

THE BELL SYSTEM TECHNICAL JOURNAL

VOLUME XXXVIII

MAY 1959

NUMBER 3

Copyright 1959, American Telephone and Telegraph Company

Probability of Error for Optimal Codes in a Gaussian Channel

By CLAUDE E. SHANNON

(Manuscript received October 17, 1958)

A study is made of coding and decoding systems for a continuous channel with an additive gaussian noise and subject to an average power limitation at the transmitter. Upper and lower bounds are found for the error probability in decoding with optimal codes and decoding systems. These bounds are close together for signaling rates near channel capacity and also for signaling rates near zero, but diverge between. Curves exhibiting these bounds are given.

I. INTRODUCTION

Consider a communication channel of the following type: Once each second a real number may be chosen at the transmitting point. This number is transmitted to the receiving point but is perturbed by an additive gaussian noise, so that the i th real number, s_i , is received as $s_i + x_i$. The x_i are assumed independent gaussian random variables all with the same variance N .

A code word of length n for such a channel is a sequence of n real numbers (s_1, s_2, \dots, s_n) . This may be thought of geometrically as a point in n -dimensional Euclidean space. The effect of noise is then to move this point to a nearby point according to a spherical gaussian distribution.

A block code of length n with M words is a mapping of the integers 1, 2, \dots , M into a set of M code words w_1, w_2, \dots, w_M (not necessarily

all distinct). Thus, geometrically, a block code consists of a collection of M (or less) points with associated integers. It may be thought of as a way of transmitting an integer from 1 to M to the receiving point (by sending the corresponding code word). A *decoding system* for such a code is a partitioning of the n -dimensional space into M subsets corresponding to the integers from 1 to M . This is a way of deciding, at the receiving point, on the transmitted integer. If the received signal is in subset S_i , the transmitted message is taken to be integer i .

We shall assume throughout that all integers from 1 to M occur as messages with equal probability $1/M$. There is, then, for a given code and decoding system, a definite probability of error for transmitting a message. This is given by

$$P_e = \frac{1}{M} \sum_{i=1}^M P_{ei},$$

where P_{ei} is the probability, if code word w_i is sent, that it will be decoded as an integer other than i . P_{ei} is, of course, the total probability under the gaussian distribution, centered on w_i in the region complementary to S_i .

An *optimal decoding system* for a code is one which minimizes the probability of error for the code. Since the gaussian density is monotone decreasing with distance, an optimal decoding system for a given code is one which decodes any received signal as the integer corresponding to the geometrically nearest code word. If there are several code words at the same minimal distance, any of these may be used without affecting the probability of error. A decoding system of this sort is called *minimum distance decoding* or *maximum likelihood decoding*. It results in a partitioning of the n -dimensional space into n -dimensional polyhedra, or polytopes, around the different signal points, each polyhedron bounded by a finite number (not more than $M - 1$) of $(n - 1)$ -dimensional hyperplanes.

We are interested in the problem of finding good codes, that is, placing M points in such a way as to minimize the probability of error P_e . If there were no conditions on the code words, it is evident that the probability of error could be made as small as desired for any M , n and N by placing the code words at sufficiently widely separated points in the n space. In normal applications, however, there will be limitations on the choice of code words that prevent this type of solution. An interesting case that has been considered in the past is that of placing some kind of *average power limitation* on the code words; the distance of the points from the origin should not be too great. We may define three different possible limitations of this sort:

i. All code words are required to have *exactly the same power* P or the same distance from the origin. Thus, we are required to choose for code words points lying on the surface of a sphere of radius \sqrt{nP} .

ii. All code words have power P or less. Here all code words are required to lie interior to or on the surface of a sphere of radius \sqrt{nP} .

iii. The *average power* of all code words is P or less. Here, individual code words may have a greater squared distance than nP but the average of the set of squared distances cannot exceed nP .

These three cases lead to quite similar results, as we shall see. The first condition is simpler and leads to somewhat sharper conclusions — we shall first analyze this case and use these results for the other two conditions. *Therefore, until the contrary is stated, we assume all code words to lie on the sphere of radius \sqrt{nP} .*

Our first problem is to estimate, as well as possible, the probability of error $P_e(M, n, \sqrt{P/N})$ for the best code of length n containing M words each of power P and perturbed by noise of variance N . This minimal or *optimal probability of error* we denote by $P_{e\text{opt}}(M, n, \sqrt{P/N})$. It is clear that, for fixed $M, n, P_{e\text{opt}}$ will be a function only of the quotient $A = \sqrt{P/N}$ by change of scale in the geometrical picture. We shall obtain upper and lower bounds on $P_{e\text{opt}}$ of several different types. Over an important range of values these bounds are reasonably close together, giving good estimates of $P_{e\text{opt}}$. Some calculated values and curves are given and the bounds are used to develop other bounds for the second and third type conditions on the code words.

The geometrical approach we use is akin to that previously used by the author¹ but carried here to a numerical conclusion. The problem is also close to that studied by Rice,² who obtained an estimate similar to but not as sharp as one of our upper bounds. The work here is also analogous to bounds given by Elias³ for the binary symmetric and binary erasure channels, and related to bounds for the general discrete memoryless channel given by the author.⁴

In a general way, our bounds, both upper and lower, vary exponentially with n for a fixed signaling rate, R , and fixed P/N . In fact, they all can be put [letting $R = (1/n) \log M$, so that R is the transmitting rate for the code] in the form

$$e^{-E(R)n+o(n)}, \quad (1)$$

where $E(R)$ is a suitable function of R (and of P/N , which we think of as a fixed parameter). [In (1), $o(n)$ is a term of order less than n ; as $n \rightarrow \infty$ it becomes small relative to $E(R)n$.]

Thus, for large n , the logarithm of the bound increases linearly with n or, more precisely, the ratio of this logarithm to n approaches a con-

stant $E(R)$. This quantity $E(R)$ gives a crude measure of how rapidly the probability of error approaches zero. We will call this type of quantity a *reliability*. More precisely, we may define the reliability for a channel as follows:

$$E(R) = \limsup_{n \rightarrow \infty} -\frac{1}{n} \log P_{e, \text{opt}}(R, n), \quad (2)$$

where $P_{e, \text{opt}}(R, n)$ is the optimal probability of error for codes of rate R and length n . We will find that our bounds determine $E(R)$ *exactly* over an important range of rates, from a certain critical rate R_c up to channel capacity. Between zero and R_c , E is not exactly determined by our bounds, but lies within a not too wide range.

In connection with the reliability E , it may be noted that, in (1) above, knowledge of $E(R)$ and n does not closely determine the probability of error, even when n is large; the term $o(n)$ can cause a large and, in fact, increasing multiplier. On the other hand, given a desired probability of error and $E(R)$, the necessary value of the code length n *will* be sharply determined when n is large; in fact, n will be asymptotic to $-(1/E) \log P_e$. This inverse problem is perhaps the more natural one in applications: given a required level of probability of error, how long must the code be?

The type of channel we are studying here is, of course, closely related to a band-limited channel (W cycles per second wide) perturbed by white gaussian noise. In a sense, such a band-limited channel can be thought of as having $2W$ coordinates per second, each independently perturbed by a gaussian variable. However, such an identification must be treated with care, since to control these degrees of freedom physically and stay strictly within the bandwidth would require an infinite delay.

It is possible to stay very closely within a bandwidth W with a large but finite delay T , for example, by using $(\sin x)/x$ pulses with one tail deleted T from the maximum point. This deletion causes a spill-over outside the band of not more than the energy of the deleted part, an amount less than $1/T$ for the unit $(\sin x)/x$ case. By making T large, we can approach the situation of staying within the allotted bandwidth and also, for example, approach zero probability of error at signaling rates close to channel capacity.

However, for the problems we are studying here, delay as related to probability of error is of fundamental importance and, in applications of our results to such band-limited channels, the additional delay involved in staying closely within the allotted channel must be remembered. This is the reason for defining the channel as we have above.

II. SUMMARY

In this section we summarize briefly the main results obtained in the paper, both for easy reference and for readers who may be interested in the results without wishing to work through the detailed analysis. It might be said that the algebra involved is in several places unusually tedious.

We use the following notations:

P = signal power (each code word is on the surface of a sphere of radius \sqrt{nP});

N = noise power (variance N in each dimension);

$A = \sqrt{P/N}$ = signal-to-noise "amplitude" ratio;

n = number of dimensions or block length of code;

M = number of code words;

$R = (1/n) \log M$ = signaling rate for a code (natural units);

$C = \frac{1}{2} \log (P + N)/N = \frac{1}{2} \log (A^2 + 1)$ = channel capacity (per degree of freedom);

θ = variable for half-angle of cones appearing in the geometrical problem which follows;

$\Omega(\theta)$ = solid angle in n space of a cone of half-angle θ , or area of unit n sphere cut out by the cone;

$\theta_0 = \cot^{-1} A$ = cone angle relating to channel capacity;

θ_1 = cone angle such that the solid angle $\Omega(\theta_1)$ of this cone is $(1/M)\Omega(\pi)$, [the solid angle of a sphere is $\Omega(\pi)$]; thus, θ_1 is a cone angle related to the rate R ;

$G = G(\theta) = \frac{1}{2}(A \cos \theta + \sqrt{A^2 \cos^2 \theta + 4})$, a quantity which appears often in the formulas;

θ_c = the solution of $2 \cos \theta_c - AG(\theta_c) \sin^2 \theta_c = 0$ (this critical angle is important in that the nature of the bounds change according as $\theta_1 > \theta_c$ or $\theta_1 < \theta_c$);

$Q(\theta) = Q(\theta, A, n)$ = probability of a point X in n space, at distance $A\sqrt{n}$ from the origin, being moved outside a circular cone of half-angle θ with vertex at the origin O and axis OX (the perturbation is assumed spherical gaussian with unit variance in all dimensions);

$E_L(\theta) = A^2/2 - \frac{1}{2}AG \cos \theta - \log (G \sin \theta)$, an exponent appearing in our bounds;

$P_{e \text{ opt}}(n, R, A)$ = Probability of error for the best code of length n , signal-to-noise ratio A and rate R ;

$\Phi(X)$ = normal distribution with zero mean and unit variance.

The results of the paper will now be summarized. $P_{e \text{ opt}}$ can be bounded as follows:

$$Q(\theta_1) \leq P_{e \text{ opt}} \leq Q(\theta_1) - \int_0^{\theta_1} \frac{\Omega(\theta)}{\Omega(\theta_1)} dQ(\theta). \quad (3)$$

[Here $dQ(\theta)$ is negative, so the right additional term is positive.] These bounds can be written in terms of rather complex integrals. To obtain more insight into their behavior, we obtain, in the first place, asymptotic expressions for these bounds when n is large and, in the second place, cruder bounds which, however, are expressed in terms of elementary functions without integrals.

The asymptotic lower bound is (asymptotically correct as $n \rightarrow \infty$)

$$\begin{aligned} Q(\theta_1) &\sim \frac{1}{\sqrt{n\pi} G \sqrt{1 + G^2 \sin^2 \theta_1} (\cos \theta_1 - AG \sin^2 \theta_1)} e^{-E_L(\theta_1)n} \\ &= \frac{\alpha(\theta_1)}{\sqrt{n}} e^{-E_L(\theta_1)n} \quad (\theta_1 > \theta_0). \end{aligned} \quad (4)$$

The asymptotic upper bound is

$$Q(\theta_1) - \int_0^{\theta_1} \frac{\Omega(\theta)}{\Omega(\theta_1)} dQ(\theta) \sim \frac{\alpha(\theta_1)}{\sqrt{n}} e^{-E_L(\theta_1)n} \left(1 - \frac{\cos \theta_1 - AG \sin^2 \theta_1}{2 \cos \theta_1 - AG \sin^2 \theta_1} \right). \quad (5)$$

This formula is valid for $\theta_0 < \theta_1 < \theta_c$. In this range the upper and lower asymptotic bounds differ only by the factor in parentheses independent of n . Thus, asymptotically, the probability of error is determined by these relations to within a multiplying factor depending on the rate. For rates near channel capacity (θ_1 near θ_0) the factor is just a little over unity; the bounds are close together. For lower rates near R_c (corresponding to θ_c), the factor becomes large. For $\theta_1 > \theta_c$ the upper bound asymptote is

$$\frac{1}{\cos \theta_c \sin^3 \theta_c G(\theta_c) \sqrt{\pi E''(\theta_c) [1 + G(\theta_c)]^2}} e^{-n[E_L(\theta_c) - R]}. \quad (6)$$

In addition to the asymptotic bound, we also obtain firm bounds, valid for all n , but poorer than the asymptotic bounds when n is large. The firm lower bound is

$$P_e \geq \frac{1}{6} \frac{\sqrt{n-1} e^{3/2}}{n(A+1)^2 e^{(A+1)^2/2}} e^{-E_L(\theta_1)n}. \quad (7)$$

It may be seen that this is equal to the asymptotic bound multiplied by a factor essentially independent of n . The firm upper bound {valid if the maximum of $G^n (\sin \theta)^{2n-3} \exp [-(n/2)(A^2 - AG \cos \theta)]$ in the range 0 to θ_1 occurs at θ_1 } is

$$P_{e \text{ opt}} \leq \theta_1 \sqrt{2n} e^{3/2} G^n(\theta_1) \sin \theta_1^{n-2} \exp \left[\frac{n}{2} (-A^2 + AG \cos \theta_1) \right] \cdot \left\{ 1 + \frac{1}{n\theta_1 \min [A, AG(\theta_1) \sin \theta_1 - \cot \theta_1]} \right\}. \tag{8}$$

For rates near channel capacity, the upper and lower asymptotic bounds are both approximately the same, giving, where n is large and $C - R$ small (but positive):

$$P_{e \text{ opt}} \doteq \Phi \left[\sqrt{n} \sqrt{\frac{2P(P+N)}{N(P+2N)}} (R - C) \right], \tag{9}$$

where Φ is the normal distribution with unit variance.

To relate the angle θ_1 in the above formulas to the rate R , inequalities are found:

$$\frac{\Gamma\left(\frac{n}{2} + 1\right) (\sin \theta_1)^{n-1}}{n \Gamma\left(\frac{n+1}{2}\right) \pi^{1/2} \cos \theta_1} \left(1 - \frac{1}{n} \tan^2 \theta_1\right) \leq e^{-nR} \tag{10}$$

$$\leq \frac{\Gamma\left(\frac{n}{2} + 1\right) (\sin \theta_1)^{n-1}}{n \Gamma\left(\frac{(n+1)}{2}\right) \pi^{1/2} \cos \theta_1}.$$

Asymptotically, it follows that:

$$e^{-nR} \sim \frac{\sin^n \theta_1}{\sqrt{2\pi n} \sin \theta_1 \cos \theta_1}. \tag{11}$$

For low rates (particularly $R < R_c$), the above bounds diverge and give less information. Two different arguments lead to other bounds useful at low rates. The *low rate upper bound* is:

$$P_{e \text{ opt}} \leq \frac{1}{\lambda A \sqrt{\pi n}} e^{n[R - (\lambda^2 A^2)/4]}, \tag{12}$$

where λ satisfies $R = [1 - (1/n)] \log (\sin 2 \sin^{-1} \lambda / \sqrt{2})$. Note that

as $R \rightarrow 0$, $\lambda \rightarrow 1$ and the upper bound is approximately

$$\frac{1}{A\sqrt{\pi n}} e^{-nA^2/4}.$$

The *low rate lower bound* may be written

$$P_{e \text{ opt}} \cong \frac{1}{2} \Phi \left[-A \left(\frac{2M}{2M-1} \frac{n}{2} \right)^{1/2} \right]. \quad (13)$$

For M large, this bound is close to $\frac{1}{2}\Phi(-A\sqrt{n/2})$ and, if n is large, this is asymptotic to $1/(A\sqrt{\pi n}) e^{-nA^2/4}$. Thus, for rates close to zero and large n we again have a situation where the bounds are close together and give a sharp evaluation of $P_{e \text{ opt}}$.

With codes of rate $R \cong C + \epsilon$, where ϵ is fixed and positive, $P_{e \text{ opt}}$ approaches unity as the code length n increases.

III. THE LOWER BOUND BY THE "SPHERE-PACKING" ARGUMENT

Suppose we have a code with M points each at distance \sqrt{nP} from the origin in n space. Since any two words are at equal distance from the origin, the $n-1$ hyperplane which bisects the connecting line passes through the origin. Thus, all of the hyperplanes which determine the polyhedra surrounding these points (for the optimal decoding system) pass through the origin. These polyhedra, therefore, are pyramids with apexes at the origin. The probability of error for the code is

$$\frac{1}{M} \sum_{i=1}^M P_{ei},$$

where P_{ei} is the probability, if code word i is used, that it will be carried by the noise outside the pyramid around the i th word. The probability of being *correct* is

$$1 - \frac{1}{M} \sum_{i=1}^M P_{ei} = \frac{1}{M} \sum_{i=1}^M (1 - P_{ei});$$

that is, the average probability of a code word being moved to a point *within* its own pyramid.

Let the i th pyramid have a solid angle Ω_i (that is, Ω_i is the area cut out by the pyramid on the unit n -dimensional spherical surface). Consider, for comparison, a right circular n -dimensional cone with the same solid angle Ω_i and having a code word on its axis at distance \sqrt{nP} . We assert that *the probability of this comparison point being moved to within its cone is greater than that of w_i being moved to within its pyramid*. This

is because of the monotone decreasing probability density with distance from the code word. The pyramid can be deformed into the cone by moving small conical elements from far distances to nearer distances, this movement continually increasing probability. This is suggested for a three-dimensional case in Fig. 1. Moving small conical elements from outside the cone to inside it increases probability, since the probability density is greater inside the cone than outside. Formally, this follows by integrating the probability density over the region R_1 in the cone but not in the pyramid, and in the region R_2 in the pyramid but not in the cone. The first is greater than the solid angle Ω of R_1 times the density at the edge of the cone. The value for the pyramid is less than the same quantity.

We have, then, a bound on the probability of error P_e for a given code:

$$P_e \geq \frac{1}{M} \sum_{i=1}^M Q^*(\Omega_i), \tag{14}$$

where Ω_i is the solid angle for the i th pyramid, and $Q^*(\Omega)$ is the probability of a point being carried outside a surrounding cone of solid angle Ω . It is also true that

$$\sum_{i=1}^M \Omega_i = \Omega_0,$$

the solid angle of an n sphere, since the original pyramids corresponded to a partitioning of the sphere. Now, using again the property that the density decreases with distance, it follows that $Q^*(\Omega)$ is a convex function of Ω . Then we may further simplify this bound by replacing each Ω_i by

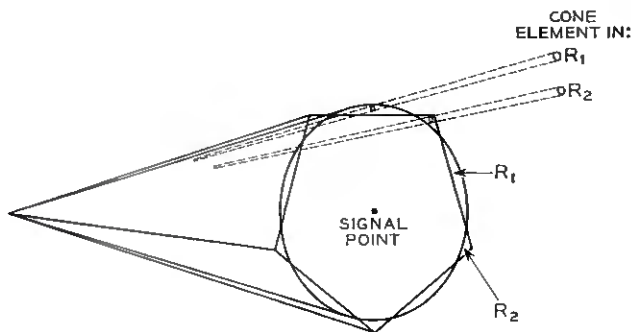


Fig. 1 — Pyramid deformed into cone by moving small conical elements from far to nearer distances.

the average Ω_0/M . In fact,

$$\frac{1}{M} \sum_{i=1}^M Q^*(\Omega_i) \geq Q^*\left(\frac{\Omega_0}{M}\right),$$

and hence

$$P_e \geq Q^*\left(\frac{\Omega_0}{M}\right).$$

It is more convenient to work in terms of the half-cone angle θ rather than solid angles Ω . We define $Q(\theta)$ to be the probability of being carried outside a cone of half-angle θ . Then, if θ_1 corresponds to the cone of solid angle Ω_0/M , the bound above may be written

$$P_e \geq Q(\theta_1). \quad (15)$$

This is our fundamental lower bound for P_e . It still needs translation into terms of P , N , M and n , and estimation in terms of simple functions.

It may be noted that this bound is exactly the probability of error that would occur if it were possible to subdivide the space into M congruent cones, one for each code word, and place the code words on the axes of these cones. It is, of course, very plausible intuitively that any actual code would have a higher probability of error than would that with such a conical partitioning. Such a partitioning clearly is possible only for $n = 1$ or 2, if $M > 2$.

The lower bound $Q(\theta_1)$ can be evaluated in terms of a distribution familiar to statisticians as the noncentral t -distribution.⁵ The noncentral t may be thought of as the probability that the ratio of a random variable $(z + \delta)$ to the root mean square of f other random variables

$$\sqrt{\frac{1}{f} \sum x_i^2}$$

does not exceed t , where all variates x_i and z are gaussian and independent with mean zero and unit variance and δ is a constant. Thus, denoting it by $P(f, \delta, t)$, we have

$$P(f, \delta, t) = \Pr \left\{ \frac{z + \delta}{\sqrt{\frac{1}{f} \sum_1^f x_i^2}} \leq t \right\}. \quad (16)$$

In terms of our geometrical picture, this amounts to a spherical gaussian distribution with unit variance about a point δ from the origin in $f + 1$ space. The probability $P(f, \delta, t)$ is the probability of being outside a

cone from the origin having the line segment to the center of the distribution as axis. The cotangent of the half-cone angle θ is t/\sqrt{j} . Thus the probability $Q(\theta)$ is seen to be given by

$$Q(\theta) = P\left(n-1, \sqrt{\frac{nP}{N}}, \sqrt{n-1} \cot \theta\right). \quad (17)$$

The noncentral t -distribution does not appear to have been very extensively tabled. Johnson and Welch⁵ give some tables, but they are aimed at other types of application and are inconvenient for the purpose at hand. Further, they do not go to large values of n . We therefore will estimate this lower bound by developing an asymptotic formula for the cumulative distribution $Q(\theta)$ and also the density distribution $dQ/d\theta$. First, however, we will find an *upper* bound on $P_{e\text{opt}}$ in terms of the same distribution $Q(\theta)$.

IV. UPPER BOUND BY A RANDOM CODE METHOD

The upper bound for $P_{e\text{opt}}$ will be found by using an argument based on random codes. Consider the ensemble of codes obtained by placing M points randomly on the surface of a sphere of radius \sqrt{nP} . More precisely, each point is placed independently of all others with probability measure proportional to surface area or, equivalently, to solid angle. Each of the codes in the ensemble is to be decoded by the minimum distance process. We wish to compute the *average probability of error* for this *ensemble of codes*.

Because of the symmetry of the code points, the probability of error averaged over the ensemble will be equal to M times the average probability of error due to any particular code point, for example, code point 1. This may be computed as follows. The probability of message number 1 being transmitted is $1/M$. The differential probability that it will be displaced by the noise into the region between a cone of half-angle θ and one of half-angle $\theta + d\theta$ (these cones having vertex at the origin and axis out to code word 1) is $-dQ(\theta)$. [Recall that $Q(\theta)$ was defined as the probability that noise would carry a point outside the cone of angle θ with axis through the signal point.] Now consider the cone of half-angle θ surrounding such a *received point* (not the cone about the message point just described). If this cone is empty of signal points, the received word will be decoded correctly as message 1. If it is not empty, other points will be nearer and the received signal will be incorrectly decoded. (The probability of two or more points at exactly the same distance is readily seen to be zero and may be ignored.)

The probability in the ensemble of codes of the cone of half-angle θ being empty is easily calculated. The probability that any particular code word, say code word 2 or code word 3, etc. is in the cone is given by $\Omega(\theta)/\Omega(\pi)$, the ratio of the solid angle in the cone to the total solid angle. The probability a particular word is *not* in the cone is $1 - \Omega(\theta)/\Omega(\pi)$. The probability that all $M - 1$ other words are not in the cone is $[1 - \Omega(\theta)/\Omega(\pi)]^{M-1}$ since these are, in the ensemble of codes, placed independently. The probability of error, then, contributed by situations where the point 1 is displaced by an angle from θ to $\theta + d\theta$ is given by $-(1/M)[1 - [1 - \Omega(\theta)/\Omega(\pi)]^{M-1}]dQ(\theta)$. The total average probability of error for all code words and all noise displacements is then given by

$$P_{er} = - \int_{\theta=0}^{\pi} \left\{ 1 - \left[1 - \frac{\Omega(\theta)}{\Omega(\pi)} \right]^{M-1} \right\} dQ(\theta). \quad (18)$$

This is an exact formula for the average probability of error P_{er} for our random ensemble of codes. Since this is an average of P_e for particular codes, there must exist particular codes in the ensemble with at least this good a probability of error, and certainly then $P_{e\text{opt}} \leq P_{er}$.

We may weaken this bound slightly but obtain a simpler formula for calculation as follows. Note first that $\{1 - [\Omega(\theta)/\Omega(\pi)]^{M-1}\} \leq 1$ and also, using the well-known inequality $(1 - x)^n \geq 1 - nx$, we have $\{1 - [1 - \Omega(\theta)/\Omega(\pi)]^{M-1}\} \leq (M - 1)[\Omega(\theta)/\Omega(\pi)] \leq M[\Omega(\theta)/\Omega(\pi)]$. Now, break the integral into two parts, $0 \leq \theta \leq \theta_1$ and $\theta_1 \leq \theta \leq \pi$. In the first range, use the inequality just given and, in the second range, bound the expression in braces by 1. Thus,

$$\begin{aligned} P_{er} &\leq - \int_0^{\theta_1} M \left[\frac{\Omega(\theta)}{\Omega(\pi)} \right] dQ(\theta) - \int_{\theta_1}^{\pi} dQ(\theta), \\ P_{er} &\leq - \frac{M}{\Omega(\pi)} \int_0^{\theta_1} \Omega(\theta) dQ(\theta) + Q(\theta_1). \end{aligned} \quad (19)$$

It is convenient to choose for θ_1 the same value as appeared in the lower bound; that is, the θ_1 such that $\Omega(\theta_1)/\Omega(\pi) = 1/M$ — in other words, the θ_1 for which one expects one point within the θ_1 cone. The second term in (19) is then the same as the lower bound on $P_{e\text{opt}}$ obtained previously. In fact, collecting these results, we have

$$Q(\theta_1) \leq P_{e\text{opt}} \leq Q(\theta_1) - \frac{M}{\Omega(\pi)} \int_0^{\theta_1} \Omega(\theta) dQ(\theta), \quad (20)$$

where $M\Omega(\theta_1) = \Omega(\pi)$. These are our fundamental lower and upper bounds on $P_{e\text{opt}}$.

We now wish to evaluate and estimate $\Omega(\theta)$ and $Q(\theta)$.

V. FORMULAS FOR RATE R AS A FUNCTION OF THE CONE ANGLE θ

Our bounds on probability of error involve the code angle θ_1 such that the solid angle of the cone is $1/M = e^{-nR}$ times the full solid angle of a sphere. To relate these quantities more explicitly we calculate the solid angle of a cone in n dimensions with half-angle θ . In Fig. 2 this means calculating the $(n - 1)$ -dimensional area of the cap cut out by the cone on the unit sphere. This is obtained by summing the contributions due to ring-shaped elements of area (spherical surfaces in $n - 1$ dimensions

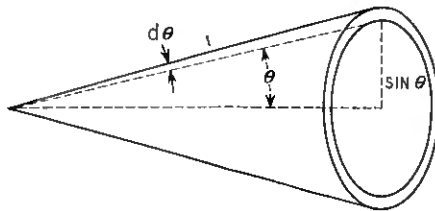


Fig. 2 — Cap cut out by the cone on the unit sphere.

of radius $\sin \theta$ and of incremental width $d\theta$). Thus, the total area of the cap is given by

$$\Omega(\theta_1) = \frac{(n - 1)\pi^{(n-1)/2}}{\Gamma\left(\frac{n + 1}{2}\right)} \int_0^{\theta_1} (\sin \theta)^{n-2} d\theta. \tag{21}$$

Here we used the formula for the surface $S_n(r)$ of a sphere of radius r in n dimensions, $S_n(r) = n\pi^{n/2}r^{n-1}/\Gamma(n/2 + 1)$.

To obtain simple inequalities and asymptotic expressions for $\Omega(\theta_1)$, make the change of variable in the integral $x = \sin \theta$, $d\theta = (1 - x^2)^{-1/2} dx$. Let $x_1 = \sin \theta_1$ and assume $\theta_1 < \pi/2$, so that $x_1 < 1$. Using the mean value theorem we obtain

$$(1 - x^2)^{-1/2} = (1 - x_1^2)^{-1/2} + \frac{\alpha}{(1 - \alpha^2)^{3/2}} (x - x_1), \tag{22}$$

where $0 \leq \alpha \leq x_1$. The term $\alpha(1 - \alpha^2)^{-3/2}$ must lie in the range from 0 to $x_1(1 - x_1^2)^{-3/2}$ since this is a monotone increasing function. Hence we have the inequalities

$$(1 - x_1^2)^{-1/2} + \frac{(x - x_1)x_1}{(1 - x_1^2)^{3/2}} \leq (1 - x^2)^{-1/2} \leq (1 - x_1^2)^{-1/2} \tag{23}$$

$$0 \leq x \leq x_1.$$

Note that $x - x_1$ is negative, so the correction term on the left is of the right sign. If we use these in the integral for $\Omega(\theta_1)$ we obtain

$$\frac{(n-1)\pi^{(n-1)/2}}{\Gamma\left(\frac{n+1}{2}\right)} \int_0^{x_1} x^{n-2} \left[(1-x_1^2)^{-1/2} + \frac{(x-x_1)x_1}{(1-x_1^2)^{3/2}} \right] dx$$

$$\cong \Omega(\theta_1) \cong \frac{(n-1)\pi^{(n-1)/2}}{\Gamma\left(\frac{n+1}{2}\right)} \int_0^{x_1} x^{n-2} \frac{dx}{\sqrt{1-x_1^2}}, \quad (24)$$

$$\frac{(n-1)\pi^{(n-1)/2}}{\Gamma\left(\frac{n+1}{2}\right) \sqrt{1-x_1^2}} \left[\frac{x_1^{n-1}}{n-1} + \frac{x_1^{n+1}}{n(1-x_1^2)} - \frac{x_1^{n+1}}{(n-1)(1-x_1^2)} \right]$$

$$\cong \Omega(\theta_1) \cong \frac{(n-1)\pi^{(n-1)/2} x_1^{n-1}}{\Gamma\left(\frac{n+1}{2}\right) (n-1) \sqrt{1-x_1^2}}, \quad (25)$$

$$\frac{\pi^{(n-1)/2} (\sin \theta_1)^{n-1}}{\Gamma\left(\frac{n+1}{2}\right) \cos \theta_1} \left(1 - \frac{1}{n} \tan^2 \theta_1 \right)$$

$$\cong \Omega(\theta_1) \cong \frac{\pi^{(n-1)/2} (\sin \theta_1)^{n-1}}{\Gamma\left(\frac{n+1}{2}\right) \cos \theta_1}. \quad (26)$$

Therefore, as $n \rightarrow \infty$, $\Omega(\theta_1)$ is asymptotic to the expression on the right.

The surface of the unit n sphere is $n\pi^{n/2}/\Gamma(n/2 + 1)$, hence,

$$\frac{\Gamma\left(\frac{n}{2} + 1\right) (\sin \theta_1)^{n-1}}{n\Gamma\left(\frac{n+1}{2}\right) \pi^{1/2} \cos \theta_1} \left(1 - \frac{1}{n} \tan^2 \theta_1 \right) \cong e^{-nR}$$

$$= \frac{\Omega(\theta_1)}{\Omega(\pi)} \cong \frac{\Gamma\left(\frac{n}{2} + 1\right) (\sin \theta_1)^{n-1}}{n\Gamma\left(\frac{n+1}{2}\right) \pi^{1/2} \cos \theta_1}. \quad (27)$$

Replacing the gamma functions by their asymptotic expressions, we obtain

$$e^{-nR} = \frac{\sin^n \theta_1}{\sqrt{2\pi n} \sin \theta_1 \cos \theta_1} \left[1 + O\left(\frac{1}{n}\right) \right]. \quad (28)$$

Thus $e^{-nR} \sim \sin^n \theta_1 / \sqrt{2\pi n} \sin \theta_1 \cos \theta_1$ and $e^{-R} \sim \sin \theta_1$. The somewhat sharper expression for e^{-nR} must be used when attempting asymptotic evaluations of P_e , since P_e is changed by a factor when θ_1 is changed by, for example, k/n . However, when only the reliability E is of interest, the simpler $R \sim -\log \sin \theta_1$ may be used.

VI. ASYMPTOTIC FORMULAS FOR $Q(\theta)$ AND $Q'(\theta)$

In Fig. 3, O is the origin, S is a signal point and the plane of the figure is a plane section in the n -dimensional space. The lines OA and OB represent a (circular) cone of angle θ about OS (that is, the intersection of this cone with the plane of the drawing.) The lines OA' and OB' correspond to a slightly larger cone of angle $\theta + d\theta$. We wish to estimate the probability $-dQ_n(\theta)$ of the signal point S being carried by noise into the region between these cones. From this, we will further calculate the probability $Q_n(\theta)$ of S being carried outside the θ cone. What is desired in both cases is an asymptotic estimate — a simple formula whose

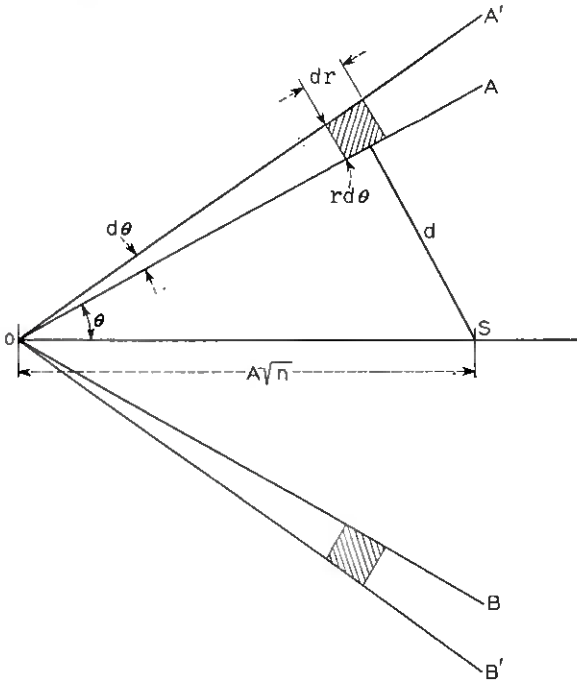


Fig. 3 — Plane of cone of half-angle θ .

ratio to the true value approaches 1 as n , the number of dimensions, increases.

The noise perturbs all coordinates normally and independently with variance 1. It produces a spherical gaussian distribution in the n -dimensional space. The probability density of its moving the signal point a distance d is given by

$$\frac{1}{(2\pi)^{n/2}} e^{-d^2/2} dV, \quad (29)$$

where dV is the element of volume. In Fig. 4 we wish to first calculate the probability density for the crosshatched ring-shaped region between

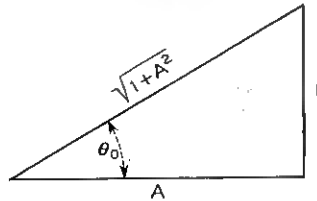


Fig. 4 — Special value θ_0 .

the two cones and between spheres about the origin of radius r and $r + dr$. The distance of this ring from the signal point is given by the cosine law as

$$d = (r^2 + A^2n - 2rA\sqrt{n} \cos \theta)^{1/2}. \quad (30)$$

The differential volume of the ring-shaped region is $r dr d\theta$ times the surface of a sphere of radius $r \sin \theta$ in $(n - 1)$ -dimensional space; that is

$$r dr d\theta \frac{(n - 1)\pi^{(n-1)/2}(r \sin \theta)^{n-2}}{\Gamma\left(\frac{n + 1}{2}\right)}. \quad (31)$$

Hence, the differential probability for the ring-shaped region is

$$\frac{1}{(\sqrt{2\pi})^n} \exp \left[\frac{-(r^2 + A^2n - 2rA\sqrt{n} \cos \theta)}{2} \right] \cdot \left[\frac{(n - 1)\pi^{(n-1)/2}(r \sin \theta)^{n-2}}{\Gamma\left(\frac{n + 1}{2}\right)} \right] r dr d\theta \quad (32)$$

The differential probability $-dQ$ of being carried between the two cones

is the integral of this expression from zero to infinity on dr :

$$\begin{aligned}
 -dQ &= \frac{1}{2^{n/2}} \frac{(n-1) d\theta}{\sqrt{\pi} \Gamma\left(\frac{n+1}{2}\right)} \\
 &\cdot \int_0^\infty \exp\left[\frac{-(r^2 + A^2n - 2rA\sqrt{n} \cos \theta)}{2}\right] (r \sin \theta)^{n-2} r dr.
 \end{aligned}
 \tag{33}$$

In the exponent we can think of A^2n as $A^2n(\sin^2\theta + \cos^2\theta)$. The \cos^2 part then combines with the other terms to give a perfect square

$$(r - A\sqrt{n} \cos \theta)^2$$

and the \sin^2 term can be taken outside the integral. Thus

$$\begin{aligned}
 -dQ &= \frac{(n-1) \exp\left[-\frac{A^2n \sin^2 \theta}{2}\right] (\sin \theta)^{n-2} d\theta}{2^{n/2} \sqrt{\pi} \Gamma\left(\frac{n+1}{2}\right)} \\
 &\cdot \int_0^\infty \exp\left[\frac{-(r - A\sqrt{n} \cos \theta)^2}{2}\right] r^{n-1} dr.
 \end{aligned}
 \tag{34}$$

We can now direct our attention to estimating the integral, which we call K . The integral can be expressed exactly as a finite, but complicated, sum involving normal distribution functions by a process of continued integration by parts. We are, however, interested in a simple formula giving the asymptotic behavior of the integral as n becomes infinite. This problem was essentially solved by David and Kruskal,⁶ who prove the following asymptotic formula as a lemma:

$$\int_0^\infty z^\nu \exp\left(-\frac{1}{2}z^2 + z\sqrt{\nu+1}w\right) dz \sim \sqrt{2\pi} \left(\frac{\bar{z}}{e}\right)^\nu \exp\left(\frac{1}{2}\bar{z}^2\right) T, \tag{35}$$

as $\nu \rightarrow \infty$, w is fixed, $T = [1 + \frac{1}{4}(\sqrt{w^2 + 4} - w)^2]^{-1/2}$ and

$$\bar{z} = \frac{1}{2}\sqrt{\nu+1}w + \sqrt{\frac{1}{4}(\nu+1)w^2 + \nu}.$$

This is proved by showing that the main contribution to the integral is essentially in the neighborhood of the point \bar{z} where the integral is a maximum. Near this point, when ν is large, the function behaves about as a normal distribution.

The integral K in (34) that we wish to evaluate is, except for a multiplying factor, of the form appearing in the lemma, with

$$z = r, \quad w = A \cos \theta, \quad \nu = n - 1.$$

The integral then becomes

$$\begin{aligned}
 K &= \exp\left(-\frac{A^2 n \cos^2 \theta}{2}\right) \int_0^\infty z^{n-1} \exp\left(-\left(\frac{z^2}{2} + zA \sqrt{n} \cos \theta\right)\right) dz \\
 &\sim \exp\left(-\frac{A^2 n \cos^2 \theta}{2}\right) \sqrt{2\pi} \left(\frac{\bar{z}}{e}\right)^{n-1} T \exp\left(\frac{\bar{z}^2}{2}\right).
 \end{aligned} \tag{36}$$

We have

$$\begin{aligned}
 \bar{z} &= \frac{1}{2} \sqrt{n} A \cos \theta + \sqrt{\frac{1}{4} n A^2 \cos^2 \theta + n - 1} \\
 &= \sqrt{n} \left[\frac{1}{2} A \cos \theta + \sqrt{\frac{A^2}{4} \cos^2 \theta + 1 - \frac{1}{n}} \right] \\
 &= \sqrt{n} \left[\frac{1}{2} A \cos \theta + \sqrt{\frac{A^2}{4} \cos^2 \theta + 1} \right. \\
 &\quad \left. - \frac{1}{2n \sqrt{\frac{A^2}{4} \cos^2 \theta + 1}} + O\left(\frac{1}{n^2}\right) \right].
 \end{aligned} \tag{37}$$

Letting

$$G = \frac{1}{2} [A \cos \theta + \sqrt{A^2 \cos^2 \theta + 4}],$$

we have

$$\bar{z} = \sqrt{n} G \left[1 - \frac{1}{nG \sqrt{A^2 \cos^2 \theta + 4}} + O\left(\frac{1}{n^2}\right) \right],$$

so

$$\begin{aligned}
 \left(\frac{\bar{z}}{e}\right)^{n-1} &= \left(\frac{\sqrt{n} G}{e}\right)^{n-1} \left[1 - \frac{1}{nG \sqrt{A^2 \cos^2 \theta + 4}} + O\left(\frac{1}{n^2}\right) \right]^{n-1} \\
 &\sim \left(\frac{\sqrt{n} G}{e}\right)^{n-1} \exp\left(-\frac{1}{G \sqrt{A^2 \cos^2 \theta + 4}}\right).
 \end{aligned} \tag{38}$$

Also,

$$\begin{aligned}
 \exp \frac{\bar{z}^2}{2} &= \exp \frac{1}{2} n G^2 \left[1 - \frac{1}{nG \sqrt{A^2 \cos^2 \theta + 4}} + O\left(\frac{1}{n^2}\right) \right]^2 \\
 &\sim \exp\left(\frac{1}{2} n G^2 - \frac{2G}{2 \sqrt{A^2 \cos^2 \theta + 4}}\right) \\
 &= \exp\left[\frac{1}{2} n (1 + AG \cos \theta) - \frac{G}{\sqrt{A^2 \cos^2 \theta + 4}}\right],
 \end{aligned} \tag{39}$$

since, on squaring G , we find $G^2 = 1 + AG \cos \theta$. Collecting terms:

$$\begin{aligned}
 K &\sim T \sqrt{2\pi} \left(\frac{\sqrt{n}G}{e} \right)^{n-1} e^{n/2} \exp \left(-\frac{1}{G \sqrt{A^2 \cos^2 \theta + 4}} \right. \\
 &\quad \left. - \frac{G}{\sqrt{A^2 \cos^2 \theta + 4}} - \frac{A^2 n}{2} \cos^2 \theta + \frac{n}{2} AG \cos \theta \right) \quad (40) \\
 &= T \sqrt{2\pi} n^{(n-1)/2} G^{n-1} e^{-n/2} \exp \left(-\frac{n}{2} A^2 \cos^2 \theta + \frac{n}{2} AG \cos \theta \right)
 \end{aligned}$$

since a little algebra shows that the terms

$$1 - \frac{1}{G \sqrt{A^2 \cos^2 \theta + 4}} - \frac{G}{\sqrt{A^2 \cos^2 \theta + 4}}$$

in the exponential cancel to zero. The coefficient of the integral (34), using the asymptotic expression for $\Gamma[(n + 1)/2]$, is asymptotic to

$$\frac{(n - 1) e^{-(\sin^2 \theta)(A^2 n)/2} \sin \theta^{n-2} e^{(n+1)/2}}{2^{n/2} \sqrt{\pi} \left(\frac{n + 1}{2} \right)^{n/2} \sqrt{2\pi}} \quad (41)$$

Combining with the above and collecting terms (we find that $T = G/\sqrt{1 + G^2}$):

$$\begin{aligned}
 -\frac{dQ}{d\theta} &\sim \\
 &\frac{n - 1}{\sqrt{\pi n}} \frac{1}{\sqrt{1 + G^2 \sin^2 \theta}} \left[G \sin \theta \exp \left(-\frac{A^2}{2} + \frac{1}{2} AG \cos \theta \right) \right]^n \quad (42)
 \end{aligned}$$

This is our desired asymptotic expression for the density $dQ/d\theta$.

As we have arranged it, the coefficient increases essentially as \sqrt{n} and there is another term of the form $e^{-E_L(\theta)n}$, where

$$E_L(\theta) = \frac{A^2}{2} - \frac{1}{2} AG \cos \theta - \log (G \sin \theta).$$

It can be shown that if we use for θ the special value $\theta_0 = \cot^{-1} A$ (see Fig. 4) then $E_L(\theta_0) = 0$ and also $E'_L(\theta_0) = 0$. In fact, for this value

$$\begin{aligned}
 G(\theta_0) &= \frac{1}{2} (A \cos \theta_0 + \sqrt{A^2 \cos^2 \theta_0 + 4}) = \frac{1}{2} \left(\frac{A^2}{\sqrt{A^2 + 1}} \right. \\
 &\quad \left. + \sqrt{\frac{A^4}{A^2 + 1} + 4} \right) = \frac{1}{2} \left(\frac{A^2}{\sqrt{A^2 + 1}} + \frac{A^2 + 2}{\sqrt{A^2 + 1}} \right) = \csc \theta_0.
 \end{aligned}$$

Hence the two terms in the logarithm cancel. Also

$$\frac{A^2}{2} - \frac{1}{2}AG \cos \theta_0 = \frac{A^2}{2} - \frac{1}{2}A \sqrt{A^2 + 1} \frac{A}{\sqrt{A^2 + 1}} = 0.$$

So $E_L(\theta_0) = 0$. We also have

$$E'_L(\theta) = \frac{1}{2}AG \sin \theta - \frac{1}{2}AG' \cos \theta - \frac{G'}{G} - \cot \theta. \quad (43)$$

When evaluated, the term $-G'/G$ simplifies, after considerable algebra, to

$$\frac{A \sin \theta}{\sqrt{A^2 \cos^2 \theta + 4}}.$$

Substituting this and the other terms we obtain

$$\begin{aligned} E'_L(\theta) &= \frac{A^2}{2} \sin \theta \cos \theta + \frac{A^3 \cos^2 \theta \sin \theta}{4\sqrt{A^2 \cos^2 \theta + 4}} \\ &+ \frac{A}{4} \frac{(A^2 \cos^2 \theta + 4)}{\sqrt{A^2 \cos^2 \theta + 4}} \sin \theta + \frac{A \sin \theta}{\sqrt{A^2 \cos^2 \theta + 4}} - \cot \theta. \end{aligned} \quad (44)$$

Adding and collecting terms, this simplifies to

$$\begin{aligned} E'_L(\theta) &= \frac{A}{2} (A \cos \theta + \sqrt{A^2 \cos^2 \theta + 4}) \sin \theta - \cot \theta \\ &= AG \sin \theta - \cot \theta \\ &= \cot \theta \left[\frac{A^2}{2} \sin^2 \theta + \frac{A}{2} \sin^2 \theta \sqrt{A^2 + \frac{4}{\cos^2 \theta}} - 1 \right]. \end{aligned} \quad (45)$$

Notice that the bracketed expression is a monotone increasing function of θ ($0 \leq \theta \leq \pi/2$) ranging from -1 at $\theta = 0$ to ∞ at $\theta = \pi/2$. Also, as mentioned above, at θ_0 , $G = \csc \theta_0$ and $A = \cot \theta_0$, so $E'_L(\theta_0) = 0$. It follows that $E'_L(\theta) < 0$ for $0 \leq \theta < \theta_0$ and $E'_L(\theta) > 0$ for $\theta_0 \leq \theta < \pi/2$.

From this, it follows that, in the range from some θ_1 to $\pi/2$ with $\theta_1 > \theta_0$, the minimum $E_L(\theta)$ will occur at the smallest value of θ in the range, that is, at θ_1 . The exponential appearing in our estimate of $Q(\theta)$, namely, $e^{-E_L(\theta)^n}$, will have its *maximum* at θ_1 , for such a range. Indeed, for sufficiently large n , the maximum of the entire expression (45) must occur at θ_1 , since the effect of the n in the exponent will eventually dominate anything due to the coefficient. For, if the coefficient is called $\alpha(\theta)$ with $y(\theta) = \alpha(\theta) e^{-nE_L(\theta)}$, then

$$y'(\theta) = e^{-nE_L(\theta)} [-\alpha(\theta)nE'_L(\theta) + \alpha'(\theta)], \quad (46)$$

and, since $\alpha(\theta) > 0$, when n is sufficiently large $y'(\theta)$ will be negative and the only maximum will occur at θ_1 . In the neighborhood of θ_1 the function goes down exponentially.

We may now find an asymptotic formula for the integral

$$Q(\theta) = \int_{\theta_1}^{\pi/2} \alpha(\theta)e^{-nE_L(\theta)}d\theta + Q(\pi/2) \tag{47}$$

by breaking the integral into two parts,

$$Q(\theta) = \int_{\theta_1}^{\theta_1+n^{-2/3}} + \int_{\theta_1+n^{-2/3}}^{\pi/2} + Q(\pi/2). \tag{48}$$

In the range of the first integral, $(1 - \epsilon)\alpha(\theta_1) \leq \alpha(\theta) \leq \alpha(\theta_1)(1 + \epsilon)$, and ϵ can be made as small as desired by taking n sufficiently large. This is because $\alpha(\theta)$ is continuous and nonvanishing in the range. Also, using a Taylor's series expansion with remainder,

$$e^{-nE_L(\theta)} = \exp \left[-nE_L(\theta_1) - n(\theta - \theta_1)E'_L(\theta_1) - n\frac{(\theta - \theta_1)^2}{2} E''_L(\theta^*) \right], \tag{49}$$

where θ^* is the interval θ_1 to θ . As n increases the maximum value of the remainder term is bounded by $n(n/2)^{-4/3} E''_{\max}$, and consequently approaches zero. Hence, our first integral is asymptotic to

$$\begin{aligned} &\alpha(\theta_1) \int_{\theta_1}^{\theta_1+n^{-2/3}} \exp [-nE_L(\theta_1) - n(\theta - \theta_1)E'_L(\theta_1)] d\theta \\ &= -\alpha(\theta_1) \exp [-nE_L(\theta_1)] \frac{\exp [-n(\theta - \theta_1)E'_L(\theta_1)]}{nE'_L(\theta_1)} \Big]_{\theta_1}^{\theta_1+n^{-2/3}} \\ &\sim \frac{\alpha(\theta_1)e^{-nE_L(\theta_1)}}{nE'_L(\theta_1)}. \end{aligned} \tag{50}$$

since, at large n , the upper limit term becomes small by comparison. The second integral from $\theta_1 + n^{-2/3}$ to $\pi/2$ can be dominated by the value of the integrand at $\theta_1 + n^{-2/3}$ multiplied by the range

$$\pi/2 - (\theta_1 + n^{-2/3}),$$

(since the integrand is monotone decreasing for large n). The value at $\theta_1 + n^{-2/3}$ is asymptotic, by the argument just given, to

$$\alpha(\theta_1) \exp [-nE_L(\theta_1) - n(n^{-2/3}) E'_L(\theta_1)].$$

This becomes small compared to the first integral [as does $Q(\pi/2) =$

$\Phi(-A)$ in (47)] and, consequently, on substituting for $\alpha(\theta_1)$ its value and writing θ for θ_1 , we obtain as an asymptotic expression for $Q(\theta)$:

$$Q(\theta) \sim \frac{1}{\sqrt{n\pi}} \frac{1}{\sqrt{1+G^2 \sin^2 \theta}} \frac{\left[G \sin \theta \exp \left(-\frac{A^2}{2} + \frac{1}{2} AG \cos \theta \right) \right]^n}{(AG \sin^2 \theta - \cos \theta)} \quad (51)$$

$$\left(\frac{\pi}{2} \geq \theta > \theta_0 = \cot^{-1} A \right).$$

This expression gives an asymptotic lower bound for $P_{e \text{ opt}}$, obtained by evaluating $Q(\theta)$ for the θ_1 such that $M\Omega(\theta_1) = \Omega(\pi)$.

Incidentally, the asymptotic expression (51) can be translated into an asymptotic expression for the noncentral t cumulative distribution by substitution of variables $\theta = \cot^{-1}(t/\sqrt{f})$ and $n-1 = f$. This may be useful in other applications of the noncentral t -distribution.

VII. ASYMPTOTIC EXPRESSIONS FOR THE RANDOM CODE BOUND

We now wish to find similar asymptotic expressions for the upper bound on $P_{e \text{ opt}}$ of (20) found by the random code method. Substituting the asymptotic expressions for $dQ(\theta)d\theta$ and for $\Omega(\theta)/\Omega(\pi)$ gives for an asymptotic upper bound the following:

$$Q(\theta_1) + e^{nR} \int_0^{\theta_1} \frac{\Gamma\left(\frac{n}{2} + 1\right) (\sin \theta)^{n-1}}{n \Gamma\left(\frac{n+1}{2}\right) \pi^{1/2} \cos \theta} \sqrt{\frac{n}{\pi}}$$

$$\frac{\left[G \sin \theta \exp \left(-\frac{P}{2N} + \frac{1}{2} \sqrt{\frac{P}{N}} G \cos \theta \right) \right]^n}{\sqrt{1+G^2 \sin^2 \theta}} d\theta. \quad (52)$$

Thus we need to estimate the integral

$$W = \int_0^{\theta_1} \frac{1}{\cos \theta \sin^3 \theta \sqrt{1+G^2}}$$

$$\cdot \exp \left\{ n \left(-\frac{P}{2N} + \frac{1}{2} \sqrt{\frac{P}{N}} G \cos \theta + \log G + 2 \log \sin \theta \right) \right\} d\theta. \quad (53)$$

The situation is very similar to that in estimating $Q(\theta)$. Let the coefficient of n in the exponent be D . Note that $D = -E_L(\theta) + \log \sin \theta$. Hence its derivative reduces to

$$\frac{dD}{d\theta} = -AG \sin \theta + 2 \cot \theta. \quad (54)$$

$dD/d\theta = 0$ has a unique root θ_c , $0 \leq \theta_c \leq \pi/2$ for any fixed $A > 0$. This follows from the same argument used in connection with (45), the only difference being a factor of 2 in the right member. Thus, for $\theta < \theta_c$, $dD/d\theta$ is positive and D is an increasing function of θ . Beyond this maximum, D is a decreasing function.

We may now divide the problem of estimating the integral W into cases according to the relative size of θ_c and θ_1 .

Case 1: $\theta_1 < \theta_c$.

In this case the maximum of the exponent within the range of integration occurs at θ_1 . Consequently, when n is sufficiently large, the maximum of the entire integrand occurs at θ_1 . The asymptotic value can be estimated exactly as we estimated $Q(\theta)$ in a similar situation. The integral is divided into two parts, a part from $\theta_1 - n^{-2/3}$ to θ_1 and a second part from 0 to $\theta_1 - n^{-2/3}$. In the first part the integrand behaves asymptotically like:

$$\frac{1}{\cos \theta_1 \sin^3 \theta_1 \sqrt{1 + G^2(\theta_1)}} \exp \left(n \left\{ -\frac{P}{2N} + \frac{1}{2} \sqrt{\frac{P}{N}} G(\theta_1) \cos \theta_1 + \log G(\theta_1) + 2 \log \sin \theta_1 - (\theta - \theta_1)[AG(\theta_1) \sin \theta_1 - 2 \cot \theta_1] \right\} \right). \tag{55}$$

This integrates asymptotically to

$$\frac{\exp \left\{ n \left[-\frac{P}{2N} + \frac{1}{2} \sqrt{\frac{P}{N}} G(\theta_1) \cos \theta_1 + \log G(\theta_1) + 2 \log \sin \theta_1 \right] \right\}}{\cos \theta_1 \sin^3 \theta_1 \sqrt{1 + G^2(\theta_1)} [-AG(\theta_1) \sin \theta_1 + 2 \cot \theta_1]n}. \tag{56}$$

The second integral becomes small in comparison to this, being dominated by an exponential with a larger negative exponent multiplied by the range $\theta_1 - n^{-2/3}$. With the coefficient

$$\frac{1}{\pi \sqrt{n}} \left[\frac{\Gamma\left(\frac{n}{2} + 1\right)}{\Gamma\left(\frac{n+1}{2}\right)} \right] e^{nR},$$

and using the fact that

$$\frac{\Gamma\left(\frac{n}{2} + 1\right)}{\Gamma\left(\frac{n+1}{2}\right)} \sim \sqrt{\frac{n}{2}},$$

our dominant term approaches

$$\frac{\left[G \sin \theta_1 \exp \left(-\frac{A^2}{2} + \frac{1}{2} AG \cos \theta_1 \right) \right]^n}{\sqrt{n\pi} \sqrt{1 + G^2 \sin^2 \theta_1 (2 \cos \theta_1 - AG \sin^2 \theta_1)}} \quad (57)$$

Combining this with the previously obtained asymptotic expression (51) for $Q(\theta_1)$ we obtain the following *asymptotic expression for the upper bound on $P_{e \text{ opt}}$ for $\theta_1 < \theta_c$* :

$$\left(1 - \frac{\cos \theta_1 - AG \sin^2 \theta_1}{2 \cos \theta_1 - AG \sin^2 \theta_1} \right) \frac{\left[G \sin \theta_1 \exp \left(-\frac{A^2}{2} + \frac{1}{2} AG \cos \theta_1 \right) \right]^n}{\sqrt{n\pi} \sqrt{1 + G^2 \sin^2 \theta_1 (AG \sin^2 \theta_1 - \cos \theta_1)}} \quad (58)$$

Since our lower bound was asymptotic to the same expression without the parenthesis in front, *the two asymptotes differ only by the factor*

$$\left(1 - \frac{\cos \theta_1 - AG \sin^2 \theta_1}{2 \cos \theta_1 - AG \sin^2 \theta_1} \right)$$

independent of n . This factor increases as θ_1 increases from the value θ_0 , corresponding to channel capacity, to the critical value θ_c , for which the denominator vanishes. Over this range the factor increases from 1 to ∞ . In other words, for large n , $P_{e \text{ opt}}$ is determined to within a factor. Furthermore, the percentage uncertainty due to this factor is smaller at rates closer to channel capacity, approaching zero as the rate approaches capacity. It is quite interesting that these seemingly weak bounds can work out to give such sharp information for certain ranges of the variables.

Case 2: $\theta_1 > \theta_c$.

For θ_1 in this range the previous argument does not hold, since the maximum of the exponent is not at the end of the range of integration but rather interior to it. This unique maximum occurs at θ_c , the root of $2 \cos \theta_c - AG \sin^2 \theta_c = 0$. We divide the range of integration into three parts: 0 to $\theta_c - n^{-2/5}$, $\theta_c - n^{-2/5}$ to $\theta_c + n^{-2/5}$ and $\theta_c + n^{-2/5}$ to θ . Proceeding by very similar means, in the neighborhood of θ_c the exponential behaves as

$$\exp \left(-n \left\{ E_L(\theta_c) + \frac{(\theta - \theta_c)^2}{2} E''_L(\theta_c) + O[(\theta - \theta_c)^3] \right\} \right).$$

The coefficient of the exponential approaches constancy in the small interval surrounding θ_c . Thus the integral (53) for this part is asymptotic to

$$\frac{1}{\cos \theta_c \sin^3 \theta_c \sqrt{1 + G^2}} \int \exp \left\{ -n \left[E_L(\theta_c) + \frac{(\theta - \theta_c)^2}{2} E''_L(\theta_c) \right] \right\} d\theta \tag{59}$$

$$\sim \frac{1}{\cos \theta_c \sin^3 \theta_c \sqrt{1 + G^2}} \exp [-nE_L(\theta_c)] \frac{\sqrt{2\pi}}{\sqrt{nE''_L(\theta_c)}}.$$

The other two integrals become small by comparison when n is large, by essentially the same arguments as before. They may be dominated by the value of the integrand at the end of the range near θ_c multiplied by the range of integration. Altogether, then, the integral (52) is asymptotic to

$$\frac{1}{\sqrt{\pi n} \cos \theta_c \sin^3 \theta_c \sqrt{1 + G^2} \sqrt{E''_L(\theta_c)}} e^{-n[E_L(\theta_c) - R]}. \tag{60}$$

The other term in (52), namely, $Q(\theta_1)$, is asymptotically small compared to this, under the present case $\theta > \theta_c$, since the coefficient of n in the exponent for $Q(\theta)$ in (51) will be smaller. Thus, all told, *the random code bound is asymptotic to*

$$\frac{1}{\cos \theta_c \sin^3 \theta_c \sqrt{n\pi E''_L(\theta_c)} [1 + G(\theta_c)^2]} e^{-n[E_L(\theta_c) - R]} \tag{61}$$

for $\theta > \theta_c$ or for rates $R < R_c$ the rate corresponds to θ_c .

Incidentally, the rate R_c is very closely one-half bit less than channel capacity when $A \geq 4$, and approaches this exactly as $A \rightarrow \infty$. For lower values of A the difference $C - R_c$ becomes smaller but the ratio $C/R_c \rightarrow 4$ as $A \rightarrow 0$.

VIII. THE FIRM UPPER BOUND ON $P_{e\text{ opt}}$

In this section we will find an upper bound, valid for all n , on the probability of error by manipulation of the upper bound (20). We first find an upper bound on $Q'(\theta)$. In Ref. 6 the integral (35) is transformed into $\bar{z}' \exp(-\frac{1}{2}\bar{z}^2 + \bar{z}\sqrt{\nu + 1}w)$ times the following integral (in their notation):

$$U = \int_{-\infty}^{\infty} \varphi_z(y) \exp \left\{ -\frac{1}{2} y^2 + \nu \left[\ln \left(1 + \frac{y}{\bar{z}} \right) - \frac{y}{\bar{z}} \right] \right\} dy.$$

It is pointed out that the integrand here can be dominated by $e^{-v^2/2}$. This occurs in the paragraph in Ref. 6 containing Equation 2.6. Therefore, this integral can be dominated by $\sqrt{2\pi}$, and our integral in (34) involved in $dQ/d\theta$ is dominated as follows:

$$\begin{aligned} \int_0^\infty \exp \left[-\frac{(r - A\sqrt{n}\cos\theta)^2}{2} \right] r^{n-1} dr \\ &= \left(\frac{\bar{z}}{e}\right)^{n-1} \exp\left(\frac{\bar{z}}{2}\right)^2 \exp\frac{-A^2n}{2} \cos^2\theta U \\ &\leq \left(\frac{\bar{z}}{e}\right)^{n-1} \exp\left(\frac{\bar{z}}{2}\right)^2 \exp\frac{-A^2n}{2} \cos^2\theta \sqrt{2\pi}. \end{aligned}$$

We have

$$\bar{z} = \frac{1}{2}\sqrt{n}(A\cos\theta + \sqrt{A^2\cos^2\theta + 4 - 4/n}) \leq \sqrt{n}G.$$

Replacing \bar{z} by this larger quantity gives

$$\left(\frac{\sqrt{n}G}{e}\right)^{n-1} \exp\left(\frac{nG^2}{2} - \frac{A^2n}{2}\cos^2\theta\right) \sqrt{2\pi}.$$

We have, then,

$$\begin{aligned} \frac{dQ}{d\theta} \leq \frac{(n-1)\exp\left(\frac{-A^2n}{2}\sin^2\theta\right)(\sin\theta)^{n-2}}{2^{n/2}\sqrt{\pi}\Gamma\frac{n+1}{2}} \left(\frac{\sqrt{n}G}{e}\right)^{n-1} \\ \cdot \exp\left(\frac{nG^2}{2} - \frac{A^2n}{2}\cos^2\theta\right) \sqrt{2\pi}. \end{aligned} \quad (62)$$

Replacing the gamma function by its Stirling expression

$$\left(\frac{n+1}{2}\right)^{n/2} \exp\left(\frac{n+1}{2}\right) \sqrt{2\pi}$$

(which is always too small), and replacing $[1 + (1/n)]^{n/2}$ by $\sqrt{2}$ (which is also too small) again increases the right member. After simplification, we get

$$\begin{aligned} \frac{dQ}{d\theta} \leq \frac{(n-1)(G\sin\theta)^n \exp\left[\left(\frac{n}{2}\right)(-A^2 + 1 + AG\cos\theta)\right]}{\sqrt{n}G\sin^2\theta \sqrt{2\pi} \exp\left(\frac{n-3}{2}\right)} \\ \leq \frac{(n-1)e^{3/2}e^{-E_L(\theta)n}}{\sqrt{2\pi n}G\sin^2\theta}. \end{aligned} \quad (63)$$

Notice that this differs from the asymptotic expression (42) only by a factor

$$\frac{e^{3/2} \sqrt{1 + G^2}}{\sqrt{2} G} \leq e^{3/2}$$

(since $G \geq 1$). A firm upper bound can now be placed on $Q(\theta)$:

$$Q(\theta_1) = \int_{\theta_1}^{\pi/2} \frac{dQ}{d\theta} d\theta + Q\left(\frac{\pi}{2}\right).$$

We use the upper bound above for $dQ/d\theta$ in the integral. The coefficient of $-n$ in the exponent of e

$$E_L(\theta) = \frac{1}{2}(A^2 - AG \cos \theta) - \log G \sin \theta$$

is positive and monotone increasing with θ for $\theta > \theta_0$, as we have seen previously. Its derivative is

$$E'_L(\theta) = AG \sin \theta - \cot \theta.$$

As a function of θ this curve is as shown in Fig. 5, either rising monotonically from $-\infty$ at $\theta = 0$ to A at $\theta = \pi/2$, or with a single maximum. In any case, the curve is concave downward. To show this analytically, take the second derivative of E'_L . This consists of a sum of negative terms.

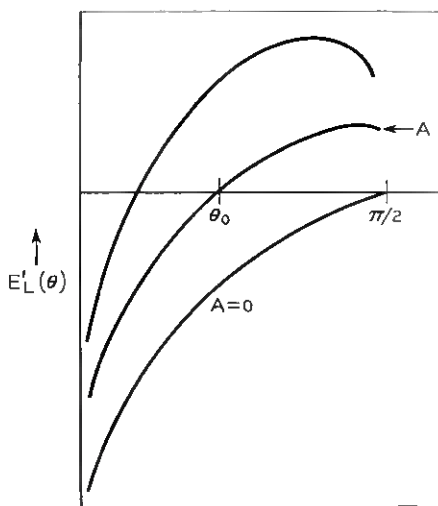


Fig. 5 — $E'_L(\theta)$ as a function of θ .

Returning to our upper bound on Q , the coefficient in (63) does not exceed

$$\frac{\sqrt{n}}{\sqrt{2\pi}} e^{3/2} \frac{1}{\sin^2 \theta_1},$$

replacing $\sin \theta$ and G by $\sin \theta_1$ and 1, their minimum values in the range. We now wish to replace $e^{-nE_L(\theta)}$ by

$$\exp -n[E_L(\theta_1) + (\theta - \theta_1)h].$$

If h is chosen equal to the minimum $E'_L(\theta)$, this replacement will increase the integral and therefore give an upper bound. From the behavior of $E'_L(\theta)$ this minimum occurs at either θ_1 or $\pi/2$. Thus, we may take $h = \min [A, AG(\theta_1) \sin \theta_1 - \cot \theta_1]$. With this replacement the integral becomes a simple exponential and can be immediately integrated.

The term $Q(\pi/2)$ is, of course,

$$\Phi(-A \sqrt{n}) \leq \frac{1}{\sqrt{2\pi n} A} e^{-A^2 n/2}.$$

If we continue the integral out to infinity instead of stopping at $\pi/2$, the extra part added will more than cover $Q(\pi/2)$. In fact, $E_L(\pi/2) = A^2/2$, so the extra contribution is at least

$$\frac{\sqrt{n} e^{3/2}}{An \sin^2 \theta_1 \sqrt{2\pi}} e^{-A^2 n/2},$$

if we integrate

$$\frac{\sqrt{n} e^{3/2}}{\sin^2 \theta_1 \sqrt{2\pi}} e^{-A^2 n/2 - n(\theta - \theta_1)A}$$

to ∞ instead of stopping at $\pi/2$. Since $e^{3/2}/\sin^2 \theta_1 \geq 1$, we may omit the $Q(\pi/2)$ term in place of the extra part of the integral.

Consequently, we can bound $Q(\theta_1)$ as follows:

$$Q(\theta_1) \leq \frac{e^{3/2} \exp \{ (n/2) [AG(\theta_1) \cos \theta_1 - A^2 + 2 \log G \sin \theta_1] \}}{\sqrt{2\pi n} \sin^2 \theta_1 \min (A, AG(\theta_1) \sin \theta_1 - \cot \theta_1)}. \quad (64)$$

In order to overbound P_{opt} by (3) it is now necessary to overbound the term

$$\int_0^{\theta_1} \frac{\Omega(\theta)}{\Omega(\theta_1)} dQ(\theta).$$

This can be done by a process very similar to that just carried out for $\int dQ(\theta)$. First, we overbound $\Omega(\theta)/\Omega(\theta_1)$ using (21). We have

$$\begin{aligned} \frac{\Omega(\theta)}{\Omega(\theta_1)} &= \frac{\int_0^\theta (\sin x)^{n-2} dx}{\int_0^{\theta_1} (\sin x)^{n-2} dx} \\ &= \frac{\int_0^\theta (\sin x)^{n-2} dx}{\int_0^\theta (\sin x)^{n-2} dx + \int_\theta^{\theta_1} (\sin x)^{n-2} dx} \\ &\cong \frac{\int_0^\theta (\sin x)^{n-2} \cos x dx}{\int_0^\theta (\sin x)^{n-2} \cos x dx + \cos \theta \int_\theta^{\theta_1} (\sin x)^{n-2} dx} \\ &\cong \frac{\int_0^\theta (\sin x)^{n-2} \cos x dx}{\int_0^\theta (\sin x)^{n-2} \cos x dx + \int_\theta^{\theta_1} (\sin x)^{n-2} \cos x dx}, \end{aligned}$$

and, finally,

$$\frac{\Omega(\theta)}{\Omega(\theta_1)} \cong \frac{(\sin \theta)^{n-1}}{(\sin \theta_1)^{n-1}}. \tag{65}$$

Here the third line follows since the first integral in the denominator is reduced by the same factor as the numerator and the second integral is reduced more, since $\cos \theta$ is decreasing. In the next line, the denominator is reduced still more by taking the cosine inside.

Using this inequality and also the upper bound (63) on $dQ/d\theta$, we have

$$\begin{aligned} \int_0^{\theta_1} \frac{\Omega(\theta)}{\Omega(\theta_1)} dQ(\theta) &\cong \int_0^{\theta_1} \frac{(\sin \theta)^{n-1}}{(\sin \theta_1)^{n-1}} \frac{(n-1)e^{3/2}(G \sin \theta)^n e^{(n/2)(-A^2+A \cos \theta G)}}{\sqrt{2\pi n G} \sin^2 \theta} d\theta \tag{66} \\ &= \frac{(n-1)e^{3/2}}{\sqrt{2\pi n} (\sin \theta_1)^{n-1}} \int_0^{\theta_1} G^n (\sin \theta)^{2n-3} e^{(n/2)(-A^2+A \cos \theta G)} d\theta. \end{aligned}$$

Near the point θ_1 the integrand here behaves like an exponential when n is large (provided $\theta_1 < \theta_c$), and it should be possible to find a firm

upper bound of the form

$$\frac{k}{\sqrt{n}} e^{-E_L(\theta_1)n},$$

where k would not depend on n . This, however, leads to considerable complexity and we have settled for a cruder formulation as follows:

The integrand may be bounded by its maximum values. If $\theta_1 < \theta_c$, the maximum of the integrand will occur at θ_1 , at least when n is large enough. In this case, the integral will certainly be bounded by

$$\theta_1 G^n(\theta_1) (\sin \theta_1)^{2n-3} e^{(n/2)[-A^2 + A \cos \theta_1 G(\theta_1)]}.$$

The entire expression for $P_{e \text{ opt}}$ may then be bounded by [adding in the bound (64) on $Q(\theta_1)$]

$$P_{e \text{ opt}} \leq \frac{\sqrt{ne^{3/2}} \theta_1 e^{-E_L(\theta_1)}}{\sqrt{2\pi} \sin^2 \theta_1} \left\{ 1 + \frac{1}{n \theta_1 \min [A, AG(\theta_1) \sin \theta_1 - \cot \theta_1]} \right\}, \quad (67)$$

It must be remembered that (67) is valid only for $\theta_1 < \theta_c$ and if n is large enough to make the maximum of the integrand above occur at θ . For $\theta_1 > \theta_c$, bounds could also be constructed based on the maximum value of the integrand.

IX. A FIRM LOWER BOUND ON $P_{e \text{ opt}}$

In this section we wish to find a lower bound on $P_{e \text{ opt}}$ that is valid for all n . To do this we first find a lower bound on $Q'(\theta)$ and from this find a lower bound on $Q(\theta)$. The procedure is quite similar to that involved in finding the firm upper bound.

In Ref. 6, the integral (35) above was reduced to the evaluation of the following integral (Equation 2.5 of Ref. 6):

$$\begin{aligned} & \int_{-\infty}^{\infty} \left(1 + \frac{y}{z}\right)^{\nu} \exp\left(-\frac{1}{2}y^2 - y\frac{\nu}{z}\right) dy \\ & \geq \int_0^{\infty} \exp\left\{-\frac{1}{2}y^2 + \nu\left[\ln\left(1 + \frac{y}{z}\right) - \frac{y}{z}\right]\right\} dy \\ & \geq \int_0^{\infty} \exp\left[-\frac{1}{2}y^2 + \nu\left(\frac{-y^2}{2z^2}\right)\right] dy \\ & = \int_0^{\infty} \exp\left[\frac{-y^2}{2}\left(1 + \frac{\nu}{z^2}\right)\right] dy = \frac{1}{2} \frac{\sqrt{2\pi}}{\sqrt{1 + \frac{\nu}{z^2}}} \\ & \geq \frac{\sqrt{2\pi}}{2\sqrt{2}} = \frac{\sqrt{\pi}}{2} \end{aligned}$$

Here we used the inequality

$$\ln \left(1 + \frac{y}{\bar{z}} \right) - \frac{y}{\bar{z}} \geq - \frac{y^2}{2\bar{z}^2} \quad \text{for} \quad \frac{y}{\bar{z}} > 0,$$

and also the fact that $\nu/\bar{z}^2 \leq 1$. This latter follows from Equation 2.3 of Ref. 6 on dividing through by \bar{z}^2 .

Using this lower bound, we obtain from (34)

$$\frac{dQ}{d\theta} \geq \frac{(n-1) \sin^{n-2} \theta \exp\left(\frac{-A^2 n}{2}\right)}{2^{n/2} \sqrt{\pi} \Gamma\left(\frac{n+1}{2}\right)} \left(\frac{\bar{z}}{e}\right)^{n-1} \exp\left(\frac{\bar{z}^2}{2}\right) \frac{\sqrt{\pi}}{2}. \quad (68)$$

Now $\bar{z} \geq \sqrt{n-1} G$ and

$$\Gamma\left(\frac{n+1}{2}\right) < \left(\frac{n+1}{2}\right)^{n/2} e^{-(n+1)/2} \sqrt{2\pi} \exp\left[\frac{1}{6(n+1)}\right]$$

and, using the fact that

$$\left(\frac{n-1}{n+1}\right)^{n/2} \geq \frac{1}{3} \quad \text{for} \quad n \geq 2,$$

we obtain

$$\frac{dQ}{d\theta} \geq \frac{1}{6\sqrt{2\pi}} \frac{\sqrt{n-1} e^{3/2} e^{-nE_L(\theta)}}{G \exp\left[\frac{G^2}{2} + \frac{1}{6(n+1)}\right] \sin^2 \theta} \quad \text{for } n \geq 2. \quad (69)$$

This is our lower bound on $dQ/d\theta$.

To obtain a lower bound on $Q(\theta)$ we may use the same device as before—here, however, replacing the coefficient by its minimum value in the range and the exponent by $-nE_L(\theta_1) - n(\theta - \theta_1)E'_L \max$:

$$\begin{aligned} E'_L &= AG \sin \theta - \cot \theta \\ &\geq AG \\ &\geq A(A+1). \end{aligned}$$

Similarly, in the coefficient, G can be dominated by $A+1$ and $\sin^2 \theta$ by 1. Thus,

$$Q(\theta_1) \geq \int_{\theta_1}^{\pi/2} \frac{\sqrt{n-1} e^{3/2} e^{-nE_L(\theta_1)} e^{-n(\theta-\theta_1)A(A+1)}}{6\sqrt{2\pi}(A+1) \exp\left[\frac{(A+1)^2}{2} + \frac{1}{6(n+1)}\right]} d\theta + Q\left(\frac{\pi}{2}\right). \quad (70)$$

Integrating and observing that the term due to the $\pi/2$ limit can be absorbed into the $Q(\pi/2) - \text{erf } A$, we arrive at the lower bound:

$$Q(\theta_1) \geq \frac{\sqrt{n-1} e^{3/2} e^{-nE_L(\theta_1)}}{6\sqrt{2\pi n} (A+1)^3 \exp\left[\frac{(A+1)^2}{2} + \frac{1}{6(n+1)}\right]} \quad (71)$$

X. BEHAVIOR NEAR CHANNEL CAPACITY

As we have seen, near channel capacity the upper and lower asymptotic bounds are substantially the same. If in the asymptotic lower bound (42) we form a Taylor expansion for θ near θ_0 , retaining terms up to $(\theta - \theta_0)^2$, we will obtain an expression applying to the neighborhood of channel capacity. Another approach is to return to the original noncentral t -distribution and use its normal approximation which will be good near the mean (see Ref. 5). Either approach gives, in this neighborhood, the approximations [since $E(\theta_0) = E'(\theta_0) = 0$]:

$$-\frac{dQ}{d\theta} \doteq \frac{\sqrt{n} (1+A^2)}{\sqrt{\pi} \sqrt{2+A^2}} \exp\left[-n \frac{(A^2+1)^2}{A^2+2} (\theta - \theta_0)^2\right] \quad (72)$$

$$Q(\theta) \doteq \Phi\left[(\theta_0 - \theta) \frac{A^2+1}{\sqrt{A^2+2}} \sqrt{2n}\right],$$

or, since near channel capacity, using $e^{-x} \doteq \sin x$,

$$\theta - \theta_0 \doteq A^{-1}(C - R)$$

$$P_{e, \text{opt}}\left(n, R, \sqrt{\frac{P}{N}}\right) \doteq \Phi\left[\sqrt{2n} A^{-1} \frac{A^2+1}{\sqrt{A^2+2}} (R - C)\right] \quad (73)$$

$$= \Phi\left[\frac{P+N}{\sqrt{P(P+2N)}} \sqrt{2n} (R - C)\right].$$

The reliability curve is approximated near C by

$$E(R) \doteq \frac{(P+N)^2}{P(P+2N)} (C - R)^2. \quad (74)$$

It is interesting that Rice² makes estimates of the behavior of what amounts to a lower bound on the exponent E near channel capacity. His exponent, translated into our notation, is

$$E^*(R) \doteq \frac{P+N}{2P} (C - R)^2,$$

a poorer value than (74); that is, it will take a larger block length to

achieve the same probability of error. This difference is evidently due to the slight difference in the manner of construction of the random codes. Rice's codes are obtained by placing points according to an n -dimensional gaussian distribution, each coordinate having variance P . In our codes the points are placed at random on a sphere of precisely fixed radius \sqrt{nP} . These are very close to the same thing when n is large, since in Rice's situation the points will, with probability approaching 1, lie between the spheres of radii $\sqrt{nP}(1 - \epsilon)$ and $\sqrt{nP}(1 + \epsilon)$, (any $\epsilon > 0$). However, we are dealing with very small probability events in any case when we are estimating probability of error, and the points within the sphere are sufficiently important to affect the exponent E . In other words, the Rice type of code is sufficient to give codes that will have a probability of error approaching zero at rates arbitrarily near channel capacity. However, they will not do so at as rapid a rate (even in the exponent) as can be achieved. To achieve the best possible E it is evidently necessary to avoid having too many of the code points interior to the \sqrt{nP} sphere.

At rates R greater than channel capacity we have $\theta_1 < \theta_0$. Since the Q distribution approaches normality with mean at θ_0 and variance $2n(A^2 + 1)^2/(A^2 + 2)$, we will have $Q(\theta_1)$ approaching 1 with increasing n for any fixed rate greater than C . Indeed, even if the rate R varies but remains always greater than C (perhaps approaching it from above with increasing n), we will still have $P_{e, \text{opt}} > \frac{1}{2} - \epsilon$ for any $\epsilon > 0$ and sufficiently large n .

XI. UPPER BOUND ON $P_{e, \text{opt}}$ BY METHOD OF EXHAUSTION

For low rates of transmission, where the upper and lower bounds diverge widely, we may obtain better estimates by other methods. For very low rates of transmission, the main contribution to the probability of error can be shown to be due to the code points that are nearest together and thus often confused with each other, rather than to the general average structure of the code. The important thing, at low rates, is to maximize the minimum distance between neighbors. Both the upper and lower bounds which we will derive for low rates are based on these considerations.

We will first show that, for $D \leq \sqrt{2nP}$, it is possible to find at least

$$M_D = \left(\sin 2 \sin^{-1} \frac{D}{2\sqrt{nP}} \right)^{1-n}$$

points on the surface of an n sphere of radius \sqrt{nP} such that no pair

of them is separated by a distance less than D . (If M_D is not an integer, take the next larger integer.) The method used will be similar to one used by E. N. Gilbert for the binary symmetric channel.

Select any point on the sphere's surface for the first point. Delete from the surface all points within D of the selected point. In Fig. 6, x is the selected point and the area to be deleted is that cut out by the cone. This area is certainly less (if $D \leq \sqrt{2nP}$) than the area of the hemisphere of radius H shown and, even more so, less than the area of the sphere of radius H . If this deletion does not exhaust the original sphere, select any point from those remaining and delete the points within D of this new point. This again will not take away more area than that of a sphere of radius H . Continue in this manner until no points remain. Note that each point chosen is at least D from each preceding point. Hence all interpoint distances are at least D . Furthermore, this can be continued at least as many times as the ratio of the surface of a sphere of radius \sqrt{nP} to that of a sphere of radius H , since each deletion takes away not more than this much surface area. This ratio is clearly

$$(\sqrt{nP}/H)^{n-1}.$$

By simple geometry in Fig. 6, we see that H and D are related as follows:

$$\begin{aligned}\sin \theta &= \frac{H}{\sqrt{nP}}, \\ \sin \frac{\theta}{2} &= \frac{D}{2\sqrt{nP}}.\end{aligned}$$

Hence

$$H = \sqrt{nP} \sin 2 \sin^{-1} \frac{D}{2\sqrt{nP}}. \quad (75)$$

Substituting, we can place at least

$$M_D = \left(\sin 2 \sin^{-1} \frac{D}{2\sqrt{nP}} \right)^{-(n-1)}$$

points at distances at least D from each other, for any $D \leq \sqrt{2nP}$.

If we have M_D points with minimum distance at least D , then the probability of error with optimal decoding will be less than or equal to

$$M_D \Phi \left(\frac{-D}{2\sqrt{N}} \right).$$

To show this we may add up pessimistically the probabilities of each

point being received as each other point. Thus the probability of point 1 being moved closer to point 2 than to the original point 1 is not greater than $\Phi[-D/(2\sqrt{N})]$, that is, the probability of the point being moved in a certain direction at least $D/2$ (half the minimum separation). The contribution to errors due to this cause cannot, therefore, exceed $(1/M_D)\Phi[-D/(2\sqrt{N})]$, (the $1/M_D$ factor being the probability of message 1 being transmitted). A similar argument occurs for each (ordered) pair of points, a total of $M_D(M_D - 1)$ contributions of this kind. Consequently, the probability of error cannot exceed $(M_D - 1)\Phi[-D/(2\sqrt{N})]$ or, more simply, $M_D\Phi[-D/(2\sqrt{N})]$.

If we set

$$e^{nR} = M_D = \left(\sin 2 \sin^{-1} \frac{D}{2\sqrt{nP}} \right)^{-(n-1)}$$

then the rate R (in natural units) is

$$R = \left(1 - \frac{1}{n} \right) \log \left(\sin 2 \sin^{-1} \frac{D}{2\sqrt{nP}} \right)^{-1}$$

with

$$P_e \leq e^{nR} \Phi \left(\frac{-D}{2\sqrt{N}} \right) \leq e^{nR} \frac{\sqrt{2N}}{D\sqrt{\pi}} e^{-(D^2/8N)}, \tag{76}$$

using the well-known upper bound $\Phi(-x) \leq (1/x\sqrt{2\pi})e^{-x^2/2}$. These are

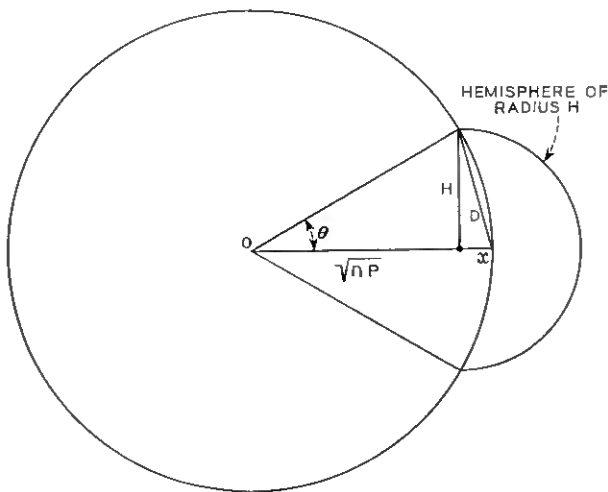


Fig. 6 — Geometry of sphere of radius \sqrt{nP} .

parametric equations in terms of D . It is more convenient to let

$$D = \lambda\sqrt{2nP}.$$

We then have

$$R = \left(1 - \frac{1}{n}\right) \log \left(\sin 2 \sin^{-1} \frac{\lambda}{\sqrt{2}}\right)^{-1},$$

$$P_e \leq \frac{1}{\lambda \sqrt{\pi n \frac{P}{N}}} e^{n[R - (\lambda^2 P)/(4N)]}. \quad (77)$$

The asymptotic reliability, that is, the coefficient of $-n$ in the exponent of P_e , is given by $(\lambda^2 P/4N) - R$. This approaches

$$\left(\sin \frac{1}{2} \sin^{-1} e^R\right)^2 \frac{P}{2N} - R \quad \text{as} \quad n \rightarrow \infty.$$

Thus our asymptotic *lower* bound for reliability is (eliminating λ):

$$E \geq \left(\sin \frac{1}{2} \sin^{-1} e^R\right)^2 \frac{P}{2N} - R. \quad (78)$$

As $R \rightarrow 0$ the right-hand expression approaches $P/(4N)$.

This lower bound on the exponent is plotted in the curves in Section XIV and it may be seen to give more information at low rates than the random code bound. It is possible, however, to improve the random coding procedure by what we have called an "expurgating" process. It then becomes the equal of the bound just derived and, in fact, is somewhat stronger over part of the range. We shall not go into this process in detail but only mention that the expurgating process consists of eliminating from the random code ensemble points which have too close neighbors, and working with the codes that then remain.

XII. LOWER BOUND ON P_e IN GAUSSIAN CHANNEL BY MINIMUM DISTANCE ARGUMENT

In a code of length n with M code words, let m_{is} ($i = 1, 2, \dots, M$; $s = 1, 2, \dots, n$) be the s th coordinate of code word i . We are here assuming an average power limitation P , so that

$$\frac{1}{nM} \sum_{i,s} m_{is}^2 \leq P. \quad (79)$$

We also assume an independent gaussian noise of power N added to each coordinate.

We now calculate the average squared distance between all the $M(M - 1)/2$ pairs of points in n -space corresponding to the M code words. The squared distance from word i to word j is

$$\sum_s (m_{is} - m_{js})^2.$$

The average $\overline{D^2}$ between all pairs will then be

$$\overline{D^2} = \frac{1}{M(M - 1)} \sum_{s,i,j} (m_{is} - m_{js})^2.$$

Note that each distance is counted twice in the sum and also that the extraneous terms included in the sum, where $i = j$, contribute zero to it. Squaring the terms in the sum,

$$\begin{aligned} \overline{D^2} &= \frac{1}{M(M - 1)} \left(\sum_{i,j,s} m_{is}^2 - 2 \sum_s \sum_{i,j} m_{is} m_{js} + \sum_{i,j,s} m_{js}^2 \right) \\ &= \frac{1}{M(M - 1)} \left[2M \sum_{i,s} m_{is}^2 - 2 \sum_s \left(\sum_i m_{is} \right)^2 \right] \\ &\leq \frac{1}{M(M - 1)} 2MPnM \\ \overline{D^2} &\leq \frac{2nMP}{M - 1}, \end{aligned} \tag{80}$$

where we obtain the third line by using the inequality on the average power (79) and by noting that the second term is necessarily non-positive.

If the *average* squared distance between pairs of points is

$$\leq (2nMP)/(M - 1),$$

there must exist a pair of points for whose distance this inequality holds. Each point in this pair is used $1/M$ of the time. The best detection for separating this pair (if no other points were present) would be by a hyperplane normal to and bisecting the joining line segment. Either point would then give rise to a probability of error equal to that of the noise carrying a point half this distance or more in a specified direction. We obtain, then, a contribution to the probability of error at least

$$\begin{aligned} \frac{1}{M} \cdot \Pr \left\{ \text{noise in a certain direction} \geq \frac{1}{2} \sqrt{\frac{2nMP}{M - 1}} \right\} \\ = \frac{1}{M} \Phi \left[-\sqrt{\frac{nMP}{(M - 1)2N}} \right]. \end{aligned}$$

This we may assign to the first of the two points in question, and the errors we have counted are those when this message is sent and is received closer to the second message (and should therefore be detected as the second or some other message).

Now delete this first message from the set of code points and consider the remaining $M - 1$ points. By the same argument there must exist among these a pair whose distance is less than or equal to

$$\sqrt{\frac{2nP(M-1)}{(M-2)2N}}$$

This pair leads to a contribution to probability of error, due to the first of these being displaced until nearer the second, of an amount

$$\frac{1}{M} \Phi \left[- \sqrt{\frac{(M-1)nP}{(M-2)2N}} \right].$$

This same argument is continued, deleting points and adding contributions to the error, until only two points are left. Thus we obtain a lower bound on $P_{e \text{ opt}}$ as follows:

$$P_{e \text{ opt}} \cong \frac{1}{M} \left[\Phi \left(- \sqrt{\frac{nP}{2N} \frac{M}{M-1}} \right) + \Phi \left(- \sqrt{\frac{nP}{2N} \frac{M-1}{M-2}} \right) + \dots + \Phi \left(- \sqrt{\frac{nP}{2N} \frac{2}{1}} \right) \right]. \quad (81)$$

To simplify this bound somewhat, one may take only the first $M/2$ terms [or $(M+1)/2$ if M is odd]. Since they are decreasing, each term would be reduced by replacing it with the last term taken. Thus we may reduce the bound by these operations and obtain

$$P_{e \text{ opt}} \cong \frac{1}{2} \Phi \left(- \sqrt{\frac{M}{M-2} \frac{nP}{2N}} \right). \quad (82)$$

For any rate $R > 0$, as n increases the term $M/(M-2)$ approaches 1 and the bound, then, behaves about as

$$\frac{1}{2} \Phi \left(- \sqrt{\frac{nP}{2N}} \right).$$

This is asymptotic to

$$\frac{1}{2\sqrt{\frac{\pi nP}{N}}} e^{-(nP)/(4N)}.$$

It follows that the reliability $E \leq P/(4N) = A^2/4$. This is the same value as the lower bound for E when $R \rightarrow 0$.

XIII. ERROR BOUNDS AND OTHER CONDITIONS ON THE SIGNAL POINTS

Up to now we have (except in the last section) assumed that all signal points were required to lie on the surface of the sphere, i.e., have a mean square value \sqrt{nP} . Consider now the problem of estimating $P'_{e \text{ opt}}(M, n, \sqrt{P/N})$, where the signal points are only required to lie on or within the spherical surface. Clearly, since this relaxes the conditions on the code, it can only improve, i.e., decrease the probability of error for the best code. Thus $P'_{e \text{ opt}} \leq P_{e \text{ opt}}$.

On the other hand, we will show that

$$P'_{e \text{ opt}} \left(M, n, \sqrt{\frac{P}{N}} \right) \geq P_{e \text{ opt}} \left(M, n + 1, \sqrt{\frac{P}{N}} \right). \quad (83)$$

In fact, suppose we have a code of length n , all points on or within the

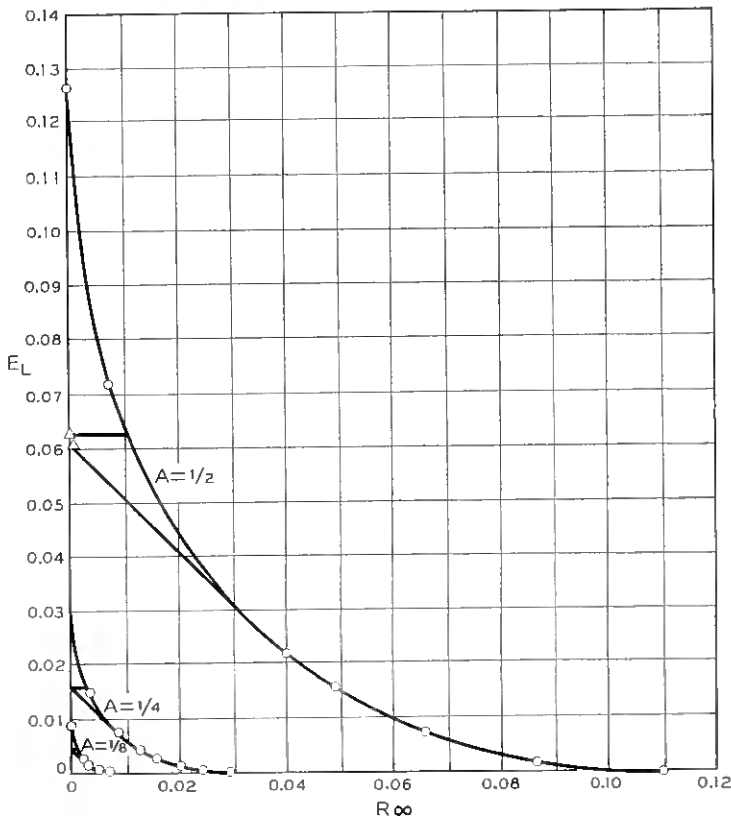


Fig. 7 — Curves showing E_L vs. R for $A = \frac{1}{8}, \frac{1}{4}$ and $\frac{1}{2}$.

n sphere. To each code word add a further coordinate of such value that in the $n + 1$ space the point thus formed lies *exactly* on the $n + 1$ sphere surface. If the first n coordinates of a point have values x_1, x_2, \dots, x_n with

$$\sum_{i=1}^n x_i^2 \leq nP,$$

the added coordinate will have the value

$$x_{n+1} = \sqrt{(n+1)P - \sum_{i=1}^n x_i^2}.$$

This gives a derived code of the first type (all points *on* the $n + 1$ sphere surface) with M words of length $n + 1$ at signal-to-noise ratio P/N . The probability of error for the given code is at least as great as that of the derived code, since the added coordinate can only improve

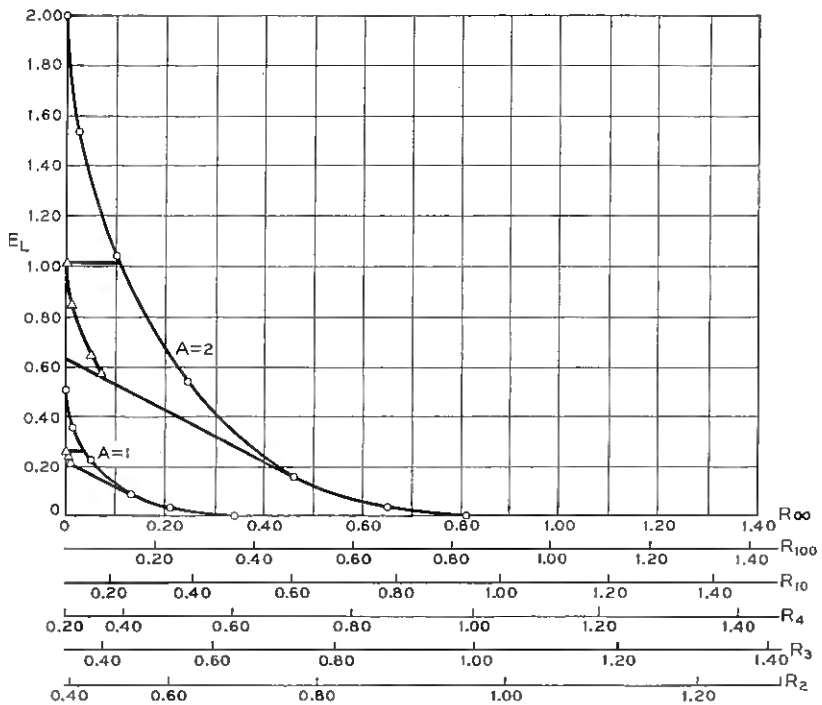


Fig. 8 — Curves showing E_L vs. different values of R for $A = 1$ and 2.

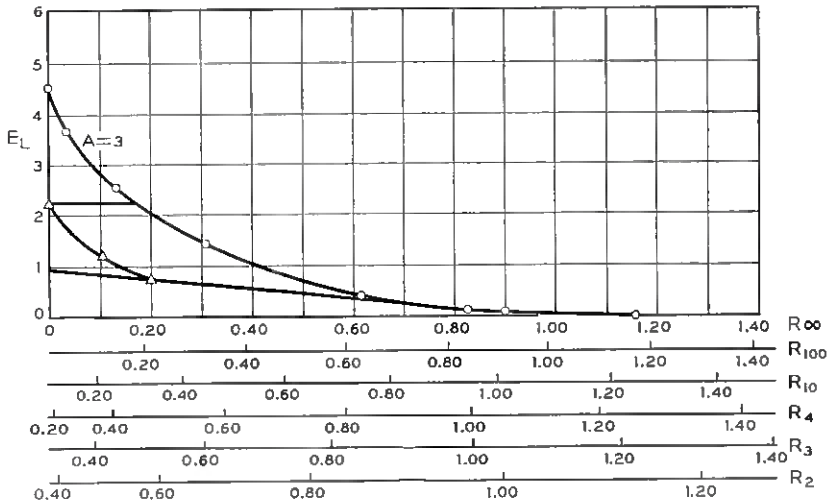


Fig. 9 — Curves showing E_L vs. different values of R for $A = 3$.

the decoding process. One might, for example, decode ignoring the last coordinate and then have the same probability of error. Using it in the best way would, in general, improve the situation.

The probability of error for the derived code of length $n + 1$ must be greater than or equal to that of the optimal code of the length $n + 1$ with all points on the surface. Consequently we have (83). Since $P_{e\text{opt}}(M, n, \sqrt{P/N})$ varies essentially exponentially with n when n is large, the effect of replacing n by $n + 1$ is essentially that of a constant multiplier. Thus, our upper bounds on $P_{e\text{opt}}$ are not changed and our lower bounds are multiplied by a quantity which does not depend much on n when n is large. The asymptotic reliability curves consequently will be the same. Thus the E curves we have plotted may be applied in either case.

Now consider the third type of condition on the points, namely, that the average squared distance from the origin of the set of points be less than or equal to nP . This again is a weakening of the previous conditions and hence the optimal probability of error, $P''_{e\text{opt}}$, is less than or equal to that of the previous cases:

$$P''_{e\text{opt}}\left(M, n, \frac{P}{N}\right) \leq P'_{e\text{opt}}\left(M, n, \frac{P}{N}\right) \leq P_{e\text{opt}}\left(M, n, \frac{P}{N}\right). \quad (84)$$

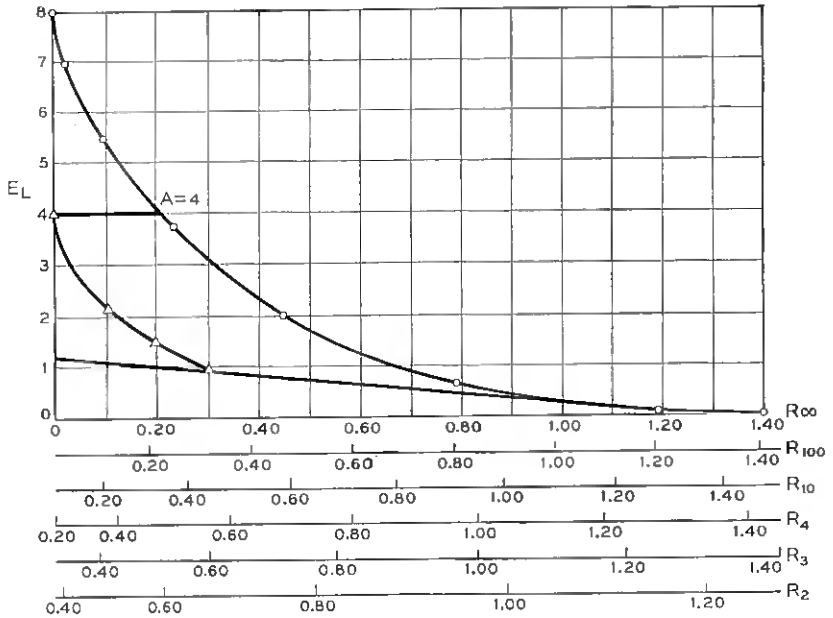


Fig. 10 — Curves showing E_L vs. different values of R for $A = 4$.

Our upper bounds on probability of error (and, consequently, lower bounds on reliability) can be used as they stand.

Lower bounds on P''_{opt} may be obtained as follows. If we have M points whose mean square distance from the origin does not exceed nP , then for any α ($0 < \alpha \leq 1$) at least αM of the points are within a sphere of squared radius $nP/(1 - \alpha)$. [For, if more than $(1 - \alpha)M$ of them were outside the sphere, these alone would contribute more than

$$(1 - \alpha)MnP/(1 - \alpha)$$

to the total squared distance, and the mean would then necessarily be greater than nP .] Given an optimal code under the third condition, we can construct from it, by taking αM points within the sphere of radius $\sqrt{nP/(1 - \alpha)}$, a code satisfying the second condition with this smaller number of points and larger radius. The probability of error for the new code cannot exceed $1/\alpha$ times that of the original code. (Each new code word is used $1/\alpha$ times as much; when used, its probability of error is at least as good as previously.) Thus:

$$\begin{aligned}
 P''_{e \text{ opt}} \left(M, n, \sqrt{\frac{P}{N}} \right) &\geq \frac{1}{\alpha} P'_{e \text{ opt}} \left(\alpha M, n, \sqrt{\frac{P}{(1-\alpha)N}} \right) \\
 &\geq \frac{1}{\alpha} P_{e \text{ opt}} \left(\alpha M, n+1, \sqrt{\frac{P}{(1-\alpha)N}} \right).
 \end{aligned}$$

XIV. CURVES FOR ASYMPTOTIC BOUNDS

Curves have been calculated to facilitate evaluation of the exponents in these asymptotic bounds. The basic curves range over values of

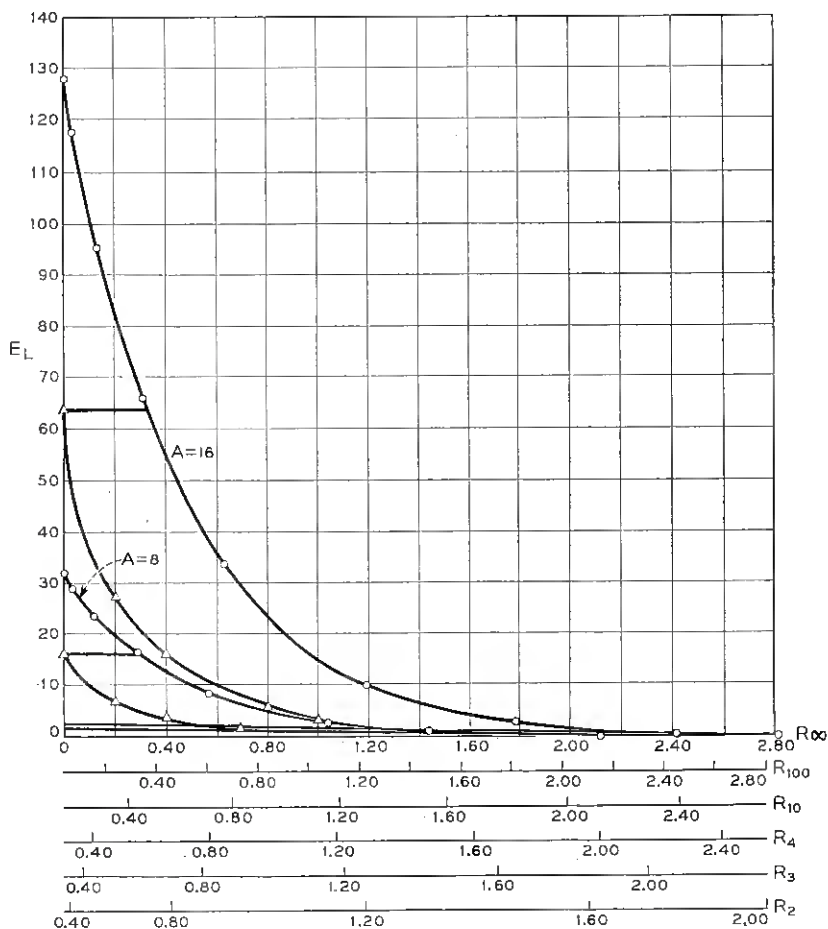


Fig. 11 — Curves showing E_L vs. different values of R for $A = 8$ and 16 .

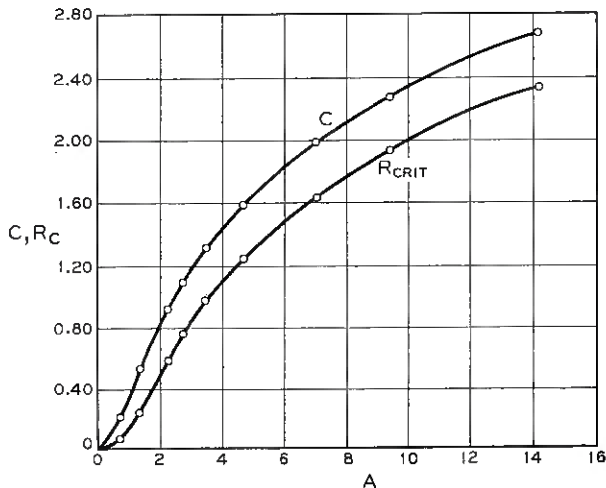


Fig. 12 — Channel capacity, C , and critical rate, R_c , as functions of θ .

$A = \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2, 3, 4, 8, 16$. Figs. 7 through 11 give the coefficients of n and E_L as functions of the rate R . Since E_L strictly is a function of θ , and the relation between θ and R depends somewhat on n , a number of slightly different R scales are required at the bottom of the curve. This, however, was considered a better means of presenting the data than the use of auxiliary curves to relate R and θ . These same curves give the coefficient of n in the upper bounds (the straight line part together with the curve to the right of the straight line segment). The point of tangency is the critical R (or critical θ). In other words, the curve and the curve plus straight line, read against the $n = \infty$ scale, give upper and lower bounds on the reliability measure. The upper and lower bounds on E for low R are also included in these curves. The upper bound is the horizontal line segment running out from $R = 0, E = A^2/4$. The lower bound is the curved line running down from this point to the tangent line. Thus, the reliability E lies in the four-sided figure defined by these lines to the left of R_c . It is equal to the curve to the right of R_c . Fig. 12 gives channel capacity C and the critical rate R_c as functions of θ . For A very small, the $E_L(R)$ curve approaches a limiting form. In fact, if $\theta = (\pi/2) - \epsilon$, with ϵ small, to a close approximation by obvious expansions we find

$$E_L(R) \doteq \frac{A^2}{2} - A\epsilon + \frac{\epsilon^2}{2} \quad \text{and} \quad R \doteq \frac{\epsilon^2}{2}.$$

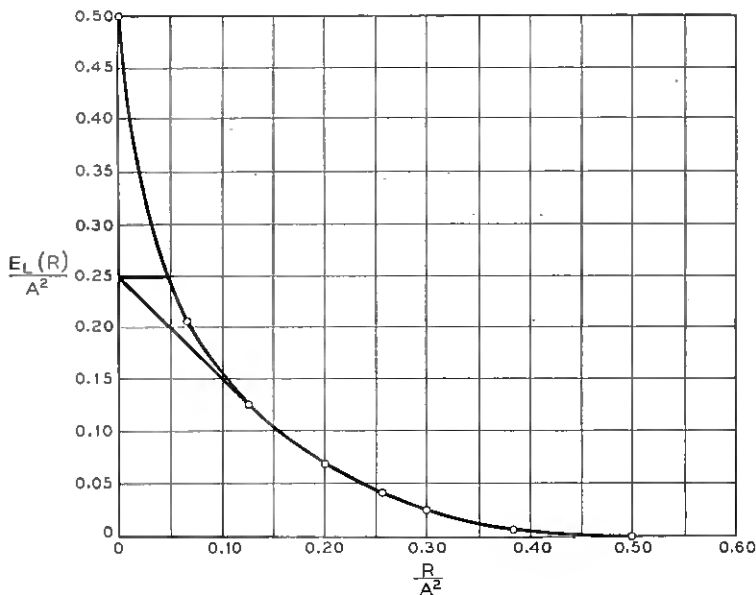


Fig. 13 — Plots of $E_L(R)/A^2$ against R/A^2 .

Eliminating ϵ , we obtain

$$\frac{E_L(R)}{A^2} \doteq \frac{1}{2} - \sqrt{\frac{2R}{A^2}}$$

Fig. 13 plots $E_L(R)/A^2$ against R/A^2 .

XV. ACKNOWLEDGMENTS

I am grateful to several people for help in preparing this paper. Mrs. Judy Frankman computed the curves of E_L and other members of the Center for Advanced Study in the Behavioral Sciences were helpful in many ways. The referee made several valuable suggestions which have been incorporated in the paper. Finally, I am particularly indebted to my wife Betty for checking much of the algebra involved in the asymptotic bounds.

REFERENCES

1. Shannon, C. E., Communication in the Presence of Noise, Proc. I.R.E., **37**, January 1949, p. 10.

2. Rice, S. O., Communication in the Presence of Noise — Probability of Error for Two Encoding Schemes, *B.S.T.J.*, **29**, January 1950, p. 60.
3. Elias, P., in *Information Theory* (Cherry, C., ed.), Academic Press, New York, 1956.
4. Shannon, C. E., Certain Results in Coding Theory for Noisy Channels, *Inform. and Cont.*, **1**, September 1957, p. 6.
5. Johnson, N. L. and Welch, B. L., Applications of the Noncentral t -Distribution, *Biometrika*, **31**, 1939, p. 362.
6. David, H. T. and Kruskal, W. H., The WAGR Sequential t -Test Reaches a Decision with Probability One, *Ann. Math. Stat.*, **27**, September 1956, p. 797

Analysis of Phonon-Drag Thermomagnetic Effects in n-Type Germanium*

By C. HERRING, T. H. GEBALLE and J. E. KUNZLER

(Manuscript received November 12, 1958)

A study has been made of the Nernst effect and the variation of thermoelectric power with magnetic field for single-crystal samples of n-type germanium of various orientations and impurity concentrations, at fields up to 18,000 gauss and temperatures from 275° to 60°K and below. Except at the highest temperatures, both effects arise predominantly from that part of the thermoelectric power which is due to phonon drag. All observations can be quantitatively accounted for by theory. They yield information about the dependence of the relaxation times for phonon-phonon scattering on the frequency of the phonons, and establish with some certainty the conclusion that four-phonon collisions are much less important than three-phonon collisions in the pertinent range of temperatures and phonon frequencies. Auxiliary investigations have shown that the quantization of electron orbits in a magnetic field has only a small effect on transport properties when the cyclotron level spacing is less than thermal energy. The mean free path of electrons is shown to be energy-independent, as acoustic-scattering theory predicts. The absolute mobility can be predicted to within 10 per cent or better from data on the fractional changes of resistance with stress and magnetic field.

A more detailed summary of the conclusions and implications of the present work is given in Section IX.

I. INTRODUCTION

In a recent paper¹ we have presented measurements of the Nernst effect and the change of thermoelectric power in a magnetic field for n-type germanium of high purity. These two effects were shown to be

* This paper, though complete in itself, constitutes part II of a study, part I of which was presented in Ref. 1. In addition to describing the phenomena and the physical principles underlying them, Ref. 1 develops the background of the present study and gives details of the measuring techniques.

compounded out of electron-diffusion and phonon-drag contributions, the latter being the predominant one in the range near liquid-air temperature. Now the phonon-drag phenomenon² — the pushing of charge carriers from hot to cold by the asymmetric phonon distribution which a thermal gradient produces — depends on the details of both the electron-phonon interaction and the processes which scatter phonons. The anisotropies of these interactions and their dependences on wave number enter in different ways into the several thermomagnetic quantities which one can measure. Thus, an analysis of thermomagnetic data should make it possible to sort out the different factors involved in phonon drag, and to obtain information not previously available about electron-phonon and phonon-phonon interactions. In our first paper we used qualitative arguments to draw a number of semiquantitative conclusions about these interactions. Our object in the present paper is to make these conclusions as quantitative as possible, by comparing the observations with explicit theoretical formulas.

1.1 Basic Concepts

Following our previous practice,¹ we shall formulate the theory in terms of the Peltier tensor $\mathbf{\Pi}$ of the electrons, which is related to the thermoelectric power tensor \mathbf{Q} in a magnetic field \mathbf{H} by

$$\mathbf{Q}_{\alpha\beta}(\mathbf{H}) = \frac{\mathbf{\Pi}_{\beta\alpha}(-\mathbf{H})}{T}. \quad (1)$$

Here, $\mathbf{\Pi}$ is defined as the tensor relating the energy flux \mathbf{F} , relative to the Fermi level, to the current density \mathbf{j} :

$$\mathbf{F} = \mathbf{\Pi} \cdot \mathbf{j}. \quad (2)$$

In terms of the \mathbf{Q} -tensor, the change of thermoelectric power in the α direction, due to \mathbf{H} , is

$$\Delta Q = Q_{\alpha\alpha}(\mathbf{H}) - Q_{\alpha\alpha}(0). \quad (3)$$

The Nernst coefficient $B(\mathbf{H})$ is defined — for the symmetrical orientations considered here — by

$$\mathbf{E}_N = -B\mathbf{H} \times \nabla T, \quad (4)$$

where \mathbf{E}_N is the open-circuit field transverse to the temperature gradient. We shall obtain B from (1) and the relation

$$BH = -Q_{yz}, \quad (5)$$

which applies if ∇T is in the x -direction and \mathbf{H} in the z -direction.

The model which we shall use as a basis for our calculations is the "electron-group" model described in our previous paper.¹ In this, the possible states of motion of the charge carriers are divided into a number of groups, g , and the assumption is made that for each such group the distribution function over the group is determined by the contribution, \mathbf{j}_g , which the states of this group make to the current density. Carriers of all groups are, of course, subject to the same electric and magnetic fields.

For a multivalley semiconductor like n germanium, a natural choice is to let a group g consist of the states in an ellipsoidal shell of the Brillouin zone, of energy range ϵ to $\epsilon + d\epsilon$, in a particular valley. Each such group g has its characteristic Peltier tensor Π_g ; this must be independent of \mathbf{H} if, as we are assuming, \mathbf{H} affects the \mathbf{j}_g 's of the different groups but not the distribution function for a given \mathbf{j}_g . The total Peltier tensor of the medium can be compounded out of those of the different groups. Explicitly,

$$\Pi(\mathbf{H}) = \Sigma_g \Pi_g \cdot \delta_g(\mathbf{H}) \cdot \rho(\mathbf{H}), \quad (6)$$

where δ_g is the conductivity tensor of group g and $\rho = (\Sigma_g \delta_g)^{-1}$ is the total resistivity tensor.

For any one of the ellipsoidal shells just mentioned, the Peltier tensor Π_g must have the symmetry of the valley, which for n germanium is axial symmetry about a [111]-type direction. Thus, Π_g is describable in terms of its two principal components, $\Pi_{\parallel}(\epsilon)$, along this direction (the high-mass direction), and $\Pi_{\perp}(\epsilon)$, normal to it. Each of these components, in turn, is a sum of an electron-diffusion term (subscript e) and a phonon-drag term (subscript p). Here, $\Pi_{e\parallel} = \Pi_{e\perp}$ and is a linear function of ϵ , while $\Pi_{p\parallel}(\epsilon) \neq \Pi_{p\perp}(\epsilon)$. Our task is to relate these functions of energy explicitly to the low- and high-field values of the Nernst coefficient B and the change ΔQ in the thermoelectric power.

1.2 Qualitative Properties of $\Pi_{p\parallel}$ and $\Pi_{p\perp}$

As we formulate the detailed theory in the following sections, it will be helpful to keep in mind two qualitative conclusions which can be drawn from our earlier study.¹ The first of these concerns the anisotropy ratio of the phonon-drag Peltier coefficient. It is that

$$\frac{\Pi_{p\parallel}(\epsilon)}{\Pi_{p\perp}(\epsilon)} \gg 1 \quad \text{but} \quad < \frac{m_{\parallel}^* \tau_{\perp}}{m_{\perp}^* \tau_{\parallel}} \approx 17, \quad (7)$$

where m_{\parallel}^* , m_{\perp}^* are the principal effective masses and τ_{\parallel} , τ_{\perp} are the

acoustic-scattering relaxation times³ for current contributions in the corresponding directions. A large anisotropy of this sort is to be expected, since, if all the phonon modes had the same velocity and relaxation time τ_{ph} , the ratio $\Pi_{p\parallel}/\Pi_{p\perp}$ would be the ratio of the rates of crystal-momentum loss to the lattice for unit currents in the high-mass and low-mass directions, and this is just $m^*_{\parallel}\tau_{\perp}/m^*_{\perp}\tau_{\parallel}$.

Our second conclusion had to do with the energy dependence of $\Pi_{p\parallel}$ and $\Pi_{p\perp}$. All the quantities which we have measured correspond to various kinds of weighted averages of Π_{\parallel} and Π_{\perp} . A typical such average of $\Pi_{p\parallel}$ and $\Pi_{p\perp}$, considered as a function of energy ϵ , decreases as ϵ increases, at a rate distinctly slower than $\epsilon^{-\frac{1}{2}}$. This was shown to correspond to a frequency dependence of the relaxation time τ_{ph} of the phonon modes intermediate between ω^{-1} and ω^{-2} and nearer to the former.

1.3 Program

In Section II and Appendix A we shall discuss critically the adequacy of the electron-group model, reporting several auxiliary experimental and theoretical investigations which bear on this question. After establishing that the model should be quite good for pure n germanium at liquid-air temperature and above, we shall describe in Section III the procedures used for deriving the formulas for B and ΔQ . The explicit derivations will be given in Appendix B in a form applicable to any multivalley medium (e.g., n silicon); the formulas we need for the present application will be summarized in Section IV. In Section V we shall present the raw experimental data and discuss the corrections which need to be applied to them before comparing them with the theory. In Sections VI and VII we shall make the comparison for temperatures in the range 60° to 131°K, with the object of deducing from it as much information as possible about the functions $\Pi_{p\parallel}(\epsilon)$ and $\Pi_{p\perp}(\epsilon)$. In Section VIII we shall discuss various further observations, such as the behavior of B at high fields and the thermomagnetic effects above 200°K. These observations, which for experimental or theoretical reasons are less amenable to precise analysis, still are in good accord with the conclusions previously drawn, and allow them to be slightly extended. Section IX will summarize the arguments and give a few remarks on the significance of the conclusions for the theories of electron-phonon and phonon-phonon interactions, topics which we hope to treat more fully in later publications.

II. CRITIQUE OF THE ELECTRON-GROUP MODEL

For a multivalley semiconductor the electron-group model, described above and in our earlier paper,¹ is based on two assumptions: (a) that the distribution function in any particular ellipsoidal energy shell can be approximated by a linear function of the crystal momentum over the shell and (b) that this distribution function is determined only by the external fields, being independent of the distribution functions of other shells. We must first examine critically the adequacy of this model; then we shall consider further specializations of it.

The adequacy of the electron-group model depends on the answers to the following questions:

- i. Is it legitimate to use a transport equation in crystal-momentum space?
- ii. How accurately can one separate the transport problems for different energy shells?
- iii. How well is the distribution function in a single shell approximated by a linear function of the crystal momentum?

As regards i, there are two things which might undermine the validity of a transport equation in crystal-momentum space alone — inhomogeneities and orbital quantization. Distortion of the current distribution by inhomogeneities in the specimen may make it necessary to consider the transport problem in position as well as in velocity space. This trivial but annoying difficulty in the interpretation of experiments will be discussed at length by one of the authors in another place;⁴ it causes the transverse magnetoresistance to fail to saturate as $H \rightarrow \infty$, but has very little effect on thermoelectric and thermomagnetic properties. The reason for the difference is simple. When a current flows, portions of the specimen with different carrier concentrations try to set up different Hall fields, and can only adjust to one another by a distortion of the current lines, which becomes large as $H \rightarrow \infty$. Thermoelectric and thermomagnetic effects, on the other hand, are measured at zero average current, and set up fields which are much less sensitive to inhomogeneities than the Hall field is. The fields also remain finite as $H \rightarrow \infty$, whereas the Hall field does not. For this reason, and because measurements on different pairs of side-arms showed the gross homogeneity of our specimens to be quite good, we believe the effect of inhomogeneity on our results to be negligible.

A less trivial objection to the use of a transport equation in crystal-momentum space is the fact of orbital quantization, which is known to make the high-field magnetoresistance qualitatively different from the

predictions of a transport theory formulated in crystal-momentum space.⁵⁻¹¹ It can be argued theoretically that this is not serious until the spacing, $\hbar\omega_c$, of the cyclotron levels becomes comparable with kT ; however, for n germanium with 18,000 gauss along a [100] axis, $\hbar\omega_c$ corresponds to kT at about 25°K, and one may wonder whether this is negligible for experiments at 77°. Fortunately, theoretical and empirical evidence, which we have gathered together in Appendix A, indicates that for our specimens the effect of orbital quantization on the resistance and thermoelectric power is probably not more than 5 per cent or so at 18,000 gauss and 77°, and that for many orientations it is probably considerably less. The main arguments can be summarized as follows:

1. The theory of longitudinal magnetoresistance,⁷ when interpreted with allowance for the energies of the phonons which do the scattering,⁸ predicts only rather small departures from the unquantized theory as long as $\hbar\omega_c < kT$. Observations¹² support this conclusion.

2. Both for longitudinal and transverse cases the phonon-drag part Π_p of the Peltier coefficient is proportional to the product of the resistivity, the average relaxation time of the phonon modes which scatter the electrons, and the average squared velocity of these modes. As the effect of orbital quantization on the average relaxation time and velocity is probably rather less than on the resistivity, the quantization effect should be of similar magnitude for electrical and for thermoelectric measurements.

3. Although the theory for transverse magnetoresistance is more complicated than for longitudinal magnetoresistance, a combination of reasonable arguments with empirical evidence yields a rough upper bound for the orbital quantization effect. Contradictory approaches to the theory have been elaborated for the case of acoustic scattering by Klinger and Voronyuk,^{9*} Argyres¹⁰ and Wolff;¹¹ the resulting formulas can be evaluated for high fields ($\hbar\omega_c \gg kT$), but become very cumbersome at lower fields. A plausible presumption is that for constant carrier concentration and in the limit of small scattering (Hall angle practically $\pi/2$) the departure from the resistivity given by the older transport theory varies at least as fast as $\hbar\omega_c/kT$, as long as this quantity is rather less than unity. The size of the factor of proportionality cannot be estimated reliably from the nonsaturation of transverse magnetoresistance,

* These authors evaluate their expressions only for the limiting case of fields so large that the energy of an average phonon causing transitions of the electrons is $\gg kT$. However, it is easy to show that, for the case $\hbar\omega_c/kT \gg 1$ and phonon energies negligible, this approach yields a magnetoresistance that goes very nearly as H^2 . (Note added in proof: We are indebted to E. N. Adams and T. Holstein for communicating to us an as-yet-unpublished analysis that shows the approach of Ref. 9 to be the correct one.)

since this can be shown to arise in large part from other effects. However, the factor can be estimated from thermoelectric data at liquid-hydrogen and intermediate temperatures.

4. The fair agreement of theory and experiment on high-field transverse magnetoresistance would be spoiled if orbital quantization increased the resistivity by any sizable amount. The resulting discrepancy would be hard to explain away, as most perturbing effects (inhomogeneities, surface conduction, etc.) also increase the magnetoresistance.

Turning to question ii at the beginning of this section, we must consider scattering processes which change the energy. In highly doped material the most important processes of this sort are electron-electron collisions. There is also another effect: at carrier concentrations high enough to make the low-frequency phonons have an effective drift velocity comparable with that of the electrons — the so-called saturation effect² on Π_p — the contribution of any electron group to Π is reduced by the unbalance of the distribution function of the phonons with which it interacts. Since this unbalance is influenced by the interactions of these phonons with all other electron groups, it is no longer possible at such concentrations to separate $\Pi \cdot \mathbf{j}$ into a sum of $\Pi_g \cdot \mathbf{j}_g$. At low carrier concentrations, however, both these effects become negligible, and the only type of energy change on scattering which must be considered is that due to the finiteness of the energies of the phonons. Only acoustical phonons need to be investigated, since scattering by intervalley or non-polar optical modes occurs with almost equal probability to all final states in an energy shell,¹³ and therefore is describable by a relaxation time dependent only on the initial state.

The most detailed study which has appeared on the effect of acoustic-phonon energies is that of Dorn.¹⁴ He estimated the effect on the distribution function for a simple-model semiconductor in an electric field, using both a high-temperature expansion method and a variational method. He found that, as one expects, the fractional alteration in the mobility due to the phonons having finite rather than zero energy is very small in the temperature range of interest to us. The alteration is of the order of m^*c^2/kT , where m^* is the effective mass and c the velocity of sound. With a reasonable average m^* for n germanium, m^*c^2/k is only a small fraction of a degree. However, Dorn's expressions show that the fractional alteration in the current of a shell of energy ϵ is of order m^*c^2/ϵ for $\epsilon \ll kT$ but $\gg m^*c^2$. Since for pure lattice scattering the low-field magnetoresistance and ΔQ effects are dominated by the behavior of low-energy carriers, the effects of the phonon energies may be expected to be considerably more serious for these than for the mobility in the ab-

sence of a magnetic field. A quantitative estimate of these effects in a magnetic field and with anisotropic masses would be of some interest. However, it seems likely that they are still small. The fractional alteration in the H^2 term in the current of a shell should still be of order m^*c^2/ϵ , although perhaps with a slightly larger coefficient than for Dorn's case. Since, for a shell ϵ to $\epsilon + d\epsilon$ containing dn carriers, this H^2 term goes as $\epsilon r^3 dn \propto d\epsilon$ for $\epsilon \ll kT$, the phonon-energy correction should involve $\int d\epsilon/\epsilon$. Thus, we guess that

$$\text{fractional correction to magnetoresistance} \approx \frac{m^*c^2}{kT} \ln \frac{kT}{m^*c^2}. \quad (8)$$

Since the phonon-drag Peltier coefficients of the different shells do not vary much with energy, (8) should apply to the thermoelectric power also. For n germanium in the range near liquid-air temperature, the value of (8) is at most a few per cent.

We can now give a prescription for answering question iii, the last of the three posed at the start of this section. This question had to do with the legitimacy of approximating the distribution function in each energy shell by a linear function of crystal momentum. Since questions i and ii have by now received favorable answers, we need only to solve the transport equation for an energy shell under the influence of arbitrary electric and magnetic fields, and see how well the current and Peltier-heat contributions from the accurate solution agree with those from the electron-group approximation. Now the latter approximation amounts to retaining only $l = 1$ spherical harmonics in the distribution function in a new set of variables which take the ellipsoidal energy surfaces into spheres³ — variables which we have used for other purposes in Appendix A. The terms in the scattering operator which mix $l = 1$ with higher l values — call these S_{ll} — can be shown³ to affect the electric current only to the second order, i.e., $|S_{ll}|^2$. They may, however, affect the Peltier flux to order $S_{ll}T_l$, where T_l measures the Peltier flux which would be produced by a distribution with the l value in question. Now the longitudinal and transverse branches individually contribute quite anisotropic scattering, but their combined effect in germanium is almost isotropic, so the S_{ll} for $l \neq 1$ are small. But since the relative contributions of the two branches to the phonon-drag Peltier effect are not known, it is quite conceivable that the $S_{ll}T_l$ for $l \neq 1$ are not negligible, even though the $|S_{ll}|^2$ are. A detailed investigation of this question using deformation-potential theory has been made for the case $l = 3$ and, to some extent, for larger l , and will be reported elsewhere.¹⁵ The result is that, for germanium, both the S_{ll} and the T_l are quite small, and that the retention

of only the $l = 1$ term in the distribution function should be a very good approximation for all conduction and thermoelectric phenomena, with or without a magnetic field.

III. PROCEDURE FOR COMPUTING B AND ΔQ

Having satisfied ourselves that the electron-group model described in Section I is a good approximation for the cases we wish to analyze, let us now use this model to evaluate explicitly the expressions for B and ΔQ which were given there. In both cases the task is to evaluate (6). This determines the Nernst coefficient B , or the Nernst field (4), via (1) and (5). Similarly, it determines ΔQ via (1) and (3). Since each group g consists of the set of states in the energy range ϵ to $\epsilon + d\epsilon$ in some valley i , it is appropriate to replace the conductivity σ_g of this group by the infinitesimal $d\sigma^{(i)}$, and to rewrite (6) in the form

$$\Pi(\mathbf{H}) = \sum_i \left[\int \Pi^{(i)}(\epsilon) \cdot d\sigma^{(i)}(\mathbf{H}) \right] \cdot \rho(\mathbf{H}), \quad (9)$$

where the tensor $\Pi^{(i)}(\epsilon)$ has principal components $\Pi_{\parallel}(\epsilon)$ and $\Pi_{\perp}(\epsilon)$ in the principal-axis system of valley i , and where the integration is really over energies ϵ . The tensor $d\sigma^{(i)}(\mathbf{H})$ can be obtained by solving the transport equation relating the isothermal current $d\mathbf{j}^{(i)}$ in the energy shell to the electric field \mathbf{E} producing it, and expressing the solution in the form

$$d\mathbf{j}^{(i)} = d\sigma^{(i)} \cdot \mathbf{E}.$$

For the case of nondegenerate statistics, which we shall assume throughout, the transport equation for $d\mathbf{j}^{(i)}$ is easily written down, in the electron-group approximation, by equating the rates of gain and loss of crystal momentum in the shell. Anisotropic scattering can be taken into account³ by assigning to the shell a relaxation-time tensor τ , with principal components $\tau_{\parallel}(\epsilon)$, $\tau_{\perp}(\epsilon)$ along and normal to the axis of the valley. The resulting transport equation — in which, for simplicity, we omit the valley suffix i — is³

$$\tau^{-1} \cdot \mathbf{m}^* \cdot d\mathbf{j} \pm \left(\frac{e}{c} \right) d\mathbf{j} \times \mathbf{H} = e^2 \left(\frac{\epsilon}{\langle \epsilon \rangle} \right) \mathbf{E} dn, \quad (10)$$

where dn is the number of carriers in the shell, \mathbf{m}^* is the effective-mass tensor and $\langle \epsilon \rangle = (3kT)/2$ is the mean energy, and where the upper sign is for electrons, the lower for holes. For small H this can be solved by iteration to give the contributions (48) and (49) of Appendix B to the

conductivity tensor $d\sigma$. As $H \rightarrow \infty$ the part of \mathbf{E} normal to \mathbf{H} becomes interpretable as the Hall field, and the asymptotic value of the part of $d\mathbf{j}$ normal to \mathbf{H} can be determined by taking the cross product of (10) with \mathbf{H} and throwing away the first term on the left. This part of $d\mathbf{j}$ is of order H^{-1} ; the H^{-1} part of $d\mathbf{j}$ parallel to \mathbf{H} can be determined from it by dotting (10) with \mathbf{H} . Higher powers in the expansion of $d\mathbf{j}$ in H^{-1} can be obtained by further iteration. The resulting contributions to $d\sigma$ are given in (68) through (72) of Appendix B.

The resistivity tensor $\rho(\mathbf{H})$ occurring in (9) is just the reciprocal of $\Sigma_i \int d\sigma^{(i)}(\mathbf{H})$. A few properties of ρ and $d\sigma^{(i)}$ are worth listing for reference. With superscripts 0, 1, 2, \dots to denote coefficients of 1, H , H^2 , \dots in the expansion of ρ or $d\sigma$ in powers of H , and $+1$, ∞ , -1 , \dots to denote coefficients of H , 1, H^{-1} , \dots in the expansion in powers of H^{-1} , we have, for a cubic crystal,

$$\rho_{\alpha\beta}^{(1)} = -R(0)\Sigma_\gamma \delta_{\alpha\beta\gamma} H_\gamma, \quad (11)$$

$$\rho_{\alpha\beta}^{(+1)} = -R(\infty)\Sigma_\gamma \delta_{\alpha\beta\gamma} H_\gamma, \quad (12)$$

where $R(0)$ and $R(\infty)$ are the limiting values of the Hall constant at $H = 0$ and $H = \infty$, respectively, and where $\delta_{\alpha\beta\gamma} = \pm 1$ when $\alpha\beta\gamma$ is an even (odd) permutation of 123, zero otherwise. For n even (odd), $\rho^{(n)}$ and $d\sigma^{(n)}$ are symmetrical (antisymmetrical). The symmetrical tensor $\rho^{(\infty)}$ is finite and positive definite, but the only nonvanishing component of $d\sigma^{(\infty)}$ is the HH component.

The procedure of Appendix B is thus to get the high- and low-field expansions of $d\sigma^{(i)}(\mathbf{H})$ by solving (10); to use these, when necessary, to get the corresponding expansions of $\rho(\mathbf{H})$; to substitute into (9) to get the expansions for various components of $\Pi(\mathbf{H})$; and thence to compute the high- and low-field B and ΔQ for various orientations. The results, containing integrals over $\Pi_{\parallel}(\epsilon)$ and $\Pi_{\perp}(\epsilon)$, can be expressed in a great variety of forms. Although in Appendix B expressions are obtained for B and ΔQ valid for arbitrary functional forms of Π_{\parallel} , Π_{\perp} , and for valleys of either the [111] or [100] types in a cubic crystal, these expressions are cumbersome in their most general forms. These forms involve Maxwellian averages of quantities $\epsilon \Pi_{\parallel, \perp} f^{(n)}(\tau_{\parallel}, \tau_{\perp})$, where $f^{(n)}$ is a homogeneous function of degree n in its arguments. These forms therefore simplify greatly if the Π 's and τ 's are each proportional to some power of the energy ϵ . Now $\Pi_{e\parallel} = \Pi_{e\perp}$ is linear in ϵ , while our previous analysis¹ indicates that $\Pi_{p\parallel}$ and $\Pi_{p\perp}$ probably vary at a rate between constancy and $\epsilon^{-\frac{1}{2}}$. Moreover, for ideal acoustic scattering and negligible phonon energies τ_{\parallel} and τ_{\perp} are $\propto \epsilon^{-\frac{1}{2}}$. Thus, specialization of the formulas to power-law dependences of the Π 's and τ 's can give a fair picture of

the behavior of Π_e and Π_p separately, although not as accurate a picture as we shall need for our final analysis. We shall give such a specialization in the next section, before discussing more complete formulas.

IV. COMPILATION OF FORMULAS

Tables I through IV give formulas for high- and low-field B_e , B_p , ΔQ_e and ΔQ_p , as obtained by the procedure just described for various cases corresponding to energy-independent anisotropies. The formulas are for [111] valleys and are specialized to the cases $\Pi_{p\parallel,\perp} \propto \epsilon^{-1}$ or independent of ϵ , and $\tau_{\parallel,\perp} \propto \epsilon^{-1}$ or independent of ϵ . Comparison of the latter two alternatives makes possible a crude estimate of the effect of impurity scattering. Although these special cases do not give as accurate a representation of the data as we shall ultimately wish to use, they do give a fair representation. They are worth listing because their comparative simplicity makes it easy to see how sensitive the various measurable quantities are to the assumed anisotropies and energy dependences. All the tables contain references to equations of Appendix B. Most of these equations are more general expressions for the quantities tabulated, valid even when the anisotropies are not energy-independent. Also, they are applicable to multivalley band structures different from that of germanium.

The formulas for the low-field Nernst coefficient (see also Equation B.4 of Ref. 1) are most conveniently expressed in terms of the dimensionless coefficients ζ_e and ζ_p , defined by

$$B_e = \zeta_e \left(\frac{k}{e} \right) \left(\frac{\mu_H}{c} \right), \quad (13)$$

$$B_p = \zeta_p |Q_p| \left(\frac{\mu_H}{c} \right),$$

where, as always, B_e and B_p represent, respectively, the electronic and phonon-drag contributions to B . Table I gives the expressions obtained in Appendix B for ζ_e and ζ_p for the case of energy-independent anisotropies. Note that $(1 + \zeta_p)$ is the product of a factor dependent only on the energy variations of the Π 's and τ 's by a factor dependent only on their anisotropies. The anisotropy factor is unity when the tensor $\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1}$ is isotropic, and likewise when the tensor $\mathbf{\Pi}_{p0}$ is isotropic ($\Pi_{p\parallel} = \Pi_{p\perp}$); the energy factor is unity if either $\boldsymbol{\tau}$ or $\mathbf{\Pi}_{p0}$ is independent of energy. Fig. 1 shows the dependence of the anisotropy factor in $(1 + \zeta_p)$ on $p = \Pi_{p\parallel}/\Pi_{p\perp}$. Curves are drawn for several values of the anisotropy w of $\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1}$, chosen to encompass the ranges likely to occur for germanium and silicon.

TABLE I

Values of the dimensionless coefficients ζ_e , ζ_p , related by (13) to the electronic and phonon-drag parts of the low-field Nernst coefficient, respectively. The expressions given apply when the ratios $w = \tau_{\perp} m_{\perp}^* / \tau_{\parallel} m_{\parallel}^*$ and $p = \Pi_{p\parallel} / \Pi_{p\perp}$ are independent of energy. They are valid for any type of valleys in a cubic crystal (e.g., n silicon or n germanium). Angular brackets represent Maxwellian averages, i.e., averages with weight $\epsilon^{1/2} \exp(-\epsilon/kT)$. For τ , one may take either τ_{\parallel} or τ_{\perp} , and, for Π_p , either $\Pi_{p\parallel}$ or $\Pi_{p\perp}$.

	ζ_e	$1 + \zeta_p$
General formula	$\frac{3}{2} \left[\frac{\langle \epsilon^2 \tau^2 \rangle}{\langle \epsilon \tau^2 \rangle \langle \epsilon \rangle} - \frac{\langle \epsilon^2 \tau \rangle}{\langle \epsilon \rangle \langle \epsilon \tau \rangle} \right]$	$\frac{\langle \epsilon \Pi_p \tau^2 \rangle \langle \epsilon \tau \rangle}{\langle \epsilon \Pi_p \tau \rangle \langle \epsilon \tau^2 \rangle} \left[\frac{(2+w)(1+w+pw)}{(1+2w)(2+pw)} \right]$
Reference, Appendix B	(57), (50)	(57), (51)
Value for: $\tau \propto \epsilon^{-1/2}$, $\Pi_p \propto \epsilon^{-1/2}$ τ independent of ϵ , any $\Pi_p(\epsilon)$ Any $\tau(\epsilon)$, Π_p independent of ϵ	$-\frac{1}{2}$ 0 —	$4/\pi$ [as above] 1 [as above] 1 [as above]

Table II gives the part of BH going as H^{-1} when $H \rightarrow \infty$, for certain special directions of \mathbf{H} and \mathbf{j} , again as obtained by specializing the formulas of Appendix B to the case of [111] valleys and energy-independent anisotropies. Note that at high fields the electronic contribution to the Nernst field is anisotropic and depends on the anisotropy w of $m^* \cdot \tau^{-1}$, whereas at low fields it does not. Fig. 2 shows the dependence of the phonon-drag contribution on $p = \Pi_{p\parallel} / \Pi_{p\perp}$, for two assumptions about the energy dependence of $\Pi_{p\parallel, \perp}$. Several features of the graphs and formulas are worth noting:

i. The high-field Nernst coefficient is very much larger when \mathbf{H} is in an [011]-type direction than in an [001]-type. This could have been anticipated from the familiar difference in rapidity of saturation of the Hall coefficient in these two directions, and the close relationship of the Nernst and Hall effects.¹ The cause in both cases is the fact that, for \mathbf{H} along [011], two of the four valleys have a high-mass direction exactly normal to \mathbf{H} , and so require large fields to obtain Hall angles near $\pi/2$.

ii. The theoretical behavior of the Nernst coefficient differs from that of the Hall coefficient in that, when \mathbf{H} is along [011], the value of B is

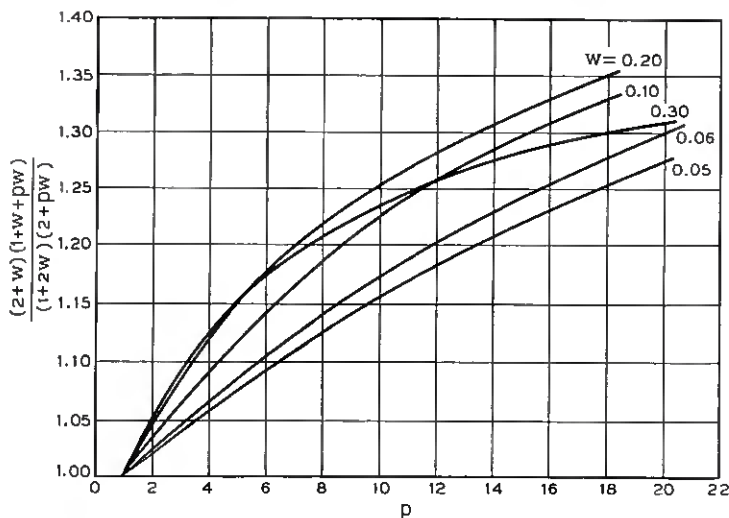


Fig. 1 — Dependence of the anisotropy factor in the expression of Table I for $1 + f_p$ on $p = \Pi_{\parallel} / \Pi_{\perp}$, for various values of $w = m_{\perp}^* \tau_{\parallel} / m_{\parallel}^* \tau_{\perp}$.

slightly different for ∇T along [100] and [01 $\bar{1}$]. Equality of Hall coefficients for these two directions of \mathbf{j} is required by the Onsager relations (principle of microscopic reversibility), but equality of the Nernst coefficients is not.¹ For the phonon-drag B_p the theoretical difference is very small when $\Pi_{\parallel} \gg \Pi_{\perp}$, but for the electronic B_e it is sizable, e.g., a factor 1.62 for acoustic scattering and $w = 0.06$.

iii. From Tables I and II it follows that the ratio of B_p to B_e is not usually the same at high fields as at low. For most specimens, therefore, there should be a range of temperatures — near that for which $B_p + B_e = 0$ — where the sign of the Nernst voltage will change with increase of H . Substitutions in the formulas give the result that, for any values of the anisotropies in the neighborhood of those obtaining for pure germanium ($p \approx 10$, $w \approx 0.06$), the ratio $|B_p/B_e|$ is greater at high fields than at low fields whenever H is along [011]. When H is along [001], this ratio is less at high fields than at low if $\Pi_{\parallel, \perp}$ increase with decreasing energy, but becomes the same at high and low fields if $\Pi_{\parallel, \perp}$ are independent of energy.

The electronic part of the thermoelectric power, $Q_e(\mathbf{H})$, behaves very simply as $H \rightarrow \infty$. When \mathbf{H} and \mathbf{j} are parallel, the energy distribution of the current is the same at $H = \infty$ and $H = 0$. So, for $\mathbf{H} \parallel \nabla T$,

$$Q_e(\infty) = Q_e(0) = \mp \frac{3}{2} \left(\frac{k}{e} \right) \frac{\langle \epsilon^2 \tau \rangle}{\langle \epsilon \rangle \langle \epsilon \tau \rangle} \quad (14)$$

TABLE II

Asymptotic behavior of the Nernst voltage as $H \rightarrow \infty$ for certain special directions of \mathbf{H} and ∇T , as derived from (86) of Appendix B. Assumptions and notations are the same as for Table I except that the expressions apply only for [111] valleys (n-germanium). In addition, μ is the drift mobility for $H = 0$.

$H, \Delta T$	Quantity	Value if $\tau \propto \epsilon^{-1/2}$, $\Pi_p \propto \epsilon^{-1/2}$	Value if $\tau \propto \epsilon^{-1/2}$, Π_p independent of ϵ	Value if τ independent of ϵ , any $\Pi_p(\epsilon)$
[001], any $\perp \mathbf{H}$	$\frac{e}{k} \left[\frac{\partial(B_x H)}{\partial(c/\mu H)} \right]_{\infty}$ $ Q_p ^{-1} \left[\frac{\partial(B_y H)}{\partial(c/\mu H)} \right]_{\infty}$	$-\frac{16}{9\pi} \left[\frac{(2+w)^2}{3(1+2w)} \right]$ $\left[\frac{256}{81\pi^2} \frac{(2+w)^2(1+w+pw)}{(1+2w)^2(2+pw)} - \frac{8}{9\pi} \frac{(2+w)^2}{(1+2w)} \right]$	same as at left	0
[011], [100]	$\frac{e}{k} \left[\frac{\partial(B_x H)}{\partial(c/\mu H)} \right]_{\infty}$ $ Q_p ^{-1} \left[\frac{\partial(B_y H)}{\partial(c/\mu H)} \right]_{\infty}$	$-\frac{16}{9\pi} \frac{(2+w)(1+2w)}{9w}$ $\left[\frac{256}{729\pi^2} \frac{(2+w)(1+2w)(5+p+w+2pw)}{w(2+pw)} - \frac{8}{27\pi} \frac{(2+w)(1+2w)}{w} \right]$	same as at left	0
[011], [011]	$\frac{e}{k} \left[\frac{\partial(B_x H)}{\partial(c/\mu H)} \right]_{\infty}$ $ Q_p ^{-1} \left[\frac{\partial(B_y H)}{\partial(c/\mu H)} \right]_{\infty}$	$-\frac{16}{9\pi} \left[\frac{w^2+7w+1}{9w} \right]$ $\left[\frac{256}{729\pi^2} \frac{(w^2+7w+1)(2+w)(p+2)}{w(2+pw)} - \frac{8}{27\pi} \frac{(2+w)(pw^2+2pw+5w+1)}{w(2+pw)} \right]$	same as at left	0

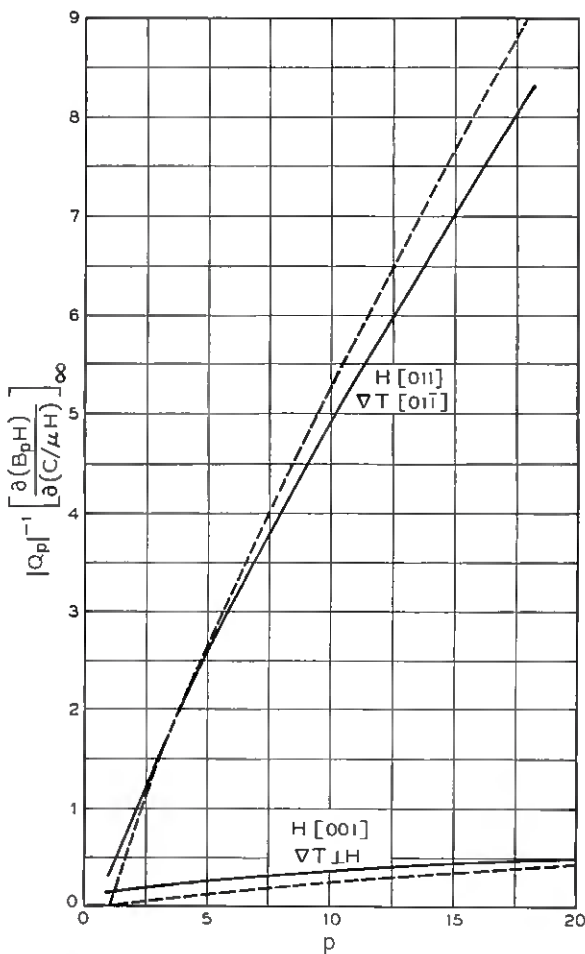


Fig. 2 — Dependence on $p = \Pi_{\parallel}/\Pi_{\perp}$ of the phonon-drag part of the Nernst coefficient at high magnetic fields, for two orientations of a cubic crystal with [111] valleys. All curves have been computed assuming $w = m_{\perp}^* \tau_0 / m_{\parallel}^* \tau_{\perp} = 0.06$, and pure acoustic scattering ($\tau \propto \epsilon^{-1/2}$). The full curves apply for $\Pi_p \propto \epsilon^{-1/2}$, the dashed for Π_p independent of energy. The curves are fairly insensitive to the energy dependence of τ , but depend strongly on w , the lower curves being for small w almost proportional to w , the upper to w^{-1} . For H along [011] and ∇T along [100] (not shown) the dashed curve would coincide with that for [011][011], while the full curve would differ by only a few per cent over most of the range shown.

TABLE III

Values of the saturation ratio $Q_p(H \rightarrow \infty)/Q_p(0)$, for special directions of \mathbf{H} and ∇T . Assumptions and notations are the same as for Table I, except that the expressions apply only for [111] valleys (n germanium). The energy factor $\langle \epsilon \Pi_p \rangle \langle \epsilon \tau \rangle / \langle \epsilon \rangle \langle \epsilon \Pi_p \tau \rangle$ has the value $8/(3\pi) = 0.849$ when $\Pi_p \propto \epsilon^{-1/2}$, and the value 1 when either Π_p or τ is independent of energy.

$\mathbf{H}, \nabla T$	Reference, Appendix B	Value, General
[001], [001]	(78)	$\frac{(2+w)(p+2)}{3(2+pw)}$
[011], [011]	(79)	$\frac{(2+w)(1+w+pw)}{(1+2w)(2+pw)}$
[111], [111]	(80)	$\frac{(2+w)(6+p+2pw)}{(7+2w)(2+pw)}$
[011], any $\perp \mathbf{H}$	(81)	$\frac{\langle \epsilon \Pi_p \rangle \langle \epsilon \tau \rangle}{\langle \epsilon \rangle \langle \epsilon \Pi_p \tau \rangle} \left[\frac{(2+w)(1+w+pw)}{(1+2w)(2+pw)} \right]$
[011], [100]	(81)	$\frac{\langle \epsilon \Pi_p \rangle \langle \epsilon \tau \rangle}{\langle \epsilon \rangle \langle \epsilon \Pi_p \tau \rangle} \left[\frac{5+w+p+2pw}{3(2+pw)} \right]$
[001], [011]	(81)	$\frac{\langle \epsilon \Pi_p \rangle \langle \epsilon \tau \rangle}{\langle \epsilon \rangle \langle \epsilon \Pi_p \tau \rangle} \left[\frac{(2+w)(p+2)}{3(2+pw)} \right]$
[011], [111]	(81)	$\frac{\langle \epsilon \Pi_p \rangle \langle \epsilon \tau \rangle}{\langle \epsilon \rangle \langle \epsilon \Pi_p \tau \rangle} \left[\frac{13+5w+5p+4pw}{9(2+pw)} \right]$

where $|\epsilon_b - \epsilon_f|$ is the distance of the Fermi level from the band edge, the upper sign is for electrons, the lower for holes, the angular brackets denote Maxwellian averages, and, for the present case of energy-independent anisotropies, τ may be taken as either τ_{\parallel} or τ_{\perp} (only the energy dependence matters). When \mathbf{H} and \mathbf{j} are perpendicular, the energy distribution of the current is the same as it would be for $H = 0$, $\tau = \text{constant}$. So, for $\mathbf{H} \perp \nabla T$,

$$Q_e(\infty) = \mp \left[\frac{5}{2} \left(\frac{k}{e} \right) + \frac{|\epsilon_b - \epsilon_f|}{eT} \right]. \quad (15)$$

Table III gives the formulas for $Q_p(\infty)$, the saturation value of the phonon-drag part of the thermoelectric power, again for [111] valleys and energy-independent anisotropies. The ratio $Q_p(\infty)/Q_p(0)$ is the product of a factor dependent on the energy variations of Π_{p0} and $m^* \cdot \tau^{-1}$ by a factor dependent on their anisotropies. Figs. 3 and 4 show how these anisotropy factors depend on $\Pi_{p\parallel}/\Pi_{p\perp}$, again for different assumptions

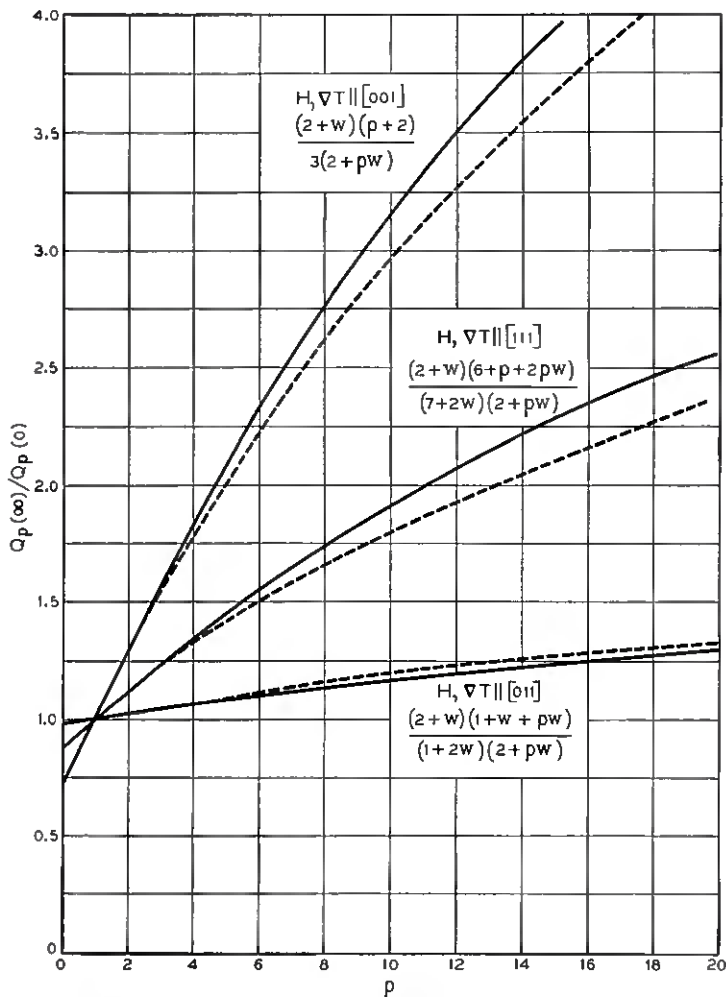


Fig. 3 — Dependence on $p = \Pi_{p\parallel}/\Pi_{p\perp}$ of the expressions of Table III for $Q_p(\infty)/Q_p(0)$, for \mathbf{H} parallel to ∇T in various directions of a cubic crystal with [111] valleys. The full curves apply for $w = m^*_{\perp}\tau_{\parallel}/m^*_{\parallel}\tau_{\perp} = 0.06$, the dashed curves for $w = 0.08$.

on the anisotropy of τ . For longitudinal effects ($\nabla T \parallel \mathbf{H}$) the energy factor is unity; thus the ratio for these cases gives information on the anisotropy of Π_{pp} independently of its energy variation. For the transverse effects the energy factor goes to unity if either Π_{p0} or τ is independent of energy.

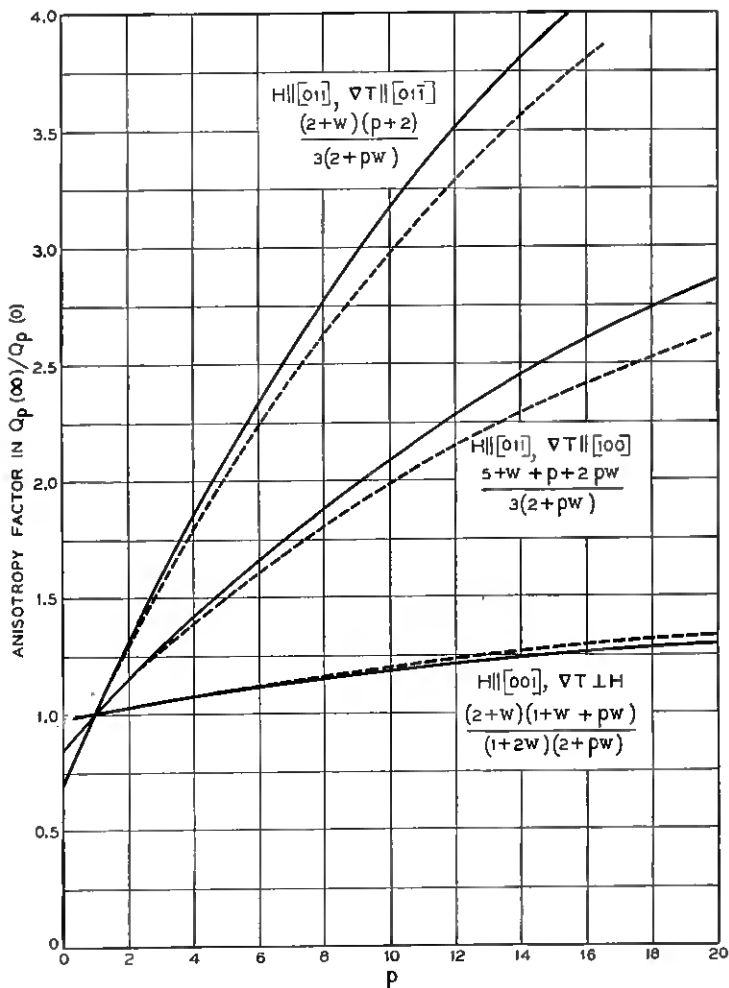


Fig. 4—Dependence on $p = \Pi_{p1}/\Pi_{p\perp}$ of the anisotropy factors in the expressions of Table III for $Q_p(\infty)/Q_p(0)$, for various orientations of a cubic crystal with [111] valleys, when \mathbf{H} is normal to ∇T . The full curves apply for $w = m^*_{\perp}\tau_{\perp}/m^*_{\parallel}\tau_{\parallel} = 0.06$, the dashed curves for $w = 0.08$. Note two of the factors are the same as for Fig. 3. No curves are given for $\mathbf{H} \parallel [011], \nabla T \parallel [1\bar{1}1]$, as this case is just one-third the way toward [011][100] from [011][011].

Table IV, finally, gives the expressions for the low-field limits of $\Delta Q_e/H^2$, $\Delta Q_p/H^2$. These have a less simple behavior than the high-field limits. However, it is noteworthy that, when \mathbf{H} is in the [001] direction and ∇T in the [100], the expression for $\Delta Q_p/Q_p$ involves the anisotropy of Π_p only through the same coefficient ζ_p as occurs in the expression (13) for B_p . The reason for this has been elucidated in our previous paper.¹

V. DATA AND CORRECTIONS

5.1 *Measurements and Extrapolations*

Measurements of electrical and thermomagnetic properties were made on a number of samples of n-type germanium having various orientations and dopings, at a number of temperatures from the liquid-hydrogen range to room temperature, and for magnetic fields up to 18,000 gauss. The samples were of the usual shape with three pairs of side-arms. Heat or electric current was fed in at the ends, the potential drop was measured between a side-arm of the first pair and one of the last pair, and the Hall or Nernst voltage was measured between the two middle side-arms. A full description of the experimental procedures has been given in our earlier paper.¹

Table V gives all the thermomagnetic quantities measured, as functions of magnetic field, for temperatures above 60°K; measurements were also made at a number of fields intermediate between those listed. Such data may be useful for comparison with future theoretical expressions for the various quantities in the intermediate-field range. In this paper, however, we have developed the theory only for the limiting cases of high and low fields, and so the most appropriate quantities for comparison with theory are the high- and low-field limits; these are given in Table VI. The limits of the various quantities as $H \rightarrow 0$ were obtained, since the quantities are even functions of H , by plotting the quantity in question against H^2 , fitting the low-field part of the plot to a parabola, and noting its intercept. The average size of the random errors involved in this procedure can be gauged from the sample plots given in Fig. 5. It will be noted that, at low fields, ΔQ is a little less accurately determinable than $\Delta\rho$, because of the smaller signals involved; similarly, the errors in B are larger than those in R . For the same reason, the random errors in the thermomagnetic quantities become larger, percentagewise, at higher temperatures. The extrapolations to $H = \infty$ were made in the same way, by plotting the quantities against $1/H^2$ and fitting with parab-

TABLE IV

Limiting low-field behavior of ΔQ_e , ΔQ_p , for various directions of \mathbf{H} and ∇T . Assumptions and notation are the same as for Table I, except that the expressions apply only for [111] valleys. The upper sign is for electrons, the lower for holes, and $\Delta\rho/\rho$ is the magnetoresistance.

Quantity	$\mathbf{H}, \nabla T$	Reference, Appendix B	Value if $\tau \propto \epsilon^{-1/2}$, $\Pi_p \propto \epsilon^{-1/2}$	Value if $\tau \propto \epsilon^{-1/2}$, Π_p independent of ϵ	Value if τ independent of ϵ , any $\Pi_p(\epsilon)$
$\frac{(e/k)\Delta Q_e}{(\mu_H H/c)^2}$	Parallel	(62)	$\mp \frac{\Delta\rho/\rho}{(\mu_H H/c)^2}$	same	0
	Perpendicular	(63)	$\mp \left[\frac{\Delta\rho/\rho}{(\mu_H H/c)^2} + \frac{1}{2} \right]$	same	0
$\frac{\Delta Q_p}{Q_p(\mu_H H/c)^2}$	[001], [001]	(67)	$\frac{8}{3\pi} \frac{(2+w)(1-w)(pw^2 + 5pw - 4w - 2)}{(1+2w)^2}$	$\frac{8}{\pi} \frac{w(2+w)(1-w)(p-1)}{(1+2w)^2}$	$\frac{2w(2+w)(1-w)(p-1)}{(1+2w)^2}$
	[001], any $\perp \mathbf{H}$	(66)	$\zeta_p - \frac{4}{3\pi} \frac{(2+w)^2}{(1+2w)}$	ζ_p	ζ_p

For any model the value of ΔQ in the direction of the unit vector \mathbf{u} is, to order H^2

$$\Delta Q = q_b H^2 + q_c (\mathbf{H} \cdot \mathbf{u})^2 + q_d (H_x^2 u_x^2 + H_y^2 u_y^2 + H_z^2 u_z^2),$$

where the coordinate axes are assumed oriented along cube-edge directions. For [111] valleys and low carrier concentration the additional relation $q_b = -q_c$ must hold whenever the electron-group approximation is valid (Ref. 1, Appendix C). Therefore, the values of ΔQ_p for any orientations can be expressed in terms of the two given above. For example,

$$\Delta Q_p \begin{matrix} 001 \\ 001 \end{matrix} = 2\Delta Q_p \begin{matrix} 011 \\ 011 \end{matrix} = 3\Delta Q_p \begin{matrix} 111 \\ 111 \end{matrix}, \quad \Delta Q_p \begin{matrix} 011 \\ 011 \end{matrix} = \Delta Q_p \begin{matrix} 001 \\ 100 \end{matrix} + \frac{1}{2} \Delta Q_p \begin{matrix} 001 \\ 001 \end{matrix}, \quad \text{etc.}$$

olas, as shown in Fig. 6. The legitimacy of assuming that ΔQ and $\Delta\rho$ are representable by power series in $1/H^2$ is borne out by the calculations of Appendix B, within the limitations imposed by the neglect of orbital quantization, inhomogeneities, etc. However, these latter factors are known to prevent $\Delta\rho$ from saturating completely as $H \rightarrow \infty$, so the extrapolations may be slightly different from the value which it is appropriate to compare with the present theory. We shall discuss this further below. Random errors in the high-field extrapolations are, however, inconsequential, as Fig. 6 shows.

A general discussion of the sources of experimental errors and the probable accuracy of the data has been given in our previous paper,¹ and will not be repeated here. However, there are some new observations which have an important bearing on the trustworthiness of the results, and in the following paragraphs we shall discuss these and some theoretical adjustments which should be made to the data in order to make them comparable with the theory. We have tried to evaluate the more predictable of these corrections for all the table entries that we shall use in the analysis in Sections VI and VII, and we have given the sum of these corrections in parentheses after the quantity to which they apply.

5.2 *Surface Damage*

In the attempt to track down empirically the magnitude of the effect of orbital quantization (an investigation which we have summarized in Appendix A) we encountered a stumbling-block in the high value of the infinite-field extrapolation of the resistivity normal to H when H is in a $[100]$ direction. At 77°K , for example, the originally measured $\Delta\rho/\rho_0$ for Sample 606 extrapolates to a value 0.5 unit above the value predicted by the electron-group model. This is many times larger than the expected effect of orbital quantization as deduced from other evidence. However, such an increase could occur if a partial short-circuiting of the large Hall field were effected by surface or dislocation conduction. A surface contribution to the conductance of only a few parts in a thousand would suffice. To check this, some of the measurements on this and other samples were repeated after etching off about a mil of the original sand-blasted surface with superoxol. For the case mentioned (Sample 606 at 77°), two-thirds of the original excess of observation over theory disappeared; the full excess reappeared on sand-blasting again. Thus, it is clear that reliable transverse magnetoresistance results cannot be expected with sand-blasted surfaces.

Theoretically, one would not expect most of the other tabulated quan-

TABLE V

Observed values of thermomagnetic properties of n germanium at various magnetic fields and temperatures for selected orientations and impurity concentrations. All entries refer to unetched samples, and none of the corrections discussed in Section V has been applied.

Sample No., (N_D and N_A , in $10^{19}/\text{cm}^3$)	Direction		T , °K	σ , $\Omega^{-1}\text{cm}^{-1}$	Q_i $\mu\text{V}/\text{deg.}$	Q_{ii} $\mu\text{V}/\text{deg.}$	Magnetic Field, H , in 10^3 gauss						
	∇T	H					0.5	1	2.5	6	10	14	18
							Increase in Seebeck voltage, ΔQ in $\mu\text{V}/\text{deg.}$, as a function of magnetic field and temperature						
606 $N_D = 1$ $N_A \sim 0.2$	100	001	60	0.0723	3196	1166	78.0	214.2	337.6	383.1	401.3	409.8	
	100	100	60	0.0723	3196	1166	346.5	1318	2938	3642	3936	4088	
	100	001	77	0.0491	2329	1199	23.7	83.1	222.5	190.6	206.7	214	
	100	100	77	0.0491	2329	1199	107.4	462.7	1265	1771	2024	2163	
	100	001	94	0.0355	1904	1224	10.4	38.8	89.1	114.1	127.6	134.2	
	100	100	94	0.0355	1904	1224	38.6	182.8	565.2	883.6	1074	1190	
	100	001	122	0.0237	1623	1258	0.96	3.62	48.2	69.9	82.3	89.6	
	100	100	122	0.0237	1623	1258	3.86	12.4	210.6	362.4	474.9	553	
	100	001	148	—	1532	1282	0.47*	1.82*	—	—	—	—	—
	100	100	148	—	1532	1282	1.16*	4.55*	—	—	—	—	—
	100	001	163	—	1504	1294	0.29*	1.16*	—	—	—	—	—
	100	100	163	—	1504	1294	0.68*	3.67*	—	—	—	—	—
580 $N_D = 3.3$ $N_A = 0.5$	100	001	79.5	0.144	2139	1094	17.3	66.8	140.6	177.4	198.7	214	
	100	100	79.5	0.144	2139	1094	71.2	344.9	1024	1487	1727	1863	
	100	001	91	0.1174	1845	1110	8.7	39.1	—	—	—	—	
	100	100	91	0.1174	1845	1110	35.9	181.0	590	922	1117	1234	

601	$N_D = 26$	100	001	94	1.143	1590	910	1.9	6.2	28.5	78.5	108.0	125.4	137.0
	$N_A = 2$	100	100	91	1.151	1590	910	6.2	22.2	123.0	457	750	945	1059
596	$N_D = 180$ $N_A = 6$	100	001	208	0.341	1152	1012	0.100*	0.40*	2.34*	8.4*	—	—	—
		100	100	208	0.341	1152	1012	0.183*	0.78*	4.75*	23.0*	—	—	—
		100	001	234	0.281	1144	1026	0.062*	0.248*	1.48*	7.2*	—	—	—
		100	100	234	0.281	1144	1026	0.132*	0.523*	3.16*	14.6*	—	—	—
		100	001	275	0.208	1142	1048	0.036*	0.144*	0.87*	4.1*	—	—	—
		100	100	275	0.208	1142	1048	0.062*	0.250*	1.53*	7.6*	—	—	—
604	$N_D = 2$ $N_A = 0.6$	100	001	91	4.86	1499	764	1.3	3.7	26.1	66.7	102.8	124.0	138.8
		100	100	91	4.86	1499	764	3.6	12.6	73.0	324	605	803	932
		100	001	235	1.36	1002	886	0.045*	0.179*	1.06	4.8	11.2	17.6	23.4
		100	100	235	1.36	1002	886	0.085	0.337	1.98	10.4	25.0	42.0	56.3
		100	011	61	0.1125	3066	1116	18.8	70.1	318	866	1298	1531	1677
		100	100	61	0.1125	3066	1116	76	274	1124	2615	3319	3624	3784
576A	$N_D = 10$ $N_A = 1$	100	011	77.6	0.0815	2277	1147	6.3*	24.6	118.5	373.3	598.6	748.8	850.0
		100	100	77.4	0.0811	2277	1147	25.1	93.5	423	1183	1676	1923	2059
		100	011	92	0.0630	1886	1171	2.5*	10.4*	54.8	187.0	319	419	491
		100	100	92	0.0630	1886	1171	11.1*	42.8	196.3	611	946	1142	1256
		100	100	235	0.0135	1406	1290	0.046*	0.184*	1.11*	—	—	—	—
		110	001	80	0.495	2019	989	4.1*	15.2	65.8	136.7	174	195	212
110	110	80	0.498	2019	989	7.1*	26.1	90.8	154.7	175.8	185.0	190.4		
110	110	93	0.394	1693	998	3.4	12.8	49.6	98.4	111.1	117.7	121.7		
110	001	133	0.224	1359	1054	0.44*	1.67	9.3	31.6	49.2	60.0	66.8		
110	110	133	0.224	1359	1054	0.77*	2.87	15.8	31.4	41.2	45.4	47.7		

TABLE V — Continued

Sample No., (N_D and N_A in $10^{19}/\text{cm}^3$)	Direction		T , °K	σ , $\Omega^{-1}\text{cm}^{-1}$	$\frac{Q_i}{\mu\text{V}/\text{deg.}}$	$\frac{Q_o}{\mu\text{V}/\text{deg.}}$	Magnetic Field, H , in 10^3 gauss						
	∇T	H					0.5	1	2.5	6	10	14	18
603 $N_D = 2$ $N_A = 0.3$	110	110	61	0.154	3043	1093	61.6	218	861	2060	2824	3237	3482
	110	110	61	0.154	3043	1093	32.8	96.7	231	316	341	351	356
	110	110	77.4	0.1136	2254	1124	19.1	74.9	333	912	1359	1643	1831
	110	110	77.4	0.1136	2254	1124	11.0	39.0	113.7	175.6	191.2	201	205
	110	110	94	0.0816	1830	1150	7.0	25.5	120.3	371	600	765	884
	110	110	94	0.0816	1830	1150	4.1	14.1	50.3	90.4	104.7	109.8	112.6
605 $N_D = 0.3 + A$ $N_A = A$	110	110	135	0.0463	1496	1196	1.44 ^s	5.6 ^s	28.3	—	—	—	—
	110	110	135	0.0463	1496	1196	0.84	3.18	14.4	—	—	—	—
	110	110	155	0.0372	1443	1213	0.71 ^s	2.9 ^s	15.5	—	—	—	—
	110	110	155	0.0372	1443	1213	0.34 ^s	1.34 ^s	7.3 ^s	—	—	—	—
	111	111	91	0.01201	2038	1303	3.4 ^s	13.3 ^s	70.9	233	369	453	499
	111	111	98	0.01145	1918	1313	2.9 ^s	11.3	60.0	204	330	407	456
610 $N_D = 1.7$ $N_A = 0.4$	111	110	92	0.0598	1894	1179	6.4 ^s	23.9	112.4	352	596	726	839
	111	111	92	0.0598	1894	1179	3.7 ^s	14.4	75.3	248	391	477	529
606 $N_D = 1$ $N_A \sim 0.2$	100	001	60	0.0723	3196	1166	123	199	234	158	107	82	66
	100	001	77.4	0.0491	2329	1199	41 ^s	72	107	93	75	61	63
	100	001	94	0.0355	1904	1224	17	31	52	54	48	44	42
	100	001	122	0.0237	1623	1258	4.7	9.4	19.2	26.8	29.0	31.3	34.5
	100	001	155	0.0157	1518	1288	1.55 ^s	3.88 ^s	7.03 ^s	—	—	—	—

Increase in Seebeck voltage, ΔQ , in $\mu\text{V}/\text{deg.}$, as a function of magnetic field and temperatureNernst voltage, BE , in $\mu\text{V}/\text{deg.}$, as a function of magnetic field and temperature

580	$N_D = 3.3$	100	001	79.5	0.144	2139	1094	31	57	94	85	64	—	—	42
	$N_A = 0.5$	100	001	91	0.1174	1845	1110	17	32	56	—	—	—	—	—
		100	001	145	0.0605	1428	1168	1.8*	3.5	7.2	—	9.5	7.9*	6.9*	6.3
601	$N_D = 26$	100	001	94	1.143	1590	910	13.8	26.4	52.3	59	47	38	38	32
	$N_A = 2$	100	001	160	0.534	1196	976	0.65*	1.27*	2.78*	—	1.8	—	—	—
		100	001	208	0.341	1152	1012	-0.39	-0.77	-1.92	0.29	—	—	—	—
		100	001	234	0.281	1144	1026	-0.54*	-1.08*	-2.65*	—	—	—	—	—
		100	001	275	0.205	1142	1048	-0.50*	-0.99*	-2.45*	—	—	—	—	—
596	$N_D = 180$	100	001	91	0.486	1499	764	11.5	22.5	49.5	71	64	53	53	43
	$N_A = 6$	100	001	235	0.136	1002	886	-0.35*	-0.68	-1.68	-4.1	-6.6	-8.4	-8.4	-9.9
604	$N_D = 2$	100	011	77.4	0.0815	2277	1147	40	81	223	485	590	590	590	540
	$N_A = 0.6$	100	011	92	0.0630	1886	1171	17.8*	37.0*	105	252	340	360	360	360
676A	$N_D = 10$	110	001	80	0.495	2019	989	30.6*	57.4	97	89	67	53	53	45
	$N_A = 1$	110	001	133	0.224	1359	1054	2.43	4.67	9.8	12.9	10.9	8.5	8.5	7.6
603	$N_D = 2$	110	110	61	0.153	3043	1093	104	221	602	1060	1080	980	980	860
	$N_A = 0.3$	110	110	77.4	0.1136	2254	1124	45	93	265	540	640	630	630	580
		110	110	94	0.0826	1830	1150	16.4*	34.7	102	245	332	357	357	340
		110	110	135	0.0463	1496	1196	2.3*	4.9	17	57	98	98	123	137
		110	110	155	0.0373	1443	1213	1.0*	2.4	9	34	62	62	84	101
610	$N_D = 1.7$	110	110	178	0.0290	1418	1238	-0.10*	-0.14	1.28*	—	—	—	—	—
	$N_A = 0.4$	111	110	92	0.598	1894	1179	19.2*	39.7*	113	271	359	385	385	380

* From analysis of Hall coefficient-temperature curves

° Value from smoothed curve

TABLE VI

Low- and high-field limits of electrical and thermomagnetic quantities, for various samples and temperatures. For all the entries used in the analyses of Sections VI and VII, and for a few other cases, the sum of the corrections discussed in Section V has been evaluated and is listed in parentheses following the table entry; these corrections are the side-arm correction (Table VII), the temperature-gradient correction (only for the less pure samples — see Table XI) and the boundary-scattering correction, (16) and (17). The estimated true value is obtained by adding the correction to the tabulated value. The symbol (---) means that the correction is zero.

Sample No. (N_D and $N_A \frac{1}{2}$ in $10^{21}/\text{cm}^3$)	Direction		$T, ^\circ\text{K}$	$\mu_H, 10^8 \text{ cm}^2/\text{vs}$	$\lim_{H \rightarrow 0} \frac{\Delta\rho}{\rho H^2}$ $10^{-3} \text{ gauss}^{-2}$	$\lim_{H \rightarrow 0} \frac{\Delta\rho}{\rho}$ $H \rightarrow 0$	$-\frac{Q_p}{\rho}$ $\mu\text{V}/\text{deg.}$	$B(H=0), 10^{-3}$ $\text{V}/\text{gauss deg.}$	$\lim_{H \rightarrow 0} \frac{\Delta Q}{H^2}$ $10^{-13} \text{ v}/\text{gauss}^2 \text{ deg.}$	$\lim_{H \rightarrow \infty} \frac{\Delta Q}{\mu^2}$ $\mu\text{V}/\text{deg.}$
	j or ∇T	H								
606 $N_D = 1$ $N_A \sim 0.2$	100	001	60	51.9 (-1.35)	10.8 (---)	1.68 (---)	2030 (+283)	—	103 (+10.6)	424 (+53)
	100	100	60	—	27.5 (-0.1)	3.94 (-0.11)	—	—	429 (+62.5)	4345 (+605)
	100	001	63	48.3* (-1.26)	8.2* (---)	0.68* (---)	1830	—	—	—
	100	100	63	—	23.3* (-0.1)	3.32* (-0.11)	—	—	—	—
	100	001	77.4	37.5* (-0.9)	5.4* (---)	0.61* (---)	1130 (+60)	80 (+0.8)	26.3 (+0.9)	225 (+10)
	100	100	77.4	36.0	5.6	0.94	—	—	120 (+8)	2415 (+127)
	100	001	94	26.7 (-0.7)	14.0* (-0.06)	3.28* (-0.12)	—	—	—	—
	100	100	94	—	13.9	3.45	—	—	—	141* (+3)
	100	001	122	—	3.3* (---)	0.56* (---)	680 (+17)	31*	9.9* (+0.1)	150
	100	100	122	—	8.1* (-0.03)	0.73	—	31	43* (+1.4)	1450* (+35)
	100	001	131	18.18 (-0.47)	8.04	3.37	—	—	42	1413
	100	100	131	—	1.58 (---)	0.625 (---)	365 (+3)	9.8 (-0.1)	3.8 (---)	104
	100	001	148	16.4* (-0.4)	3.76 (-0.01)	3.04 (-0.10)	315 (+2)	4.2* (-0.08)	12.8 (+0.2)	721 (+6)
	100	100	148	—	1.23* (---)	0.54* (---)	250	—	2.4* (---)	87* (+0.3)
	100	001	155	—	2.99* (-0.01)	2.98* (-0.10)	210	—	8.5* (+0.14)	597* (+4)
	100	100	155	12.06	0.88	1.87	—	3.1	1.89	—
100	001	163	—	0.72	—	230	—	4.6	—	
100	100	163	—	1.43	—	210	—	1.16	—	
									2.7	

TABLE VI — Continued

Sample No. (N_D and N_A † in $10^{19}/\text{cm}^3$)	Direction		T, K	$\mu H, 10^3 \text{ cm}^2/\text{vs}$	$\lim_{H \rightarrow 0} \frac{\Delta \rho}{\rho H^2},$ $10^{-3} \text{ gauss}^{-2}$	$\lim_{H \rightarrow \infty} \frac{\Delta \rho}{\rho}$ $H \rightarrow \infty$	$-Q_p,$ $\mu\text{v}/\text{deg.}$	$B(HA=0), 10^{-9}$ $\text{v}/\text{gauss deg.}$	$\lim_{H \rightarrow 0} \frac{\Delta O}{H^2},$ $10^{-12} \text{ v}/\text{gauss}^2 \text{ deg.}$	$\lim_{H \rightarrow \infty} \Delta O,$ $\mu\text{v}/\text{deg.}$
	j or V_T	H								
576A $N_D = 10$ $N_A = 1$	100	100	77.4	—	$12.7^* \dagger (-0.02)$ 12.4	$3.28^* \dagger (-0.06)$ 3.24	—	—	—	—
	100	011	92	$26.4 (-0.7)$	$3.2 (-0.04)$	$1.94 (-0.02)$	715	$36 (+0.4)$	$103 (+5)$	$2265 (+120)$
	100	100	92	—	$7.86 (-0.03)$	$3.23 (-0.11)$	—	—	$10.7 (-0.2)$	$638 (+16)$
	100	011	195	7.96	0.32	0.68	155	—	$45 (+1.5)$	$1502 (+40)$
	100	100	195	—	0.747	1.72	—	—	—	—
	100	011	235	$5.70 (-0.15)$	$0.183 (-0.002)$	—	116	—	—	—
	100	100	235	—	$0.388 (-0.002)$	—	—	—	0.18	—
	110	001	63	40.9	—	—	1830	—	—	—
	110	001	77.4	32.1	—	—	$1130 (+60)$	—	—	—
	110	110	77.4	—	4.6	0.334	$1130 (+60)$	—	—	—
603 $N_D = 2$ $N_A = 0.3$	110	001	80	—	—	—	1030	62.5	16.7	252
	110	110	80	—	3.1	0.319	—	—	29.2	200
	110	110	91	—	—	—	735	—	—	—
	110	001	93	25.8	—	—	695	—	—	130
	110	110	93	—	—	—	—	—	14.4	80.3
	110	001	133	14.8	—	—	305	4.92	1.74	52.2
	110	110	133	—	—	—	—	—	3.13	—
	110	001	195	—	0.30	—	155	—	—	—
	110	110	61	47.0	—	—	1950	200	257	3960
	110	110	61	$35.8^* (-0.9)$	$10.4^* (+0.06)$	$3.72^* (+0.04)$	$1130 (+60)$	—	150	365
110	110	110	77.4	34.6	6.20* (-0.22)	0.328* (-0.035)	1130 (+60)	87 (+0.9)	79.5 (+4.2)	2264 (+118)
	110	110	77.4	—	5.95	0.328	—	—	—	—
	110	110	94	$26.0 (-0.7)$	$5.95 (+0.04)$	$3.32 (+0.03)$	680 (+17)	—	47 (+2.8)	214 (+11)
	110	110	94	—	$3.40 (-0.12)$	$0.313 (-0.034)$	—	—	27.9 (+0.75)	1163 (+28)
	110	110	135	15.21	2.25	—	300	4.6	15.9 (+0.5)	116.5 (+2.9)
	110	110	135	—	1.20	—	—	—	5.8	—
	110	110	135	—	—	—	—	—	3.4	—

	110	110 155	12.31	1.53	—	230	2.1	2.9	—
	110	110 155	—	0.82	—	—	—	1.4	—
	110	110 178	9.38	1.02	—	180	-0.21	1.7	—
	110	110 178	—	0.50	—	—	—	0.80	—
	110	110 200	7.77	—	—	150	—	—	—
	110	110 204	7.49	0.565	—	146	—	0.66	—
	110	110 204	—	0.305	—	—	—	0.39	—
595	111	110 77.4	30.4	—	—	1130(+60)	—	—	—
$N_D = 0.3 + A$	111	111 91	—	—	—	785	—	13.8	593
$N_A = A$	111	111 98	—	—	—	605	—	11.7	546
610	111	110 77.4	36.6*(-0.9)	8.65*(\approx)	3.07*(+0.02)	1130(+60)	—	—	—
$N_D = 1.7$	111	111	35.7	8.6	3.08	—	—	—	—
$N_A = 0.4$	111	111 77.4	—	4.31*(-0.15)	1.42*(-0.06)	—	—	—	—
	111	110 92	28.3(-0.7)	4.35	1.44	715	38(+0.3)	26.0(+0.5)	1079(+28)
	111	111 92	—	5.5(\approx 0)	2.90(+0.02)	—	—	14.9(+0.5)	641(+17)
	111	111 92	—	2.78(-0.09)	1.39(-0.04)	—	—	—	—

* Etched samples

† Side-arms removed

‡ From analysis of Hall coefficient-temperature curves

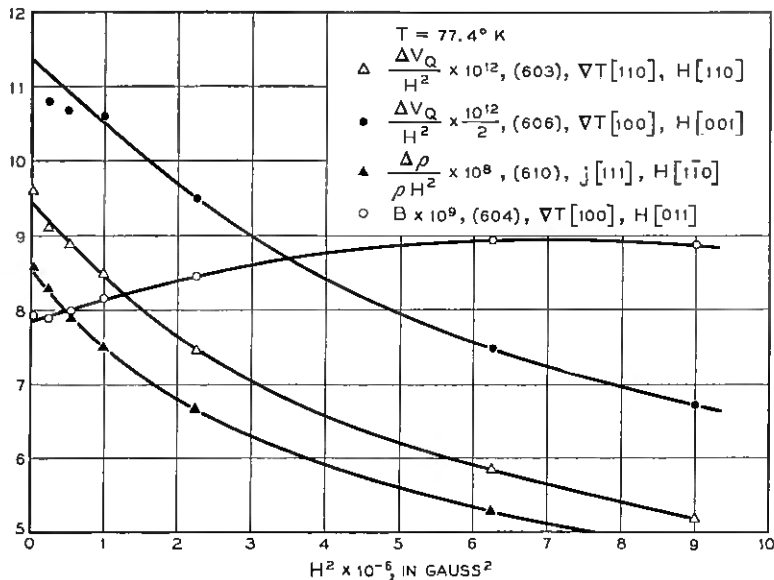


Fig. 5— Sample plots of low-field Nernst, magnetoresistance and ΔQ data against the square of the magnetic field H showing extrapolations to zero field. Here ΔV_Q is the change in the measured thermoelectric emf in volts between the end-arms, due to H , at constant temperature gradient. The field H is measured in gauss, and B is in volts/degree gauss.

ties to be nearly as sensitive to surface effects as is the high-field transverse magnetoresistance, since only in the latter case is a large Hall field present. The high-field Nernst coefficient (see Section VIII) is the one exception, since its approach to zero makes a small Hall field due to longitudinal counter-currents amount to a large percentage error. The measurements confirm the expectation that the effects on the other quantities should be small. However, though small, the changes found in the apparent mobility, etc., were often appreciable. The effects seem to involve a complex interplay of reduced mobility in the damaged surface layer and conduction with little or no Hall effect; moreover, the latter conduction seems to be extremely anisotropic, since the large excess of magnetoresistance was found only when there were [100] surfaces normal to the magnetic field.

Table VI lists, along with the original values for sand-blasted surfaces, the values (identified by asterisks) appropriate to etched surfaces, for all cases for which the latter were measured.

5.3 Side-Arm Corrections

The raw data given in Tables V and VI were obtained from the measurements on the assumption that the flow of heat or electric current down the sample was purely one-dimensional and that the side-arms acted merely as infinitesimal probes with which potentials could be measured without disturbing this flow. This assumption is not quite correct, for the width of the side-arms (0.06 cm) was not negligible compared with the width (0.15 cm) of the current-carrying portion. Currents passed through such a sample must bulge out a little into the side-arms, in a way which will be altered by a magnetic field. Moreover, the anisotropy which a magnetic field introduces into the thermoelectric power can cause circulating currents to flow near the side-arms in the presence of an electric field, even though no net current flows along the specimen. These effects can become very serious when a large magnetic field is ap-

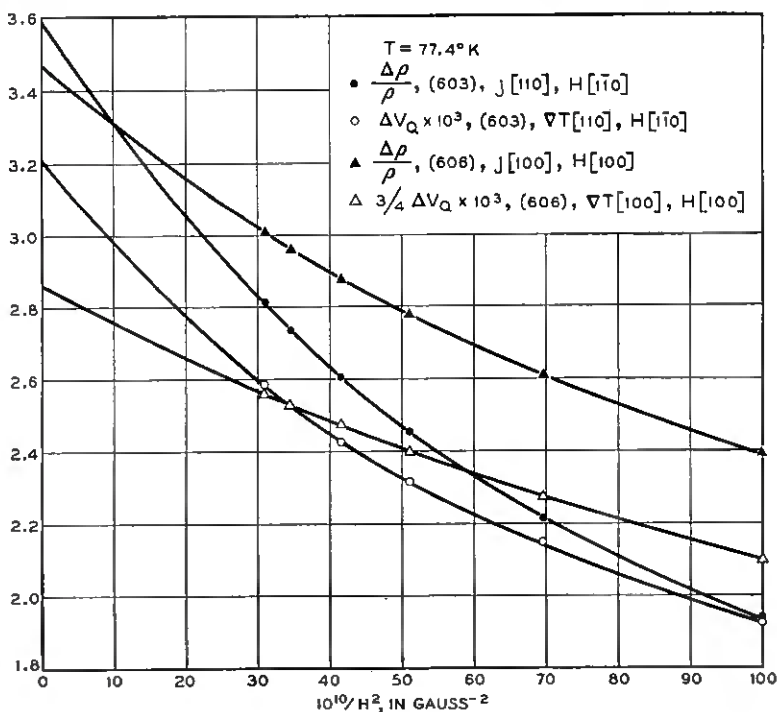


Fig. 6 — Sample plots of high-field magnetoresistance and ΔQ data against the inverse square of the magnetic field H , showing extrapolations to infinite field. Notation and units are the same as for Fig. 5.

plied in the direction of the side-arms. For example, in indium antimonide magnetoresistances of several times the ideal value have been observed for this geometry.¹⁶ For this reason, all the transverse measurements reported here were made with the magnetic field normal to the plane of the arms, and we shall show that, for this case and for all longitudinal and zero-field cases, the side-arm errors are small and can be approximately corrected for.

A number of statements about the effects of the side-arms can be made from simple theoretical considerations. For example, it is obvious that the current distortions have no effect on the measured values of the zero-field thermoelectric power or of the Hall coefficient. The conductance of the sample is, of course, higher with the side-arms than without; measurements on scaled cutouts of stainless-steel sheet showed the difference of true and apparent resistivities in the absence of a magnetic field to be about 2.6 per cent for our geometry. However, it can be shown fairly easily that the measured transverse magnetoresistance $\Delta\rho/\rho$, and likewise the transverse ΔQ , are unaffected by the presence of the side-arms when \mathbf{H} is in a [100] direction; this invariance arises from the two-dimensional nature of the potential distribution and the isotropy of ρ and Q in the plane normal to \mathbf{H} . When \mathbf{H} is in a [110] direction the anisotropies of $\Delta\rho$ and ΔQ cause a difference between the apparent and true $\Delta\rho/\rho$ and ΔQ , which for transverse cases can be expressed in terms of the anisotropies and the correction to the zero-field resistivity. When \mathbf{H} is longitudinal, more complicated effects arise, due to the Hall fields where the current lines depart from longitudinality at the bases of the side-arms. However, for almost all the cases listed in Table VI there exist scaling relations which allow the side-arm corrections to be expressed, at least approximately, in terms of the correction to the apparent conductivity at $H = 0$.

Several experiments were performed to verify that these effects were small and to evaluate the corrections for the two cases (high-field transverse ΔQ with \mathbf{H} in a [110] direction) for which the scaling relations do not suffice for a theoretical prediction. These experiments consisted in removing the middle pair of side-arms of some of the specimens and noting the changes in the apparent values of $\Delta\rho/\rho$ and ΔQ . The changes due to removal of all side-arms would, of course, be twice as great. Table VII summarizes the theoretical and experimental results for the corrections to be added to the entries of Table VI. In no case are these greater than about 4 per cent.

5.4 Gross Inhomogeneities

Gradients of impurity concentrations can greatly falsify the results. For example, the apparent Hall constant depends on the carrier concen-

TABLE VII

Approximate correction factors by which the entries of Table VI should be multiplied to give the "true" quantities as they would be measured on a specimen without side-arms. Except for the ΔQ 's, the factors are independent of temperature, or nearly so. In the last two rows, the factors given apply in the neighborhood of liquid-air temperature.

Direction of \mathbf{j} or ∇T Direction of \mathbf{H}	[100] [001]	[100] [011]	[110] [1 $\bar{1}$ 0]	[110] [110]	[100] [100]	[111] [111]
Quantity:						
μ_H	1.026	1.026	1.026	---	---	---
$\lim_{H \rightarrow 0} \left(\frac{\Delta \rho}{\rho H^2} \right)$	1.000	0.987	1.006	0.965	0.996	---
$\lim_{H \rightarrow \infty} \frac{\rho(H)}{\rho(0)}$	1.000	0.994	1.008	0.974	0.974	0.974
$B(0)$	0.974	0.974	0.974	---	---	---
$\lim_{H \rightarrow 0} \frac{\Delta Q}{H^2}$	1.000	0.975	1.009	1.014	1.013	---
$\lim_{H \rightarrow \infty} \Delta Q$	≈ 1.000	≈ 1.000	≈ 1.000	1.000	1.000	1.000

tration between the middle side-arms, while the apparent conductivity depends on the carrier concentration over the whole length; if the concentration varies, the apparent Hall mobility will be incorrect. Most of the samples were checked for homogeneity by measuring Hall constant on the end arms as well as on the middle pair. All three measurements usually agreed to within 1 per cent. Inhomogeneities can also cause measurements designed to give quantities even in the magnetic field ($\Delta \rho$, ΔQ) to contain contributions from effects odd in the field (R , B), and vice versa. This effect occasionally showed up as an asymmetry in the measured values when the direction of \mathbf{H} was reversed. This asymmetry was small for all cases listed in the tables and was consistent with the present interpretation, e.g., the magnetoresistance asymmetry could be calculated from the Hall inhomogeneity. The entries represent averages of measurements for \mathbf{H} and $-\mathbf{H}$.

5.5 Value of the Temperature Gradient

The measurements tabulated in Tables V and VI were all taken with the special-purpose apparatus described in our earlier paper.¹ With this

apparatus the temperature gradient used for the thermomagnetic effects was not measured directly, but was computed from the thermoelectric voltage developed between the potential leads at $H = 0$, together with the thermoelectric power Q of the specimen. The latter was computed on the assumption that, for a given temperature, the only source of variation of Q from specimen to specimen is the variation in the height of the Fermi level associated with the different carrier concentrations. In other words, the Q for any sample was computed from values previously measured on other samples by adding a term equal to $86 \mu\text{v}/\text{degree}$ times the logarithm of the ratio of carrier concentrations. This will not be correct if impurity scattering affects Q appreciably.

Previous empirical evidence suggested that, for fairly pure material, Q is almost unaffected by impurity scattering. The same conclusion follows from a theoretical analysis which we shall carry out in Section VII. When $T \geq 77^\circ\text{K}$, it turns out that no correction is needed to the data for Samples 603, 604, 606 and 610, on which most of the analysis of the next two sections will be based, nor for Sample 595. The corrections should be small for these samples at $T \approx 60^\circ\text{K}$, and for Samples 576A and 580 at $T \geq 77^\circ$. For Sample 601 the correction begins to be perceptible at 94° , and for Sample 596 at 87° it amounts to about 8 per cent of the quantities tabulated.

5.6 *Boundary Scattering*

The final—and largest—correction we shall consider has to do with the scattering of phonons from the boundary of the specimen. We wish to interpret the data in terms of a theory of conduction in an infinite homogeneous medium. The side-arm correction has eliminated effects of the shape of the sample, but we have yet to correct for the finiteness of its dimensions.

We may start from the observation that phonon-drag effects are due to the anisotropy of the low-frequency part of the phonon distribution. This anisotropy, due to the temperature gradient, has a certain value at depths below the surface that are large compared with the mean free path of a phonon. However, immediately below the surface the anisotropy is only about half as great, since that half of the phonons that comes from the surface has no tangential anisotropy if the surface scatters diffusely. Thus, one may say, approximately, that the tensor \mathbf{Q}_p has one value at depths greater than some value Δ , of the order of a phonon mean free path, and has half this value at depths less than Δ . The effect of this alteration is easily seen to be to make the effective $\mathbf{Q}_p(\mathbf{H})$ less than the in-

terior value by the ratio

$$\frac{Q_p}{Q_p \text{ (ideal)}} = 1 - \frac{2\Delta}{A^{\frac{1}{2}}}, \quad (16)$$

independently of H , where A is the cross-sectional area of the sample. The effect on the low-field Nernst coefficient can be calculated similarly; it is

$$\frac{B_p}{B_p \text{ (ideal)}} = 1 - \frac{\Delta}{t}, \quad (17)$$

where t is the thickness of the sample normal to its plane. We shall take Δ as roughly independent of the orientation of the surface and equal, in millimeters, to $0.037 (77^\circ/T)^{3.9}$, the best value currently available from experiments on the variation of Q_p with diameter near liquid-air temperature.²

In parentheses after each entry in Table VI we have given the sum of the side-arm correction calculated from Table VII and (for the thermoelectric quantities) the temperature-gradient correction and the boundary-scattering correction computed from (16) or (17). This total correction is to be added to the table entry to get the ideal value appropriate for the comparison with theory.

VI. ANALYSIS

At each temperature and for each impurity content we have 10 or 11 measured quantities which depend on the partial Peltier coefficients $\Pi_{\parallel}(\epsilon)$, $\Pi_{\perp}(\epsilon)$ of the energy shells, and among which no relations derivable from phenomenological theory exist. These are: Q and B at $H = 0$, $\Delta Q(H = \infty)$ for [100], [110] and [111] longitudinal orientations, one transverse $\Delta Q(H = \infty)$ for $\mathbf{H} \parallel$ [100] and two such for $\mathbf{H} \parallel$ [110], and two (or three) constants describing the low-field ΔQ behavior. (There are three phenomenologically independent constants of the latter type, but there is one relation between them which follows merely from the band structure and the electron-group approximation.¹) We shall now derive what information we can about $\Pi_{p\parallel}(\epsilon)$ and $\Pi_{p\perp}(\epsilon)$ by fitting the theory of Section IV and Appendix B to the data.

We shall first note briefly that all these data can be approximately fitted by the specialized formulas of Tables I, III and IV, with three adjustable constants A_{\parallel} , A_{\perp} and n , where

$$\Pi_{p\parallel, \perp} = A_{\parallel, \perp} \left(\frac{\epsilon}{kT} \right)^n$$

TABLE VIII

Comparison of thermomagnetic data observed at 94°K with values computed from the formulas of Tables I and III and from the mean moment curves of Fig. 7. For the former calculation $w = 0.061$ and the approximately optimum values $p \equiv \Pi_{p\parallel}/\Pi_{p\perp} = 9.5$, $Q_p = 697$, $\Pi_p \propto \epsilon^{-0.25}$ were used. The corrections in parentheses in Table VI have been included in the observed values, and the data for Samples 603 and 610 have been corrected to 94°. The data for Sample 606 are starred to indicate that this sample had been etched. The units of the ΔQ 's are $\mu\text{V}/\text{degree}$, those of $B(0)$ are $10^{-9}\text{V}/\text{gauss degree}$. Acoustic scattering has been assumed; correction for impurity scattering would lower the computed $B(0)$ slightly. In the first row, $Q(0) - Q^*$ is T^{-1} times the Peltier heat relative to the band edge, i.e., the sum $Q_p + 172\mu\text{V}/\text{degree}$.

Quantity	Sample No.	Directions		Value Computed from Tables I and III	Observed Value	Value Computed from Fig. 7
		H	∇T			
$Q^* - Q(0)$ $-\Delta Q(\infty)$				869	869	863
				1440		1459
	606	[001]	[001]		1484*	
	604				1460	
	610	[111]	[111]	605	621	615
	603	[011]	[011]	118	119	124
	606	[001]	[100]	91	144*	93
	603	[011]	[011]	1325	1191	1187
$B(0)$	604	[011]	[100]	655	613	610
				22		34
	606				31*	
	603				32	
	604				36	

and with the value of $w = m^*_{\perp} \tau_{\parallel} / m^*_{\parallel} \tau_{\perp}$ given by, say, the high-field magnetoresistance measurements. Table VIII gives, in its fifth and sixth columns, a sample comparison of observed values with the predictions of this simplified theory; low-field ΔQ values have been omitted from the table because, as we shall see in detail in Section VII, they need to be corrected for impurity scattering. The extent of the agreement is a fairly good measure of the adequacy of the assumed behavior of Π_p , since for most of the quantities listed the electron-diffusion contribution is small enough so that uncertainties in its value are unimportant. The observations are fitted well enough to engender some confidence in the theory, but there are systematic discrepancies which suggest that $\Pi_{p\parallel}$ and $\Pi_{p\perp}$ may have different energy dependences that a less restrictive analysis could reveal.

6.1 Analysis in Terms of Moments

As we have mentioned in Section III, the most general formulas of the electron-group theory, as derived in Appendix B, give Q or ΔQ and B as linear combinations of Maxwellian averages of quantities of the form $\epsilon \Pi_{\parallel, \perp} f^{(n)}(\tau_{\parallel}, \tau_{\perp})$, where $f^{(n)}$ is homogeneous of degree n . The coefficients contain in their denominators linear combinations of similar Maxwellian averages, without the Π 's. Now when n is 0 or 1, corresponding respectively to the high-field transverse Q and to $Q(0)$, or the high-field longitudinal Q , it should be quite a good approximation to take τ_{\parallel} and τ_{\perp} both proportional to $\epsilon^{-\frac{1}{2}}$ in the purest of our samples (ideal acoustic scattering). The value of $w = m^*_{\perp} \tau_{\parallel} / m^*_{\parallel} \tau_{\perp}$ appropriate to these cases can be determined from the high-field longitudinal magnetoresistance. However, for $n = 2$ and especially for $n = 3$ (low-field effects), the Maxwellian averages in question become very sensitive to small amounts of impurity scattering, and one might question the legitimacy of taking the ratio $\tau_{\parallel} / \tau_{\perp}$ independent of energy. But since the term containing $\tau_{\parallel}^m \tau_{\perp}^{n-m}$ also contains $m^*_{\parallel}{}^{-m} m^*_{\perp}{}^{m-n}$, and since $\Pi_{\parallel, \perp}$ vary only slowly with energy, it will always be a good approximation to assume that, with angular brackets as usual denoting Maxwellian averages,

$$\left\langle \epsilon \Pi_{\parallel, \perp} \left(\frac{\tau_{\parallel}}{m^*_{\parallel}} \right)^m \left(\frac{\tau_{\perp}}{m^*_{\perp}} \right)^{n-m} \right\rangle \approx w_n^m \left\langle \epsilon \Pi_{\parallel, \perp} \left(\frac{\tau_{\perp}}{m^*_{\perp}} \right)^n \right\rangle, \quad (18)$$

$$\left\langle \epsilon \left(\frac{\tau_{\parallel}}{m^*_{\parallel}} \right)^m \left(\frac{\tau_{\perp}}{m^*_{\perp}} \right)^{n-m} \right\rangle \approx w_n^m \left\langle \epsilon \left(\frac{\tau_{\perp}}{m^*_{\perp}} \right)^n \right\rangle, \quad (19)$$

where, for $n = 1, 2, 3$, w_n is determined from purely electrical data. Specifically, we can define w_n so that (19) is exact for $m = 1$; (18) will also then be nearly exact for $m = 1$. For $m > 1$, (19) and (18) may be less accurate, but this inaccuracy will have very little effect on any electrical or thermomagnetic quantity, because the terms with $m > 1$, being of order w_n^m , will be very small.

We shall use the approximations (18) and (19) to express the various thermomagnetic quantities in terms of "moments" $\Pi_{\parallel}^{(n)}$ and $\Pi_{\perp}^{(n)}$, where, for $n = 0, 1, 2, 3$,

$$\Pi_{\parallel, \perp}^{(n)} \equiv \frac{\langle \epsilon \tau_{\perp}^n \Pi_{\parallel, \perp} \rangle}{\langle \epsilon \tau_{\perp}^n \rangle}. \quad (20)$$

From the formulas of Appendix B we obtain the expressions listed in Table IX. Since the low-field Nernst coefficient involves $\Pi_{\parallel, \perp}^{(1)}$ and $\Pi_{\parallel, \perp}^{(2)}$, we have listed as the second entry in the first column that combination of B and Q which involves only the $\Pi_{\parallel, \perp}^{(2)}$. Similarly, we have

TABLE IX

Thermomagnetic quantities in terms of the moments $\Pi_{\parallel, \perp}^{(n)}$ defined by (20), as obtained from the formulas of Appendix B using the approximations (18) and (19) and the assumption of n-type material and [111] valleys.

Quantity	Directions of H and ∇T	Expression	Reference, Appendix B
$TQ(\infty)$	[011] [100]	$\frac{1}{3} \frac{1+2w}{2+w} \Pi_{\parallel}^{(0)} + \frac{1}{3} \frac{(5+w)}{(2+w)} \Pi_{\perp}^{(0)}$	(81)
	[011] [011]	$\frac{1}{3} \Pi_{\parallel}^{(0)} + \frac{2}{3} \Pi_{\perp}^{(0)}$	(81)
	[001] any \perp H	$\frac{w}{1+2w} \Pi_{\parallel}^{(0)} + \frac{1+w}{1+2w} \Pi_{\perp}^{(0)}$	(81)
	[001] [001]	$\frac{1}{3} \Pi_{\parallel}^{(1)} + \frac{2}{3} \Pi_{\perp}^{(1)}$	(78)
	[111] [111]	$\frac{1+2w}{7+2w} \Pi_{\parallel}^{(1)} + \frac{6}{7+2w} \Pi_{\perp}^{(1)}$	(80)
	[011] [011]	$\frac{w}{1+2w} \Pi_{\parallel}^{(1)} + \frac{1+w}{1+2w} \Pi_{\perp}^{(1)}$	(79)
$TQ(0)$	— —	$\frac{w}{2+w} \Pi_{\parallel}^{(1)} + \frac{2}{2+w} \Pi_{\perp}^{(1)}$	(50)
$-T \left[\frac{B(0)}{(\mu_H/c)} + Q \right]$	— —	$\frac{w_2}{1+2w_2} \Pi_{\parallel}^{(2)} + \frac{1+w_2}{1+2w_2} \Pi_{\perp}^{(2)}$	(57)
$T \left[-\lim_{H \rightarrow 0} \frac{\Delta Q + (B\mu_H H^2/c)}{(\Delta\rho/\rho) + (\mu_H H/c)^2} + Q \right]$	[001] [100]	$\frac{w_3}{2+w_3} \Pi_{\parallel}^{(3)} + \frac{2}{2+w_3} \Pi_{\perp}^{(3)}$	(58), (65)
$T \left[-\lim_{H \rightarrow 0} \frac{\Delta Q}{(\Delta\rho/\rho)} + Q \right]$	[001] [001]	$-\frac{w_3}{1-w_3} \Pi_{\parallel}^{(3)} + \frac{1}{1-w_3} \Pi_{\perp}^{(3)}$	(58), modified (65)

listed in the last two rows those combinations of the low-field ΔQ , B , and Q which involve only the $\Pi_{\parallel, \perp}^{(3)}$. These equations can in most cases be inverted to give the $\Pi_{\parallel, \perp}^{(n)}$ in terms of the observed quantities. Thus, the first two rows just suffice to determine $\Pi_{\parallel, \perp}^{(0)}$, the next four rows overdetermine $\Pi_{\parallel, \perp}^{(1)}$, the next row gives one relation between $\Pi_{\parallel}^{(2)}$ and $\Pi_{\perp}^{(2)}$, and the last two rows just determine $\Pi_{\parallel, \perp}^{(3)}$.

Our procedure will now be to determine the moments $\Pi_{\parallel, \perp}^{(n)}$ from these equations and the observational data, and then to discuss the separation of each such empirically determined $\Pi_{\parallel, \perp}^{(n)}$ into electron-diffusion and phonon-drag contributions and to interpret the variation of the latter part with n in terms of the energy dependence of $\Pi_{p\parallel, \perp}$.

TABLE X

Values of w_1 , w_2 and w_3 used in the moment analysis. The values of w_1 were obtained independently for each case by fitting the electron-group theory to the observed high-field limit of the longitudinal magnetoresistance. The side-arm corrections of Table VI have been included and, whenever possible, data obtained after etching have been used. The latter cases are identified by asterisks. The values of w_3 were estimated from these w_1 's (assumed for the first four samples equal to 0.061) and the approximate dependence of w_3 on impurity concentration found by Goldberg.¹⁷ The values of w_2 were interpolated between w_1 and w_3 .

Sample	T, °K	Empirical w_1	Assumed w_2	Assumed w_3
606	60°		0.0625	0.064
	63°	0.061*		
	77.4°	0.0615*	0.0615	0.062
	94°	0.061*	0.061	0.061
604	63°	0.062		
	77.4°	0.061*	0.0625	0.064
	92°	0.0625	0.0615	0.062
610	77.4°	0.061*	0.0625	0.064
	92°	0.061	0.0615	0.062
603	77.4°		0.0625	0.064
	94°		0.0615	0.062
Value adopted for all above cases		0.0610		
601	94°	0.069	0.076	0.084
596	87°	0.081		
606	131°	0.067* (0.0595 adopted, see text)	0.0595	0.0595

6.2 Choice of the w_n

To begin with, we need values of w_1 , w_2 and w_3 . Magnetoresistance evidence on w ($\equiv w_1$) for the present samples is summarized in the column labeled "Empirical w_1 " of Table X. The good agreement of the w_1 values from Samples 606 and 604 with that from Sample 610 speaks well for the validity of the electron-group model. There seems to be no significant change of w_1 with temperature in the range 60° to 94°K. No reliable value of w_1 can be obtained from Sample 603, since in this orientation a change of w_1 from 0.04 to 0.08 only changes the theoretical high-field resistance by 0.05 of the zero-field value. The data for 603 do, however, agree fairly well with the theory, in that the $\Delta\rho/\rho_0$ predicted from the w_1 of the other specimens is within 0.02 or 0.03 of the value in Table VI.

The empirical w_1 value for Sample 606 at 131° , though listed, is not to be trusted, as a field of 18,000 gauss does not give a very good approach to saturation. We have preferred to rely on the value 0.061 found at the lower temperatures, decreasing it to 0.0595 to allow for the gradual decrease of w with temperature found by Goldberg,^{17,18} a decrease which also shows up, though less clearly, in our own low-field data at higher temperatures.

The estimation of w_3 is somewhat less satisfactory. In principle w_3 can be estimated from the longitudinal and transverse low-field magnetoresistances, since, as we noted at the start of this section, two independent phenomenological constants should suffice for the complete description of the low-field magnetoresistance in reasonably pure material. However, the w_3 values computed in this way for Samples 603, 604 and 610 turn out to be quite different — ranging, for example, from 0.056 to 0.073 at 77°K , all for etched samples — whereas the similarity in impurity content between these three samples would lead one to expect similar w_3 's. This sort of discrepancy might have been anticipated from the fact (Ref. 1, Fig. 13) that the four data for 603 and 604 depart perceptibly from the predictions of the phenomenological theory with two adjustable constants. The trouble can arise from any of several causes. To compute w_3 from the data, one needs the longitudinal and transverse magnetoconductances, and the second of these equals the corresponding magnetoresistance plus $(\mu_H/c)^2$. Now the empirical value of μ_H , unlike magnetoresistance and ΔQ values, depends on the measured thickness of the sample and is affected by errors in this measurement. Also, it can be falsified if the carrier concentration between the Hall arms differs slightly from the average over the length of the sample. Since a change of μ_H by 1 per cent affects w_3 by 0.003, this source of scatter from specimen to specimen could well be serious. Moreover, our random errors are large enough to affect w_3 perceptibly, and comparison of etched and unetched samples has shown that the surface damage effect can alter the empirical w_3 by 0.007 or so.

Fortunately, a study of the formulas of Table IX, which we shall describe below, shows that, once w_3 is specified, the $\Pi_{\parallel,\perp}$ ⁽³⁾ computed from the measured quantities is not nearly so sensitive to small errors in the latter as we have just found the empirical w_3 to be. This suggests that one can get a fairly reliable analysis by disregarding the apparent w_3 's of the specimens and estimating the true w_3 for each temperature and impurity concentration in some other way. Now for pure acoustic scattering w_3 should be practically the same as w_1 , for which we have obtained the presumably trustworthy value 0.061 (Table X). Several studies of the effect of impurity scattering on magnetoresistance have

been published,^{17,18,19} and, while the absolute values of w_3 obtained in these studies may possibly have been affected by various kinds of systematic errors, their picture of the variation with impurity content and temperature is undoubtedly correct. From the work of Goldberg¹⁷ and of Goldberg and Howard,¹⁸ it appears that, at 77°K, a sample with impurity concentration $7 \times 10^{13} \text{ cm}^{-3}$ has a w_3 differing from that of an ideally pure sample by only about 10 per cent, or at the very most 20 per cent. † Thus it is a reasonable guess to say that the corresponding deviation at a concentration of $2 \times 10^{13} \text{ cm}^{-3}$ is perhaps 5 per cent, i.e., 0.003 in w_3 . While the error in this estimate may be as large as the estimate itself, it is surely much smaller than the scatter mentioned in the preceding paragraph, and we shall see that this accuracy suffices for our present purpose. The last column of Table X gives the values we shall adopt for our analysis, obtained by adding an estimate of this sort to our w_1 . The next to the last column gives the corresponding w_2 's, obtained by interpolating between w_1 and w_3 .

6.3 Results for $\Pi_{\parallel}^{(n)}$ and $\Pi_{\perp}^{(n)}$

It is now a straightforward matter to insert the w_n of Table X into the formulas of Table IX, equate to the empirical quantities of Table VI (with the corrections included), and solve for the moments $\Pi_{\parallel,\perp}^{(n)}$ defined by (20). To make a comparison of the values for different specimens more meaningful, it is best to compare not the $\Pi_{\parallel,\perp}^{(n)}$, which depend on the location of the Fermi level through Π_c , but the quantities $\Pi^* - \Pi_{\parallel,\perp}^{(n)}$, where

$$\Pi^* = \frac{\epsilon_F - \epsilon_b}{e} \tag{21}$$

represents the part of Π_c due to the difference between the Fermi level ϵ_F and the band edge ϵ_b . The quantities $\Pi^* - \Pi_{\parallel,\perp}^{(n)}$ thus are positive for n-type material and measure moments of the Peltier heat relative to the band edge.

Fig. 7 shows the results for $\Pi^* - \Pi_{\perp}^{(n)}$ and $\Pi^* - \Pi_{\parallel}^{(n)}$, respectively, for a number of the purest samples at several temperatures from 60° to 131°K. Here the moments for $n = 1$ were calculated from Q and the high-field limit of ΔQ for H longitudinal. Those for $n = 3$ were calculated from the longitudinal and transverse values of

$$\lim_{H \rightarrow 0} \frac{\Delta Q}{H^2},$$

† We are indebted to C. Goldberg for having informed us of some of these results prior to publication.

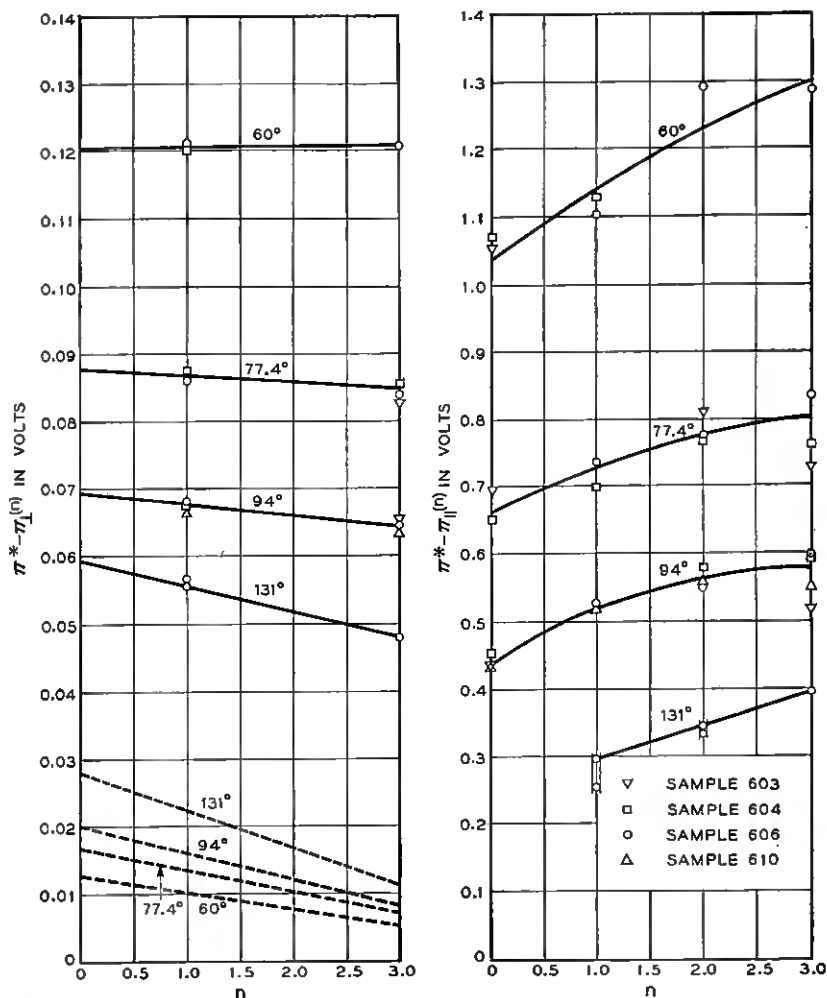


Fig. 7 — Empirical values of the moments $\Pi_{\perp}^{(n)}$ (left) and $\Pi_{\parallel}^{(n)}$ (right) defined by (20), as obtained from the method of analysis described in the text. Note that the two scales differ by a factor of ten. The dashed curves at the bottom of the chart for $\Pi_{\perp}^{(n)}$ show the contributions which Π_a would make to $\Pi_{\perp}^{(n)}$ or $\Pi_{\parallel}^{(n)}$ at the various temperatures, if the scattering were purely acoustic. The full curves are reference means which we shall use in Figs. 11 and 12.

using, if necessary, the phenomenological identities mentioned at the bottom of Table IV to relate the ΔQ values for the directions going with a particular sample to those for the directions listed in Table IX. For $n = 2$, only one moment can be determined; since $\Pi_{\perp}^{(n)}$ is very reproducible and varies very little with n , we have chosen to determine $\Pi_{\perp}^{(2)}$ by linear interpolation between $\Pi_{\perp}^{(1)}$ and $\Pi_{\perp}^{(3)}$, and then to determine $\Pi_{\parallel}^{(2)}$ for each specimen by the equation specified in Table IX. For $n = 0$, a similar procedure was followed, using a linearly extrapolated value of $\Pi_{\perp}^{(0)}$; the alternative procedure of combining transverse ΔQ data on two different samples would lead to larger random errors.

For the 131° data the method just described for determining the moments for $n = 1$ starts to break down, since, as we noted above, the high-field limit is not reliable. However, it is probably safe to say that $\Delta Q(\infty)$ lies in between the extrapolated value of Table VI and a value increased from this in the ratio of the $\Delta\rho(\infty)$ calculated for the assumed w_1 to the $\Delta\rho(\infty)$ of the table. The first moments obtained from both these limits are plotted in Fig. 7 and connected by vertical lines to emphasize the uncertainty. A corresponding though smaller ambiguity results for the second moment ($\Pi^* - \Pi_{\parallel}^{(2)}$). The second and probably more nearly correct assumption gives the higher ($\Pi^* - \Pi_{\parallel}^{(1)}$) and the lower ($\Pi^* - \Pi_{\perp}^{(1)}$).

Fig. 7 also shows, for comparison, the electronic contribution

$$\Pi^* - \Pi_r^{(n)},$$

as given by an equation analogous to (20) on the assumption $\tau_{\perp} \propto \epsilon^{-\frac{1}{2}}$, i.e., without correction for impurity scattering. The latter correction should be very small for these specimens; some estimates of it will be given in Section VII. It will be seen that, with our neglect of impurity scattering, the electronic contribution is linear in n , and that it accounts for the major part of the variation of $\Pi_{\perp}^{(n)}$ with n ; this makes our use of linear interpolations and extrapolations seem reasonable. Subtraction of the electronic and total moments shows that $|\Pi_{p\perp}|$ is almost independent of energy, decreasing slightly with increasing energy, and that $|\Pi_{p\parallel}|$ has a more rapid, though still modest, decrease. This is a refinement of the conclusion of our first paper to the effect that some average of $|\Pi_{p\parallel}|$ and $|\Pi_{p\perp}|$ decreases with increasing energy. Since an accurate quantitative estimate of the energy dependences of $\Pi_{p\parallel}$ and $\Pi_{p\perp}$ requires a correction for impurity scattering, we shall postpone this topic to Section VII.

The last two columns of Table VIII give a sample comparison of the measured thermomagnetic quantities with values predicted from the mean moment curves of Fig. 7. The fit is, indeed, significantly better

than was obtained in the middle column on the assumption of energy-independent anisotropies.

6.4 Accuracy of the Results

A few words are in order regarding the sensitivity of these results to random or possible systematic errors in the inputs. As one might infer from the consistency of the values of $\Pi_{\perp}^{(1)}$ and $\Pi_{\perp}^{(3)}$ for different samples and from their regular trend with temperature, these quantities are rather insensitive. For germanium with its small w , Table IX shows that the quantity $\Pi_{\perp}^{(1)}$ equals $TQ(0)$ minus a small fraction of the longitudinal high-field limit of ΔQ ; since the latter term is only about one-fifth of $TQ_p(0)$, errors in it have little effect. This term is almost proportional to w . Thus, for example, it turns out that, at 77°K, an alteration of even as much as 10 per cent in the assumed w , or in the high-field ΔQ , would affect the plotted quantity $\Pi^* - \Pi_{\perp}^{(1)}$ by only 2 per cent of its value. The makeup of $\Pi_{\perp}^{(3)}$ is similar. It equals $TQ(0)$ plus a number of terms involving the Hall mobility and the low-field $\Delta\rho$'s and ΔQ 's; the latter terms are again only a small fraction of $Q_p(0)$, and they depend hardly at all on w_3 . For example, for Sample 606 at 77°K, a change of the assumed w_3 by 10 per cent would alter $\Pi^* - \Pi_{\perp}^{(3)}$ by only one-third of 1 per cent; changes of 10 per cent in any one of the quantities B , μ_H or the transverse $\Delta\rho/\rho H^2$ or $\Delta Q/H^2$ would affect $\Pi^* - \Pi_{\perp}^{(3)}$ by 1 per cent or less; a 10 per cent change in the longitudinal $\Delta\rho/\rho H^2$ or $\Delta Q/H^2$ would have a 3 per cent effect. The sensitivity for the other cases is similar. Thus, the values of $\Pi_{\perp}^{(3)}$ should be quite reliable.

Turning to the $\Pi_{\parallel}^{(n)}$, we find these quantities to be much more sensitive. The values of $\Pi_{\parallel}^{(1)}$, which depend primarily on the longitudinal high-field limit of ΔQ , should be more reliable than those for other values of n . Thus, for example, for Sample 606 at 77°K an error of 5 per cent in the high-field ΔQ would affect $\Pi^* - \Pi_{\parallel}^{(1)}$ by a little over 4 per cent. A 10 per cent change in the assumed w would affect this quantity by only about one-half of 1 per cent. On the other hand, $\Pi_{\parallel}^{(3)}$ is much more sensitive to w_3 , being in fact very nearly proportional to w_3^{-1} . It is also fairly sensitive to the [100] longitudinal $\Delta Q/H^2$ and $\Delta\rho/\rho H^2$. For example, at 77°K a 10 per cent change in either of the latter quantities would change $\Pi^* - \Pi_{\parallel}^{(3)}$ by about 9 per cent. However, the sensitivity to the other inputs mentioned in the preceding paragraph is small; a 10 per cent change in any one of them at 77°K would affect $\Pi^* - \Pi_{\parallel}^{(3)}$ by no more than 2 per cent. The quantity $\Pi^* - \Pi_{\parallel}^{(2)}$ behaves very similarly. It is very nearly proportional to w_2^{-1} ; however, the uncertainties in w_2

should be rather less than those in w_3 . Changes of 10 per cent in B or μ_H would affect $\Pi^* - \Pi_{\parallel}^{(2)}$ at 77° by about 5 per cent, and a change of 3 per cent in the interpolated $\Pi_{\perp}^{(2)}$ would affect it by 6 per cent. The values of $\Pi^* - \Pi_{\parallel}^{(0)}$, finally, are insensitive to w when derived from measurements with H in a [110] direction. The sensitivity to the assumed high-field transverse ΔQ is moderate; for Sample 603 at 77° , for example, a 10 per cent change in the latter quantity would change $\Pi^* - \Pi_{\parallel}^{(0)}$ by perhaps 9 per cent.

6.5 Other Orientations

Several conceivable calculations have been omitted from the list used in the construction of Fig. 7; these correspond to cases where the observed quantity is very insensitive to the moment which one might wish to calculate from it. Thus, as Fig. 3 shows, the longitudinal ΔQ in the [011] direction is very small and insensitive to p (i.e., to $\Pi_{\parallel}^{(1)}$). Similarly, Fig. 4 shows that the transverse ΔQ with H in an [001] direction is small and very insensitive to p (or $\Pi_{\parallel}^{(0)}$). Although no reliable points for Fig. 7 can be computed from these data, the entries in the last two columns of Table VIII show that the ΔQ 's for these orientations are in fair agreement with the predictions of the theory. (It is more significant to compare values of $Q(\infty) = Q(0) + \Delta Q$ than values of ΔQ itself, since the latter quantity could conceivably be of either sign.) The small discrepancies remaining are of the order of magnitude of the expected effect of orbital quantization.

VII. DEPENDENCE OF $\Pi_{p\parallel,\perp}$ ON ENERGY

The obvious next step is to convert the empirical moments of $\Pi^* - \Pi_{\parallel,\perp}$ plotted in Fig. 7 into moments of $-\Pi_{p\parallel,\perp}$ by subtracting the theoretical moments of $\Pi^* - \Pi_c$. One may then hope to determine the rate of variation of $\Pi_{p\parallel,\perp}$ with energy ϵ by comparing the way in which these moments vary with n with the n -variation to be expected for, say, $\Pi_{p\parallel,\perp} \propto \epsilon^l$. For a specimen completely free of impurity scattering, all this could be done very simply indeed. However, as we shall see, the second and third moments, i.e., those with weight factors τ_{\perp}^2 and τ_{\perp}^3 in (20), are extremely sensitive to impurity scattering, since this modifies τ_{\perp} greatly in the low-energy region, where it is large, and since $\Pi_{p\parallel,\perp}$ also tend to be largest at low energies. Thus, in spite of our efforts to avoid the morass of impurity-scattering theory by using very pure specimens, we must still resort to it if we are to make a complete analysis of our data. However, the smallness of the impurity scattering in the sam-

ples used for Fig. 7 will turn out to simplify greatly the task of correcting for it. We shall therefore digress for a moment to discuss the nature of impurity scattering and to establish the validity of the procedure we wish to use to take account of it.

7.1 Impurity Scattering

The Conwell-Weisskopf theory of ionized-impurity scattering,²⁰ combined with various modifications,^{21,22,23,24} has had some success in correlating mobility data and other properties of semiconductors.^{25,26*} However, a theory of this form can be rigorously justified only in the limit of carrier and ion densities rather lower than those normally encountered in semiconductors; even for this case, a rigorous treatment requires inclusion of the effects of electron-electron collisions,²⁷ a refinement which has usually been ignored in the semiconductor literature. Under conditions such as those of our experiments, there are several aspects of the theory which become rather dubious, e.g., use of the Born approximation in scattering,^{23,24} the justification for treating the screened coulomb potential of an ion as a static scattering potential, etc. Thus, about all one can safely say about ionized-impurity scattering in this range is that it is something which modifies the distribution function more and more as the energy becomes lower. Fortunately, this fact seems to be about all one needs to know in order to correlate the various effects of impurity scattering with one another, in the range where impurity scattering is slight. We shall now try to demonstrate this empirically.

Although the fact of electron-electron scattering prevents a relaxation-time model from being strictly valid for impurity scattering, we shall follow custom and assume that, in the range of interest, the effects of impurity scattering can be adequately allowed for by modifying the form of the functions $\tau_{\parallel}(\epsilon)$, $\tau_{\perp}(\epsilon)$. The usual assumption takes the form

$$\tau^{-1} = a\epsilon^{\frac{1}{2}} + b\epsilon^{-\frac{1}{2}}, \quad (22)$$

where the first term represents acoustic scattering, and the second impurity scattering. This leads to difficult integrals in the various transport expressions. Since no real justification can be given for the exact form of (22) — except at densities far below ours — we shall try using instead another expression, which has the same property of reducing τ at low

* For a general review with further references, see Ref. 26.

energies but leads to more tractable integrals, namely,

$$\tau \propto \epsilon^{-3} \left[1 - \exp\left(\frac{-\epsilon}{k\theta}\right) \right]. \quad (23)$$

Here θ is a parameter which measures the amount of impurity scattering.

7.2 Legitimacy of (23)

Our argument that (23) is justified will consist of two parts. The first will be to show that (23) gives practically the same results as (22) for a variety of problems, i.e., that when impurity scattering is just beginning to be appreciable the exact form of the impurity-scattering law is of less consequence than the mere fact that the τ 's of low-energy electrons are greatly reduced. The second argument will be empirical and will consist in a rough quantitative correlation of magnetoresistance, Hall, and mobility data, based on (23).

The curves of Fig. 8 show the first of the comparisons just mentioned. The abscissas are values of the ratio μ_H/μ_{Ha} , where μ_H is the Hall mobility computed with (22) or (23) for a simple-model semiconductor, and μ_{Ha} is the corresponding quantity for pure acoustic scattering ($b = 0$ or $\theta = 0$). The upper two sets of curves show the values of the ratio of Hall to drift mobility and of the electronic Peltier heat relative to the band edge. The full curves were computed from (23), while the dashed ones were taken from the literature^{28,29,30} based on (22).^{*} Note that the agreement is very good near $\mu_H/\mu_{Ha} = 1$, but that, as one must expect, it becomes poor for $\mu_H/\mu_{Ha} < \frac{1}{2}$. For the magnetoconductivity (bottom curve) the agreement with the dashed curve, calculated from (22) by means of the tables of Beer, Armstrong and Greenberg,³⁰ is even better. It is worth noting, incidentally, that Mansfield³¹ has also found μ_H/μ to be insensitive to the form of the impurity-scattering law for $\mu_H/\mu_{Ha} > 0.7$ or so.

Although it is encouraging to find such insensitivity of the predictions of impurity-scattering theory to the exact form of the scattering law, a comparison of the predictions of (23) with experiment is desirable as a test of its reliability. In applying (23) to n germanium we must, of course, take account of the anisotropy of the valleys, and of possible anisotropy in the impurity scattering. We shall do this in the same way as was described for the moment analysis of Section VI. Thus, we shall describe

^{*} Ref. 28 gives a general review and references to the earlier literature, Ref. 29 an accurate table of integrals and Ref. 30 the most complete tables of Hall and magnetoresistance integrals.

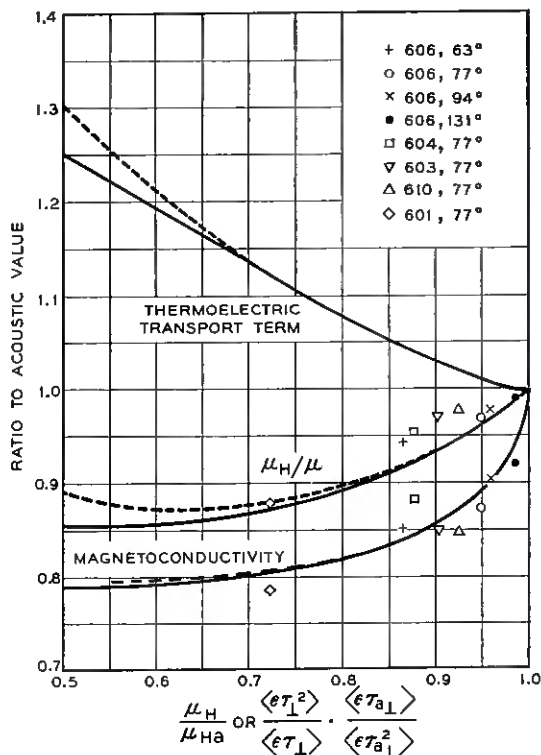


Fig. 8 — Tests of the predictions of the impurity-scattering assumption (23). The full curves represent the predictions of (23), the dashed curves those of (22), for a semiconductor with isotropic effective mass. Abscissas are ratios of Hall mobility to its value for pure acoustic scattering. Ordinates are the ratios of the values computed for the designated quantities with impurity scattering to the values with pure acoustic scattering. The top curves pertain to $\Pi^* - \Pi_s = \langle \epsilon^2 \tau \rangle / e \langle \epsilon \tau \rangle$, where Π^* is defined by (21); the acoustic value of $\Pi^* - \Pi_s$ is $2kT/e$. The middle curves pertain to $\mu_H/\mu = \langle \epsilon \tau^2 \rangle \langle \epsilon \rangle / \langle \epsilon \tau \rangle^2$; the acoustic value is $3\pi/8$. The bottom curves pertain to the low-field transverse magnetoconductivity ratio $(c/\mu_H H^2)(\Delta\sigma/\sigma_0) = \langle \epsilon \tau^3 \rangle \langle \epsilon \tau \rangle / \langle \epsilon \tau^2 \rangle^2$; the acoustic value is $4/\pi$. The points represent empirical values for the various etched specimens as obtained from the observed μ_H/μ and the magnetoresistance constant $b = (\Delta\rho/\rho H^2)_{100}^{001}$ with abscissas based on the assumption; $\mu_{Ha} = 54.3, 38.6, 28.0, \text{ and } 16.2 \times 10^3 \text{ cm}^2/\text{vs}$ at $63^\circ, 77^\circ, 94^\circ,$ and 131°K , respectively; both ordinates and abscissas represent ratios of observed quantities to values which would obtain if the w_1, w_2 and w_3 of Table X were combined with the acoustic value of τ_{\perp} . The ratios of μ_H/μ and of b to their ideal values are rather insensitive to the anisotropy of m^*/τ when this is allowed for in this manner.

the anisotropy by the w_3 , w_2 and w_1 values given in Table X, and shall take the transverse magnetoconductivity of a [100] specimen to be given by

$$-\frac{\Delta\sigma}{\sigma(\mu_H H/c)^2} = \frac{\Delta\rho}{\rho(\mu_H H/c)^2} + 1 \approx \left[\frac{(2 + w_3)(1 + 2w_3)(2 + w_1)}{3(1 + 2w_2)^2} \right] \left[\frac{\langle \epsilon\tau_{\perp}^3 \rangle \langle \epsilon\tau_{\perp} \rangle}{\langle \epsilon\tau_{\perp}^2 \rangle^2} \right], \tag{24}$$

an expression which follows from the known³ forms for $\Delta\sigma$, σ and μ_H in the electron-group approximation, and the further assumption (19). Similarly, we shall take for the Hall and drift mobilities

$$\mu_H \propto \left[\frac{(1 + 2w_2)}{(2 + w_1)} \right] \left[\frac{\langle \epsilon\tau_{\perp}^2 \rangle}{\langle \epsilon\tau_{\perp} \rangle} \right], \tag{25}$$

$$\mu \propto (2 + w_1) \left[\frac{\langle \epsilon\tau_{\perp} \rangle}{\langle \epsilon \rangle} \right], \tag{26}$$

so that the ratio of μ_H to its ideal acoustic-scattering value μ_{Ha} can be obtained from the ratio of the right of (25) to its ideal value. If now we assume (23) for τ_{\perp} , and the w_n 's of Table X, we can use the curves of Fig. 8, since the ratio of the magnetoconductivity or the Hall mobility to the factor involving w_n 's in (24) or (25) is just the value for the isotropic-mass case.

As we have no direct measurement of μ_{Ha} from which to compare (25) with (24) for a single specimen, we have adjusted this one parameter to fit the average behavior of the magnetoresistance of the purer specimens at 77.4°K. The μ_{Ha} values for other temperatures were determined from the assumption $\mu_{Ha} \propto T^{-1.65}$, this exponent having been chosen to make the fractional departure of μ_H from μ_{Ha} scale a little more slowly than T^{-3} , as one expects from the Conwell-Weisskopf law. Using the ratio of the observed μ_H to this μ_{Ha} and the w_1 and w_2 of Table X, we can determine the ratio $\langle \epsilon\tau_{\perp}^2 \rangle \langle \epsilon\tau_{a\perp} \rangle / \langle \epsilon\tau_{\perp} \rangle \langle \epsilon\tau_{a\perp}^2 \rangle$ of the last factor in (25) to the value it has when $\tau_{\perp} = \tau_{a\perp}$, the value for pure acoustic scattering. With this ratio as abscissa, we have plotted, on the graphs of Fig. 8, the ratios of empirical value to acoustic value for the last factor in (24) and for the ratio of the last factors in (25) and (26), as obtained from the ratio of the low- and high-field Hall constants. Only data from etched specimens have been used, as extreme accuracy is required for a significant comparison; even so, there is a fair amount of scatter in the points. Much of this may be due to errors in μ_H resulting from errors — perhaps 1 per cent — in the assumed dimensions, or from the fluctuations in impurity

density; such errors will displace the μ_H/μ points horizontally and the magnetoconductivity points along a line sloping at about 45° from upper left to lower right. However, the curves account at least roughly for the at first sight surprisingly low magnetoconductivity of the pure specimens, and for the further reduction of magnetoconductivity and μ_H/μ found for the less pure specimen (Sample 601).

In concluding this digression, a few words are in order regarding the validity of the assumption that nothing except ionized-impurity scattering interferes appreciably with the ideal acoustic scattering law $\tau \propto \epsilon^{\frac{1}{2}}$. Optical-mode or intervalley scattering has sometimes been suggested as one of several possible causes of the departure of the mobility from a $T^{-\frac{1}{2}}$ law. Though such scattering freezes out at low temperatures, it is easily verified¹³ that the amount of it which would be needed to account entirely for the mobility exponent could raise μ_H/μ by 3 per cent or so at 131°K , depress it by perhaps 1 per cent at 77° , and lower the absolute value of the magnetoconductivity by 1 per cent or so. The fact that, for Sample 606 at 131° , the ratio μ_H/μ is 0.992 of the acoustic-scattering value suggests that this and other departures from ideal acoustic-scattering behavior are small.

7.3 Behavior of the Moments of $\Pi_{p\parallel, \perp}$

Fig. 9 shows the way in which the various moments of Π_p , as defined by (20) for either component, would be affected by impurity scattering of the type (23), if Π_p were $\propto \epsilon^{-\frac{1}{2}}$ or $\epsilon^{-\frac{1}{4}}$. Impurity scattering, of course, does not affect the zeroth moment, and if Π_p is independent of energy it does not affect any moment. We have chosen as abscissa the same ratio of Hall mobility to ideal acoustic mobility which was used in Fig. 8. Fig. 10 shows, analogously, the effect of impurity scattering on the second and third moments of the quantity $(\Pi^* - \Pi_e) \propto \epsilon$; the effect on the first moment has already been shown in Fig. 8. Comparison of Figs. 9 and 10 with Fig. 8 shows that, for even our purest samples at 77°K , the effects of impurity scattering on the moments with $n = 2$ and 3 can be considerable if Π_p varies as rapidly as $\epsilon^{-\frac{1}{2}}$. This is why we have undertaken such an elaborate discussion of these effects.

The first step in analyzing the data of Fig. 7 is to correct $\Pi^* - \Pi_e$ for the effects of impurity scattering and subtract it from the curves of Fig. 7 to get curves representing the moments of $\Pi_{p\parallel}$ and $\Pi_{p\perp}$, averaged over the purer specimens. To describe the average behavior of the purer specimens, we have chosen the values $\langle \epsilon \tau_{\perp}^2 \rangle \langle \epsilon \tau_{\perp} \rangle / \langle \epsilon \tau_{\perp} \rangle \langle \epsilon \tau_{\perp}^2 \rangle = 0.89$ at 60° , 0.93 at 77.4° , 0.96 at 94° and 0.98 at 131° . These choices, though

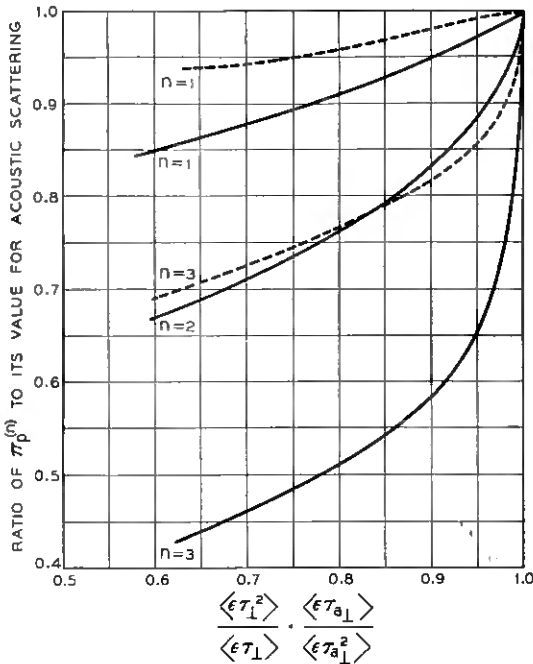


Fig. 9 — Effect of impurity scattering on the moments $\Pi_p^{(n)}$ if Π_p is assumed $\propto \epsilon^{-1/2}$ (full curves) or $\propto \epsilon^{-1/4}$ (dashed curves). Abscissas are the ratio of Hall mobility to its value for same anisotropy parameters w_1 and w_2 , but $\tau_{\perp} \propto \epsilon^{-1/2}$.

reasonable enough, differ slightly from other equally reasonable choices which one might make from Fig. 8; we shall discuss the sensitivity of the results to the choice in the next paragraph. From these abscissas in Fig. 10 and the top curve of Fig. 8 we have determined the corrections to $\Pi^* - \Pi_e^{(n)}$. The resulting values of $\Pi_{p\perp}^{(n)}$ and $\Pi_{p\parallel}^{(n)}$ are shown in Fig. 11 as full curves. The dashed curves are the moments to be expected for $\Pi_{p\parallel,\perp} \propto \epsilon^0, \epsilon^{-1/4},$ or $\epsilon^{-1/2}$, as determined from the values of μ_H/μ_{H0} just given and the curves of Fig. 9; the dashed curves have been fitted to the full curves at $n = 1$.

A word about the sensitivity of the results to the choice of abscissas in Figs. 8 to 10 is in order. Shifting the abscissa 0.01 unit at 77° has a negligible effect on the moments with $n = 0$ or 1. It shifts the full and dashed curves of Fig. 11 in the same direction for $n > 1$, the dashed curves being shifted much the greater amount; for $n = 3$, it shifts the ϵ^{-3} dashed curves by about 0.02 volt for $\Pi_{p\parallel}$, or 0.002 volt for $\Pi_{p\perp}$. At 131°K , the effect is more serious; a decrease of the assumed abscissa from 0.98 to 0.97, though affecting the $n = 1$ moments and $\Pi_{p\parallel}^{(3)}$ very

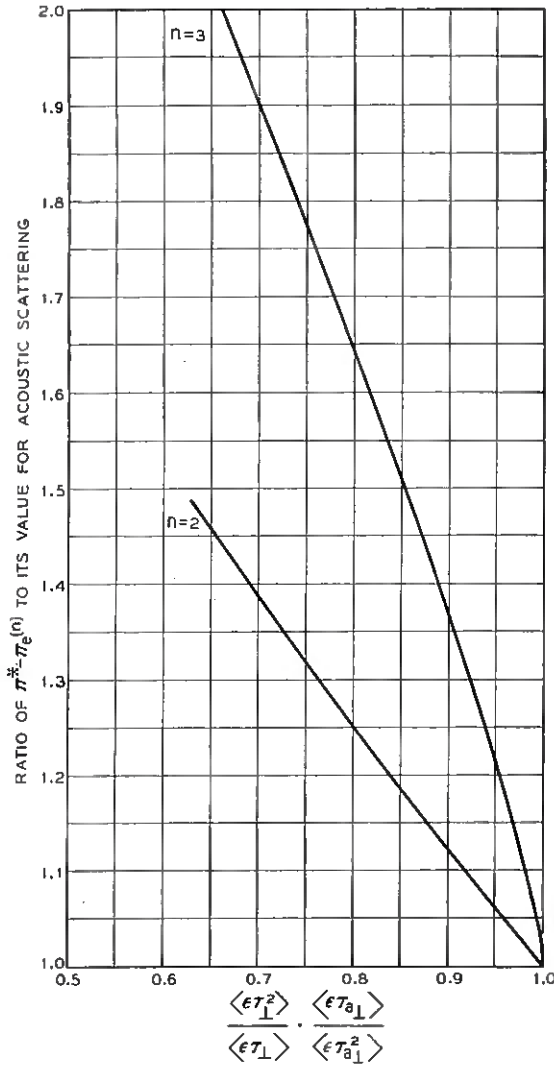


Fig. 10 — Effect of impurity scattering on the moments $\Pi^* - \Pi_{\perp}^{(n)}$. Abscissas are the ratio of Hall mobility to its ideal acoustic value for the same anisotropy parameters w_1 and w_2 , but $\tau_{\perp} \propto \epsilon^{-1/2}$.

little, would lower the empirical $-\Pi_{p\perp}^{(3)}$ by about 0.001 volt, and would lower the corresponding ordinate of the $\epsilon^{-1/2}$ curve for $-\Pi_{p\parallel}^{(3)}$ by 0.09 volt.

As is to be expected from our experience in constructing Fig. 7, the results for $\Pi_{p\perp}^{(n)}$ are the more consistent. At all four temperatures the

empirical $\Pi_{p\perp}$ curves lie between those for $\Pi_{p\perp} \propto \epsilon^0$ and ϵ^{-1} , and nearer to the former. By interpolation we estimate $\Pi_{p\perp} \propto \epsilon^{-0.08}$, the exponent being within 0.02 or so of this value at all four temperatures. The curves for $\Pi_{p\parallel}^{(n)}$ have, for $n > 1$, roughly the behavior to be expected for $\Pi_{p\parallel} \propto \epsilon^{-0.25}$ (77° and 94°K), $\epsilon^{-0.35}$ (60°) and $\epsilon^{-0.50}$ (131°K); between $n = 0$ and $n = 1$ the slope is steeper, corresponding on the average to,

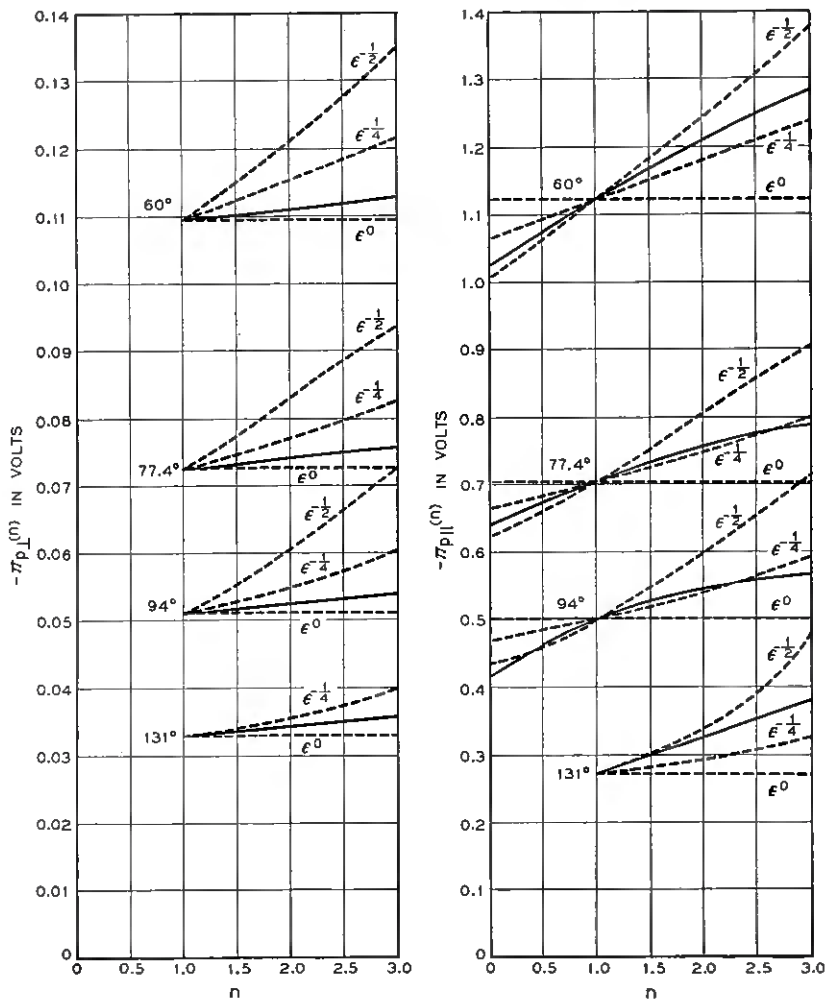


Fig. 11 — Moments of $\Pi_{p\perp}$ and $\Pi_{p\parallel}$ (full curves), as determined from Fig. 7 and the values of $\Pi^* - \Pi_e$ corrected for impurity scattering. The dashed curves show the behavior which would be expected if $\Pi_{p\perp}$ were $\propto \epsilon^0$, $\epsilon^{-1/4}$ or $\epsilon^{-1/2}$, with allowance for impurity scattering.

TABLE XI

Computed fractional errors in the assumption that Q is the sum of the ideal Q_e and Q_p values which would occur in the absence of impurity scattering. Here $\delta |Q| = |Q|_{\text{true}} - |Q|_{\text{ideal}} = \delta |Q_e| + \delta |Q_p|$ and $\delta |Q_p| = \delta_w |Q_p| + \delta_1 |Q_p|$, where $\delta_w |Q_p|$ arises from the departure of $w_1 = m^*_{\perp}(\epsilon_{T\parallel})/m^*_{\parallel}(\epsilon_{T\perp})$ from its acoustic-scattering value, and $\delta_1 |Q_p|$ arises from the corresponding departures of the $\Pi_{p\parallel}^{(1)}$, $\Pi_{p\perp}^{(1)}$ defined by (20).

Sample	Estimated Abscissa for Figs. 8 & 9	T , °K	$\delta Q_e $, $\mu\text{V}/\text{deg.}$	$\delta_w Q_p $, $\mu\text{V}/\text{deg.}$	$\delta_1 Q_p $, $\mu\text{V}/\text{deg.}$	$\frac{\delta Q }{ Q }$
606	0.91	60°	+5	$\leq +5$	-25	≈ -0.005
	0.99	122°	0	0	0	0.000
603	0.93	77°	+3	$\leq +5$	-4	$\leq +0.002$
601	0.75	94°	+18	+16	-17	+0.010
596	0.56	87°	+38	+57	-30	+0.041

say, $\epsilon^{-0.4}$ at the three lower temperatures. Considering these variations and the random scatter of the points in Fig. 7, we can state as our best estimate that $\Pi_{p\parallel} \propto \epsilon^{-0.3 \pm 0.05}$ at low energies and that its energy dependence probably becomes more rapid at higher energies. There is no reliable evidence that the energy dependence of either component changes appreciably in the temperature range 60° to 131°K, although the accuracy of the determination of energy dependence is not very good for $\Pi_{p\parallel}$. The ratio $\Pi_{p\parallel}^{(1)}/\Pi_{p\perp}^{(1)}$ seems to be decreasing slightly with increasing temperature, although here again it is hard to be sure that the effect is real.

7.4 Effect of Impurity Scattering on the Assumed Q

As we have mentioned in Section V, all our thermomagnetic quantities were measured relative to the value of Q at $H = 0$, and the values originally tabulated for them were based on the value of Q characteristic of the observed carrier concentration and pure acoustic scattering. If impurity scattering causes the total Q to depart appreciably from this assumed value, these values will require correction. The results obtained in this section allow us to estimate how much effect impurity scattering has on Q_e and Q_p , and so to check this point.

According to Table IX, impurity scattering can affect Q_p in two ways: through changing the anisotropy $w (=w_1)$ of $m^* \cdot \tau^{-1}$, and through its effect on the moments $\Pi_{p\parallel}^{(1)}$, $\Pi_{p\perp}^{(1)}$. The effect on Q_e is, of course, simply obtained from the top curve of Fig. 8. Table XI shows some sample eval-

uations of the two kinds of effects on Q_p and of the effect on Q_e , computed from Figs. 8 and 9, the w 's of Table X, and the assumptions $\Pi_{p\parallel} \propto \epsilon^{-0.30}$, $\Pi_{p\perp} \propto \epsilon^{-0.08}$. It will be noted that the change in w largely compensates the changes in $\Pi_{p\parallel, \perp}$ ⁽¹⁾; this explains the previously puzzling lack of sensitivity of Q_p to impurity scattering. The fractional errors $\delta Q/Q$ in our assumed Q values are seen to be at most a fraction of a per cent for the purer specimens; this implies errors in the ordinate of Fig. 7 of rather less than a per cent. For Samples 601 and 596, however, at the temperatures shown, $\delta Q/Q$ is 1 per cent and 4 per cent, respectively, and $\delta\Pi/(\Pi^* - \Pi)$ is about twice as great.

7.5 Confirmation of Impurity-Scattering Theory on Sample 601

We have based our corrections for impurity scattering in the purer samples on plausible theoretical reasoning and on favorable though not ideally consistent evidence from isothermal electrical measurements. It is therefore of some interest to see if the thermomagnetic measurements on a sample with a sizable amount of impurity scattering differ from those on the pure samples, in the way in which we would predict on the basis of (23) and Figs. 8 to 10, assuming the energy dependences of $\Pi_{p\parallel}$ and $\Pi_{p\perp}$ which we inferred from Fig. 11.

Fig. 12 shows the comparison just mentioned for Sample 601 at 94°K. This sample was chosen for the comparison because good isothermal data were available for it and because it had about the maximum amount of impurity scattering for which the predictions of (23) are likely to be reliable. (For some of the other samples, the electrical and thermomagnetic data were at different temperatures, and for 596 at 87° the divergence between the full and dashed curves of Fig. 7 is beginning to be appreciable.) The dashed curves in Fig. 12 are the curves of Fig. 7, and represent the mean moments of the total $\Pi^* - \Pi_{p\perp}$ for the purer samples. The plotted points represent the empirical moments for Sample 601, computed from the data of Table VI in the same way as for the pure samples of Fig. 7. The small temperature-gradient correction just discussed is, of course, included in the corrections of Table VI. These points are to be compared with the full curves, which were constructed by computing the differences $\Pi_{p\perp}^{(n)}(601) - \Pi_{p\perp}^{(n)}(\text{ideal})$ theoretically and combining them with the dashed curves. The differences in question were computed by the following steps: (a) From (25), the observed μ_H of 601, the empirical w_1 and w_2 of Table X, and the assumed acoustic value $\mu_{H\alpha}(94^\circ) = 28,200 \text{ cm}^2/\text{vs}$ (probably good to within a few per cent), we computed the abscissa appropriate to Sample 601 in Figs. 8 to

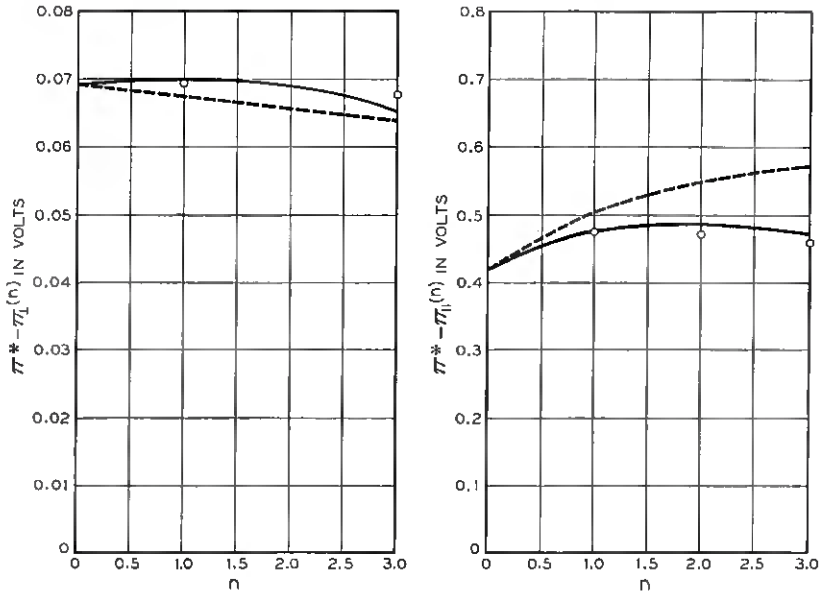


Fig. 12 — Comparison of the empirical moments $\Pi^* - \Pi_{\parallel, \perp}^{(n)}$ for Sample 601 (points) with those of the purer samples of Fig. 7 (dashed curves) at 94°K. The full curves are moments computed from the latter, with allowance for the amount of impurity scattering present for Sample 601.

10. (b) From these figures and the assumptions $\Pi_{p\parallel} \propto \epsilon^{-0.30}$, $\Pi_{p\perp} \propto \epsilon^{-0.08}$, the ratios of $\Pi_{p\parallel, \perp}^{(n)}$ (601) to $\Pi_{p\parallel, \perp}^{(n)}$ (ideal) were obtained. (c) The same ratios were obtained for the average of the pure specimens, assuming an abscissa of 0.97 in the figures. (This is larger than corresponds to the above μ_{H0} and the mobilities measured on the unetched samples, but is more reasonable in the light of the 77° data and the effect of etching at that temperature.) (d) From steps (b) and (c) and the empirical $\Pi_{p\parallel, \perp}^{(n)}$ of Fig. 11 the $\Pi_{p\parallel, \perp}^{(n)}$ to be expected for Sample 601 were calculated. (e) The theoretical $\Pi^* - \Pi_e^{(n)}$ were calculated for Sample 601 from Figs. 8 and 10, and combined with the $\Pi_{p\parallel, \perp}^{(n)}$ just obtained.

It will be seen that, for $\Pi_{\parallel}^{(n)}$ especially, the empirical points fall much closer to the curve computed with allowance for impurity scattering than to the curve for the pure samples. Since the theoretical difference between the two curves is proportional to the rate of variation of Π_{\parallel} with energy and to the departures of the ordinates of Figs. 8 to 10 from unity, the agreement is a significant bulwark for the analysis of this section. However, it is of more qualitative than quantitative significance, for the inputs to the calculation were numerous and some of them a little

uncertain. The discrepancy for $\Pi_{p\perp}^{(3)}$ may well be due to the great sensitivity of the $\Pi^* - \Pi_c^{(3)}$ curve of Fig. 10 to the abscissa chosen.

VIII. FURTHER OBSERVATIONS SUBSTANTIATING AND EXTENDING THE PRESENT MODEL

In this section we shall discuss several parts of our data which are not adapted to as precise an analysis as that of Sections VI and VII, but which can still be semiquantitatively explained and give indications of how well our conclusions can be extended to higher energies and higher temperatures.

8.1 *The High-Field Asymptote of BH*

So far, we have made no use of our observations of the Nernst coefficient at large magnetic fields. If the asymptotic high-field behavior were accurately known, it could provide a valuable extension of the moment analysis of Section VI. For, as (86) of Appendix B shows, this asymptotic behavior involves the moments (20) for $n = 0$ and -1 , so that measurements in two orientations could be used to determine $\Pi_{\parallel}^{(-1)}$ and $\Pi_{\perp}^{(-1)}$. In practice, however, the accuracy of such a determination would be very questionable, as it would depend on the already rather uncertain zeroth moments and on accurate knowledge of the asymptotic BH , which, as we shall presently see, is rather uncertain for various experimental reasons. We have therefore chosen merely to compare the high-field behavior of BH with the asymptotic behavior predicted by the special-case formulas of Table II, with the object of showing that there is at least fair agreement with the model previously deduced.

As we have noted in Section V, the small value of the Nernst constant at high fields makes it very sensitive to falsification by the Hall fields of the longitudinal counter-currents which will be set up if surface conduction partially short-circuits the thermoelectric field; a surface conductance of a fraction of one per cent can be serious. We did not realize this soon enough to carry out extensive Nernst measurements on etched specimens. The lowest temperature at which such measurements were made was 94°K. At this temperature the difference between the behavior of Sample 606 etched and unetched was profound, the unetched having about twice the Nernst field of the etched at 18 kilogauss. At higher temperatures, the surface effect gets even worse, because the counter-currents are driven by the full Q , including the term proportional to the distance of the Fermi level below the band edge, whereas the Nernst effect arises only from Q_p and the small transport term in Q_s .

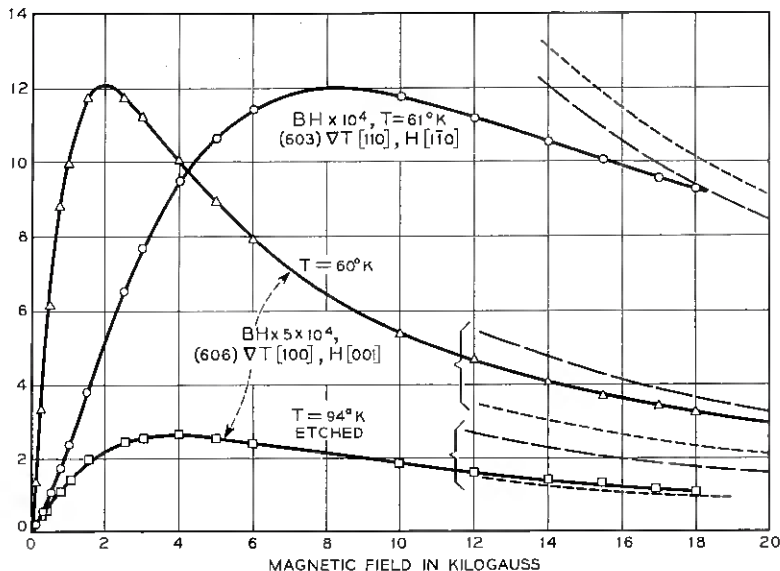


Fig. 13 — Comparison of the asymptotic behavior of the Nernst field constant, BH , at high fields with the predicted $(B_s + B_p)H$ of the simplified asymptotic formulas of Table II. The three experimental curves are for two samples of different orientations; note that the two scales are different. The short-dashed curves are computed from the formulas of Table III for Π_p independent of energy, the long-dashed curves for $\Pi_p \propto \epsilon^{-1/2}$. For both sets, the values $w = 0.061$, $p = 9.8$ were assumed. The unit of BH is volts/degree.

Fig. 13 shows a comparison of three observed curves of BH against H with the leading (H^{-1}) term in the asymptotic expression calculable from Table II. The data for Sample 606 at 94°K are the most reliable, since this sample has been etched; however, these data do not extend to as high a value of $\mu H/c$ as is obtainable at lower temperatures, so there may be some doubt as to how well they approach asymptotic behavior at 18 kilogauss. The data for Sample 606 at 60° , though on an unetched specimen, appear not to have been very seriously falsified by surface conduction, since they match the 94° data pretty well — as far as they go — if ordinates are scaled proportionally to Q_p and abscissas to μ , and if a little allowance is made for the electronic contribution. However, at the high-field end, the 60° data lie closer to the theoretical asymptote for $\Pi_p \propto \epsilon^{-1/2}$ than to that for Π_p independent of energy, whereas the reverse is the case for the 94° data; this may be due to a slight raising of the experimental curve by surface conduction.

The data for Sample 603 — with H in a $[110]$ direction — show the

vastly slower saturation which is to be expected from the fact that, for this orientation, half the valleys have a very large cyclotron mass. Since the curves for the two samples are plotted on different scales, the values of BH at 18 kilogauss and 60° – 61°K actually differ by a factor of 14 in the two orientations, roughly the factor predicted by the theory. As nearly as one can estimate from the range available, the experimental curve seems to be heading for an asymptote closer to that for Π_p independent of energy than to that for $\Pi_p \propto \epsilon^{-1}$ — perhaps even the other side of it. This is a slight discrepancy with theory which might again be due to the presence of a little surface conduction.

Thus, it appears that the qualitative features of the high-field Nernst curves are accounted for by the theory. The curves also give some information on the value of the anisotropy $p = \Pi_{p\parallel}/\Pi_{p\perp}$. Since the entries in the next to the last column of Table II are almost proportional to $(p - 1)$, it is probably safe to say that the present curves indicate a p within 20 per cent or so of the value derived in the last section. Regarding the variation of the weighted average of $\Pi_{p\parallel}$ and $\Pi_{p\perp}$ with energy, about all that can be said is that this variation again appears not to be as rapid a decrease as ϵ^{-1} .

8.2 Reversal of the Sign of B with H

All our analysis so far has been based on data at temperatures in the range 60° to 131°K . No comparably accurate analysis can be made at higher temperatures, partly because the limits of the various quantities as $H \rightarrow \infty$ can no longer be determined, and partly because the value of Q_p is much less accurately known than in the lower range. What evidence there is, however, suggests that, at temperatures in the range 150° to 235°K , the anisotropy and energy dependence of Π_p are similar to those which we have found in the liquid-air range. This evidence is provided partly by the magnitudes of the low-field B and $\Delta Q/H^2$, which we shall discuss presently, and partly by an interesting reversal of the sign of B with increasing magnetic field. The latter effect occurs at temperatures at which the negative $B_e(0)$ and the positive $B_p(0)$ almost cancel each other. This cancellation occurs near 175°K , $B(0)$ being greater than zero below this temperature, less than zero above it¹ (although easily falsified by surface conduction). Fig. 14 shows some curves of B versus H at temperatures in this range, taken on etched samples. The data previously reported¹ on 606 and 604 unetched showed no region of negative $B(0)$, but now both show sign reversals at 171° to 174° . What we wish to discuss now, however, is the fact that, as H increases,

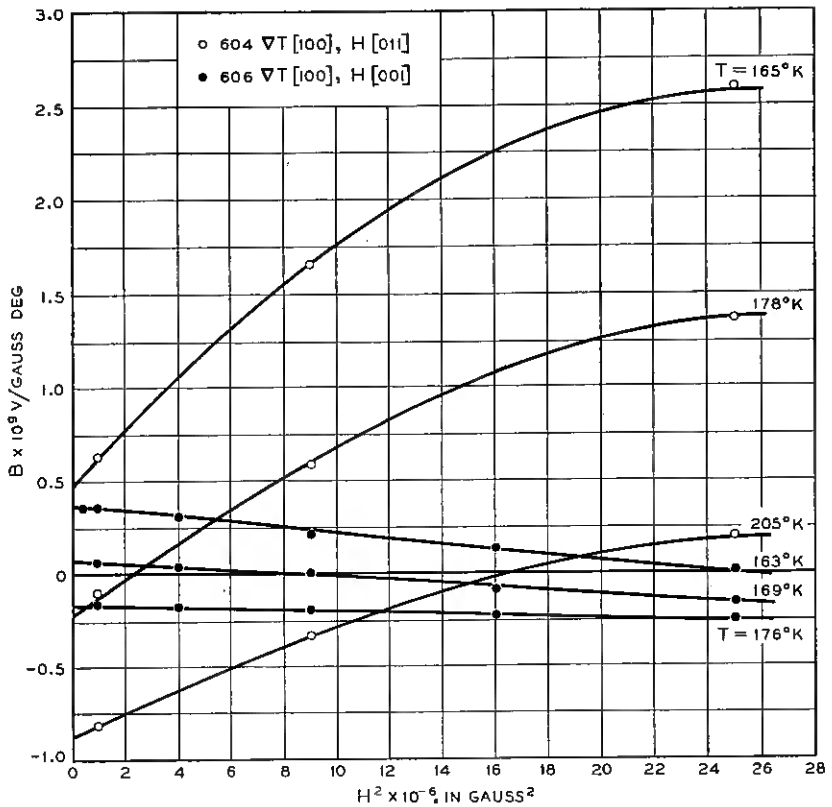


Fig. 14 — Variation of the Nernst coefficient B with magnetic field H , for two orientations, showing the sign reversal which takes place at temperatures near where B_s and B_p compensate each other. For comparison, the computed $B_s(178^\circ)$ is $-4.0 \times 10^{-9} \text{ v/deg gauss}$.

the B of Sample 604 goes from negative to positive, and that of 606 from positive to negative.

Although we have not elaborated the theory of B at intermediate fields, it is doubtless safe to assume that, for all field strengths, $B_p > 0$ and $B_s < 0$. We can get a fair idea how the ratio of B_p to B_s varies with H by comparing the low- and high-field limits of $|B_p/B_s|$ as given by the formulas of Tables I and II. Table XII gives the comparison, for several sets of assumptions about the anisotropy and energy dependence of Π_p . We have constructed schematic curves of B_s and B_p against H , shown in Fig. 15, from this information and the knowledge¹ that B_p decreases monotonically with H for $\mathbf{H} \parallel [001]$, while it has an initial

TABLE XII

Comparison of high- and low-field ratios of B_p (positive) to B_e (negative), for several orientations and for different assumptions regarding the anisotropy $p = \Pi_{p\parallel}/\Pi_{p\perp}$ and the dependence of Π_p on energy ϵ . The numerical values given are for $w = m^*_{\perp}\tau_{\parallel}/m^*\tau_{\perp} = 0.06$.

Direction of		Ratio of $\left \frac{B_p(\infty)}{B_e(\infty)} \right $ to $\left \frac{B_p(0)}{B_e(0)} \right $, for			
\mathbf{H}	νT	Π_p independent of ϵ , any p	$\Pi_p \propto \epsilon^{-1/2}$		
			$p = 5$	$p = 10$	$p = 20$
[001]	$\perp \mathbf{H}$	1	0.45	0.50	0.54
[011]	[100]	$\frac{1 + 2w}{3w} = 6.2$	1.46	2.06	2.58
[011]	[011]	$\frac{(2 + w)(1 + 2w^2)}{3w(1 + 7w + w^2)} = 8.1$	2.27	3.31	4.17

rise for $\mathbf{H} \parallel [011]$, if $p = \Pi_{p\parallel}/\Pi_{p\perp}$ is large. From these curves it is clear that, if the anisotropy and energy dependence of Π_p are similar at the temperatures of Fig. 14 to what we found in Fig. 11, then B should become more negative with increasing H for Sample 606, more positive for 604. This is exactly what is observed.

Only the crudest sort of quantitative limits can be derived, from this comparison, for p and the rate of variation of Π_p with energy. From an inspection of Table XII, Fig. 15 and the curves of Fig. 9 of Ref. 1, it seems clear that the sizes of the shifts of B with field in Fig. 14 are eminently consistent with the previously found anisotropy and energy dependence, and that they would be hard to reconcile with a Π_p that increased with increasing energy, or with a p close to unity. Specifically, Table XII suggests that the theoretical curve for 606 is nearly independent of p , so that the downward slopes in Fig. 14 must arise from an inverse energy dependence of Π_p . The swing in B between zero and 5000 gauss is of the order of a twentieth of $|B_e(0)|$; this would not be unreasonable for $\Pi_p \propto \epsilon^{-1/2}$. The swing for 603, on the other hand, is probably due to the fact that B_p rises with H while B_e falls; if so, it should be of the order of one or two times the rise in B_p , a rise due mainly to the anisotropy of $\Pi_p (p \gg 1)$. For the observed anisotropy at liquid-air temperature ($p \approx 10$) this rise (see Fig. 9 of Ref. 1), at the $\mu_H H/c$ corresponding to 5000 gauss at 178° , is perhaps 15 per cent of $B_p(0)$. The observed swing is about twice this.

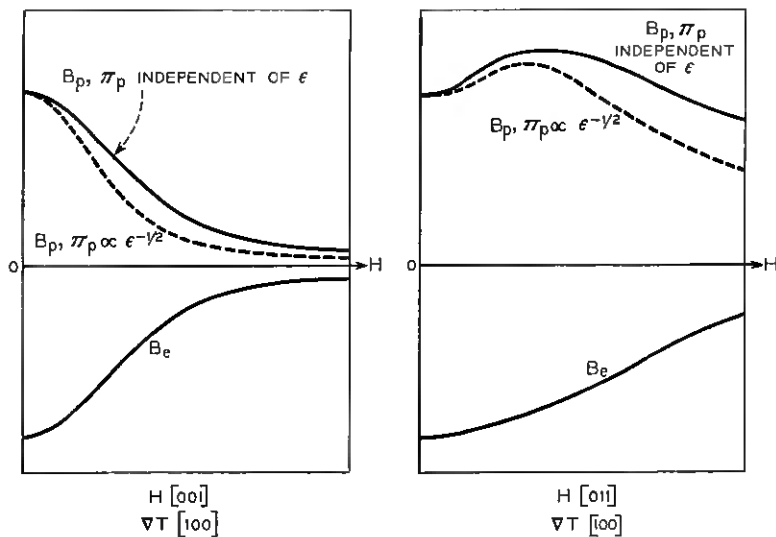


Fig. 15 — Schematic variation of B_e and B_p with magnetic field, for two orientations, assuming the temperature such that $B_e(0) = -B_p(0)$. Two B_p curves are shown for each case, corresponding to $\Pi_p \propto \epsilon^0$ and $\epsilon^{-1/2}$, respectively. It is assumed that $\Pi_{p\parallel}/\Pi_{p\perp} \gg 1$.

8.3 Analysis of High-Temperature Data

The sort of moment analysis which we made in Figs. 7 and 11 gets increasingly difficult as T increases, for several reasons: high-field limits become unavailable; the phonon-drag effects become a much smaller part of the total; because of this, it is easier for such things as surface inversion layers to falsify measurements of bulk thermomagnetic effects. We shall nevertheless try to draw what inferences we can from the temperature at which $B(0)$ changes sign, and from the totality of measurements on Sample 601 at 234°K, a sample pure enough to be fairly free of impurity scattering at this temperature, yet sufficiently highly conducting to be insensitive to surface conduction.

As Fig. 14 shows, $B(0)$ reverses sign in the pure, etched specimens at 174°K (603) or 172°K (606). From the average of these and interpolations of the Q_p values of Table VI we can deduce a value for the ζ_p of (13). With the acoustic-scattering value of B_e , we obtain $\zeta_p = 0.226$, a value close to the average ζ_p of 0.229 found for the pure specimens at 77°K, corrected for the effect of impurity scattering on B_e . This suggests either that the anisotropy and energy variation of Π_p are very similar at 173° and 77° or that the difference in anisotropy just happens to be such as to compensate for the difference in energy dependence. The 131°

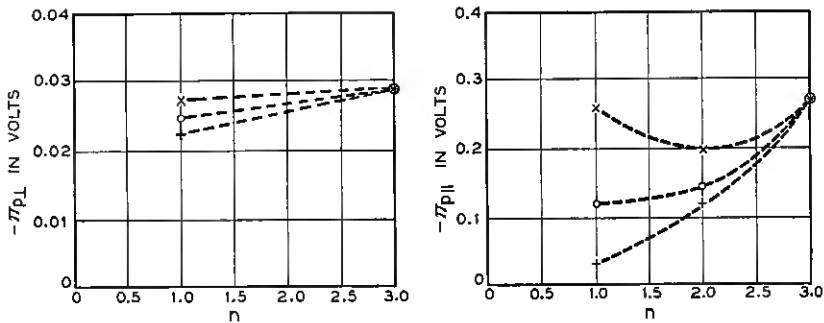


Fig. 16 — Results of an attempt to determine empirically the moments $\Pi_{p||,\perp}^{(n)}$ defined by (20) for Sample 601 at 234°K. The points O, + and X correspond respectively to three choices of the assumed $\Pi_{p\perp}^{(1)}$; it will be seen that the O choice gives the most nearly plausible behavior for $\Pi_{p||}^{(n)}$.

data of Fig. 11 suggested a possible decrease of the anisotropy with increasing temperature, and a possible increase in the rapidity of variation of $\Pi_{p||}$ with energy. These factors would affect ζ_p in opposite ways, so they could cancel.

We have attempted a moment analysis like that of Sections VI and VII for Sample 601 at 234°K. We adopted the plausible anisotropy $w_1 = w_2 = w_3 = 0.056$, (see Table X and the discussion of Section VI, especially the remarks about change of w with temperature). This is a value which happens to fit the anisotropy of the magnetoresistance. The data give $\Pi^* - \Pi_{\perp}^{(3)} = 0.0507$ volt, $\Pi^* - \Pi_{||}^{(3)} = 0.298$ volt. The only additional knowledge obtainable, however, is one relation between $\Pi_{||}^{(1)}$ and $\Pi_{\perp}^{(1)}$ (to fit Q) and one relation between $\Pi_{||}^{(2)}$ and $\Pi_{\perp}^{(2)}$ (to fit B). If, as before, we assume $\Pi_{\perp}^{(n)}$ to be linear in n , we can adjust $\Pi_{\perp}^{(1)}$ until the curve of $\Pi_{||}^{(n)}$ against n is reasonably smooth. Fig. 16 shows the results for $\Pi_{p||,\perp}^{(n)}$ after subtracting off the values of $\Pi^* - \Pi_e^{(n)}$, slightly corrected for impurity scattering. Several choices of $\Pi_{p\perp}^{(1)}$ have been made, and their consequences are shown in the figure. There are strong suggestions that $\Pi_{p||}^{(1)}/\Pi_{p\perp}^{(1)}$ is still smaller than at 131°K, that $\Pi_{p||}$ varies more rapidly with energy than at 77° and that $\Pi_{p\perp}$ continues to vary slowly. However, the data are so limited that only the last of these points can be considered reasonably established.

IX. SUMMARY AND IMPLICATIONS OF THE RESULTS

9.1 Evidence on $\Pi_{p||}(\epsilon)$ and $\Pi_{p\perp}(\epsilon)$

We have concluded that, near liquid-air temperature, the dependence of the phonon-drag Peltier tensor of an energy shell on the energy ϵ of

the shell is given, at least roughly, by $\Pi_{p\parallel}(\epsilon) \propto \epsilon^{-0.3}$, $\Pi_{p\perp}(\epsilon) \propto \epsilon^{-0.08}$, while the first-moment ratio $\langle \epsilon^{\frac{1}{2}} \Pi_{p\parallel} \rangle / \langle \epsilon^{\frac{1}{2}} \Pi_{p\perp} \rangle$ is about 9.6. These conclusions were based primarily on the analysis of high- and low-field ΔQ data and low-field Nernst data on the several purest samples, with auxiliary use of mobility and magnetoresistance data. The use made of the low-field data depended to some extent on corrections for impurity scattering, the method of making the corrections having been roughly checked both against purely electrical measurements and against alternative theoretical assumptions. The data for Sample 601, which showed a sizable effect of impurity scattering, were consistent with the assumed characteristics of $\Pi_p(\epsilon)$ combined with the assumed impurity-scattering law. The conclusions regarding $\Pi_{p\perp}(\epsilon)$ were very insensitive to the more uncertain of the inputs and were consistent from sample to sample and from one temperature to another. Those regarding $\Pi_{p\parallel}(\epsilon)$ were more sensitive and showed more fluctuation, but always showed roughly the same behavior. High-field Nernst data, though less suited to accurate analysis, independently indicated the ratio of $\Pi_{p\parallel}$ to $\Pi_{p\perp}$ to be of the order of 10, and also indicated a decrease of the weighted average of these quantities with ϵ , slower than $\epsilon^{-\frac{1}{2}}$.

The high-field limit of the transverse ΔQ 's gave strong but not entirely certain indications that the decrease of $\Pi_{p\parallel}$ with increasing energy becomes more rapid than the above-mentioned $\epsilon^{-0.3}$ at high energies. The variation of the anisotropy and energy dependence of Π_p with temperature was too small to be detected over the range 60° to 94°K. However, we found some rather uncertain indications that $\Pi_{p\parallel}/\Pi_{p\perp}$ decreases slowly with increasing T , and that the energy dependence of $\Pi_{p\parallel}$ becomes steeper. These conclusions came from a comparison of 131° thermomagnetic data with data at lower temperatures, and from studies of low-field thermomagnetic data at 234°. The low- and intermediate-field behavior of the Nernst coefficient near 175°K was found to be consistent with the anisotropy and energy dependences found from the previous sources, but did not suffice to check the temperature variation of these two factors individually. The data at all temperatures are clearly inconsistent with a law $\Pi_{p\parallel} \propto \epsilon^{\frac{1}{2}}$, such as would result (see below) if the relaxation times of the phonons were independent of frequency.

9.2 Implications for Phonon-Phonon Scattering

The energy dependence of $\Pi_p(\epsilon)$ reflects the dependence of the relaxation times τ_{ph} of the different phonon modes on their frequency ω . For example, if all branches of the phonon spectrum had

$$\tau_{ph} = \omega^i \times (\text{function of direction}), \quad (27)$$

then, since the wave numbers, or frequencies, of the modes coupled to electrons of energy ϵ scale proportionally to $\epsilon^{\frac{1}{2}}$, we would have, as ϵ varies

$$\begin{aligned} \Pi_p &\propto \tau^{-1} \times (\text{a mean } \tau_{ph}) \\ &\propto \epsilon^{\frac{3}{2} + t}. \end{aligned} \quad (28)$$

If the longitudinal and transverse branches were to obey (27) individually, but with different exponents t_l and t_t , then the effective energy exponents for $\Pi_{p\parallel}$ and $\Pi_{p\perp}$ would be intermediate between $\frac{1}{2} + \frac{1}{2}t_l$ and $\frac{1}{2} + \frac{1}{2}t_t$, and the exponents would in general be different for $\Pi_{p\parallel}$ and $\Pi_{p\perp}$.

Our liquid-air-temperature result is thus that the phonons interacting with a current parallel to the valley axis have an average relaxation time $\tau_{ph} \propto \omega^{-1.6}$, while those interacting with a current in the perpendicular direction have an average $\tau_{ph} \propto \omega^{-1.15}$. These are to be compared with the "ideal" phonon-phonon scattering laws^{2,32}

$$\text{longitudinal branch: } \tau_{ph} \propto \omega^{-2}, \quad (29)$$

$$\text{transverse branch: } \tau_{ph} \propto \omega^{-1}, \quad (30)$$

valid if $\omega \ll$ the frequency of the average thermal phonon and if only three-phonon noncollinear collisions contribute to the relaxation of the modes concerned. Estimates³² of the frequency at which appreciable departures from (29) or (30) set in have suggested that, for germanium at liquid-air temperature, these laws should be approximately obeyed, but that departures from them may well be appreciable. As it can be shown that τ_{ph} goes as a higher power of ω when ω gets large, it seems likely that the finiteness of ω will lead to higher, rather than lower, negative exponents. Thus, at present a more plausible explanation of the observed exponents is that the longitudinal and transverse branches contribute comparably to the phonon-drag phenomenon. Since it is known³ that the longitudinal and transverse branches contribute comparable amounts to the scattering of the electrons, this would require that, in the frequency range most important for the scattering of electrons, the relaxation times of the transverse modes be comparable with or rather larger than those of the longitudinal modes. This is contrary to what one would anticipate from (29) and (30), but cannot yet be considered unreasonable, in view of our present ignorance about scattering matrix elements, etc.

There are two further possibilities which would lead to departures from the asymptotic laws (29) and (30). One of these is three-phonon collisions in which the three phonons belong to the same branch and have

almost collinear wave vectors. If such collisions are possible, they will change (29) into the same form as (30). For this to occur for a low-frequency longitudinal mode, it is necessary that some higher-frequency mode with the same direction of propagation have a higher group velocity than a mode of infinitesimal frequency.³² This could happen, for example, if the curve of frequency against wave number in the longitudinal branch curved upward instead of downward. Although this does not seem impossible *a priori*, no real crystal or mathematical model has to our knowledge been found with this property. Moreover, by now our knowledge of the elastic spectrum of germanium³³ is good enough to make it fairly certain that this does not occur for the present material.

More intriguing is the possibility of higher-order processes, i.e., collisions involving four or more phonons,³⁴ or possibly even processes more appropriately describable as relaxation effects.³⁵ From the standpoint of collision theory, a calculation of, say, three-phonon collisions becomes invalid when the frequency of one of the modes involved becomes less than the natural widths τ_{pk}^{-1} of the other modes, or when the failure of energy conservation in an "energetically impossible" process gets less than some relevant $\hbar\tau_{pk}^{-1}$. Moreover, even when neither of these catastrophes occurs, and the three-phonon process that is being considered is perfectly meaningful, it may happen that the number of such possible three-phonon processes is so small that the relaxation of the low-frequency mode involved occurs predominantly by higher-order processes.

However, the rates of four-phonon processes vary much more rapidly with temperature than do those of three-phonon processes.³⁴ At the frequencies of the phonon-drag modes, which are high compared to the reciprocal thermal relaxation time of the majority modes, relaxation can be described as due to some, though not all, of the collisions involving four or more protons. Since the temperature dependence of Q_p in the range 60° to 90°K is not far from that to be expected for three-phonon processes,² we must conclude that higher-order processes are at most a small correction in this range. But if they are at all appreciable at one temperature, they must become completely dominant at a higher temperature. The changeover from dominance of three-phonon processes to dominance of higher-order processes would almost certainly change the anisotropy or the energy dependence of $\Pi_p(\epsilon)$. In particular, if the relaxation time of the phonon-drag modes became independent of frequency — as it would for the simplest type of higher-order processes — Π_p would go as $\epsilon^{\frac{1}{2}}$, in violent disagreement with our observations at all temperatures. The absence of any abrupt change of the energy dependence or anisotropy of Π_p with temperature thus shows that, for the phonon-drag

modes, three-phonon processes are more important than higher-order processes, up to at least 234°K, and that over most of this range the dominance of the former must be extreme.

It is not possible to follow the phonon-phonon scattering to temperatures much below 60°K without a theory that takes account of boundary scattering, an effect which spoils the crystalline symmetry of the problem. In fact, some of the apparent changes which Fig. 11 shows in the behavior of $\Pi_{p\parallel}$ at 60° may well be due to the fact that boundary scattering is beginning to be too large to be corrected for by the crude method of Section V.

9.3 *Transport Theory and Deformation-Potential Theory*

Although the study of the purely electrical properties of n-type germanium was only an incidental feature of the present work, our results have led to a number of new conclusions in this field, which we shall summarize here.

To begin with, we have found considerable evidence (presented in Appendix A and summarized in Section II) that in a magnetic field one need not abandon conventional transport theory in crystal-momentum space in favor of a quantum treatment until the cyclotron spacing $\hbar\omega_c$ has become close to kT . Specifically, we found that, at $\hbar\omega_c/kT \approx \frac{1}{3}$ and for acoustic scattering, the longitudinal and transverse resistivities are correctly given by conventional theory to within 5 per cent or so, and in most cases probably rather closer than this. Even at $\hbar\omega_c/kT \approx 1$ our ΔQ evidence suggests, again for predominantly acoustic scattering, that the errors may be only of the order of 10 per cent or so.

We have found the magnetoresistance and Hall data at high and low magnetic fields to be fairly consistent with the predictions of the electron-group theory, i.e., with the picture which approximates the effect of the scattering processes by relaxation times $\tau_{\parallel}(\epsilon)$ and $\tau_{\perp}(\epsilon)$. For the most accurately measured cases (etched samples at 77.4°K), the accuracy of the fit can be described as follows: The single parameter $w_1 = m_{\perp}^* \langle \epsilon \tau_{\parallel} \rangle / m_{\parallel}^* \langle \epsilon \tau_{\perp} \rangle$ can be chosen to reproduce the longitudinal high-field magnetoresistance limits $\rho(\infty)/\rho(0)$ for the four purest samples in [100], [110] and [111] orientations to within the reproducibility of measurements from sample to sample (2 per cent or so). The values of $\rho(\infty)/\rho(0)$ for three transverse orientations also agree with predictions based on this anisotropy and the assumption $\tau_{\parallel, \perp} \propto \epsilon^{-\frac{1}{2}}$, to within a somewhat larger uncertainty (worst discrepancy 13 per cent), due in part to surface conduction, etc., and in part to quantization. The low-field magneto-

resistance constants obey the symmetry relation characteristic of the electron-group approximation (the relations at the bottom of Table IV with $\Delta\rho$ substituted for ΔQ_p) about as accurately as they obey the requirements of macroscopic symmetry (consistency to a few per cent).

Concerning the energy dependence of τ_{\parallel} and τ_{\perp} , the most accurate evidence we have is that from the ratio of low- to high-field Hall constant, i.e., the ratio of Hall to drift mobility. As was shown in Section VII and Fig. 8, this ratio, when measured on etched samples, gets closer and closer to the value predicted for $\tau_{\parallel, \perp} \propto \epsilon^{-\frac{1}{2}}$ as the importance of impurity scattering is reduced. We were able to conclude that only a very small amount of optical-mode or intervalley scattering, or of any other departure from the law $\tau_{\parallel, \perp} \propto \epsilon^{-\frac{1}{2}}$, can be present at temperatures in the range 77° to 131°K.

In the calculation of the absolute value of the acoustic mobility, μ_a , of n-type germanium from deformation-potential theory,³ two deformation-potential constants are required. One of these can be obtained with some accuracy from piezoresistance measurements, and the other can be inferred if magnetoresistance or other data give a reliable value of the acoustic-scattering anisotropy $\tau_{\parallel}/\tau_{\perp}$. Our present measurements have yielded much better values, both for μ_a and for $w = m_{\perp}^* \tau_{\parallel} / m_{\parallel}^* \tau_{\perp}$, than were available when Ref. 3 was written, so it is of interest to see how well the deformation-potential calculation of μ_a now works out.

The first step is to obtain $\tau_{\parallel}/\tau_{\perp}$ from the cyclotron masses³⁶ $m_{\perp}^* = 0.082m$, $m_{\parallel}^* = 1.58m$ (assumed for the moment to be valid near liquid-air temperature) and the high-purity limit $w = 0.061_0$, presumably characteristic of acoustic scattering. The result is $\tau_{\parallel}/\tau_{\perp} = 1.18$. From this and Equations (52) and (53) or Fig. 5 of Ref. 3 (also based on the masses mentioned), we get the value -0.40 for the ratio Ξ_a/Ξ_u of the deformation-potential constant for stretch along a direction normal to the valley axis to that for the combination of a compression along such a direction with an equal stretch along the valley axis. From this value of Ξ_a/Ξ_u and the piezoresistance constant^{3,37}

$$m_{44}^{(P)} = \frac{1}{3} \frac{\Xi_u}{kT} \frac{w - 1}{w + 2} = 30,400/T \quad (31)$$

we obtain Ξ_a , Ξ_u , and thence the acoustic mobility. The value predicted in this way for 77.4°K is

$$\mu_a (77.4^\circ) = 45,100 \text{ cm}^2/\text{vs.} \quad (32)$$

The empirical value is obtained by dividing the ideal value $\mu_{Ha}/\mu_a = 0.933$ into the observed Hall mobility, after correction for impurity

scattering by the estimates described in Section VII. With the value $38,600 \text{ cm}^2/\text{vs}$ there adopted for μ_{Ha} (see Fig. 8), we get

$$\mu_a(77.4^\circ) = 41,400 \text{ cm}^2/\text{vs}. \quad (33)$$

The agreement between (32) and (33) is gratifying. Although the 8 or 9 per cent discrepancy between them is larger than the uncertainty in either due to uncertainties in the experimental inputs (perhaps 2 or 3 per cent), it is of the right sort to be attributable to departures of the masses or deformation-potential constants at 77° from the low-temperature values we have assumed; moreover, the inherent inaccuracy of the deformation-potential method could conceivably be of this same order. As regards the former point, we may notice that the difference between mobilities varying as $T^{-1.60}$ and as $T^{-1.65}$ amounts to 8 per cent between 46° and 77°K . If the latter variation, which we found to be closer to the truth over the range 60° to 94°K , were due to something like a temperature variation of the masses which disappeared below 40° or so, the discrepancy would be eliminated.

Regarding impurity scattering, finally, we have found in Section VII that it is possible to correlate with each other the effects of ionized-impurity scattering on mobility, Hall-to-drift-mobility ratio, magneto-resistance, etc. — at least in the range where impurity scattering is small to moderate. In this correlation the anisotropy of the impurity scattering was treated as an empirical parameter; beyond this, no specific assumption was necessary about the mechanism of impurity scattering, except that it affected low-energy carriers much more than high-energy ones. The correlations we were able to make in Fig. 8 and Fig. 12 seem to be quantitatively significant, but still leave something to be desired.

All these items deserve more careful study. Our Table XV (Appendix B), incidentally, which gives Maxwellian averages of powers of the energy, may prove handy for various types of calculations.

9.4 Surface Effects

We have encountered serious effects of surface conduction in a number of places in the present study. Our experience justifies the statement that, at all temperatures, one must eliminate or allow for surface effects before one can get really accurate measurements of any bulk electrical property of high-purity material. For some properties, even the qualitative behavior can be completely falsified by surface effects. Specifically, we have found that a surface contribution of less than one per cent to the conductance can greatly affect the apparent high-field transverse magneto-

resistance and the high-field Nernst coefficient near liquid-air temperature and below, and that it can change the sign of the low-field Nernst coefficient in the range 175° to 250°K. Near liquid-air temperature, changes in surface treatment can alter the apparent conductivity and Hall coefficient by amounts ranging up to several per cent.

X. ACKNOWLEDGMENTS

We are indebted to G. W. Hull for his skill in preparing the samples and in mounting them in the various pieces of apparatus used. A large amount of careful observation was required in order to obtain the data presented here. Much of this work was entrusted to J. J. Byrnes, G. W. Hull and P. A. Liberti, and we are indebted to them for their most helpful and sustained efforts. R. G. Treuting and J. N. Hobstetter provided us with a useful program that enabled much of the data to be processed on the I.B.M. 704 computer. Thanks are also due T. T. Wu for discussions of current flow near side-arms.

APPENDIX A. ORBITAL QUANTIZATION

We give here the reasoning behind a number of statements which one can make, including those made in the text. We shall take up first some quantitative statements about the relation of the magnetoresistance of a multivalley semiconductor to that of the simple model assumed in the existing theoretical literature (items 2 and 3 below). Next we shall present empirical and theoretical evidence on the smallness of the effect of quantization on longitudinal magnetoresistance (item 4). Finally, in items 5, 6 and 7, we shall discuss the relation of the effect of quantization on ΔQ to that on magnetoresistance, and shall present evidence that these effects are fairly small even for transverse orientations. As the basis of the whole treatment we shall take the fairly obvious statement:

1. *For the region $\mu H/c \gg 1$, the orbital-quantization correction functions for $\rho(\mathbf{H})/\rho(0)$ and $Q_p(\mathbf{H})/Q_p(0)$ are practically the same functions of H as they would be if the electron-phonon coupling were allowed to approach zero.*

In other words, the difference between the true $\rho(\mathbf{H})/\rho(0)$ or $Q_p(\mathbf{H})/Q_p(0)$ and the value calculated without quantization is a function of the cyclotron spacing $\hbar\omega_c$ (primarily of $\hbar\omega_c/kT$), which depends, of course, on the anisotropy of the electron-phonon interaction, but which depends on the magnitude of the latter (i.e., on $\mu H/c$) in such a way as to approach a limit uniformly as the scattering goes to zero ($\mu H/c \rightarrow \infty$). This is, of course, what one expects intuitively; it follows from any trans-

port theory in which $\rho(0)$, $\rho(H)$, etc. are representable by expansions in powers of the strength of the scattering.

We shall show first that, if the scattering processes are energy-conserving collisions with a squared matrix element that is the same for all transitions within a valley, the following statement holds:

2. *In the approximation of isotropic scattering over an energy shell (i.e., $\tau_{\parallel} = \tau_{\perp}$), the longitudinal magnetoresistance ratio $\rho(\mathbf{H})/\rho(0)$ for any multivalley case equals the product of the value calculated in the absence of quantization by a weighted average of the values calculated for isotropic m^* 's equal to the cyclotron masses of the various valleys for the given direction of \mathbf{H} , provided $\mu H/c \gg$ a value which for germanium is of the order of a fraction of $m_{\parallel}^*/m_{\perp}^*$.*

Thus, for longitudinal magnetoresistance, the orbital-quantization correction can be deduced from the simple-model calculations in the literature.^{5,6,7,8}

For the proof we note first that when $\mathbf{E} \parallel \mathbf{H}$ in, say, the z direction, $\rho_{zz} = 1/\sigma_{zz}$; this is exact for all \mathbf{H} in symmetry directions, and asymptotically exact as $H \rightarrow \infty$ in any direction. Therefore we need consider only σ_{zz} , which is a sum of contributions $\sigma_{zz}^{(i)}$ from the different valleys i . For any such valley we can reduce the transport problem to that for the simple model by means of a transformation which projects the energy surfaces into spheres.³ In considering this transformation we must remember that it is not an orthogonal one, and that it therefore carries with it a need to distinguish between covariant and contravariant vectors. For covariant vectors, \mathbf{F} , the transformed vector, $\bar{\mathbf{F}}$, is related to the original vector \mathbf{F} by

$$\bar{\mathbf{F}}_{\parallel} = r^{\frac{1}{2}}\mathbf{F}_{\parallel}, \quad \bar{\mathbf{F}}_{\perp} = \mathbf{F}_{\perp}, \quad (34)$$

where \parallel and \perp refer to the valley axis and $r = m_{\perp}^*/m_{\parallel}^*$. For contravariant vectors \mathbf{G} the relation is

$$\bar{\mathbf{G}}_{\parallel} = r^{-\frac{1}{2}}\mathbf{G}_{\parallel}, \quad \bar{\mathbf{G}}_{\perp} = \mathbf{G}. \quad (35)$$

The scalar product of a covariant and a contravariant vector is unchanged by the transformation. It is easily verified that, if we regard \mathbf{E} , the crystal momentum \mathbf{P} , and the vector-potential \mathbf{A} as covariant vectors, while regarding velocity (or \mathbf{j}) and \mathbf{H} as contravariant vectors, all the usual relations involved in the Schrödinger equation and the transport equation are the same for the transformed and original vectors, except that in the transformed space the effective-mass tensor becomes isotropic.

Now let $\bar{\sigma}_l(\bar{\mathbf{H}})$ be the conductivity in the transformed system when

$\bar{\mathbf{E}} \parallel \bar{\mathbf{H}}$, and $\bar{\sigma}_i(\bar{\mathbf{H}})$ be that for $\bar{\mathbf{E}} \perp \bar{\mathbf{H}}$. These conductivities will be independent of the orientation of $\bar{\mathbf{H}}$ or of the azimuth of $\bar{\mathbf{E}}$ about $\bar{\mathbf{H}}$, if the scattering for $H = 0$ is isotropic over each energy shell, as is known to be very nearly the case for acoustic scattering in n-type germanium, and if at the same time any effects associated with finiteness of phonon energies are also isotropic. Since it is known that the latter effects are important for incipient quantization effects in longitudinal magnetoresistance,⁸ this last assumption deserves further investigation; however, we shall make it and proceed. Then, for the conductivity along H of all the valleys i in parallel, we have

$$\sigma_{zz}(\mathbf{H}) = \sum_i [\lambda_i^{(i)} \bar{\sigma}_i(\bar{H}^{(i)}) + \lambda_i^{(i)} \bar{\sigma}_i(\bar{H}^{(i)})], \quad (36)$$

where the coefficients $\lambda_i^{(i)}$, $\lambda_i^{(i)}$ depend only on the orientation of \mathbf{H} with respect to the axis of valley i . Therefore, since $\bar{\sigma}_i = \bar{\sigma}_i = \bar{\sigma}_0$ at $\mathbf{H} = 0$,

$$\frac{\sigma_{zz}(\mathbf{H})}{\sigma(0)} = \sum_i \left[\frac{\sigma_i(\bar{H}^{(i)})}{\bar{\sigma}_0} \frac{\lambda_i^{(i)}}{\sum_i (\lambda_i^{(i)} + \lambda_i^{(i)})} + \frac{\bar{\sigma}_i(\bar{H}^{(i)})}{\bar{\sigma}_0} \frac{\lambda_i^{(i)}}{\sum_i (\lambda_i^{(i)} + \lambda_i^{(i)})} \right]. \quad (37)$$

Now, in the absence of quantization, $\bar{\sigma}_i$ is of order \bar{H}^{-2} or less as $\bar{H} \rightarrow \infty$, while $\bar{\sigma}_i$ remains finite; when quantization is just starting to be important and $\mu H/c \gg 1$, $\bar{\sigma}_i$ will still be $\ll \bar{\sigma}_i$. Under our assumptions, the dependence of both these quantities on \bar{H} is given by the theory for the simple model with isotropic effective mass. If we keep only the first term of (37), we have the result italicized above. For some orientations, however, the coefficient involving λ 's in the second term of (37) is larger than that in the first term by a factor of the order of a fraction of $n_{\parallel}^*/m_{\perp}^*$. This explains the qualification in the statement of the theorem.

3. *In the approximation of isotropic scattering over an energy shell (i.e., $\tau_{\parallel} = \tau_{\perp}$), any contribution to the transverse magnetoresistance ratio $\rho(\mathbf{H})/\rho(0)$ which goes as H^2 must for a cubic crystal be independent of the orientations of \mathbf{E} and \mathbf{H} , in the range $\mu H/c \gg 1$, and any contribution which goes as H has only the modest orientation dependence shown in Table XIII.*

We interpret our assumption in the same way as for item 2 above, and with the same degree of risk. The proof is again quite simple if we use the transformations (34) and (35). Since \mathbf{E} is covariant, \mathbf{H} contravariant, $\mathbf{E} \cdot \mathbf{H} = 0$ implies $\bar{\mathbf{E}} \cdot \bar{\mathbf{H}} = 0$ in every valley. When $\mu H/c \gg 1$,

TABLE XIII

Orientation dependence of a hypothetical orbital-quantization term in the transverse magnetoresistance which is proportional to H . The value $r = m^*_\perp/m^*_\parallel = 0.052$ has been assumed.

Direction of		Value of (40) = relative value of δ in $\rho(H)/\rho(0) = a + bH$
\mathbf{H}	\mathbf{j}	
[001]	$\perp \mathbf{H}$	1.03
[011]	[100]	1.74
[011]	[011]	1.28

$E \sim RHj$, and so, with \mathbf{j} along x and \mathbf{H} along z ,

$$\rho_{xx}(H) = j^{-2} \sigma_{yy} E^2 = R^2 H^2 \sigma_{yy} = R^2 H^2 \sum_i \frac{\bar{\sigma}_i(\bar{H}^{(i)}) \bar{E}^{(i)2}}{E^2}, \quad (38)$$

where, as before, $\bar{\sigma}_i(\bar{H})$ is the conductivity of the simple-model sphericalized valley transverse to a magnetic field $\bar{\mathbf{H}}$, and the summation is over the different valleys. In the absence of orbital quantization, $\bar{\sigma}_i$ would go as \bar{H}^{-2} , and so (38) would be independent of the magnitude of H . A term in ρ_{xx} going as H^2 can arise only from a term in $\bar{\sigma}_i$ independent of \bar{H} . If $\delta\bar{\sigma}_i$ is such a term, the corresponding increment in ρ_{xx} is

$$\delta\rho_{xx}(\mathbf{H}) = R^2 H^2 \delta\bar{\sigma}_i \sum_i \frac{\bar{E}^{(i)2}}{E^2}. \quad (39)$$

But the summation is a quadratic function of the components of \mathbf{E} with the symmetry of the crystal, hence it is isotropic for any cubic case. This proves the first statement. A term in ρ_{xx} going as H must come from a term in $\bar{\sigma}_i$ going as \bar{H}^{-1} , so the relative magnitudes of this term in various orientations must be proportional to the values of the quantity

$$\sum_i \left(\frac{H}{\bar{H}^{(i)}} \right) \left(\frac{\bar{E}^{(i)}}{E} \right)^2. \quad (40)$$

These values are shown in Table XIII.

4. *Theory and experiment agree that the effect of orbital quantization on longitudinal magnetoresistance in pure germanium is very small as long as $\hbar\omega_c/kT \lesssim 1$ or so.*

Calculations⁸ based on the simple-model theory of Argyres⁷ give $\rho(H)$ actually about 15 per cent below $\rho(0)$ for $\hbar\omega_c/kT = 1$, rising to $\rho(0)$ for $\hbar\omega_c/kT \approx 2$ and increasing linearly at high fields. The initial depression, however, is due to the fact that, in this theory, the scattering is

describable by a relaxation time $\tau(\epsilon)$ which goes to zero when $\epsilon = (n + \frac{1}{2})\hbar\omega_c$. Since τ^{-1} is proportional to the density of states in energy, the mean value of τ^{-1} in an energy range several times $\hbar\omega_c$ in width is close to the value calculated without quantization. Because of the fluctuations, the mean value of τ over such a range is greater than the reciprocal of the mean of τ^{-1} ; hence the conductivity is greater than in the absence of quantization, the current being carried preferentially in the energy ranges with the highest τ 's. Now these fluctuations in τ will be considerably smoothed out if the energies of the phonons responsible for the scattering are even a fraction of $\hbar\omega_c$, and crude numerical calculations⁸ show that a reasonable amount of such smoothing will change the maximum reduction of ρ from 15 per cent to a few per cent. Thus, we expect theoretically that the effect of quantization will be small for at least the range $\hbar\omega_c/kT \lesssim 1$.

The highest values of $\hbar\omega_c/kT$ which have so far been obtained experimentally for material with essentially acoustic scattering have been those of Furth and Waniek¹² for germanium in pulsed fields at room temperature. Their maximum field of 460,000 gauss corresponds to $\hbar\omega_c/kT = 2.1$. Up to this field, their $\rho(H)$ shows nice saturation, agreeing fairly well with the predictions of the theory ignoring quantization. By contrast, a combination of the latter curve with the quantization correction computed for zero phonon energies gives a curve which is rising conspicuously at the highest fields, and this is in marked disagreement with the observations. Moreover, the saturation values of longitudinal magnetoresistance in both the [100] and the [111] directions agree with the predictions of the unquantized theory with $w = m_{\perp}^* \tau_{\parallel} / m_{\parallel}^* \tau_{\perp}$ equal to about 0.057, a fairly reasonable value.

5. *In the region $\mu H/c \gg 1$, $\hbar\omega_c/kT < 1$ and, for acoustic scattering, the fractional changes in $Q_p(\mathbf{H})$ and in $\rho(\mathbf{H})$ attributable to orbital quantization either should both be $\ll 1$, or else should be roughly equal.*

Since the contribution of a phonon mode \mathbf{q} to the Peltier flux is proportional to $[c(\mathbf{q})]^2 \tau_{ph}(\mathbf{q})$, where $c(\mathbf{q})$ is its velocity and $\tau_{ph}(\mathbf{q})$ is its relaxation time, the statement just made will be valid if we can show that the crystal momentum lost by a current to the phonons is distributed in nearly the same manner among the different modes, regardless of whether quantization is assumed or ignored.

It is easy to see qualitatively that, as $\hbar\omega_c/kT$ becomes $\gg 1$, the scattering of the electrons must involve phonons whose wave-vector components $\perp \mathbf{H}$ become very large. This will undoubtedly cause the mean relaxation time of the phonons to become much shorter. Also, the change in the average direction of the wave vectors of the contributing phonons

will alter their mean relaxation time if the individual $\tau_{ph}(\mathbf{q})$ are anisotropic. Finally, if quantization changes the relative contributions of transverse and longitudinal branches, the average $[c(\mathbf{q})]^2 \tau_{ph}(\mathbf{q})$ will be changed. We shall try to show, however, that all these effects are rather small when $\hbar\omega_c/kT < 1$.

Since the theory of quantum magnetoresistance is rather complicated when $\hbar\omega_c/kT < 1$, we shall base our estimates on calculations of the distribution of crystal momentum among the different phonon modes for the case $\hbar\omega_c/kT \gg 1$. These calculations give the asymptote of the curve of $(c^2 \tau_{ph})_{Av}$ against $\hbar\omega_c/kT$, an asymptote which should be very nearly valid for $\hbar\omega_c/kT \approx 3$. Since the curve almost certainly is monotonic and comes in to its low-field limit with a horizontal tangent, one can estimate the order of magnitude of its deviation from the low-field value for $\hbar\omega_c/kT < 1$, or at least a rough upper bound to this. It will suffice for our purposes to give the argument for a semiconductor with isotropic effective mass, since we have shown above that, when the scattering is nearly isotropic over an energy shell, the transport problem for an anisotropic valley can be reduced to that for the isotropic case.

For \mathbf{E} parallel to \mathbf{H} the appropriate quantum transport theory is well understood. A brief inspection of the formulas given by Argyres⁷ shows that, for $\hbar\omega_c/kT \gg 1$, the crystal momentum in the direction of the current which is delivered to the phonon modes of various wave numbers \mathbf{q} is proportional to

$$q_z^2 \exp\left(\frac{-\hbar^2 q_z^2}{8m^*kT}\right) \exp\left[\frac{-\hbar(q_x^2 + q_y^2)}{4m^*\omega_c}\right], \tag{41}$$

where we have taken the z direction to be that of \mathbf{E} and \mathbf{H} . From this, one can derive any kind of average which may be desired. Although our main interest is in an average τ_{ph} , hence in an average of something like q^{-2} or q^{-1} with the weight (41), it will suffice to compute the more easily evaluated averages of q_x^2 , q_y^2 , q_z^2 , noting that the average of q^{-2} will differ by a somewhat smaller percentage from its zero-field value than will that of q^2 . It is clear from (41) that the averages of q_x^2 and q_y^2 will go proportionally to ω_c , while that of q_z^2 will be independent of ω_c . The most convenient quantities to tabulate are the ratios of these averages to the corresponding averages for $H = 0$; such ratios are given for $\hbar\omega_c/kT = 3$ in the first row of Table XIV. They differ by so little from unity that it seems safe to assume that, for $\hbar\omega_c/kT < 1$, the average $c^2 \tau_{ph}$ is very close indeed to the value at $H = 0$.

It is worth remarking that the decrease of longitudinal resistance below the zero-field value, which the theory predicts^{7,8} in the region

TABLE XIV

Ratios of the average q_x^2 , q_y^2 , q_z^2 of the phonon modes, weighted in proportion to the crystal momentum they receive in the current direction, to the corresponding averages in the absence of a magnetic field. The ratios were evaluated from the high-field asymptotes for the value $\hbar\omega_c/kT = 3$. For the transverse case the predictions of two rival theories are given. Acoustic scattering has been assumed.

Direction		Theory	$\frac{(q_x^2)_{Av. nt 3}}{(q_x^2)_{Av. nt 0}}$	$\frac{(q_y^2)_{Av. nt 3}}{(q_y^2)_{Av. nt 0}}$	$\frac{(q_z^2)_{Av. nt 3}}{(q_z^2)_{Av. nt 0}}$
E	H		(linear in H)	(linear in H)	(constant)
z	z	Argyres ⁷	1.12	1.12	1.29
x	z	Argyres ¹⁰	1.12	1.12	$\rightarrow 0$
x	z	Klinger-Voronyuk ⁹	1.77	3.37	$\rightarrow 0$

$\hbar\omega_c/kT \lesssim 2$, probably has no very noticeable counterpart in the behavior of the $(q_x^2)_{AV}$ etc. This decrease arises from the current being carried preferentially in the regions of energy just below the values $(n + \frac{1}{2})\hbar\omega_c$, where the relaxation time is longer than average. But the changes in the average q_x^2 etc., for transitions at energy ϵ as ϵ passes through $(n + \frac{1}{2})\hbar\omega_c$ are much less pronounced than those in the relaxation time.

Similar calculations can be made from the theory of transverse quantum magnetoresistance, except that the correct form of the theory is still in dispute.^{9,10} The alternative versions agree in predicting that for $\hbar\omega_c/kT \gg 1$ the current should be carried predominantly in states of low k_z , hence should involve transitions of low q_z . If the levels had zero natural width and the phonons had zero energy, the integral for the current in the **E**-direction would diverge. In the last two rows of Table XIV, where we have tabulated the predictions of the two types of theory, we have therefore used the entry " $\rightarrow 0$ " in the last column. The distribution of crystal momentum in the current direction among the various modes turns out to be proportional, as far as q_x and q_y are concerned, to the last factor of (41) in the Argyres theory,¹⁰ and to q_y^2 times this factor in the theory of Klinger and Voronyuk.⁹ While the latter theory spreads the q_y -distribution rapidly as H increases, the over-all average $c^2\tau_{ph}$ reflects the combined behavior of all three components of \mathbf{q} , and it seems likely that, for $\hbar\omega_c/kT < 1$, the average $c^2\tau_{ph}$ will not deviate from its zero-field value by more than a small fraction of the latter.

By contrast, the values which these theories predict for the resistance $\rho(\mathbf{H})$ are at least several times $\rho(0)$ for $\hbar\omega_c/kT = 3$, for a substance like germanium.

6. *The high-field transverse magnetoresistance is so sensitive to inhomogeneities, surface conduction and other perturbing influences that it cannot serve as a reliable indicator of the size of the orbital-quantization effect.*

At least, before it can be so used, many auxiliary investigations will have to be made to ensure the elimination of such perturbing influences. For the present purpose it will suffice to consider the case where the normal bulk conduction is in parallel with conduction by some other mechanism (e.g., surface conduction) which has a widely different Hall effect, say zero. Slight inhomogeneities in the bulk conduction itself can be shown to have a similar effect,⁴ though they are probably not of major importance for our specimens.

Let the x -axis be along the length of the specimen, and let the magnetic field \mathbf{H} be in the z -direction. Let ρ_b and ρ_s be the effective resistivities due respectively to normal bulk conduction alone and to the perturbing conduction alone, and let R_b , $R_s = 0$ be the corresponding Hall constants. Then we have, assuming for simplicity that ρ_b and ρ_s are isotropic,

$$j_{sy} = -j_{by}, \quad (42)$$

$$E_y = \rho_b j_{by} - R_b H j_{bx} = -\rho_s j_{by}, \quad (43)$$

$$E_x = \rho_b j_{bx} + R_b H j_{by} = \rho_s j_{sx}. \quad (44)$$

From (43) and (44) we get the effective resistivity

$$\begin{aligned} \rho_{\text{eff}} &= \frac{E_x}{(j_{bx} + j_{sx})} = \frac{\rho_s [\rho_b (\rho_b + \rho_s) + (R_b H)^2]}{(\rho_b + \rho_s)^2 + (R_b H)^2} \\ &\approx \rho_b \left[1 + \frac{(R_b H)^2}{\rho_b \rho_s} \right] \end{aligned} \quad (45)$$

if $\rho_s \gg R_b H \gg \rho_b$. The physical meaning of (45) is that the Hall field due to the b -conduction is partially shorted by the s -conduction, with consequent increase of the total dissipation. If ρ_s is, say, 200 ρ_b but $R_b H / \rho_b$ is 7, the second term in the bracket in (45) will be 0.25, and the effect on a parabolic extrapolation to infinite field, such as we have made in Fig. 6, will be between two and three times this.

The etching experiments described in Section V have shown that surface treatment can alter the observed transverse magnetoresistance by a very sizable amount at liquid-air temperature, especially for $\mathbf{H} \parallel [001]$, $\mathbf{j} \parallel [100]$. The subject merits a much more thorough investigation. However, it should be emphasized that the high-field ΔQ effect is not com-

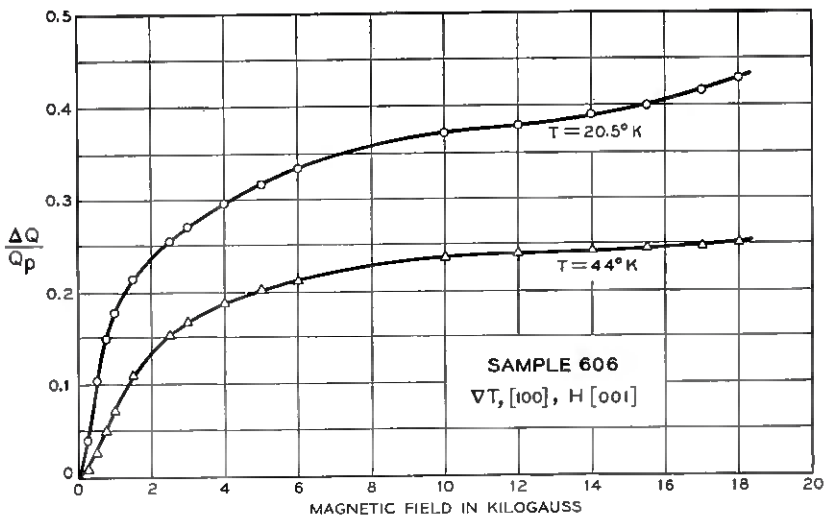


Fig. 17 — Values of the ratio of the transverse $\Delta Q(H)$ to $Q_p(0)$ for Sample 606 at 20.5° and 44°K, showing the presumed effect of orbital quantization. The values of $|Q_p(0)|$ are 8600 $\mu\text{V}/\text{deg}$ at 20.5° and 3730 at 44°.

parably sensitive, either experimentally or theoretically; this is, of course, because there is no large Hall field in the ΔQ experiment.

7. At 20° and 44°K the effect of orbital quantization on the transverse Q_p can be identified in the data, and it is small. Since the mechanisms of scattering of phonons are quite different at these two temperatures, these small values can hardly be due to an accidental cancellation of effects of quantization on $\rho(\mathbf{H})$ and on $c^2\tau_{ph}$, and it is probably safe to conclude that at liquid-air temperature quantization modifies ρ and Q_p by only a few per cent.

Fig. 17 shows the ΔQ data for Sample 606 at these two temperatures. Both curves start to saturate as H increases, then turn upward again. These upturns are undoubtedly due to orbital quantization; one may also suspect the presence of a quantization term varying less rapidly with H than the conspicuous part of the upturn, say linearly. One would like to separate ΔQ into a nonquantization part that saturates when $\mu H/c \gg 1$ and a quantization correction dependent on $\hbar\omega_c/kT$. If, as the data suggest, a curve of the quantization correction alone against H would be concave upward, then we must surely underestimate the infinite-field limit without quantization if we assume it to equal the intercept on the $H = 0$ axis of the tangent to the $Q(H)$ curve at its most nearly horizontal portion. Thus, the difference between this intercept and the or-

dinate of the curve at 18,000 gauss ought to be an overestimate of the quantization correction. In this way we estimate that at 20.5°K, the orbital quantization correction is $\lesssim 10$ per cent of $Q_p(0)$, and that at 44° it is $\lesssim 3$ per cent. The reasonableness of our partition of ΔQ into its two parts is attested by the values obtained for what $\Delta Q/Q_p$ would be in the absence of quantization. Our intercept at 44° corresponds (if $|\Delta Q_c| \approx 43 \mu\text{v}/\text{deg}$) to a $\Delta Q_p/Q_p$ of 0.21, only slightly larger than the values found in the liquid-air range. The intercept at 20.5° corresponds to a $\Delta Q_p/Q_p$ of about 0.33; this exceeds the high-temperature value by about the right amount to be attributed to the alteration of the phonon-phonon scattering law (see next paragraph).

Now, at 20.5°K, the scattering of the phonons involved in Q_p is almost entirely due to the boundaries of the specimen, while at 44° the phonon-phonon scattering exceeds the boundary scattering. Since these two types of scattering have quite different dependences on frequency and direction, it would not be reasonable to postulate that the effect of quantization on the average relaxation time of the participating phonons at both temperatures happens to be such as almost to cancel the effect of quantization on the magnetoresistance. The scattering of the electrons, incidentally, can be taken as approximately acoustic at both temperatures, since at 20.5° the observed Hall mobility, 215,000 cm^2/vs , is about 0.8 of the estimated ideal value. We therefore conclude that, with acoustic scattering, the effect of quantization on the transverse $\rho(H)/\rho(0)$ is probably no more than 10 per cent or possibly 20 per cent at 20.5°K. This conclusion is in violent contrast to the result of direct magnetoresistance measurements at this temperature, which for this sample (unetched!) gave $\rho(\mathbf{H})/\rho(0) = 4.8$ at 18,000 gauss. We believe the discrepancy is to be attributed to the great sensitivity of high-field transverse magnetoresistance to perturbing influences, especially surface conduction, as discussed in the Section V of the text.

From the figures given above and the plausible assumption that above 44°K the orbital-quantization term decreases at least as fast as T^{-1} , we conclude that, at liquid-air temperature and 18,000 gauss, this term is probably no more than a few per cent of the resistivity $\rho(H)$, and almost certainly is no more than this for the thermoelectric power $Q(H)$ or for $Q_p(H)$. It should be remembered, of course, that the effect of the quantization term on the extrapolation of ρ or Q to infinite field, made on a plot such as those of Fig. 6, will be proportionally larger than the effect at the field at which the extrapolation commences; for our parabolic extrapolation procedure, the factor is a little under 2 for a term proportional to H and a little over 2 for a term proportional to H^2 .

APPENDIX B. DERIVATION OF FORMULAS

In this appendix we shall give some of the intermediate steps involved in applying the procedure outlined in Section III for the calculation of B and ΔQ at low and high magnetic fields. Our objective in doing so is partly to clarify the procedure used, but primarily to list a number of intermediate formulas which are especially useful or illuminating, and to give the final formulas in a form sufficiently general to be applied to other cubic many-valley structures (e.g., silicon), and to cases where the anisotropies are energy-dependent.

Our starting point is the electron-group picture, with each ellipsoidal energy shell constituting a group g . The quantities we wish to calculate are the components of the Peltier tensor $\mathbf{\Pi}$, or, rather, the terms of order 1, H , H^2 in $\mathbf{\Pi}$ at low H , and those of order 1, H^{-1} at high H . For both cases we shall base our calculations on (6) or (9), which relate $\mathbf{\Pi}$ to the partial conductivity tensors δ_θ of the groups. These δ_θ , functions of \mathbf{H} , are derivable from the transport equation (10) for group g , namely,

$$\boldsymbol{\tau}^{-1} \cdot \mathbf{m}^* \cdot d\mathbf{j} \pm \left(\frac{e}{c}\right) d\mathbf{j} \times \mathbf{H} = e^2 \left(\frac{\epsilon}{\langle \epsilon \rangle}\right) \mathbf{E} dn, \quad (46)$$

where the upper and lower signs, it will be remembered, are for electrons and holes, respectively.

B.1 *Small H, Basic Formulas*

The conductivity tensor $d\delta (\equiv \delta_\theta)$ of an energy shell is the coefficient of \mathbf{E} in the solution of (46) for $d\mathbf{j}$. For $H \rightarrow 0$ the usual way to expand $d\delta$ in powers of H is by iteration. In a coordinate system oriented along the principal axes of the valley we find for the parts containing H^0 , H^1 , H^2 respectively:

$$d\sigma_{\alpha\alpha}^{(0)} = e^2 \left(\frac{\epsilon}{\langle \epsilon \rangle}\right) \left(\frac{\tau_\alpha}{m_\alpha^*}\right) dn, \quad d\sigma_{\alpha\beta}^{(0)} = 0 \quad \text{if } \beta \neq \alpha; \quad (47)$$

$$d\sigma_{\alpha\beta}^{(1)} = \mp \left(\frac{e^3}{c}\right) \left(\frac{\epsilon}{\langle \epsilon \rangle}\right) \left(\frac{\tau_\alpha \tau_\beta}{m_\alpha^* m_\beta^*}\right) \sum_\gamma \delta_{\alpha\beta\gamma} H_\gamma dn, \quad (48)$$

where $\delta_{\alpha\beta\gamma}$ is ± 1 when α, β, γ are even (odd) permutations of 123, and is 0 otherwise; and

$$d\sigma_{\alpha\beta}^{(2)} = - \left(\frac{e^4}{c^2}\right) \left(\frac{\epsilon}{\langle \epsilon \rangle}\right) \sum_{\lambda\gamma\delta} \left(\frac{\tau_\alpha \tau_\beta \tau_\lambda}{m_\alpha^* m_\beta^* m_\lambda^*}\right) \delta_{\alpha\lambda\gamma} \delta_{\beta\lambda\delta} H_\lambda H_\delta. \quad (49)$$

B.2 Value of Q at $H = 0$

Let us start by computing Q for $H = 0$ in terms of the principal components $\Pi_{\parallel}(\epsilon)$, $\Pi_{\perp}(\epsilon)$ of $\mathbf{\Pi}_g$. From (6), (1) and (47) we have, for a cubic crystal with axially symmetrical valleys,

$$\begin{aligned}
 TQ(0) &= \frac{1}{3} \sum_{\alpha} \Pi_{\alpha\alpha} = \frac{\frac{1}{3} \sum_{\alpha} \int \Pi_{g\alpha\alpha} d\sigma_{\alpha\alpha}^{(0)}}{\frac{1}{3} \sum_{\alpha} \int d\sigma_{\alpha\alpha}^{(0)}} \\
 &= \frac{\langle \epsilon \Pi_{\parallel} \tau_{\parallel} / m_{\parallel}^* \rangle + 2 \langle \epsilon \Pi_{\perp} \tau_{\perp} / m_{\perp}^* \rangle}{\langle \epsilon \tau_{\parallel} / m_{\parallel}^* \rangle + 2 \langle \epsilon \tau_{\perp} / m_{\perp}^* \rangle},
 \end{aligned}
 \tag{50}$$

where, as always, the angular brackets denote Maxwellian averages. This equation, of course, also holds for Q_e and Q_p separately, with $\Pi_{\parallel, \perp} \rightarrow \Pi_e$ or $\Pi_{p\parallel, \perp}$. When the anisotropies $p = \Pi_{p\parallel} / \Pi_{p\perp}$ and $w = m_{\perp}^* \tau_{\parallel} / m_{\parallel}^* \tau_{\perp}$ are energy-independent, the expression for Q_p reduces to

$$TQ_p(0) = \frac{\langle \epsilon \Pi_{\perp} \tau_{\perp} \rangle}{\langle \epsilon \tau_{\perp} \rangle} \frac{(2 + pw)}{(2 + w)}.
 \tag{51}$$

When $\mathbf{\Pi}_p$ and $\boldsymbol{\tau}$ depend on energy according to simple power laws, the Maxwellian averages occurring in these and subsequent expressions can be evaluated from Table XV at the end of this appendix.

B.3 The Low-Field Value of B

From (1), (5) and (6) we have, if \mathbf{H} is in the z direction

$$TBH = \Pi_{xy}^{(1)} = \left[\sum_g \Pi_g \cdot \sigma_g^{(1)} \right]_{xy} \rho^{(0)} + \frac{\Pi^{(0)} \rho_{xy}^{(1)}}{\rho^{(0)}}.
 \tag{52}$$

If the crystal has over-all cubic symmetry, the evaluation of the first term on the right of (52) can be greatly simplified by noting that it is isotropic and may therefore be replaced by its average over all permutations of the directions of x , y and z (keeping the system right-handed, of course). If this averaging is done before the summing on valleys, the result will be the same for all valleys. For a single valley it can be evaluated by taking the coordinate axes along the principal axes of the valley. We find in this way, using (48),

$$\left[\sum_g \mathbf{\Pi}_g \cdot \boldsymbol{\sigma}_g^{(1)} \right]_{xy} = \mp \frac{ne^3 \rho^{(0)} H}{3c\langle \epsilon \rangle} \left[\frac{\langle \epsilon \tau_{\perp}^2 \Pi_{\perp} \rangle}{m_{\perp}^{*2}} + \frac{\langle \epsilon \tau_{\parallel} \tau_{\perp} \Pi_{\perp} \rangle}{m_{\parallel}^* m_{\perp}^*} + \frac{\langle \epsilon \tau_{\parallel} \tau_{\perp} \Pi_{\parallel} \rangle}{m_{\parallel}^* m_{\perp}^*} \right],
 \tag{53}$$

where n is the total carrier density and, as usual, the angular brackets

represent Maxwellian averages. In the last term of (52) we have, from (11),

$$\frac{\rho_{zu}^{(1)}}{\rho^{(0)}} = \pm \frac{\mu_H H}{c}, \quad (54)$$

where μ_H is the low-field Hall mobility. We may use these and the expressions of Ref. 3 — also obtainable from (47) and (48) —

$$\sigma^{(0)} = \rho^{(0)-1} = \frac{ne^2}{3\langle\epsilon\rangle} \left[\frac{2\langle\epsilon\tau_{\perp}\rangle}{m_{\perp}^*} + \frac{\langle\epsilon\tau_{\parallel}\rangle}{m_{\parallel}^*} \right], \quad (55)$$

$$\mu_H = e \left[\frac{\langle\epsilon\tau_{\perp}^2\rangle}{m_{\perp}^{*2}} + \frac{2\langle\epsilon\tau_{\parallel}\tau_{\perp}\rangle}{m_{\parallel}^*m_{\perp}^*} \right] \left[\frac{2\langle\epsilon\tau_{\perp}\rangle}{m_{\perp}^*} + \frac{\langle\epsilon\tau_{\parallel}\rangle}{m_{\parallel}^*} \right]^{-1}, \quad (56)$$

to obtain finally

$$TB = \left\{ \mp \left[\frac{\langle\epsilon\tau_{\perp}^2\Pi_{\perp}\rangle}{m_{\perp}^{*2}} + \frac{\langle\epsilon\tau_{\parallel}\tau_{\perp}\Pi_{\perp}\rangle}{m_{\parallel}^*m_{\perp}^*} + \frac{\langle\epsilon\tau_{\parallel}\tau_{\perp}\Pi_{\parallel}\rangle}{m_{\parallel}^*m_{\perp}^*} \right] \cdot \left[\frac{\langle\epsilon\tau_{\perp}^2\rangle}{m_{\perp}^{*2}} + \frac{2\langle\epsilon\tau_{\parallel}\tau_{\perp}\rangle}{m_{\parallel}^*m_{\perp}^*} \right]^{-1} \pm TQ \right\} \frac{\mu_H}{c}. \quad (57)$$

This again applies to B_p and B_e separately, as well as to their sum, if $\Pi_{\parallel,\perp}$ and Q are given the appropriate subscript. When the anisotropies are independent of energy the formulas of Table I are easily derived from this with the use of (51). Note that the sign of B is independent of the sign of the carriers, since for electrons the upper sign is to be used in (57) with negative Π 's and Q , while the reverse obtains for holes.

All the formulas we have written thus far apply to any cubic material with axially symmetrical valleys.

B.4 General Formulas for ΔQ at Small H

For the low-field ΔQ we must take the second-order part of (9) and use (48) and (49). Alternatively, we can eliminate the $d\sigma^{(1)}$'s in favor of B . For the second-order part of the thermoelectric power in the α direction we find, for extrinsic material,

$$Q_{\alpha\alpha}^{(2)} = \left(\frac{\rho^{(0)}}{T} \right) \sum_{\theta} [\Pi_{\theta} \cdot \sigma_{\theta}^{(2)}]_{\alpha\alpha} + Q^{(0)} \left[\frac{\rho_{\alpha\alpha}^{(2)}}{\rho^{(0)}} \right] + \left[Q^{(0)} \left(\frac{\mu_H}{c} \right) \mp B \left(\frac{\mu_H}{c} \right) \right] (H^2 - H_a^2), \quad (58)$$

with the usual convention of upper sign for n type, lower sign for p . The quantity $\rho_{\alpha\alpha}^{(2)}/\rho^{(0)}$ is just the magnetoresistance. Our task is now to compute the first term of (58) for various orientations.

Since in cubic material three phenomenological constants suffice to describe $Q^{(2)}$ for all possible orientations of \mathbf{H} and ∇T , only three evaluations of (58) need be made. We can show further than when the valleys have [100] or [111] orientations in the Brillouin zone, their axial symmetry alone gives a further relation between the three constants, at least to the accuracy of the electron-group approximation, so that only two orientations need be computed. The argument, which is a generalization of the corresponding one for magnetoresistance, was given in Appendix C of our first paper.¹ In terms of the q_b, q_c, q_d defined under Table IV the results are

$$q_c = -q_b \quad \text{for [111] valleys,} \tag{59}$$

$$q_b + q_c + q_d = 0 \quad \text{for [100] valleys.} \tag{60}$$

The evaluation of (58) can take a variety of forms, depending on whether $\rho^{(2)}$ and B are expressed explicitly in terms of the τ 's and Π 's, or are left as themselves. These forms can be converted into one another by use of (57) etc., and of

$$\frac{\rho_{\alpha\alpha}^{(2)}}{\rho^{(0)}} = - \frac{\sum_g \sigma_{g\alpha\alpha}^{(2)}}{\sigma^{(0)}} \quad (\alpha \parallel \mathbf{H}) \tag{61}$$

or

$$= - \left[\left(\frac{\mu_H}{c} \right)^2 + \frac{\sum_g \sigma_{g\alpha\alpha}^{(2)}}{\sigma^{(0)}} \right] \quad (\alpha \perp \mathbf{H}).$$

The expressions for $\sum_g \sigma_{g\alpha\alpha}^{(2)}$ can be taken from the magnetoresistance literature, or (for [111] valleys) from (65) below or its longitudinal modification, by setting $\mathbf{\Pi} = \mathbf{1}$.

B.5 ΔQ_e at Small H with Energy-Independent Anisotropies

As with the Nernst coefficient, the evaluation of $Q_e^{(2)}$ for the case of an energy-independent anisotropy of τ is a little easier than that of $Q_p^{(2)}$. When \mathbf{H} is parallel to the direction of measurement (the α direction), the last term of (58) drops out and, with $\mathbf{\Pi}_g$ replaced by $\mp \epsilon/e$, the ratio of the first term of (58) to the second is a function only of the energy dependence of τ . Using (61) we find easily, with $\rho^{(2)} = \Delta\rho$ and with the upper sign for electrons, the lower for holes,

$$\Delta Q_e = Q_{e\alpha\alpha}^{(2)} = \mp \frac{3}{2} \frac{k}{e} \left[\frac{\langle \epsilon^2 \tau \rangle}{\langle \epsilon \tau \rangle \langle \epsilon \rangle} - \frac{\langle \epsilon^2 \tau^2 \rangle}{\langle \epsilon \tau^2 \rangle \langle \epsilon \rangle} \right] \frac{\Delta\rho}{\rho^{(0)}} \quad (\alpha \text{ longitudinal}). \tag{62}$$

For transverse or general orientations of \mathbf{H} we must use (57) to evaluate

the last term of (58) and must modify the relation of $\mathbf{d}_\sigma^{(2)}$ to $\mathbf{g}^{(2)}$; a little algebra gives

$$\Delta Q_\sigma = Q_{\sigma\sigma}^{(2)} = \mp \frac{3k}{2e} \left[\left(\frac{\langle \epsilon^2 \tau \rangle}{\langle \epsilon \tau \rangle \langle \epsilon \rangle} - \frac{\langle \epsilon^2 \tau^3 \rangle}{\langle \epsilon \tau^3 \rangle \langle \epsilon \rangle} \right) \frac{\Delta \rho}{\rho(0)} + \left(\frac{\langle \epsilon^2 \tau^2 \rangle}{\langle \epsilon \tau^2 \rangle \langle \epsilon \rangle} - \frac{\langle \epsilon^2 \tau^3 \rangle}{\langle \epsilon \tau^3 \rangle \langle \epsilon \rangle} \right) \left(\frac{\mu_H}{c} \right)^2 (H^2 - H_\alpha^2) \right]. \quad (63)$$

These formulas are valid for any multivalley band structure in a cubic crystal. Still more general expressions for ΔQ_σ , valid for [111] valleys if the anisotropy of τ is not energy-independent, could be obtained from (65) below.

B.6 ΔQ and ΔQ_p at Small H

Because of the anisotropy of $\mathbf{\Pi}_p$, the first term of (58) is a little more difficult to evaluate for the phonon-drag contribution or for the total ΔQ . We may again eliminate the sum over different valleys by noting that ΔQ is unchanged if we apply any operation D of the cubic symmetry group both to \mathbf{H} and to ∇T , the α -direction in (58). If we average the contribution of each valley over all such D , the result will be the same for all valleys. If $\mathbf{u}^{(T)}$ is the unit vector in the direction of ∇T , the result can be written, in the principal-axis system of some one valley,

$$\mathbf{u}^{(T)} \cdot \sum_{\nu} [\mathbf{\Pi}_\sigma \cdot \mathbf{d}_\sigma^{(2)}] \cdot \mathbf{u}^{(T)} = -\frac{ne^4}{c^2} \sum_{\alpha\beta\lambda\gamma\delta} \frac{\langle \epsilon \mathbf{\Pi}_\alpha \tau_\alpha \tau_\beta \tau_\lambda \rangle}{\langle \epsilon \rangle m_\alpha^* m_\beta^* m_\lambda^*} \delta_{\lambda\beta\gamma} \delta_{\lambda\alpha\delta} \cdot \langle H_\gamma H_\delta u_\alpha^{(T)} u_\beta^{(T)} \rangle_D, \quad (64)$$

where $\langle \rangle_D$ means the average of the components shown over all sets, $D\mathbf{H}$, $D\mathbf{u}^{(T)}$.

So far, the equations are valid for any cubic multivalley material. We shall now specialize to valleys along the [111] axes of the Brillouin zone. For this case and for \mathbf{H} along a [100] axis, $\nabla T \perp \mathbf{H}$, (64) can be evaluated to

$$\mathbf{u}^{(T)} \cdot \sum_{\sigma} [\mathbf{\Pi}_\sigma \cdot \mathbf{d}_\sigma^{(2)}] \cdot \mathbf{u}^{(T)} = -\frac{ne^4 H^2}{c^2} \left[\frac{2}{9} \frac{\langle \epsilon \mathbf{\Pi}_{\parallel} \tau_{\parallel}^2 \tau_{\perp} \rangle}{\langle \epsilon \rangle m_{\parallel}^* m_{\perp}^*} + \frac{1}{9} \frac{\langle \epsilon \mathbf{\Pi}_{\parallel} \tau_{\parallel} \tau_{\perp}^2 \rangle}{\langle \epsilon \rangle m_{\parallel}^* m_{\perp}^*} + \frac{4}{9} \frac{\langle \epsilon \mathbf{\Pi}_{\perp} \tau_{\parallel} \tau_{\perp}^2 \rangle}{\langle \epsilon \rangle m_{\parallel}^* m_{\perp}^*} + \frac{2}{9} \frac{\langle \epsilon \mathbf{\Pi}_{\perp} \tau_{\perp}^3 \rangle}{\langle \epsilon \rangle m_{\perp}^*} \right]. \quad (65)$$

When \mathbf{H} and ∇T are both along a [100] axis we find an expression of similar form to (65) but with the coefficients $\frac{2}{9}$, $\frac{1}{9}$, $\frac{4}{9}$, $\frac{2}{9}$ replaced respec-

tively by $\frac{2}{3}$, $-\frac{2}{3}$, $-\frac{2}{3}$, $\frac{2}{3}$. These expressions apply, of course, to either the electron or the phonon-drag contribution, or to their sum; one need only affix the proper subscript to all Π 's.

For the most general case the expression for ΔQ in terms of moments of $\Pi_{\parallel, \perp}$ is obtained by inserting (65), or its longitudinal modification, into (58); one may or may not wish to express the $Q^{(0)}$ and B occurring there by (50) and (57). The magnetoresistance can be expressed if desired by (61) and (65) with $\mathbf{\Pi} = \mathbf{1}$, and $\rho^{(0)}$ and μ_H by (55) and (56). The last entries in Table IX were obtained in this way.

If we assume that $p = \Pi_{p\parallel}/\Pi_{p\perp}$ and $w = \tau_{\parallel}m_{\perp}^*/\tau_{\perp}m_{\parallel}^*$ are energy-independent, we get, after a little reduction,

$$\Delta Q_p \Big|_{100}^{001} = Q_p \left(\frac{\mu_H H}{c} \right)^2 \left[\zeta_p + \frac{(2+w)^2}{3(1+2w)} \left(\frac{\langle \epsilon \tau^3 \rangle \langle \epsilon \tau \rangle}{\langle \epsilon \tau^2 \rangle^2} - \frac{\langle \epsilon \tau \rangle^2 \langle \epsilon \Pi \tau^3 \rangle}{\langle \epsilon \Pi \tau \rangle \langle \epsilon \tau^2 \rangle^2} \right) \right], \quad (66)$$

$$\begin{aligned} \Delta Q_p \Big|_{001}^{001} &= Q_p \frac{\Delta \rho}{\rho} - Q_p \left(\frac{\mu_H H}{c} \right)^2 \\ &\quad \cdot \frac{2(2+w)^2(1-w)(1-pw)}{3(2+pw)(1+2w)^2} \frac{\langle \epsilon \tau \rangle^2 \langle \epsilon \Pi \tau^3 \rangle}{\langle \epsilon \Pi \tau \rangle \langle \epsilon \tau^2 \rangle^2} \\ &= Q_p \left(\frac{\mu_H H}{c} \right)^2 \frac{2(2+w)(1-w)}{3(1+2w)^2} \left[\frac{(1-w)\langle \epsilon \tau \rangle \langle \epsilon \tau^3 \rangle}{\langle \epsilon \tau^2 \rangle^2} \right. \\ &\quad \left. - \frac{(2+w)(1-pw)}{(2+pw)} \frac{\langle \epsilon \tau \rangle^2 \langle \epsilon \Pi \tau^3 \rangle}{\langle \epsilon \Pi \tau \rangle \langle \epsilon \tau^2 \rangle^2} \right]. \end{aligned} \quad (67)$$

It is perhaps worth remarking that there exist a variety of checks which can be applied to calculations of this sort. For example, ΔQ should always vanish if all $\mathbf{\Pi}_\sigma$ are the same multiple of the unit tensor. Again, setting $\mathbf{\Pi} \propto \mathbf{m} \cdot \boldsymbol{\tau}^{-1}$ should give the magnetoresistance.¹ These and other tests of the same sort are very helpful in uncovering algebraic errors.

B.7 Large H , Basic Formulas

The leading term in the expansion of the solution of (46) in powers of H^{-1} is especially simple when the geometry is such that \mathbf{E} is either exactly parallel or exactly perpendicular to \mathbf{H} . As these cases include all the orientations dealt with in this paper, we shall give first some formulas based on these specializations.

When \mathbf{E} and \mathbf{H} are parallel, the part of $d\mathbf{j}$ normal to \mathbf{H} must be of order EH^{-1} , since, if it were of lower order in H^{-1} , the term of (46) in $d\mathbf{j} \times \mathbf{H}$ could not be compensated by anything else. Therefore, only the part $d\mathbf{j}_H$ of $d\mathbf{j}$ parallel to \mathbf{H} can be of order EH^0 . It is obtained by taking

the component of (46) along \mathbf{H} :

$$dj_H = e^2 \left(\frac{\epsilon}{\langle \epsilon \rangle} \right) \frac{E_H dn}{(\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1})_{HH}}, \quad \text{if } \mathbf{E} \parallel \mathbf{H}, \quad (68)$$

with neglect of terms of order H^{-1} .

When \mathbf{E} and \mathbf{H} are perpendicular, as they always are when \mathbf{H} is perpendicular to a symmetry plane and the total \mathbf{j} is in this plane, we can start the solution of (46) by taking its cross product with the unit vector $\mathbf{u}^{(H)}$ in the direction of \mathbf{H} . The Hall field now dominates \mathbf{E} , so all components of $d\mathbf{j}$ are of order EH^{-1} , and the component normal to \mathbf{H} measures in essence the Hall velocity cEH^{-1} . We find, to the first order in H^{-1} ,

$$(\text{part of } d\mathbf{j} \perp \mathbf{H}) = \pm ec \left(\frac{\epsilon}{\langle \epsilon \rangle} \right) \mathbf{u}^{(H)} \times \mathbf{E} H^{-1} dn \quad \text{if } \mathbf{E} \perp \mathbf{H}, \quad (69)$$

where the upper sign is for electrons, the lower for holes. The part of $d\mathbf{j}$ parallel to \mathbf{H} is now obtained by taking the dot product of (46) with $\mathbf{u}^{(H)}$ and using (69). The result of combining the two is, to the first order,

$$d\mathbf{j} = \pm ec \left(\frac{\epsilon}{\langle \epsilon \rangle} \right) \left[\mathbf{u}^{(H \times E)} - \frac{(\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1})_{H \times E, H} \mathbf{u}^{(H)}}{(\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1})_{HH}} \right] \frac{E}{H} \quad \text{if } \mathbf{E} \perp \mathbf{H}, \quad (70)$$

where now we have used the label $H \times E$ in subscripts and superscripts to designate the corresponding coordinate direction, just as we used H above; $\mathbf{u}^{(H \times E)}$ is the unit vector in this direction.

The expressions (68) and (70) will suffice for the calculation of the high-field limit of Q ; to get the high-field behavior of the Nernst effect, however, for Table II, we need the conductivity tensor to one higher order. Rather than writing this down explicitly, we shall outline the result of a systematic solution of (46) by iteration of the procedures just used. Let us write for the conductivity tensor of the shell

$$d\sigma(\mathbf{H}) = d\sigma^{(\infty)} + d\sigma^{(-1)} + d\sigma^{(-2)} + \dots, \quad (71)$$

where the superscripts ∞ , -1 , -2 , etc. denote respectively terms of zeroth, first, second, etc. orders in H^{-1} . With a corresponding notation for the components of $d\mathbf{j}$, we find that, for $k > 1$, the part of $d\mathbf{j}^{(-k)}$ perpendicular to \mathbf{H} can be obtained in terms of $d\mathbf{j}^{(-k+1)}$ by taking the cross product of (46) with $\mathbf{u}^{(H)}$, and the part parallel to \mathbf{H} from both this result and the dot product of (46) with $\mathbf{u}^{(H)}$. The general term can be written concisely but abstractly in the form

$$d\delta^{(-k)} = - \left(\frac{\pm c}{eH} \right)^k e^2 \left(\frac{\epsilon}{\langle \epsilon \rangle} \right) (\mathbf{G} \cdot \mathbf{U} \cdot \mathbf{S})^{k-1} \cdot \mathbf{G} \cdot \mathbf{U} \cdot \mathbf{G}^+ dn, \quad (72)$$

if $k \geq 1$; the direction of \mathbf{E} is now unrestricted. Here \mathbf{S} stands for the tensor $\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1}$, while

$$G_{\alpha\beta} = \delta_{\alpha\beta} - S_{HH}^{-1} u_\alpha^{(H)} \sum_\lambda u_\lambda^{(H)} S_{\lambda\beta}, \quad (73)$$

\mathbf{G}^+ is the transposed tensor, and

$$U_{\beta\gamma} = \sum_\mu \delta_{\beta\gamma\mu} u_\mu^{(H)}. \quad (74)$$

By use of the identity $G^+S = SG$ one can easily show that (72) is symmetrical for even k , antisymmetrical for odd, as it should be.

b.s Limits of Q_e and Q_p at Large H

For either Q_p or Q_e , or their sum, we have, from (6) and the Kelvin relation, the high-field limit $Q^{(\infty)}$ given by

$$TQ_{\alpha\alpha}^{(\infty)} = \sum_g (\Pi_g \cdot \boldsymbol{\rho}_g^{(\infty)} \cdot \boldsymbol{\rho}^{(\infty)})_{\alpha\alpha} + \sum_g (\Pi_g \cdot \boldsymbol{\rho}_g^{(-1)} \cdot \boldsymbol{\rho}^{(+1)})_{\alpha\alpha}, \quad (75)$$

where α labels the rectangular direction in which Q is measured, and where, by analogy with the notation already used, $\boldsymbol{\rho}^{(\infty)}$ is the limiting resistivity tensor at infinite H and $\boldsymbol{\rho}^{(+1)}$ is the high-field Hall tensor (12).

When the α direction is parallel to \mathbf{H} , the only nonvanishing component of $\boldsymbol{\rho}_g^{(\infty)}$ is, by (68), the $\alpha\alpha$ one. Moreover, in this case $\rho_{\beta\alpha}^{(+1)}$ vanishes for all β , so only the first term of (75) remains. In this term we can replace $\rho_{\alpha\alpha}^{(\infty)}$ by $(\sum_g \sigma_{g\alpha\alpha}^{(\infty)})^{-1}$, so we have

$$TQ_{\alpha\alpha}^{(\infty)} = \frac{\sum_i \langle \epsilon \Pi_{\alpha\alpha}^{(i)} / (\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1})_{\alpha\alpha}^{(i)} \rangle}{\sum_i \langle \epsilon / (\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1})_{\alpha\alpha}^{(i)} \rangle} \quad \text{if } \mathbf{H} \parallel \alpha, \quad (76)$$

where i labels the different valleys.

If the different components of $\boldsymbol{\tau}$ all have the same energy dependence, the relative magnitudes of the currents (68) for shells of different energies will be the same as for $H = 0$, so, as (76) and (50) show explicitly,

$$Q_e^{(\infty)} = Q_e(0) \text{ parallel to } \mathbf{H}. \quad (77)$$

If the valleys have [111] orientations in a cubic crystal, the general form (76) becomes, for symmetry directions,

$$TQ_p^{(\infty)} \Big|_{001}^{001} = \frac{\langle \epsilon[\Pi_{p\parallel} + 2\Pi_{p\perp}][(m_{\parallel}^*/\tau_{\parallel}) + 2(m_{\perp}^*/\tau_{\perp})]^{-1} \rangle}{3\langle \epsilon[(m_{\parallel}^*/\tau_{\parallel}) + 2(m_{\perp}^*/\tau_{\perp})]^{-1} \rangle}, \quad (78)$$

$$TQ_p^{(\infty)} \Big|_{011}^{011} = \frac{\langle \epsilon\Pi_{p\perp}\tau_{\perp}/m_{\perp}^* \rangle + \langle \epsilon[\Pi_{p\perp} + 2\Pi_{p\parallel}][(m_{\perp}^*/\tau_{\perp}) + 2(m_{\parallel}^*/\tau_{\parallel})]^{-1} \rangle}{\langle \epsilon\tau_{\perp}/m_{\perp}^* \rangle + 3\langle \epsilon[(m_{\perp}^*/\tau_{\perp}) + 2(m_{\parallel}^*/\tau_{\parallel})]^{-1} \rangle}, \quad (79)$$

$$TQ_p^{(\infty)} \Big|_{111}^{111} = \frac{\langle \epsilon\Pi_{p\parallel}\tau_{\parallel}/m_{\parallel}^* \rangle + 3\langle \epsilon[\Pi_{p\parallel} + 8\Pi_{p\perp}][(m_{\parallel}^*/\tau_{\parallel}) + 8(m_{\perp}^*/\tau_{\perp})]^{-1} \rangle}{\langle \epsilon\tau_{\parallel}/m_{\parallel}^* \rangle + 27\langle \epsilon[(m_{\parallel}^*/\tau_{\parallel}) + 8(m_{\perp}^*/\tau_{\perp})]^{-1} \rangle}. \quad (80)$$

As usual, these equations also hold for Q_e and Q_p individually, if Π_{\parallel} and Π_{\perp} are given the corresponding subscript. Further simplifications are possible if the anisotropies are energy-independent, as assumed in Table III, or if the τ 's are $\propto \epsilon^{-1/2}$, etc.

When the α direction is perpendicular to \mathbf{H} and \mathbf{H} is normal to a symmetry plane, (70) suffices for the evaluation of (75). For this case, the first term of (75) vanishes, since $\sigma_{\sigma}^{(\infty)}$ has only an HH component and $\rho^{(\infty)}$ has no αH component. We find easily

$$TQ_{\alpha\alpha}^{(\infty)} = N_{\nu}^{-1} \langle \epsilon \rangle^{-1} \sum_i \left\langle \left(\Pi_{\alpha\alpha}^{(i)} - \frac{S_{\alpha H}^{(i)}}{S_{HH}^{(i)}} \right) \Pi_{\alpha H}^{(i)} \right\rangle \quad (\alpha \perp \mathbf{H}), \quad (81)$$

where i runs over the different valleys, N_{ν} in number, and where, as above, $\mathbf{S}^{(i)}$ is the tensor $\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1}$ for the i th valley. This holds for any multivalley structure when the symmetry is such that \mathbf{E} is exactly $\perp \mathbf{H}$. The entries of Tables III and IX are specializations of (81).

B.9 Asymptotic Behavior of B as $H \rightarrow \infty$

Let \mathbf{H} be in the z -direction, ∇T in the x -direction. From (1), (5), (6) and the symmetries which we proved above from (72), we find

$$T \left[\frac{\partial(BH)}{\partial(H^{-1})} \right]_{\infty} = H \left[\sum_{\sigma} \Pi_{\sigma} \cdot (\sigma_{\sigma}^{(-2)} \cdot \rho^{(+1)} + \sigma_{\sigma}^{(-1)} \cdot \rho^{(\infty)}) \right]_{xy}, \quad (82)$$

where, it will be remembered, the superscripts -2 , -1 , ∞ , $+1$ denote quantities going as H^{-2} , H^{-1} , H^0 , H as $H \rightarrow \infty$. Here we may use (12) for $\rho^{(+1)}$ and get $\sigma^{(-1)}$ from (70). Since $\rho^{(+1)}$ has only xy and yx components, the first term on the right of (82) involves just $(\sum_{\sigma} \Pi_{\sigma} \cdot \sigma_{\sigma}^{(-2)})_{xx}$. Use of the explicit evaluation of $\delta_{\sigma}^{(-2)} = d\delta^{(-2)}$ from (72) gives

$$\begin{aligned}
 (\Pi_g \cdot d\delta^{(-2)})_{xx} = & -\frac{c^2 \epsilon}{H^2 \langle \epsilon \rangle} dn \left[\left(\Pi_{gxy} - \frac{\Pi_{gzz} S_{zy}}{S_{zz}} \right) \left(S_{xy} - \frac{S_{xz} S_{yz}}{S_{zz}} \right) \right. \\
 & \left. - \left(\Pi_{gxx} - \frac{\Pi_{gzz} S_{zz}}{S_{zz}} \right) \left(S_{yy} - \frac{S_{yz}^2}{S_{zz}} \right) \right], \tag{83}
 \end{aligned}$$

where, as above, \mathbf{S} is the tensor $\mathbf{m}^* \cdot \boldsymbol{\tau}^{-1}$ for group g . If the directions of \mathbf{H} and ∇T are such that $\rho^{(\infty)}$ is diagonal in our coordinate system, the second term on the right of (82) involves $(\Pi_g \cdot \sigma_g^{(-1)})_{xy}$, for which we find

$$(\Pi_g \cdot d\delta^{(-1)})_{xy} = \pm \frac{ce\epsilon}{H \langle \epsilon \rangle} dn \left(\Pi_{gxx} - \frac{\Pi_{gzz} S_{zz}}{S_{zz}} \right), \tag{84}$$

where, as always, the upper sign is for electrons, the lower for holes. By the use of (50) and the relation

$$\frac{c}{\mu} = \mp \frac{\rho^{(0)}}{R(\infty)} \tag{85}$$

we can evaluate (82) in the form,

$$\begin{aligned}
 \frac{1}{Q(0)} \left[\frac{\partial(BH)}{\partial(c/\mu H)} \right]_{\infty} = & \mp \frac{\langle \epsilon (S_{\parallel}^{-1} + 2S_{\perp}^{-1})^2 \rangle}{3 \langle \epsilon \left(\frac{\Pi_{\parallel}}{S_{\parallel}} + \frac{2\Pi_{\perp}}{S_{\perp}} \right) \rangle \langle \epsilon \rangle^2 N_v} \\
 & \cdot \sum_i \left\langle \epsilon \left(\Pi_{xy}^{(i)} - \frac{\Pi_{zz}^{(i)} S_{zy}^{(i)}}{S_{zz}^{(i)}} \right) \left(S_{xy}^{(i)} - \frac{S_{xz}^{(i)} S_{yz}^{(i)}}{S_{zz}^{(i)}} \right) \right. \\
 & \left. - \epsilon \left(\Pi_{xx}^{(i)} - \frac{\Pi_{zz}^{(i)} S_{zz}^{(i)}}{S_{zz}^{(i)}} \right) \left(S_{yy}^{(i)} - \frac{S_{yz}^{(i)2}}{S_{zz}^{(i)}} \right) \right\rangle \\
 & \mp \frac{\langle \epsilon (S_{\parallel}^{-1} + 2S_{\perp}^{-1}) \rangle}{\langle \epsilon \left(\frac{\Pi_{\parallel}}{S_{\parallel}} + \frac{2\Pi_{\perp}}{S_{\perp}} \right) \rangle \langle \epsilon \rangle N_v} \sum_i \left\langle \epsilon \left(\Pi_{xx}^{(i)} - \frac{\Pi_{zz}^{(i)} S_{zz}^{(i)}}{S_{zz}^{(i)}} \right) \right\rangle, \tag{86}
 \end{aligned}$$

where the superscript i labels the different valleys, N_v in number. This equation applies to B_e or B_p or their sum, if all Π 's, B and Q are given the appropriate subscript; it is valid for any multivalley structure with over-all cubic symmetry if z is normal to a symmetry plane and x is along a two-fold or four-fold axis in the plane. The entries in Table II result from tedious but straightforward specializations of (86) to [111] valleys and energy-independent anisotropies.

B.10 Maxwellian Averages

As an appendix to this appendix we list here the Maxwellian averages of various powers of the energy, which are needed for the evaluation of

TABLE XV — MAXWELLIAN AVERAGES OF POWERS OF ENERGY

n	$\langle (\frac{\epsilon}{kT})^n \rangle$
-1	2
$-\frac{3}{4}$	1.38
$-\frac{1}{2}$	$2\pi^{-1/2}$
$-\frac{1}{4}$	1.02
1	1
$\frac{1}{4}$	1.04
$\frac{1}{2}$	$2\pi^{-1/2}$
$\frac{3}{4}$	1.28
1	$\frac{3}{2}$
$\frac{5}{4}$	1.81
$\frac{3}{2}$	$4\pi^{-1/2}$
$\frac{7}{4}$	2.88
2	$\frac{15}{4}$
$\frac{5}{2}$	4.99
$\frac{3}{2}$	$12\pi^{-1/2}$

the various expressions we have derived when the τ 's and Π 's are assumed to have power-law dependences on energy. Since the number of Maxwellian carriers in the energy range ϵ to $\epsilon + d\epsilon$ is proportional to $\epsilon^{\frac{3}{2}} \exp(-\epsilon/kT)d\epsilon$, we have for any quantity Z

$$\langle Z \rangle \equiv \frac{\int_0^{\infty} Z \epsilon^{\frac{3}{2}} \exp(-\epsilon/kT) d\epsilon}{\int_0^{\infty} \epsilon^{\frac{3}{2}} \exp(-\epsilon/kT) d\epsilon} \quad (87)$$

The results for various powers of the energy are given in Table XV.

REFERENCES

- Herring, C., Geballe, T. H. and Kunzler, J. E., Phys. Rev., **111**, 1958, p. 36.
- Herring, C., in *Halbleiter und Phosphore*, F. Vieweg & Sohn, Braunschweig, Germany, 1958, p. 184.
- Herring, C. and Vogt, E., Phys. Rev., **101**, 1956, p. 944.
- Herring, C., to be published.
- Argyres, P. N. and Adams, E. N., Phys. Rev., **104**, 1956, p. 900.
- Appel, J., Z. Naturforsch., **11a**, 1956, p. 892.
- Argyres, P. N., J. Phys. Chem. Solids, **4**, 1958, p. 19.
- Pollak, H. O. and Herring, C., unpublished calculations.
- Klinger, M. I. and Voronyuk, P. I., J. Exp. Theor. Phys. (U.S.S.R.), **33**, 1957, p. 77. Translation: Soviet Phys. JETP, **6**, 1958, p. 59.
- Argyres, P. N., Phys. Rev., **109**, 1958, p. 1115.
- Wolff, P. A., unpublished work.
- Furth, H. P. and Wanick, R. W., Phys. Rev., **104**, 1956, p. 343.

13. Herring, C., B.S.T.J., **34**, 1955, p. 237.
14. Dorn, D., Z. Naturforsch., **12a**, 1957, p. 18.
15. Herring, C., to be published.
16. Tanenbaum, M. and Hrostowski, H. J., unpublished measurements.
17. Goldberg, C., Phys. Rev., **109**, 1958, p. 331.
18. Goldberg, C. and Howard, W. E., Phys. Rev., **110**, 1958, p. 1035.
19. Glicksman, M., Phys. Rev., **108**, 1957, p. 264; *Progress in Semiconductors*, Vol. 3, Heywood & Co., London, 1958.
20. Conwell, E. M. and Weisskopf, V. F., Phys. Rev., **77**, 1950, p. 388.
21. Brooks, H., Phys. Rev., **83**, 1951, p. 879; *Advances in Electronics and Electron Physics*, Vol. 7, Academic Press, New York, 1955, p. 128.
22. Dingle, R. B., Phil. Mag., **46**, 1955, p. 831.
23. Sclar, Phys. Rev., **104**, 1956, p. 1548.
24. Blatt, F. J., J. Phys. Chem. Solids, **1**, 1957, p. 262.
25. Debye, P. P. and Conwell, E. M., Phys. Rev., **93**, 1954, p. 693.
26. Fan, H. Y., *Solid State Physics*, Vol. 1, Academic Press, New York, 1955, p. 283; Conwell, E. M., Proc. I.R.E., **46**, 1958, p. 1281.
27. Spitzer, L., Jr. and Härm, R., Phys. Rev., **89**, 1953, p. 977.
28. Madelung, O., *Handbuch der Physik*, Vol. 20, Springer-Verlag, Berlin, 1957, p. 1.
29. Dingle, R. B., Arndt, D. and Roy, S. K., Appl. Sci. Res., **B6**, 1956, p. 155.
30. Beer, A. C., Armstrong, J. A. and Greenberg, I. N., Phys. Rev., **107**, 1957, p. 1506.
31. Mansfield, R., Proc. Phys. Soc., **B69**, 1956, p. 862.
32. Herring, C., Phys. Rev., **95**, 1954, p. 954.
33. Brockhouse, B. N. and Iyengar, P. K., Phys. Rev., **111**, 1958, p. 747.
34. Pomeranchuk, I., J. Phys. (U.S.S.R.), **4**, 1941, p. 259; **6**, 1942, p. 237; Phys. Rev., **60**, 1941, p. 820.
35. Pomeranchuk, I., J. Phys. (U.S.S.R.), **4**, 1941, p. 529.
36. Dresselhaus, G., Kip, A. F. and Kittel, C., Phys. Rev., **98**, 1955, p. 368.
37. Morin, F. S., Geballe, T. H. and Herring, C., Phys. Rev., **105**, 1957, p. 525.

Stabilization of Silicon Surfaces by Thermally Grown Oxides*

By M. M. ATALLA, E. TANNENBAUM and E. J. SCHEIBNER

(Manuscript received January 7, 1959)

A study has been carried out of the stability of silicon surfaces when they are provided with a chemically bound solid-solid interface. Stable surfaces have been obtained with the system silicon-silicon dioxide when the oxide is thermally grown. This latter system has been studied in some detail. In this paper the following phases of our investigation are presented: (i) some aspects of the thermal oxidation process and properties of the oxide; (ii) the electronic properties of the resulting silicon-silicon dioxide interface; (iii) the application of the process to devices and resulting device characteristics.

I. INTRODUCTION

In this introduction we will give a qualitative discussion of the problem of stabilization of semiconductor surfaces.

1.1 *Atomically Clean Surfaces*

Fig. 1(a) is a schematic diagram of a clean [100] silicon surface showing the surface dangling bonds or unfilled orbitals. Since two electrons can occupy a free orbital, surface atoms may become negatively charged, thus acting as acceptor surface states [Fig. 1(b)]. The existence of these states in crystals was first proposed by Tamm,¹ based on a special one-dimensional model. A more general treatment by Shockley² showed that surface states can occur only if there is a separate potential trough at the surface or if the energy bands arising from separate atomic levels overlap. He further concluded that surface states should occur in the forbidden gap and their number should be approximately equal to the number of surface atoms.

Experimentally, therefore, one expects to find that an atomically clean

* This work was supported in part by the Department of the Army under Contract DA 36-039 sc-64618.

surface would be strongly p-type, corresponding to a surface state density roughly equal to the density of surface atoms [Fig. 1(c)], and strongly sensitive to surface effects such as chemisorption. Measurements on germanium surfaces prepared by the Farnsworth³ technique (argon bombardment, then annealing at 500°C in high vacuum), of work function,^{4,5,6} photoconductance,^{7,8,9,10} surface conductivity^{9,10,11} and field effect^{7,9,11} have indicated that there is a large density of surface states and that these surface states are the dangling bonds which act as acceptor-type states by trapping bulk electrons and forming a p-type surface. The data also indicate that the adsorption of gases which can bond co-

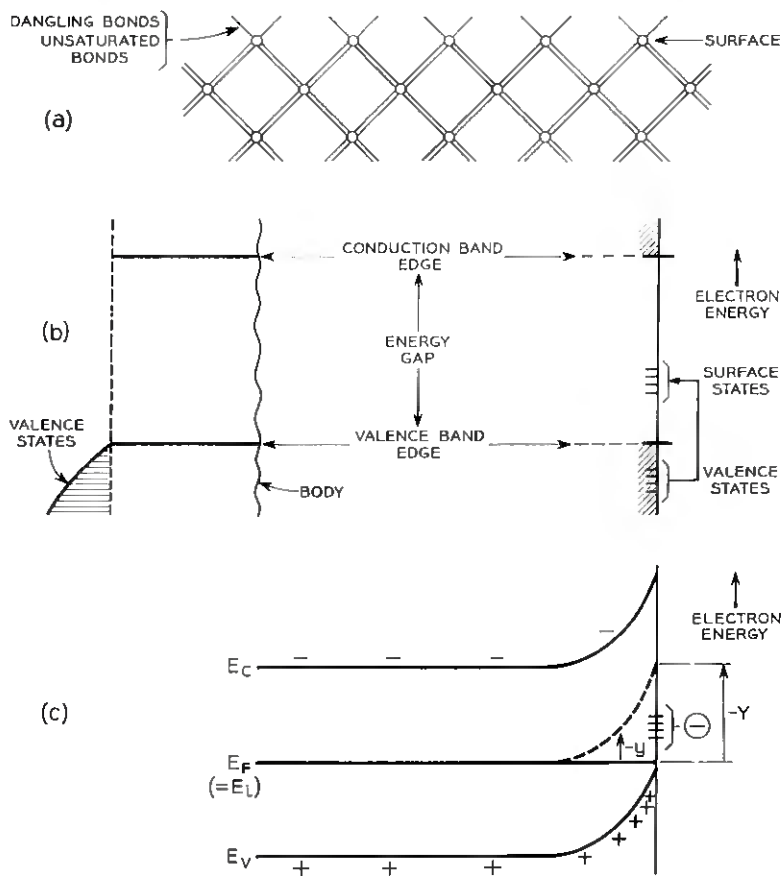


Fig. 1 — (a) Atomic diagram of [100] surface; (b) surface states; (c) resulting band bending near the surface of an atomically clean intrinsic material.

valently with these unfilled orbitals affects the density of surface states and the electrical properties of the surface. More recent measurements,¹² still only preliminary, on cleaved germanium surfaces in high vacuum, appear to be in agreement with those obtained on bombarded and annealed surfaces.

Measurements on atomically clean silicon surfaces are not yet as complete as are those on germanium surfaces. There is strong evidence that clean silicon surfaces are obtainable by high-temperature, high-vacuum heating, as indicated by field emission patterns,^{13,14} secondary electron emission¹⁵ and low energy electron diffraction.¹⁶ From the nature of the surface states on germanium surfaces as outlined above, one would expect considerable resemblance between clean germanium surfaces and silicon surfaces.

From the standpoint of device technology, one may expect that truly clean surfaces are not desirable. A junction device would not be operative if its surfaces were atomically cleaned, since the resulting surfaces would be strongly p-type (as present evidence indicates), regardless of the body doping, and a low-resistance surface path would shunt the device.

1.2 *Practical Surfaces*

It is obvious from the above discussion that the usual processes of "surface cleaning" as used in device preparation, such as etches or rinses, do not provide "surface cleaning" in its true sense. A more appropriate term is "surface doping," since such processes generally provide surface structures, complex and unknown, which are found empirically to produce surface properties compatible with the device body properties (through the attainment of a suitable density, distribution and type of surface states, and hence the term surface doping).

Fig. 2(a) shows a possible distribution of surface states on an atomically clean surface. These are all acceptor-type states (neutral when unoccupied, negative when occupied). Chemisorption of some foreign species on the surface may produce the following changes [see Fig. 2(b)]: (a) removal of some of the original acceptor states, (b) introduction of new acceptor and donor states (neutral or positive) and (c) introduction of outer states located beyond the interface, within or on the surface of a grown film. Due to their usually longer relaxation times, these are called "slow" states, in contrast to the interface or "fast" states [see Fig. 2(c)].

For a surface in equilibrium with the body, the potential distribution near the surface is determined by the density, distribution and type of all the surface states. In contrast to the atomically clean surface, which develops a strongly p-type surface conductance [Fig. 1(c)], the surface

conductance of a practical surface may be either p-type or n-type, depending on the existing surface states as determined by surface treatment. This important concept was first appreciated by Bardeen.¹⁷ Fig. 3 illustrates potential distributions at the surfaces of n-type and p-type materials when *donor*-type surface states predominate, hence producing a net positive surface charge which must be compensated by an equal negative charge in a surface space-charge layer.

A point of particular interest and basic significance follows from the above. If a certain surface treatment produces a predominance of either donor- or acceptor-type states (irrespective of body doping), it is evident that the surfaces near the various regions of a device (p and n) can, at best, provide a compromise compatible with the device. For instance, for a p-n junction with predominantly donor-type surface states (net positive charge), the tendency is to make the p-region less p (and possibly inverted), while the n-region becomes more n. For a multiple junction

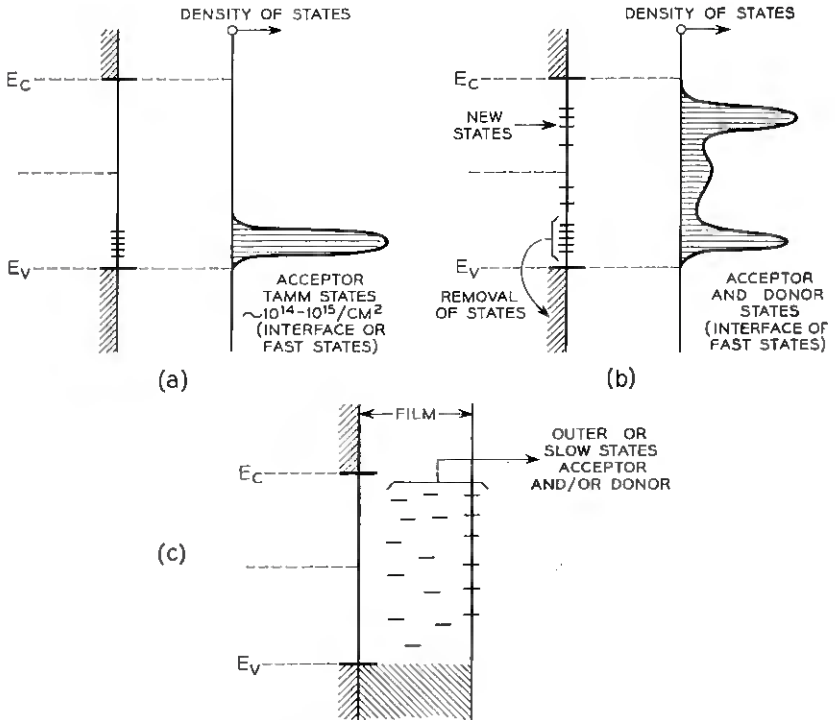


Fig. 2 — (a) Surface states on an atomsically clean surface; (b) interface or fast states on a practical surface; (c) outer or slow states on a practical surface.

device it is obvious that the problem of simultaneous compatibility of the surface with the various regions may become quite critical, particularly with high-resistivity materials. This point will be considered in more detail later.

1.3 Nature of the Surface Stability Problem

Practical surfaces obtained by the usual empirical techniques involving etches, rinses, etc., are known to show various degrees of instability. They are generally sensitive to various ambients producing drifts in surface properties that are usually not completely reversible. There are at least three main reasons for surface instability:

i. Replacement of interface impurities where an impurity *B* may replace a previously existing surface impurity *A*, depending on their relative affinities for the semiconductor. This, in general, will modify the distribution and density of the interface or fast states and, accordingly, the surface potential.

ii. Ionization of adsorbed neutral impurities. This requires both the availability of active sites at the free surface and electronic exchange across a surface film between the adsorbed species and the semiconductor. Depending on its relative electronic affinity, a neutral molecule located at an active surface site (constituting an electronic state) will either re-

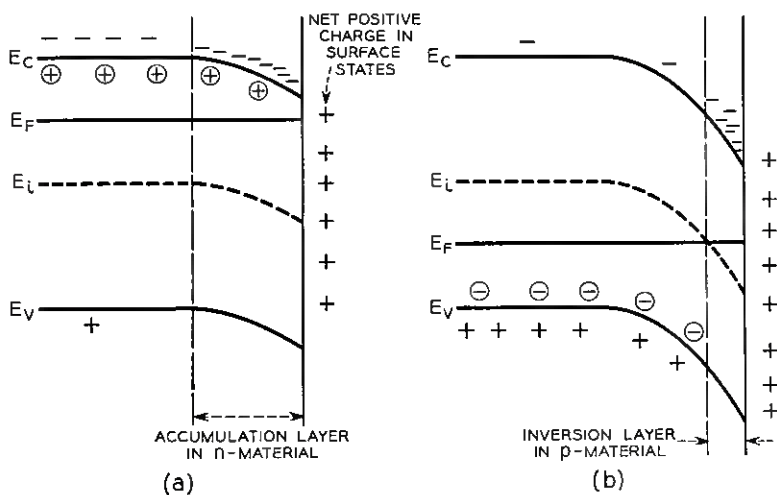


Fig. 3 — Band bending at the surfaces of p and n semiconductors due to surface states that are predominantly donor-type.

ceive an electron from the semiconductor, producing a negative surface charge, or give up an electron to the semiconductor, producing a positive surface charge. Tunneling or direct conduction are probably the mechanisms responsible for the electronic exchange. These (i and ii) are the outer or slow states.

iii. Replacement of these outer impurities may also occur, producing a change in the density, distribution and type of the slow states.

In addition to the above mechanisms, one might also consider a distinctly different process of instability which is generally operative in the vicinity of a junction. Here, surface ions produced by processes other than electronic exchange with the semiconductor, such as hydrolysis, etc., can migrate on the surface under the influence of the fringing field of the junction. A second paper to be published is entirely devoted to a detailed study of this process.

II. TECHNIQUES FOR MEASUREMENT OF SURFACE STATES

In the course of this work several techniques have been used for measuring the detailed surface properties of silicon surfaces, particularly as related to surface states. These included the large-signal ac field effect technique,^{18,19} the large-signal ac field effect technique with superimposed dc bias, a zero-frequency four-point probe field effect technique, the surface photovoltage technique²⁰ and the n-p-n or p-n-p channel technique.^{21,22} The majority of the results presented here, however, have been obtained by the field effect techniques. It seems in order, therefore, to present here a brief review of the basic principles and concepts underlying field effect work.

2.1 *Basic Concepts of the Field Effect*

The object of this technique is to determine the density and distribution of surface states. Application of a field transverse to the semiconductor surface produces a change in the band bending near the surface to maintain charge neutrality. This, in effect, alters the location of the Fermi level at the surface and, accordingly, changes the occupancy of the surface states. One measures changes in surface conductance with field (or plate voltage), from which the density and distribution of states are obtained.

2.1.1 *Field Effect with No Surface States*

Fig. 4(a) shows the bands near the surface of an intrinsic semiconductor which are straight in the absence of surface states or charge on the

field plate. The conductance of a region near the surface is due to electrons in the conduction band and an equal number of holes in the valence band. When a voltage is applied to the plate with respect to the semiconductor the bands near the surface must adjust to the point where the net charge in the semiconductor is equal and opposite to that on the plate, Fig. 4(b). If ΔN and ΔP are the changes in the number of electrons and holes per unit area, due to plate charge Q_P per unit area, the corresponding change in surface conductance ΔG is given by

$$\Delta G = q(\mu_n \Delta N + \mu_p \Delta P). \tag{1}$$

One defines the derivative dG/dQ_P , as the "field effect mobility" $\mu_{F.E.}$. It can be shown¹⁸ that ΔG exhibits a minimum value that corresponds to a unique potential distribution near the surface for each body resistivity. For intrinsic material, the bands are bent by an amount $Y = -\ln b$, where $b = \mu_n/\mu_p$. If the bands are so strongly bent in either direction that the number of one type of carrier is negligibly small as compared with the number of the other, then a change in field-plate charge will induce an equal and opposite change in the majority carrier and the field effect mobility approaches the surface mobility of that carrier. Fig. 4(c)

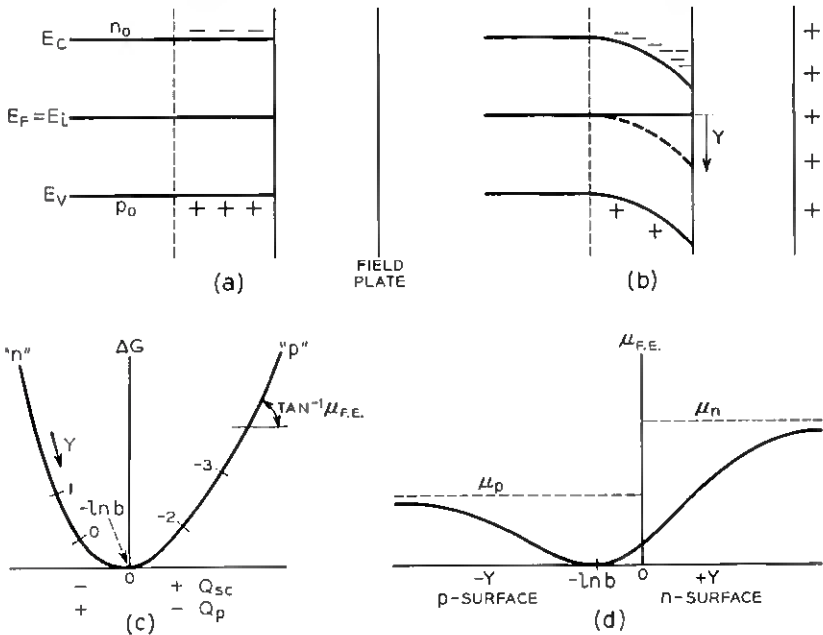


Fig. 4 — Field effect in the absence of surface states.

shows the variation of conductance with plate charge (equal and opposite to space charge Q_{sc}) with the corresponding values of surface potential Y [in $(kt)/q$ units]. In Fig. 4(d) the changes in field effect mobility with surface potential are shown. For any other resistivity one can calculate the variation in surface conductance with charge and the corresponding surface potential using the space-charge equations derived by Garrett and Brattain,²³ which are also conveniently available in graph form.²⁴ In many cases, one needs to include the effect of surface scattering on carrier mobility. No direct measurements are available, however, and the common practice is to apply Schrieffer's theoretical results.²⁵

2.1.2 Field Effect with Surface States

In the presence of surface states, part of the charge on the field plate will terminate on charges in surface states. If one now changes the charge by some amount, both the semiconductor space charge and charges in the surface states will change. The resulting change in surface conductance, however, will always be *smaller* than one which would be obtained in the absence of the surface states. The experimentally obtained curve of ΔG versus Q_p is generally wider than that calculated with no surface states. Experimentally,^{18,19} one varies Q_p at some frequency and observes on an oscilloscope the variation of ΔG versus Q_p . Only states with relaxation times shorter than the period of the frequency used will be effective. Slower states will only contribute a fixed charge.

The point of minimum surface conductance must be obtained in order to provide the reference point for the quantitative evaluation of the results.¹⁸ It has been pointed out in the preceding section that the minimum conductance corresponds to a unique value of the surface potential Y for a given resistivity. Furthermore, the conductance in general is a unique function of Y for each resistivity. Therefore, one adjusts the experimental curve of ΔG versus Q_p with the theoretical space charge curve along the ΔG -axis until the two minima are matched, as shown in Fig. 5. It must be noted that their lateral relative position is arbitrary. Now one transfers the points of various values of Y on the theoretical curve horizontally to the experimental curve, thus providing the value of surface potential at each point on the curve. Furthermore, the horizontal segments between the two curves at each value of Y represent essentially the net charge in surface states Q_{ss} (difference between total charge or plate charge and charge in space charge) with an unknown additive constant. From plots of dQ_{ss}/dY versus Y , one may then fit the results with an approximate distribution of surface states. Rigorously

speaking, an exact distribution is only obtainable if the experiment is carried out throughout the range of $Y = \pm \infty$.²⁶

2.1.3 Effect of Slow States

If one applies a step voltage to the plate, one will first get a corresponding change in surface conductance resulting from a change in surface potential, in conformity with the fast state occupancy requirements discussed above. If no other states are present, the new value of surface conductance would be maintained indefinitely. However, if slow states are present, their occupancy, which is determined by their location with respect to the Fermi level, must gradually change to correspond to the new equilibrium condition. Since the total charge in the surface states and the space charge must remain constant (equal and opposite to the plate charge), a change in occupancy of the slow states will cause a gradual redistribution of charge in the space-charge region and in the fast surface states. It is obvious that this will always correspond to a *decay* in the initial change in conductance. If the density of the slow states exceeds the plate charge density, the decay is almost complete, with a time constant of the order of the relaxation time of the slow states. *This is evidence of electronic exchange between the semiconductor and the slow or outer states.* The same phenomenon is also observed in three different ways in ac field effect measurements. First, if one superimposes a dc bias on the ac field signal, one observes an immediate shift of the field effect curve towards p or n, depending on the polarity of the dc bias. This is followed by a slow decay towards the initial curve. Release of the dc bias will simply reverse the effect, with a final return to the initial condition. Second, by changing the ambient, one may observe a shift to-

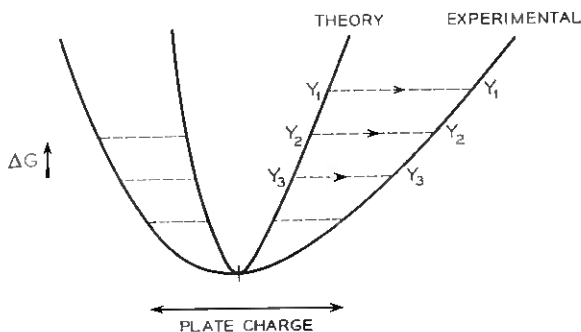


Fig. 5 — Determination of charges in surface states from field effect experiments.

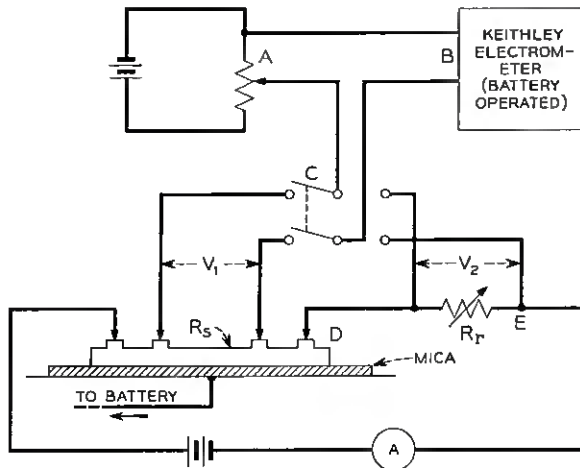


Fig. 6 — Circuit diagram of the zero-frequency four-point probe technique.

wards p or n, depending on the ambient.²⁷ This is an indication of a change in the density of the slow states and in the magnitude of the fixed charge they represent. It is a further indication of the ionization of neutral adsorbed impurities by electronic exchange with the semiconductor through the surface film. Third, if the frequency of the measurement is reduced to the range of frequencies of the slow states, one observes smaller changes in conductance corresponding to very low field effect mobilities.

2.2 Zero-Frequency Four-Point Probe Field Effect Technique

In many cases, particularly for device purposes, one requires a measurement of all the states, fast and slow. This may be done by the ac method, by repeating the measurements at various frequencies. Another alternative is the zero-frequency* four-point probe technique, shown in Fig. 6. The outer electrodes are used to feed in the current and the voltage drop is measured across the inner electrodes. The voltage drop across the specimen is initially balanced against a reference resistor R_r in series with the specimen. The change in conductance ΔG produced by a field is obtained from the voltage difference ΔV between the specimen and the reference resistor, $|V_2 - V_1|$:

$$\Delta G = \frac{\Delta V(L/w)}{I_0 R_s^2} \text{ mhos}/\square, \quad (2)$$

* The term "zero-frequency" is used to distinguish between this method and the conventional dc field effect method used to observe conductivity decay.

where L and w are the length and width of the specimen, I_0 is the current and R_s is the sample resistance. The plate charge is obtained from the voltage applied and the measured capacity of the plate specimen.

2.3 Characterization of a Surface Treatment

A complete characterization of a surface is a detailed description of the type, density and distribution of all the surface states. This allows one to predict the surface potential obtained for any resistivity material as a result of these same surface states. In many cases, however, particularly in device applications, it may suffice to know whether a certain surface treatment will or will not produce surface inversion. To extrapolate field effect results obtained on a material of some resistivity to other resistivities, one must show that the surface treatment and the resulting fast surface states are not dependent on concentration and type of body doping.

Fig. 7(a) shows a field effect curve produced at zero-frequency and hence includes all the states. Point o corresponds to zero charge on the field plate. Point I corresponds to an intrinsic surface or $Y_i = u_b$, where $u_b = \ln \lambda$ and $\lambda = p_0/n_i$, with p_0 the hole density in the body and n_i the intrinsic density. Fig. 7(b) corresponds to the "onset of surface inversion". To change the surface potential from its initial value Y_0 to the onset of inversion Y_i one must apply a charge $-Q_{0-I}$ to the plate. Invoking the condition of charge neutrality at point I:

$$[Q_{ss}]_I = Q_{0-I} - [Q_{sc}]_I, \tag{3}$$

where Q_{0-I} is measured and $[Q_{sc}]_I$ is calculated for the particular resistivity of the specimen used. This gives $[N_{ss}]_I = [Q_{ss}]_I/q$, the absolute

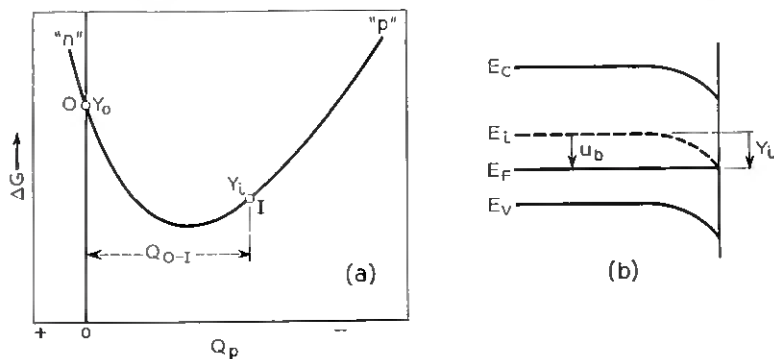


Fig. 7 — Determination of type and strength of surface treatment.

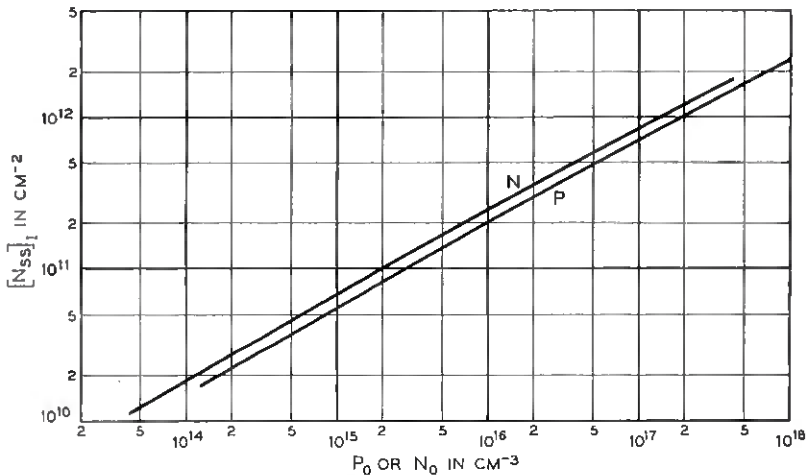


Fig. 8 — Net surface charge required for surface inversion.

number of charges in all the surface states when the Fermi level coincides with the intrinsic level at the surface corresponding to the onset of inversion.

Now, for any doping, the number of charges in the space charge at the onset of inversion is given by²³

$$(N_{sc})_I = n_i L \left[\lambda (e^{-u_b} - 1) + \frac{1}{\lambda} (e^{u_b} - 1) + \left(\lambda - \frac{1}{\lambda} \right) u_b \right]^{\frac{1}{2}}. \quad (4)$$

Hence, for any resistivity material, the surface states will just bring about the onset of inversion if $(N_{ss})_I$, obtained experimentally, is equal to $(N_{sc})_I$ in (4). For higher resistivity, inversion will occur and, for lower resistivity, inversion will not occur. Fig. 8 is a plot of (4) for both p- and n-type silicon.

The procedure we have used for characterizing the effect of a surface treatment is simply as follows:

i. Determine $[N_{ss}]_I$ from the field effect data as outlined above. If this number is positive, it indicates that inversion will occur on n-type material.

ii. From Fig. 8 determine p_0 or n_0 corresponding to $[N_{ss}]_I$. This determines the limiting resistivity below which inversion will not occur; all materials with resistivities above this value will form inversion layers. Now we will introduce some terminology which concisely describes the properties of an oxidation treatment as related to inversion: *p-type oxides* or *n-type oxides* — this indicates the tendency of the oxidation process to *induce* a p-type surface or an n-type surface, respectively, or to induce

inversion in n-type material or in p-type material, respectively. The *strength of an oxide* is defined as the limiting resistivity above which inversion will occur and below which no inversion will occur. To illustrate, let the oxide be characterized as n-type oxide of strength 10 ohm-cm. This means that this oxide produces sufficient surface states to invert the surface of all p-type material of resistivity above 10 ohm-cms.

III. THE OXIDATION PROCESS

3.1 *Introductory Remarks*

In the search for a coating for stabilization of silicon surfaces, several possibilities were considered. These included two general categories: coatings that are deposited by methods such as dipping, painting or evaporation followed by various processes such as baking, and coatings that are grown from the silicon surface itself. The first category generally suffers from such difficulties as: (a) ionic impurities (remembering that 1/10,000th of a monolayer is sufficient to invert the surface of 1 ohm-cm silicon) which will drift with field, producing device instability; (b) lack of uniformity of chemical bonding (if there is any bonding) at the silicon-film interface, particularly with an ill-defined initial silicon surface; (c) poor structures of the coatings from the standpoint of permeability to various ambients and (d) various incompatibilities such as dielectric strength or coefficients of expansion.

For coatings grown out of the silicon surface, we have restricted our efforts to silicon dioxide. We considered two alternatives: oxides grown anodically in electrolytes and thermally grown oxides. Measurements of their transverse conduction indicated substantial differences. The resistivity of the anodic oxide is of the order of 10^{12} ohm-cm, while that of the thermal oxide (grown at about 1000°C) is of the order of 10^{16} . Furthermore, in the anodic oxides, part of the conduction is due to mobile ionic impurities of the order of 10^{18} per cm^3 , while the ionic content of the thermal oxide is less than 10^{14} ions per cm^3 .

3.2 *Preoxidation Treatment of Silicon Surfaces*

Preliminary studies, using various silicon diodes, indicated a strong dependence of reverse characteristics on the surface treatment preceding the oxidation process. Impurities left on the surface can be partially maintained at the Si—SiO₂ interface through the oxidation process, and may affect the interface characteristics and the device characteristics. To illustrate, by purposely varying the preoxidation treatment on the

same junctions, it was possible to obtain a range of reverse current of 10^{-10} to 10^{-3} amperes. This work provided an extensive empirical background on the effect of preoxidation treatment and indicated the desirability of obtaining as clean a Si—SiO₂ interface as possible. Accordingly, all the details of the procedure prior to and during the oxidation process were chosen to achieve this objective. A suitable preoxidation process must be capable of: (a) removing surface imperfections, damage, etc.; (b) removing the bulk of organic surface residues (weakly bound); (c) removing chemically bound organic substances; (d) removing metallic impurities and (e) providing in a controlled fashion a lightly oxidized surface prior to thermal oxidation.

3.3 *Monitoring the Preoxidation Process*

In exploring the effects of preoxidation treatments, it appeared desirable to have some simple monitoring technique to indicate qualitatively certain effects of the various treatments on some property of the surface. A technique was used which determined whether the surface was hydrophobic or hydrophilic. The technique is based on the water break and water spray tests, whereby the contact angle of water droplets determines whether the surface is wettable (hydrophilic) or not (hydrophobic). The surface to be examined is dipped in liquid nitrogen for about 10 seconds, then mounted under a high-power microscope (400X) in a closed chamber where wet nitrogen is circulated. One first observes a thin sheet of ice forming uniformly on the surface. After the ice melts, a distinctive pattern of water droplets forms which remains for a few minutes before final evaporation. It has been observed that the shape and size of the droplets and the uniformity of the pattern is sensitive to the surface treatment. This process may be repeated (without re-treating the surface) about five to ten times before observable changes in the pattern start to occur, indicating surface deterioration.

The following is a summary of our observations: (i) A silicon surface that has been freshly etched in HNO₃—HF mixtures is strongly hydrophobic. A typical pattern of this surface as obtained by the "liquid nitrogen test" is shown in Fig. 9(a). The droplet size is of the order of 0.0001 inch. (ii) After the etched surface has been boiled in deionized water for as long as one hour, the surface remains hydrophobic. The size of the droplets, however, may slightly increase. (iii) After having been boiled in organic solvents the surface remains hydrophobic, although the shape and size of the droplets may change, [Fig. 9(b)]. (iv) Boiling in water containing detergents does not produce a hydrophilic surface. (v) When the surface is heated in hot oxidizing agents such as HNO₃ for

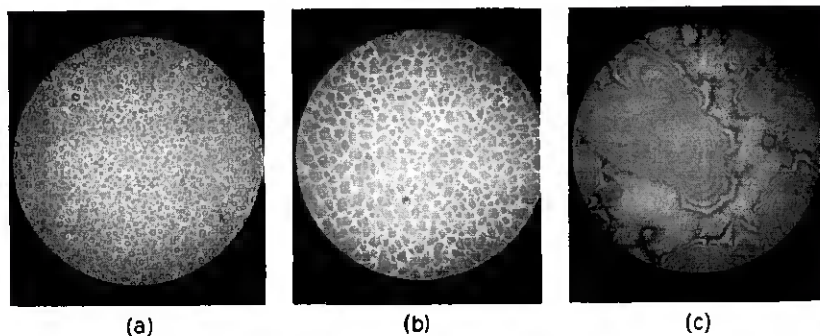


Fig. 9 — Silicon surface examination by the “liquid nitrogen test”: (a) after exposure to HF; (b) after subsequent boiling in xylene; (c) after boiling in nitric acid.

10 minutes or more the surface becomes uniformly hydrophilic. Fig. 9(c) shows the optical interference pattern of such a surface just before complete evaporation of the once continuous film of water. (vi) After having been boiled in water this hydrophilic surface remains hydrophilic. (vii) If, however, the hydrophilic surface is exposed for a few seconds to HF vapor, it reverts to a strongly hydrophobic surface. (viii) If the hydrophilic surface obtained by the HNO_3 treatment is exposed to room air, it becomes gradually and nonuniformly hydrophobic. Boiling the exposed surface in water or organic solvents will not restore the surface to its former hydrophilic state. (ix) A thermally oxidized surface is hydrophilic, and shows similar behavior to the nonoxidized one when treated with HNO_3 exposed to room air. A major difference, however, is that its original hydrophilic state can be largely restored simply by rinsing in an organic solvent. This suggests the relative inertness of an oxidized surface compared to a freshly etched surface.

3.4 *The Thermal Oxidation Process and Properties of the Thermal Oxide*

Silicon surfaces were oxidized at temperatures in the vicinity of 1000°C and at atmospheric pressure in both dry and wet oxygen. The oxidation rates in oxygen and in water vapor have been measured²⁸ using a vacuum microbalance technique. For the range of film thicknesses of about a hundred to thousands of angstroms the oxidation process follows the following parabolic laws:

$$\text{O}_2 : X^2 = 8.4 \times 10^{10} p^{4/5} t e^{(-1.7q)/kT}, \quad (5)$$

$$\text{H}_2\text{O} : X^2 = 2.54 \times 10^{13} p^{8/5} t e^{(-1.7q)/kT}, \quad (6)$$

where k is the Boltzman constant, T the absolute temperature, X the film thickness in angstroms, p the pressure in atmospheres, t the time in minutes and q the electronic charge. It is of particular interest to note that the activation energies are equal for both O_2 and H_2O oxidation, suggesting that the two types of oxides have the same structure. The high activation energy of 1.7 electron volts is also significant, since it presumably corresponds to extraction of silicon from the lattice and its diffusion in the oxide network.

The following are some properties of the thermal oxide:

i. Electron diffraction studies of thermally grown silicon oxides have indicated no crystalline structure. The film is essentially continuous and amorphous.

ii. The dielectric constant of the film is about 4. It has a resistivity of about 10^{16} ohm-cm, measured in a transverse direction to the oxidized surface. The dielectric strength of the film is between 5 and 10×10^6 volt/cm.

iii. If the surface impurities are not carefully removed prior to oxidation, oxide crystallites will form on the surface that become visible under the microscope if the film grown is of sufficient thickness. This provides rather serious discontinuities in the film.

iv. Discontinuities in the film are also observed at some surface imperfections, due either to the nonuniform growth of the oxide at the imperfection or to impurities trapped at the imperfection and subsequently



MAG 500 X

Fig. 10 — Oxide imperfections as revealed by the hot chlorine technique.

affecting the local oxidation. These local imperfections are made beautifully visible, even for films a few atomic layers thick, by the "chlorine etching" technique. At a temperature of about 900°C , chlorine gas will etch silicon at a rate of the order of 0.001 inch per minute but will not react with silicon dioxide. To test for film imperfections, the oxidized specimen is placed in the chlorine oven for a few minutes, allowing the chlorine to etch through discontinuities in the oxide. Fig. 10 shows a chlorine-treated oxide surface indicating some imperfections along well-defined surface damage.

v. Uniform oxides have been obtained consistently. They showed few imperfections as indicated by the chlorine test. Oxidized surfaces stored for more than 12 months before exposure to the chlorine test have indicated the permanence of film uniformity. Furthermore, oxidized devices given the 900°C chlorine test indicated no change in their electrical characteristics.

IV. FIELD EFFECT MEASUREMENTS ON THERMALLY OXIDIZED SILICON:

Both the large-signal ac field effect technique, using a frequency of 40 to 500 cps, and the zero-frequency four-point probe field effect technique were used. The sample resistivities ranged from 80 to several hundred ohm-cms, both for p-type and n-type silicon. Pulled crystals, rotated and nonrotated, as well as floating-zone crystals were used.

4.1 Evidence for the Absence of Slow States

Fig. 11(a) shows the ac field effect pattern obtained at 70 cps for a 300-ohm-cm p-type pulled crystal with an oxide thickness of about 500

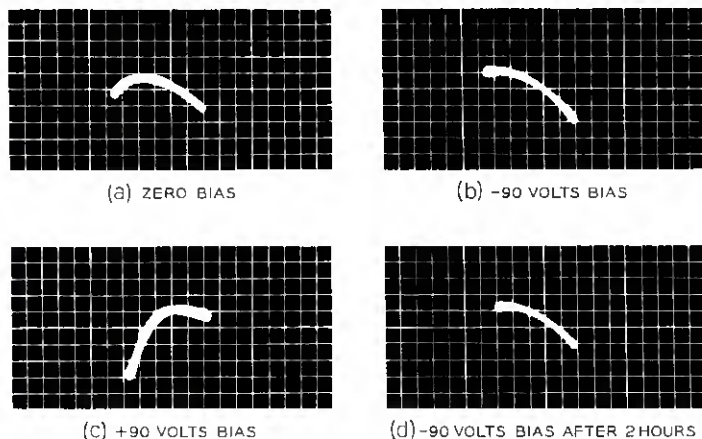


Fig. 11 — Large-signal ac field effect — effect of a superimposed dc bias.

angstroms. The field-plate signal was ± 280 volts. For an etched silicon surface, if one superimposed a dc bias on the ac signal, the pattern would shift to one side or the other, then drift back within about a minute to the initial pattern. For the oxidized silicon, however, an application of -90 dc volts produced the shift indicated in Fig. 11(b) and $+90$ volts produced the shift indicated in Fig. 11(c), all figures having the same calibrations. If one superimposes all three curves, one finds that the actual measured shift of the curves corresponds closely to the 90 volts which have been applied. *This indicates that the usual slow states observed on unoxidized surfaces have been virtually eliminated.* Furthermore, the -90 volts dc bias was left on for two hours with no change in the field effect pattern, which is shown in Fig. 11(d) at the end of two hours. To determine the possible existence of slower states, we have also maintained a dc field in a four-point probe experiment (discussed in a later section) for over 3000 hours with no measurable decay in surface conductance.

To explore the possibility of states with relaxation times corresponding to higher frequencies than the 70 cps used above, a frequency run was performed. This was done on a 300-ohm-cm p-type rotated crystal with an oxide thickness of 180 angstroms. The field effect pattern was ob-

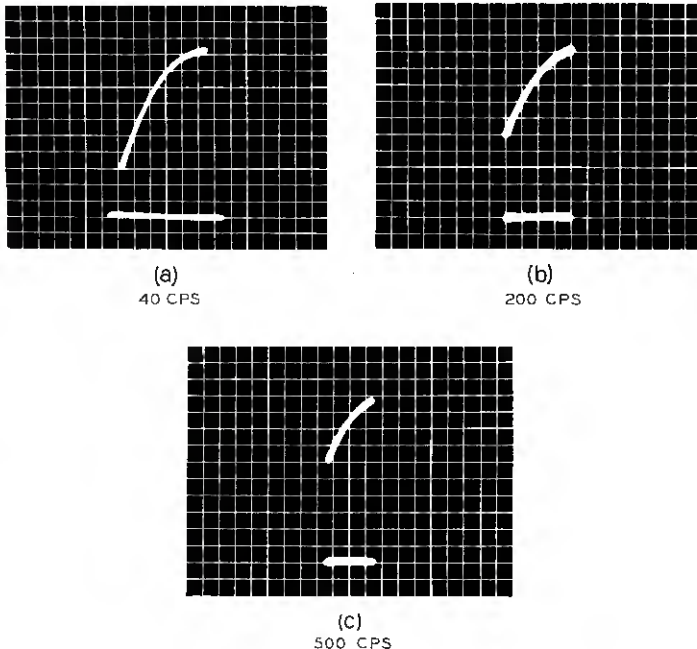


Fig. 12 — Large-signal ac field effect — effect of frequency.

served at 11 different frequencies from 40 to 500 cps. Fig. 12 shows the patterns obtained at 40, 200 and 500 cps, all drawn to the same scales. Due to distortions at the higher frequencies, lower plate voltages had to be applied at 200 and 500 cps. Here again, if one superimposes all three patterns by making the midpoints coincide, one obtains *one identical curve*. This shows that no slow states (or not-so-slow states) are observable up to 500 cps. Unfortunately, our measuring set is limited to 500 cps and, accordingly, we have no knowledge of the lower limit of the relaxation time of the surface states.

4.2 *Effect of Ambients*

The experiments reported in the above section were performed in dry oxygen. To study the effects of ambients, the following gases have been circulated: dry oxygen, dry nitrogen, wet oxygen and wet nitrogen (~ 40 per cent relative humidity), ammonia and ozone. No shifts in the field effect pattern were observed over several hours. On replacing the oxidized sample with etched silicon and germanium samples, shifts in the field effect patterns were observed in accordance with the usual experience.

4.3 *Effect of Presence of External Impurities During Oxidation — Not-So-Fast States*

This effect was first observed accidentally when the oxidation tube was cracked and impurities from the heating elements of the oven were present during the oxidation process. The general effect is to produce "*not-so-fast states*," as indicated by both the zero-frequency technique and the ac technique. The zero-frequency technique indicates a very small change in conductance with applied field. The ac field effect technique, however, by simple changes of the frequency, shows the effects of the not-so-fast states. Fig. 13 shows the results obtained from a sample

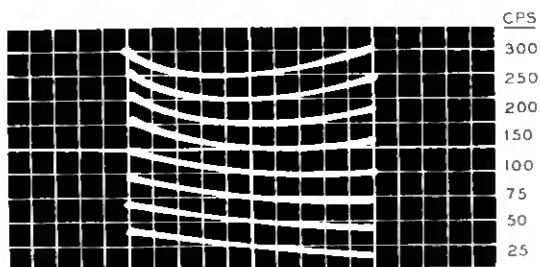


Fig. 13 — Large-signal ac field effect — effect of frequency when not-so-fast states are present.

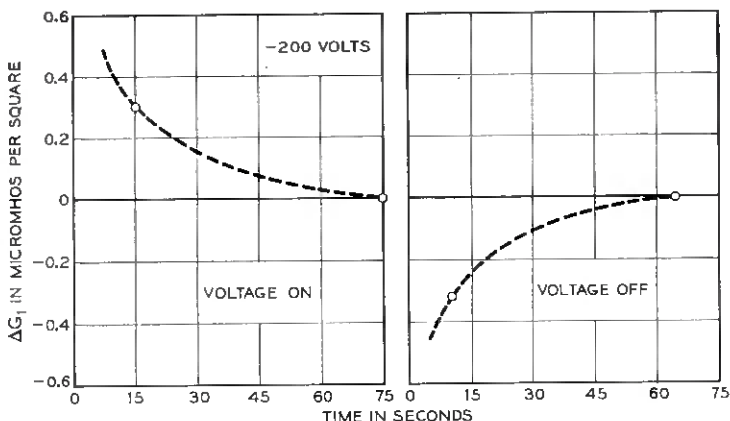


Fig. 14 — Zero-frequency four-point probe field effect — etched silicon surface.

showing this effect. From bottom to top the frequencies are 25, 50, 75, 100, 150, 200, 250 and 300 cps. As shown, the tendency is for the surface to become more p-type with increased frequency, indicating that the responsible states are mainly donor type (neutral or positive). It is of further interest to note that most of the change in the pattern occurred between 50 and 300 cps, setting a range of relaxation times for these states between 3 and 20 milliseconds.

4.4 Zero-Frequency Field Effect Measurements on Oxidized Floating-Zone Silicon

First we will demonstrate the typical behavior of etched silicon surfaces when they are examined by the zero-frequency four-point probe technique. Fig. 14 shows, for a p-type silicon specimen, a typical decay with time of surface conductance of 180 ohm-cm produced by the floating-zone technique.* The change in conductance decays to practically zero in a little over one minute. The plate voltage was -200 volts, corresponding to a plate electronic charge density of $7 \times 10^{10}/\text{cm}^2$. On removal of the plate voltage, the reverse effect is obtained.

Samples of the same crystal were oxidized to various oxide thicknesses. The field effect results were quite reproducible for each sample and from sample to sample. In all cases, no decay in the change in conductance was observed. Fig. 15 shows the results from one unit, plotted as the change in surface conductance versus plate charge. The calculated space

* These samples had 22 zone passes in a hydrogen atmosphere. It is believed that most fast-diffusing impurities, as well as oxygen, have been removed.

charge plot is also shown, from which the values of surface potential were obtained. They are shown on the experimental curve. Note that, with this technique, the *lateral location of the theoretical space-charge curve is absolute*. The horizontal segments between the theoretical and the experimental curves represent the total net charge in surface states (and also fixed charges if present — which is not the case, as will be shown shortly from data on the effect of time). The point of intersection corresponds to zero net charge in the surface states; to the right, there is a net positive charge and, to the left, a net negative charge. From this, it is concluded that *the surface states present are of two types, donor and acceptor*. By analyzing the data in the usual way, the following approximate distribution of states was obtained:

- i. An acceptor level located at about 0.4 electron volt below the intrinsic level of density of the order of $10^{11}/\text{cm}^2$.
- ii. A donor level located higher than 0.3 electron volt above the intrinsic level of density of the order of $10^{10}/\text{cm}^2$.
- iii. A nonuniform and less significant distribution near the center of the band of average density of the order of $10^9/\text{cm}^2$ volt.

For this high-purity floating-zone material, *the general tendency after*

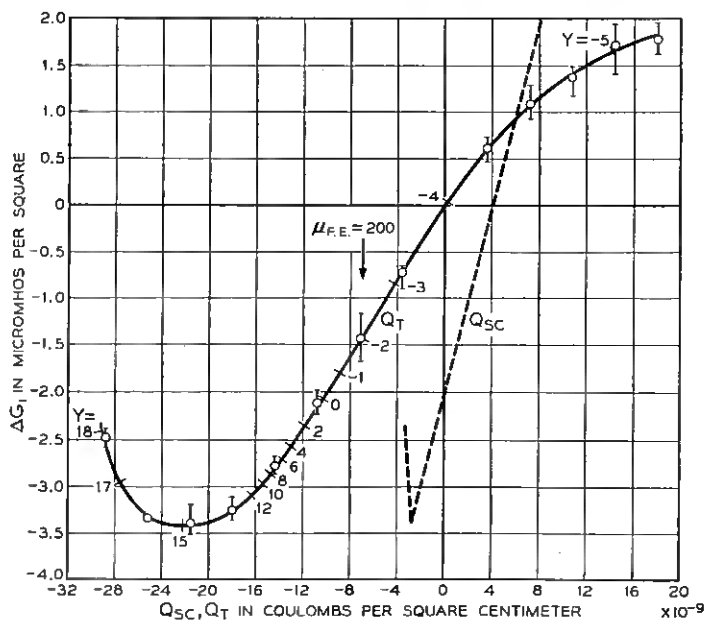


Fig. 15 — Zero-frequency four-point probe field effect — thermally oxidized high purity p-type crystal.

oxidation is to produce a p-type surface, indicating the predominance of the acceptor-type surface states. In terms of oxide type and strength, discussed previously, these results indicate a p-type oxide of strength of about 3 ohm-cm; i.e., the same surface states would invert all n-type silicon surfaces with resistivities above 3 ohm-cms.

To determine whether there is any significant number of *very slow states* within or on the surface of the oxide, a sample was given a life test in which, after its field effect characteristics were measured, a dc bias of 400 volts was left on (corresponding to a plate electronic charge of about $1.5 \times 10^{11}/\text{cm}^2$). Periodically, the surface conductance was remeasured and the field effect curve was re-obtained. The test has so far exceeded 3000 hours with no detectable change in the conductance or the field effect curve. Since the film dielectric constant and resistivity are 4 and 10^{16} ohm-cms, respectively, its relaxation time $\kappa\rho/4\pi$ (in electrostatic units) is of the order of one hour. One concludes, therefore, that *the number of outer or very slow states is insignificantly small.*

4.5 Surface Properties of Oxidized Silicon Pulled Crystals; Effect of Body Impurities

Field effect experiments have also been performed on more than 20 samples from a number of high-resistivity pulled crystals, nonrotated and rotated at different rates during growth. They generally exhibited, when oxidized, the same general behavior as far as the absence of slow states, high field effect mobility, and insensitivity to ambients. One very significant difference that was observed, however, was that *the surfaces obtained were generally n-type; i.e., the oxides were n-type and had a wide range of strength from 0.5 to 20 ohm cms.* This indicates that, in contrast to high-purity floating-zone crystals, the donor-type surface states are more predominant than the acceptor-type states. The variations in film strength were observed both from crystal to crystal and from sample to sample in the same crystal. In many cases, the surface was so strongly n-type that no conductivity minimum was observed even when a plate signal corresponding to about 2×10^{11} electron charges per cm^2 was applied, indicating a density of donor-type states in excess of this figure. To show that the n-type surface was *induced* by charges in surface states rather than by some body doping in a small region near the surface, the oxides were removed by exposure to HF vapor for 30 seconds, followed by a water rinse. Thermal probing of the surface indicated p-type material, and ac field effect observations indicated the disappearance of the initial n-type surface.

In view of the above results on the high-purity crystal, it was further

concluded that a strong n-type surface, as is obtained with pulled crystals, is *not* a characteristic behavior of thermally grown silicon oxide. To explain this behavior of pulled crystals, the main remaining possibility is crystal impurities, particularly the fast diffusants. This possibility has been supported by two main observations:

i. Pulled n-type crystals have shown, in the majority of cases, rather large increases in body resistivities during the oxidation process (as high as two orders of magnitude). These were observed when the samples were quenched from the oxidation temperature to room temperature. This change, however, will decay at room temperature, and the decay is nearly complete in the order of an hour. At 200°C the decay is complete in less than 10 minutes. This phenomenon is not understood, but it is presumably due to some body impurities producing donor-type levels,²⁹ possibly through formation of complexes with oxygen in the crystal or precipitation of acceptors.³⁰ In our work with high-purity floating-zone crystals, on the other hand, these changes in body resistivity were either absent or comparatively small, the largest observed increase being about 15 per cent. Therefore, it seems quite plausible that the source of the high density of donor-type surface states observed with pulled crystals is fast-diffusing body impurities which are getterd by the oxide.

ii. The second observation was a result of the following experiment. A 10,000-angstrom gold film was evaporated on the same high-purity floating-zone sample for which results were given in Fig. 15. The sample was then heated in an argon oven for 30 minutes at 950°C. The excess gold was then removed and the sample re-etched. From the increase in resistivity it was estimated that approximately 3×10^{14} atoms/cm³ of gold were added. This is based on published data^{31,32} on the donor and acceptor levels introduced by gold in silicon (donor level 0.35 electron volt from valence band; acceptor level 0.54 electron volt from conduction band), and on the assumption that each gold atom is electrically active. The sample was oxidized (300 angstroms) and the zero-frequency field effect measurement was repeated. The results are shown in Fig. 16. One now finds that the surface is *n-type*, the conductivity minimum is barely observed and the field effect mobility has decreased markedly to about 80 cm²/volt-sec. The oxide was then removed by HF vapor, followed by a water rinse, and the surface reverted to p-type. From this, one concludes that added body impurities such as gold (and possibly others that were unknowingly introduced into the crystal along with the gold), can diffuse to the surface and produce donor-type surface states that will induce a strong n-type surface. Finally, to determine the effectiveness of the oxide in gettering body impurities, the same sample was re-oxidized

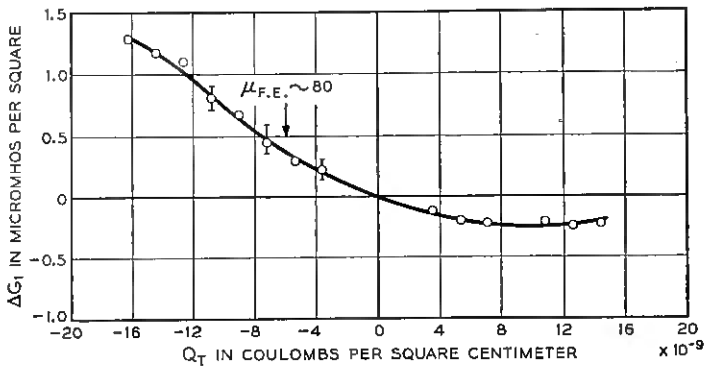


Fig. 16 — Zero-frequency four-point probe field effect — effect of adding gold to crystal.

for 65 hours at 920°C (4000 angstroms), the oxide was removed and a one-half-hour oxide was regrown (300 angstroms). Field effect measurements were made and are presented in Fig. 17. Comparing this with Figs. 15 and 16, one finds the resulting surface *still n-type*, yet the field effect mobility is slightly higher and the conductance minimum is more pronounced. From body resistivity measurements, one calculates a decrease in gold content from the initial concentration of 3×10^{14} to 1×10^{14} atoms/cm³.

4.6 Comment on the Nature of Surface States at an Oxidized Surface

It was pointed out in the introduction that surface states exist on an atomically clean surface. These are associated with the dangling bonds

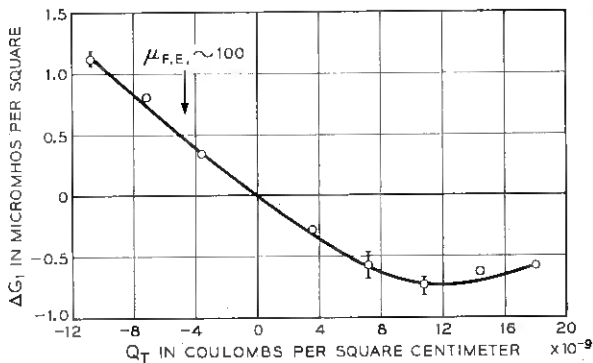


Fig. 17 — Zero-frequency four-point probe field effect — gold gettering by prolonged oxidation.

or free orbitals. These "Tamm" states are acceptor-type states which are neutral when unoccupied and negatively charged when occupied. It has further been shown, theoretically,³³ that, at the interface between two dissimilar crystals, Tamm-type states will also exist. In the case of a thermally grown oxide on a silicon surface, one may visualize a *gradual* transition in structure (within a few atomic layers) from an all-crystalline structure to an all-amorphous structure. The composition must also undergo a gradual transition from all silicon to all silicon dioxide. One intuitively expects, therefore, that such a transition should give rise to some Tamm-type surface states, as well as states associated with vacancies due to lattice mismatch. For the surface states observed in our experiments with thermally oxidized surfaces, we propose the following tentative model. The acceptor-type states observed are Tamm-type states and vacancy states, while the donor-type states observed are associated with impurities that may be introduced in a variety of ways. From this, the concepts result of surface doping and surface-state compensation, in exact analogy to their counterparts in the body of a semiconductor. For further examination of this model, it appears quite important to study the surface states on an atomically clean silicon surface which has been thermally oxidized in pure oxygen. This study is currently in progress.

V. CHARACTERISTICS OF THERMALLY OXIDIZED DEVICES

5.1 *Characteristics of Oxidized Diodes*

The oxidation process has been applied to several thousand diffused junctions. These included n^+p and p^+n diodes (graded junctions) in the range of breakdown voltages of 20 to 400 volts, $n-p-n$ and $p-n-p$ structures and $p-n-p-n$ switching transistor structures. Crystals used were all pulled crystals, with phosphorus and boron the doping impurities.

The following results were all obtained on junctions with thin oxide films, 150 to 300 angstroms, prepared at 920°C in dry oxygen (10 to 30 minutes oxidation time). After oxidation the junctions were quenched in room air. No special precautions were taken for their storage; they were usually stored in plastic boxes in room air. Units stored in this fashion for as long as 15 months have shown no change in electrical characteristics. Most of the measurements given below were obtained in room air, and contacts were made by a point pressed on the top surface of the device. By applying a sufficiently high voltage it is possible to break through the oxide film (only for the thin films) and obtain satisfactory contact for reverse bias current measurements.

Fig. 18 shows the reverse V - I characteristics for two different boron-doped crystals A and B. The junctions were obtained by diffusing phosphorus into a mechanically polished surface, giving a junction which is about 0.001 inch deep. The structure is essentially n^+p . Each slice had 12 individual diodes 0.025×0.025 inch, and slices A and B had body breakdown voltages of 36.5 and 43 volts, respectively. The circles indicate the average current for each voltage; the range is also indicated. It is seen that the spread is within a factor of 2 to 3, which has been typical for all our results within one slice and within various slices with the same history. From crystal to crystal, however, and sometimes from one diffusion run to another, the reverse currents may vary by one to two orders of magnitude. This is illustrated in Fig. 18. This figure also indicates another significant effect that has been consistently observed with n^+p

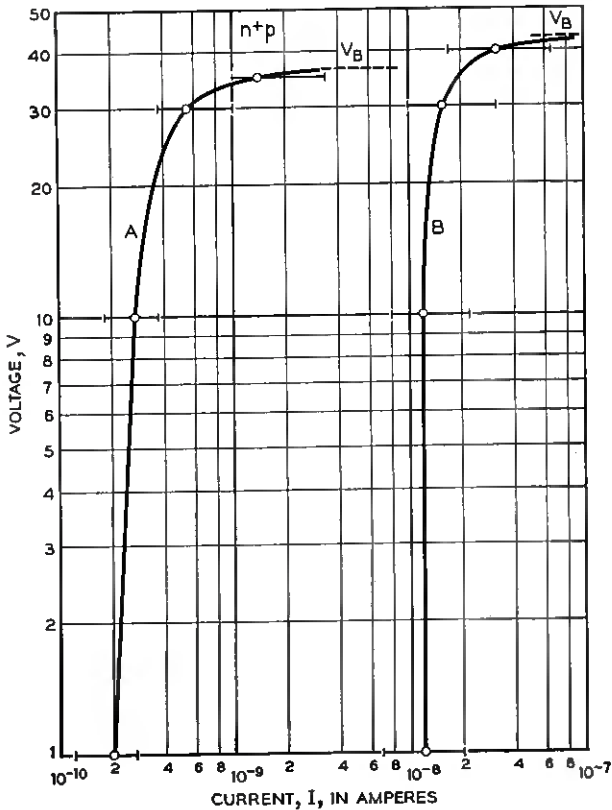


Fig. 18 — Typical characteristics of oxidized n^+p graded junction silicon diodes.

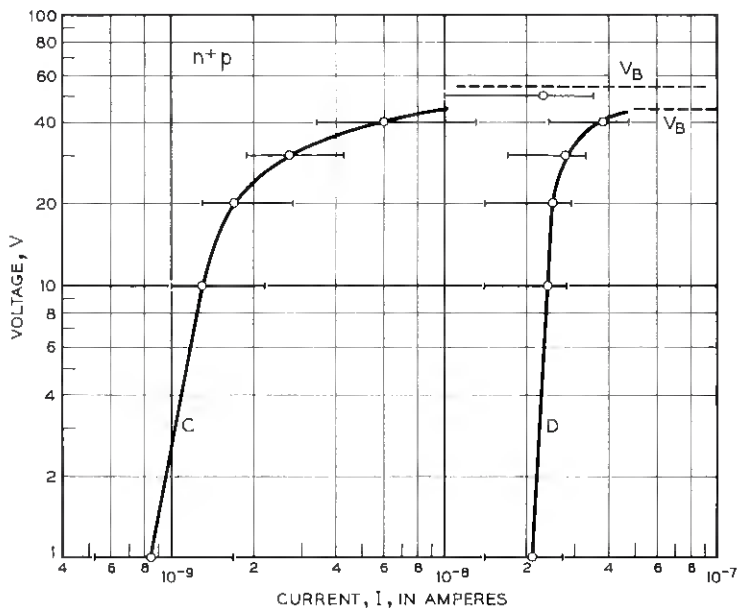


Fig. 19 — Typical characteristics of oxidized n^+p graded junction silicon diodes.

diodes. For junctions giving higher reverse currents, there is less dependence of current on voltage. In general, we have observed that for n^+p graded junctions the reverse current is proportional to the applied voltage raised to a power which varies from nearly zero to one-third. Crystal A in Fig. 18 gave a power of 0.11 to 0.18 (in the range 1 to 10 volts). Fig. 19 shows the same effects for two other sets of n^+p diodes made from two different crystals, with higher body breakdown voltage as indicated. For set c the power dependence varied from 0.14 to 0.28.

From measurements of lifetime by the injection-extraction technique,³⁴ it was shown that the body saturation current (due to diffusion according to the simple diode theory³⁵) is negligible. The contribution of body current to the observed currents must, therefore, be due to space-charge generation.^{36,37,38} For graded junctions as used above, however, the space-charge generation current is proportional to the applied voltage (plus a built-in voltage) to the one-third power. Occasionally, this dependence has been observed with oxidized units, but, in most cases, the dependence observed is weaker for n^+p junctions, as discussed above. This strongly suggests a separate contribution of the surface to the observed currents.

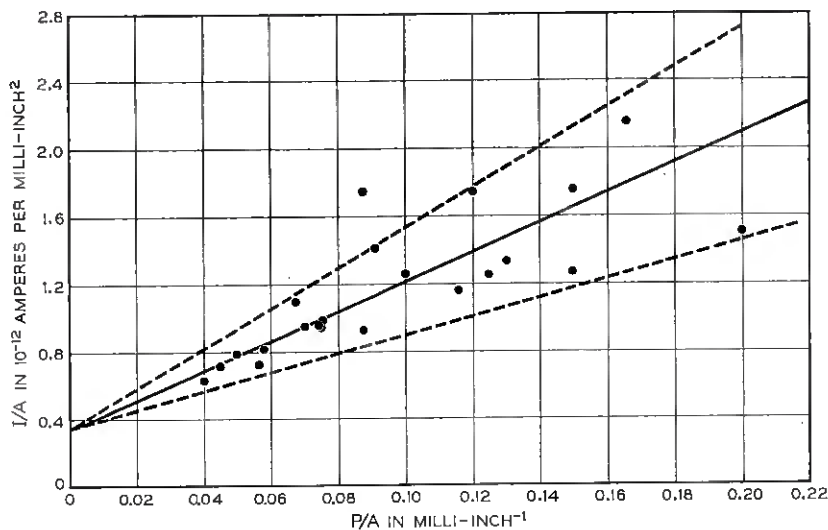


Fig. 20 — Separation of body and surface currents for oxidized diodes.

To illustrate this, the following experiment was performed. On one slice of diffused n^+p material, 22 rectangular or square units of different sizes were etched and oxidized to a thickness of about 200 angstroms. The range of perimeter to area ratio obtained was 40 to 200 inch^{-1} . This is based on the nominal dimensions rather than the actual dimensions, including perimeter irregularities, resulting from etching, and accounts for some of the spread of the data shown in Fig. 20. This figure is a plot of the I/A versus P/A , where I is the reverse current measured at 10 volts, A is the nominal area of the diode and P is its nominal perimeter. If the currents observed were only body currents, this plot would be a horizontal line, indicating that the current density is constant, regardless of shape or size of the diode. One observes, however, that the experimental points have a definite trend towards an increase in I/A with an increase in the P/A ratio. The median line is shown, together with two other lines that contain all the points except one. From the intercepts and slopes, one finds that the currents observed can be split into body and surface components. From Fig. 20, the body component is 5.3×10^{-8} amperes/cm² and the surface component is $3.5 \times 10^{-9} \pm 40$ per cent amperes/cm, both measured at 10 volts (for a 40-volt breakdown voltage n^+p diode) and at 25°C.

For further identification of the mechanisms of the body and surface currents, one must repeat the above experiments at various tempera-

tures and examine the activation energy in different temperature ranges. Preliminary results indicate that, between room temperature and 250°C, the activation energy increases with temperature, approaching approximately 1 electron volt near 160°C and *in some cases* being as low as 0.55 electron volt at room temperature. Tentatively, this suggests that reverse currents near room temperature are due to space-charge generation by traps located approximately at the middle of the energy gap, and that, at higher temperatures, the usual saturation currents due to diffusion will predominate (activation energy equal to the energy gap).

Now, returning to Figs. 18 and 19, one must explain why the dependence of current on voltage is to a power less than one-third. This may be explained in terms of our results from field effect measurements on pulled crystals, as reported in previous sections. It was pointed out that, for pulled crystals, the general tendency is to form n-type oxides, i.e., oxides with a predominance of donor-type states which may invert the surface of a p-type crystal. This means that a "channel" can form on the p-side of an n⁺-p diode. The surface component of the observed current, therefore, corresponds to this channel current. Now, if sufficiently high voltage is applied, the channel will "pinch off" and its current will essentially *saturate*. This current, added to the body current which varies as $V^{1/3}$, gives the total current, which will appear to have a weaker dependence on voltage. The stronger the channel formed, the higher its pinch-off current will be and the voltage dependence of the total current above pinch-off voltage will get weaker. For the results of Figs. 18 and 19, this suggests that the pinch-off voltages were below one volt, the lowest voltage used in these experiments. More recent work on n-p-n structures, whereby similar channels were formed by oxidation and the channel characteristics obtained both *before* and *after* pinch-off, substantiates the above explanation.

The above reasoning suggests that one should obtain a markedly different behavior for p⁺-n structures. This is, indeed, the case, and Figs. 21 and 22 show typical reverse characteristics of two different samples, E and F, of p⁺-n diodes. The junction depth is about 0.001 inch, and the body breakdown voltages are 64 and 32.5 volts for samples E and F, respectively. Sample E consists of 12 diodes and sample F consists of 7 diodes. In the range of 1 to 10 volts, the current is proportional to the voltage raised to a power of 0.4 to 0.64, which is *in excess* of the one-third power expected for a graded junction on the basis of space-charge generation. This, again, is due to surface contribution, which is explainable in the same way we have explained the behavior of n⁺-p junctions. The oxide is n-type, which tends to form an *enhancement* or *enrichment* layer

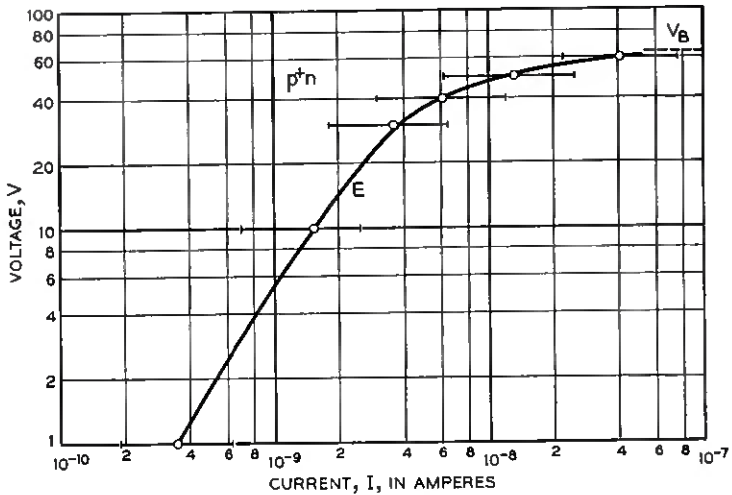


Fig. 21 — Typical characteristics of oxidized p⁺-n graded junction silicon diodes.

on the n-region of the p⁺-n junction. If this enhancement is not sufficiently strong, one obtains softer reverse characteristics, as shown in Figs. 21 and 22. It is believed that this soft characteristic is made up of a succession of small localized surface breakdowns which occur at various voltages below body breakdown. We have seen cases where, instead of a gradual softness, the reverse characteristic is made up of a succession of distinct segments starting at some low voltage and proceeding until body breakdown is reached. The slopes of these lines correspond to resistances of the order of several thousand ohms. This suggests localized surface breakdowns. These effects are now being studied for shallow junctions (of the order of a micron), using the light emission technique.^{39,40} Preliminary results indicate that the occurrence of each of these segments (or knees on the V-I curve) corresponds exactly to the onset of light emission (white light from one point on the junction surface).

5.2 Characteristics of Oxidized p-n-p-n Diodes

Another device that was oxidized was the two-terminal p-n-p-n structure transistor switch, which is an all-diffused diode with junction spacings of the order of 0.0001 inch. Fig. 23 shows typical characteristics of a set of 20 units processed on one slice, including the high impedance characteristics in both directions. (All these units had a 300-angstrom oxide.) The circles represent the average readings and the corresponding spreads are indicated.

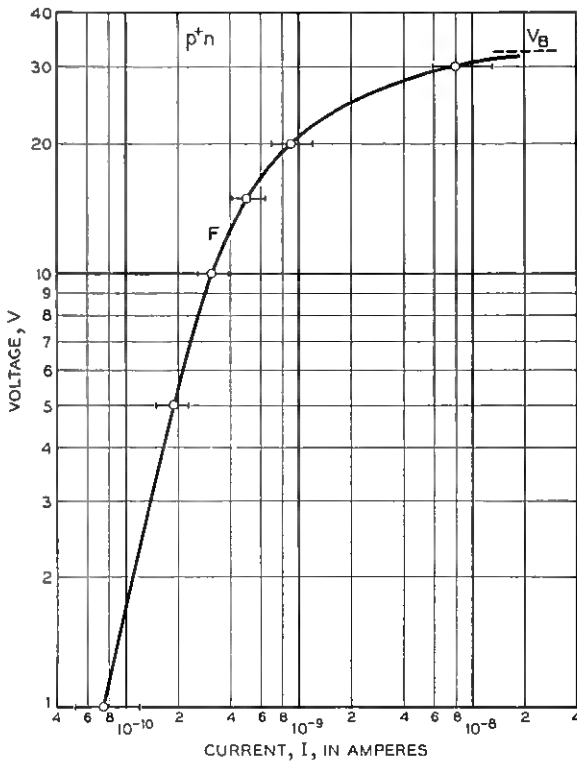


Fig. 22 — Typical characteristics of oxidized p^+-n graded junction silicon diodes.

In this device, as well as in various others, it is important to maintain the lifetime in the oxidation process. For a switching diode, the lifetime will be mainly reflected in the turn-on current of the device. This requirement is compatible with the oxidation process, provided that the cooling after oxidation is controlled. It was possible to maintain the lifetime, or even show some improvement, by cooling the oxidation oven at a rate of a few degrees per minute. This was checked on both diodes, using the injection-extraction technique for lifetime measurement, and on the p-n-p-n switch, as reflected in its turn-on current.

VI. NOISE IN OXIDIZED DIODES

Excess noise in single-crystal filaments^{41,42} and diodes⁴³ usually exhibits a $1/f^\alpha$ spectrum ($\alpha \sim 1$) and is sensitive to surface conditions. Several theories^{43,44,45,46} have been proposed to explain this frequency dependence. McWhorter,⁴³ for instance, suggests that the $1/f$ noise is caused

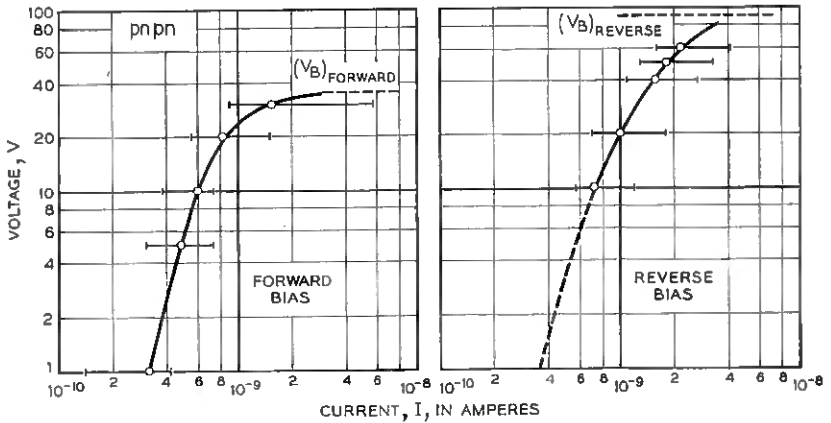


Fig. 23. Typical characteristics of oxidized p-n-p-n silicon switching diodes.

primarily by fluctuations in the occupancy of slow surface states. From field effect measurements as a function of frequency he has deduced a distribution of the relaxation times for the slow states on germanium surfaces, and he has shown that the same distribution can account for the $1/f$ dependence of excess noise. For single-crystal filaments, the fluctuating occupancy of slow states produces conductivity modulation in the bulk. In diodes and transistors, it also modulates the recombination rate at the surface, because of the relation between surface recombination velocity and surface potential.⁴⁷

As discussed previously, thermally oxidized silicon surfaces show no

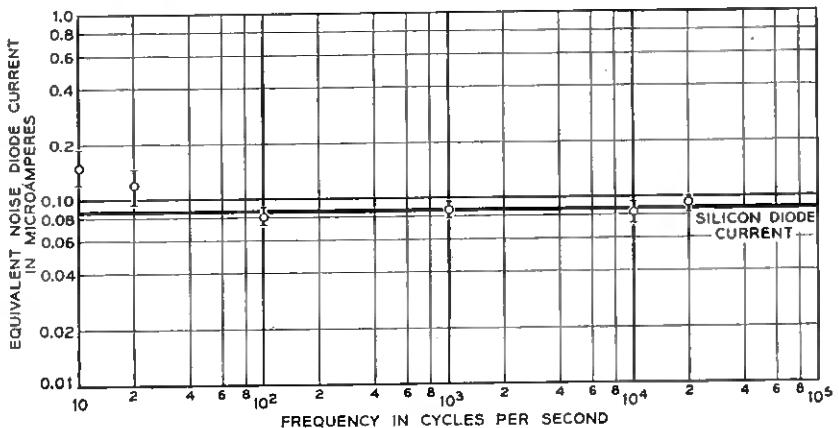


Fig. 24 — Noise in a thermally oxidized silicon diode.

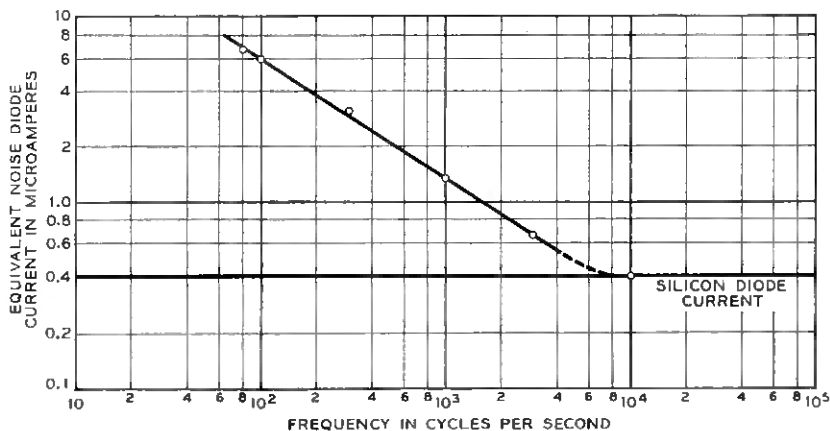


Fig. 25 — Noise in silicon diode after removal of oxide and light etching.

effects of slow states, so that it is of interest to determine whether the $1/f$ noise also is absent. Preliminary measurements of excess noise in oxidized silicon diodes have shown the following:

- i. Reverse-biased diodes operated at 75°C in dry air show no noise in excess of shot noise to frequencies as low as 30 cps (Fig. 24).
- ii. Removal of the oxide in HF vapor generally does not introduce any additional noise.
- iii. After a light etching in an HF-HNO₃ solution a $1/f^{0.66}$ dependence of excess noise was found up to a frequency of about 10 kc (Fig. 25).

VII. SUMMARY

The problem has been re-examined of the stabilization of silicon surfaces when the silicon surface is chemically bound to an appropriate solid film that is well defined in composition and structure. We have studied one such system, namely, the silicon-silicon dioxide system when the oxide film was produced by high temperature oxidation of the silicon surfaces.

i. *Film properties.* The films produced were amorphous, continuous and uniform. Various thicknesses were studied from a few hundred angstroms to tens of thousands of angstroms. The films had a resistivity of about 10¹⁶ ohm-cm, dielectric constant of about 4, and dielectric strength of about 10⁷ volts/cm. The concentration of impurities which were detected as mobile ionic charges were less than 10¹⁴/cm³.

ii. *Interface properties.* A completely new kind of surface was produced by thermal oxidation, with its own characteristic distribution of surface states. These were fast or interface states; slow or outer states were not

observed. The $1/f$ noise commonly associated with fluctuations in slow states was also not observed down to a frequency of 30 cps. The fast states consisted of both donor-type and acceptor-type states. Their densities and distributions, which determine the surface potential (surface inversion or accumulation), were affected by both the pre-oxidation treatment and the concentration of certain impurities in the silicon crystal. Furthermore, the resulting surface potential was virtually locked and showed no sensitivity to wet ambients, ammonia vapor or ozone.

iii. *Compatibility.* We have demonstrated the applicability of the oxidation process to certain device structures. We must stress, however, that it cannot be considered as a universal process which can be applied to any arbitrary device structure and maintain some desired surface properties. The oxidation process provides not only a stable surface but a completely new kind of surface with its own characteristic properties. In applying the process to devices, it cannot be considered simply as a stabilization of surface properties but must, instead, be considered as an integral part of the device structure. From this evolves the important concept of "surface design" which should be incorporated in the early stages of device design.

VIII. ACKNOWLEDGMENT

Among the many people who have contributed to this work we wish to thank J. M. Goldey, B. T. Howard and C. A. Bittmann for providing all the diffused silicon junctions; M. Tanenbaum and E. Kolb for the high-purity crystals; A. V. Voinier and G. Reich for the design and construction of the equipment and E. I. Povilonis for processing. We also wish to thank V. O. Mowery for his contributions to the field effect studies, H. Gummel for the noise measurements and Mrs. M. H. Read for the diffraction studies.

REFERENCES

1. Tamm, I., *Physik. Z. Sowjetunion*, **1**, 1932, p. 733.
2. Shockley, W., *Phys. Rev.*, **56**, 1939, p. 317.
3. Farnsworth, H. E., Schlier, R. E., George, T. H. and Burger, R. M., *J. Appl. Phys.*, **26**, 1955, p. 252.
4. Allen, F., Tech. Rep. 236, Cruft Lab., Harvard Univ., Cambridge, Mass., 1955.
5. Allen, F., Tech. Rep. 237, Cruft Lab., Harvard Univ., Cambridge, Mass., 1955.
6. Dillon, J. A., *Bull. Amer. Phys. Soc.*, **1**, 1956, p. 53; Dillon, J. A. and Farnsworth, H. E., *Phys. Rev.*, **99**, 1955, p. 1643.
7. Wallis, G. and Wang, S., *Bull. Amer. Phys. Soc.*, **1**, 1956, p. 52.
8. Madden, H. H. and Farnsworth, H. E., *Bull. Amer. Phys. Soc.*, **1**, 1956, p. 53.
9. Autler, S. H., McWhorter, A. L. and Gebbie, H. A., *Bull. Amer. Phys. Soc.*, **1**, 1956, p. 145.
10. Law, J. T. and Garrett, C. G. B., *J. Appl. Phys.*, **27**, 1956, p. 656.
11. Handler, P., *Bull. Amer. Phys. Soc.*, **1**, 1956, p. 144.

12. Palmer, D. R. and Davenbough, C. E., *Bull. Amer. Phys. Soc.*, **3**, 1958, p. 138.
13. D'Asaro, L. A., *J. Appl. Phys.*, **29**, 1958, p. 33.
14. Allen, F., private communication.
15. Hagstrum, H. D., *Bull. Amer. Phys. Soc.*, **2**, 1957, p. 133.
16. Dillon, J. A., *Bull. Amer. Phys. Soc.*, **3**, 1958, p. 31.
17. Bardeen, J., *Phys. Rev.*, **71**, 1947, p. 717.
18. Brown, W. L., *Phys. Rev.*, **100**, 1955, p. 590.
19. Montgomery, H. C. and Brown, W. L., *Phys. Rev.*, **103**, 1956, p. 865.
20. Garrett, C. G. B. and Brattain, W. H., *B.S.T.J.*, **35**, 1956, p. 1041.
21. Brown, W. L., *Phys. Rev.*, **91**, 1953, p. 518.
22. Stutz, H., DeMars, G. A., Davis, L., Jr., Adams, A., Jr., *Phys. Rev.*, **101**, 1956, p. 1272.
23. Garrett, C. G. B. and Brattain, W. H., *Phys. Rev.*, **99**, 1955, p. 376.
24. Kingston, R. H. and Neustadter, S. F., *J. Appl. Phys.*, **26**, 1955, p. 718.
25. Schrieffer, J. R., *Phys. Rev.*, **97**, 1955, p. 641.
26. Wolff, P., private communication.
27. Brattain, W. H., and Bardeen, J., *B.S.T.J.*, **32**, 1953, p. 1.
28. Ligenza, J. R., private communication (to be published).
29. Fuller, C. S. and Logan, R. A., *J. Appl. Phys.*, **28**, 1957, p. 1427.
30. Fuller, C. S., private communication.
31. Taft, E. A. and Horn, F. H., *Phys. Rev.*, **93**, 1954, p. 64.
32. Collins, E. D., Carlson, R. O. and Gallagher, C. J., *Phys. Rev.*, **106**, 1957, p. 1168.
33. Pratt, G. W., Jr., *Lincoln Lab. Quart. Prog. Rep.*, **12**, February 1, 1955.
34. Lax, B. and Neustadter, S., *J. Appl. Phys.*, **26**, 1954, p. 1148.
35. Shockley, W., *Electrons and Holes in Semiconductors*, D. Van Nostrand and Co., New York, 1950.
36. Shockley, W. and Read, W. T., Jr., *Phys. Rev.*, **87**, 1952, p. 835.
37. Pell, E. M., *J. Appl. Phys.*, **26**, 1955, p. 659.
38. Sah, C. T., Noyce, R. N. and Shockley, W., *Proc. I.R.E.*, **45**, 1957, p. 1228.
39. Chynoweth, A. G. and McKay, K. G., *Phys. Rev.*, **102**, 1956, p. 369.
40. Newman, R., *Phys. Rev.*, **100**, 1955, p. 700.
41. Montgomery, H. C., *B.S.T.J.*, **31**, 1952, p. 950.
42. Maple, T. G., Bess, L. and Gebbie, H. A., *J. Appl. Phys.*, **26**, 1955, p. 490.
43. McWhorter, A. L., *Semiconductor Surface Physics*, Univ. of Pennsylvania Press, Philadelphia, 1957.
44. North, D. O., *Bull. Am. Phys. Soc.*, **2**, 1957, p. 319.
45. Bess, L., *Phys. Rev.*, **91**, 1953, p. 1569.
46. Bess, L., *Phys. Rev.*, **103**, 1956, p. 72.
47. Stevenson, D. T. and Keyes, R. J., *Physica*, **20**, 1954, p. 1041.

Analysis and Design of a Transistor Blocking Oscillator Including Inherent Nonlinearities

By J. A. NARUD and M. R. AARON

(Manuscript received July 17, 1958)

An analysis of transistor blocking oscillators is presented which differs significantly from previous approaches in that the nonlinear dependence of collector current and base voltage upon base current is considered. First, the nonlinear differential equations governing circuit performance are derived considering the effects of alpha cutoff, collector capacitance and leakage and magnetizing inductance. The relative importance of the various terms in this equation during transition, relaxation and recovery intervals is discussed from a physical viewpoint.

Analog and digital computer solutions of the nonlinear differential equation yield pulse responses which give excellent agreement with experimental results for various signals, transistor characteristics and values of passive circuit elements. The paper concludes with a design example in which the analysis is confirmed by experiment.

I. INTRODUCTION

Transistor blocking oscillators are finding application in digital computers¹ and transmission systems² where it is necessary to reconstruct a pulse train that has been dispersed in a noisy medium. In fact, several papers^{3,4,5} have been published concerning the design and analysis of such circuits. All of the previous analytical approaches have been based on the use of a piecewise linear analysis. Linear circuit approximants have been used which are approximately valid during portions of the circuit operation, and formulas have been derived for the rise time and pulse width of the output pulse. The piecewise linear approach does not deal with the inherent nonlinear character of the circuit, and the approximations involved in linearization are often difficult to justify. The most serious shortcoming of the linearization technique is the fact that

it does not give information concerning trigger requirements for reliable operation.

In view of the inherent restrictions of the approximate analysis, a more detailed analysis which deals directly with the nonlinear nature of the circuit has been undertaken. The following sections of this paper will be concerned with the development of this procedure. First, the equivalent circuit, including the nonlinear dependencies, will be displayed and the assumptions on which the analysis is based will be discussed. The nonlinear integro-differential equations governing the performance of the circuit throughout its *entire cycle* of operation will be derived. A clear qualitative picture of the circuit operation is achieved with the nonlinearities included. Physical reasoning permits the dissection of the nonlinear equations into equations governing the circuit operation during the transition intervals ("off and on"), "on" interval and "recovery" interval. Trigger requirements for reliable regenerative action are derived. These are, of course, dependent on the nonlinear characteristics of the transistor as well as on other circuit parameters. Requirements on the circuit parameters to minimize ringing, achieve fast rise time, and recover quickly for repetitive operation are developed. Application of the results of the analysis to a specific example is given and confirmed by experiment.

Both analog and digital computers have been used as adjuncts to the strictly analytic approach to search out relationships between circuit parameters to achieve desired performance. It is believed that the combined analytical and computer approach used in this paper is more general and more easily applied than the piecewise linear analysis and, furthermore, that it is a valuable design philosophy for exploring other nonlinear circuits. Further work on other nonlinear transistor circuits using this approach is being pursued.

II. THE EQUIVALENT CIRCUIT

In selecting the equivalent circuit for an active device one is usually confronted with two conflicting requirements. On the one hand, the equivalent circuit must accurately exhibit the same properties as the physical device, so that computations based on the equivalent circuit give substantial agreement with experiment. On the other hand, there is a twofold requirement for simplicity. First, the number of parameters involved should be minimized to make the analysis tractable. Second, it should be possible to draw qualitative conclusions regarding the circuit operations from the equivalent circuit. In the case of the transistor blocking oscillator, these requirements are particularly important. Since the oscillator is highly nonlinear in its operation the number of coupled

equations involved is closely related to the complexity of the equivalent circuit.

With these factors in mind, the equivalent circuit shown in Fig. 1(a) was chosen to represent the transistor in the grounded emitter blocking oscillator. It takes into account the following properties of the transistor: (a) the nonlinear relationship between base voltage and base current; (b) the nonlinear dependence of collector current on base current and frequency; (c) partially, the effect of minority carrier storage during "turnoff" and (d) the combined effect of collector and stray capacitance.

However, the circuit neglects, among other things, the effect of coupling between the base and the collector circuit, emitter junction capacitance, base spread resistance and recombination. In the following we will discuss the quantities listed above, and how they can be obtained from measurements on the transistor. The surface barrier transistor 2N128 will be used throughout this paper to exemplify the procedure; the results are applicable to other types of junction transistors with broadly similar characteristics. Modifications in the analysis to accommodate drastic differences from this model can be made, with an attendant increase in the complexity of the analysis.

2.1 *Nonlinear Relationship Between Base Voltage and Base Current*

In the equivalent circuit of Fig. 1(a) the nonlinear relationship between base voltage and current is denoted by;

$$v_b = V_b(i_b). \quad (1)$$

The dependence on collector current has been neglected since, as shown later, the terms involving v_b in the equations governing the blocking oscillator are negligible during the "turn-on" and "on" intervals. Thus, v_b may be obtained by measuring the base voltage as a function of base current, with the collector circuit connected to the load and bias with which it will normally operate. The result of such a measurement is shown in Fig. 1(b). It is seen that this result is almost a straight line above a certain value of i_b (saturation) and almost vertically straight when $i_b < 0$. Thus, it can sometimes be approximated by two straight-line segments through the origin with different slopes.⁵ Another way of approximating it is to use the relationship between current and voltage of p-n junction diodes:⁶

$$I = I_s(e^{qV_b/kT} - 1), \quad (2)$$

where $kT/q = 0.026$ volt at room temperature (25°C) and I_s is the reverse saturation current.

2.2 Nonlinear Dependence of Collector Current on Base Current and Frequency

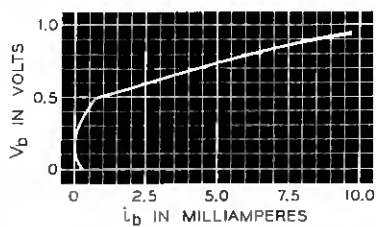
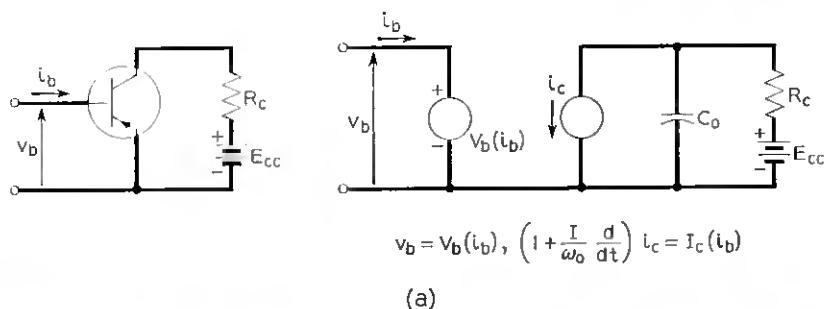
The dependence of the collector current upon base current and frequency is assumed in the equivalent circuit to be of the form

$$\left(1 + \frac{1}{\omega_0} \frac{d}{dt}\right) i_c = I_c(i_b), \quad (3)$$

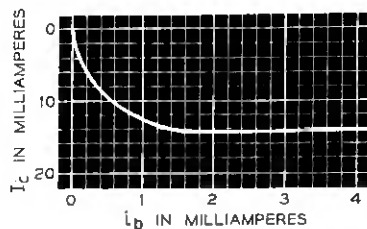
where $I_c(i_b)$ represents the static collector-current-base-current characteristic, i_c the actual collector current, and ω_0 the "large-signal" cutoff frequency or inverse time constant of the transistor. This expression is the "large-signal" equivalent of the linear relationship

$$\beta(s) = \frac{\beta_0}{1 + \frac{s}{\omega_\alpha(1 - \alpha_0)}} = \frac{\frac{dI_c}{di_b}}{1 + \frac{s}{\omega_\alpha(1 - \alpha_0)}}. \quad (4)$$

The derivative of (3) with respect to i_b reduces to (4) when the circuit operation is linear. The function $I_c(i_b)$ may be determined experi-



(b)



(c)

Fig. 1 — (a) The grounded emitter configuration and its large-signal equivalent circuit; (b) v_b vs. i_b ; (c) i_c vs. i_b .

TABLE I

2N128 Unit No.	$I_{b\beta}$, ma	$I_{c\beta}$, ma	$\bar{\beta}$	f_{α} , mc	f_0 , mc	Rise Time with 100-ohm Load, μs		Fall Time, Meas- ured, μs	f_{01} Fall, mc	$C_c(V_{cc})$, μμf	T_{R} observed/ f_{R} computed, 400-ohm Load
						Cal- culated	Meas- ured				
92	2.25	14	6.22	106.5	14.75	23.8	21	145	2.41	2.3	0.85
150	2.0	14	7.0	66.5	8.32	42.1	46	240	1.46	1.6	1.2
151	1.0	14	14.0	120	8.0	43.8	42	190	1.84	1.6	1.06
176	1.88	14	7.45	63.3	7.5	47	53	260	1.35	2.1	1.15
273	1.25	14	11.2	99	8.1	43.2	39	140	2.50	3.0	0.94
369	1.63	14	8.59	36	3.75	93.4	88	260	1.35	3.0	0.87
Average Values						48.8	48.2				1.01

mentally by measuring the collector direct current for various values of base current, with the collector connected to the load and bias with which it will normally operate. Such a characteristic is shown in Fig. 1(c). It is seen that $I_c(i_b)$ saturates very sharply, the saturation value being approximately E_{cc}/R_c . In some cases in the subsequent analysis it will be convenient to approximate this function with an analytic expression.

From (4) it is seen that, when the circuit operation is linear, the effective cutoff frequency, ω_0 , for the grounded emitter configuration is

$$\omega_0 = \omega_{\alpha}(1 - a) = \frac{\omega_{\alpha}}{1 + \beta}, \quad (5)$$

and therefore is dependent upon ω_{α} and β . Thus, as the operating point is changed, ω_0 will vary, having a minimum when β (or α) is a maximum and becoming equal to ω_{α} when $\beta = 0$ as, for example, at cutoff or saturation. When a large signal is applied to the transistor, ω_0 will, therefore, change continuously as the various points of the i_c/i_b characteristics are traversed and the "effective cutoff frequency" will then be some kind of average of (5). Although ω_{α} is also dependent on the location of the operating point,⁷ it has been found that the ω_{α} determined at standard bias (V_{cc}) yields transient results which agree closely with experiment. These results are shown in Table I and discussed further below.

Moll⁸ has suggested the following calculation of the large-signal cutoff frequency, ω_0 : If an average β is defined as the ratio of the total excursion of the collector and base currents; that is,

$$\bar{\beta} = \frac{I_{c \max} - I_{c \min}}{I_{b \max} - I_{b \min}}, \quad (6)$$

then the large-signal cutoff frequency is given by:

$$\omega_0 \doteq \frac{\omega_\alpha}{1 + \beta}. \quad (7)$$

Accordingly, in a circuit like the blocking oscillator, in which the operation extends from cutoff to saturation, ω_0 will be given by (neglecting I_{co})

$$\omega_0 \doteq \frac{\omega_\alpha}{1 + \frac{I_{cs}}{I_{bs} - I_{bc}}}. \quad (8)$$

The large-signal cutoff frequency may also be determined experimentally by applying a current step function of magnitude I_{bs} to the base and then measuring the 10 to 90 per cent rise time of the voltage across the load resistance. Since the output response can be described by a single exponential, provided the load resistance, R_c , is selected to be so small that the effect of the collector capacitance may be neglected, ω_0 is given by:

$$\omega_0 \doteq \frac{2.2}{T_R}, \quad (9)$$

where T_R is the 10 to 90 per cent rise time. In Fig. 2(a) the output response is shown for steps of various magnitudes in base current. It is seen that the rise time continually decreases as the size of the step increases, and that the response is very nearly exponential in shape, thereby justifying the above approximation.

With two ways of calculating ω_0 , the natural question that arises is how well these two methods agree. To provide a partial answer to this question, the calculated and measured rise times of six transistors were compared. The rise time was computed by use of (8) and (9) and measured across a 100-ohm load when a current step function of magnitude I_{bs} was applied to the base. The results are shown in Table I, from which it is seen that the calculated and measured values agree within ± 10 per cent and the average values agree within 2 per cent. Taking into account the error involved in measuring rise times from the face of an oscilloscope, the agreement is indeed quite good.

2.3 Effect of Minority Carrier Storage During Turnoff

The effect of minority carrier storage in the collector-base region is not very easily related to any of the commonly measured transistor param-

ters. The most convenient way to take it into account is to determine its effective time-constant experimentally. Fig. 2(b) shows the voltage across a 400-ohm load resistor in series with the collector of a junction surface barrier transistor (2N128) when various square current pulses are applied to the base. The rise time, fall time and width of the applied current pulse were 1, 7 and 100 milimicroseconds, respectively, and its magnitude was varied in steps from 0.2 to 10 ma. The base saturation current for that particular transistor (Unit 151) was 1 ma. From the

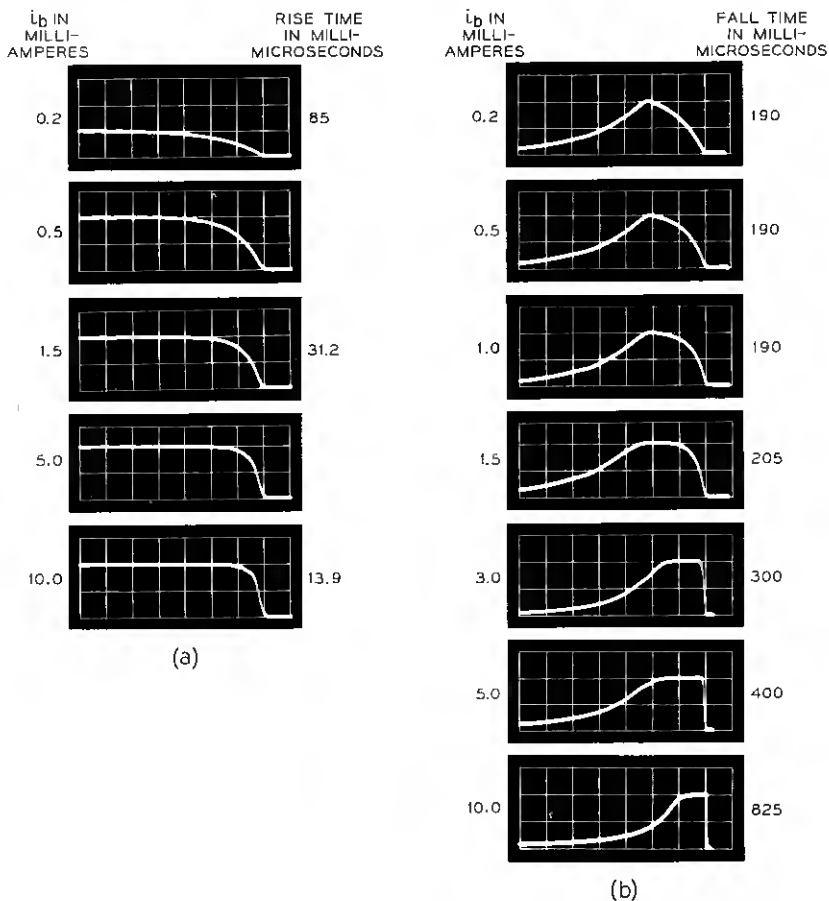


Fig. 2 — Collector current in response to (a) a step of base current, i_b ; (b) a pulse in base current of magnitude i_b , with rise time of 1 μ sec, fall time of 7 μ sec and width of 100 μ sec. Time scales: (a) 25 μ sec/division; (b) first four curves, 50 μ sec/division; next two curves, 100 μ sec/division; bottom curve, 250 μ sec/division.

figures it is seen that, in contrast to the rising edge of the output pulse, the falling part becomes slower as the magnitude of the input pulse is increased. Also, the 10 to 90 per cent fall time is almost constant up to saturation, after which it increases rapidly. Again, since the falling response has roughly an exponential shape the fall time can be used to estimate the equivalent ω_0 during turnoff by the relationship:

$$\omega_{01} \sim \frac{2.2}{T_f}. \quad (10)$$

Columns 8 and 9 of Table I list the fall times and the corresponding ω_{01} 's for the six transistors.

2.4 Combined Effect of Collector and Stray Capacitance

In the small-signal case the output impedance of the grounded emitter configuration is:

$$Z_0 = Z_c(1 - \alpha) + r_e \frac{r_b + \alpha Z_c + Z_g}{r_b + r_e + Z_g}, \quad (11)$$

where Z_g is the impedance of the generator driving it. Since, in the blocking oscillator, Z_g is large and r_e is very small, the output impedance consists essentially of a capacitance of magnitude

$$C_0 = \frac{C_c}{1 - \alpha} = C_c(1 + \beta). \quad (12)$$

Under large-signal operation, both β and C_c will vary over the range, the collector capacitance varying with the collector voltage as

$$C_c = kV^{-r}, \quad (13)$$

where the exponent r is $\frac{1}{2}$ for step junctions and $\frac{1}{3}$ for uniformly graded junctions. Measurements made on a dozen 2N128 transistors yield an average value for r very close to $\frac{1}{3}$.

Bashkow⁹ has shown that, as far as transient analysis is concerned, this nonlinear capacitance may be replaced by an average capacity C_{av} . The average capacity is defined as that which displaces the same charge as the nonlinear capacity for the same voltage change. A minor generalization of Bashkow's results yields

$$C_{av} = \frac{C_c(V_{cc})}{1 - r} \frac{1 - \left(\frac{V_1}{V_{cc}}\right)^{1-r}}{1 - \frac{V_1}{V_{cc}}}, \quad (14)$$

where V_{cc} is the collector voltage at cutoff and V_1/V_{cc} is the ratio of the

voltage swing on the collector. For a swing from V_{cc} to 0, and $r = \frac{1}{3}$, (12) becomes

$$C_0 = 1.5C_c(1 + \bar{\beta}), \quad (15)$$

where, as in the case of ω_0 , β has been replaced by $\bar{\beta}$. In a physical circuit, stray capacity from collector to ground should be added to C_0 .

With the transistor characterized in this manner it is readily shown that the Laplace transform of the collector current of the resistively loaded (R_c) common emitter stage is given by

$$I_c(s) = \frac{I_{cs}}{s \left(\frac{s}{\omega_0} + 1 \right) (R_c C_0 s + 1)} \quad (16)$$

for a step function in base current sufficient to saturate the transistor. To check the validity of the equivalent circuit, the 10 to 90 per cent rise time was computed for six transistors from the inverse transform of (16), using (7) and (15).^{*} The rise time was also determined experimentally and the ratio of observed rise time to computed rise time is given in the last column of Table I. It can be seen that agreement within ± 15 per cent has been achieved on individual units, and that the average is within 2 per cent. These figures are within both the experimental error in reading a scope face and the error inherent in determining C_c , ω_α and $\bar{\beta}$.

In passing, it should be noted that a good approximation to the 10 to 90 per cent rise time is given by

$$T_R = \frac{2.2}{\omega_0} \sqrt{1 + \omega_0^2 R_c^2 C_0^2}. \quad (17)$$

The above expression results from (16) by an application of the rise time definition given by Elmore.¹⁰ This will be mentioned again in the section devoted to the rise time of the blocking oscillator. It can also be shown by using Elmore's expression that the delay time (time to traverse from zero to 50 per cent of the final output) is given by

$$T_0 = \frac{1 + \omega_0 R_c C_0}{\omega_0}. \quad (18)$$

Over a wide range of parameters, two times the time delay, as given above, is a very close approximation to the actual rise time. In this case, the rise time [two times (18)] is similar in form to that given by Easley⁷ and can be justified by other considerations.¹¹† Equation (18) will be

^{*} 5 micromicrofarads were added to C_0 to account for stray socket and wiring capacity.

† Several definitions of rise time are given in Ref. 11, one of which, when applied to the grounded emitter stage, yields (18).

useful in the following sections in indicating how the factors in the design of the transistor blocking oscillator can be related to the rise time or delay time of the resistively loaded transistor alone.

Before proceeding to the blocking oscillator, it should be mentioned that the equivalent circuit of Fig. 1(a) may be used for the grounded base configuration. In this case, $\omega_0 = \omega_\alpha$ and $C_0 = 1.5C_c$. Therefore, the following discussion is also applicable to the grounded-base case.

III. THE GROUNDED EMITTER TRANSISTOR BLOCKING OSCILLATOR

Fig. 3 shows the circuit diagram as well as the equivalent circuit for the grounded emitter blocking oscillator with series feedback. The primary of the transformer is connected in series with the collector and the

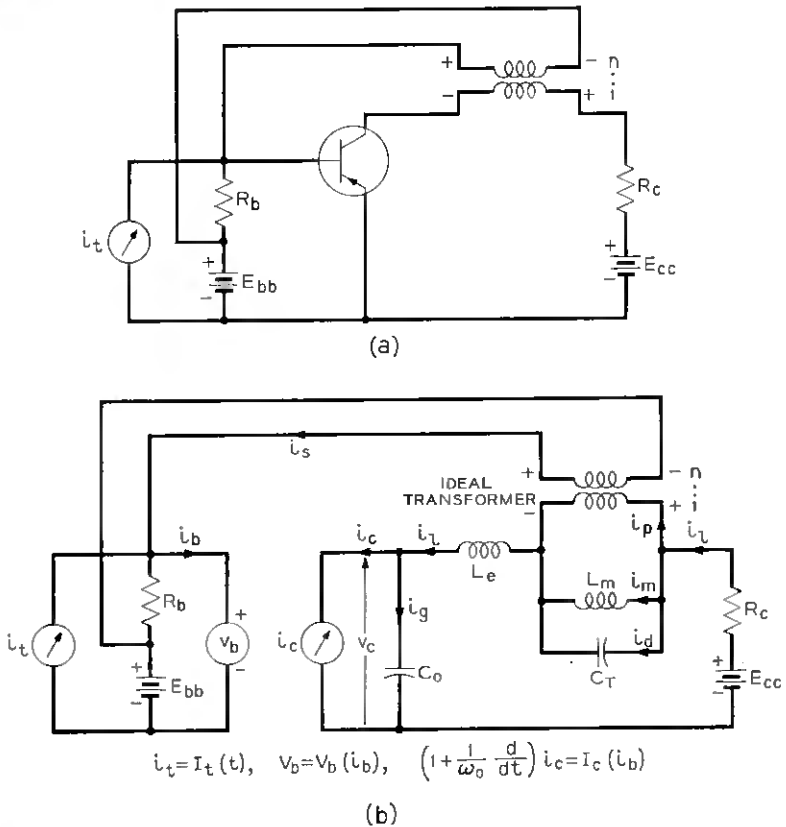


Fig. 3 — (a) Grounded emitter transistor blocking oscillator; (b) nonlinear equivalent circuit.

load resistance, R_c , and its secondary is connected between the base and the bias source E_{bb} . The polarity of the coupling is such that positive feedback exists around the loop when the transistor is conducting. The current generator i_t represents the trigger source and R_b the total external load in the base circuit.* As discussed in the previous section, ω_0 represents the large-signal cutoff frequency of the transistor and is related to ω_α through (8). The capacitance C_0 , is the sum of the total stray capacity of the circuit and that of the transistor given by (15). Finally, L_e , L_m and C_T represent, respectively, the leakage and magnetizing inductances and the parasitic capacity of the transformer.

Briefly, the circuit operation is as follows. For monostable operation the bias E_{bb} is selected in such a way that the transistor is cut off at its stable operating point. If a trigger signal of sufficient magnitude and width to bring the circuit into its regenerative region is applied, the collector and base currents will build up until a point in the saturation region is reached. (This point will hereafter be referred to as the "quasi-stable operating point.") Then the magnetizing inductance becomes charged slowly, shorting out more and more of the current fed back to the base. This process will go on until the gain around the loop again becomes equal to unity, at which point the circuit regenerates back into a state of cutoff. Then the energy stored in the transformer slowly dissipates and the circuit recovers to its stable operating point.

The equation governing this complete cycle of operation is derived in Appendix A and normalized in Appendix B.

One point is immediately obvious from these equations: they are extremely long, complicated and relatively useless as they stand. In order to restore the physical feel for the circuit, it will be shown in the remainder of the paper that the equations can be simplified, depending on the region of operation. The primary basis for reducing the more exact equations will be whether or not the base emitter junction is forward- or back-biased. In principle, this approach is similar to that used in piecewise linear analysis. There are several important advantages, however, in starting with the more exact description of the circuit given here. First, retention of the nonlinear terms enables one to determine the operating points of the circuit, which, in turn, give a clear definition of the trigger requirements for reliable operation. Second, the approximations involved in taking the general equations apart can be easily drawn. Finally, the simulation of the blocking oscillator from the exact equations on an analog or digital computer can be readily achieved for purposes of com-

* For convenience, R_b has been placed between base and E_{bb} rather than from base to ground. This approximation serves to eliminate a few small constant terms from the governing equation in Appendix A, and in no way restricts the analysis.

plete analysis or evaluation of the approximations that are made in tearing the more complete description apart. The computer approach also permits the investigation of the effects of the various parameters of the circuit on its response under conditions that can be more closely controlled than can be those in the experimental circuit.

The equations can be broken down in accordance with the importance of the various terms during the following four intervals: (a) "Transition On", (b) "On", (c) "Transition Off" and (d) "Recovery." These will be discussed separately in the following sections.*

IV. THE "TRANSITION ON" INTERVAL

This interval is defined as the time from the initiation of the "On" trigger pulse to the time at which the circuit has apparently come to rest at its quasistable operating point, that is, the time when all transients associated with switching have ceased. Specifically,

$$I_{bc} \approx i_b \approx I_{br}. \quad (19)$$

By the time triggering occurs, the input impedance to the transistor is low compared to R_b and

$$\frac{V_b(i_b) - E_{bb}}{R_b} \ll i_b, \quad (20)$$

where $[V_b(i_b) - E_{bb}]/n$ represents the voltage drop that must be maintained across the primary of the feedback transformer to overcome the bias and forward drop. This voltage is unavailable in series with the load and should be made as small as possible, implying the use of a transistor with a low forward drop. Fig. 1(b) shows that this voltage is essentially constant during the on interval; therefore very little feedback current will be diverted by the transformer capacity C_T . Subject to condition (20) and the constancy of $[V_b(i_b) - E_{bb}]/n$, this term and derivatives thereof may be neglected. During the switching interval, the contribution from the integral

$$\frac{R_b}{n^2 L_m} \int_0^t \frac{v_b - E_{bb}}{R_b} dt$$

will be negligible, since the magnetizing inductance is usually large in

* The equations in Appendix A can be readily modified to deal with the blocking oscillator with shunt feedback. For example, $R_c = 0$ and the load resistance is reflected through the transformer in parallel with R_b .

comparison with the other parameters of the circuit.* Finally, the magnetizing current is equal to nI_{bc} . Therefore, in this interval, (109) of Appendix A simplifies to

$$\begin{aligned} & \frac{L_e C_0}{\omega_0} \frac{d^3 i_b}{dt^3} + \left(L_e C_0 + \frac{R_e C_0}{\omega_0} \right) \frac{d^2 i_b}{dt^2} + \left(\frac{1}{\omega_0} + R_e C_0 \right) \frac{d i_b}{dt} - (I(i_b) - i_b + I_{bc}) \\ & = \left[1 + \left(\frac{1}{\omega_0} + R_e C_0 \right) \frac{d}{dt} + \left(L_e C_0 + \frac{R_e C_0}{\omega_0} \right) \frac{d^2}{dt^2} + \frac{L_e C_0}{\omega_0} \frac{d^3}{dt^3} \right] I_t(t). \end{aligned} \quad (21)$$

The equation governing the equilibrium points (or the "steady-state equation," as it has been called) for this interval is obtained by setting $I_t(t)$ and all the derivatives of i_b equal to zero; that is, we obtain:

$$I(i_b) - i_b + I_{bc} = 0, \quad (22)$$

where

$$I(i_b) = \frac{I_c(i_b)}{n} - \left[1 + \frac{R_b}{n^2 L_m} \left(\frac{1}{\omega_0} + R_e C_0 \right) \right] \frac{V_b(i_b) - E_{bb}}{R_b} \quad (23)$$

represents the nonlinear or active part of this equation. In this expression, the term $I_c(i_b)/n$ represents the contribution to $I(i_b)$ by the collector current and the term $[V_b(i_b) - E_{bb}]/R_b$, as noted, is the amount of current flowing down through the resistance R_b to overcome the bias. The term

$$\frac{R_b}{n^2 L_m} \left(\frac{1}{\omega_0} + R_e C_0 \right) \frac{V_b(i_b) - E_{bb}}{R_b}$$

represents the modification in $I(i_b)$ due to finite magnetizing inductance, or the equivalent drain of this inductance. In other words, the magnetizing inductance may, as far as the steady-state equation is concerned, be thought of as modifying the resistance R_b to:

$$R_b' = \frac{R_b}{1 + \frac{R_b}{n^2 L_m} \left(\frac{1}{\omega_0} + R_e C_0 \right)}. \quad (24)$$

It is instructive to study the graphical solution of the steady-state equation, since this reveals the locations of the operating points and gives a qualitative idea of the magnitude of the trigger signal needed to trigger the circuit. In Fig. 4 the various components of $I(i_b)$ and the graphical solution of the steady-state equation are shown. For the sake

* The principal effects of neglecting this term are: (a) the trigger current required for reliable triggering will be somewhat larger than indicated in the following for slow trigger signals; (b) the rise time will also increase slightly above that computed in the following, due again to diversion of feedback current through L_m .

of clarity the contribution due to the term containing $V(i_b)$ has been somewhat exaggerated.

In this figure the long-dashed curve represents the dynamic $I_c(i_b)/n$ characteristic for the transistor and the short-dashed curve the contribution due to the $[V(i_b) - E_{bb}]/R_b$ term. The solid curve constitutes the combined effect of these two curves, while the three intersections with the straight line $i_b - I_{bc}$ represent the three operating (or equilibrium) points of the blocking oscillator. Since the loop gain of the circuit

$$T = \frac{dI(i_b)}{di_b} \quad (25)$$

is obviously less than unity at the two extreme points (s and q.s.), these are operating points of the stable kind. In the following, these two points will be referred to as the "stable" and the "quasistable" operating points, respectively. The circuit cannot stay for an indefinite period of time at the q.s. point, due to the accumulation of current in the magnetizing

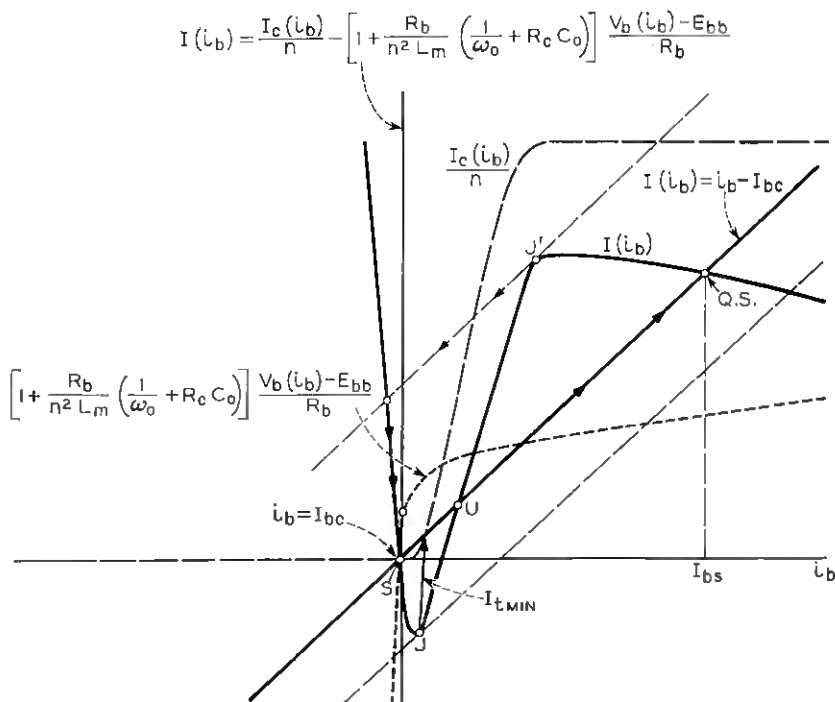


Fig. 4 — Graphical solution of "steady-state equation".

inductance. Finally, the point U is unstable (saddle-point) since the loop gain here is larger than unity, giving rise to responses in the neighborhood of this point directed away from it.

In general, the location of the operating points will depend upon I_{bc} , the shape and magnitude of $I(i_b)$ and the turns ratio n . However, since the contribution to $I(i_b)$ of the term containing $V(i_b)$ is very small when $i_b > 0$ and the backward input resistance $dV(i_b)/di_b$ is very large, the bias current I_{bc} or the bias voltage E_{bb} will not significantly affect the location of any of the operating points, provided the magnetizing inductance is relatively large. In addition, this means that the stable operating point is more or less fixed to a position just to the left of the origin on the i_b axis. Hence, only the locations of the unstable and the quasistable operating points are subject to change; such changes can only be effected by varying $I_c(i_b)$ and n . The magnitudes or saturating values, I_{cs} , of $I_c(i_b)$ and n affect the location of these two points in an opposite manner. For example, as I_{cs} increases or n decreases, the unstable point moves closer to the stable one and the quasistable operating point is displaced further out, provided the value of i_b at which $I_c(i_b)$ saturates remains unchanged. Finally, if $I_c(i_b)$ is made steeper close to the origin, this will also result in a smaller distance between the unstable and the stable operating points.

On the other hand, if the magnetizing inductance, or n , is made very small such that

$$\left[1 + \frac{R_b}{n^2 L_m} \left(\frac{1}{\omega_0} + R_c C_0 \right) \right] \frac{V_b(i_b) - E_{bb}}{R_b} \tag{26}$$

is significant with respect to $[I_c(i_b)]/n$, the term containing $V_b(i_b)$ will also affect the location of the unstable and the quasistable operating points. For instance, by reducing L_m (but letting n remain fixed) the $I(i_b)$ curve is moved downward, causing the unstable operating point to move further away and the quasistable one to move closer to the stable operating point, thereby reducing the range that i_b and $I(i_b)$ traverse in going from one operating point to the other.

The dependence upon n is more complicated, however, since it affects both the saturation value of $I(i_b)$ and $n^2 L_m$. If n is very small, $n^2 L_m$ will also be small, tending to make the distance between the stable and unstable points large. On the other hand, a small n will also increase the saturation value of $I(i_b)$, thereby tending to decrease the distance between these two points. However, since $n^2 L_m$ is proportional to the square of n , while I_s varies as $1/n$, the net effect at first will be that the unstable operating point will move closer to the stable one when n is increased

from a small initial value. Then, as n becomes larger, the effect of the $V(i_b)$ term becomes negligible and $I_c(i_b)/n$ will predominate, with the result that the distance between these two operating points will increase. Hence, as a function of n , this distance will go through a minimum. Therefore, in order to make sure that the magnetizing inductance does not affect the location of the operating points and does not become significantly charged during the "Transition On" interval, L_m and n should be selected well beyond this minimum; that is,

$$n^2 - \frac{I_c(i_b)}{V_b(i_b) - E_{bb}} n + \frac{R_b}{L_m} \left(\frac{1}{\omega_0} + R_c C_0 \right) \ll 0 \quad (27)$$

In such cases, it is an excellent approximation to neglect the effect of the term containing $V_b(i_b)$ entirely, and use only $I_c(i_b)/n$ for $I(i_b)$. The graphical determination of the operating points then simplifies to that of Fig. 5. It is seen that n is now simply equal to the ratio, $I_{cs}/(I_{bs} - I_{bc})$, which, by virtue of (8), means that ω_0 may be expressed in terms of ω_α and n only; that is,

$$\omega_0 = \frac{\omega_\alpha}{1 + n}. \quad (28)$$

In the following, the requirements for triggering the blocking oscillator, the effect of trigger magnitude, the shape of the collector current charac-

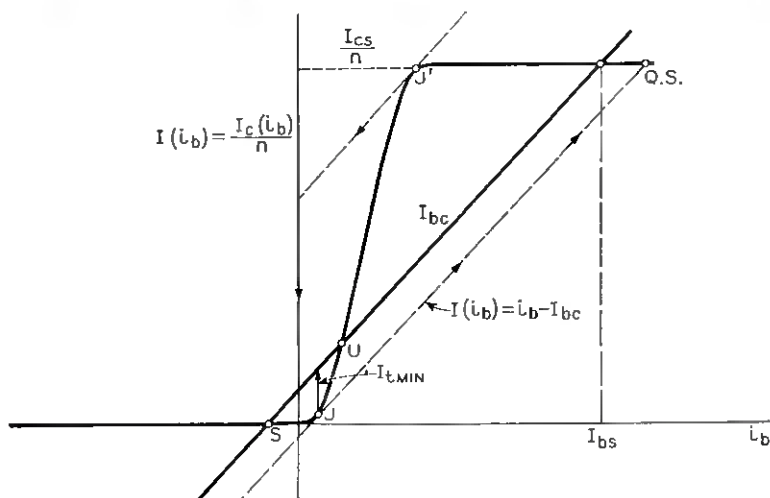


Fig. 5 — Approximate location of operating points.

teristics and the coefficients of the differential equation upon rise time and delay and the conditions for a nonoscillatory response will all be discussed. It will be assumed that condition (27) is satisfied, so that the effect of the magnetizing inductance may be neglected. Also, for most of the following, $I_c(i_b)$ will be approximated by an analytic function.

These items are most easily discussed if (21) is written in normalized form. From Appendix B, it can be shown that the normalized equation governing the "Transition On" interval is

$$\begin{aligned}
 a^2 \frac{d^3 x}{d\tau^3} + (a^2 + b) \frac{d^2 x}{d\tau^2} + (1 + b) \frac{dx}{d\tau} + x - f(x) \\
 = \left[1 + (1 + b) \frac{d}{d\tau} + (a^2 + b) \frac{d^2}{d\tau^2} + a^2 \frac{d^3}{d\tau^3} \right] g(\tau),
 \end{aligned}
 \tag{29}$$

where all symbols are defined in Appendix B. The normalized steady-state equation is now

$$f(x) - x = 0, \tag{30}$$

where both x and $f(x)$ attain the values 0 and 1 at the stable and the quasistable operating points, respectively. Finally, the normalized output voltage for this interval may be written

$$y = x - g(\tau). \tag{31}$$

V. TRIGGER SIGNAL REQUIREMENTS

If the blocking oscillator is released from an initial position that is different from any of its operating points, it will proceed, depending upon the initial conditions, either to the stable or the quasistable operating point. For instance, if the initial values of all the derivatives in (29) are equal to zero, the circuit will simply come to rest at the operating point that is on the same side of the unstable point as is the initial position. In general, however, the circuit may cross the unstable operating point, provided the initial values of some of the derivatives in (29) are large enough. Hence, in order to make the blocking oscillator flip over from the stable to the quasistable operating point, the applied trigger signal must have sufficient magnitude and duration either to bring the circuit past its unstable operating point or to impart energy into it so that the initial values of the derivatives in (29) are large enough to make the circuit traverse the unstable point by itself. In this section the conditions this requirement imposes upon the trigger signal will be discussed.

The simplest possible case is that of a trigger signal with a very large width. In such cases, the trigger pulse will vary very slowly or be nearly

constant during the transition interval, and the derivatives of $g(\tau)$ in (29) may therefore be neglected. Up to the point of triggering, the base current will also vary very slowly and its derivatives will consequently be small.* Therefore, prior to the time regenerative action takes place, (29) reduces to:

$$x - f(x) = g(\tau). \quad (32)$$

Under these conditions, the trigger signal acts mainly as a variable bias causing the load line in Fig. 4 or Fig. 5 to be displaced to the right. This, in turn, makes the stable and the unstable operating points move closer and closer together until they merge into one point. At this point (j) the load line is tangential to $f(x)$ [or $I(i_b)$] and an infinitesimal additional displacement of the load line will cause the circuit to have only one possible operating point. Hence, a "jump" to the quasistable operating point will take place at this point, and the minimum required trigger magnitude is given by:

$$G_{\min} = x_j - f(x_j), \quad (33)$$

where x_j represents the value the normalized base current has at the point of tangency, j ; that is, x_j is found by solving the equation:

$$\frac{df(x_j)}{dx_j} = 1. \quad (34)$$

In terms of actual trigger current, this minimum corresponds to the vertical distance, $I_{t \min}$, between the point of tangency and the load line, as shown in Figs. 4 or 5. The variation of $I_{t \min}$ with n will proceed in the same manner as for the distance between the stable and unstable operating points. This effect will be described more fully in the following.

When more rapidly varying trigger signals are applied to the blocking oscillator, the terms containing derivatives of $g(\tau)$ and x become so large that they must be taken into account. Since (29) is both nonlinear and inhomogeneous, no analytical solution exists, and it is therefore impossible to derive in closed form conditions that the trigger signal must satisfy for reliable operation. Hence, numerical methods or analog computer solutions must be resorted to. In the particular case, where the magnitude required to trigger the circuit for a particular fixed width of the trigger pulse must be found, the latter method was chosen, because repeated adjustments of parameters until the desired solution is obtained are much more readily accomplished on an analog computer.

* The quantitative definition of "small" is, of course, dependent on the coefficients (circuit parameters) multiplying the various derivatives.

To prevent overload of amplifiers and integrators and also to simplify the analog computer setup, (29) was first rewritten in terms of the normalized output voltage, y , which simplifies it to:

$$a^2 \frac{d^3 y}{d\tau^3} + (a^2 + b) \frac{d^2 y}{d\tau^2} + (1 + b) \frac{dy}{d\tau} + y = f[y + g(\tau)]. \quad (35)$$

The Gaussian pulse was used as a representative trigger signal; that is,

$$g(\tau) = Ge^{-[(\tau-\tau_0)/0.6\tau_w]^2}, \quad (36)$$

where G is the magnitude, τ_w the 50 per cent width and τ_0 the time at which $g(\tau)$ attains its maximum value, G . For $f(x) = f[y + g(\tau)]$, representing the normalized relationship between collector current and base current, the analytical function

$$f(x) = \frac{x^2}{2x^2 - 2x + 1} \quad (37)$$

was chosen. This choice was dictated strictly on the basis of mathematical simplicity and tractability of the solution of (35) when g , a and b are zero. This function is depicted by the solid line in Fig. 6, from which it can be seen that it has a minimum at $x = 0$ and a maximum $x = 1$, where it is equal to zero and unity, respectively. Also, $f(x)$ has an inflection point at $x = \frac{1}{2}$ around which it is symmetrical. The slope of $f(x)$ at this point is equal to 2. For values of x less than zero and larger than one it is slowly asymptotic to the horizontal line $\frac{1}{2}$. The normalized

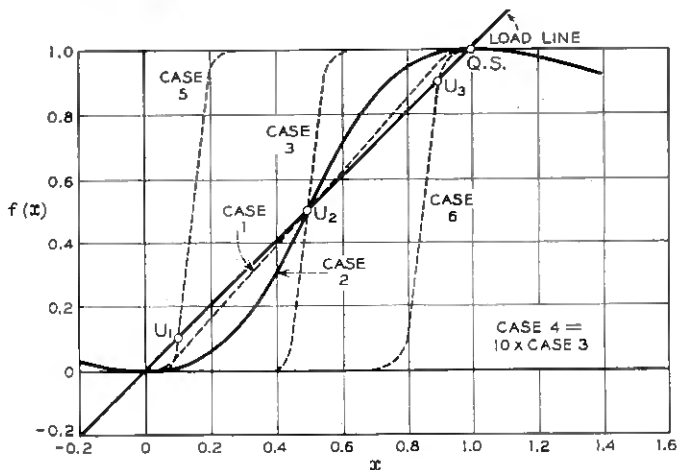


Fig. 6 — Various nonlinear characteristics.

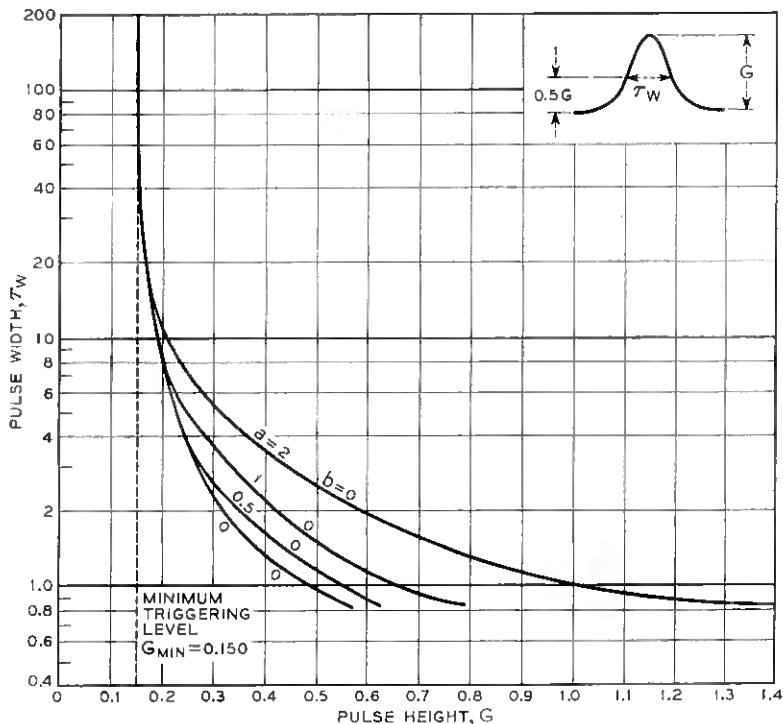


Fig. 7 — Relationships between magnitude and width of applied signal required to trigger the blocking oscillator with $b = 0$.

load line of the circuit is also shown in Fig. 6. It intersects $f(x)$ at the points x equal to 0, $\frac{1}{2}$ and 1, and these are therefore the stable, unstable and quasistable operating points, respectively. Although the function $f(x)$ differs from a real i_c/i_b characteristic in that it is not constant beyond the stable and the quasistable operating points, it resembles such a characteristic sufficiently closely within the region of interest. Hence, solutions of (35) having reasonable overshoots should be expected to approximate exact ones very closely and exhibit the main features of the circuit accurately.

On the analog computer the following procedure was used. The circuit was assumed to be at rest at the stable operating point prior to the application of the trigger pulse, and the initial value of all the derivatives was set equal to zero. Then, for each value of the width τ_w , the magnitude G was adjusted so that the response barely flipped over to the quasistable operating point. The resulting relationship between the

magnitude and width of the trigger pulse is shown in Figs. 7 and 8 for various values of the parameters a and b . It is seen that the curves are somewhat hyperbolic in shape and that all of them are asymptotic to the vertical line $G_{\min} = 0.15$ at large pulse widths. This asymptote corresponds exactly to the minimum trigger magnitude given by (33) and, indeed, if the expression for $f(x)$ from (37) is substituted into (33) and (34), the result is $G_{\min} = 0.1504$, which agrees quite closely with the computed value. From Figs. 7 and 8 it can be seen that, at large pulse widths, the critical trigger magnitude is nearly independent of a and b . As the width becomes smaller, a larger magnitude is required to trigger

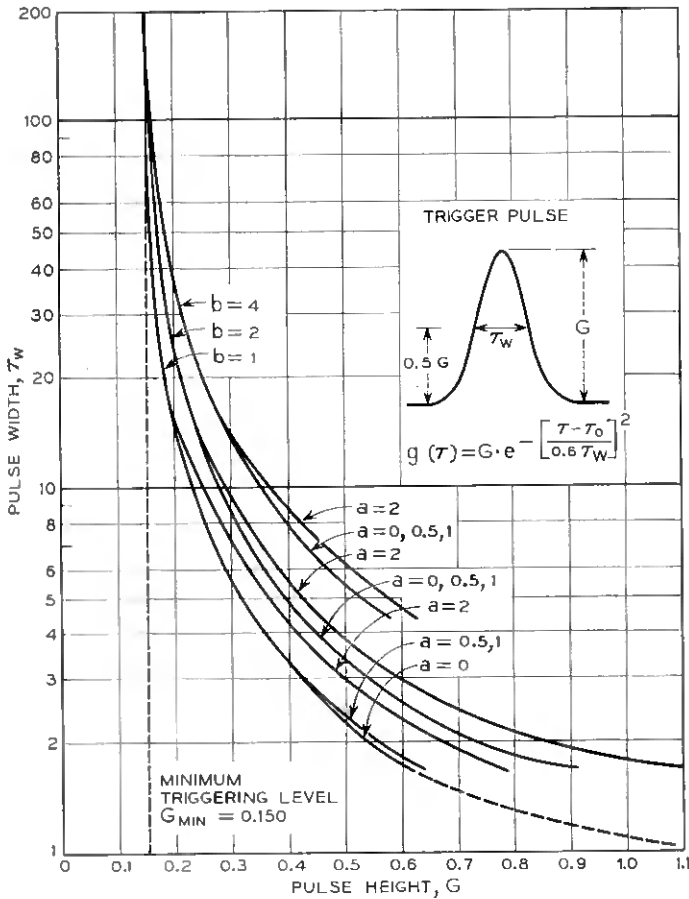


Fig. 8 — Relationships between magnitude and width of applied signal required to trigger the blocking oscillator with $b = 1, 2$ and 4 .

TABLE II

τ_w	$b = 0$			$b = 1$			$b = 2$			$b = 4$		
	$a = 0$	$a = 0.5$	$a = 2.0$	$a = 0$	$a = 0.5$	$a = 2.0$	$a = 0$	$a = 0.5$	$a = 2.0$	$a = 0$	$a = 0.5$	$a = 2.0$
	$G =$	$G =$	$G =$	$G =$	$G =$	$G =$	$G =$	$G =$	$G =$	$G =$	$G =$	$G =$
0.835	0.57	0.925	1.4	1.7	1.7	1.7	0.91	0.91	0.91	0.58	0.58	0.605
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
	0.3507	0.5344	0.8884	1.2943	1.2943	1.2943	1.2692	1.2692	1.2692	1.7953	1.7953	1.8996
1.67	G =	G =	G =	G =	G =	G =	1.2692	1.2692	1.2692	1.7953	1.7953	1.8996
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
	0.334	0.5344	0.8884	1.2943	1.2943	1.2943	1.2692	1.2692	1.2692	1.7953	1.7953	1.8996
4.175	G =	G =	G =	G =	G =	G =	1.1398	1.1398	1.1398	0.433	0.433	0.47
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
	0.24	0.34	0.335	0.34	0.34	0.34	0.328	0.328	0.328	1.8004	1.8004	2.1376
6.08	0.3758	0.5428	0.7724	0.7933	0.7933	0.7933	1.1890	1.1890	1.1890	1.8004	1.8004	2.1376
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
8.35	0.215	0.215	0.27	0.28	0.28	0.28	0.31	0.31	0.31	0.4	0.4	0.43
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
16.70	0.4342	0.4342	0.8016	0.8684	0.8684	0.8684	1.3360	1.3360	1.3360	2.0875	2.0875	2.8380
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
41.75	0.2	0.2	0.235	0.225	0.225	0.225	0.224	0.224	0.224	0.27	0.27	0.27
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
83.50	0.4175	0.4175	0.7098	0.6263	0.6263	0.6263	1.2358	1.2358	1.2358	2.0040	2.0040	2.8380
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
167.0	0.17	0.17	0.17	0.19	0.19	0.19	0.224	0.224	0.224	0.27	0.27	0.27
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
41.75	0.3340	0.3340	0.3340	0.6680	0.6680	0.6680	1.2358	1.2358	1.2358	2.0040	2.0040	2.8380
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
83.50	0.159	0.159	0.159	0.167	0.167	0.167	0.175	0.175	0.175	0.195	0.195	0.195
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
167.0	0.3758	0.3758	0.3758	0.7098	0.7098	0.7098	1.0438	1.0438	1.0438	1.8788	1.8788	2.8380
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
167.0	0.151	0.151	0.151	0.152	0.152	0.152	0.157	0.157	0.157	0.165	0.165	0.165
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
167.0	0.0835	0.0835	0.0835	0.1670	0.1670	0.1670	0.5845	0.5845	0.5845	1.2525	1.2525	1.2525
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
167.0	0.151	0.151	0.151	0.151	0.151	0.151	0.151	0.151	0.151	0.151	0.151	0.151
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
167.0	0.1670	0.1670	0.1670	0.1670	0.1670	0.1670	0.1670	0.1670	0.1670	0.1670	0.1670	0.1670
	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =	G =
	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =	N =
\bar{N}	0.3746	0.3931	0.4533	0.8183	0.8255	0.8255	1.2023	1.2023	1.2023	1.9312	1.9312	2.0516

the circuit and this magnitude increases with a and b . The curves indicate a stronger dependence of G upon b than upon a , meaning that the output capacitance and the load resistance affect the trigger sensitivity of the circuit more than the leakage inductance does.

The fact that the curves in Figs. 7 and 8 resemble hyperbolas suggests that the product of the pulse width and the difference between the magnitude and its minimum value is nearly constant; that is,

$$(G - G_{\min})\tau_w \doteq N(a, b), \quad (38)$$

where N is a function of a , b and the nonlinear characteristic. To investigate this further, the above product was formed for all the measurements taken, and the results are shown in Table II. It is seen that the value that N assumes for each measurement varies randomly around its mean value with no consistent deviation in any direction, except at large pulse widths, where it is quite small, and at short pulse widths, where it is quite large. Taking into account the subjectiveness involved in determining the point at which the response barely flips over to the quasi-stable operating point, the accuracy of the computer and the fact that at large pulse widths the value of G is very nearly equal to G_{\min} , it must be concluded that (38) is indeed a good approximation over particular regions of τ_w . Since the area of the Gaussian pulse in (36) is

$$A = \int_{-\infty}^{\infty} g(\tau)d\tau = \sqrt{\pi}G\tau_w, \quad (39)$$

it can be concluded that the area of an equivalent trigger pulse of magnitude $(G - G_{\min})$, not the area of the actual trigger pulse, stays approximately constant for given values of a and b . This brings out the fact that the blocking oscillator is not strictly an amplitude threshold device, but an energy threshold device.

The data given in Table II can be closely approximated by a function of the form

$$N(a, b) = K_0 + K_1a + K_2b. \quad (40)$$

For $1.67 \leq \tau_w \leq 41.75$, and $b \neq 0$, these values are: $K_0 = 0.38$, $K_1 = 0.06$, $K_2 = 0.42$. With these values in (40), the computed $N(a, b)$ agrees with the tabulation in Table II to within better than 30 per cent. In addition to the inaccuracies introduced by the computer and operator when G is close to G_{\min} , it is to be expected that, at the largest pulse widths, $N(a, b)$ will be independent of both a and b . For $\tau_w = 167$, Table II confirms the contention that the very slow trigger signal acts strictly to vary the bias. When τ_w is between 16.70 and 83.5, $N(a, b)$ is independent

of a for the range of parameters covered. Only the coefficients of the second and third derivatives with respect to the base current are functions of a . Therefore it can be concluded that this range of τ_w covers a transition region where only the first derivative term (ω_0 and collector capacity) is important. At the other extreme — the very narrow trigger pulse — the dependence on a is larger than indicated by the value of K_1 derived for the region $1.67 \leq \tau_w \leq 41.75$. Here the dependence on the higher derivatives (leakage inductance) is more pronounced, again confirming our physical reasoning.

The case $b = 0$ is of little practical concern, since this corresponds to an undamped output circuit working into a short-circuited load.

Other implications of the relationship between trigger magnitude and width can be brought out more clearly if we examine the unnormalized version of (40). If we substitute (40) in (38), neglect the weak dependence on leakage inductance, neglect parasitic capacity from collector to ground, and substitute from the normalizing equations, we get

$$(I_t - I_{t \min})T_w = \left(\frac{I_{cs}}{n} - I_{bc} \right) (1 + n) \left(\frac{K_0}{\omega_\alpha} + K_2 1.5C_c R_c \right). \quad (41)$$

From this relationship it can be seen that the magnitude-width product decreases as ω_α is increased until K_0/ω_α can be neglected with respect to $K_2 1.5C_c R_c$. This dependence on ω_α can also be viewed in an equivalent manner. Since K_0 and K_2 are both about equal to 0.4, (41) can be written in the following form

$$(I_t - I_{t \min})T_w \doteq \left(\frac{I_{cs}}{n} - I_{bc} \right) K_0 \frac{(1 + b)}{\omega_0}. \quad (42)$$

It will be recalled, from the discussions on the equivalent circuit for the transistor, specifically (18), that $(1 + b)/\omega_0$ is a good approximation to the large-signal (saturating step) delay time of the resistively loaded common emitter stage. Therefore (42) can be written

$$(I_t - I_{t \min})T_w \doteq K_0 \left(\frac{I_{cs}}{n} - I_{bc} \right) T_0. \quad (43)$$

As ω_α is increased, the delay time of the basic transistor stage (T_0) without external feedback decreases. In the blocking oscillator this is equivalent to a reduction of one of the time constants in the positive feedback loop, thereby permitting a more rapid build-up of the regenerative loop, or, alternately, regeneration can occur for a lower trigger signal for a given trigger width.

The dependence of this product on n is more difficult to visualize, since K_0 , K_1 and K_2 are functions of the circuit operating points, which in turn are dependent on n . As mentioned previously, $I_{t \min}$ will vary with n in much the same way that the distance between the stable and unstable operating point varies with n . For small n , $I_{t \min}$ will be large and decrease to a minimum as n increases, and then begin to increase with n . This can be seen from the graphical construction of Figs. 4 and 5, and it has also been observed experimentally. It is reasonable to expect that the product given by (42) will vary in a similar manner, though no rigorous proof can be given here.

For trigger signals that are nonzero over only a finite time interval, a relationship between magnitude and width similar to that shown in Figs. 7 and 8 exists. The principal difference is that the curves corresponding to those of Fig. 7 or 8 are now asymptotic to a finite value of width as well as of magnitude. In other words, for such trigger signals there exists a minimum width, $\tau_{w \min}$, below which the blocking oscillator cannot be triggered even when the magnitude of the trigger pulse is made infinite. That such a minimum width exists can be easily realized by considering (35) when a square pulse of infinite magnitude is applied to the circuit. In such a case, the collector current generator will saturate immediately, making $f[y + g(\tau)] = 1$. Hence, the actual forcing function is finite even if the magnitude of $g(\tau)$ is infinite, and the applied trigger pulse must therefore have a finite width in order to make the circuit flip over to its quasistable operating point. The reason $\tau_{w \min}$ is not finite in the case of the Gaussian pulse is the fact that this function is nonzero except at infinity. Thus, the width of $g(\tau)$ must be zero in order to make the duration of the actual forcing function finite.

These considerations may be used to calculate $\tau_{w \min}$. As an example, let us consider the case when $a = b = 0$. The solution to (35) for $\tau \leq \tau_{w \min}$ is then simply

$$y(\tau) = 1 - e^{-\tau}, \quad (43)$$

and the minimum width is now the time sufficient to bring y up to the unstable operating point; that is, $y = \frac{1}{2}$. This gives $\tau_{w \min} = 0.693$. In the more general case, when a and b are nonzero, the procedure is somewhat more complicated. It involves finding first the relationship between y and its first and second derivatives necessary to slide the circuit past its unstable operating point when $g(\tau)$ is zero, and then determining the width that will make these quantities satisfy this relationship. Since this procedure is usually quite laborious, it is more practical to estimate $\tau_{w \min}$ from two points on the magnitude-width curves.

VI. DELAY AND RISE TIME

The delay and rise time of the blocking oscillator are, in general, dependent upon all parameters in the circuit. The most important of these are the magnitude and width of the trigger signal, the shape of the i_c/i_b characteristic, the cutoff frequency ω_0 (or ω_α and n) and, finally, the collector capacitance in conjunction with the load resistance and the leakage inductance. To see clearly how each particular parameter affects the delay and rise time, these will be discussed separately, although it may not be possible to vary some of them independently in the actual physical circuit. In the following, the delay time is defined as the time difference between the maximum of the trigger pulse and the 50 per cent point of the normalized output. The rise time is the time required for the output to traverse 10 to 90 per cent of its final value.

6.1 *The Effect of Trigger Magnitude and Width*

It is instructive to consider first the transient responses when the blocking oscillator is simply released from various initial positions. In such cases these responses are obtained from (35) by setting $g(\tau) = 0$ and assigning initial values to y and its derivatives. In particular, if the collector capacitance and the leakage inductance are so small that they may be neglected, an analytical expression for these responses may be obtained. If we limit ourselves to this particular case by setting $a = b = 0$ in (35) and using the analytical approximation to the i_c/i_b characteristic, the solution to this equation becomes (by simple separation of variables)

$$\frac{y(y-1)}{y-0.5} = \frac{y(0+)[y(0+)-1]}{y(0+)-0.5} e^{-\tau}, \quad (44)$$

where $y(0+)$ represents the initial value of y . The bistable nature of the blocking oscillator is reflected in (44), because, as $\tau \rightarrow \infty$, y approaches either the stable ($y = 0$) or the quasistable ($y = 1$) operating point. Which one it will attain in an actual case depends upon whether the initial position is below or above the unstable operating point. In the neighborhood of the point ($y = 0.5$) the output voltage becomes an exponential ascending function of time, thereby exhibiting the unstable characteristic of this point. Fig. 9 shows the responses of (44) for six initial values located between the unstable and the quasistable operating points.¹² It is seen that, as the initial value of y is decreased toward the unstable operating point, the delay and rise time of the circuit become larger, increasing rapidly as $y(0+)$ approaches this point. To minimize

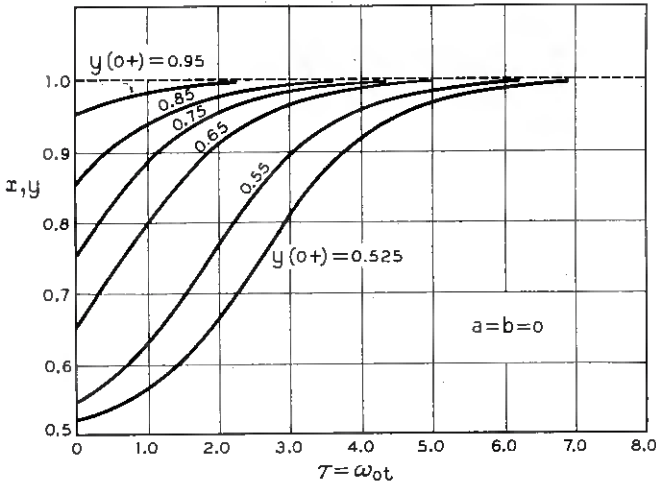


Fig. 9 — Normalized time response of the blocking oscillator when the circuit is released from various initial positions.

rise time, the magnitude and width of the trigger signal should at least be such that the circuit has been brought well beyond the unstable point when the trigger terminates.

When the coefficients a and b (or the trigger function) are not equal to zero, (44) cannot be solved analytically and numerical methods must be used. In Fig. 10 the results of solving this equation by the Runge-Kutta method on a digital computer are shown. Again, the analytical approximation to the i_c/i_b characteristics was used and the Gaussian function represented the trigger signal. To conform with values encountered in a particular circuit, a and b were, in all of these cases, made equal to 0.447 and 0.8, respectively. In the first of these figures, Fig. 10(a), the normalized response of the base current, x , and the output voltage, y , to a slow trigger pulse, are shown. It is seen that the output voltage does not change appreciably until $g(\tau)$ reaches the critical triggering level, which, from Table II, is equal to 0.167 for $\tau_w = 41.75$. Subsequently, the output voltage switches rapidly in comparison to the trigger signal to the quasistable operating point. Hence, one should expect, in the case of a slow trigger signal, that the rise time of the blocking oscillator is more or less independent of the trigger magnitude as long as this magnitude is larger than the critical value necessary to activate the circuit. The upper curve in Fig. 11 bears out this statement. Here, the rise time is shown as a function of pulse magnitude, and it is seen that, except for a trigger magnitude close to the critical triggering level, the rise time is

virtually constant over the whole range. Thus, one may conclude that in this case the rise time is, practically speaking, dependent only upon the parameters of the blocking oscillator itself. The delay time, on the other hand, will of course vary with trigger magnitude, but only to the extent of the time it takes the trigger signal to reach the critical triggering level.

Figs. 10(b) through 10(d) show the response of the blocking oscillator to a short trigger pulse of successively increasing magnitude. Actually, in all figures except Fig. 10(c) two trigger pulses were applied, one that

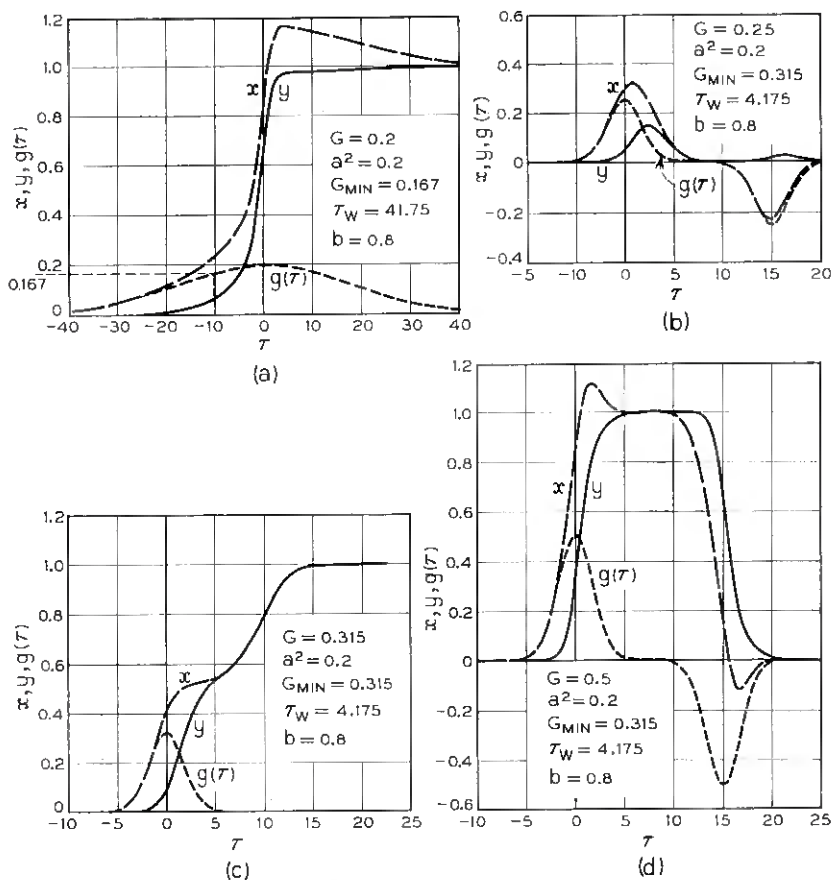


Fig. 10—Computed responses of the blocking oscillator to trigger pulses of various widths and magnitudes: (a) slow trigger pulse; (b) inadequate trigger pulse; (c) trigger pulse having just limiting magnitude; (d) adequate trigger pulse.

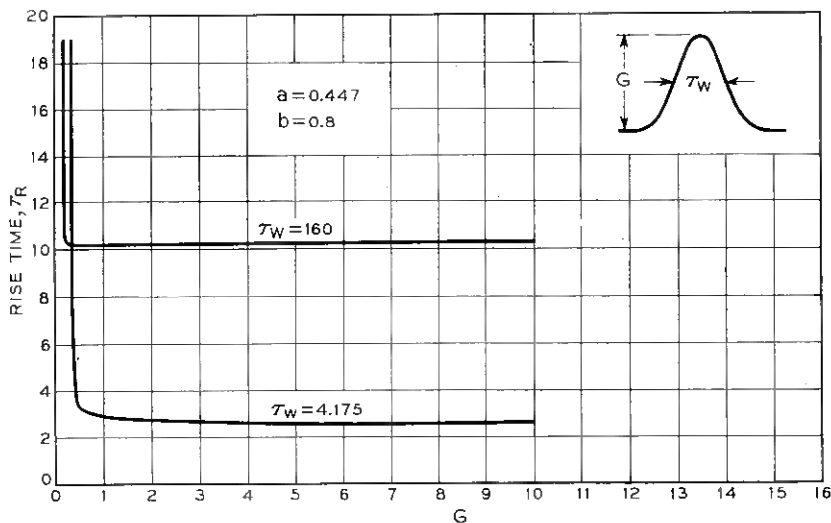
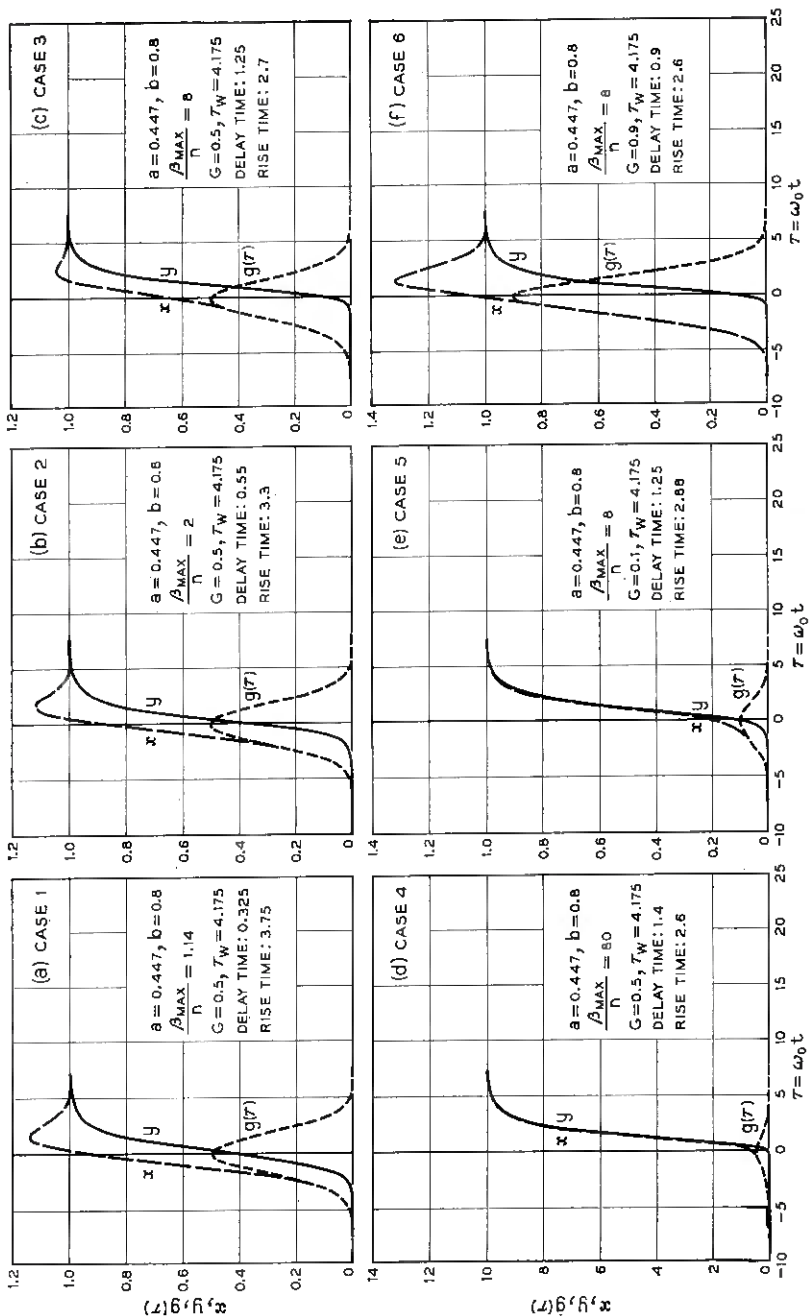


Fig. 11 — Rise time as a function of trigger pulse magnitude.

turns the blocking oscillator “on”, followed by another that turns it “off”. This was done in order to simulate the conditions for a blocking oscillator that might be used in a regenerative repeater. In Fig. 10(b) the magnitude of the trigger signal was insufficient to activate the blocking oscillator, while in Fig. 10(c) it was just barely large enough to make the circuit flip over to the quasistable operating point. The rather large rise time and the inflection point at $y = 0.5$ should be noted in the latter case, which, as explained previously, is due to the fact that the circuit is now essentially released from a position close to the unstable operating point. As the magnitude is increased further the rise time at first decreases rapidly and then levels off, approaching a limit for large values of G . This is depicted by the lower curve in Fig. 11. In this case the rise time is largely independent of trigger magnitude as long as the magnitude is at least 50 per cent larger than the critical triggering level. It should be pointed out, however, that the apparent leveling off of the rise time in the fast pulse case is only true when ω_0 may be considered to be independent of trigger magnitude. If a fast and large trigger pulse is applied such that the base current significantly exceeds the value I_{bs} during transition, the maximum value of this current should, according to (7), be used to calculate ω_0 rather than I_{bs} . This results in a larger ω_0 and therefore a smaller rise time than that indicated by Fig. 11. The correction to be applied to ω_0 in such cases is, by virtue of (6) and (8), equal to:

Fig. 12 — Transient response of the blocking oscillator for the normalized i_c/i_b characteristics shown in Fig. 6.

$$\frac{\omega_0'}{\omega_0} = \frac{1 + \frac{I_{cs}}{I_{be} - I_{bc}}}{1 + \frac{I_{cs}}{I_{b\max} - I_{bc}}} = \frac{1 + n}{1 + \frac{n}{x_{\max}}} \tag{45}$$

In most practical blocking oscillators, however, the stable and unstable operating points are located very close together, so that such correction is rarely needed. Also, the desirability of a large power gain favors the use of a design requiring only small trigger magnitudes. Finally, in the case of a slow trigger signal, such correction is never needed, since here the trigger signal does not change appreciably during the relatively fast transition from the stable to the quasistable operating point.

6.2 *The Effect of the Shape of the i_c/i_b Characteristic*

To see how the shape of the i_c/i_b characteristic affected the delay and the rise time of the blocking oscillator, transient responses were computed from (35) using for $f(x)$ the functions depicted in Fig. 6. In all cases, the location of the stable operating point was kept fixed and the circuit was assumed to be initially at rest at this point. To trigger the circuit a Gaussian pulse was used which had a 50 per cent width of 4.175 and a magnitude in each case equal to the distance between the stable and unstable operating points. In all cases, $a = 0.447$ and $b = 0.8$.

In the first set of computations the location of the unstable operating point was kept fixed while the slope of the characteristic through this point was varied by using four different functions for $f(x)$. These are the cases numbered 1 through 4 in Fig. 6, the fourth one being essentially Case 3 multiplied by ten. The slope at the unstable operating point at each of these cases, which on an unnormalized scale represents β_{\max}/n of the circuit, are, respectively, 1.14, 2, 8 and 80. However, the average value of β/n is the same in all cases, since the quasistable operating point is located at (1, 1) in the first three cases and at (10, 10) in the last one. Figs. 12(a) through 12(d) give the results of computing the transient response from (44) in each of these cases. In Fig. 13 the rise time is plotted as a function of β_{\max}/n . It is seen that the delay and rise time change only slightly over the whole range of β_{\max}/n ; the most rapid variation being confined to the region where β_{\max}/n is nearly equal to unity. The slope must be greater than this value for $f(x)$ to intersect the load line at three points. Hence, it can be concluded that the rise time decreases very little when β_{\max}/n is increased, as long as this quantity is larger than unity. In fact, Fig. 13 shows that a 70-to-1 change in loop gain results in only a 30 per cent change in rise time, and that the bulk

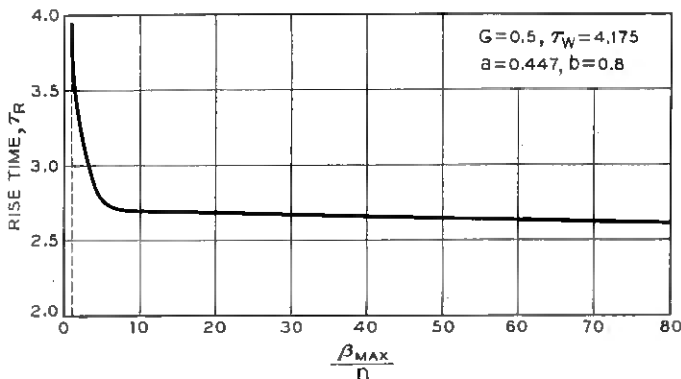


Fig. 13 — Rise time as a function of β_{MAX}/n .

of this change occurs for the first 7-to-1 increase in loop gain. The reason for this somewhat surprising result is, of course, the fact that the average of β/n , and therefore the average of the loop gain, is the same in all these cases.

Finally, it should also be noted in Figs. 12(a) through 12(d) that, while the rise time decreases as the slope of $f(x)$ becomes steeper, the delay increases. This is so because the collector current starts to flow much earlier in the case of a gentle slope than in the case where $f(x)$ is steep.

To investigate the effect of moving the unstable operating point, the transient response was also computed for Case 5 and Case 6 in Fig. 6. The unstable operating points in these two cases are, respectively, located in $x = 0.1$ and $x = 0.9$ and both functions have slopes identical to that of Case 3. Comparing the resulting responses of Figs. 12(e) and 12(f) with that of Fig. 12(c), it is seen that, as far as the delay and rise time are concerned, the location of the unstable operating point has even less effect than the slope has. A slightly smaller rise time is obtained when the unstable operating point is located farther away from the stable one, since the trigger is larger in such cases.

Before concluding this section, it should again be pointed out that the dependence of ω_0 on the average value of β must be taken into account when the above results are applied to actual circuits. The fact that $1/\omega_0$ and C_0 are nearly proportional to $\bar{\beta}$ [(7) and (15)] tends to override the somewhat weak dependencies discussed above. Transistors should therefore be selected for large ω_0 and small C_0 rather than a high β to obtain the best rise time. On the other hand, a high-low current β is desirable to increase the trigger sensitivity of the circuit. Finally, the

above results for a single pair of values of a and b and a single value of τ_w apply to other values of these parameters. The basic shapes of the curves in Figs. 11 and 13 remain practically unaltered. The principal effect of other values of these constants is to displace these curves both horizontally and vertically, with the vertical displacement the most pronounced effect.

6.3 *The Effect of Leakage Inductance and Collector Capacitance upon the Delay and Rise Time*

To investigate how the leakage inductance and collector capacitance affect the delay and rise time, the transient response of the blocking oscillator was computed from (35) for a large number of values of the constants a and b . Four different types of trigger signals were used throughout this investigation. In the first case, the effect of a very slow trigger signal was simulated by employing a Gaussian pulse for $g(\tau)$ of magnitude and width equal to $G = 0.3$ and $\tau_w = 166.7$, respectively. As pointed out in a previous section, this represents the case in which the rise time is solely dependent upon the properties of the feedback loop of the blocking oscillator and not upon the detailed behavior of the trigger signal itself. In the next two cases, Gaussian pulses of medium ($\tau_w = 25$) and extremely small ($\tau_w = 4.175$) widths were used, illustrating situations where both the trigger signal and the circuit properties affect the rise time. In all three cases, however, a magnitude of the trigger signals was selected such that the circuit operated on the flat portions of the curves in Fig. 11. In the fourth and last case, a step function of magnitude sufficient to saturate the transistor immediately ($G = 1$) was applied. Since the entire driving capability of the transistor is now employed from the very beginning, this represents the situation in which the delay and rise time are minimized. Also, the loop gain is zero during the entire transition, so that the degenerative effect of positive feedback on rise time is nonexistent in this case. Finally, in all cases, the blocking oscillator was assumed to be initially at rest at its stable operating point, and the analytical function in (37) was used to represent the i_c/i_b characteristic.

Figs. 14(a) through 14(d) give the delay time as a function of the constants, a and b , in the four cases. It is seen that in nearly all cases the delay is a linear function of b when either this constant or a is considerably larger than one. The only noticeable deviation from this rule is the $a = 8$ curve in Fig. 14(b), which starts to bend upwards for large b . The reason for this is that here the critical triggering level approaches the actual trigger magnitude as b is increased; in other words, the circuit is

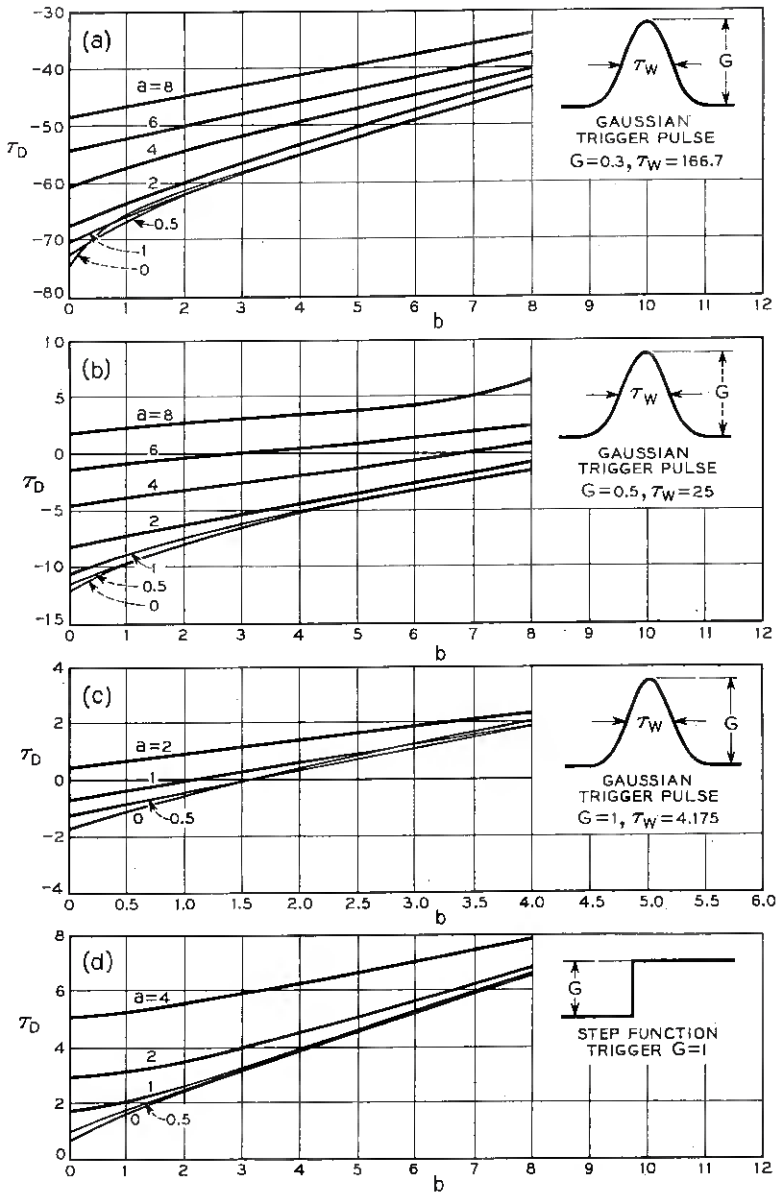


Fig. 14 — The delay as a function of a and b when the trigger signal is (a) a very slow Gaussian pulse; (b) a moderately fast Gaussian pulse; (c) a short Gaussian pulse and (d) a step function.

no longer operating on the flat portions of the curves in Fig. 11. From Fig. 14(a) it should be noted that, for values of b less than 4, the delay goes through a minimum when a is varied.

Figs. 15(a) through 15(g) show the rise time as a function of a and b for the four trigger functions. Comparing these figures, it is seen that the rise time is largest in the case of a very slow trigger signal and smallest in the case of the saturating step function. However, the variation in rise time with b is much less in the former case than in the latter one, while with a it is somewhat the other way around. Also, the dependence is stronger upon b than upon a . Finally, in all cases, the rise time goes through a minimum when a is varied and b kept fixed, this minimum being most pronounced in the case of a moderately fast trigger pulse [Fig. 15(d)].

Since it is usually not possible to obtain an analytical solution to (35), an expression for the rise time that is valid under all conditions cannot be found either. However, in two of the cases discussed above, approximate expressions for the rise time were found to fit the computed data over a limited range. The first of these pertains to the case of a very slow trigger signal and was obtained by a trial-and-error procedure. It was discovered that the expression

$$\tau_R \approx 7.85 \sqrt{1 + 2b + a} \tag{46}$$

approximates the computed data in Fig. 15(a) within ± 10 per cent. The second expression for the rise time is concerned with the case where a unit step function is applied to the circuit. Since the transistor is saturated immediately in this case, (35) is linear with $f[y + g(\tau)] = 1$, and the transistor can be characterized by the transfer function

$$\frac{Y(s)}{U(s)} = \frac{1}{(s + 1)(a^2s^2 + bs + 1)}, \tag{47}$$

where $U(s)$ is the Laplace transform of the unit step function. The poles of the transfer function are given by $s = -1, -\alpha, -\beta$, where

$$\alpha = \frac{b + \sqrt{b^2 - 4a^2}}{2a^2} \quad \text{and} \quad \beta = \frac{b - \sqrt{b^2 - 4a^2}}{2a^2}. \tag{48}$$

It can be shown that a system with a transfer function whose only singularities in the finite part of the complex plane are negative real poles will have monotonic step response.¹³ For the blocking oscillator, α and β are real when $b \geq 2a$ and a monotonic step response results. Under these conditions, Elmore's definition¹⁰ of delay and rise time gives results in good agreement with those obtained by the exact calculations

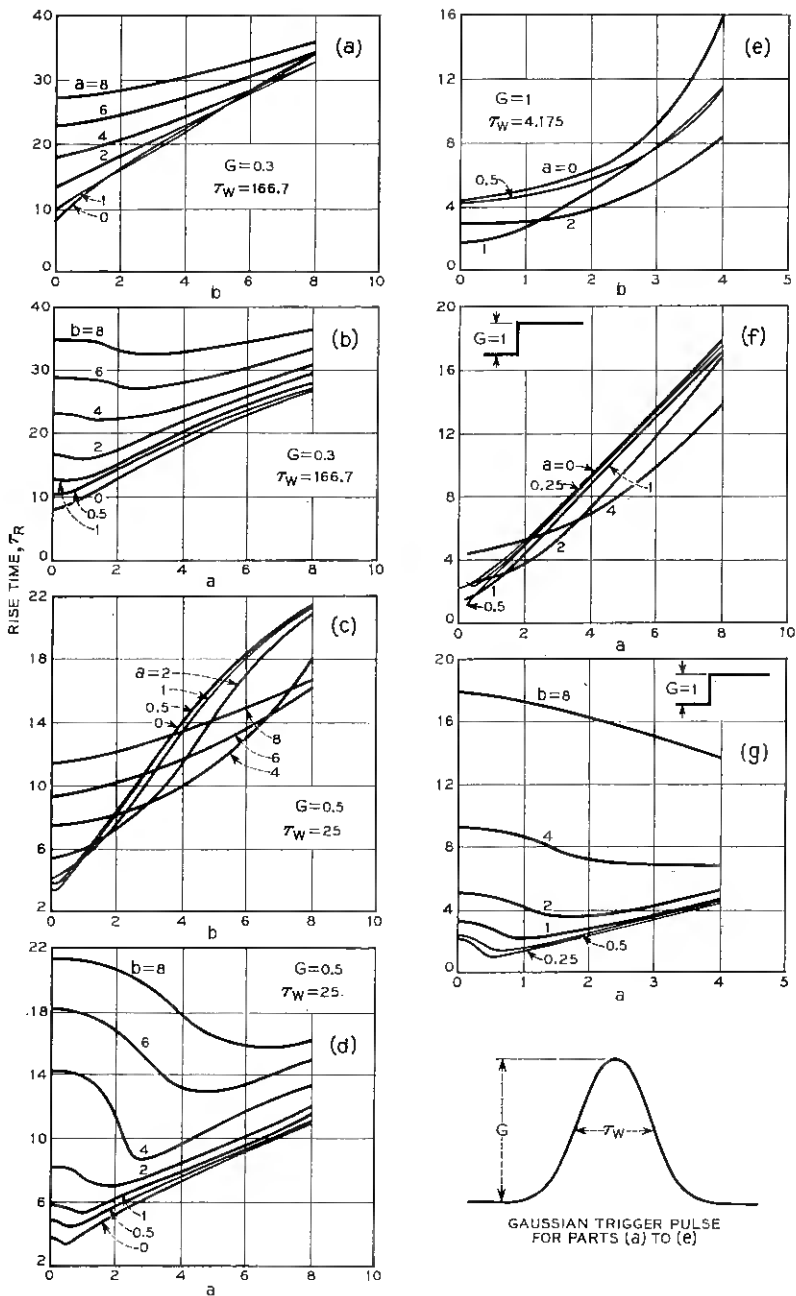


Fig. 15 — Rise time of the blocking oscillator as a function of the constants a and b when the trigger signal is: (a) and (b) a very slow Gaussian pulse; (c) and (d) a moderately fast Gaussian pulse; (e) a very short Gaussian pulse; (f) and (g) a step function of sufficient magnitude to saturate the transistor immediately.

used to obtain the results of Fig. 15. Application of Elmore's definitions yield the normalized delay

$$\tau_D = 1 + b \quad b \geq 2a \quad (49)$$

and the normalized rise time

$$\tau_R = 2.2 \sqrt{1 + b^2 - 2a^2} \quad b \geq 2a \quad (50)$$

or, in unnormalized form (neglecting parasitic capacity):

$$T_D = \frac{1}{\omega_0} + R_c C_0 = (1 + n) \left(\frac{1}{\omega_\alpha} + 1.5 R_c C_c \right), \quad (51)$$

$$\begin{aligned} T_R &= 2.2 \sqrt{\frac{1}{\omega_0^2} + R_c^2 C_0^2 - 1.5 L_c C_0} \\ &= 2.2(1 + n) \sqrt{\frac{1}{\omega_\alpha^2} + 2.25 R_c^2 C_c^2 - 2.25 \frac{L_c C_c}{1 + n}}. \end{aligned} \quad (52)$$

As expected, n and $R_c C_c$ should be selected as small as possible to reduce the rise time while ω_α should be large. Equation (52) suggests that there exists a value of the leakage inductance that will minimize the rise time. The exact value of the leakage inductance that accomplishes this cannot be determined from (52), since it is strictly valid only when the response is monotonic.

Finally, in Figs. 16(a) through 16(c) the overshoot is shown as a function of a , with $b/2a$ as a parameter. It is seen that, as $b/2a$ increases, the overshoot decreases, becoming zero when $b/2a \geq 1$. It should be noted that, for fixed values of $b/2a$, the overshoot approaches a constant when a is large. Spot checks of these analog computer results were made on a digital computer, with excellent agreement between the two results.

6.4 Conditions for Monotonic Response

In the preceding section it has been shown that the blocking oscillator will have monotonic or oscillatory response when it is driven by a saturating step function, depending on whether b is greater than or less than $2a$ respectively. For example, the response will be monotonic if

$$b > 2a \quad \text{or} \quad R_c > 2 \sqrt{\frac{L_c}{C_0}}, \quad (53)$$

and be oscillatory if

$$b < 2a \quad \text{or} \quad R_c < 2 \sqrt{\frac{L_c}{C_0}}. \quad (54)$$

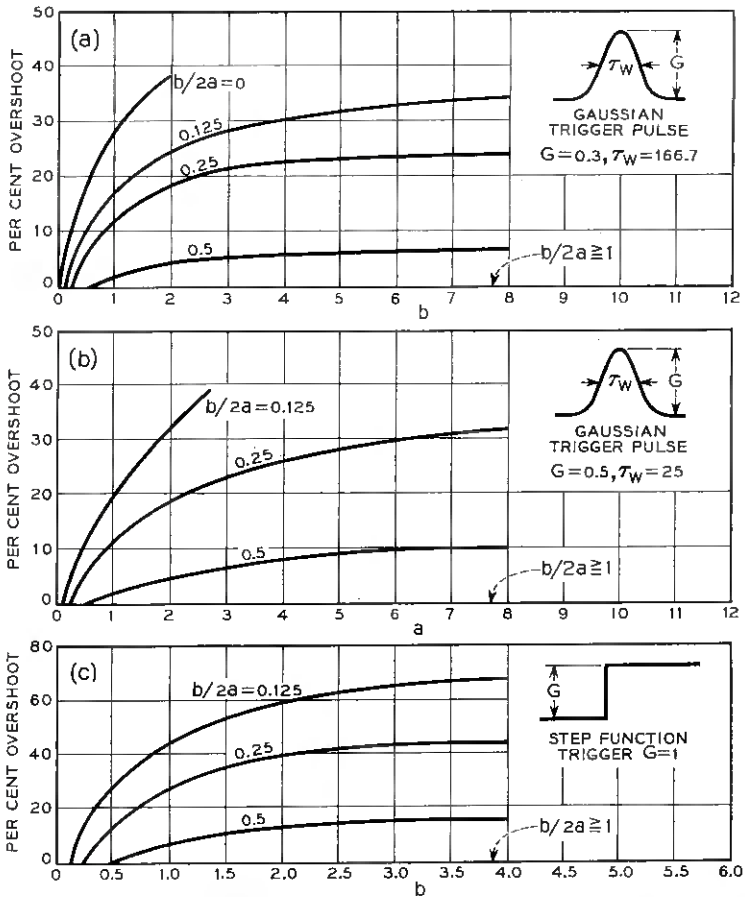


Fig. 16 — The percentage overshoot of the blocking oscillator when the trigger signal is (a) a very slow Gaussian pulse; (b) a moderately fast Gaussian pulse and (c) a step function.

The condition $b = 2a$ yields two coincident real roots and divides the a - b plane into regions of monotonic and oscillatory responses, as shown in Fig. 17. Two real roots can also result when either α or β is equal to -1 . This leads to the parabola

$$b = 1 + a^2 \quad (55)$$

in the a - b plane. At the point $(1, 2)$ this parabola is tangent to the straight line $b = 2a$ and the blocking oscillator has three coincident real poles at $s = -1$.

There is one exception to the above division in the a - b plane. When $b = 2a^2$ and $0 < a < 1$, a pair of complex conjugate poles whose real part = -1 arises. The poles all lie on a line that is parallel to the imaginary axis in the complex plane and spaced one unit to the left. In this case, the impulse response is tangent to the zero line, so that the step response is monotonic. It can be shown that in the region $b \geq 2a^2$, $0 \leq a \leq 1$ the response will be monotonic. The rise time in this case with $b = 2a^2$ is given by

$$\tau_r = 2.2 \sqrt{1 + 4a^4 - 2a^2}. \tag{56}$$

The above rise time is a minimum when $a = 0.5$; the minimum rise time then is 1.87. Actual transient calculations for this condition yield a minimum rise time of 1.45 at $a = 0.45$. In this case, the 10 to 90 per cent rise time is faster than it is in any of the other cases of monotonic response. However, the response proceeds quite slowly from the 90 per cent point to its final unit asymptote. This behavior explains the discrepancy between the results of (56) and the exact transient calculations.

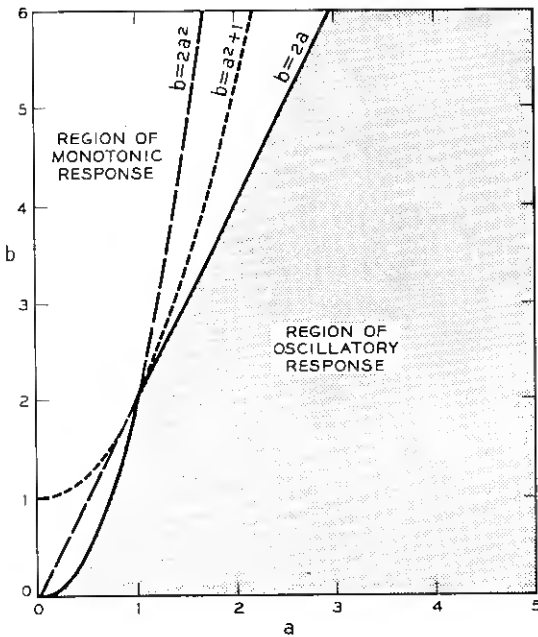


Fig. 17 — Conditions for monotonic step response.

Ideally, the blocking oscillator should have a response that has both as small a rise time as possible and a flat, smooth top. The best physical compromise to this requirement is to make the circuit nearly critically damped during this interval. In other words, the coefficients a and b should be selected such that they lie in a region close to the straight line, $b = 2a$, in Fig. 17, i.e.,

$$b \doteq 2a \quad \text{or} \quad R_c \doteq 2 \sqrt{\frac{L_c}{C_0}}. \quad (57)$$

Furthermore, Figs. 15(f) and 15(g) show that a minimum rise time can be achieved for small values of a and b . In fact, a rise time faster than the minimum rise time of the transistor itself ($2.2/\omega_0$) can be obtained. Some overshoot is associated with these cases. This result is similar to that achieved with series peaking in a conventional interstage, and confirms the conclusions of Linvill and Mattson.³ However, in this case the blocking oscillator is acting as an overdriven amplifier with zero loop gain. The larger the drive the faster the response, since ω_0 increases with drive. Smallest rise time can only occur when the power gain required of the blocking oscillator is low. This requirement may be satisfied in some computer applications where the primary function of the blocking oscillator is to retime a relatively sharp pulse train. In applications where the blocking oscillator must reshape as well as retime a pulse train, such as in pulse code modulation, considerable power gain will be demanded from the circuit and the minimum rise time cannot be achieved. This is the analog of the familiar gain-bandwidth product for linear systems.

When signals other than the saturating step function are used to trigger the blocking oscillator, the manner in which the quasistable operating point is approached is dependent on the parameters a and b , the nonlinear characteristic and the trigger signal. When the transistor saturates after suitable triggering $f[y + g(\tau)] = 1$. At this point, as noted previously, (35) is linear and the Laplace transform of the output can be expressed as:

$$y(s) = \frac{1}{s[a^2s^3 + (a^2 + b)s^2 + (1 + b)s + 1]} + \frac{[a^2s^2 + (a^2 + b)s + (1 + b)]y(t_s) + [a^2s + (a^2 + b)]y'(t_s) + a^2y''(t_s)}{a^2s^3 + (a^2 + b)s^2 + (1 + b)s + 1}, \quad (57)$$

where $y(t_s)$, $y'(t_s)$ and $y''(t_s)$ are, respectively, the values that y , its first and its second derivatives attain when the transistor becomes

saturated. The kind of response $y(t)$ will have is therefore dependent upon both the initial conditions at saturation and the roots of the characteristic equation

$$a^2 s^3 + (a^2 + b)s^2 + (1 + b)s + 1 = (s + 1)(a^2 s^2 + bs + 1) = 0. \tag{58}$$

Relationships between the initial conditions and a and b can be established such that the quasistable operating point is approached monotonically.¹³ This will not be pursued here, since the values of the derivatives at saturation are not normally measured quantities. On purely physical grounds, it would be expected that the conditions given previously for monotonic response in the case of the saturating step function should also apply here. This follows when it is realized that the saturating step function brings the full capability of the transistor into play immediately, thereby imparting maximum energy into the output circuit. For other trigger signals this will not be true, so that the energy transferred to the output circuit will be smaller than it will be in the case of the saturating step. Since the values of the output function and its derivatives when the transistor goes into saturation are a measure of this energy, it would appear that, if no overshoot occurs with the saturating step, none should occur for other trigger signals. This heuristic argument is bolstered by the results displayed in Figs. 18(a) and 18(b). These figures give the analog computer solutions to (35) for various values of a and b when either a fast or slow trigger is applied to the circuit. In both cases the response adheres to the conditions given on Fig. 17 for monotonic response.

VII. THE "ON" INTERVAL

After the blocking oscillator has reached the quasistable operating point the magnetizing inductance becomes charged slowly, reducing the current fed back to the base. This causes the load line, and therefore the quasistable operating point, to move to the left, in Fig. 4 or 5, until the load line is tangential to the $I(i_b)$ curve. Beyond this point (J'), the load line intersects $I(i_b)$ only in the cutoff region, and the circuit now has only one possible operating point. Hence, a rapid "jump" toward cutoff takes place at J' , and it is therefore natural to define the "On" interval as the time it takes the base current to decrease from the quasistable operating point to the value it has at the point of tangency, J' , that is:

$$I_{bJ} \leq i_b \leq I_{b0} \tag{59}$$

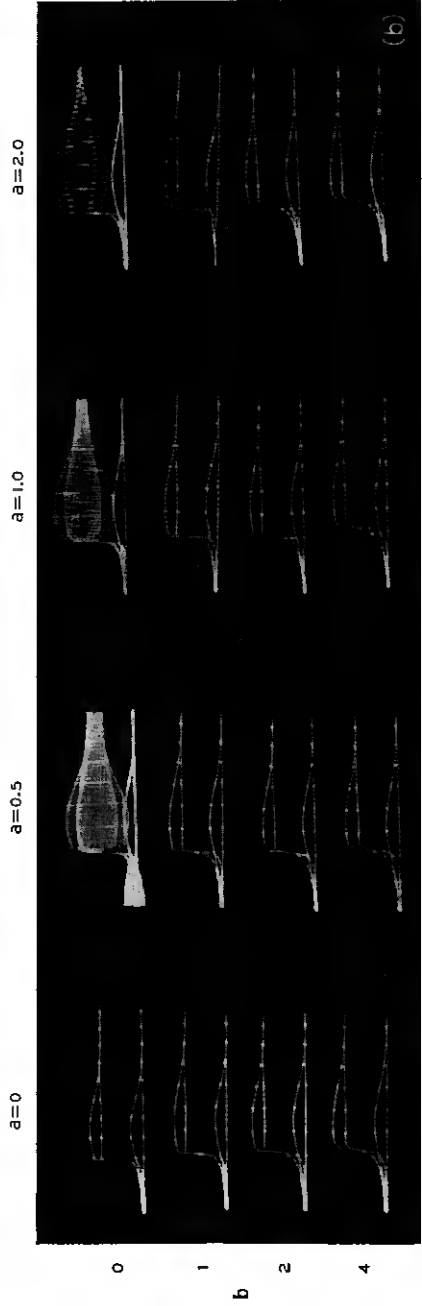
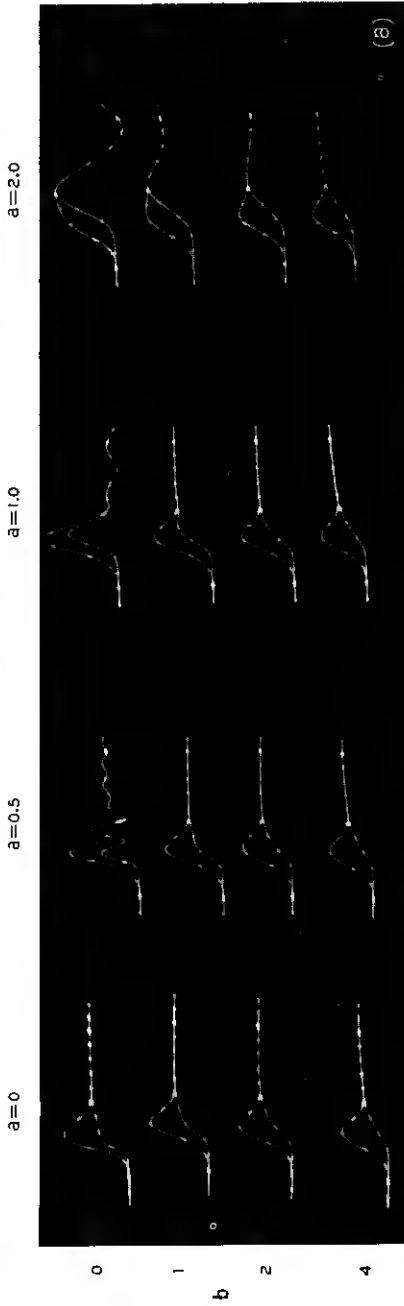


Fig. 18 — (a) Transient responses when $G = 1$ and $T_w = 4.175$, with the upper curves being base current, x , and the lower curves output voltage, y ; (b) transient responses when $G = 0.3$ and $T_w = 160$, with the upper curves being base current, x , the middle curves output voltage, y , and the lower curves the linear function

In Fig. 19 the results of solving the complete equation [(116) of Appendix B] of the blocking oscillator on an analog computer are shown. The analytical representation (37) of the i_c/i_b characteristic was used and $V(x)$ was approximated by two straight-line segments. Also, the constants a, b, d and e were kept fixed while the constant c , which is proportional to the magnetizing inductance, was varied. It is seen that, except when c is small and the circuit oscillates, the base current and the base voltage vary rather slowly during the "On" interval. Hence, during this interval terms involving derivatives of x and $V(x)$ may be neglected. This also follows from the fact that, in all cases of practical interest, the magnetizing inductance is very large compared with the other storage elements in the circuit. For the same reason, the amount of current collected by this inductance during the "Transition On" interval is also assumed to be negligible, i.e., $m(\tau_f) \doteq 0$. Assuming the trigger signal to have terminated by this time, the equation governing the normalized base current during this interval therefore reduces to

$$x - h(x) + \frac{1}{c} \int_0^{\tau} V(x) d\tau = 0$$

(60)

or

$$c \frac{dx}{d\tau} + \frac{V(x)}{1 - \frac{dh(x)}{dx}} = 0.$$

The unnormalized version of this equation becomes

$$n^2 L_m \frac{di_b}{dt} + \frac{V_b(i_b) - E_{bb}}{1 - \frac{dI(i_b)}{di_b}} = 0. \quad (61)$$

The above equation reflects the operation of the circuit during this interval exactly as it was described at the beginning of this section. First, it indicates that the base current decreases rather slowly, since L_m is large and $V_b(i_b) - E_{bb}$ small. Then, when $dI(i_b)/di_b$ becomes equal to unity [which corresponds to the point of tangency (J') in Fig. 4], it predicts that a jump takes place, by virtue of the fact that $di_b/dt = \infty$ at this point. Since

$$1 - \frac{dI(i_b)}{di_b} = \frac{R_b' + \frac{dV_b(i_b)}{di_b}}{R_b'} \left[1 - \frac{1}{n} \frac{R_b'}{R_b' + \frac{dV_b(i_b)}{di_b}} \frac{dI_c(i_b)}{di_b} \right], \quad (62)$$

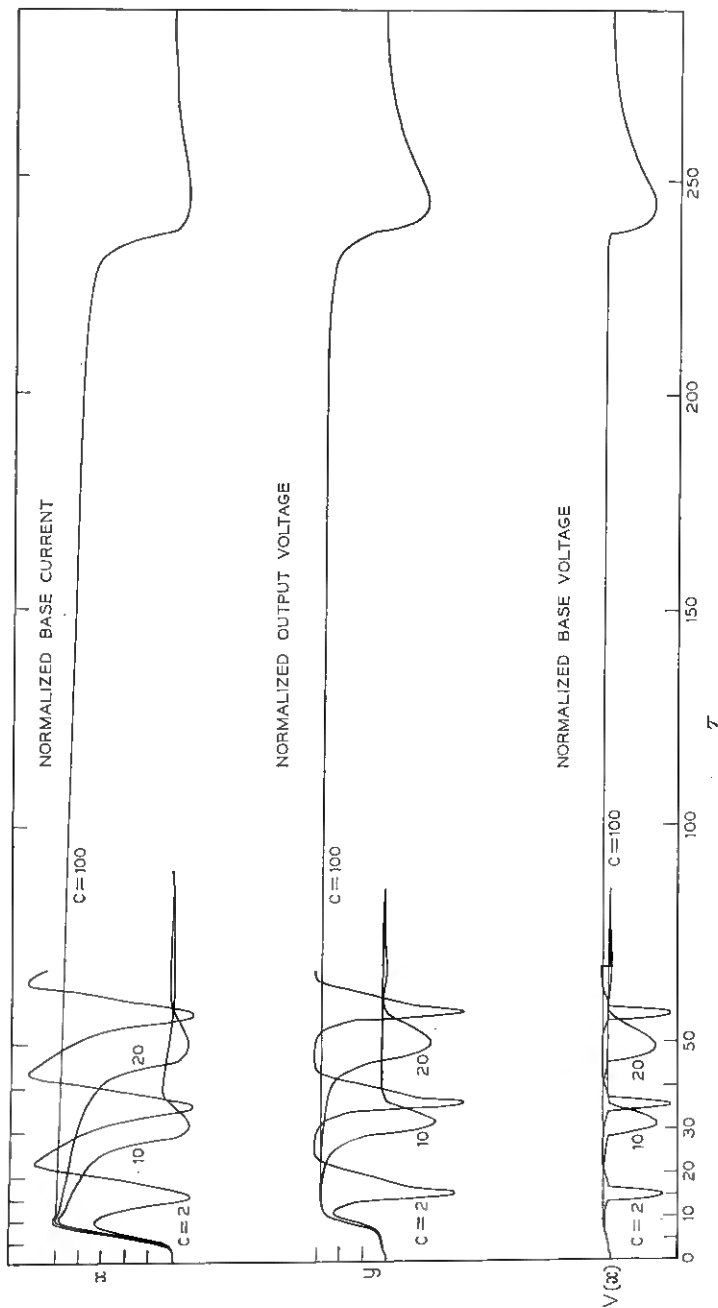


Fig. 19 — Natural pulse width of the blocking oscillator when $c = 2, 10$ and 100 . An analytical $f(x)$ was used and $\sigma^2 = 0.2, b = 0.8, d = 0.28, G = 0.5$ and $\tau_w = 4.175$.

it can be seen that the “midband” loop gain of the circuit is equal to unity at the point where the jump takes place.

As a first step toward deriving a formula for the natural pulse width of the blocking oscillator, let us introduce the above expression into (61) and also take advantage of the fact that $V_b(i_b)$ is very nearly linear in the saturation region, so that $V_b(i_b) = r_{bs}i_b$. Thus, (61) becomes:

$$n^2 L_m \frac{R_b' + r_{bs}}{R_b' r_{bs}} \frac{di_b}{dt} + \frac{i_b - \frac{E_{bb}}{r_{bs}}}{1 - \frac{1}{n} \frac{R_b'}{R_b' + r_{bs}} \frac{dI_c(i_b)}{di_b}} = 0. \quad (63)$$

By substituting the analytical approximation for $I_c(i_b)$, the above equation may be solved exactly.¹⁴ However, in cases of practical interest, the quasistable operating point extends quite far into the saturation region and, furthermore, most i_c/i_b characteristics bend quite sharply before going into saturation. Therefore the jump point, j' , is positioned very close to the saturation level. Hence, the transistor is saturated during almost the entire “On” interval with the result that $dI_c(i_b)/di_b$ is equal to zero except in the near vicinity of the point j' . For the purpose of deriving an expression for the pulse width, it is therefore an excellent approximation to neglect the loop gain term in (63). Solving this equation under these conditions, introducing i_b as equal to I_{bs} and I_{bj} at the beginning and end of the “On” interval and substituting for R_b' from (24), the expression for the natural pulse width may be written

$$T_w \sim n^2 L_m \frac{R_b + r_{bs} \left[1 + \frac{R_b}{n^2 L_m} \left(\frac{1}{\omega_0} + R_c C_0 \right) \right]}{R_b r_{bs}} \ln \frac{I_{bs} - \frac{E_{bb}}{r_{bs}}}{I_{bj} - \frac{E_{bb}}{r_{bs}}}. \quad (64)$$

In most circuits

$$R_b \gg r_{bs} \left[1 + \frac{R_b}{n^2 L_m} \left(\frac{1}{\omega_0} + R_c C_0 \right) \right] \quad \text{and} \quad I_{bs} \doteq \frac{I_{cs}}{n}, \quad (65)$$

so the above equation reduces to

$$\frac{T_w}{L_m / r_{bs}} \doteq n^2 \ln \frac{\frac{1}{n} - \frac{E_{bb}}{I_{cs} r_{bs}}}{\frac{I_{bj}}{I_{cs}} - \frac{E_{bb}}{I_{cs} r_{bs}}}. \quad (66)$$

Figs. 20 and 21 show the normalized pulse width as a function of n for various values of I_{bj}/I_{cs} and $E_{bb}/I_{cs} r_{bs}$ as computed from (66). It is seen

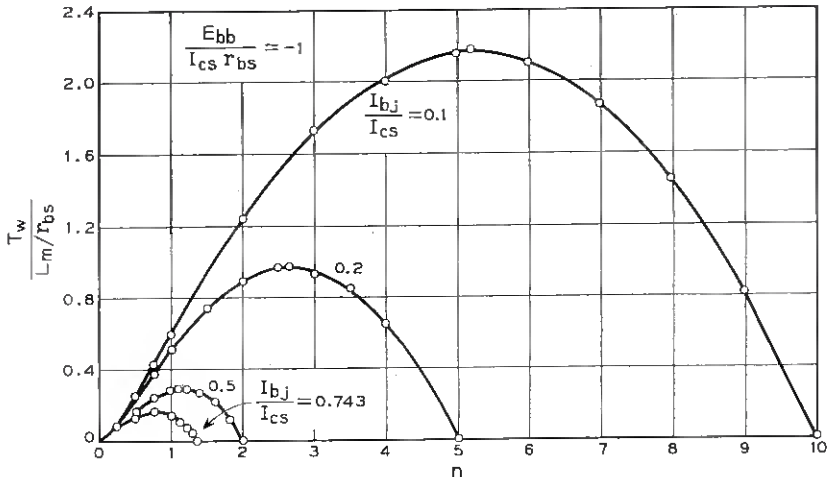


Fig. 20 — Pulse width as a function of turns ratio, with $E_{bb}/I_{cs}r_{bs} = -1$.

that the pulse width goes through a maximum and is nonzero over only a finite range of n , the limiting values of n being 0 and I_{cs}/I_{bj} . Physically, the upper limit corresponds to the coincidence of the quasistable and the unity loop gain points, making it impossible for the circuit to remain in a conducting state for a finite time. It should be noted from these figures that, for a fixed turns ratio, the natural pulse width exhibits a strong dependence upon both $E_{bb}/I_{cs}r_{bs}$ and I_{cs}/I_{bj} , i.e., on both the nonlinearities and the bias.

Before concluding this section, two things should be pointed out. First, the quantity I_{bj} in (66) is, strictly speaking, also a function of n , because, as n is varied, the point of tangency in Fig. 4 changes. However, since the i_c/i_b characteristic usually bends quite sharply before going into saturation, I_{bj} will vary only slightly with n . Second, in the above analysis the effect of minority carrier storage has been neglected. The presence of such carriers in the collector junction introduces a delay before the transistor can leave the saturation region; as if the $I(i_b)$ curve had been temporarily moved somewhat to the left in Fig. 4, thereby causing the pulse width to be slightly larger than that predicted by (60). However, this effect may be partially accounted for by calculating the pulse width from (64) rather than from (66), using, as explained in the section on the equivalent circuit, the value of ω_0 appropriate to turnoff.

VIII. THE "TRANSITION OFF" INTERVAL

In discussing the "Transition Off" interval one must discriminate between two cases: one in which the blocking oscillator turns off by

itself, and the other when it is triggered off by virtue of an externally applied signal. In the following, these two cases will be discussed separately.

8.1 Internal "Turn Off"

As explained in the previous section, when the magnetizing inductance becomes sufficiently charged so that the load line is tangent to the $I(i_b)$ curve at the unity loop gain point, J' , a rapid transition toward the other point of intersection between the load line and the characteristic takes place. Since the input impedance to the transistor is very high in the reverse direction, the bias current, I_{bc} , is very small and the $I(i_b)$ curve is almost vertical in the cutoff region. Thus, the end point of the transition is essentially located on the $i_b = 0$ axis. Hence, it is natural in this case to define the "Transition Off" interval as the time it takes the base current to vary between I_{bj} and $I_{bc} \doteq 0$. The transistor is therefore conducting over almost this entire interval and the same approximations can be made here as for the "Transition On" interval. Accordingly, the equations governing the responses in the two intervals will be alike, the only differences being that the trigger signal is absent and the magnetizing current is

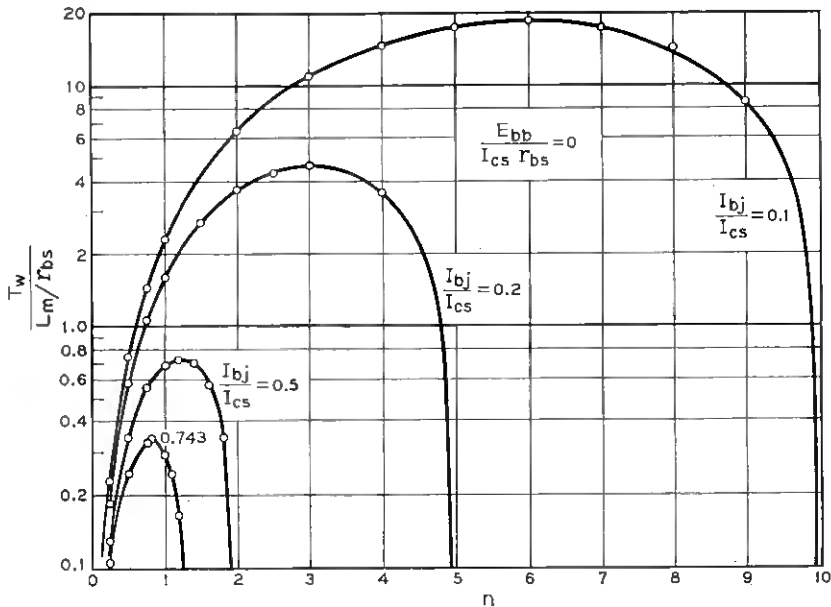


Fig. 21 — Pulse width as a function of turns ratio, with $E_{bb}/I_{cs}r_{bs} = 0$.

$$I_m(t_j) = I_c(I_{bj}) - n \left\{ \left[1 + \frac{R_b}{n^2 L_m} \left(\frac{1}{\omega_0} + R_c C_0 \right) \right] \frac{V_b(I_{bj}) - E_{bb}}{R_b} + I_{bj} \right\} \quad (67)$$

$$\doteq I_c(I_{bj}) - n I_{bj}$$

or, in normalized form,

$$m(\tau_j) = f(x_j) - \left(1 + \frac{1+b}{c} \right) V(x_j) - x_j \doteq f(x_j) - x_j. \quad (68)$$

As explained in the section on the equivalent circuit, a different ω_0 that takes into account the effect of minority carrier storage must now be used. This latter modification usually results in a fall time that is from two to ten times that of the rise time.

The fall time in this case may be calculated in an alternate way which permits one to use the results in Section IV. The accumulation of current in the magnetizing inductance acts like a slowly varying bias that slides the load line toward the point of tangency, at which point rapid regeneration takes place. This fact suggests that the mode of operation now is very similar to the case in which a slow trigger signal is applied to the circuit. Hence, if one defines the fall time as the time taken for the response to traverse 90 to 10 per cent of the interval between the stable and quasistable operating points, one should expect that the fall time can be obtained directly from Fig. 15(a) or (46), provided the ω_0 valid in this interval is used. Indeed, comparisons of rise and fall times between Figs. 15(a) and 19 bear this out.

8.2 External "Turn Off"

When the blocking oscillator is turned off by an external signal the natural pulse width is always selected to be much larger than the interval between the "Turn On" and the "Turn Off" pulses, in order to secure reliable operation. The magnetizing inductance will therefore not become significantly charged during the time the circuit is on and conditions similar to those during the "Transition On" interval exist. By using the value ω_0 , which takes into account the effect of minority carrier storage, the fall time in this case may be calculated from the results in Section IV.

IX. THE "RECOVERY" INTERVAL

This interval is defined as the difference in time between cutoff and the time the circuit comes to rest at the stable operating point. During this interval, assuming that the recovery is overdamped or slightly

oscillatory, it is reasonable to assume that the transistor is essentially cut off. This means that the β of the transistor is zero, $f(x)$ is essentially zero and C_0 reduces to

$$C_0 = C_c + C_s . \tag{69}$$

In addition, the terms involving the normalized base voltage will be considerably larger than those containing the normalized base current in the complete equation governing the circuit operation. Since the terms containing derivatives of the normalized base voltage are unimportant in the preceding operating intervals, their coefficients can be modified in terms of the values they assume during recovery without affecting performance in the other intervals. In this case, we do not operate on the base voltage terms with the expression

$$\left(1 + \frac{1}{\omega_0} \frac{d}{dt}\right) i_c = I_c(i_b), \tag{70}$$

since the transistor, except for its output capacity C_c , is effectively disconnected from the output circuit while it rings out. In addition, the parameter e which arises from the capacity across the transformer is lumped in with the constant d , to give a value of d modified to

$$d_1 = \frac{\omega_0 R_{cb}}{n^2} \left[(C_c + C_s) \left(1 + \frac{L_e}{L_m}\right) + C_\tau \right]. \tag{71}$$

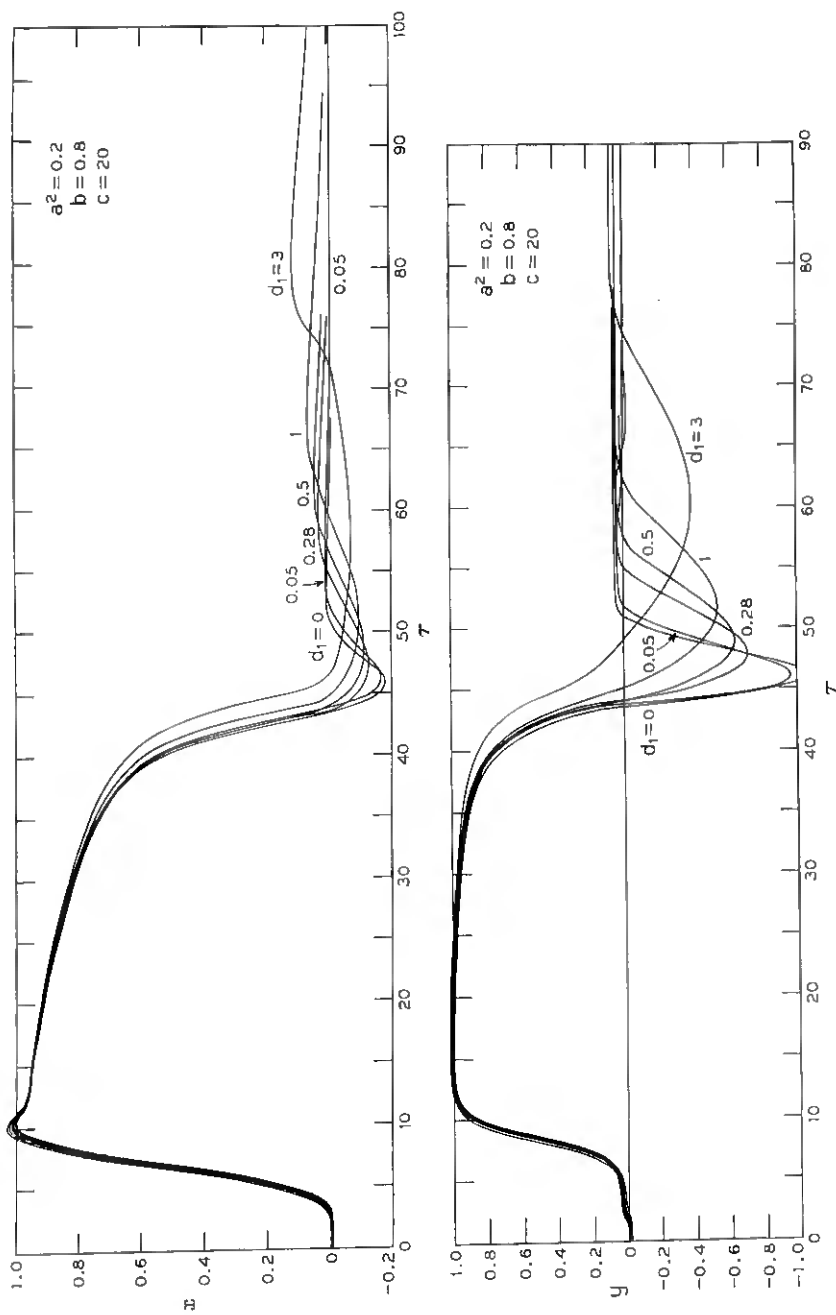
Finally, derivatives of v_b higher than the second were neglected. With this reasoning, the terms that remain involving $V(x)$ in the normalized equation of the blocking oscillator are

$$(a_1^2 + be) \frac{d^2 V}{d\tau^2} + (b + d_1) \frac{dV}{d\tau} + \left(1 + \frac{b}{c}\right) V + \frac{1}{c} \int_{\tau_i}^{\tau} V d\tau, \tag{72}$$

where

$$\left. \begin{aligned} a_1 &= a \\ b_1 &= b \end{aligned} \right\} \text{with } C_0 = C_c + C_s . \tag{73}$$

With this change, the equation governing the performance of the blocking oscillator over its complete cycle of operation was simulated on an analog computer to verify the contention that the transistor is cut off during recovery, and to give quantitative substance to the assumption that only the $V(x)$ terms need be considered. Again, the analytical function was used to represent the i_c - i_b characteristic, while the normalized base voltage function was approximated by two straight line segments.

Fig. 22 — Response of the blocking oscillator for various values of d_1 .

The circuit was turned on by a Gaussian trigger pulse of width $\tau_w = 4.175$ and magnitude $G = 0.5$, and left to its own devices to turn off. The computer results when the constant d_1 was varied and the remaining constants fixed are shown in Fig. 22. The main effect of varying d_1 is to change the response after the blocking oscillator has turned off. Recovery is rather slow, and decreases as d is increased. It should be noted that, as long as d is not too large, the overshoot of the base current is small, confirming the assumption that the transistor remains essentially cut off during recovery. The undershoot of the output voltage is rather large, becoming larger and narrower as d_1 is increased. This is largely due to the fact that the transformer capacity has been assumed negligible compared to C_0 (i.e., $e \neq 0$). When e is greater than C_0 , it will be found that the undershoot of the output voltage is quite small. This will be discussed more fully in a succeeding paragraph. From Fig. 19 it can be seen that recovery is very much dependent on the parameter c . This would be expected physically since c is associated with the energy stored in the magnetizing inductance.

From the above results it can be seen that the transistor is essentially cut off during recovery and that the base-emitter junction of the transistor is back-biased, assuring that $V(x) \gg x$. Therefore, all the assumptions of the beginning of this section are valid and the equation governing recovery can be written from the Appendix and (72) as

$$(a_1^2 + b_1e) \frac{d^2V}{d\tau^2} + (b_1 + d + e) \frac{dV}{d\tau} + \left(1 + \frac{b_1}{c}\right)V + \frac{1}{c} \int_{\tau_j}^{\tau} V d\tau + m(\tau_j) = \left(1 + b_1 \frac{d}{d\tau} + a_1^2 \frac{d^2}{d\tau^2}\right)g(\tau). \tag{74}$$

If we differentiate (74) we get, after some rearranging,

$$c(a_1^2 + b_1e) \frac{d^3V}{d\tau^3} + c(b_1 + d + e) \frac{d^2V}{d\tau^2} + (b_1 + c) \frac{dV}{d\tau} + V = c \left(\frac{d}{d\tau} + b_1 \frac{d^2}{d\tau^2} + a^2 \frac{d^3}{d\tau^3}\right)g(\tau). \tag{75}$$

Whether the response of (75) is monotonic or not depends upon the initial conditions when the recovery interval is entered, as well as the parameters a through e . In the absence of information concerning the initial conditions, we can only determine whether or not the response will be nonoscillatory. This can be determined by using the discriminant for the cubic on the homogeneous equation. Instead, we assume that the third derivative term can be neglected and deal with the second-or-

der equation in V . This approximation appears to be justified in practice. In this case, it is readily shown that the recovery response will be non-oscillatory if

$$d + e \leq \frac{(b + c)^2}{4c} - b = \frac{(c - b)^2}{4c} \doteq \frac{c}{4}. \quad (76)$$

In unnormalized form, this becomes

$$R_b \leq \frac{n^2 L_m}{R_c(C_c + C_s) + 2 \sqrt{L_m C_T + L_m(C_c + C_s) \left(1 + \frac{L_c}{L_m}\right)}} \quad (77)$$

$$\doteq \frac{n^2}{2} \left(\frac{L_m}{C_c + C_s + C_T} \right)^{1/2}.$$

During this interval, the output voltage is given by

$$y = x + e \frac{dV}{d\tau} + \frac{1}{c} \int_{\tau_j}^{\tau} V d\tau + V + m(\tau_j) - g(\tau). \quad (78)$$

When the transformer capacity (directly proportional to e) is much larger than $C_c + C_s$, and the derivative terms of $g(\tau)$ and the highest derivative of V are neglected, (74) becomes

$$e \frac{dV}{d\tau} + V + \frac{1}{c} \int_{\tau_j}^{\tau} V d\tau + m(\tau_j) - g(\tau) \doteq 0. \quad (79)$$

Comparing (79) and (78), it can be seen that they are identical when $x = 0 = y$. During recovery, x is small compared to V . The fact that $y \doteq 0$ implies that the output voltage is essentially zero during recovery when the transformer capacity $C_T \gg C_c + C_s$. This must be so from physical considerations, since the current through the load must of necessity be small when the output impedance of the transistor becomes quite high.

X. MULTIPULSE RESPONSE OF THE BLOCKING OSCILLATOR

To investigate how the blocking oscillator should be designed in order to work reliably as a regenerator in a PCM repeater, (116) of Appendix B was solved on an analog computer with $g(\tau)$ consisting of a sequence of alternating positive and negative Gaussian pulses. Each "Turn On" and "Turn Off" trigger pulse was spaced 15 time units apart, and the magnitude and width of these were $G = 0.5$ and $\tau_w = 4.175$, respectively. The analytical expression (37) was again used to represent the i_c/i_b characteristic and $V(x)$ was approximated by two straight line segments. The constants a and b were selected such that the rise time of the circuit was small

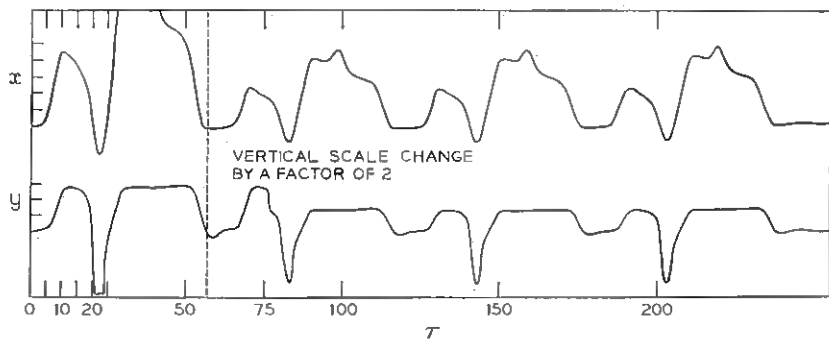


Fig. 23 — Response of the blocking oscillator to a sequence of “Turn On” and “Turn Off” pulses when $a^2 = 0.2$, $b = 0.8$, $c = 5$ and $d = 0.28$. This value of c corresponds to a natural pulse width of $\tau_w = 10$. Each “Turn On” and “Turn Off” pulse was spaced 15 units apart and had a magnitude and width of $G = 0.5$ and $\tau_w = 4.175$. The maximum of the first “On” pulse occurs when $\tau = 7.5$.

compared to the spacing between the trigger pulses and d was picked such that the recovery was reasonably fast. Finally, it was assumed that the blocking oscillator was initially at rest at the stable operating point so that $m(0) = 0$.

The results are shown in Figs. 23 through 26 for various values of the constant c ; that is, the natural pulse width of the circuit was altered in these cases. In Fig. 23 the pulse width was equal to 10, and therefore it was less than half the period of the trigger signal. It is seen that, not only do the height and width of the output pulses vary considerably, but the timing information contained in the trigger signal is almost completely destroyed as well. In Fig. 24, in which the natural width

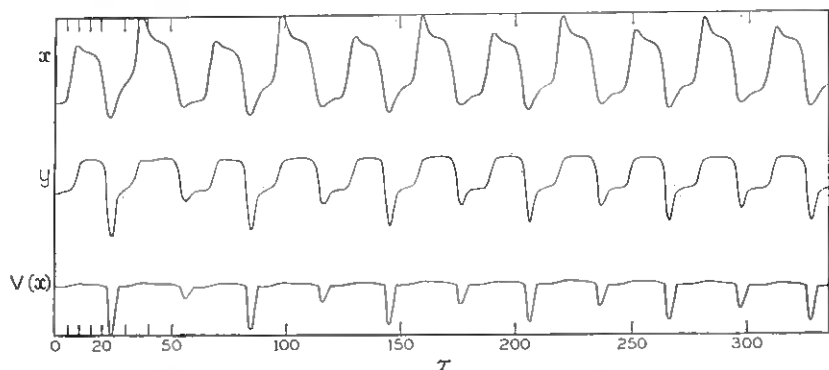


Fig. 24 — Response of the blocking oscillator to a sequence of “Turn On” and “Turn Off” pulses when $c = 10$.

was almost doubled, conditions have improved considerably. However, variations in width and timing are still intolerable. By increasing the width further by almost 40 per cent, as shown in Fig. 25, the result becomes acceptable, although the undershoot of y varies considerably. Finally, in Fig. 26, where the natural pulse width is about 1.4 times the period of the trigger on signal, the train of output pulses is almost completely uniform, with scarcely any change from one pulse to another. Hence, for the blocking oscillator to perform reliably as a regenerator, its natural pulse width should be at least equal to 1.5 times the repetition period or equivalent, three times the time interval, T_T , between each "Turn On" and "Turn Off" pulse for a 50 per cent duty cycle; that is,

$$T_w \cong 3T_T \quad (80)$$

In addition, the rise time should, of course, be considerably less than T_T and the recovery should be such that it is overdamped or only slightly oscillatory.

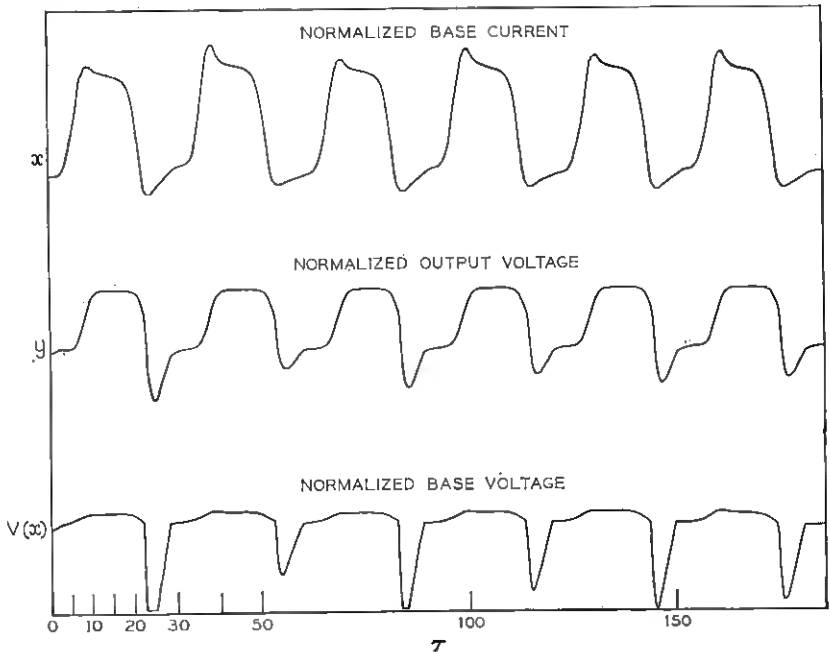


Fig. 25 — Response of the blocking oscillator to a sequence of "Turn On" and "Turn Off" pulses when $c = 14.3$.

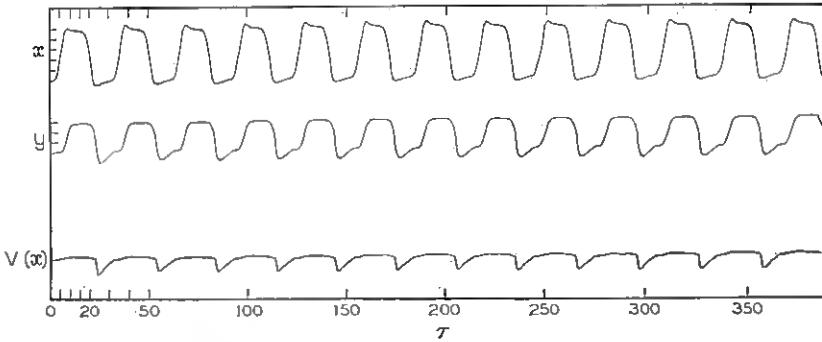


Fig. 26 — Response of the blocking oscillator to a sequence of 13 “Turn On” and “Turn Off” pulses when $c = 25$. This value of c corresponds to a natural pulse width for the circuit of $\tau_w = 41.5$.

Physically, the above results are what one should qualitatively expect. When the natural width is about equal to or less than T_T , the turnoff mechanism in the blocking oscillator itself will interfere with the trigger signal, while in the case of a large width the circuit works more or less like a bi-stable device, thus reproducing the incoming signal faithfully.

XI. DESIGN EXAMPLE

To illustrate how the results of the foregoing discussion may be used to develop a design procedure for the blocking oscillator, let us consider the following typical problem.

A grounded emitter blocking oscillator is to be used as a pulse regenerator working at a maximum repetition rate of 1.5 mc. It should be triggered “on” and “off” by a circuit having an output impedance of $R_b = 1000$ ohms and it must work into a load of $R_c = 400$ ohms. The available bias voltages for the base and collector circuits are $E_{bb} = -0.8$ volt and $E_{cc} = -6$ volts, respectively, and the stray and wiring capacity was estimated to be about 5 micromicrofarads. Available transformers for this application have parameters referred to the primary in the neighborhood of the following:

Leakage inductance: $L_e \doteq 1 \mu\text{h}$,

Magnetizing inductance: $L_m \doteq 100 \mu\text{h}$,

Stray capacitance: $C_T \doteq 20 \mu\mu\text{f}$,

Possible turns ratios: $n = 1.5, 3, 4.5$.

A surface barrier transistor, 2N128, is to be used, for which typical

$I_c(i_b)$ and $[V_b(i_b) - E_{bb}]/R_b$ characteristics are shown in Fig. 27(a). Its other specifications are as follows:

Collector saturation current with $R_c = 400$ -ohm load: $I_{cs} = 14.5$ ma,

Minimum base current to saturate transistor: $I_{bs \text{ min}} = 2$ ma,

Base voltage at the point of saturation: $\doteq 0.5$ volt,

Base input impedance in the conduction region: $r_{bs} = 60$ ohms,

α cutoff frequency: $f_\alpha = 60$ mc,

Collector capacitance: $C_c = 1.6 \mu\mu\text{f}$,

$\beta_{\text{max}} = 37$.

The "Turn On" and "Turn Off" trigger pulses have roughly a Gaussian shape and are spaced 0.33 microsecond apart. Their 50 per cent width is about 0.2 microsecond and their magnitude should be determined. The rise time and fall time of the blocking oscillator should be less than 0.15 microsecond.

The first step to be taken in designing the blocking oscillator is to determine what turns ratio should be used. However, the turns ratio affects a number of things, such as the location of operating points, natural pulse width, etc. in the circuit, so that no single criterion is sufficient to determine n . The range of values imposed upon n by the various requirements in the circuit must be determined first, then a turns ratio can be chosen compatible with all these requirements:

i. The first of these requirements comes from the fact that the turns ratio must be such that the load line intersects the $I(i_b)$ curve in three points. This restricts n to values:

$$n < \frac{I_{cs}}{I_{bs \text{ min}} - I_{bc}} \doteq \frac{I_{cs}}{I_{bs \text{ min}}} = \frac{14.5}{2} = 7.25, \quad (81)$$

where $I_{bs \text{ min}}$ is the minimum base current that will saturate the transistor.

ii. The second requirement pertains to the desirability that n is such that the magnetizing inductance does not significantly affect the location of the quasistable operating point. According to inequality (27), this requires that:

$$n^2 - \frac{I_{cs}}{\frac{V_{bs} - E_{bb}}{R_b}} n + \frac{R_b}{L_m} \left(\frac{1}{\omega_0} + R_c C_0 \right) \ll 0. \quad (82)$$

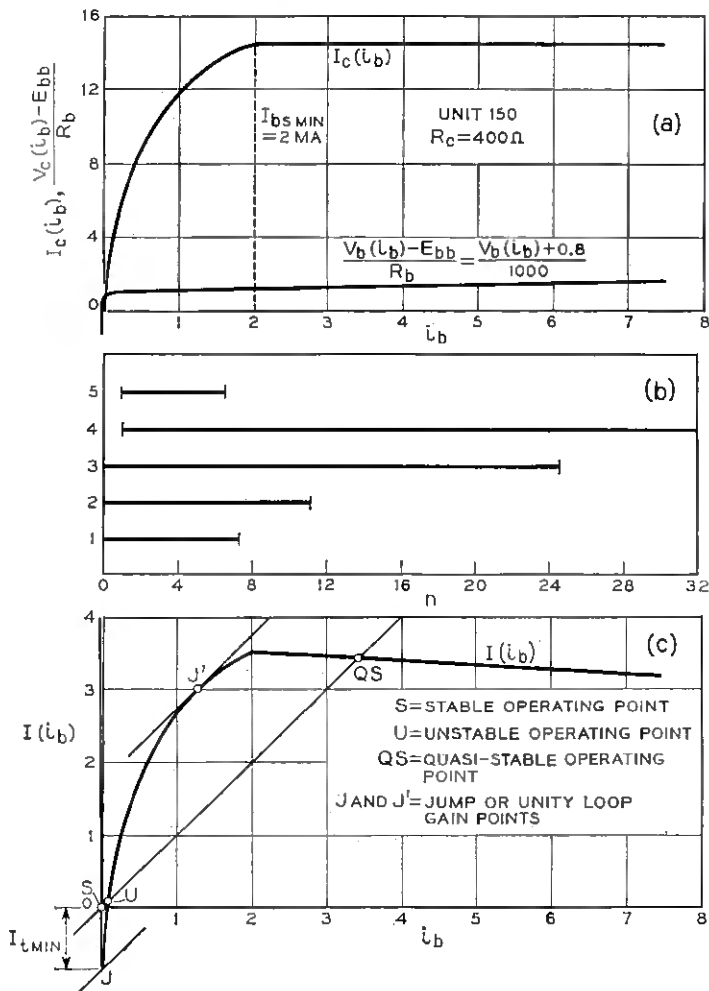


Fig. 27 — Design example: (a) $[V_b(i_b) - E_{bb}]/R_b$ and $I_c(i_b)$ characteristics of 2N128 transistor; (b) range of values imposed upon n ; (c) $I(i_b)$ for $n = 3$.

Substituting for ω_0 and C_0 from (28) and (15), respectively, this condition may also be written in the form:

$$n^2 - R_b \left[\frac{I_{cs}}{V_{bs} - E_{bb}} - \frac{1}{L_m} \left(\frac{1}{\omega_\alpha} + 1.5R_c C_c \right) \right] n + \frac{R_b}{I_m} \left[\frac{1}{\omega_\alpha} + R_c(C_s + 1.5C_c) \right] \ll 0. \tag{83}$$

However, in most cases of practical interest,

$$\frac{1}{L_m} \left[\frac{1}{\omega_\alpha} + R_c(C_s + 1.5C_c) \right] \ll \frac{I_{cs}}{V_{bs} - E_{bb}}. \quad (84)$$

Thus, inequality (83) reduces to

$$0 < n < \frac{I_{cs}}{\frac{V_{bs} - E_{bb}}{R_b}}. \quad (85)$$

In our numerical example, condition (84) leads to $5.9 \times 10^{-3} \ll 1.11$. So as far as this requirement is concerned, n should be selected within the range:

$$0 < n < 11.1. \quad (86)$$

iii. In order to assure that the rise time is reasonably independent of β_{\max} of the transistor, n should, according to Fig. 13, be selected such that β_{\max}/n is not too close to unity, say,

$$n < \frac{\beta_{\max}}{1.5} = \frac{37}{1.5} = 24.6. \quad (87)$$

iv. The fourth requirement on n concerns itself with the desirability of a non-oscillatory recovery, which, according to (77), requires that:

$$n > \sqrt{\frac{2R_b}{\frac{L_m}{C_c + C_s + C_t}}} = \sqrt{\frac{2 \times 10^3}{\frac{10^{-4}}{3 \times 10^{-11}}}} = 1.05. \quad (88)$$

v. Finally, in the section on the multipulse response of the blocking oscillator it was found that the natural width should be equal to or larger than three times the time interval, T_r , between the "Turn On" and "Turn Off" pulses in order for the circuit to reproduce the incoming information faithfully. According to (80), this requires that n must be such that:

$$\frac{n^2 L_m}{r_{bs}} \ln \frac{\frac{1}{n} - \frac{E_{bb}}{I_{cs} r_{bs}}}{\frac{I_{bj}}{I_{cs}} - \frac{E_{bb}}{I_{cs} r_{bs}}} \geq 3T_r. \quad (89)$$

From Figs. 20 or 21 it is seen that this condition will limit n to a range of values between a minimum and a maximum. All quantities are known in the above equation except I_{bj} , which, if n was known, could be determined from Fig. 27(a). However, since the lower and upper limit of n

defined by the above condition approach each other when I_{bj} is increased, and since I_{bj} increases as n is decreased, we only have to consider the case where I_{bj} is determined by the lowest possible value of n . Examining the previous restrictions on n , we see that its smallest possible value is 1.05, which determines the point of tangency in Fig. 27(a) to be $I_{bj} = 2$. Hence, condition (89) becomes

$$n^2 \ln \frac{\frac{1}{n} + 1}{1.138} \geq 0.564, \tag{90}$$

which gives:

$$0.996 < n < 6.53. \tag{91}$$

The range of values imposed upon n by the above requirements is plotted in Fig. 27(b). It is seen that any value of n between 1.05 and 6.53 satisfies all these requirements. However, remembering that the turns ratio should be selected small in order to reduce the rise time as much as possible, and taking into account that, for reasons of reliability, n should not be chosen close to the edges of the allowable range, the value $n = 3$, seems to be a good compromise. With this value of n the resulting $I(i_b)$ curve is as shown in Fig. 27(c), from which it is seen that the base current at the quasistable operating point is $I_{bs} = 3.42$ ma. Hence, the various circuit parameters become:

$$\omega_0 \sim \frac{\omega_a}{1 + \frac{I_{cs}}{I_{bs}}} = \frac{2\pi 6.08 \times 10^7}{1 + \frac{14.5}{3.42}} = 7.3 \times 10^7, \tag{92}$$

$$C_0 \sim C_s + 1.5C_c \left(1 + \frac{I_{cs}}{I_{bs}}\right) = 5 + 1.5(1.6) \left(1 + \frac{14.5}{3.42}\right) = 17.6 \mu\mu f,$$

$$a = \omega_0 \sqrt{L_c C_0} = 7.3 \times 10^7 \sqrt{10^{-6}(17.8 \times 10^{-12})} = 0.31, \tag{93}$$

$$b = \omega_0 R_c C_0 = (7.3 \times 10^7)(1.76 \times 10^{-11}) = 0.52,$$

$$\frac{b}{2a} = 0.84.$$

The normalized trigger pulse width is

$$\tau_w = \omega_0 T_w = (7.3 \times 10^7)(2 \times 10^{-7}) = 14.6. \tag{94}$$

From Fig. 15(c) it is seen that with these values of a and b the normalized rise time is about 4.6, which means that the real rise time is:

$$T_R = \frac{4.6}{7.3 \times 10^7} = 0.063 \mu\text{s}. \quad (95)$$

Also, from Fig. 16(b) it can be concluded that the overshoot is completely negligible in this case. This is, of course, to be expected, since the ratio $b/2a$ is close to unity.

By separate measurements on the transistor it was determined that the equivalent ω_0 during turnoff was equal to 3×10^7 , giving, as before, $a = 0.13$ and $b = 0.21$. According to Fig. 15(c), the fall time is therefore

$$T_F = \frac{3.8}{3.10^7} = 0.127 \mu\text{s}. \quad (96)$$

Hence, both the rise and fall times are well within the specified limit. If this had not been the case, it would have been necessary to repeat the above procedure with a different trigger source, transformer or load, or with a different transistor.

From Fig. 27(c) it is seen that the minimum trigger magnitude in the case of a very slow trigger is $I_{t \min} = 0.85$ and, by interpolation between Figs. 7 and 8, the normalized magnitude for a width of $\tau_w \sim 15$ is $G = 0.175$. Hence, the magnitude of the applied trigger should be larger than:

$$\begin{aligned} I_t &> (I_{b_s} - I_{b_c})(G - G_{\min}) + I_{t \min} \\ &\sim 3.42(0.175 - 0.150) + 0.85 \\ &\sim 0.95 \text{ ma}. \end{aligned} \quad (97)$$

If the blocking oscillator had only been turned "on" externally but left to its own devices afterwards, its natural pulse width would be:

$$T_w = \frac{3^2 \times 10^{-4}}{56.4} \ln \frac{\frac{1}{3} + 1}{\frac{1.25}{14.5} + 1} = 3.35 \mu\text{s}, \quad (98)$$

where I_{b_j} was found from Fig. 27(c) to be equal to 1.25 ma.

The undershoot of the output voltage will be very small since the transformer capacitance is much larger than the output capacitance of the transistor during the recovery interval.

To check how well the above calculations and the rise times predicted by Fig. 15 agreed with experimental results, a blocking oscillator having the same specifications as the one above was built and measured. In the case of the rise times, trigger pulses having magnitudes and widths corresponding exactly to those of Fig. 15 were applied to the circuit and

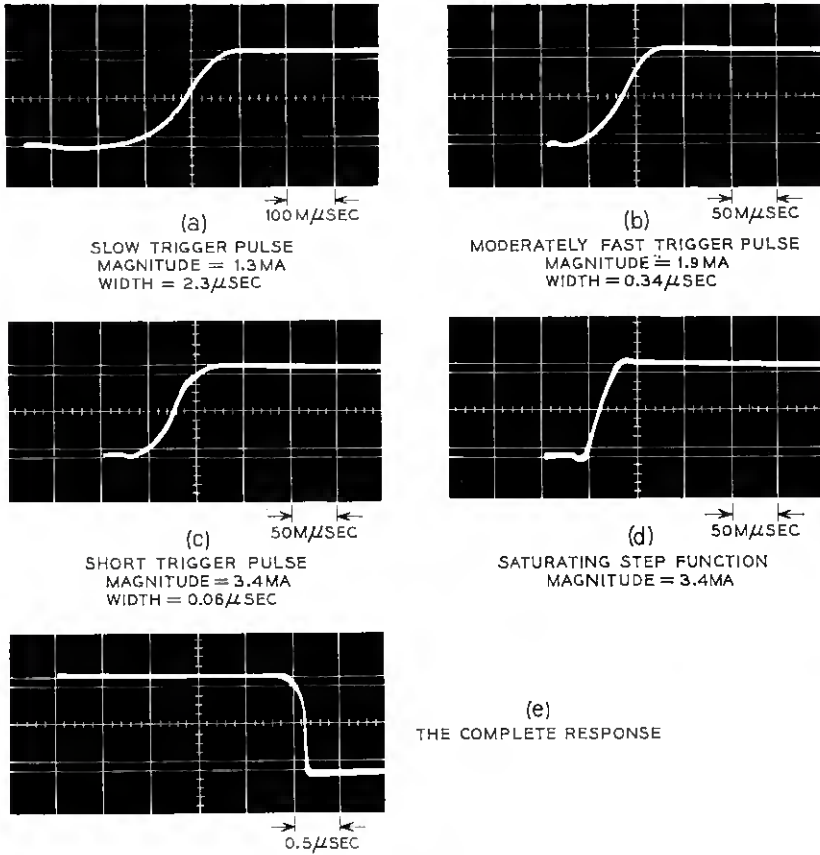


Fig. 28 — (a) through (d) output voltage of the blocking oscillator in response to various types of trigger signals; (e) the complete response.

then compared to the values predicted by these figures. The results are shown in Figs. 28(a) through (d) and in Table III. In all cases the overshoot was not observable, which is what one should expect from Fig. 16. Also, the critical triggering level in the case when a 0.2-microsecond pulse was applied was found to be equal to $I_{l \text{ min}} = 0.9 \text{ ma}$. Fig. 28(e) shows the complete response of the blocking oscillator, from which it is seen that its natural pulse width was about 2.5 microseconds and the undershoot less than 5 per cent.

XII. DISCUSSION OF RESULTS

An approach to the design of transistor blocking oscillators with emphasis on the nonlinearities inherent in their operation has been pre-

TABLE III

Trigger Signal	Actual Magnitude, ma	Actual Width, μ s	Rise Time, T_R , in μ s	
			Calculated	Measured
Slow Trigger Pulse $G = 0.3$, $\tau_w = 166.7$ [Fig. 14(a)]	1.3	2.3	$\frac{11.5}{7.3 \times 10^7} = 0.157$	0.15
Moderately Fast Trigger Pulse $G = 0.5$, $\tau_w = 25$ [Fig. 14(c)]	1.9	0.34	$\frac{4.6}{7.3 \times 10^7} = 0.063$	0.06
Short Trigger Pulse $G = 1$, $\tau_w = 4.175$ [Fig. 14(e)]	3.4	0.06	$\frac{4.65}{7.3 \times 10^7} = 0.064$	0.06
Saturating Step Function $G = 1$ [Fig. 14(f)]	3.4	—	$\frac{2.3}{7.3 \times 10^7} = 0.031$	0.03

sented. A graphical picture of the operating points of the circuit has been drawn which gives a useful feel for the circuit behavior. In addition, both the minimum trigger signal needed to cause the circuit to regenerate, and the requirements for reliable triggering can be inferred from the static characteristic (steady-state equation). An approximate relationship between trigger magnitude and width was determined from analog computations. This relationship, combined with other analysis, shows that the minimum energy for reliable triggering is attained when the transistor has a high-low current β , large ω_c , small collector capacity, and low forward drop (a good switch) and the transformer is of high quality, i.e., has large magnetizing inductance and high coupling coefficients. This is hardly surprising, but it satisfies our intuition. Most importantly, an indication is given as to how these factors affect trigger sensitivity.

The rise time of the circuit is relatively independent of the maximum loop gain β_{\max}/n , and therefore also of the details of the nonlinear i_c-i_b characteristic for trigger signals about 50 per cent or more larger than the minimum. On the other hand, the rise time is very much dependent upon the rise time of the trigger pulse. Slowly varying triggers give, as expected, the poorest rise time, while the minimum rise time is achieved with fast trigger signals that saturate the transistor immediately. The implications of this result are obvious: power gain and considerable pulse reshaping, as required in a regenerative repeater for PCM, can be purchased at the expense of a rise time significantly slower than that attainable with a sharp trigger signal. This follows from the fact that the positive feedback loop, which is active over a greater portion of the transition on interval with slow trigger signals, has a deleterious effect on rise time. The results of the rise time calculations and the conditions for

monotonic response permit the designer to establish specifications on the circuit to obtain the desired leading edge of the output pulse with small overshoot. The analysis of the "On" interval shows clearly the instabilities in pulse width that result when this width is determined by the natural pulse width of the circuit. It can be seen from Figs. 20 and 21 that the natural pulse width is critically dependent upon the circuit operating points, the nonlinearities and the bias. This is further emphasized by the fact that the largest discrepancy between measured and computed results in the example occurs in the natural pulse width. When the pulse width must be accurately controlled, it is highly desirable to have this accomplished by an external trigger-off pulse. This implies a large magnetizing inductance for the feedback transformer and justifies the initial approximations that were made in the derivation of the equation governing the transition on interval.

During recovery the energy stored in the magnetizing inductance rings out. Quick recovery can be accomplished without retriggering if the transformer circuit is critically damped. This may be accomplished with the existing circuit elements, or it may be necessary to add a damping resistor and diode combination across the transformer.⁴ This has the effect of modifying the value of R_b used in the recovery interval. When the circuit is clocked off, analog computer results show that the natural pulse width should be at least three times the desired pulse width for reliable performance.

Several points regarding this analysis and its application should be clarified before concluding. As pointed out, the equivalent circuit chosen for the transistor is a necessity for a tractable analysis. Gross deviations from the assumed configuration may partially invalidate some of the results. In most cases, physical reasoning can be used to account for the effects of these deviations. Furthermore, this equivalent circuit is a good one for transistor designs to shoot for.

Second, the design example clearly shows that the analysis presented does not lead to a set of design equations that can be solved immediately for the values of all the circuit elements. In particular, the parameter n , the turns ratio, is woven into all of the normalized design parameters, as evidenced by the design example. Furthermore, both the magnetizing inductance of the transformer and its parasitic capacity are dependent on n . The fact that there is no simple, immediate means of arriving at a unique set of parameters is due to this interdependence. This is typical of most engineering problems. However, this analysis does bring out the compromises required in the design.

Finally, we believe that the ingredients employed in the present ap-

proach are useful in the successful execution of nonlinear circuit design. The recipe is as follows: apply some experiment; some intuition; some analysis, including minor deviations from superposition, and, finally, a healthy sprinkling of planned computing.

XIII. ACKNOWLEDGMENTS

During the development of this approach to the design of transistor blocking oscillators several stimulating discussions were held with the authors' colleagues. In particular, discussions with F. T. Andrews and R. C. Chapman were most helpful in clarifying some of our concepts. The extensive use of both analog and digital computers is a tribute to the skill and speed with which Mrs. W. L. Mammel, Mrs. G. J. Hansen and Miss E. G. Cheatham coded the problems and nursed them through the computers. S. H. Rothrock made the measurements of the transistor parameters.

APPENDIX A

Derivation of the Differential Equation

From the equivalent circuit in Fig. 2(b) the equations governing blocking oscillator behavior are:

$$i_i = I_i(t), \quad (99)$$

$$v_b = V_b(i_b), \quad (100)$$

$$\left(1 + \frac{1}{\omega_0} \frac{d}{dt}\right) i_c = I_c(i_b), \quad (101)$$

$$i_p = ni_s = n \left(i_b + \frac{v_b - E_{bb}}{R_b} - i_i \right), \quad (102)$$

$$i_i = i_p + i_m + i_d = i_p + \frac{1}{L_m} \int_{t_f} v_b - \frac{E_{bb}}{n} dt + I_m(t_j) + \frac{C_T}{n} \frac{dv_b}{dt}, \quad (103)$$

$$i_c = i_i - i_g = i_i - C_0 \frac{dv_c}{dt}, \quad (104)$$

$$v_c = E_{cc} - i_l R_c - \frac{v_b - E_{bb}}{n} - L_c \frac{di_l}{dt}. \quad (105)$$

Substituting for v_c from (105) in (104):

$$i_c = i_i + C_0 R_c \frac{di_l}{dt} + \frac{C_0}{n} \frac{dv_b}{dt} + L_c C_0 \frac{d^2 i_l}{dt^2}, \quad (106)$$

and combining 102 and 103:

$$i_i = ni_b + n \frac{v_b - E_{bb}}{R_b} + \frac{1}{L_m} \int_{t_j}^t \frac{v_b - E_{bb}}{n} dt + I_m(t_j) + \frac{C_T}{n} \frac{dv_b}{dt} - i_i. \tag{107}$$

Substituting for i_i from (107) into (106) and combining terms yields:

$$\begin{aligned} i_c = & L_c C_0 n \frac{d^2 i_b}{dt^2} + C_0 R_c n \frac{di_b}{dt} + ni_b + \frac{L_c C_0 C_T}{n} \frac{d^3 v_b}{dt^3} \\ & + \frac{L_c C_0 n}{R_b} + \frac{C_0 C_T R_c}{n} \frac{d^2 v_b}{dt^2} + \frac{C_0 R_c n}{R_b} + \frac{C_0 + C_T}{n} + \frac{L_c C_0}{L_m n} \frac{dv_b}{dt} \\ & + \frac{n}{R_b} + \frac{C_0 R_c}{L_m n} (v_b - E_{bb}) + \frac{1}{L_m} \int_{t_j}^t \frac{v_b - E_{bb}}{n} dt + I_m(t_j) \\ & - n \left(1 + C_0 R_c \frac{d}{dt} + C_0 L_c \frac{d^2}{dt^2} \right) i_i. \end{aligned} \tag{108}$$

Operating upon (108) in accordance with (101), grouping like terms, dividing through by n and rewriting in terms of $[v_b(i_b) - E_{bb}]/R_b$ gives the final form of the defining equation as:

$$\begin{aligned} & \frac{L_c C_0}{\omega_0} \frac{d^3 i_b}{dt^3} + \left(L_c C_0 + \frac{R_c C_0}{\omega_0} \right) \frac{d^2 i_b}{dt^2} + \left(\frac{1}{\omega_0} + R_c C_0 \right) \frac{di_b}{dt} \\ & - \left[I(i_b) - i_b - \frac{I_m(t_j)}{n} \right] + \frac{R_b}{n^2 L_m} \int_{t_j}^t \frac{V_b(i_b) - E_{bb}}{R_b} dt \\ & + \left\{ \frac{C_T R_b}{n^2} + \frac{1}{\omega_0} + C_0 \left[R_c \left(1 + \frac{R_b}{\omega_0 n^2 L_m} \right) + \frac{R_b}{n^2} \left(1 + \frac{L_c}{L_m} \right) \right] \right\} \\ & \cdot \frac{d}{dt} \frac{V_b(i_b) - E_{bb}}{R_b} + \left\{ \frac{C_0 C_T R_c R_b}{n^2} + \frac{R_b C_T}{n^2 \omega_0} + L_c C_0 \right. \\ & \left. + \frac{C_0}{\omega_0} \left[R_c + \frac{R_b}{n^2} \left(1 + \frac{L_c}{L_m} \right) \right] \right\} \frac{d^2}{dt^2} \frac{V_b(i_b) - E_{bb}}{R_b} \\ & + \left(\frac{C_0 C_T R_c R_b}{n^2 \omega_0} + \frac{L_c C_0}{\omega_0} + \frac{L_c C_0 C_T R_b}{n^2} \right) \frac{d^3}{dt^3} \frac{V_b(i_b) - E_{bb}}{R_b} \\ & + \frac{L_c C_0 C_T R_b}{n^2 \omega_0} \frac{d^4}{dt^4} \frac{V_b(i_b) - E_{bb}}{R_b} \\ & = \left[1 + \left(\frac{1}{\omega_0} + R_c C_0 \right) \frac{d}{dt} + \left(L_c C_0 + \frac{R_c C_0}{\omega_0} \right) \frac{d^2}{dt^2} + \frac{L_c C_0}{\omega_0} \frac{d^3}{dt^3} \right] I_i(t). \end{aligned} \tag{109}$$

In the above, $I(i_b)$ is defined as

$$I(i_b) = \frac{I_c(i_b)}{n} \left[1 + \frac{R_b}{n^2 L_m} \left(\frac{1}{\omega_0} + R_c C_0 \right) \right] \frac{V_b(i_b) - E_{bb}}{R_b}. \tag{110}$$

The "steady-state equation" which determines the operating points of the circuit is given by (109) when all derivatives and integrals are zero. It is

$$I(\dot{i}_b) - \dot{i}_b - \frac{I_m(t_j)}{n} = 0. \quad (111)$$

From (107), the output voltage becomes:

$$v_0 = R_c i_t = nR_c \left[\dot{i}_b + \frac{V_b(\dot{i}_b) - E_{bb}}{R_b} + \frac{R_b}{n^2 L_m} \int_{t_j}^t \frac{V_b(\dot{i}_b) - E_{bb}}{R_b} dt \right. \\ \left. + R_b \frac{C_T}{n^2} \frac{d}{dt} \frac{V_b(\dot{i}_b) - E_{bb}}{R_b} + \frac{I_m(t_j)}{n} - I_t(t) \right]. \quad (112)$$

Combining (105) and (112), the collector voltage may be written

$$v_c = E_{cc} = \frac{V_b(\dot{i}_b) - E_{bb}}{n} - \left(1 + \frac{L_m}{R_c} \frac{d}{dt} \right) v_0. \quad (113)$$

APPENDIX B

Normalization

Considerable simplification and savings in space can be achieved when the dependent and independent variables are normalized in the following manner:

$$\begin{aligned} \tau &= \omega_0 t, \\ x &= \frac{\dot{i}_b - I_{bc}}{I_{bs} - I_{bc}} = \text{normalized base current}, \\ V(x) &= \frac{V_b(\dot{i}_b) - E_{bb}}{R_b(I_{bs} - I_{bc})} = \text{normalized base voltage}, \\ g(\tau) &= \frac{I_t(t)}{I_{bs} - I_{bc}} = \text{normalized trigger function}, \\ f(x) &= \frac{I_c(\dot{i}_b)}{n(I_{bs} - I_{bc})} \\ &= \text{normalized collector current base current function}, \\ m(\tau_j) &= \frac{I_m(t_j) - nI_{bc}}{n(I_{bs} - I_{bc})} = \text{normalized initial current in } L_m, \\ h(x) &= \frac{I(\dot{i}_b)}{I_{bs} - I_{bc}} = \text{normalized nonlinear characteristics.} \end{aligned} \quad (114)$$

Further, we define

$$\begin{aligned}
 a^2 &= \omega_0^2 L_e C_0, \\
 b &= \omega_0 R_c C_0, \\
 c &= \omega_0 \frac{n^2 L_m}{R_b}, \\
 d &= \frac{\omega_0 R_b C_0}{n^2} \left(1 + \frac{L_e}{L_m} \right), \\
 e &= \frac{\omega_0 R_b C_T}{n^2}, \\
 k &= \frac{R_b}{n^2 R_c}.
 \end{aligned} \tag{115}$$

Using the above definitions in (109) through (111) of Appendix A gives the normalized defining equations:

$$\begin{aligned}
 a^2 \frac{d^3 x}{d\tau^3} + (a^2 + b) \frac{d^2 x}{d\tau^2} + (1 + b) \frac{dx}{d\tau} + x + m - h(x) \\
 + \frac{1}{c} \int_{\tau_i}^{\tau} V(x) d\tau + \left[1 + b \left(1 + \frac{1}{c} \right) + d + e \right] \frac{dV(x)}{d\tau} \\
 + (a^2 + b + d + e + be) \frac{d^2 V(x)}{d\tau^2} \\
 + (a^2 + be + a^2 e) \frac{d^3 V(x)}{d\tau^3} + a^2 e \frac{d^4 V(x)}{d\tau^4} \\
 = \left[1 + (1 + b) \frac{d}{d\tau} + (a^2 + b) \frac{d^2}{d\tau^2} + a^2 \frac{d^3}{d\tau^3} \right] g(\tau)
 \end{aligned} \tag{116}$$

and

$$h(x) = f(x) - \left(1 + \frac{1 + b}{c} \right) V(x). \tag{117}$$

The normalized “steady-state equation” which determines the operating points of the circuit is given by (22) in Section IV when all derivatives and integrals are zero. It is

$$h(x) - x - m(\tau_i) = f(x) - \left(1 + \frac{1 + b}{c} \right) V(x) - x - m(\tau_i) = 0. \tag{118}$$

In terms of the normalized base current, x , the stable operating point is

now located at $x = 0$ and the quasistable one at $x = 1$, and $f(x)$ attains the same values at these two points.

If the normalized output voltage is defined as

$$y = \frac{v_o}{nR_c(I_{bs} - I_{bc})} = \frac{i_i}{n(I_{bs} - I_{bc})},$$

the normalized version of (112) of Appendix A becomes

$$y = x + V(x) + \frac{1}{c} \int_{\tau_j}^{\tau} V(x) d\tau + e \frac{dV(x)}{d\tau} + m(\tau_j) - g(\tau). \quad (119)$$

In a similar way, the normalized collector voltage is defined as

$$z = \frac{E_{cc} - v_c}{nR_c(I_{bs} - I_{bc})},$$

and (113) of Appendix A becomes

$$z = kV(x) + \left(1 + \frac{a^2}{b} \frac{d}{d\tau}\right) y. \quad (120)$$

REFERENCES

1. Tendick, F. H., Jr., Transistor Pulse Regenerative Amplifiers, B.S.T.J., **35**, September 1956, p. 1085.
2. Wrathall, L. R., Transistorized Binary Pulse Regenerators, B.S.T.J., **35**, September 1956, p. 1059.
3. Linvill, J. G. and Mattson, R. N., Junction Transistor Blocking Oscillators, Proc. I.R.E., **43**, November 1955, p. 1632.
4. Bowers, F. K., unpublished manuscript.
5. Senatorov, K. J. and Guzhov, V. P., On the Analysis of Processes in Transistor Blocking Oscillators, Radio. and Elect. (U.S.S.R.), **2**, 1957, p. 1119.
6. Ebers, J. J. and Moll, J. L., Large-Signal Behavior of Junction Transistors, Proc. I.R.E., **42**, December 1954, p. 1761.
7. Easley, J. W., Effect of Collector Capacity on the Transient Response of Function Transistors, I.R.E. Trans., **ED-4**, January 1957, p. 6.
8. Moll, J. L., Large-Signal Transient Response of Junction Transistors, Proc. I.R.E., **42**, December 1953, p. 1773.
9. Bashkow, T. R., Effect of Nonlinear Collector Capacitance on Collector Current Rise Time, I.R.E. Trans., **ED-3**, October 1956, p. 167.
10. Elmore, W. C. and Sands, M., *Electronics, Experimental Techniques*, McGraw-Hill, New York, 1949, p. 136.
11. Mathers, G. W. C., The Synthesis of Lumped Element Circuits for Optimum Transient Response, Stanford Research Lab. of Electronics Report No. 38, Stanford Univ., Palo Alto, Calif.
12. Narud, J. A., Theory of Nonlinear Feedback Systems Having a Multiple Number of First Order Operating Points and Its Application to Milli-Microsecond Techniques, Atomic Energy Commission Report HEPL-34, February 1955.
13. Aaron, M. R. and Segers, R. G., A Necessary and Sufficient Condition for a Bounded Nondecreasing Step Response, I.R.E. Trans., **CT-5**, September 1958, p. 226.
14. Narud, J. A., The Secondary Emission Pulse Circuit, Its Analysis and Application, Cruft Lab. Tech. Report No. 245, Harvard Univ., Cambridge, Mass., April 5, 1957, p. 6.

Hall Effect Devices

By W. J. GRUBBS

(Manuscript received October 21, 1958)

A wealth of devices which depend on the Hall effect for their operation have been proposed in the last decade. This paper gives the results of a survey of these devices. Original work in this field is included in those sections which describe the circulator, one-piece gyrator, switch, frequency spectrum analyzer, phase discriminator and digital-to-analog encoder. Semiconductor materials are discussed in terms of what type of material is most desirable and how currently available materials limit the usefulness of Hall effect devices.

I. INTRODUCTION

If a magnetic field is applied perpendicular to a current flow in any conductor, the moving charges (which constitute the current) are deflected sidewise and build up a potential difference between the two sides of the conductor. The creation of this transverse electric field (perpendicular to both the magnetic field and the original current flow) is called the Hall effect. During recent years, interest in this effect has increased tremendously. Before semiconductors and their capabilities were understood, the Hall effect in solids was little more than a laboratory curiosity. Now it is not only an important tool in metallurgy and semiconductor device development, but it has been the mode of operation of many proposed devices. This article describes how 20 or so of these devices operate. In each case, the major advantages or disadvantages are mentioned, but no attempt is made actually to determine the usefulness of the device.

The devices to be discussed have been arbitrarily divided into two groups: devices which use a constant magnetic field and devices in which a signal or an oscillator produces at least a part of the magnetic field. Such a division is not entirely arbitrary, because the first group inherently has a very high limit on the operating frequency and the second group has a considerably lower limit. The devices are listed on the following page.

Constant magnetic field

Gyrator
 Isolator
 Negative-resistance amplifier
 Circulator

Signal-produced magnetic field

Switch
 Transducer
 Magnetic field meter
 Electrical compass
 Magnetic field variation meter
 Ammeter
 Wattmeter
 Amplifier
 Modulator (and nondrift dc amplifier)
 Demodulator
 Frequency spectrum analyzer
 Phase discriminator
 Digital-to-analog encoder
 Analog multiplier

II. CONSTANT MAGNETIC FIELD DEVICES

2.1 *Gyrator*

The gyrator has received more attention^{1,2,3,4,5} than has any other Hall effect device — undoubtedly because it was the first nonreciprocal circuit element to which the electrical world was exposed. Casimir⁶ was the first to point out in an English publication that, in a conducting solid, if $R_{ik}(H)$ is the transfer impedance from terminal pair i to pair k in the presence of an orthogonal magnetic field H , then

$$R_{ik}(H) \neq R_{ki}(H),$$

but

$$R_{ik}(H) = R_{ki}(-H). \quad (1)$$

This is simply a statement of the fact that the presence of the magnetic field causes the reciprocity theorem to be violated. With $H = 0$,

$$R_{ik} = R_{ki}.$$

Casimir credited Meixner⁷ with being the first to prove (1).

McMillan⁸ pointed out that many transducers (such as the crystal or condenser types) are reciprocal, but that electrodynamic or magnetic transducers are antireciprocal. An antireciprocal four-pole is a transducer in which the transfer impedance from left-to-right is opposite in sign and equal in magnitude to the transfer from right-to-left. McMillan suggested that by combining the two kinds of transducers (reciprocal and antireciprocal) it should be possible to produce a nonreciprocal transducer — one in which the two transfers have unequal magnitude.

It may be noted that McMillan's antireciprocal four-pole corresponds to at least one definition of the gyrator and his nonreciprocal four-pole bears a resemblance to an isolator. The gyrator may be defined as a four-pole in which the two transfer impedances are equal in magnitude but have phase angles which differ by 180° (i.e., a gyrator is an antireciprocal four-pole).

In a letter to the editor, McMillan⁹ suggested the Hall effect as one

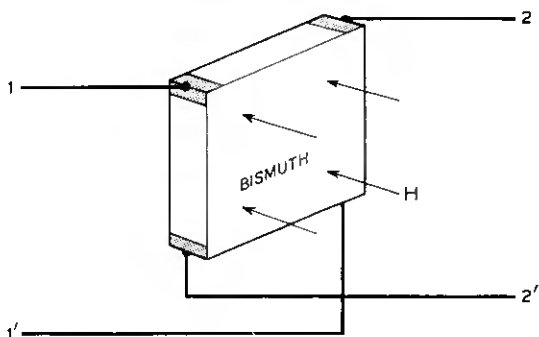


Fig. 1 — McMillan's antireciprocal four-pole (proposed).

possible means of achieving a nonreciprocal electrical system. His scheme was to place a slab of bismuth with contacts in a magnetic field, as shown in Fig. 1. He selected bismuth because it exhibits a very large Hall effect.

Tellegen⁴ referred to McMillan's work and went on to give the gyrator its name,* propose in general terms possible means of making it and suggest many ways in which it would be helpful in electrical network synthesis.

Hogan¹⁰ used the Faraday effect to make a gyrator at microwave frequencies. Mason *et al.*² reported making gyrators which use the Hall effect for their operation, and Wick⁵ showed that the minimum possible power loss of a Hall effect gyrator is 7.66 db.

The type of Hall effect gyrator which has been studied most thoroughly

* The gyrator is the electrical analog of a mechanical gyroscope.

consists of a square slab of semiconductor with an ohmic contact in the middle of each edge face and a constant magnetic field H perpendicular to the plane of the slab (see Fig. 2). There are no junctions and the reader may readily convince himself that such a device is antireciprocal.

If the Hall angle is small, it is given by

$$\theta_H = \tan^{-1} (\mu_H H \times 10^{-8}), \quad (2)$$

where μ_H is the Hall mobility of the semiconductor's majority carrier in $\text{cm}^2/\text{volt-second}$ and H is the magnetic field intensity in oersteds. (This approximation is valid as long as the product of the majority carrier mobility and density is much greater than the same product for minority carriers.) It can be shown (again, for small θ_H) that, in the Hall effect gyrator with the output leads open-circuited, the output voltage is

$$V_{\text{out}} = V_{\text{in}}(\mu_H H \times 10^{-8}). \quad (3)$$

An n-type germanium gyrator has a loss of about 14 db with $H = 17,500$. If a material with higher mobility were used, the same loss would be observed with a smaller field or a smaller loss with the same field. Hogan¹⁰ refers to a vacuum tube Hall effect gyrator proposed by R. O. Grisdale which had four electrodes that could both emit and collect electrons. Such a device could have very high electron mobility. It was built and found to have a loss of "about 7 db". Apparently $\tan \theta_H$ became so very high that the minimum possible loss was approached. Of course, this is more complicated than most users prefer, so the device was not developed for practical use.

It is worth mentioning that a semiconductor Hall effect gyrator operates at direct current and, theoretically, at any frequency up to the dielectric relaxation frequency of the material used.

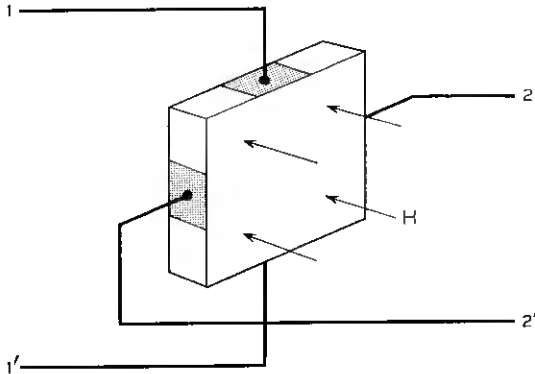


Fig. 2 — Hall effect gyrator (built).

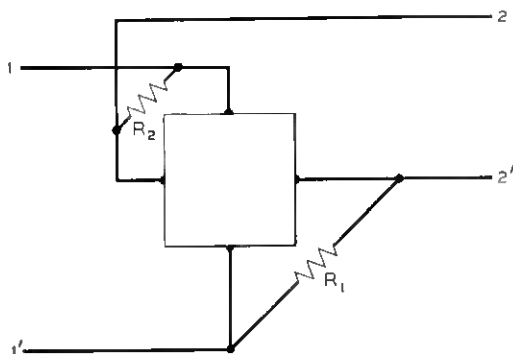


Fig. 3 — Hall effect isolator made from gyrator.

An ideal gyrator is also an ideal impedance inverter, but the loss inherent in the Hall effect gyrator makes it useless as an inverter. This loss is of such a nature that it could be effectively removed by placing negative resistance of the proper value in series with the gyrator. This lossless or ideal gyrator should then find application in network synthesis.

The gyrator has been discussed here in some detail not only because of its historical interest but, primarily, because it has been the basis of the major portion of the devices to be described in the remainder of this article.

2.2 Isolator

An isolator^{2,3,5} is a nonreciprocal four-pole in which one of the two transfer impedances is zero. Thus it transmits signals in one direction only — say, from terminal pair 1 to pair 2. Such a device can be made from a Hall effect gyrator by the addition of two shunting resistors (see Fig. 3). Assuming the gyrator is symmetrical, its two self impedances may be called z_s and its transfer impedances z_{T1} and z_{T2} . Furthermore,

$$z_{T1} = -z_{T2} = z_T > 0.$$

Then it can be shown that the parallel resistance values R_1 and R_2 must be such that

$$R_1 + R_2 = \frac{z_s^2 + z_T^2}{z_T}$$

in order to obtain an isolator.

Since $z_T/z_s = \tan \theta_H$, it may be seen that, with R_1 and R_2 fixed in value, to maintain isolation one must maintain a constant H , a constant μ_H and, therefore, a constant temperature. Also, if the resistance of the

semiconductor contacts is current-sensitive, isolation might be possible only in a limited range of signals. The forward loss of such an isolator is just slightly less than that of the gyrator from which it is made. The minimum possible forward loss is 6 db, and the loss in experimental germanium isolators has been approximately 14 db. Thus, to transmit 1 mw of power, one must apply 26 mw to the input because 25 mw will be dissipated in the isolator, most of it in the germanium. This is another source of current sensitivity — changing the current level changes the amount of power dissipated in the isolator, changing its temperature and therefore z_s , z_T/z_s , R_1 and R_2 . The higher the reverse loss is (i.e., the better the balance) the more sensitive the isolator will be to all such variations. The reverse loss in germanium isolators has been made 75 db, and this loss is quite sensitive to small variations.

When indium antimonide (which has a much higher μ_H) was used, the forward loss was reduced to about 7.5 db, even with a smaller H . With carefully adjusted shunt resistance values, a reverse loss of “the order of 100 db” was obtained. Notice that the forward loss was very nearly the minimum possible loss. Unfortunately, InSb has a quite small energy gap and consequently a very high np product. This fact, in conjunction with its very high mobility, inevitably gives rise to a very low resistivity. Thus, the InSb isolator had an impedance of around 1 ohm. Germanium isolators can easily be made with impedance levels ranging from 10 to 1000 ohms.

Typically, μ_H in InSb is a much stronger function of temperature than is μ_H in germanium. Thus, this low-loss isolator has some serious drawbacks.

There is another way of making a Hall effect isolator, and it has been suggested that this second method should remove the temperature sensitivity. Unfortunately, this is not the case. This form of the isolator

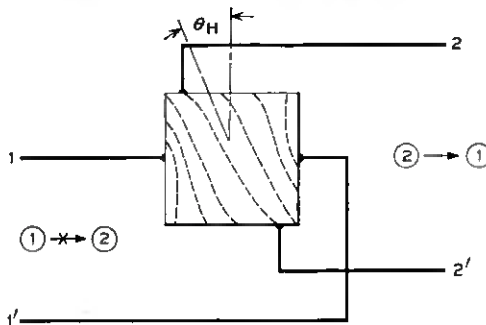


Fig. 4 — Hall effect isolator (or skew gyrator).

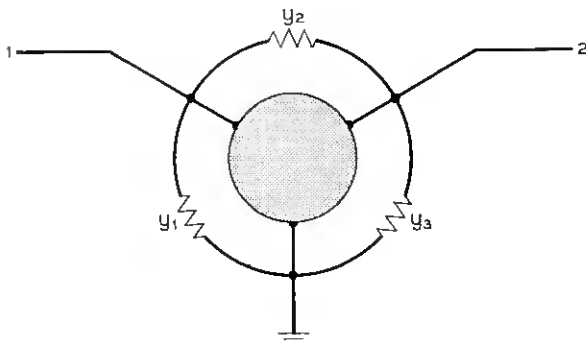


Fig. 5 — Negative resistance amplifier.

uses no resistors and has been called a “skew gyrator,” in which the sample is square but the leads are not symmetrically spaced. Fig. 4 shows the sample with the dotted equipotential lines which would occur with a signal applied to side 1 if the magnetic field were the proper value. Obviously, no voltage appears between leads 2 and 2'. If a signal is applied to side 2, there will be an output at side 1. This is stated symbolically as

$$1 \leftrightarrow 2 \rightarrow 1.$$

This type of isolator is still sensitive to variations in H , μ_H and temperature (and, therefore, current level). This sensitivity arises from the need for a constant θ_H to maintain the slope of the equipotential lines. Of course, if the semiconductor has the proper doping level, so that μ_H is independent of temperature near room temperature, this device could be quite stable. For that matter, most of these devices could be stable with such a constant mobility material.

As for the shunt-resistor isolator, it is conceivable that resistors could be used that had the proper temperature coefficient for this isolator also to be made stable — perhaps even more stable than the skew gyrator.

2.3 Negative-Resistance Amplifier

The negative-resistance amplifier^{5,11} uses a slice of semiconductor with three equally spaced edge contacts, a perpendicular magnetic field of constant value H and three negative resistances (see Fig. 5). For the moment, assume the parallel conductances are not in the circuit. The short-circuit admittances may be defined as

$$y_{11} = \left. \frac{\dot{i}_1}{v_1} \right|_{v_2=0} = y_0, \quad y_{22} = \left. \frac{\dot{i}_2}{v_2} \right|_{v_1=0} = y_0,$$

$$y_{12} = \left. \frac{\dot{i}_1}{v_2} \right|_{v_1=0}, \quad y_{21} = \left. \frac{\dot{i}_2}{v_1} \right|_{v_2=0}.$$

A parameter α (an odd function of θ_H) may be such that

$$y_{12} = -\frac{y_0}{2}(1 - \alpha), \quad y_{21} = -\frac{y_0}{2}(1 + \alpha) \quad |\alpha| < 1.$$

However, if the parallel admittances are connected and given the values

$$y_1 = y_3 = y_{21},$$

$$y_2 = y_{12},$$

it can be shown that $i_1 = 0$ and $i_2 = \alpha y_0 v_1$. Therefore, with the negative conductances in parallel, the device becomes an isolator which can have gain in the forward direction. One of its advantages is that it permits the construction of negative-resistance amplifiers which are unidirectional, so that higher gain with the same degree of stability is possible with this type of device than with the usual two-terminal negative resistance.

Negative resistance obtained from gas tubes was used in this way to give a forward gain of 6 db and a reverse loss of 46 db.

This type of device was proposed as a high-frequency amplifier because of the inherent insensitivity to frequency of Hall effect gyrators and related devices. Just how good they are at very high frequencies is a question which will be discussed a little later.

The same principle may be applied to the shunt-resistor isolator. Still another form of negative resistance amplifier employing Hall effect will be mentioned in the next section.

2.4 Circulator

A circulator is a nonreciprocal n -port device ($n > 2$) in which a signal applied at one port is transmitted only to an adjacent port (e.g., the adjacent clockwise port). No output appears at any other port. A three-port circulator can be made which uses the Hall effect. As one may observe from the skew gyrator, the Hall effect merely enables one to prevent any signal from appearing at one of the outputs. Therefore, if we had four or more ports, a portion of the signal would appear at all ports except one, and the device would not really be a circulator. See Fig. 6 for the symbol of a three-port circulator.

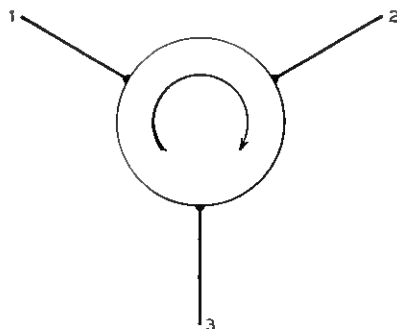


Fig. 6 — General three-port circulator.

The Hall effect circulator^{12,13} consists of a circular slab of semiconductor (n-type germanium has been used to date) with six equally spaced edge contacts and a constant magnetic field H applied perpendicular to the plane of the slab. Such a sample is shown in Fig. 7 with equipotential lines dotted in. With the proper value of H , no output appears at 3 but there is an output signal at 2. The results would be similar for an input at 2 or 3 as well. So we may say

$$1 \rightarrow 2 \rightarrow 3 \rightarrow 1,$$

$$3 \leftrightarrow 2 \leftrightarrow 1 \leftrightarrow 3.$$

When a load is attached so that a current may flow at 2 (with the signal applied at 1) this secondary current flow produces an additional Hall electric field which has the net effect of reducing θ_H . To maintain

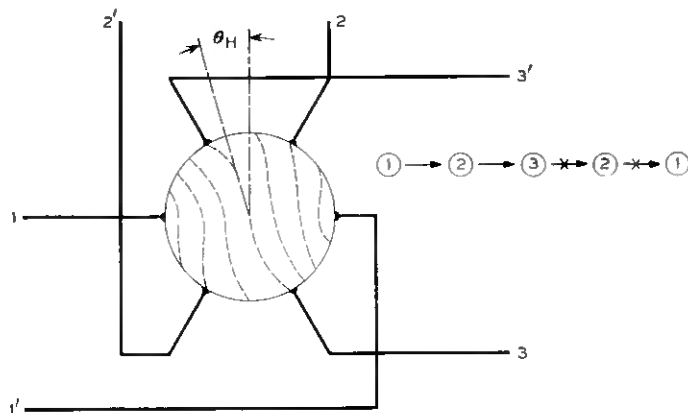


Fig. 7 — Hall effect circulator.

circulation, H must be increased. Thus, it is obvious that the circulator must be terminated in the proper impedances if its circulation property is to be retained. At a particular magnetic field, H_0 , the load impedance necessary to "balance" the circulator is the same as the circulator's input impedance, so H_0 is the logical field to use to avoid impedance mismatches.

When n-type germanium is used, $H_0 \approx 14,500$ oersteds, and the matching load impedance, Z_L , is a function of the sample's thickness and resistivity, and might vary from 10 to 1000 ohms. If n-type InSb were used, H_0 would be in the order of 1000 oersteds, but (as with the isolator) Z_L would have to be quite low — 1 ohm or even less.

It is almost impossible to make a circulator which has the proper impedance level and is symmetrical so that it will circulate a signal applied at any input. Fortunately there is a means of overcoming this difficulty so that a circulator of the type which is quite easily fabricated may be made to appear symmetrical and to operate at the proper impedance level. All that is involved is placing a network of six resistors in parallel with the circulator (see Fig. 8). Use C_P with $H < H_0$ and C_S with $H > H_0$. Apparently the three-resistor networks shown in Fig. 9 can be used equally effectively but, since C_P and C_S have six variables, they should offer greater flexibility than does C_P' or either form of C_S' .

The parallel networks permit Z_L to be adjusted to any value between about one tenth its normal value and infinity (open-circuit), using $H = H_0$ all the while. One can also use the normal value of Z_L with any H from zero to infinity. Another possible function of the networks is to permit an approximately symmetrical circulator to be used with unequal loads at the three-ports.

The forward loss of a circulator is about 17 db if it is used without parallel networks, and reverse losses of about 65 db have been achieved. With n-type germanium, H must be about 14,500 oersteds. If InSb were

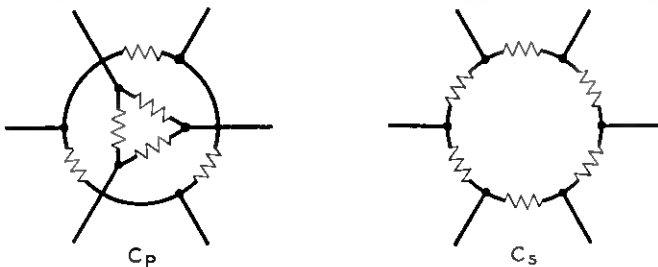


Fig. 8 — Parallel networks for circulator.

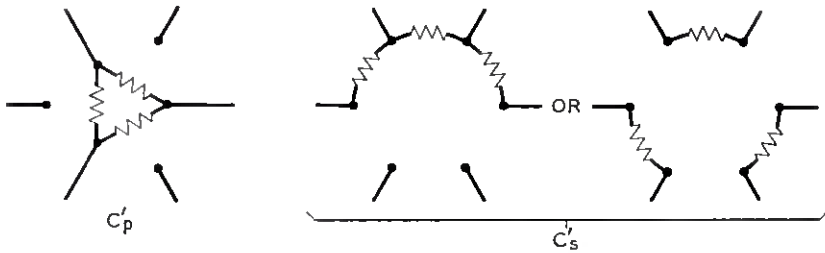


Fig. 9 — Alternative parallel networks.

used, H would be only about 1000. However, the forward loss would still be 17 db, since the loss is determined by θ_H , and θ_H must be the same in any Hall effect circulator.

On the other hand, if one uses network C_s and increases H , the forward loss is reduced. With n-type germanium, it has been reduced to about 15 db. Judging by the low loss achieved with an InSb isolator, one would expect that an InSb circulator could be made with a forward loss within one db or so of the minimum possible loss. The minimum loss has been calculated to be 8.4 db, so an InSb circulator probably could be made with a 9- or 10-db forward loss. The reverse loss would depend solely on how precisely the resistance values of C_s were adjusted. Of course, InSb still has the disadvantage of very low resistivity, and the impedance level would be about 1 ohm. As previously stated, such a circulator could be operated with any higher impedance load, but the necessary mismatch would increase the insertion loss tremendously.

Perhaps another means of reducing the loss should be considered — that is, using negative resistance in C_p or C_s . By this means, the loss could be removed, or gain could even be achieved, and one would have an amplifying circulator. Negative resistance is still hard to get, but if someone knew of a good means of producing it, he could make the circulator an even more interesting device.

III. GENERAL REMARKS

All the devices described in the preceding pages could theoretically transmit dc signals as well as ac signals of any frequency up to the dielectric relaxation frequency of the semiconductor material. In order to realize even a major portion of this huge bandwidth (\approx thousands of megacycles) special care must be taken to shield all input and output leads properly. The samples should be small, so that only a small magnet will be required. Consequently, all the leads must be fairly close to one

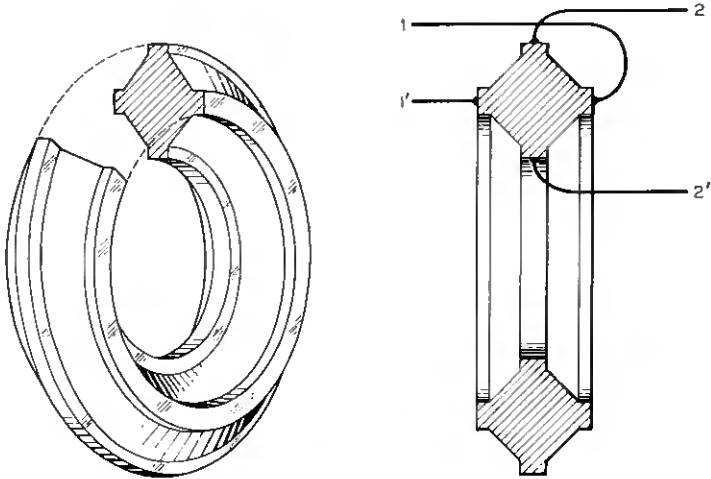


Fig. 10 — One-piece gyrator.

another where they attach to the sample. If the lead problem can be overcome, the devices will remain resistive (and essentially insensitive to frequency) until the relaxation frequency is approached. If leads are treated carelessly, the top frequency may be less than one megacycle.

For linear operation, the contacts must be as nearly ohmic as possible and contact resistance must be minimized. Nonlinearity might be a serious problem with any of the above devices if a wide dynamic range of operation were required.

If a good permanent magnet material were known (which also had a respectable Hall mobility) the above devices could be fabricated from a single piece of this material (see Fig. 10). The torus would be made of the unknown material, and it would be magnetized. The figure shows four contacts equally spaced about the torus' body, so the pictured device would be a gyrator. It could also be an isolator or a circulator. Two advantages of this structure would be that it would be all in one piece (for mechanical simplicity and rigidity) and that there would be no air-gap in the path of the magnet.

IV. SIGNAL-PRODUCED MAGNETIC FIELD DEVICES

4.1 *Switch*

Hall effect switches⁵ for the most part are not simply single-pole single-throw switches. The simplest Hall effect switch might be termed

a double-pole single-throw switch. This might be a gyrator sample in the gap of an electromagnet. Let us call the two terminal pairs of the gyrator 1 and 2. When there is no magnet current flowing, 1 is not connected to 2. But, with the magnet energized, 1 is connected to 2 and the loss introduced by the switch "contacts" is the insertion loss of the gyrator.

One can make a more intriguing switch by using an isolator instead of a gyrator. With no magnetic field, this is not very interesting, because 1 and 2 are then merely reciprocally connected very inefficiently. But with the proper value of H , 1 is connected to 2 but 2 is isolated from 1 ($1 \rightarrow 2 \leftrightarrow 1$). Reversing the field reverses the direction of isolation. Furthermore, if the isolator is constructed from a gyrator plus negative resistances, so that the isolator is a unidirectional amplifier, switching the magnetic field changes the direction of amplification. It is as though a broadband amplifier could be physically turned around by energizing a coil in the opposite sense.

An even more elaborate switch involves switching a circulator. It was first proposed that a signal be applied to terminals 1 of a circulator and that it be switched between output 2 and output 3 by reversing the direction of H . However, when this is done, no less than three "isolator switches" are reversed at the same time because any two terminal pairs in the circulator make up an isolator.

Unfortunately, none of the above-mentioned switches are extremely fast in their operation. They all involve collapsing a rather sizable magnetic field and reproducing it in the opposite sense. This requires a particular amount of energy and, therefore, either a lot of time or a lot of power. There is a faster, more elegant means of switching, but it presents certain difficulties.

This second means permits the magnetic field to remain constant, but requires the resistance values in the parallel networks to change. The network values must be negative for at least one of the conditions and could be either positive or negative for the other. This might be achieved by using a fixed negative resistance in series with a positive resistance whose value could be shifted or switched from one desired value to another desired value. For example, one might use three resistors, one negative and two positive, and have a solid-state switch across one of the positive resistances. By shorting or opening the switch, the Hall effect switch could be operated.

Thus, an isolator switch could be reversed by opening or closing two electronic switches. A circulator switch could be operated by opening or closing three electronic switches. One calculation on a particular circu-

lator indicates that, by changing the resistance values in C_p from -304 ohms to -168 ohms, the circulator could be switched. In this instance the forward loss in one condition is 4.5 db and in the other condition there is a forward gain of 3 db.

The reverse loss in each condition should still be of the order of 50 to 60 db. The only real difficulty is in obtaining the negative resistances. It should be noted that amplifier connections could be reversed with electronic switches, and the Hall effect devices would not be needed. Only the circulator switch performs a new function. Perhaps it would be useful.

4.2 *Transducer*

Hall effect transducers³ can be constructed which convert mechanical motion into electrical signals. The mechanical motion moves a gyrator sample in and out of the air gap of a permanent magnet. Such a device could be useful in measuring strain or other displacements. It has been estimated (admittedly optimistically) that an InSb transducer could detect a displacement of about 1 angstrom if the gradient of H were 10 oersteds/micron.

The same device could be used as a phonograph pickup or as a microphone. It should be interesting here because it responds to very low frequency — even direct current.

A microphone and strain gauge have been constructed, and they behaved more or less as expected.

4.3 *Magnetic Field Meter*

The Hall effect has been used to measure magnetic field strengths for about ten years.^{2,3,14,15,16} Since the output voltage of a Hall effect sample is proportional to the product of the sample input current I and the applied H , if I is constant the output voltage is proportional to H . For some years now, there has been a commercial instrument based on this principle on the market. An alternating current is applied to the sample so that the output is this ac carrier modulated by the magnetic field. Thus, either slowly varying direct current or low-frequency ac fields can be measured. The probe is about 25 mils thick (and could be made thinner) so that the field even in a thin air gap can be measured. Field strengths from a few oersteds up to 30,000 oersteds can be readily measured with a germanium probe. With care, even smaller fields are measurable. InSb can supposedly detect fields as low as 10^{-3} oersteds, but that sensitivity has not yet been achieved. However, an InSb sample can easily measure the earth's field.

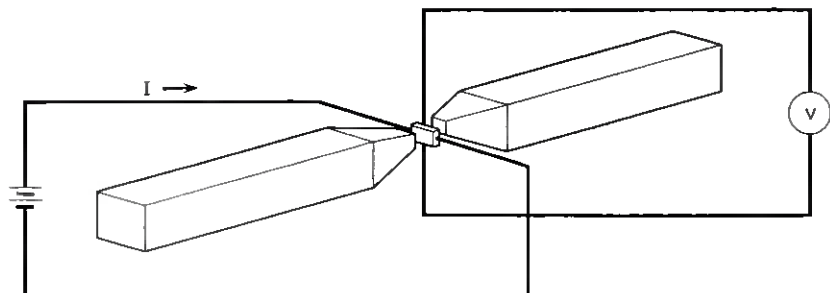


Fig. 11 — Electrical compass.

The Hall effect magnetic field meter is useful in measuring air-gap flux densities, flux distributions (because of its small size), demagnetization and hysteresis curves, and several other parameters. Its accuracy is inversely proportional to the temperature coefficient of carrier mobility in the semiconductor, so something like germanium should be of more general use than InSb. Nevertheless, InSb definitely gives a greater output for the same field, and its use is indicated where small fields or small variations are to be determined. It can be shown that

$$\frac{\text{power out}}{\text{power in}} \propto H^2 \mu_H^2.$$

Thus, InSb has an obvious advantage when H is small.

4.4 *Electrical Compass*

A Hall effect electrical compass³ can be constructed as shown in Fig. 11. The rods are made of high permeability material and serve to concentrate the earth's field on the sample, thus increasing the sensitivity. Since the output is proportional to $I \times H$, the reading will be maximum (say positive) when the rods are in line with the earth's field in one direction, and maximum negative when the rods are in line in the opposite direction. When the rods are perpendicular to H , there will be a null. Using an InSb sample, with no amplification, a rotation of 1° from the null position is detectable. With amplification, the sensitivity could be increased.

In this device, a change in temperature would only affect its sensitivity slightly if it were used at a maximum or minimum position.

4.5 *Magnetic Field Variation Meter*

If two Hall effect samples are connected in series as shown in Fig. 12, and if they are identical samples, there will be no output voltage except

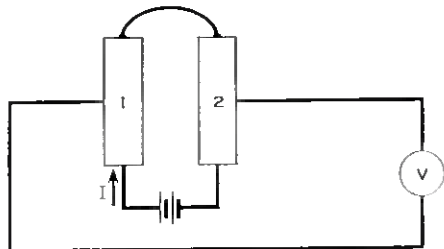


Fig. 12 — Magnetic field variation meter.

when the average field in sample 1 is different from the average field in sample 2.¹⁷ If the field is uniform, the two outputs will cancel. Such a device is useful in determining magnetic field gradients, and has been used to detect cracks in metals. To do this, a permanent magnet produces a field in the metal, and where there are surface cracks (perhaps too small to be seen with the naked eye), there will be a stray magnetic field at the metal surface. Such a localized field can be readily detected with this type of device.

If a single Hall effect sample were used, a positive and then a negative pulse would appear in the output for each crack. With the two samples in series and properly spaced, the two pulses can be additive.

4.6 Ammeter

A Hall effect ammeter³ can be constructed in at least a couple of ways. If one uses an ordinary gyrotator (with a permanent magnet field) and connects the input leads in series with a current-carrying line, the gyrotator's output voltage will be proportional to the line's current. If the current is too high for the gyrotator sample, the gyrotator input can be shunted with a calibrated resistor. One advantage of this is that it permits one to measure very high frequency currents. This hardly seems preferable to inserting a simple resistor in the line and measuring the voltage developed across it, but this Hall effect ammeter is mentioned because it will be referred to in the next section.

Another and perhaps more useful form of ammeter is shown in Fig. 13. A yoke of ferromagnetic material is hinged at the bottom so that it can be clipped around a current-carrying conductor. A Hall effect sample is in the magnetic path, with a direct current applied to two of its terminals. If there is either a direct or alternating current flowing in the conductor, there will be a corresponding output voltage proportional to the current, since the magnetic field is proportional to the conductor

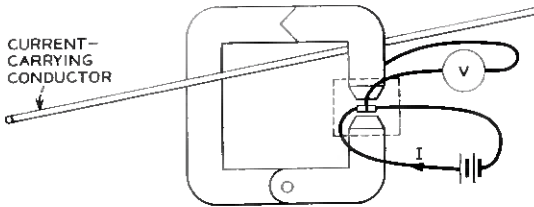


Fig. 13 — Ammeter.

current. This device should be useful when it is undesirable or impossible to break a line to insert an ordinary ammeter. It should be particularly useful for direct currents, because there is no similar instrument which measures direct currents.

4.7 Wattmeter

A Hall effect wattmeter^{3,18,19} is easily obtained by using the first ammeter described above and employing the voltage between the two conductors to energize a coil wound on the magnetic core. Thus, the sample's input current is proportional to the circuit's voltage, and the output voltage is proportional to the power transmitted through the conductors. This is a simple wattmeter, but it is not useful at high frequencies.

A wattmeter can be used on a coaxial line by connecting the Hall effect sample from the center conductor to the shell so that the plane of the sample includes the axis of the line (see Fig. 14). Thus, the magnetic field is supplied by current in the central conductor, and the output voltage is proportional to the power.

Due to its high-frequency capabilities, the Hall effect has been proposed as a means of measuring power in waveguide circuits. In this case, the E field provides the current to the sample and the H field supplies

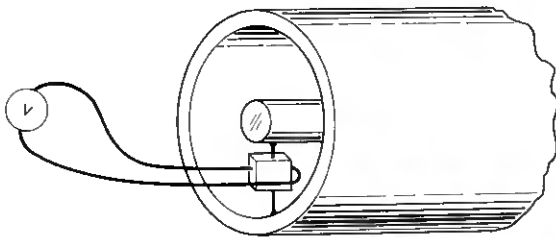


Fig. 14 — Coaxial line wattmeter.

the H field. The output leads are brought out on equipotentials through holes in the sides of the waveguide. It is theoretically possible to meter power in this manner to quite high frequencies.

4.8 *Amplifier*

If a direct current is supplied to the input of a Hall effect sample which is in the gap of an electromagnet, the setup can be an amplifier.^{3,20} If the magnet is efficient enough, and if the sample's input current is large enough, the output power from the sample can be greater than the power applied to the magnet. With InSb, a gain of 5 has been achieved. Greater gains are probably attainable.

Such an amplifier could be used for direct or alternating current. The upper frequency limitation would be determined by the quality of the magnetic core and the stray capacitance of the core winding. This type of amplifier seems to offer little or nothing that a more conventional amplifier would not offer, but it is interesting in that it requires a low-voltage, high-current power supply — perhaps the type of thing thermoelectric generators could furnish.

4.9 *Modulator*

The Hall effect can be used to make a product modulator.² If one signal is applied to the input leads of a Hall effect sample and another signal is applied to a winding on a magnetic core to produce the magnetic field, then the output voltage is proportional to the product of the two signals. This type of modulator has been proposed as a substitute for the mechanical chopper found in most sensitive dc amplifiers. The output voltage of a straight dc amplifier drifts with time and temperature because of instability of the components. Mechanical choppers are used to convert the dc input to an ac signal, which in turn can be amplified with much less drift. The Hall effect modulator has an inherent drift-free zero set and can convert dc to ac with no mechanical motion. Thus, it should be longer lived than relay-type choppers.

Furthermore, mechanical choppers usually operate at 60 cps, whereas the Hall effect modulator frequency may easily be 1000 cps. This permits a broader band dc amplifier to be realized. Direct current signals as low as 20 microvolts (across 400 ohms input impedance) have been successfully amplified, and a dynamic range of 70 db has been observed. If additional effort were applied to reducing pickup in the output circuit, the 20-microvolt minimum could presumably be further reduced.

4.10 Demodulator

A Hall effect square-law demodulator or detector²¹ is precisely the same as a device referred to as a Hall effect full-wave rectifier elsewhere.³ If an incoming signal represented by $I = I_0 \cos \omega t$ is applied to both the Hall effect sample input leads and also to the magnet winding, the output voltage will be

$$V = \frac{k_1 I_0^2}{2} + \frac{k_1 I_0^2}{2} \cos 2\omega t, \quad (4)$$

where k_1 is a constant involving the Hall constant R_H and the various parameters of the magnetic circuit. The first term in (4) is a dc term proportional to the square of the amplitude of the input signal. The second is a double-frequency component which can be easily filtered out. Such a detector should give an accurate square conversion for a wide range of amplitudes.

A linear Hall effect detector²¹ may be constructed on the same principle if a sine wave of constant amplitude (and the same frequency ω) is applied to the magnet winding. Then (4) becomes

$$V = \frac{k_2 I_0}{2} + \frac{k_2 I_0}{2} \cos 2\omega t. \quad (5)$$

The first term in (5) is proportional to the first power of the signal amplitude, and so this device should give a dc output voltage proportional to the input signal amplitude. Again, the conversion ratio should be linear for a wide range of amplitudes.

One difficulty with both these demodulators is that they are operable only up to frequencies at which magnetic fields of the order of 1000 oersteds or greater can be sinusoidally reversed with the power that happens to be available.

4.11 Frequency Spectrum Analyzer

A highly selective frequency spectrum analyzer²¹ can be made from the linear Hall effect demodulator. If the signal applied to the sample is $I_s \cos \omega_s t$ and the signal applied to the magnet winding is of frequency ω_H , the usual sum and difference frequency components appear in the output:

$$V = \frac{k_2 I_0}{2} \cos (\omega_s - \omega_H)t + \frac{k_2 I_0}{2} \cos (\omega_s + \omega_H)t. \quad (6)$$

It may be noted that (6) reduces to (5) when $\omega_s = \omega_H = \omega$.

Suppose the signal applied to the sample is composed of many frequency components and that the amplitude of the magnetic field is held constant as its frequency is changed. If the output is applied to a dc voltmeter, there will be a reading only when ω_H coincides with some ω_s , say ω_{si} . Actually, the dc output is given by

$$V_{dc} = \frac{k_3 I_{si}}{2} \cos(\Phi_H - \Phi_{si}), \quad (7)$$

where k_3 is a constant and $(\Phi_H - \Phi_{si})$ is the difference between the phases of the magnetic field and the i th component of the input. Thus, unless the phase difference is known, it is impossible to determine I_{si} . Therefore, it is recommended that ω_H be adjusted to within about 1 cps of ω_{si} so that the "dc" output is a 1-cps cosine wave whose argument $(\Phi_H - \Phi_{si})$ either increases or decreases 2π radians per second. If this output is applied to a zero-centered dc voltmeter with a response time $\tau < 0.25$ second, the meter needle follows the voltage variation, and the amplitude of the i th component is proportional to the maximum swing of the needle.

Suppose the frequency of the next component is ω_{sk} . If $(\omega_H - \omega_{sk})$ is large enough so that

$$\frac{2\pi}{\omega_H - \omega_{sk}} < \frac{\tau}{10}$$

then the meter will not respond to this component's contribution to the output. Thus, by slowly sweeping ω_H , a frequency spectrum analysis of any input wave can be obtained.

This device responds with equal accuracy to both small- and large-amplitude components. Its upper frequency is limited, because a sinusoidal magnetic field must be obtained at that frequency. Fortunately, the signal does not have to supply this field; a high-power oscillator may be used. Since the available output level is determined by the component amplitude and the value of magnetic field, the upper frequency will be determined by the required sensitivity. Perhaps its upper limit in a typical use would lie between 0.5 and 5 mc.

Of course, it operates very well down to zero frequency, and this is probably where it will find its greatest use — at audio frequencies and even lower. Conventional analyzers require huge tuned filters at these frequencies to obtain the Q required for high selectivity. With the Hall effect analyzer, all that one requires is a single low-pass filter. Notice that the amplitude of the output is independent of the frequency difference $(\omega_H - \omega_{si})$. Thus, a single low-pass filter of adjustable upper fre-

quency gives one a very simple means of adjusting selectivity without affecting the sensitivity.

With a crude device like this, it was quite easy to measure the first ten harmonics of a square wave of 25 cps fundamental frequency. With care, 18 harmonics were detected. This device does not represent the ultimate in accuracy, but it may do so in small size, simplicity and ease of operation.

4.12 *Phase Discriminator*

The Hall effect phase discriminator simply makes use of the phase sensitivity of the spectrum analyzer. It is sometimes necessary to obtain an electrical signal which is a measure of the difference between the phases of two signals of the same frequency. If these two signals are applied to the spectrum analyzer, the output is a dc component plus a double-frequency component [see (5)]. The dc term is given more completely by (7). If the two amplitudes are constant (not necessarily equal), the dc output can be calibrated to give the phase difference directly. Like the analyzer described above, the frequency limit will be determined largely by the required sensitivity. Perhaps this limit will be somewhere between 1 and 10 mc. This is estimated to be higher than that for the analyzer, on the assumption that the signal level can be higher in the discriminator (since it needs to operate at only one level).

Incidentally, this phase sensitivity must be considered if one constructs the wattmeter or demodulator previously referred to.

4.13 *Digital-to-Analog Encoder*

The Hall effect digital-to-analog encoder is obtained by applying the maximum allowable dc voltage to the sample input leads and using n separate windings (to encode n binary digits) on the magnetic core (see Fig. 15). Possibly one extra winding would be used for zeroing purposes. Current limiters assure that the current input levels are fixed, all at the same level; the digits are weighted by different numbers of turns in the windings. If the current limiters do not pass precisely the desired current the situation might be improved by adjustment of turns. All the input circuits may be completely isolated in a dc sense. A disadvantage is that the minimum time in which the magnetic field can be changed is limited to something on the order of a millisecond, so that only about 1000 binary numbers per second may be applied. Since this is the same setup as the Hall effect amplifier, it is conceivable that the encoder might introduce gain.

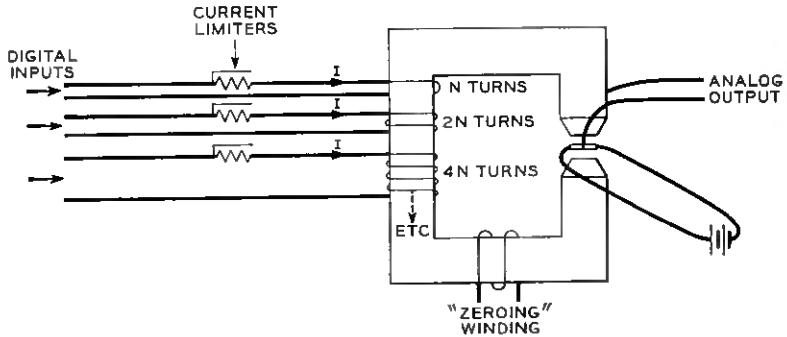


Fig. 15 — Digital-to-analog encoder.

4.14 Analog Multiplier

The Hall effect analog multiplier^{3,20} is essentially a Hall effect demodulator. A direct current is driven through the sample input leads and another direct current is applied to a winding on the core. The output voltage is proportional to the product of these two signals. Thus the device is an analog multiplier. Ref. 22 concluded that silicon offers the best combination of properties for this device. In order to realize a reasonable accuracy in the output (about 0.3 per cent) it was necessary to have a response time of approximately 1 millisecond. The error was reduced slightly (to about 0.1 per cent) when the input repetition rate was decreased to 20 cps. Ref. 22 also presents a scheme for minimizing temperature effects.

V. MATERIALS

Strange as it may seem, all of the above devices place approximately the same requirements on the semiconductor material used. Generally speaking, the ideal material would have a high-electron Hall mobility and a large energy gap, and both these parameters would be constant.

If a material with these properties were available, minimum losses would be realizable with rather small magnetic field strength — say, 1000 oersteds. Furthermore, sample resistances would be fixed by their geometry and doping level. Of course, other requirements could be laid down, but they are certainly of less importance than the two given above.

Now, these two requirements are mutually incompatible. That is, a material with a high mobility μ_H is almost certain to have a small energy gap, E_g , and, *vice versa*, a material with a large E_g will have a small μ_H . If this relationship did not exist, many Hall effect devices would find much wider application than they now enjoy.

Intermetallic semiconductors offer the best hope of a partial solution to this problem. They show a wide range of properties but still follow

the trend of combining small Eg with high μ_H . One material however, gallium arsenide, has a mobility of 6000 cm²/volt-second (higher than that of germanium) and an energy gap of 1.35 electron volts (higher than that of silicon). Let us hope that other materials will be found which will violate the ($\mu_H - Eg$) trend even more definitely and that they will be relatively easy to fabricate in highly pure single-crystal form. Refs. 23, 24 and 25 describe a large number of semiconductors (both elemental and intermetallic) and their properties and refer to the apparent relationship between μ_H and Eg . Ref. 26 gives experimental information on mercury selenide, a II-VI compound.

Before leaving materials it is worth mentioning that Ref. 5 develops an elaborate mathematical solution to the problem of finding the electric field distribution in a Hall effect sample. Several different shapes are treated, including the square gyrator, the skew isolator and the circulator. Ref. 27 suggests a method of measuring resistivity and Hall effect in flat samples of any arbitrary shape. Ref. 28 gives an analog method for obtaining the same two parameters; this last method involves no analytical computation.

VI. CONCLUSIONS

The Hall effect has thus far furnished a tremendous variety of devices. All of these devices have at least one important advantage over most semiconductor devices — Hall effect devices require the use of majority current carriers only and therefore involve no junctions. There are no areas of concentrated electric fields, so the surfaces need not be elaborately protected. Most or perhaps even all of the devices can be operated indefinitely in ordinary indoor atmospheres with no protection whatever.

Because these are majority carrier devices, the lifetimes of minority carriers are of no consequence. When an electric field is induced in a semiconductor, a brief time elapses while the current carriers redistribute themselves. This relaxation time (ordinarily less than a millimicrosecond) is the only factor which limits the frequency of signals which can be satisfactorily transmitted through a Hall effect sample. Interlead capacitance will tend to shunt part of the signal around the sample and will thus limit the frequency range of operation more severely. Of course, where a signal supplies the magnetic field, the frequency limitations will depend on the core used, the winding, the signal level and the sensitivity and linearity required.

The Hall effect is a small effect. As a result, the devices it makes possible are inefficient or lossy. In order for Hall effect devices to be stable, the majority carrier mobility must be held constant, as must the magnetic field (or core permeability as the case may be). In order that these devices have linear characteristics, all current-carrying contacts must be ohmic.

Despite the disadvantages outlined in the preceding paragraph, Hall effect devices are being studied with steadily increasing interest.

REFERENCES

1. Carlin, H. J., Synthesis of Nonreciprocal Networks, Proc. Symp. on Modern Networks, Polytechnic Inst. of Brooklyn, Brooklyn, N. Y., April 1955.
2. Mason, W. P., Hewitt, W. H. and Wick, R. F., Hall Effect and "Gyrators" Employing Magnetic Field Independent Orientations in Germanium, *J. Appl. Phys.*, **24**, 1953, p. 166.
3. Ross, I. M. and Saker, E. W., Applications of Indium Antimonide, *J. Elect.*, **1**, 1955, p. 223.
4. Tellegen, B. D. H., The Gyrator, A New Network Element, *Phil. Res. Rep.*, **3**, 1948, p. 81.
5. Wick, R. F., Solution of the Field Problem of the Germanium Gyrator, *J. Appl. Phys.*, **25**, 1954, p. 741.
6. Casimir, H. B. G., On Onsager's Principle of Microscopic Reversibility, *Rev. Mod. Phys.*, **17**, 1945, p. 343.
7. Meixner, J., Zur Theorie der elektrischen Transporterscheinungen im Magnetfeld, *Ann. der Phys.*, **40**, 1941, p. 165.
8. McMillan, E. M., Violation of the Reciprocity Theorem in Linear Passive Electromechanical Systems, *J. Acoust. Soc. Amer.*, **18**, 1946, p. 344.
9. McMillan, E. M., Further Remarks on Reciprocity, *J. Acoust. Soc. Amer.*, **19**, 1947, p. 922.
10. Hogan, C. L., The Ferromagnetic Faraday Effect at Microwave Frequencies and Its Applications, *B. S. T. J.*, **31**, 1952, p. 1.
11. Shockley, W. and Mason, W. P., Dissected Amplifiers Using Negative Resistance, *J. Appl. Phys.*, **25**, 1954, p. 677.
12. Semmelman, C. L., U. S. Patent No. 2,774,890.
13. Grubbs, W. J., The Hall Effect Circulator — A Passive Transmission Device, *Proc. I.R.E.*, **47**, 1959, p. 528.
14. Hennig, G. R., Applying the Hall Effect to Practical Magnet Testing, *Elect. Mfg.*, **59**, 1958, p. 132.
15. Pearson, G. L., A Magnetic Field Strength Meter Employing the Hall Effect in Germanium, *Rev. Sci. Inst.*, **19**, 1948, p. 263.
16. Saker, E. W., Cunnell, F. A. and Edmond, J. T., Indium Antimonide as a Fluxmeter Material, *Brit. J. Appl. Phys.*, **6**, 1955, p. 217.
17. Wolfe, R., private communication.
18. Barlow, H. E. M., The Application of the Hall Effect in a Semiconductor to the Measurement of Power in an Electromagnetic Field, *Proc. Inst. Elect. Eng.*, **102B**, 1955, p. 179.
19. Barlow, H. E. M., The Design of Semiconductor Wattmeters for Power-Frequency and Audio-Frequency Applications, *Proc. Inst. Elect. Eng.*, **102B**, 1955, p. 186.
20. Ross, I. M. and Thompson, N. A. C., An Amplifier Based on the Hall Effect, *Nature*, **175**, 1955, p. 518.
21. Bogomolov, V. N., Some New Semiconductor Devices (New Uses of the Hall Effect), *Zh. Tekh. Fiz.*, **26**, 1956, p. 693.
22. Löfgren, L., Analog Multiplier Based on the Hall Effect, *J. Appl. Phys.*, **29**, 1958, p. 158.
23. Coblenz, A., Semiconductor Compounds, *Electronics*, **30**, 1957, p. 144.
24. Pincherle, L. and Radcliffe, J. M., Semiconducting Intermetallic Compounds, *Adv. in Phys.*, **5**, 1956, p. 271.
25. Stello, P. E., Van Winkle, D. M. and Turner, J. D., Semiconductors, Materials and Properties, WESCON Conf., August 1956.
26. Blum, A. I. and Regel, A. R., Electrical Properties of Solid Solutions of Mercury Selenide and Selenium, *Zh. Tekh. Fiz.*, **21**, 1951, p. 316.
27. Van der Pauw, L. J., A Method of Measuring Specific Resistivity and Hall Effect of Discs of Arbitrary Shape, *Phil. Res. Rep.*, **13**, 1958, p. 1.
28. Broudy, R. M., Galvanomagnetic Coefficients for Arbitrary Geometry, *J. Appl. Phys.*, **29**, 1958, p. 853.

An Appraisal of Received Telephone Speech Volume

By O. H. COOLIDGE and G. C. REIER

(Manuscript received November 3, 1958)

One of the attributes of telephone service which is of importance to a telephone user is the loudness with which he hears the voice of a distant talker. Related to this loudness is an objective measurement of "received volume."

This paper represents the results of subjective tests made to determine a relationship between received volume and the satisfaction of telephone listeners. The results are shown as statistical distributions of listeners' opinion, which may be combined with estimated distributions of received volume in the telephone plant to give "grade of service."

Grade of service objectives have been stated as 95 per cent of connections rated "Good", 5 per cent "Fair" and a negligible percentage "Poor." Increased use of the more efficient 500-type telephone sets, and planned improvements in the circuits which interconnect them, will make it possible to meet these objectives.

I. INTRODUCTION

Outside of his monthly telephone bills, an occasional visit to the Telephone Company commercial office and now and then a brief telephone conversation with the telephone operator, the customer's only contact with the telephone system is his telephone set. He knows that the wires and cables he sees strung on poles and the telephone building he occasionally passes on his way downtown bear some relation to the telephone in his home or office, but this is not very important to him. The important thing is his telephone set. Through it he can communicate with anyone else who has access to a telephone, no matter how far away that person happens to be.

Now what does the user expect from his telephone? We might say that he wants it to be reasonably pleasing in appearance, comfortable to use and simple to operate. He expects accuracy; that is, he wants to be connected to the party he calls and not to some stranger, and he does

not want to be annoyed by answering "wrong numbers". Our typical customer is impatient when his call is unduly delayed without an explanation as to the reason for the delay. He wants a telephone bell or other signal which can be heard in all parts of his home but which is not objectionably loud if he happens to be near the telephone when it rings. He will object to "clicks" or other loud or unpleasant noises from his telephone. He is annoyed if he can hear "crosstalk" which is intelligible or nearly intelligible, because it implies violation of the privacy he expects in his own telephone conversation. He wants to hear in his telephone receiver a reasonably faithful, undistorted reproduction of the voice of the speaker. Finally, the most important attribute a customer expects of his telephone set — important because, without it, he might find no reason to have a telephone at all — is that it permit him to hear and to be heard with a minimum of effort on his part.

It is with this last attribute that the present paper is concerned. Briefly, if a customer hears the voice of his partner in telephone conversation with sufficient loudness and freedom from distortion and noise, little effort will be required to hear and understand what is said. If the distant voice is weak or distorted, the strain of listening requires effort. Likewise, if the distant party hears only with effort he, in turn, may request our typical customer to talk more loudly, which again calls for effort. In telephone systems of modern design, distortion and noise are no longer troublesome factors, so we shall consider only the loudness aspect. Thus, this paper is concerned with the determination of the magnitudes of received volume, or loudness, which give varying degrees of satisfaction to users, and with obtaining an estimate of how well the grade of transmission service provided by the Bell System meets the objective of satisfying customers in this respect.

What the customer would like to have and what it is economically feasible to give him may not be entirely compatible. Giving him what he wants has not always been feasible in the past, but, with the more general application of the improved telephone facilities that have become available over the past decade, this now appears within the realm of possibility.

What do we mean by "what he wants"? We must bear in mind that people vary widely in their judgment of preferred loudness, just as they do in their judgments involving their senses of feeling, sight, taste or smell. This is therefore a statistical problem involving the likes and dislikes of a large number of people. Their combined judgment might be expected to approximate (except at the tails) some distribution related to the normal law. A small number of telephone users will be quite toler-

ant of low received volumes, and another small number at the opposite end of the distribution will desire very loud volumes. But the large bulk of the telephone users will prefer some loudness level about midway between these extremes. It is the job of the engineer to find out by experiment that area of received volumes wherein satisfaction for the greatest number lies.

To carry out such an experiment with the telephone-using public would be a prodigious and impractical undertaking. However, by sampling methods and by closely controlled experimentation one can obtain an answer in the laboratory. The first thing to be done is to decide on a yardstick for expressing degree of satisfaction. Testing procedures are then set up whereby people can listen to different grades of transmission and then express their opinions in terms of the selected yardstick. Measurements of this type are known as "subjective" measurements. The results obtained depend wholly on the judgments of those taking the test. For this reason, large numbers of participants are required to obtain reliable answers.

Laboratory tests of this nature were conducted at Bell Telephone Laboratories over different portions of the total range of received telephone speech volumes at various times from 1947 to 1954. The results of all these tests are combined and summarized here.

II. RESULTS

Two types of subjective transmission tests commonly made in the laboratory are here called "appraisal tests" and "comparison tests". During the appraisal tests observers listen to speech over a telephone connection with specified transmission parameters such as volume, noise interference, etc. The observers are then asked to give their opinions of the transmission qualities of the connection by assigning it to one of several specified categories such as "Good", "Fair", "Poor", or the like. One of the transmission parameters is then changed and the process repeated. Each condition is thus rated on its own merits without direct comparison to any reference condition. In comparison tests, on the other hand, each transmission condition is compared immediately with some other condition, usually a known reference condition, and observers are asked which of the two they prefer in each case. The two types of subjective tests are useful under different circumstances that need not be pursued further here. The point to be noted is that results shown in this paper, except in one case, are based on appraisal tests.

Distributions of appraisal categories in a very wide range of received telephone speech volumes are shown in Fig. 1. The range shown is con-

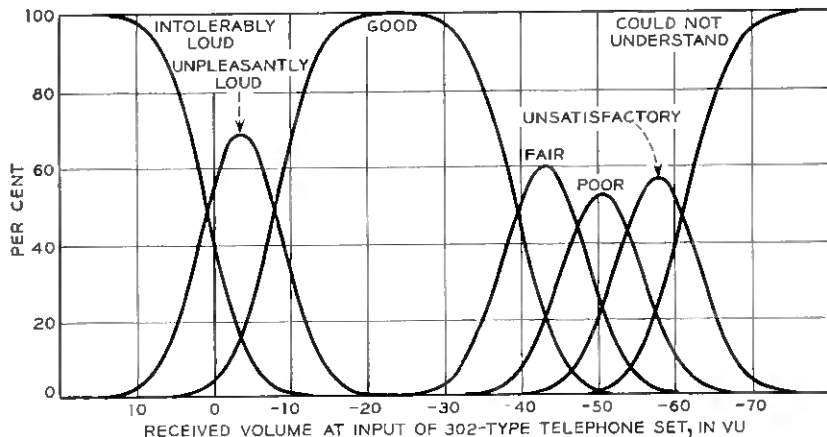


Fig. 1 — Opinion distributions — percentage of observations in which various values of received volume are assigned to the indicated basic categories.

siderably greater than any to be expected in commercial telephone service, since it extends from received volumes which are unbearably loud to those which are practically inaudible. Within this range there are seven natural categories descriptive of the transmission performance of a telephone connection at different levels of received volume in units known as VU. They are here called basic categories, and are as follows:

- Intolerably Loud,
- Unpleasantly Loud,
- Good,
- Fair,
- Poor,
- Unsatisfactory,
- Could Not Understand.

The distribution curves for these categories must not be confused with probability density curves, under which the total area must always equal unity. Instead, each one shows the percentage of total observations at each volume level in which that level is assigned to the indicated category. Since any volume level must be assigned to *some* category, the sum of the percentages of the various categories at any one level must be 100. For example, for a received volume of -10 VU at the line terminals of a 302-type telephone set, 33 per cent of the listeners would say the volume is "Unpleasantly Loud" and 67 per cent would say that it is "Good". For a received volume of -25 VU, everybody would rate the call as "Good". For a received volume of -40 VU, 43 per cent would vote "Good", 50 per cent "Fair" and 7 per cent "Poor".

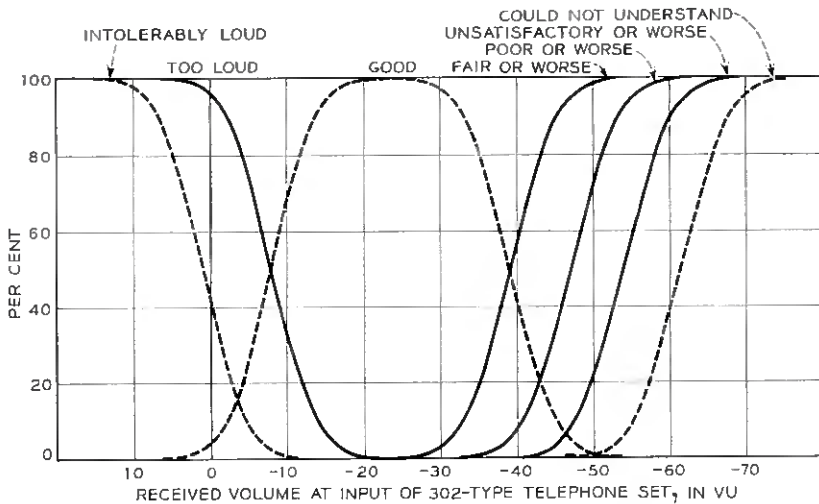


Fig. 2 — Opinion distributions — percentage of observations in which various values of received volume are assigned to the indicated categories. Solid curves are cumulative categories; dashed curves are basic categories.

The various appraisal categories in Fig. 1 are centered about the “Good” category. Curves to the left of “Good” indicate categories which are progressively less acceptable because of undesirably higher received volumes. The categories to the right are progressively less acceptable because of lower and lower received volumes. The two end categories, “Intolerably Loud” and “Could Not Understand” act as barriers because there is no further category on one side of each of them. Thus, each must eventually reach 100 per cent at sufficiently high (or low) volume levels. The distribution of the “Good” category also reaches a maximum of 100 per cent, not because it is a barrier, but because it covers such a wide range of received volumes that there is practically unanimous opinion in the middle of that range, which is around -23 or -24 VU.

The data shown in Fig. 1 can also be plotted in cumulative categories, as shown by the solid curves in Fig. 2. The three dashed curves, “Intolerably Loud”, “Good” and “Could Not Understand” are individual (basic) rather than cumulative categories, copied from Fig. 1 for purposes of orientation. The cumulative categories at the high-volume end include all individual categories to their left, while the cumulative categories at the low-volume end include all individual categories to their right, as may be seen by comparing Figs. 1 and 2. Thus, the percentage value at any volume level for “Too Loud” on Fig. 2 is the sum of the percentage values for “Unpleasantly Loud” and “Intolerably Loud” on

Fig. 1. For example, at -5 VU, 67 per cent of the listeners vote "Unpleasantly Loud" and 8 per cent vote "Intolerably Loud". The percentage in the cumulative category "Too Loud" is therefore 75, as shown in Fig. 2. Likewise, at -45 VU 56 per cent vote "Fair", 30 per cent "Poor" and 3 per cent "Unsatisfactory". Thus, on Fig. 2 "Fair or Worse" is 89 per cent and "Poor or Worse" is 33 per cent. It will be noticed that each of the cumulative categories in turn takes over the functions of "Intolerably Loud" or "Could Not Understand" in acting as an ultimate barrier, and thus eventually reaches a value of 100 per cent.

It would be too much to expect reliable results from tests at one sitting, if observers were required to carry in their heads the large number of basic categories shown by the curves of Figs. 1 and 2. Consequently, the actual tests were made over reduced ranges of received volume, one series in the lower range embracing the five categories from "Good" to "Could Not Understand", inclusive, and another series in the upper range including the three categories from "Good" through "Intolerably Loud".

One is struck also by the tolerance implied by the great width of the "Good" distribution. In addition to being a fortunate fact, it suggests that there must be some smaller range of volumes included within the "Good" category which would be found to be most acceptable of all. This preferred range, or "Excellent" category as it has sometimes been called, has been determined by two methods and is shown in Fig. 3. It was originally determined by appraisal tests similar to those employed in determining the results shown in Figs. 1 and 2. Some years later, the earlier determination was checked by a different technique employing comparison tests. In the comparison tests ten different volume levels (covering the "Good" range of Fig. 1) were compared, each with every other level. Each volume level was then rated according to the percentage of times it was preferred over the other levels with which it was compared.

Fig. 3 shows that the two methods give results which check each other satisfactorily, since the modal points differ by only 1 db. Using the appraisal distribution (to be consistent with the major tests on Figs. 1 and 2), it appears that a received volume level of -19 VU (the modal point of the curve) is the value preferred over all others. It is interesting to look back at Fig. 1 and note that -19 VU occurs, not at the center of the "Good" distribution, but well toward the high-volume side of its table top. This indicates that listeners prefer about the highest volume they can get just short of the point where too much loudness becomes annoying.

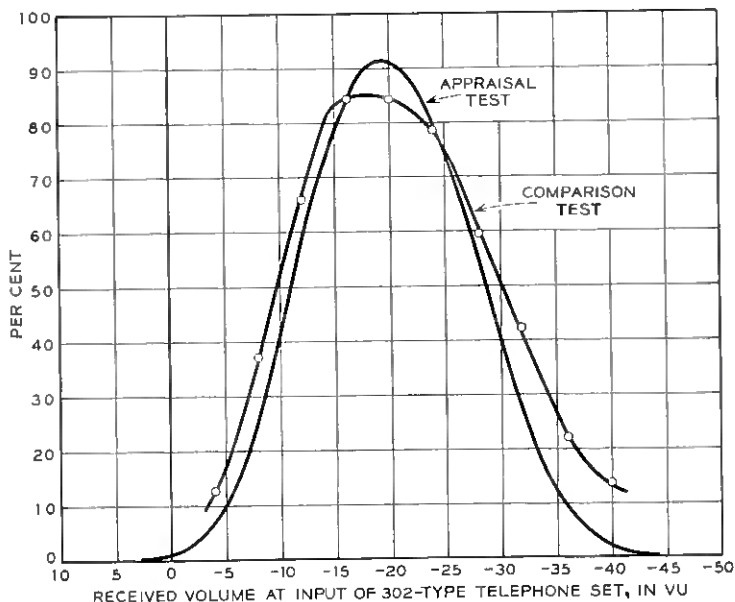


Fig. 3 — Preferred level of received volume, appraisal test versus comparison test. Appraisal test curve shows percentage of observations in which the indicated volume level is assigned to the preferred category; modal point is -19 VU. Comparison test curve shows percentage of observations in which the indicated level is preferred in comparison with other levels tested; modal point is -18 VU.

III. DESCRIPTION OF TESTS

It was recognized that the intelligibility of received telephone speech, and the listener's satisfaction, are affected by transmission factors other than just speech volume, or loudness. Consequently, wherever the effect of other factors might have been significant, they were controlled during the tests at values which cause little or no transmission impairment. The factors which were controlled and their values are as follows:

- line noise: 17 dba at telephone receiver;*
- room noise: 50 db RAP (reference acoustic pressure.);
- speech transmission band: 150 to 3200 cps.

The test circuit is shown schematically in Fig. 4. It is patterned on a typical telephone circuit employing 302-type telephone sets, but with means provided for varying the loss of the trunk so that different levels

* Dba is the unit employed in the Bell System for measurements of line noise with the 2B Noise Measuring Set. This unit was adopted when the 302-type telephone set came into use, in order to provide equal numerical readings of the 2B set when noise of equal impairment was encountered in circuits with different types of telephone set.

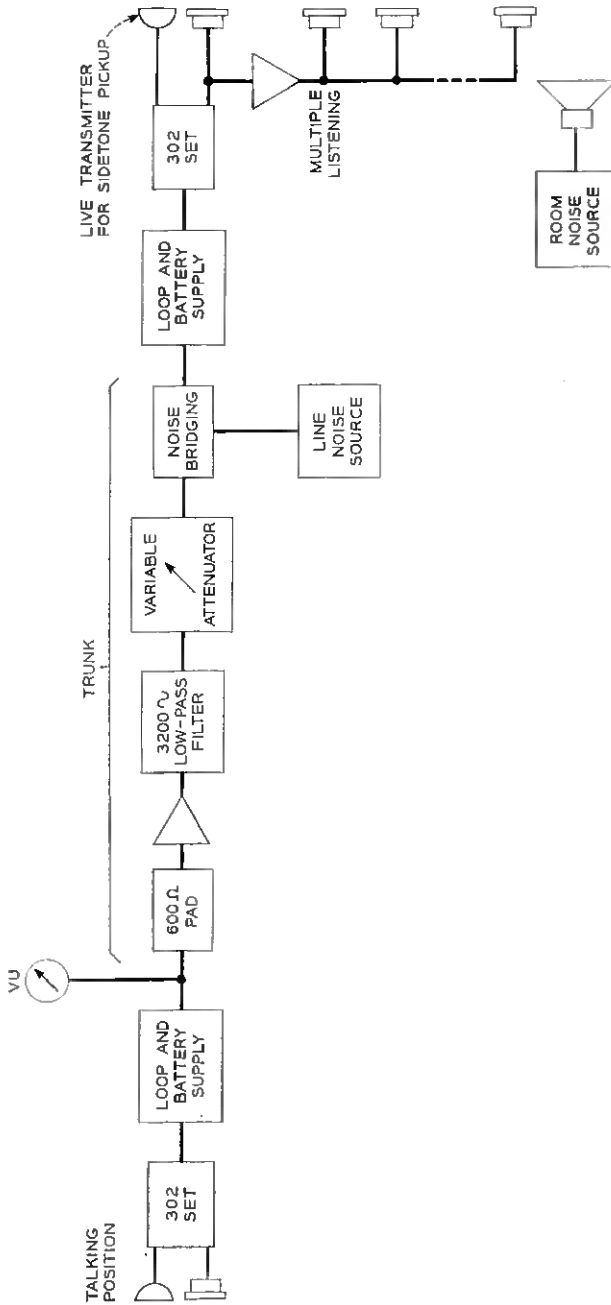


Fig. 4 — Circuit used for appraisal tests.



Fig. 5 — Observers at multiple listening positions during appraisal tests.

of speech volume may be obtained at the receiving end. A volume indicator at the transmitting end of the test circuit enables the talker to keep the level of his own speech constant, thus providing means for determining, from the known losses in between, the level of speech volume at the receiving telephone. The transmitter of the receiving telephone set is kept operative, so that a normal amount of room noise will be present in the receiver through the sidetone path. A multiple listening arrangement is provided so that many observers may listen simultaneously, each hearing the same speech volume and room noise through the sidetone path that are present at the main receiving station. This permits gathering a large amount of data with a minimum expenditure of time and effort on the part of those conducting the tests. A group of observers is shown taking the test at the multiple positions in the photograph in Fig. 5. Fig. 6 shows the board at which the test circuit is set up and controlled, and the main receiving station to which the observers' handsets are multiplied. The talking position does not show in the photographs since it is located in an adjacent soundproofed room.

At first, three men and one woman were used as talkers. It was found that differences in the type of voice (even in the case of the woman talker)



Fig. 6 — Circuit control board and main receiving station for appraisal tests.

had almost no effect on the opinions of listeners, provided that each talker impressed the same speech volume on the line. Consequently, in the later tests almost all of the talking was done by one of the original male talkers.

The test itself was a very simple and straightforward procedure. Observers were allowed to listen to speech over the test circuit at a particular level of received volume. They were then asked to assign that sample to one of the categories previously listed for them. In the case of the tests at the lower levels the categories, as already stated, were "Good", "Fair", "Poor", "Unsatisfactory" and "Could Not Understand", while those in the high-level tests were "Good", "Unpleasantly Loud" and "Intolerably Loud". No attempt was made to define these categories for the observers, each one deciding for himself what the terms meant. When observers had recorded their opinions of the first speech sample, the received volume was changed by adjusting the variable attenuator shown in Fig. 4, and the process was repeated until the entire range had been covered three times. The specific values of received volume used and the ranges covered in the different series of appraisal tests are shown in Table I. The different volume levels were presented to listeners in random order. This is necessary because it has been found that sequential orders (high to low or low to high) result in displacements of the entire opinion distribution, apparently due to conditioning of observers by volume levels

to which they have become accustomed. No disclosure of the magnitudes of the volume levels was made until the completion of a test series.

The process is illustrated by the test sheet shown in Fig. 7, which represents the judgments of one individual observer in one of the tests in the lower range of volumes. This observer recorded his opinions of the 30 different conditions to which he listened in the left-hand column under "TEST 1". Later, the analyst recorded the received volumes corresponding to the 30 conditions and arranged the observer's opinions in descending order of received volumes, as shown in the right-hand portion of the sheet. It will be noticed that this particular observer was very consistent in his judgments (as most were), and that he was among those fairly tolerant of the lowest levels.

Those who participated in the tests were drawn from personnel of the American Telephone and Telegraph Co., Bell Telephone Laboratories and several of the operating telephone companies. Naturally, they represented many facets of the telephone business, but only a small minority could be rated as experts in transmission matters. One group of more than 150 men consisted of young engineers just hired by Bell Laboratories who had practically no telephone experience except as users. The re-

TABLE I — VOLUME LEVELS USED IN VARIOUS APPRAISAL TESTS

VU at Input of Receiving 302-Type Telephone Set		
High-Range Tests	Preferred-Range Tests	Low-Range Tests
+6		
+3	+2	
0		
-3		
-6	-6	
-9	-9	
-12	-12	
-15	-15	
-18	-18	
-21	-21	
	-24	
	-27	
	-30	-30
	-33	
	-36	
	-39	
		-38
		-41
		-43
		-45
	-48	-47
		-49
		-51
		-53
		-55

TEST SUBJECT _____		OBSERVER <u>J. A. O.</u>					
RECORD JUDGMENT AS FOLLOWS:-							
G - GOOD							
F - FAIR							
P - POOR							
U - UNSATISFACTORY							
N - SENTENCE NOT UNDERSTOOD							
CONDITION	TEST 1	ANALYSIS					
		TEST 2	TEST 3	TEST 4			
		REC'D COND. VOL. %	REC'D VOL. VO	1ST 10	2ND 10	3RD 10	
1	G	1	-38		G	G	G
2	G	2	-43	-38	G	F	G
3	P	3	-55	-41	G	F	G
4	F	4	-45	-43	G	F	F
5	P	5	-53	-45	F	F	F
6	P	6	-51	-47	F	F	F
7	G	7	-30	-44	F	F	F
8	F	8	-47	-51	P	P	F
9	F	9	-49	-53	P	P	P
10	G	10	-41	-55	P	P	U
11	P	11	-53				
12	F	12	-49				
13	G	13	-30				
14	P	14	-55				
15	F	15	-47	-30	3	-	-
16	F	16	-41	-38	2	1	-
17	P	17	-51	-41	2	1	-
18	F	18	-38	-43	1	2	-
19	F	19	-45	-45	-	3	-
20	F	20	-43	-47	-	3	-
21	P	21	-53	-49	-	3	-
22	G	22	-30	-51	-	1	2
23	F	23	-49	-53	-	-	3
24	F	24	-45	-55	-	-	2
25	F	25	-51				
26	G	26	-41				
27	U	27	-55				
28	F	28	-47				
29	G	29	-38				
30	F	30	-49				

Fig. 7 — Sample test sheet for appraisal tests.

sults obtained from these various groups were analyzed separately to determine whether different age groups (with corresponding differences in hearing acuity) or differences in transmission background might be a factor in their appraisals of received volume. In general, it was found that there was no significant difference between the opinions of different

groups. For example, the data obtained from the student engineer group alone differed from the data for the entire group by no more than ± 0.5 db.

A sample of the data in cumulative categories is shown in Fig. 8. By plotting the experimental points on arithmetic-probability paper one is able to smooth the data by passing a straight line through the data points. Except for some divergence at the tails of the distributions, the data are a reasonable approximation to the straight lines, and therefore to the normal law of error. The distributions of cumulative categories in Fig. 2 are taken directly from the straight lines of Fig. 8 and from similar distributions for the categories not shown here. The distributions of basic categories in Fig. 1 are obtained by taking vertical differences between the straight lines of Fig. 8 at each level of received volume. Thus, "Unpleasantly Loud" is the difference between "Too Loud" and "Intolerably Loud"

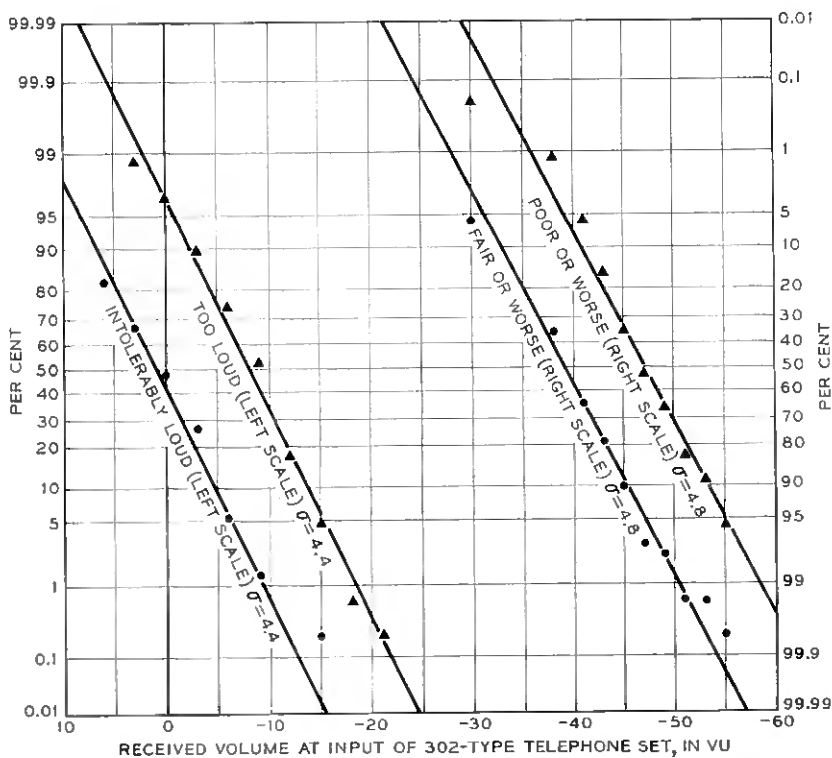


Fig. 8 — Sample of data for opinion distributions — percentage of observations in which various received volumes are assigned to indicated cumulative categories.

Loud", "Fair" is the difference between "Fair or Worse" and "Poor or Worse", and so on.

The reliability of the data is attested by the fact that, for the low-volume categories, they are based on approximately 1400 observations by 470 different listeners at each of the values of received volume indicated in Table I. The data for the high-volume categories are the result of 528 observations by 176 observers at each indicated level.

These two series have been spoken of here as the "major" tests to distinguish them from the tests of the preferred range of received volume levels, the results of which are shown in Fig. 3. The appraisal tests of Fig. 3 were made in the same manner as the major tests (at an earlier date), but were based on fewer observations, 56 at each of the levels indicated on Table I. It was decided that results based on so few observations might be questionable, and that further tests to confirm the earlier results were warranted. Instead of additional appraisal tests, the type of comparison tests which has already been described seemed applicable to the objective of obtaining the preferred volume level, and was adopted.

The talking circuit for the comparison tests was the same as that used in the appraisal tests (Fig. 4). In these tests each of the ten selected volume levels was compared with every other level twice, once with the level A preceding level B in the presentation, and again in the opposite order, B preceding A. However, the two comparisons of like levels were not made consecutively, but were interspersed among other combinations. The comparisons were thus presented to observers in an order that was random both as to the volume levels involved in the comparison and

TABLE II — ANALYSIS OF PREFERRED-RANGE COMPARISON TESTS

Level Preferred	Times That Volume Level in Left Column Was Preferred Over Level in Headings Below										W	L	%
	-4	-8	-12	-16	-20	-24	-28	-32	-36	-40			
-4	—	1	0	0	0	0	0	2	7	6	16	110	12.7
-8	13	—	0	0	1	4	5	6	7	11	47	79	37.3
-12	14	14	—	1	5	7	7	9	14	12	83	43	65.8
-16	14	14	13	—	7	7	13	13	13	12	106	20	84.2
-20	14	13	9	7	—	8	13	14	14	14	106	20	84.2
-24	14	10	7	7	6	—	13	14	14	14	99	27	78.5
-28	14	9	7	1	1	1	—	14	14	14	75	51	59.5
-32	12	8	5	1	0	0	0	—	14	13	53	73	42.1
-36	7	7	0	1	0	0	0	0	—	13	28	98	22.2
-40	8	3	2	2	0	0	0	1	1	—	17	109	13.5
	110	79	43	20	20	27	51	73	98	109			

W (Won) denotes that the level in left column was preferred.

L (Lost) denotes that the level in left column was not preferred.

as to the precedence of one value over the other in a combination. As in the appraisal tests, the magnitudes of received volumes were not disclosed to observers until after the tests were completed.

The method of analysis of the comparison test data is somewhat unusual. If each observed comparison between two volume levels is considered as a game between two teams, the "games won" and "games lost" may be charted in the same manner as the results usually published for baseball leagues for a season in which each team plays the same number of games with every other team. This analysis is shown in Table II, the right-hand column of which shows the percentage of games "won" by each of the ten volume levels, that is, the percentage of observed comparisons in which each of the indicated levels is preferred over the other nine levels with which it is compared. The values from this column are plotted as data points for the comparison test distribution shown in Fig. 3.

IV. APPLICATIONS

The tests which have been described and the curves which show the results of these tests provide fundamental data which, in the authors' opinion, are indicative of a telephone customer's expectancy with regard to hearing and being heard with a minimum of effort on his part. In this respect, they are a measure of what is required to satisfy telephone users. By themselves, the results of the tests indicate observers' opinions (as to category) of any specific value of received volume, such as might be encountered on any one individual telephone connection.

While this is of interest, it is more important to the management of an operating telephone company to know customers' reactions to the grade of service provided over that company's telephone plant, since this information may affect decisions on such questions as spending money for plant improvements. Specifically, the management should know the percentage of telephone connections over the plant which customers consider good, the percentage fair, the percentage poor, etc. Such information may be obtained by combining the opinion distributions described here with the estimated distribution of received volumes which customers actually obtain in their daily use of the telephone. This is done by integrating, over the range of received volumes, the compound probability: (a) that the telephone user considers a particular volume "Good" (or "Fair" or "Poor"), and (b) that he actually receives that volume. It should be pointed out that the distribution of volumes actually received is, in turn, a combination of the distribution of talking volumes, which vary over a range of some 30 VU, and the distribution of plant losses from the point of the talking volume measurement to the input of a listener's telephone

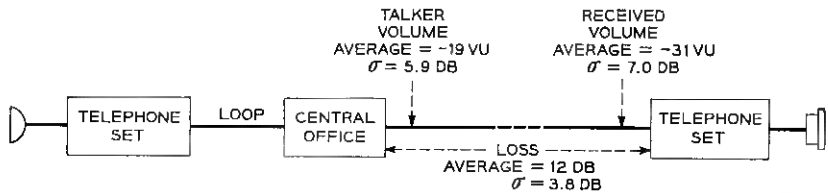


Fig. 9 — Circuit showing points of volume measurement and intervening plant losses.

set. Fig. 9 is a simple diagram which illustrates this point. It shows a transmitting telephone and a loop to the central office. Field measurements of telephone users' talking volume indicated that the average volume input into the line leaving the central office was -19 VU, with a standard deviation of 5.9 db at the time the survey was made.¹ An analysis made several years ago of the large variety of telephone connections encountered in the telephone plant indicated that the average transmission loss between the talker's central office and the distant listener's telephone set was 12 db, with a standard deviation of 3.8 db. The average received volume at the telephone terminals was therefore -31 VU, with a standard deviation of 7.0 db.*

The results of combining the distributions of observers' opinions and actual received volumes have been called "grade of service". "Grade of service" curves may be plotted showing the percentage of telephone connections considered "Good", "Fair", "Poor", etc., against the *average* of various received volume distributions, assuming that these distributions are allowed to vary in average value, but not in standard deviation. In this form they have been found useful in setting objectives for received volumes to be provided and for plant losses which will permit attainment of those received volumes. An example of grade of service curves in cumulative categories is shown in Fig. 10. The same data in another form, which includes individual categories, appear in Fig. 11. For an over-all picture of grade of service provided in the plant just described, assuming that all telephone sets were of the 302 type, the average received volume may be taken as -31 VU, with a standard deviation of 7.0 db, as indicated in Fig. 9. This average value constitutes one point on the abscissae of Figs. 10 and 11. It will be noted in Fig. 11 that, if the average of the distribution of received volumes is -31 VU, the tele-

* In this calculation it is assumed that talker volume and plant loss are independent variables. It is recognized that there may be some small correlation between them. For instance, Subrizi¹ found a small correlation between talker volume and distance between talker and listener on very long distance calls. Other tests to determine correlation have given conflicting results so the correlation is probably small enough to justify the assumption.

-28 VU, rather than at the preferred volume level, -23 VU (for the future plant with 500 sets.) Once we have increased the average value of received volumes to this objective, no further improvement in grade of service may be expected (should it become desirable) except by lowering the standard deviation still further.

V. ACKNOWLEDGMENT

The authors wish to express their gratitude to C. W. Carter of Bell Telephone Laboratories for making available to them his 1947-48 data covering appraisal tests in the preferred range of received volumes. They also extend their thanks to H. R. Huntley, now Chief Engineer of the American Telephone and Telegraph Company; to W. E. Bloecker and the late L. B. Bogan of Mr. Huntley's former department for their support of this program; and, finally, to the many others in the American Telephone and Telegraph Company and Associated Companies, and in Bell Telephone Laboratories, for their conscientious efforts in forming the unbiased judgments which were so necessary for the success of these tests.

REFERENCES

1. Subrizi, V., A Speech Volume Survey on Telephone Message Circuits, Bell Labs. Record, **31**, August 1953, p. 292.
2. MacAdam, W. K., A Basis for Transmission Performance Objectives in a Telephone Communication System, Comm. & Elect., No. 36, May 1958, p. 205.

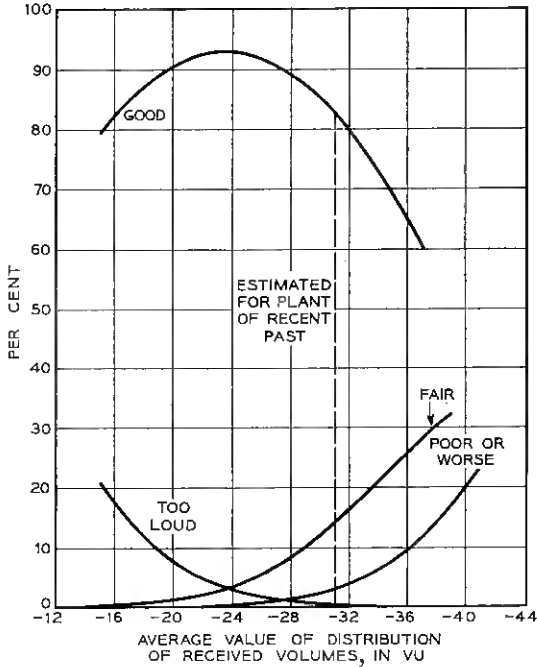


Fig. 11 — Grade of service distributions provided by plant of recent past — percentage of telephone connections assignable to various categories when the average of the received volume distribution has the indicated value.

is the same as that of the corresponding opinion curve on Fig. 8. By using arithmetic-probability paper, this simple construction of Fig. 10 replaces laborious summations of the compound probabilities. The summations cannot be avoided, however, if the distribution of received volumes is other than normal law.

For applications such as these, some reminders concerning the nature of the fundamental data presented here will not be out of place. Observers' opinions as to the category of specific values of received volume have been found to vary with (a) the amount of interference present in the form of line noise and ambient room noise at the listener's location, (b) the upper cutoff frequency of the telephone message channel and (c) the sensitivity of the listening telephone set and receiver, both for speech and for noise. As indicated early in this paper, data shown in Figs. 1 and 2 apply only to conditions of line noise, room noise and transmission bandwidth which have been found to cause little or no transmission impairment. Line noise and transmission bandwidth are controllable fac-

tors, and the values used are those which have been worked to as objectives in the Bell System for many years. Room noise is not under control of the telephone companies, but the value used is one that has been found representative of room noise in the average residence or fairly quiet office. In addition, the data of Figs. 1 and 2 apply only to receiving telephone sets of the 302 type, the type in greatest use in the telephone plant of the early 1950's, but now being largely supplanted by the 500 type. While all the conditions mentioned were reasonably normal a few years ago, it must be borne in mind that changes are being made over the years which necessitate adjustment of the opinion distributions with respect to the VU scale.

Long-range objectives for received loudness in the Bell System have been the subject of much study in recent years. W. K. MacAdam, Transmission Engineer of the American Telephone and Telegraph Company, has stated² that the design objectives of the Bell System might be about as follows:

- i. A negligible number of calls rated "Poor".
- ii. No more than 5 per cent rated "Fair".
- iii. The balance rated "Good" or "Excellent".

It is evident that the telephone plant of a few years ago did not fully meet this objective. However, the picture is altered in the case of the future plant. If modern sets of the 500 type having higher receiving sensitivity are employed, the opinion curves will be found to shift with respect to the VU scale by approximately the amount of the sensitivity difference.* Thus, a value around -23 VU, rather than -19 , would represent the volume level preferred over all others. In addition, it has been found that, with the 500-type telephone set, although the average talker volume output for a given level of acoustic input is increased, the spread in volume is decreased (smaller standard deviation) because there is a certain amount of speech compression on higher volumes. Furthermore, planned improvements in telephone lines (more precisely, the transmission medium between telephone sets) are expected to result in lower losses and lower standard deviation. It should be noted that, in the region of received volumes in which we are interested, lower standard deviation contributes to improved grade of service, just as higher average values of received volume do.

Based on these and other considerations, estimates have been made of the grade of service which will be provided by the future telephone plant, assuming that it includes the 500-type telephone set preponderantly,

* Assuming the total noise reaching the user's ear remains unchanged. Because of improvements in the sidetone circuit this is roughly the case.

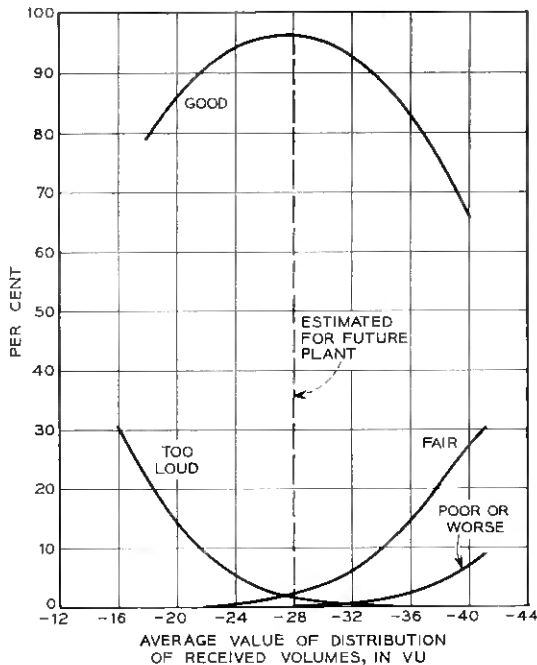


Fig. 12 — Grade of service distributions estimated for the future plant — percentage of telephone connections assignable to various categories when the average of the received volume distribution has the indicated value.

and the contemplated line improvements. Distribution curves based on these estimates are shown in Fig. 12. It is estimated that the average received volume resulting from the improvements mentioned will be -28 VU, as indicated in this figure. The percentage of calls rated "Good" would then exceed the objective of 95, while the calls rated "Fair" would come well within the 5 per cent objective. This is a satisfactory outlook for the future.

We then shall have gone about as far in the direction of increasing received volumes as we should. Fig. 12 shows that, at an average received volume of -28 VU, the percentage of calls considered "Too Loud" is about 1.5. Any further increase in received volume might raise the percentage "Too Loud" to undesirable values, with a corresponding decrease in the percentage of calls rated "Good". The fact that the "Good" distribution turns down at higher volume levels is the reason why it is advisable to set an objective for the average value of the received volume distribution at a volume level near the peak of the "Good" distribution, around

-28 VU, rather than at the preferred volume level, -23 VU (for the future plant with 500 sets.) Once we have increased the average value of received volumes to this objective, no further improvement in grade of service may be expected (should it become desirable) except by lowering the standard deviation still further.

V. ACKNOWLEDGMENT

The authors wish to express their gratitude to C. W. Carter of Bell Telephone Laboratories for making available to them his 1947-48 data covering appraisal tests in the preferred range of received volumes. They also extend their thanks to H. R. Huntley, now Chief Engineer of the American Telephone and Telegraph Company; to W. E. Bloecker and the late L. B. Bogan of Mr. Huntley's former department for their support of this program; and, finally, to the many others in the American Telephone and Telegraph Company and Associated Companies, and in Bell Telephone Laboratories, for their conscientious efforts in forming the unbiased judgments which were so necessary for the success of these tests.

REFERENCES

1. Subrizi, V., A Speech Volume Survey on Telephone Message Circuits, Bell Labs. Record, **31**, August 1953, p. 292.
2. MacAdam, W. K., A Basis for Transmission Performance Objectives in a Telephone Communication System, Comm. & Elect., No. 36, May 1958, p. 205.

Recent Monographs of Bell System Technical Papers Not Published in This Journal*

BALASHEK, S., see Dudley, H.

BEMSKI, G. and STRUTHERS, J. D.
Gold in Silicon, Monograph 3184.

BOYLE, W. S. and NOZIÈRES, P.
Band Structure and Infrared Absorption of Graphite, Monograph
3177

BRADY, G. W.
Structure of Tellurium Oxide Glass, Monograph 2878.

BRESLIN, J., see Kaiser, W.

CHYNOWETH, A. G.
Barkhausen Pulses in Barium Titanate, Monograph 3165.

CHYNOWETH, A. G. and PEARSON, G. L.
Effect of Dislocations on Breakdown in Silicon p-n Junctions, Mono-
graph 3170.

CLOSSON, H. T., DANIELSON, W. E. and NIELSEN, R. J.
Automatic Measurement of Small Deviations in Periodic Structures,
Monograph 3187.

DANIELSON, W. E., see Closson, H. T.

* Copies of these monographs may be obtained on request to the Publication Department, Bell Telephone Laboratories, Inc., 463 West Street, New York 14, N. Y. The numbers of the monographs should be given in all requests.

- DECOSTE, J. B., HOWARD, J. B., WALLDER, V. T. and ZUPKO, H. M.
Plasticized Poly (Vinyl Chloride) for Retractable Cords, Monograph 3186.
- DEVLIN, G. E., see Schwalow, A. L.
- DILLON, J. F., JR.
Observation of Domains in the Ferrimagnetic Garnets by Transmitted Light, Monograph 3183.
- DORSI, D., see Wernick, J. H.
- DUDLEY, H. and BALASHEK, S.
Automatic Recognition of Phonetic Patterns in Speech; Phonetic Vocoder, Monograph 3172.
- FRISCH, H. L., see Liehr, A. D.
- FULLER, C. S. and WHEELAN, J. M.
Diffusion, Solubility, and Electrical Behavior of Copper in Gallium Arsenide, Monograph 3176.
- GEBALLE, T. H., see Herring, C.
- GIBBONS, D. F.
Acoustic Relaxations in Ferrite Single Crystals, Monograph 2882.
- GILLES, M. A.
Superexchange Interaction Energy for Fe^{3+} — O^{2-} — Fe^{3+} Linkages, Monograph 2951.
- GLASS, M. S.
Distribution of Leakage Flux Around a TWT-Focusing Magnet, Monograph 3078.
- GORDON, J. P. and WHITE, L. D.
Noise in Maser Amplifiers — Theory and Experiment, Monograph 3178.

GROSSMAN, A. J.

Synthesis of Tchebycheff Parameter Symmetrical Filters, Monograph 2964.

HEIDENREICH, R. D., see Miller, R. C.

HEISS, J. H. and LANZA, V. L.

Thermal Embrittlement of Stressed Polyethylene, Monograph 3188

HERRING, C., GEBALLE, T. H. AND KUNZLER, J. E.

Phonon-Drag Thermomagnetic Effects in n-Type Germanium, Monograph 3171.

HOBSTETTER, J. N., see Wernick, J. H.

HORTON, A. W., JR.

An Introduction to Computer Binary Arithmetic, Monograph 3050.

HOWARD, J. B., see DeCoste, J. B.

JAHN, A. P. and VACCA, G. N.

Accelerated Aging Tests and Service Performance of Neoprene Jacketed Drop Wire, Monograph 3181.

KAISER, W. and BRESLIN, J.

Factors Determining Oxygen Content of Liquid Silicon at Its Melting Point, Monograph 3179.

KAPLAN, E. L. and MEIER, P.

Nonparametric Estimation from Incomplete Observations, Monograph 3160.

KUNZLER, J. E., see Herring, C.

LANZA, V. L., see Heiss, J. H.

LAX, M., see Levitas, A.

LEVITAS, A. and LAX, M.

Statistics of the Ising Ferromagnet, Monograph 3167.

LIEHR, A. D. and FRISCH, H. L.

Dynamical Stability Criteria for Molecular Motions, Monograph 3161.

LOVELL, L. C., see Wernick, J. H.

MAY, J. E., JR.

Precise Measurement of Time Delay, Monograph 3084.

MEIER, P., see Kaplan, E. L.

MEITZLER, A. H.

Temperature and Frequency Dependence of Insertion Loss in Delay Lines, Monograph 3083.

MILLER, R. C.

Some Experiments on the Motion of 180° Domain Walls in BaTiO_3 , Monograph 3174.

MILLER, R. C. and HEIDENREICH, R. D.

Interaction of Low-Energy Electrons with Ferroelectric Materials, Monograph 3164.

MITRA, S. S., see Wood, D. L.

MOLL, J. L., see Senitzky, B.

MORGAN, S. P.

General Solution of the Luneberg Lens Problem, Monograph 3182.

NIELSEN, R. J., see Closson, H. T.

NOZIÈRES, P., see Boyle, W. S.

PEARSON, G. L., see Chynoweth, A. G.

PEARSON, G. L. and TREUTING, R. G.

Surface Melt Patterns on Silicon, Monograph 3133.

SCHAWLOW, A. L. and DEVLIN, G. E.

Intermediate State of Superconductors: Influence of Crystal Structure, Monograph 3162.

SENITZKY, B. and MOLL, J. L.

Breakdown in Silicon, Monograph 3189.

SLEPIAN, D.

Some Comments on the Detection of Gaussian Signals in Gaussian Noise, Monograph 3163.

STADLER, H. L.

Ferroelectric Switching Time of BaTiO₃ Crystals at High Voltages, Monograph 3185.

STEPHENS, S. J.

Chemisorption and Surface Reactions of Ethylene on Evaporated Palladium Films, Monograph 3159.

STRUTHERS, J. D., see Bemski, G.

SYKES, R. A.

New Approach to the Design of High Frequency Crystal Filters, Monograph 3180.

THURSTON, R. N. and TORNILLO, L. M.

Coiled Wire Torsional Wave Delay Line, Monograph 3085.

TIEN, P. K.

Parametric Amplification and Frequency Mixing in Propagating Circuits, Monograph 3106.

TORNILLO, L. M., see Thurston, R. N.

TREUTING, R. G., see Pearson, G. L.

TURNER, E. H.

A Fast Ferrite Switch for Use at 70 KMC, Monograph 3169.

UNGER, S. H.

A New Type of Computer Oriented Towards Spatial Problems, Monograph 3080.

VACCA, G. N., see Jahn, A. P.

WAHL, A. J.

Three-Dimensional Analytic Solution for Alpha of Alloy Junction Transistors, Monograph 2928.

WALLDER, V. T., see DeCoste, J. B.

WARNER, R. M., JR.

A New Passive Semiconductor Component, Monograph 3082.

WERNICK J. H., HOBSTETTER, J. N., LOVELL, L. C. and DORSI, D.

Dislocation Etch Pits in Antimony, Monograph 3168.

WERTHEIM, G. K.

Electron Bombardment Damage in Silicon, Monograph 3166.

WHEATLEY, G. H., see Whelan, J. M.

WHELAN, J. M., see Fuller, C. S.

WHELAN, J. M. and WHEATLEY, G. H.

Preparation and Properties of Gallium Arsenide Single Crystals, Monograph 3175.

WHITE, L. D., see Gordon, J. P.

WOOD, D. L., and MITRA, S. S.

Effect of Convergence on the Infrared Spectra of Anisotropic Substances, Monograph 3173.

ZUPKO, H. M., see DeCoste, J. B.

Contributors to This Issue

M. R. AARON, B.S.E.E., 1949, and M.S.E.E., 1951, University of Pennsylvania; Bell Telephone Laboratories, 1951—. He first worked on analysis, design and synthesis of transmission networks for L3 and submarine cable systems. From 1954 to 1956 he supervised a group concerned with design of networks for the L3 system. Since 1956 he has been in charge of a group engaged in systems analysis of PCM. Member I.R.E., Eta Kappa Nu, Tau Beta Pi, Sigma Tau.

M. M. ATALLA, B.S., 1945, Cairo University (Egypt); M.S., 1947, and Ph.D., 1949, Purdue University; Bell Telephone Laboratories, 1950—. For five years Mr. Atalla headed a group engaged in basic studies in contact physics. He is now in charge of a group carrying out fundamental studies in surface physics of semiconductors and applications to devices. Member American Physical Society, Sigma Xi, Sigma Pi Sigma, Pi Tau Sigma.

OLIVER H. COOLIDGE, A.B., 1922, Harvard College; New York Telephone Company, 1921-27; American Telephone and Telegraph Company, 1927-34; Bell Telephone Laboratories, 1934—. After six years in plant maintenance methods work with the New York Telephone Company, he joined A.T.&T.'s Development and Research Department, which was later transferred to Bell Laboratories. At that time he was engaged in experimental and theoretical work in the field of transmission engineering, especially in problems of crosstalk and noise interference. During World War II he served as a radar maintenance instructor in Bell Laboratories' School for War Training. Since 1949 he has been concerned with problems of quality and standards of local transmission, and more recently with general transmission objectives.

T. H. GEBALLE, B.S., 1940, and Ph.D., 1950, University of California; Bell Telephone Laboratories, 1952—. Mr. Geballe has specialized in solid state research, with special interest in the study of mechanisms involving the transport of heat and electricity by crystalline semiconductors. At present he is in charge of a group in the Physical Research Department which is studying a variety of fundamental properties. Member American Physical Society, American Chemical Society, Phi Beta Kappa, Sigma Xi.

W. J. GRUBBS, B.S.E.E., 1951, University of Kentucky; Bell Telephone Laboratories, 1951—. After completing rotational assignments in the C.D.T. Program, Mr. Grubbs was engaged in design and development of ferrite core inductors, especially for a rural carrier telephone system. Recently he has been engaged in fundamental development of solid state devices and applications of the Hall effect. Member Eta Kappa Nu, Tau Beta Pi.

CONYERS HERRING, A.B., 1933, University of Kansas; Ph.D., 1937, Princeton University; National Research Council Fellow, Massachusetts Institute of Technology, 1937-39; research associate, Princeton, 1939-40; instructor in physics, University of Missouri, 1940-41; Division of War Research, Columbia University, 1941-45; professor of applied mathematics, University of Texas, 1946; Bell Telephone Laboratories, 1945—. Mr. Herring has specialized in theoretical physics of the solid state. He was a member of the Institute for Advanced Study at Princeton in 1952-53, and was awarded the 1959 Oliver E. Buckley Solid State Physics Prize for "his interpretation of the transport properties of semiconductors". Fellow American Physical Society; member American Association for the Advancement of Science.

J. E. KUNZLER, B.S., 1945, University of Utah; Ph.D., 1950, University of California; research associate, University of California, 1950-52; Bell Telephone Laboratories, 1952—. Mr. Kunzler has been engaged in low-temperature solid state research. He was concerned with the design and establishment of a thermodynamics and cryogenics laboratory and has been engaged in the investigation of electrical, thermal and magnetic properties of solids. Member American Physical Society, American Chemical Society, Sigma Xi, Tau Beta Pi, Alpha Chi Sigma.

J. A. NARUD, B.S. and M.S., 1951, California Institute of Technology; Ph.D., 1955, Stanford University. Mr. Narud is an assistant professor at Harvard University and, since 1957, has been a consultant to Bell Telephone Laboratories, where he has been working on various problems in connection with PCM. He is also engaged in studies of properties of nonlinear feedback networks and applications to control systems and pulse circuitry. Member I.R.E., Sigma Xi.

G. C. REIER, A.B., 1913, Washington College; B.S. in Engineering, 1916, Johns Hopkins University; American Telephone and Telegraph Company, 1916-34; Bell Telephone Laboratories, 1934-58. As an engineer with the A.T.&T. Co., Mr. Reier took part in fundamental studies of electrical wave filters and studies of speech and transmission quality. He was also concerned with transmission problems relating to

central office equipment, loading systems and telephone sets. After transferring to Bell Laboratories he was concerned with general transmission problems and later with air raid warning systems and other civil defense projects. From 1952 until his retirement in 1958, he was project engineer on a systems study for the U.S. Navy.

E. J. SCHEIBNER, B.S., 1950, Georgia Institute of Technology; M.S., 1952, and Ph.D., 1955, Illinois Institute of Technology; Bell Telephone Laboratories, 1955-59; research associate professor of physics, Georgia Institute of Technology, 1959—. At Bell Laboratories Mr. Scheibner was engaged in physical measurements of semiconductor surface properties and studies of single-crystal surfaces by electron diffraction. Member American Physical Society, Sigma Xi, Tau Beta Pi.

CLAUDE E. SHANNON, B.S.E.E., 1936, University of Michigan; Ph.D., 1940, Massachusetts Institute of Technology; Bell Telephone Laboratories, 1941—; professor of communications sciences and mathematics, M.I.T., 1956—. At Bell Laboratories Mr. Shannon has specialized in mathematical research on communication theory and computing machines and automata. He has made outstanding contributions to the communications field, especially in the mathematical theory of communication. In 1956 he was granted a leave of absence from Bell Laboratories to return to M.I.T. as visiting professor in electrical communications. He became a permanent member of the M.I.T. faculty in 1957, while continuing his association with the Laboratories as mathematical consultant. In 1957-58 he was a fellow at the Center for Advanced Study in the Behavioral Sciences. In October 1958 he was appointed to the newly established Donner Chair of Science at M.I.T. Mr. Shannon has been awarded the Alfred Noble Prize of the American Institute of Electrical Engineers, the Morris Liebmann Award of the Institute of Radio Engineers, the Stuart Ballantine Medal of the Franklin Institute and the Research Corporation Award. Member National Academy of Sciences, American Academy of Arts and Sciences, American Mathematical Society, Institute of Radio Engineers, Sigma Xi, Phi Kappa Phi, Eta Kappa Nu, Tau Beta Pi.

EILEEN TANNENBAUM, B.A., 1950, and M.A., 1952, Mount Holyoke College; Ph.D., 1955, University of California; Bell Telephone Laboratories, 1956—. Miss Tannenbaum has been engaged in work on solid state devices and basic research in surface studies of semiconductors, especially with regard to transistors. She has been awarded several fellowships for graduate and postgraduate work. Member American Physical Society.

