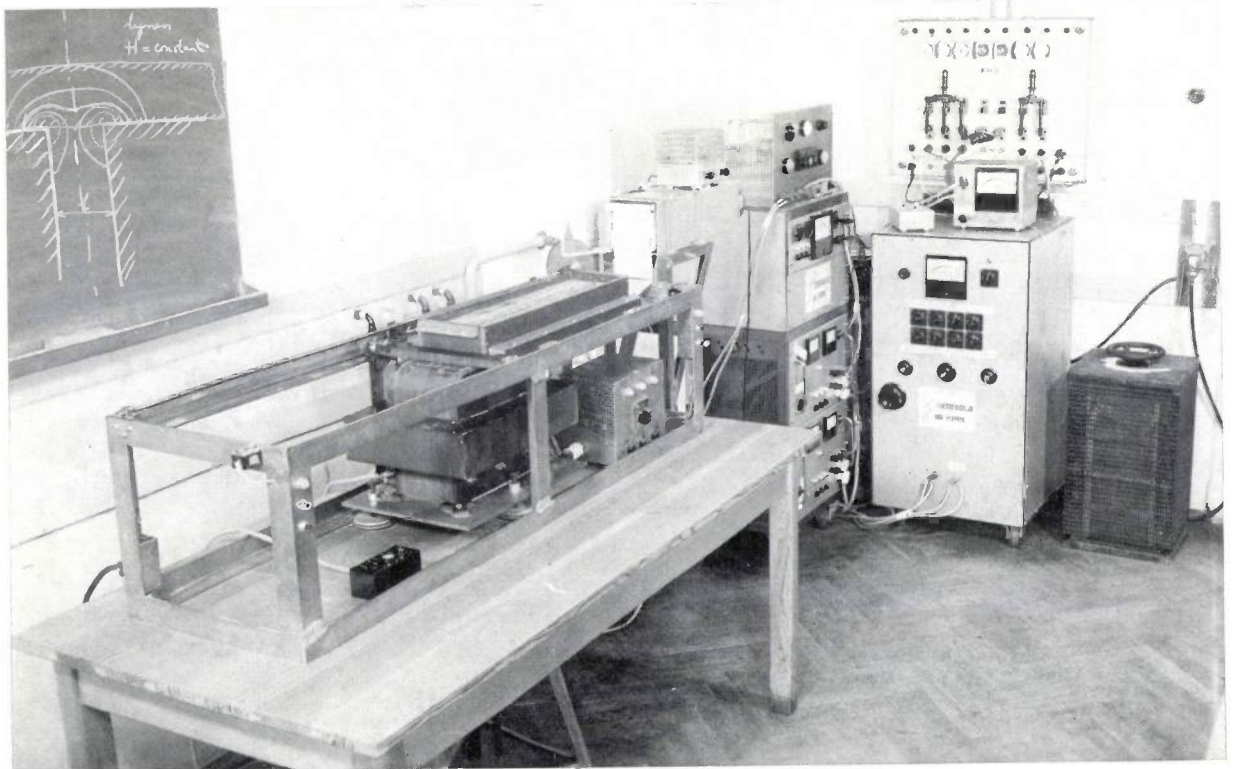


and of the Q -factor of a silver-coated ferrite body with H . From the variation it can be deduced that the permeability varies between 25 and 6 with a loss factor varying from less than 1% to about 2%.

The photograph shows the ferrite body (a, diameter 2 cm). Its resonance frequency can be magnetically varied between the resonance frequencies (in the same mode) of the cavities (b) and (c).

B6. Large scale model of the magnetic recording process, by D. L. A. Tjaden and J. Leyten.

For a better understanding of the magnetic recording process we need more detailed experimental information concerning magnitude and direction of the magnetization as a function of location in the magnetic coating of the tape. This information is not available by direct measurement.



To overcome this difficulty a model of the tape and recording head has been built in which the relevant dimensions are enlarged by a scale factor of about 5×10^3 (see photograph on the preceding page). Disc-shaped samples are taken from the "tape" which are measured with the help of a specially designed magnetometer.

The recorded magnetization can be studied as a function of several parameters, e.g. signal current, bias current, wavelength of the recorded signal. Also the recording of step functions can be investigated.

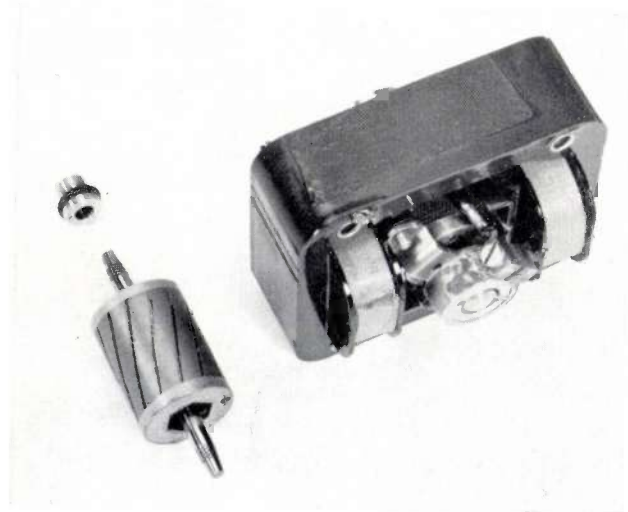
After a demonstration of the model some of the results are discussed.

B7. Self-acting and pressurized bearings, by E. A. Muijlderman.

Lubrication is necessary where forces are transmitted from moving to stationary (motionless) machine parts. Our investigations deal with that type of lubrication which is known as full-film lubrication. This implies that the bearing surfaces are fully separated; the film between the bearing surfaces can be any viscous fluid, e.g. air, water, oil.

There are essentially two classes of solutions to the problem of achieving the above-mentioned separation in practice; one is based on the principle of self-acting lubrication, the other on the principle of pressurized lubrication. A short explanation of the physical mechanism underlying these two principles is given with special emphasis on that of the self-acting bearing. In fact, on this principle the spiral-groove bearing¹⁾ is based which has been investigated in rather great detail in this laboratory. This type of bearing seems to be especially promising for application for small-thrust bearings.

At the end of the paper an apparatus is demonstrated in which the two classes of full-film lubrication are applied. Shown are air and oil bearings of both types. Again special emphasis is laid on the grease-lubricated spiral-groove bearings which are



¹⁾ This type of bearing was first mentioned in the literature by R.T.P. Whipple: Report T/R 622, Atomic Energy Research Establishment, Harwell, rev. Oct. 1958.

likely to go far towards the practical realization of ideal full-film lubrication for small diameters. The photo shows such a bearing as applied in a small electromotor (diameter of shaft 5 mm).

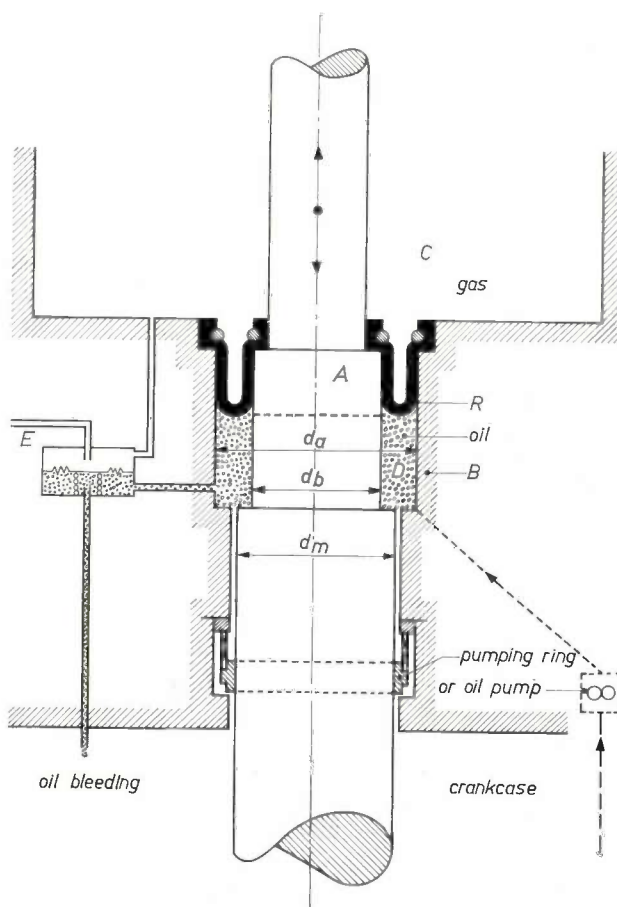
B8. Positive seal for pistons and moving rods, by J. A. Rietdijk.

A short survey is given of the conventional seals for pistons and axially moving rods (piston ring, O-ring, small gap). These seals have to meet the following requirements:

- no gas leakage,
- no oil penetration into the cylinder,
- sufficient lubrication of the piston or the rod.

The conventional methods mentioned above constitute a compromise between these requirements and we can call them only "reducers of leakage".

In Philips Stirling machines the demands on the seals are rather severe. A new method was developed, which completely satisfies the above-mentioned requirements. In this method (see figure) a rolling



diaphragm *R*, fixed to the rod or piston *A* as well as to the cylinder *B*, is used as sealing element. However, rolling diaphragms cannot stand high pressures and pressure variations. This difficulty has been solved by supporting the diaphragm by oil in space *D*. The oil for *D* is supplied from the crankcase either by a small oil pump or by a

pumping ring. The diameters d_a , d_b and d_m are chosen in such a way that the volume of D remains constant during the entire stroke of A . A constant difference is maintained between the pressures in C and D by means of the regulating and safety valve E . In practice $d_a - d_b$ is of the order of 4 mm.

The basic principle is shown by means of a simple experiment. Experience with this new method is discussed: operation times of over 10 000 hours with 1500 r.p.m., stroke 65 mm and gas pressures up to 100 atm have been realized.

Some applications — still in the experimental stage — are mentioned. The diaphragms are manufactured in simple moulds. The material used is a special polyurethane rubber.

B9. Regenerators for Stirling machines, by G. Vonk.

The regenerator is one of the most important parts of machines based on the Stirling cycle. In this cycle gas flows to and fro between two spaces which have different temperatures, and in order to reduce the heat transport due to this alternating gas flow the regenerator is introduced as a heat barrier between these two spaces.

After the regeneration process is explained, the great influence of the regenerator losses is discussed. This influence is indicated by the following data referring to the Philips liquid-air machine: a regeneration loss of only 1% involves a loss of about 20% in refrigerating power of this machine. In view of the importance of the regenerator extensive research is being performed on this subject. In the course of this research a special regenerator-measuring machine has been developed. This machine has proved extremely useful in tracing the laws governing the regenerator process and in determining the properties of newly devised regenerator matrix forms. Some results are discussed, while special attention is paid to the local rate of balance of the gas flow in both directions. The experiments seem to lead to rather surprising conclusions concerning this rate.

B10. The P.D. process, a new photographic tool in electronics, by H. Jonker and C. J. Dippel.

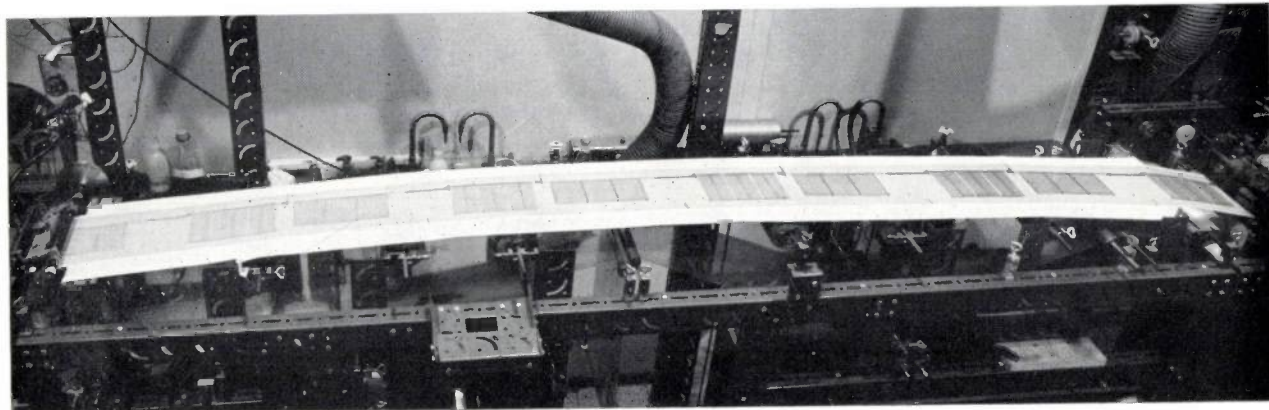
The P.D. process is based on a molecularly dis-

persed photosensitive system, taking up an intermediate position between the expensive Ag-halide emulsion process and the inexpensive diazotype sensitizing process. After exposure metallic nuclei are introduced by a special nucleation reaction and the resulting latent image is selectively intensified to a visible metallic image by the method of physical development or by one of the "electroless" methods. In this way Ag, Hg, Au, Pt, Pd, Cu, Ni and Co may be deposited at the surface of the latent image nuclei. Visible dye-images may be formed by applying physical colour development. Physical developers are in principle thermodynamically unstable solutions, giving rise to spontaneous (non-selective) precipitation of metal. We have succeeded in slowing down this spontaneous deterioration considerably by adding to the developer a small amount of ionic surfactants (by preference cationics). Furthermore, by making use of a reversible redox couple (ferrous/ferric) in composing the developer, which further contains a surfactant-stabilizer, we are able to realize physical developers having a shelf-life of 3 months and capable of developing a latent image within 1 minute, while on the other hand sufficiently stable high-speed solutions may be composed which are capable of developing within 1 second.

The most striking properties of a high-contrast modification of the P.D. process, making it an attractive tool in electronics are:

- extremely high contrast between opaque and clear areas: gamma about 10; maximum optical density > 3 ;
- extremely high resolving power: 1000 lines per mm are resolved with excellent definition;
- extremely high acuity: edge sharpness from max. to min. transmission $< 1 \mu\text{m}$ (the 20th generation of a 100-diameter reduction of a book page on P.D. film is still legible);
- the possibility to provide a wide diversity of flexible and stiff materials with external (adherent, strippable or transferable) patterns, being dimensionally accurate to $\pm 0.5 \mu\text{m}$, ranging in thickness from 300-2000 Å and having resistivities from 200 to 1 ohm per square.

Properties (a)-(c) are of the utmost importance in using the process for making micro-photomasks, which are needed in the photofabrication of electronic devices and microelectronic circuits (e.g.



thin film and solid-state circuits), as well as for making gratitudes for measurement or reference in optical instrumentation and microforms for data storage (storage capacity of P.D. film: 10^8 bits of information per cm^2).

As a consequence of (d) the P.D. process has added quite new interesting possibilities to photofabrication methods by ruling out the use of photosensitive resists in a number of cases. In fact, external, electrically conductive patterns made according to the P.D. process may be used as a base upon which metal, semiconductive or insulating material can be selectively deposited by electroplating, chemical reduction plating, simple immersion plating or electrophoretic coating. P.D.-photoplatting, as this modification of the process is called, may be applied in — among other things — the production of micro metal-masks for evaporation processes, fine metal grids for TV tubes, printed components, flexible printed wiring and static screens. The process lends itself very well to a continuous method, which is much more difficult to realize with the methods based on the use of photosensitive resists. We have built experimental machines for continuous imaging and copper-plating on which trial production of static screens for TV tubes (see photograph on the preceding page) and of flexible wiring is carried out.

B11. Spark doping of epitaxial silicon, by J. Goorissen, H. G. Bruijning and M. Knobbe.

A simple technique for the addition of dope during epitaxial deposition of silicon was studied, which makes use of the tendency of a silicetetra-chloride-hydrogen mixture to react, under proper conditions, with the commonly used impurity elements. These conditions are realized in a spark discharge: ions and radicals are produced which readily react with the electrode material owing to the heat evolved. Doping is accomplished by controlled sparking in the halide-hydrogen mixture between electrodes consisting partly or entirely of the doping element (see illustration). Variability is obtained by controlling the frequency of the spark, which makes the method extremely flexible. The amount of impurity-compound formed per spark depends on

the electrode material and on the amount of energy dissipated.

Reproducible results are obtained with LaB_6 and 0.1% phosphorus-doped silicon as electrode materials, in order to supply boron and phosphorus respectively. Preliminary experiments yield similar results for an analogous germanium system. In the paper experimental details and a discussion of results on silicon epitaxial doping are presented.

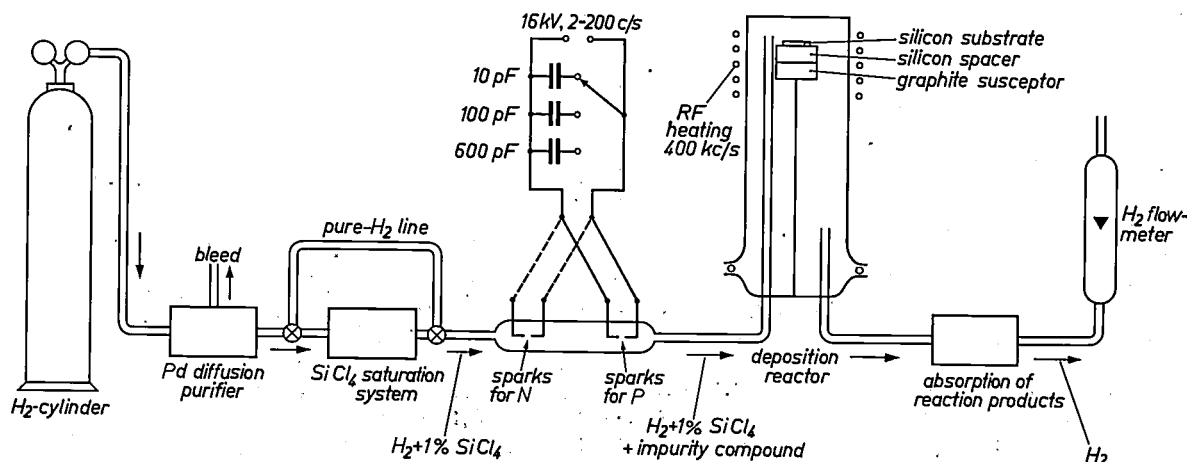
B12. Opto-electronic circuits, by J. G. van Santen.

Opto-electronic circuits consisting of a combination of photoconductors and gas-discharge tubes or of photoconductors and incandescent lamps have already found practical application, for instance in telephone exchanges, in the read-out of counters and in digital voltmeters. A big advantage of this type of circuit is that two types of signal power are present in these circuits so that crosstalk can easily be avoided.

Opto-electronic circuits consisting of electroluminescent (e.l.) and photoconductive (p.c.) elements have some additional advantages. Some of these are:

- 1) E.l. and p.c. layers are well suited to the construction of integrated circuits.
- 2) The circuits can be built compactly.
- 3) The dissipation can be small.
- 4) Low cost.

An e.l.-p.c. integrated circuit consists of two plates, one containing the e.l. elements and the other the p.c. elements. The e.l. substrate consists of a glass plate with an electrode system. On this electrode system is sintered an e.l. layer in glass enamel and on the e.l. layer is sprayed a transparent conductive tin oxide layer. The p.c. substrate consists of an enamelled ceramic plate with an electrode system. On this is sintered a layer of photoconductive material. A good light-coupling between the e.l. and p.c. layers is obtained by placing the e.l. layer directly on top of the p.c. layer with the tin oxide facing the p.c. layer. A thin mask is put between the two layers to avoid short circuits and unwanted light-couplings. The electrode systems under both the e.l. and p.c. layers, which



divide these layers into a number of discrete elements, are obtained by silk-screening a gold containing lacquer. In the paper the properties of opto-electronic circuits are discussed and some applications are shown.

B13. Germanium *N-P-N* transistor for UHF applications, by P. J. W. Jochems.

In the production of transistors for high and very high frequencies several methods are used to reproducibly obtain thin base layers and small distances between the emitter and base contact in order to render capacities and resistances as small as possible. The distance between base contact and emitter region of such transistors is generally determined

by the dimensions of jigs or photomasks used in the process, resulting in distances of at least several microns.

A new technique applied in the fabrication of a germanium double-diffused *N-P-N* transistor for UHF application has been developed which enables us to make the distance between the emitter and the base contact smaller than one micron. This distance is achieved without using mechanical means but only by a temperature- and time-controlled process.

Typical figures for transistors made by this technique are 10-12 dB power gain and 6 dB noise at 800 Mc/s. Combination of this technique with photolithographic methods results in a very small spread of these figures.

SECTION C, SYSTEMS AND MEASURING

C1. Magnetic core computer, by H. J. Heyn.

The design of a medium-speed computer using core logic and a flexible microprogramme is described. The idea is an extension of the possibilities of a word-organized core memory. By appropriate wiring of the registers it is possible to perform micro-operations such as: shift right, shift left, and form the one's-complement. Several groups of registers can be read out at the same time, thereby linearly combining their output. In this way, with the aid of non-linear amplifiers, it is possible to realize threshold functions. The operation of the registers is controlled by a microprogramme, stored in a word-organized memory. Information can stay "written" in this storage permanently by placing small permanent magnets on top of the cores.

The chief goal was to construct a very flexible, inexpensive and very rugged computer. Because of the fact that every micro-operation needs a complete core cycle, and because every macro-operation is built up from several of these micro-operations, it is not an extremely fast computer. The adding of, e.g., two numbers of 42 bits will take 10 core cycles on the average, the times necessary to extract numbers and the instruction itself from the memory not being included.

In a laboratory model (see illustration), in which cores of type 6C1 are used, a read-write cycle-time of about 2.5 μ s can be reached.

It certainly seems possible to apply the same ideas to corresponding types of stores such as e.g. magnetic film memories.

C2. Motional feedback in loudspeakers, by J. A. Klaassen and S. H. de Koning.

In order to cope with the problems arising from the non-linear and frequency dependent behaviour of moving-coil loudspeakers, it seems rather natural

and desirable to try to apply the principle of feedback which is so often used for the reduction of linear and non-linear distortion in electrical systems. For the realization of such a programme ("motional feedback"), it is necessary to have an electrical signal available which is a faithful image of the mechanical cone movement. In principle there are several ways of accomplishing this. In practice, however, all of them have their specific difficulties. A method was found which seems relatively easy to realize. It is based on a measurement of the cone acceleration by means of a piezo-electric element which is acted upon by inertial forces only. Motional feedback not only reduces distortion but provides an adequate means to modify loudspeaker characteristics and adapt a loudspeaker to all kinds of desired applications, among which the use of small enclosures must be mentioned. By the nature of the moving-coil loudspeaker, improvements by motional feedback are restricted to the low frequency range.

C3. Producing "natural" musical sound by electronic means, by N. V. Franssen.

Most electronic musical instruments produce their sound by the switching of AC signals. Conventional musical instruments, on the other hand, may be divided roughly into two classes: firstly, instruments in which pulses of energy are applied to filtering circuits (tympany, piano), and secondly, instruments in which a flow of DC energy is switched and converted into sound by means of generating circuits (violin, reed instruments). In designing electronic instruments according to these latter concepts it appears, however, that in order to obtain a "naturally" sounding result more sophisticated measures have to be taken. Among these, controlled instability of the generators and complex forms of the transients appear to be especially important.

C4. Metrological application of diffraction gratings,
by H. de Lang.

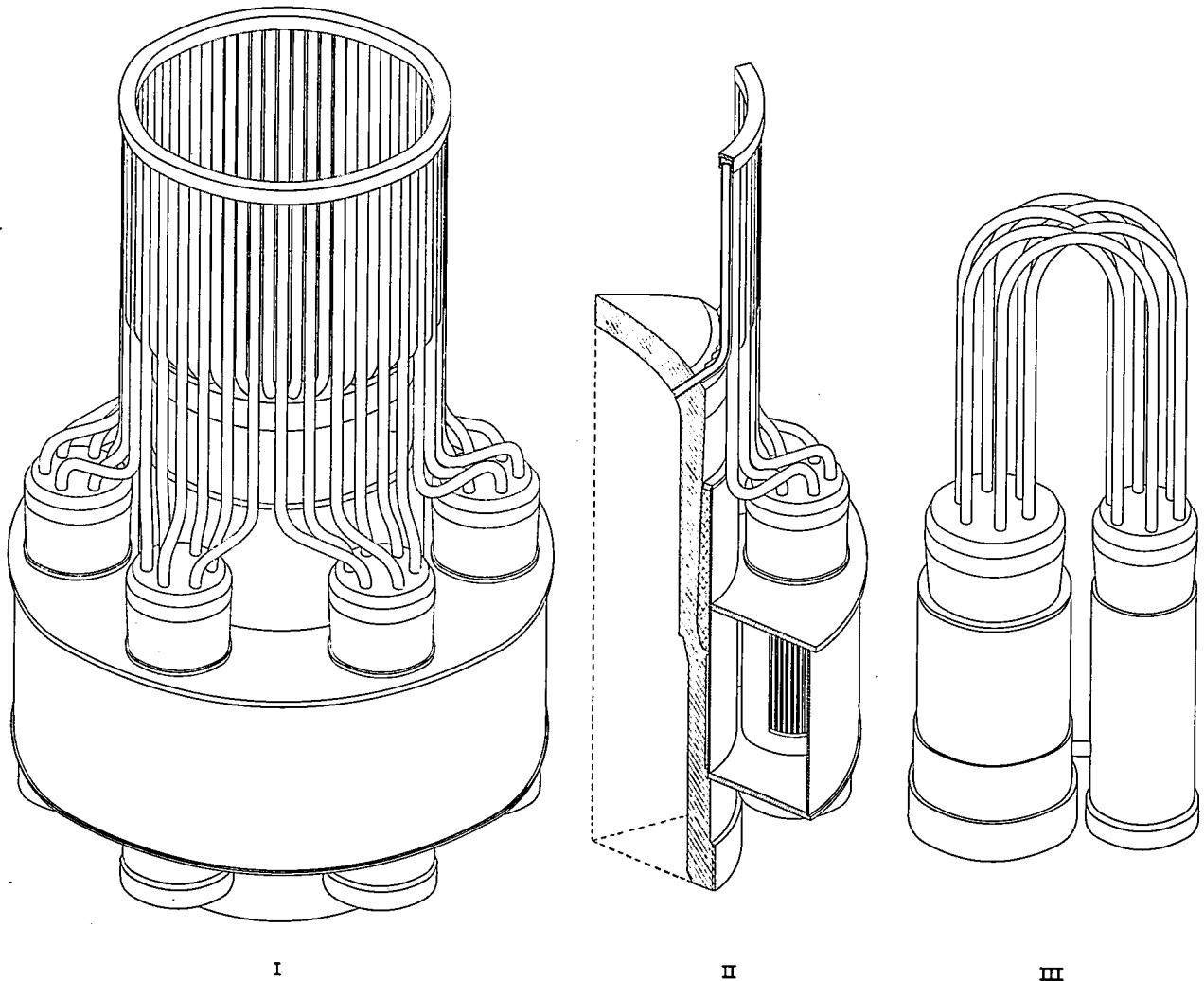
The principles of metrology based on the use of diffraction gratings are discussed and demonstrated. In particular, an optical photoelectric method for the scanning of gratings is described, which makes use of a highly corrected optical system with magnification -1 and which gives an output of two sinusoidal photoelectric signals with a stable phase difference of 90° . In order to obtain the possibility of AC coupling both signals have been given a time-harmonic phase modulation. With this method fast digital directional measurement of displacement in steps of one micron can be achieved. A working model of such a scanning device is shown.

In addition a method for the detection of extremely small translations or rotations, making use of one phase-modulated signal, is discussed and demonstrated.

C5. Experiments on high power hot gas engines,
by R. A. J. O. van Witteveen.

In 1816 the Scotsman Stirling invented the hot air engine and occasional applications of this engine have appeared throughout the 19th century. In recent years extended research has been performed in the Philips laboratories on the Stirling process and due to the use of modern materials, the application of new design possibilities and better knowledge of flow and heat-transfer phenomena the Stirling process can now be realized in a very efficient way.

After a short explanation of the principles of the Stirling process and the basic design of the hot gas engine some problems connected with the design of high power engines are discussed. In 1957 an overall efficiency of 38% was reached and with the help of an electronic computer work is now under way to optimize some twenty independent design parameters in the hope of reaching an efficiency of 40 to 45% for large engines.



- I* The top part of a 400 hp engine, showing circular symmetry; the crankcase has been omitted.
II A sector of one sixth of the engine has been cut out.
III The "part-engine" is obtained by transforming the sector-shaped cylinder-part to a cylinder with one sixth of the original swept volume and keeping the inside gas circuit the same. The phenomena in the circuit of the "part-engine" will be exactly the same as those prevailing in the original 400 hp engine.

Engines have been made so far with up to 80 hp per cylinder. However, the power output is by no means limited to this figure and for the near future plans exist for 400 hp per cylinder. Of course, the design of such an engine requires a far extrapolation of the present theory which may lead to unexpected disturbing effects; therefore, additional experimental work on an engine of this type is required. This, however, would be very expensive and time-consuming owing to the large size. Thus it is desirable to have a kind of model which is simpler, easier to modify and well fitted for the performance of exact measurements. It should be noted that a model with an identical thermodynamical behaviour is needed. Fortunately, the hot gas engine offers the possibility for such a model owing to its circular symmetry. This model, called the "part-engine", is explained and demonstrated (see illustration on the preceding page).

C6. Electronic precision versus tolerance of components, by J. J. Zaalberg van Zelst.

There is a discrepancy between the tolerance of components, especially of the active ones, and the precision often required in electronics. Nevertheless in many cases it is possible to achieve the required accuracy, even without using digital methods, by the application of stabilizing principles or by the use of the components in somewhat less obvious ways. A number of these procedures and their fundamental limits are briefly discussed and their possibilities are illustrated by some examples.

C7. Numerical processing of measured data, by N. F. Verster.

The processing of measured data with a computer saves tedious routine work, especially if the measurements are directly recorded on tape or cards. Special measures are necessary, however, to deal with such a blindfold operation. In this paper numerical methods are described which can be used to detect erroneous measurements and deviations from the expected behaviour, provided that the measurements are redundant.

C8. Measurements and waveguides in the millimetre wave range, by M. Gevers.

During the last five years a large number of microwave components and measuring instruments for the 4- and 2-mm wave bands have been developed in our laboratory. To test the feasibility of a new component it is usually designed first for use at 3-cm waves, for which the necessary measuring equipment is readily available and no stringent mechanical tolerances are required. However, for the measurement of the overall performance of the finally obtained mm-wave components direct measurements of their properties at the mm-wave frequencies are a necessity. The latest addition to the

facilities available for this purpose is an automatic swept-frequency Smith-chart plotter which covers the range between 72 and 77 Gc/s (see paper No. C9).

While conventional hollow metallic waveguides are widely used in the cm-wave range without any difficulty, these types of transmission lines for sub-mm waves are difficult to make with reasonable accuracy and have high attenuation. Other types of transmission lines, such as dielectric lines, image lines and quasi-optical transmission lines have been investigated at 4 mm wavelength in this laboratory with a view to application at sub-mm waves in due time. The use of these less conventional transmission lines is shown in a demonstration where some optical phenomena, which are well known in principle, can be displayed much simpler than with normal optical instruments (see paper No. C10).

C9. An automatic swept-frequency Smith-Chart plotter for 72-77 Gc/s, by F. C. de Ronde.

An impedance plotter has been made for a frequency automatically swept from 72 to 77 Gc/s which allows the instantaneous measurement of phenomena in the order of microseconds or even shorter. The well-known four-probe method is applied. Instead of electric coupling by probes, branch-guide coupling with quarter-wave transformers is used. In this way a rather simple, stable device could be made containing broadband detectors having identical frequency characteristics.

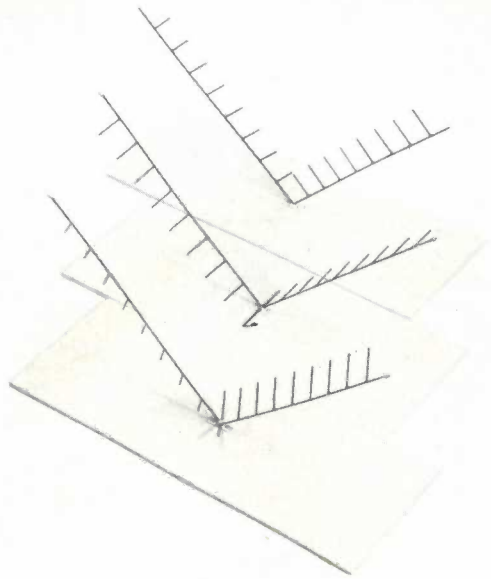
One broadband detector was used as a leveller in order to make the output of the sweep oscillator (C.S.F. carcinotron) as seen by the other detectors constant as a function of frequency.

Decoupling between the load and the sweep oscillator has been obtained by using a broadband 25 dB pad.

C10. Demonstrations of optical behaviour of 4-mm waves, by A. G. van Nie.

At a wavelength of 4.2 mm some optical phenomena are shown. By means of a horn and a lens a beam is formed. The phenomena demonstrated are:

- 1) The polarization direction of the beam can be changed by reflection at a metallic mirror. This change of the polarization direction is caused by a phase jump, equal to π radians, of the tangential component of the electric field on the surface (three cases are illustrated by the models in the photograph on next page).
- 2) In case of total reflection at the boundary between two dielectrics — one of which in our case is air — a wave, which decreases exponentially, exists behind the boundary. Disturbances in this wave, caused for instance by the application of some other dielectric, induce changes in the reflected wave.
- 3) For a plane wave incident on a lens the highest intensity does not occur in the focal plane when the beam is diffracted. In that case the highest intensity is found in a plane closer to the lens.



C11. Compatible single-sideband modulation, by T. J. van Kessel.

Single-sideband (SSB) modulation has a number of advantages over amplitude modulation (AM), which is double-sideband: Its bandwidth is only half that of AM, and for the same transmitter power it leads to a better signal-to-noise ratio. A disadvantage is that its envelope is distorted so that demodulation is considerably more complicated than in the case of AM, where a peak detector — the radio receiver — can be used.

It appears to be possible to combine the advantages of AM and SSB in one system that is compatible with normal broadcasting receivers. This type of modulation has been named compatible single-sideband (CSSB) modulation. The principle of the CSSB-system discussed in this paper is based on squaring of a full-carrier SSB signal and suppression of all the frequency components except those around the double-carrier frequency. This method implies that in order to convert an AM- to a CSSB-transmitter only a simple unit has to be added to the broadcasting side.

C12. Continuous delta modulation, by F. de Jager and J. A. Greefkes.

In delta modulation the signal to be transmitted is converted into a series of "1" and "0" pulses. Upon reception these pulses are applied to an integrating network and a low-pass filter, such that an approximation of the original input signal can be obtained. In order to achieve the correct sequence of the pulses at the transmitting end an auxiliary receiver is applied here and the difference between the input signal and the approximating signal determines whether the next pulse is "1" or "0". This implies a feedback loop in which the generated pulses are "quantized" in time and in amplitude.

All systems in which continuous information signals are coded into a sequence of binary pulses require the input signal to have the correct level. Otherwise either overloading occurs, or the distortion due to the process of quantization will be too large with respect to the signal amplitude. The object of "continuous delta modulation" is to obtain an automatic adaptation of the quantization steps to the level of the input signal.

To this end a voltage corresponding to the level of the signal is derived at the transmitting end, and its information is also transmitted to the receiver by means of the generated pulse series. In fact, the same circuit of delta modulation can be used for this purpose by introducing an extra modulator in the feedback loop. The system is adjusted so that for low-level input signals one third of the mean number of transmitted pulses is a "1" whereas for input signals of higher amplitude the ratio between the mean number of "1" and "0" pulses is gradually changed, in such a way that a 30 dB change in the level of the input signal corresponds to a change in the ratio of the mean numbers of the two kinds of pulses from $\frac{1}{3}$ to $\frac{2}{3}$. A signal, derived from this change in the generated pulse series, is used for controlling the size of the quantization steps in transmitter and receiver.

As this control mechanism is effected simultaneously in transmitter and receiver, the overall characteristic remains linear at any time. In this way the quantization noise is made nearly proportional to the amplitude of the information signal, which enables the process of coding to be applied to a much wider range of information signals. This is demonstrated by experiments, showing the difference between normal delta modulation and continuous delta modulation.

C13. Delta modulation for video transmission, by J. C. Balder.

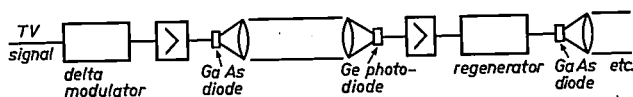
A method is described which enables video signals to be transmitted by the form of pulse-code modulation known as delta modulation. A balanced pair of tunnel-diodes is used for converting the video signal into a binary signal. With a system consisting of two tunnel-diodes, one coil and several resistors, operating at a bit rate of 100 Mc/s a ratio of signal to quantizing noise of 42 dB is obtained. An advantage over the normal P.C.M.-system is the simplicity of the apparatus involved. Also P.C.M. has to use rather high bit rates in order to avoid isophots. With the delta modulation system no isophots are visible because there are no fixed quantizing levels. With this system a colour picture can be transmitted which shows no degradation.

The system will be used to reduce the cumulative effect of interference in communication systems having many links, as e.g. optical systems or systems using circular wave guides in the TE_{01} mode.

C14. Infrared transmission of delta-modulated TV signals, by C. Kramer.

A system is described in which the output of a

tunnel-diode delta modulator (see paper No. C13) modulates the current through a GaAs light-emitting diode (see illustration). This diode is imaged by a pair of lenses on a Ge photo-diode which detects the infrared light emitted by the GaAs diode. The detected pulse signal is then regenerated and supplied to a second GaAs diode etc. In this



way the signal can be transmitted over a much larger distance than would be possible with a linear modulation system such as AM or FM. It is shown that with a delta-modulator clock frequency of 40 Mc/s, which assures a good picture quality, and a current of 100 mA through the GaAs diode, the distance between the repeaters can be made 500 m with a negligible deterioration of the signal. Because only one in 10^{20} pulses will be received incorrectly an almost unlimited number of repeaters can be used.

SECTION D, BIOCHEMISTRY AND PERCEPTION

D1. Investigations on an isolated enzyme, by L. A. Æ. Sluyterman.

Enzymes, as true catalysts, accelerate not only some forward reaction but also the corresponding backward reaction. Which reaction prevails, depends upon the thermodynamic conditions. Therefore, the product of the forward reaction will be the substrate of the backward reaction. Thus, there should be certain similarities in the enzyme-substrate and the enzyme-product interaction. This principle was tested with the protein-splitting enzyme papain using a simple synthetic substrate and its corresponding product. Kinetic experiments showed five points of similarity, which are discussed in the paper. It is pointed out how the results may be used for more precise investigations on the enzyme-substrate interaction.

D2. Radiation inactivation and genetic transformation in bacteria, by J. H. Stuy.

Ultraviolet and ionizing radiation are lethal to bacteria, i.e. they inhibit their indefinite multiplication. There are strong indications that the vital cell structure that is damaged by U.V. is deoxyribonucleic acid (DNA), the bearer of the genetic information of the cells. In order to obtain further evidence for the destruction of DNA as the cause of bacterial death by radiation, and to investigate the nature of the damage to the DNA molecule, the phenomenon of bacterial transformation has been used. This means that genetic characteristics can be transferred from one kind of cell to another by extracted DNA. The transforming properties of DNA, extracted from irradiated cells, were studied and compared with DNA which has been damaged in vitro e.g. by shear, ultrasound or the enzyme DNA-ase. In this way some conclusions could be drawn regarding the relative role of chain breakage and chemical degradation of the bases in DNA inactivation induced by U.V. or ionizing radiation.

D3. The chloroplast as an enzyme system, by J. S. C. Wessels.

The determination of the biochemical function

of specific cytoplasmic structures provides a meeting ground for biochemistry and physiology and, through their joint contributions, a better insight into the nature of cellular organization and metabolism. Chloroplasts are found in green plants and are indispensable to photosynthesis. They vary in size but are commonly discs or flat ellipsoids, 5 micron in diameter on the average. Chloroplasts contain all the chlorophyll and accessory pigments required for the absorption of light and its conversion into chemical energy. Their role in the primary photo-chemical reactions of photosynthesis is beyond dispute.

In order to shed some light on the still open problems concerning the role of chloroplasts in the cellular metabolism, experiments have been performed in this laboratory on the mechanism of photoreduction of biological oxidation-reduction systems and of photophosphorylation (synthesis of adenosine triphosphate). The capacity of chloroplasts for carrying out these reactions outside the living cell supports the conclusion that these structures are autonomous cytoplasmic bodies in which are localized all the enzyme systems needed for photosynthesis.

D4. Introduction to the visit to the Instituut voor Perceptie Onderzoek (Institute for Perception Research), by J. F. Schouten.

The Institute is a joint undertaking of the Philips Research Laboratories and the Technological University at Eindhoven. It was founded in 1957. Its scope is formed by the human "informative" cycle consisting of the senses, human perception and perceptual skills.

The main subjects under investigation are:

- Vision: - subjective brightness in relation to retinal adaptation,
- visual delay,
- pupillary contraction,
- subjective stroboscopy.
- Hearing: - subjective sound analysis,
- perception of pitch (harmonic and quasi-harmonic sounds).
- Phonetics: - relevant physical parameters of phonemes and morphemes,
- subjective analysis,
- speech synthesis.

Perception and perceptual skills:

- recognition,
- reaction times,
- industrial tasks like pinmounting etc.

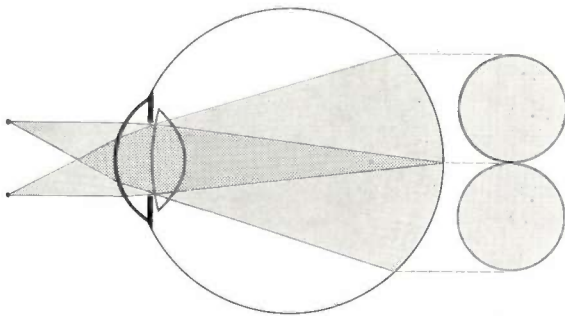
For these psychophysical measurements considerable attention is given to adequate instrumentation and in particular to automation of tedious data-reading and data-processing activities. As an example the DONDERS, a 20-channel reaction-time and reaction-quality recorder, may be mentioned.

The laboratory is also engaged in help for the disabled. Several aids to the deaf, the blind and the spastic are under development.

The activities of the laboratory are closely related to information theory, cybernetics and ergonomics (human engineering).

D5. Measurement of the efficiency spectrum of the pupillary reaction by an entoptical method, by H. Bouma.

When having to perform measurements of pupil size, one may choose between several objective methods. In a restricted area of application, however, the subjective entoptical projection method offers some definite advantages. It makes use of a twofold projection of the pupil edge (see figure) on the retina and provides a simple, rapid and accurate determination of static or semi-static pupil size by the subject himself. In addition, spontaneous pupillary movements may be detected easily.



As the subject performs these measurements himself, the experimenter's task is restricted to the necessary readings. This laborious task may well be taken over by automatic equipment. In our apparatus the settings of several continuous and discontinuous variables are registered automatically by pressing a single button. In general such an automatic procedure greatly facilitates the collection of a great amount of psychophysical data.

The entoptical method has been used to study the influence on pupil size of several parameters of retinal illumination. In particular the effect of wavelength is considered here. Generally the efficiency spectrum of a reaction to light reflects the absorption spectrum of the retinal pigments involved. It is well known that under bright lighting conditions (photopic or cone vision) the maximum

efficiency of the human eye is located at 555 nm. Under dim lighting conditions (scotopic or rod vision) this maximum lies considerably lower, at 505 nm.

The static pupillary contraction, even under photopic conditions, is found to follow roughly the scotopic efficiency curve. Thus under photopic conditions brightness and colour are determined by the cone mechanism. The static pupillary contraction, however, is determined by the rod mechanism.

D6. Short-term characteristics of pitch perception, by B. Lopes Cardozo.

In a sinusoid of sufficiently long duration the ear can distinguish the tonal part and the clicks arising from transients at the beginning and at the end of the signal. As the sinusoid is made shorter the tonal part wanes and the clicks become prevalent at the expense of the tonal quality.

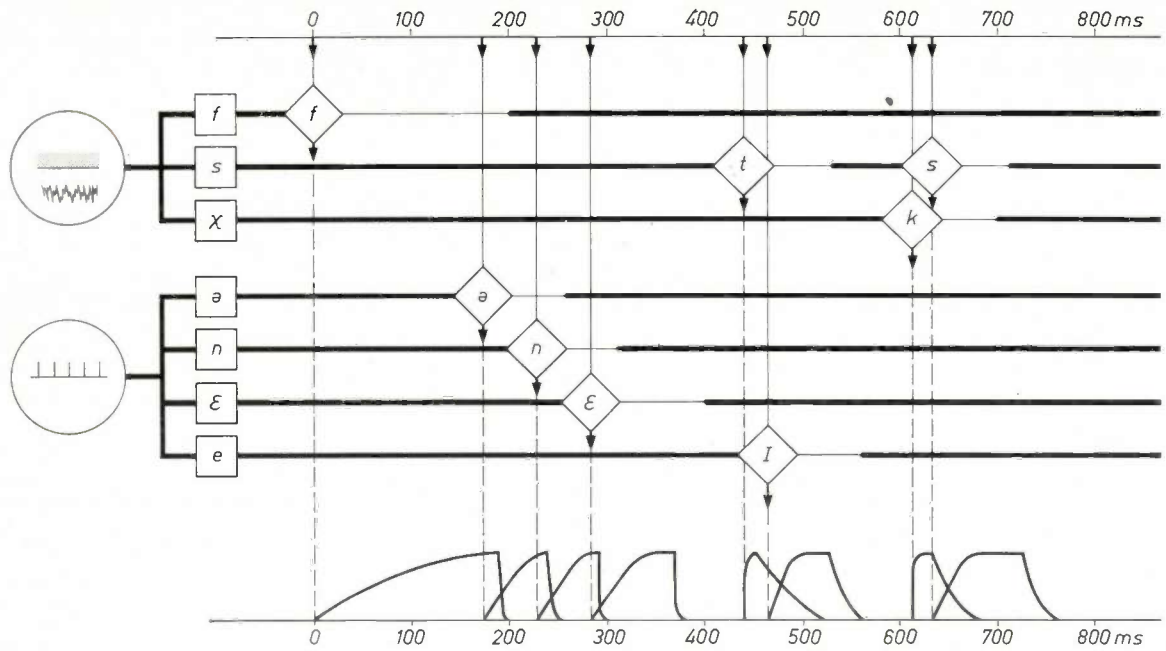
It is possible to obtain quantitative data about this phenomenon by measuring the just noticeable frequency difference between two pure tones sounding one after the other. Experiments for this measurement are shown and discussed. As an outcome of the experiments it was found that the just noticeable frequency difference is roughly inversely proportional to the duration of the signal. This relation holds for signals shorter than about 60 milliseconds and does not seem to depend very much on the frequency of the signals. The possible implications of this relation are discussed in terms of the signal-to-noise ratio of a hypothetical model of the human ear.

D7. Speech Perception, by A. Cohen.

In phonetic analysis, apart from orthodox instruments such as the oscillograph and spectrograph, use is made of an electronic gating circuit for analysing purposes.

This device enables the investigator to run perceptual tests by gating out any desired portion from a speech message that can be used as a stimulus. In this way a means is provided for analytical listening. It turns out to be possible to decide on perceptual criteria where one speech segment ends and the next one begins, resulting in an inventory of perceptual segments corresponding roughly to the number of phonemes of the language investigated (Dutch). The same technique is applied for intonation (pitch) measurements. To check the hypotheses inherent in the analytic technique a system for generating synthetic speech has been devised, the IPOVOX. The main feature of this instrument is that it generates words as a series of discrete segments, providing accurate means of shaping the amplitude envelopes of individual segments, which prove to convey as much information to the listener as the frequencies involved, known as formants.

The figure gives a block diagram of synthesizing



the word *phonetics* as a series of discrete segments. The speech sounds *f*, *s*, *t*, and *k* are obtained by filtering the output of a noise generator. A second source, generating periodic impulses, is used for the remaining speech sounds. These vowel-like sounds are obtained by a process of selective filtering involving two frequency regions (formants) per vowel. The required amplitude envelopes of individual segments are given by electronic gating circuits. In principle the falling-off of one segment is a cue for the onset of the next segment.

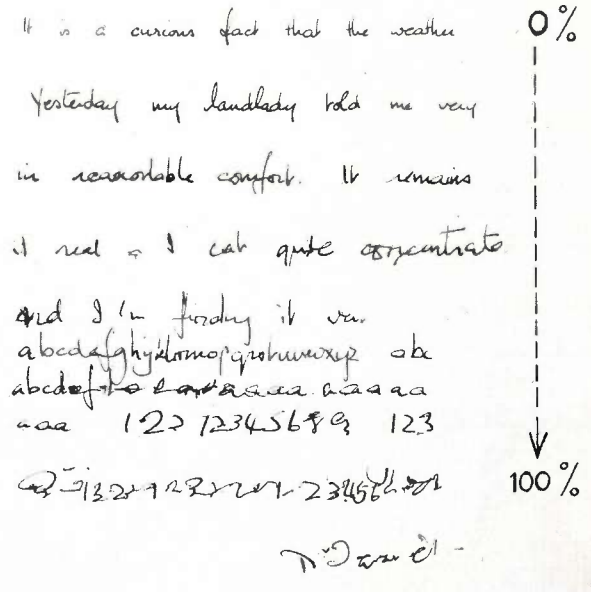
formed. With increasing rate of the standard task the quality of the handwriting strongly deteriorates, the content becomes more pedestrian and errors crop up. Under severe strain the handwriting assumes the character of that of mentally disturbed persons.

D8. Measurement of perceptual load, by H. W. Horeman and W. G. Koster.

In modern industry human labour cannot be adequately expressed in terms of physical energy only, because of the almost negligible low energy consumption of the human actions. Analysis shows that human labour exerts a load on the perceptual and mental faculties of the human being, rather than constituting a physical load. To measure a perceptual load quantitatively, a method has been developed in which the subject is put in a situation of working on two tasks simultaneously. By means of this method the amount of attention a subject can pay to a specific task while having his attention capacity partly blocked by simultaneous presentation of a standard task. From the attention required to perform the specific task a measure of the perceptual load can be gauged.

Results of work carried out in these dual-task situations are presented.

The figure shows a specimen of a subject's manuscript written extempore, used as a task to be performed simultaneously with a standard task, viz, discriminating between tones of high and low pitch. Each line is a sample representing different rates at which the standard task was per-



Brief mention is made of results of the dual-task method applied to measure the influence of lighting conditions on task performance as well as to establish the perceptual load constituted by driving a car. In the latter application use has been made of a simple rhythmical standard task, where the magnitude of the fluctuations in the rhythm serves as a measure of the load. The rationale of this choice is to have the smallest amount of interference with the performance of the specific task.