# Transport Properties of a Many-Valley Semiconductor

### By CONYERS HERRING

*The simple model of a semiconductor, based on a single effective mass for the charge carriers and a spherical shape for the surfaces of constant energy, is now known to be inadequate for most of the semiconductors which have been extensively studied experimentally. However, some of these do correspond to what may be called the "many-valley" model, a model for which the band edge occurs at a number of equivalent points $\mathbf{K}^{(i)}$ in wave number space, and for which the surfaces of constant energy are multiple ellipsoids, one centered on each of these points. This paper develops, for models of this type, the theory for: mobility (Section 2) and its temperature dependence (Section 3); thermoelectric power (Section 4); piezoresistance (Section 5); Hall effect (Sections 6 and 9); high-frequency dielectric constant (Section 7); and magnetoresistance (Sections 8 and 9). These phenomena are treated, for cases to which Maxwellian statistics apply, on the assumption that the scattering of the charge carriers is describable by a relaxation time which depends on energy only, but is otherwise unrestricted. This assumption can be shown to be justified in a large class of cases, although for some cases it fails, notably when ionized impurity scattering predominates and at the same time the effective mass is very anisotropic. Special attention is given to the role of inter-valley lattice scattering, i.e., to processes whereby a charge carrier is scattered from the neighborhood of one of the band edge points*

237

$\mathbf{K}^{(i)}$ *to the neighborhood of a different one. Numerical calculations are presented which show the effects of such processes on the magnitudes and temperature variations of the effects listed above.*

### 1. THE MANY-VALLEY MODEL

Most of the literature of semiconductor theory has been based on what we shall call the simple model. This model is based on the assumption that the minimum energy in the conduction band, or the maximum energy in the valence band, is possessed by only one quantum state of either spin. This state has the form of a Bloch wave with wave number $\mathbf{K} = 0$.* States with energies near the band edge value therefore have small $K$ values, and, since their energies $\epsilon(\mathbf{K})$ must vary continuously with $\mathbf{K}$ in this region, $\epsilon(\mathbf{K})$ for small $K$ must be a quadratic form in $K_x$, $K_y$, $K_z$. If the crystal structure is cubic, $\epsilon(\mathbf{K}) \propto K^2$, and the surfaces of constant energy are spheres in $\mathbf{K}$-space.

It has long been known that other models are possible, and indeed likely in many cases. In recent years it has become clear that the simple model does not apply to *any* of the four cases corresponding to n- and p-type germanium and silicon. The evidence for this includes magnetoresistance[1, 2, 3] and piezoresistance[4] effects, cyclotron resonances,[5] and many other phenomena. Now the possible alternatives to the simple model are the various models for which there is more than one state, apart from spin degeneracy, with the band edge energy. These models fall into two general categories.

(A) Models for which the band edge energy occurs for several wave number vectors $\mathbf{K}^{(i)}$, but for which there is only one state of each spin having this energy and a given $\mathbf{K}^{(i)}$. For a conduction band model of this sort the energy $\epsilon$, considered as a function of $\mathbf{K}$, has a number of minima or "valleys", hence we shall call these "many-valley" or "simple

---

* For the convenience of the reader the notations defined in the text are recapitulated on page 288.

[1] I. Estermann and A. Foner, Phys. Rev., **79**, p. 365, 1950; G. L. Pearson and H. Suhl, Phys. Rev., **83**, p. 768, 1951; and G. L. Pearson and C. Herring, Physica, to appear.

[2] W. Shockley, Phys. Rev., **78**, p. 173, 1950, and unpublished work.

[3] S. Meiboom and B. Abeles, Phys. Rev., **93**, p. 1121, 1954; B. Abeles and S. Meiboom, Phys. Rev., **95**, p. 31, 1954; and M. Shibuya, J. Phys. Soc., Japan, **9**, p. 134, 1954 and Phys. Rev., **95**, 1385, 1954.

[4] C. S. Smith, Phys. Rev., **94**, p. 42, 1954.

[5] G. Dresselhaus, A. F. Kip, and C. Kittel, Phys. Rev., **92**, p. 827, 1953; B. Lax, H. J. Zeiger, R. N. Dexter, and E. S. Rosenblum, Phys. Rev., **93**, p. 1418, 1954; R. N. Dexter, H. J. Zeiger, and B. Lax, Phys. Rev., **95**, p. 557, 1954; R. N. Dexter, B. Lax, A. F. Kip, and G. Dresselhaus, Phys. Rev., **96**, p. 222, 1954; and R. N. Dexter and B. Lax, Phys. Rev., **96**, p. 223, 1954.

many-valley" models. For a valence band the situation is similar, but inverted.

(B) Models for which, apart from spin degeneracy, there are two or more states with the band edge energy and the same wave number vector. These we shall call "degenerate" models. We may subdivide them further into "degenerate single-valley" and "degenerate many-valley" cases according to whether the band edge energy occurs for only one wave vector, or for several.

This paper is concerned with the transport properties of the simple many-valley models defined under (A). These models are much simpler to handle than the degenerate types, for reasons which are illustrated in Fig. 1. This illustration shows schematically the form in wave number space of the surfaces of constant energy, near the band edge energy, for four models. For the simple model, shown in (a), the locus of a given energy is, as already stated, a sphere. For a simple many-valley model, shown in (b), the locus is a set of ellipsoids centered about the band edge points $\mathbf{K}^{(i)}$. The ellipsoidal shape is required by the facts that energy must depend continuously and differentiably on $\mathbf{K}$ and have an extremum at each $\mathbf{K}^{(i)}$. For a degenerate model, however, the dependence of energy on the components of $\mathbf{K}$ is singular at the band edge point,[2] in that unique second derivatives do not exist: energy varies quadratically with $\mathbf{K}$ in any given direction from this point, but the coefficients going with different directions are determined by a secular equation. The result is that the contours of constant energy may look as shown in Fig. 1(c) (degenerate single-valley case). Degenerate many-valley cases are of course similar, but with the surfaces multiplied, as in Fig. 1(d). Such situations are obviously harder to handle mathematically than those of Fig. 1(b).

Besides the irregularity of the energy surfaces, there is another difference between these two types of cases which greatly complicates theoretical work with degenerate models. This is that in most cases the energies of the two or more states going with a given band-edge $\mathbf{K}$ will be split by spin orbit coupling. If this splitting is $\ll kT$ it can usually be ignored, and if it is $\gg kT$ it may effectively convert the degenerate model into a simple or simple many-valley model. Unfortunately, it usually happens that neither of these extremes applies, and for such intermediate cases not only do we have to deal with energy surfaces like those of Fig. 1(a), but, much worse, the variation of energy with $\mathbf{K}$ is not a simple quadratic dependence even in a fixed direction from the band edge point.

In view of all these complications of the degenerate models, it is for-

tunate that the simple many-valley case does seem to occur for n-type silicon and n-type germanium. One of the best ways of finding out whether it occurs for any given semiconductor is to compare observations on this material with theoretical predictions for the various possible models of the simple many-valley type. We proceed now to derive these predictions, assuming for simplicity that the charge carriers have Maxwellian statistics and an effective relaxation time dependent only on energy. We shall take up the simplest properties first, the more complicated ones later. Sections 2 to 5 will consider perturbation of the distribution function of the carriers by a static electric field, Section 6
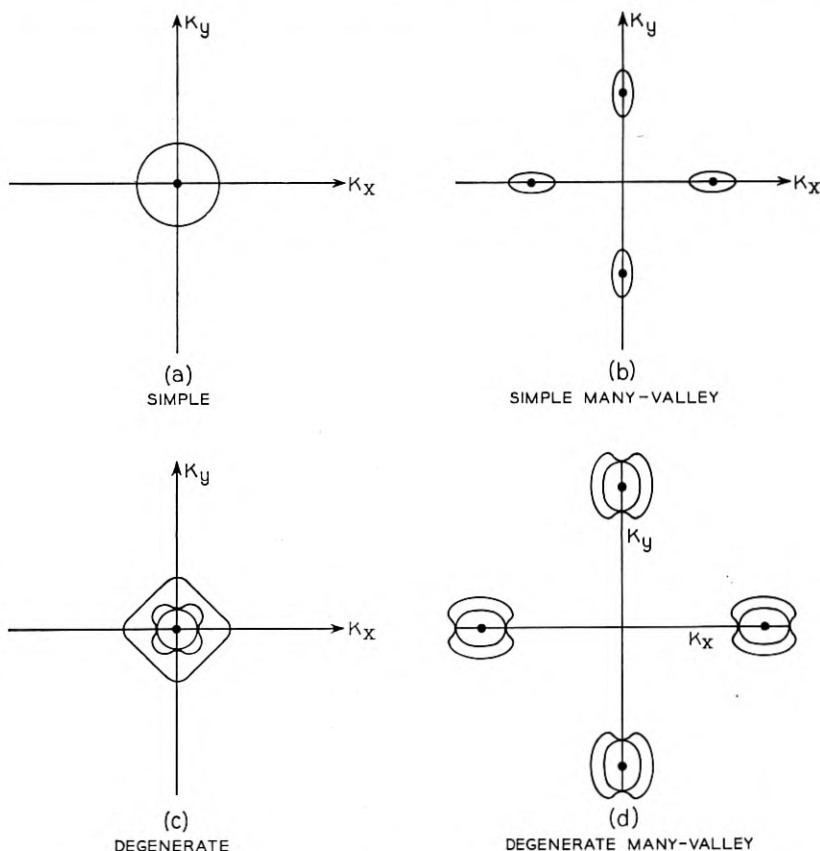


Fig. 1 — Different types of band structure for a semiconductor, illustrated by the forms of the surfaces of constant energy in wave number space. The band edge points are represented by heavy dots.

the Hall effect, Section 7 the perturbation by an oscillating electric field, Section 8 magnetoresistance, and Section 9 effects at high magnetic fields.

In presenting this material our primary objective will be to provide a coherent treatment of all the effects in language as simple and physical as possible. Thus, for example, the Hall and magnetoresistance effects will be discussed *ab initio*, although many of the details presented here have been derived and published independently by several workers.[2, 3] Nor is the theory of this paper the ultimate in refinement: at cost of a little more mathematical complication, the assumption of a relaxation time dependent only on energy can be dispensed with, and anisotropy in the scattering processes acting on the carriers can be taken into account.[6] However, the present simpler treatment illustrates most of the physical principles involved in the various phenomena, and turns out to be quantitatively adequate in a large class of cases.

## 2. CONDUCTIVITY

In this section we shall solve the Boltzmann equation for the effect of a constant electric field $\mathbf{E}$ on the motion of charge carriers in a simple many-valley band. Maxwellian statistics will be assumed. Thus if $\Delta \mathbf{P} = \hbar \, (\mathbf{K} - \mathbf{K}^{(i)})$ measures the deviation in crystal momentum space from one of the band edge points $\mathbf{K}^{(i)}$, then for $\mathbf{E} = 0$ the probability of occupation of the state described by $\Delta \mathbf{P}$ (by an electron or hole, depending on the sign of the carriers) is

$$f^{(0)} = \exp \left[ \frac{ -| \, \epsilon_F - \epsilon_b \, | - \dfrac{\Delta P_1^{\,2}}{2m_1^*} - \dfrac{\Delta P_2^{\,2}}{2m_2^*} - \dfrac{\Delta P_3^{\,2}}{2m_3^*} }{ kT } \right] \tag{1}$$

where $\epsilon_F$ is the Fermi level, $\epsilon_b$ the band edge energy and $m_1^*$, $m_2^*$, $m_3^*$ are the effective masses in the three coordinate directions 1, 2, 3 which are principal axes for the energy surfaces of the valley in question. When $\mathbf{E} = \neq 0$ the distribution function $f$ is determined by the competition between the perturbing effect of $\mathbf{E}$ and the restoring effect of scattering processes which try to restore the form (1). To make the problem tractable we shall assume that the scattering processes which the charge carriers undergo are described by a relaxation time which is a function of energy $\epsilon$ only. In other words, we shall assume that for any slight

---

[6] C. Herring and E. Vogt, to be published.

departure of $f$ from $f^{(0)}$ the time rate of change of $f$ due to collisions is

$$\left(\frac{df}{dt}\right)_c = -\frac{(f - f^{(0)})}{\tau(\epsilon)} \tag{2}$$

The legitimacy of this assumption is analyzed in Appendix A. It is shown there that the assumption should be rather accurately valid for all kinds of inter-valley scattering — defined as scattering from the neighborhood of one band edge point $\mathbf{K}^{(i)}$ to the neighborhood of another $\mathbf{K}^{(j)}$ — and for intra-valley lattice scattering due to optical modes or to neutral impurities, provided, in the latter case, that the temperature is low enough. For intra-valley lattice scattering due to acoustical modes the assumption $\tau = \tau(\epsilon)$ is not necessarily valid, but the arguments of Appendix A suggest tht it will often be a good approximation. For scattering by ionized impurities, however, this assumption will usually be a poor approximation. There is a good prospect that in the near future the adequacy of this approximation for lattice scattering can be quantitatively estimated for some substances. If it should turn out to be inadequate, the necessary generalization of the calculations of this paper can probably be made without great effort.

With the assumptions just stated in (1) and (2), the Boltzmann equation for a steady state in the presence of a field $\mathbf{E}$ takes the form

$$0 = \frac{\partial f}{\partial t} = \pm e\mathbf{E} \cdot \nabla_P f - \frac{(f - f^{(0)})}{\tau} \tag{3}$$

where the upper sign is for electrons in a conduction band, the lower for holes in a filled band. If, as is customary, we set

$$f = f^{(0)} + \mathbf{E} \cdot \mathbf{f}^{(1)} + O(E^2), \tag{4}$$

(3) gives, just as in the simple theory,

$$\mathbf{f}^{(1)} = \pm e\tau \nabla_P f^{(0)} \tag{5}$$

Having obtained the solution of the Boltzmann equation in the form (4), (5), we shall now evaluate the electron current density $\mathbf{j}$ from it. If $f^{(0)}$ is Maxwellian,

$$\nabla_P f^{(0)} = \frac{df^{(0)}}{d\Delta\epsilon} \nabla_P \Delta\epsilon = -\frac{\mathbf{v}}{kT} f^{(0)} \tag{6}$$

where $\mathbf{v}$ is the group velocity and $\Delta\epsilon = |\epsilon - \epsilon_b|$ is the distance from the band edge. The contribution of carriers in the $i$th valley to the current

density is then

$$\mathbf{j}^{(i)} = \sum_{\Delta \mathbf{P}^{(i)}, s} (\pm e) \mathbf{v}(\Delta \mathbf{P}^{(i)}) f(\Delta \mathbf{P}^{(i)})$$

$$= \frac{e^2}{kT} \sum_{\Delta \mathbf{P}^{(i)}, s} f^{(0)} \tau(\Delta \epsilon) \mathbf{E} \cdot \mathbf{v} \mathbf{v} \qquad (7)$$

where the summations are over all $\Delta \mathbf{P}^{(i)}$ occurring in the $i$th valley in unit volume of material, and over both states of spin.

The expression (7) states that any single valley $i$ possesses an anisotropic conductivity tensor $\sigma_{\alpha\beta}{}^{(i)}$, or an anisotropic mobility tensor $\mu_{\alpha\beta}{}^{(i)}$, i.e.,

$$j_\alpha{}^{(i)} = \sum_\beta \sigma_{\alpha\beta}{}^{(i)} E_\beta = \sum_\beta (n^{(i)} e \mu_{\alpha\beta}{}^{(i)}) E_\beta \qquad (8)$$

where

$$n^{(i)} = \sum_{\Delta \mathbf{P}^{(i)}, s} f^{(0)} \qquad (9)$$

is the number of carriers in valley $i$ per unit volume. If we choose the $x$, $y$, and $z$ axes to be along the principal axes of the ellipsoidal energy surfaces of valley $i$, (7) shows that $\sigma_{\alpha\beta}{}^{(i)}$ and $\mu_{\alpha\beta}{}^{(i)}$ will be diagonal. Each diagonal element $\mu_{\lambda\lambda}{}^{(i)}$ will involve a Maxwellian average of $v_\lambda{}^2 \tau(\Delta \epsilon)$. Now the equipartition principle leads us to expect that the average, over an energy shell $\Delta \epsilon$ to $\Delta \epsilon + d\Delta \epsilon$ in $\Delta \mathbf{P}$-space, of the kinetic energy associated with the $x$-component of velocity should be the same as that associated with the $y$- or $z$-component. This is easily demonstrated explicitly (Appendix B). Thus, if $m_1{}^*$, $m_2{}^*$, $m_3{}^*$ are the effective masses in the three principal directions,

$$\tfrac{1}{2} m_1{}^* v_1{}^2 = \tfrac{1}{2} m_2{}^* v_2{}^2 = \tfrac{1}{2} m_3{}^* v_3{}^2 \qquad (10)$$

Therefore (7) and (8) give, in our system of axes,

$$\mu_{\alpha\alpha}{}^{(i)} = \frac{e}{m_\alpha{}^*} \frac{\langle \Delta \epsilon \tau \rangle}{\langle \Delta \epsilon \rangle} \qquad (11)$$

$$\mu_{\alpha\beta}{}^{(i)} = 0 \qquad (\alpha \neq \beta) \qquad (12)$$

where the angular brackets denote Maxwellian averages:

$$\langle \Delta \epsilon \tau \rangle = \sum_{\Delta \mathbf{P}} \Delta \epsilon \tau f^{(0)} \Big/ \sum_{\Delta \mathbf{P}} f^{(0)}, \text{ etc.} \qquad (13)$$

The denominator of (11), $\langle \Delta \epsilon \rangle$, of course equals $\tfrac{3}{2} kT$ by equipartition.

The formula (11), it will be noted, is the same as that of the simple theory[7] with $m^*$ replaced by $m_\alpha^*$.

The overall conductivity tensor due to the carriers in all the valleys is of course $\sum_i \sigma_{\alpha\beta}^{(i)}$, and the overall mobility tensor is the average of $\mu_{\alpha\beta}^{(i)}$ over the different valleys. For a cubic crystal the mobility is the same in all directions, so we have

$$\mu = \mu_{\lambda\lambda} = \tfrac{1}{3} \sum_\alpha \mu_{\alpha\alpha} = \frac{1}{3N_V} \sum_i \sum_\alpha \mu_{\alpha\alpha}^{(i)}$$

$$= \frac{e}{m^{(I)}} \frac{\langle\Delta\epsilon\tau\rangle}{\langle\Delta\epsilon\rangle} \tag{14}$$

where $N_V$ is the number of valleys and $m^{(I)}$ is an average inertial mass defined by

$$\frac{1}{m^{(I)}} = \tfrac{1}{3}\left[\frac{1}{m_1^*} + \frac{1}{m_2^*} + \frac{1}{m_3^*}\right] \tag{15}$$

This mass, as we shall see in Section 7, is the one most directly measured by the Benedict-Shockley experiment on high-frequency dielectric constant.

## 3. TEMPERATURE VARIATION OF LATTICE MOBILITY

The $\tau$ occurring in the mobility expression (14) differs from the $\tau$ of the simple model in that it contains the effect of inter-valley scattering in addition to intra-valley and impurity scattering. Inter- and intra-valley scattering differ in that most of the phonons emitted or absorbed in intra-valley scattering have energies $\ll$ the energies of the charge carriers, while those involved in inter-valley scattering usually do not. If $\mathbf{K}^{(i)}$ and $\mathbf{K}^{(j)}$ are two different band edge points, scattering of a carrier from valley $i$ to valley $j$ must involve emission or absorption of a phonon of wave number close to $\pm\mathbf{q}_{ij}$, where $\mathbf{q}_{ij} = \mathbf{K}^{(i)} - \mathbf{K}^{(j)}$. If $\mathbf{q}_{ij}$ has a magnitude of the order of the radius of the Brillouin zone, as is likely in most cases, the energy $\hbar\omega_{ij}$ of this phonon will be a major fraction of of $k\Theta$, where $\Theta$ is the Debye temperature. This is usually $\geq kT$ in the extrinsic range. One must therefore use the Planck, rather than the Rayleigh-Jeans, distribution function for these phonons. At very low temperatures, inter-valley scattering is negligible: absorption of an $ij$ phonon is rare because few such phonons are present; emission is com-

---

[7] See, for example, W. Shockley, *Electrons and Holes in Semiconductors*, (Van Nostrand 1951) p. 276.

parably rare because few carriers have energy enough to create such a phonon. With rising temperature inter-valley scattering becomes more important. This causes $\tau$ (hence $\mu$) to decrease more rapidly with increasing $T$ than it would if there were no inter-valley scattering. In this section we shall develop this idea quantitatively.

The matrix element for scattering of a carrier from some state in valley $i$ to another state in valley $j$, by absorption or emission of a phonon $\hbar\omega$, has the form common to all one-phonon scattering processes[8]

$$M_{ij} = \left.\begin{array}{c} N^{1/2} \\ (N+1)^{1/2} \end{array}\right\} D_{ij} \text{ for } \left\{\begin{array}{l} \text{absorption} \\ \text{emission} \end{array}\right. \tag{16}$$

where $N$ is the number of phonons of the given type present in the initial state and where $D_{ij}$ is independent of the occupation of the phonon states. For inter-valley scattering $D_{ij}$ is practically independent of the locations of the initial and final states in their respective valleys. In general, of course, $D_{ij}$ will be different for the different branches of the vibrational spectrum. The transition probability from a state of energy $\epsilon$ in valley $i$ to a state in valley $j$ of energy $\epsilon' = \epsilon + \hbar\omega$ (absorption) or $\epsilon - \hbar\omega$ (emission) is proportional to $|M_{ij}|^2$ times the density of states at energy $\epsilon'$ in the $j$th valley, provided the variation of $\hbar\omega$ with position in the valley is negligible, as is the case for most transitions. Since the number of states between $\epsilon'$ and $\epsilon' + d\epsilon'$ is proportional to $\Delta\epsilon'^{1/2} d\Delta\epsilon'$, where $\Delta\epsilon'$ is the distance from the band edge, this transition probability has the form

$$\text{absorption: } W_a \propto \frac{(\epsilon + \hbar\omega)^{1/2}}{\exp(\hbar\omega/kT) - 1} \tag{17}$$

$$\text{emission: } W_e \propto \frac{(\Delta\epsilon - \hbar\omega)^{1/2}}{1 - \exp(-\hbar\omega/kT)} \text{ for } \Delta\epsilon > \hbar\omega$$

$$0 \qquad \text{for } \Delta\epsilon < \hbar\omega \tag{18}$$

Since either of the processes (17) and (18) randomizes the initial velocity of the charge carriers, and since in this paper we are assuming the existence of an effective relaxation time $\tau_{ii}(\epsilon)$ for randomization of velocity by the intra-valley scattering of acoustic modes, the total relaxation time for lattice scattering is given by

$$\frac{1}{\tau} = \frac{1}{\tau_{ii}} + \sum_{j,\alpha}' [W_a(ij,\alpha) + W_e(ij,\alpha)] \tag{19}$$

where $\alpha$ labels the branches of the vibrational spectrum and $W_a$ and $W_e$

_____
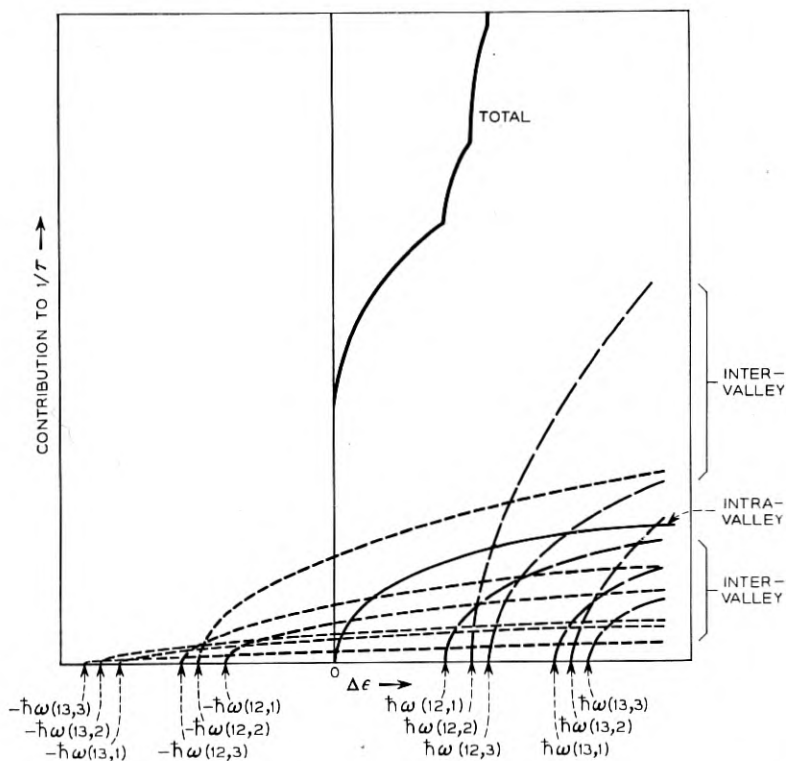[8] See, for example, Reference 7, p. 520.

Fig. 2 — Contributions to the reciprocal relaxation time of a charge carrier, due to inter-valley and intra-valley lattice scattering. The dot-dash curves are the inter-valley scattering contributions $W_e(ij, \alpha)$ for emission of a phonon, the dashed curves are the corresponding quantities $W_a(ij, \alpha)$ for absorption of a phonon. There are many transitions from valley 1 to different ones of the other valleys, due to different branches $\alpha$ of the phonon spectrum.

are given for each type of transition by (17) and (18) respectively, with $\hbar\omega = \hbar\omega(ij, \alpha)$. The prime on the summation means that when $\alpha$ is an acoustic branch, the term $j = i$ is to be omitted. However, since (17) and (18) apply to intra-valley scattering by modes of the optical branches, such scattering is included in (19) as the terms with $j = i$. Fig. 2(a) shows the various contributions to $1/\tau$ as functions of the initial energy $\Delta\epsilon$ of the carrier being scattered: $1/\tau_{ii}$ is proportional to $\Delta\epsilon^{1/2}$, as in the simple theory (this corresponds to a mean free path independent of energy for any given direction of motion), and each of the other terms is proportional to some $(\Delta\epsilon \pm \hbar\omega)^{1/2}$.

We shall try to estimate the order of magnitude of the steepness of

the parabolas describing the various inter-valley terms, relative to that of the intra-valley term $1/\tau_{ii}$. For the low-frequency acoustic modes involved in intra-valley scattering the factor $D_{ij}$ in the matrix element (16) is proportional to $q/(\hbar\omega)^{1/2}$, and since $\omega \propto q$ and $N \approx kT/\hbar\omega \gg 1$ for such modes $|M_{ij}|^2 \propto T$ and is independent of $q$. The $q$'s involved in inter-valley scattering will usually be too large to satisfy $kT/\hbar\omega \gg 1$, at least in the extrinsic ranges of Ge and Si, but we may hope to estimate a plausible order of magnitude for their $W_a$'s and $W_e$'s by assuming their $D_{ij}$'s to be $\propto q/\hbar\omega$ with a factor of proportionality of the same order as for intra-valley scattering. With this assumption the steepness of a typical $(ij)$ parabola corresponding to phonon emission ($W_e$) should be of the same order as the steepness of the intra-valley parabola when $kT \geqq \hbar\omega(ij)$, while for $kT < \hbar\omega(ij)$ the $W_e(ij)$ parabola should become nearly independent of $T$, as contrasted with $1/\tau_{ii}(\epsilon) \propto T$. The parabolas corresponding to phonon absorption are of course always less steep, the ratio of the steepness of $W_a(ij)$ to that of $W_e(ij)$ being

$$\frac{W_a/(\Delta\epsilon + \hbar\omega)^{1/2}}{W_e/(\Delta\epsilon - \hbar\omega)^{1/2}} = \frac{1 - \exp(-\hbar\omega/kT)}{\exp(\hbar\omega/kT) - 1} = \exp(-\hbar\omega/kT) \qquad (20)$$
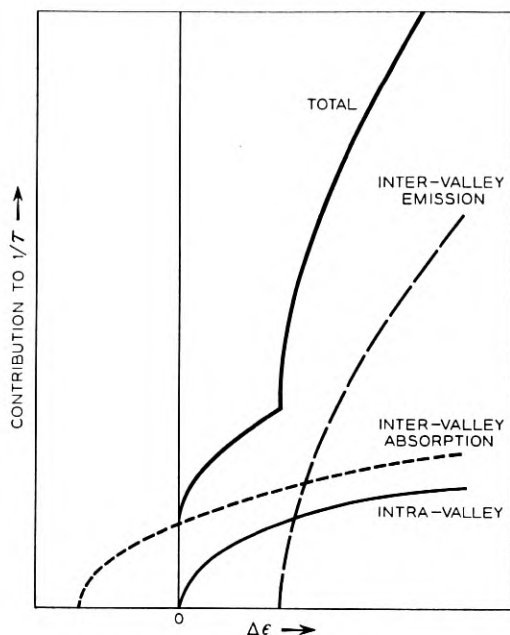


Fig. 3 — Same as Fig. 2, but for the simplified model of Equation (24), on which the numerical calculations of Figs. 4, 5, 8, 9 and 10 are based.

Because of the large number of terms $W_{a,e}(ij, \alpha)$ in (19) or Fig. 2, each with an at present unknown amplitude and critical frequency, it would be pointless to undertake calculations taking individual account of all the possible types of transitions. However, it is reasonable to hope that the behavior of the inter-valley terms can be roughly approximated by a model which considers absorption and emission of just a single type of inter-valley phonon. This model is illustrated in Fig. 3. It contains three adjustable parameters $w_1$, $w_2$, and $\hbar\omega$, defined by

$$1/\tau_{ii} = w_1 \left(\frac{\Delta\epsilon}{\hbar\omega}\right)^{1/2} \left(\frac{kT}{\hbar\omega}\right) \tag{21}$$

$$W_a = \frac{w_2 \left(\dfrac{\Delta\epsilon}{\hbar\omega} + 1\right)^{1/2}}{\exp\left(\hbar\omega/kT\right) - 1} \tag{22}$$

$$W_e = \frac{w_2 \left(\dfrac{\Delta\epsilon}{\hbar\omega} - 1\right)^{1/2}}{1 - \exp\left(-\hbar\omega/kT\right)} \text{ or } 0 \tag{23}$$

Equation (19) becomes

$$w_1\tau =$$

$$\left\{\left(\frac{\Delta\epsilon}{\hbar\omega}\right)^{1/2}\left(\frac{kT}{\hbar\omega}\right) + \frac{w_2}{w_1}\left[\frac{\left(\dfrac{\Delta\epsilon}{\hbar\omega} + 1\right)^{1/2}}{\exp\left(\hbar\omega/kT\right) - 1} + \frac{\left(\dfrac{\Delta\epsilon}{\hbar\omega} - 1\right)^{1/2} \text{ or } 0}{1 - \exp\left(-\hbar\omega/kT\right)}\right]\right\}^{-1} \tag{24}$$

Thus $w_1\tau$ is a function of the two variables $\Delta\epsilon/\hbar\omega$ and $kT/\hbar\omega$, and the single parameter $w_2/w_1$. The behavior of the mobility as a function if $kT/\hbar\omega$ therefore depends, apart from the constant scale factor $w_1$, only on $w_2/w_1$.

Fig. 4 shows the results of some calculations of this mobility-temperature relation, made by numerical evaluation of (24) and (14). It is evident that with reasonable values of $w_2/w_1$, the negative exponent describing the temperature variation of the mobility can be increased to a value considerably above the $\frac{3}{2}$ of the simple theory, over a considerable range of temperature. This is often what is needed to explain the observed mobility behavior. In Sections 4, 6, and 8 we shall see the extent to which this mobility exponent is correlated with, respectively, the electronic part of the thermoelectric power, the ratio of Hall to drift mobility, and the magnitude of the magnetoresistance.
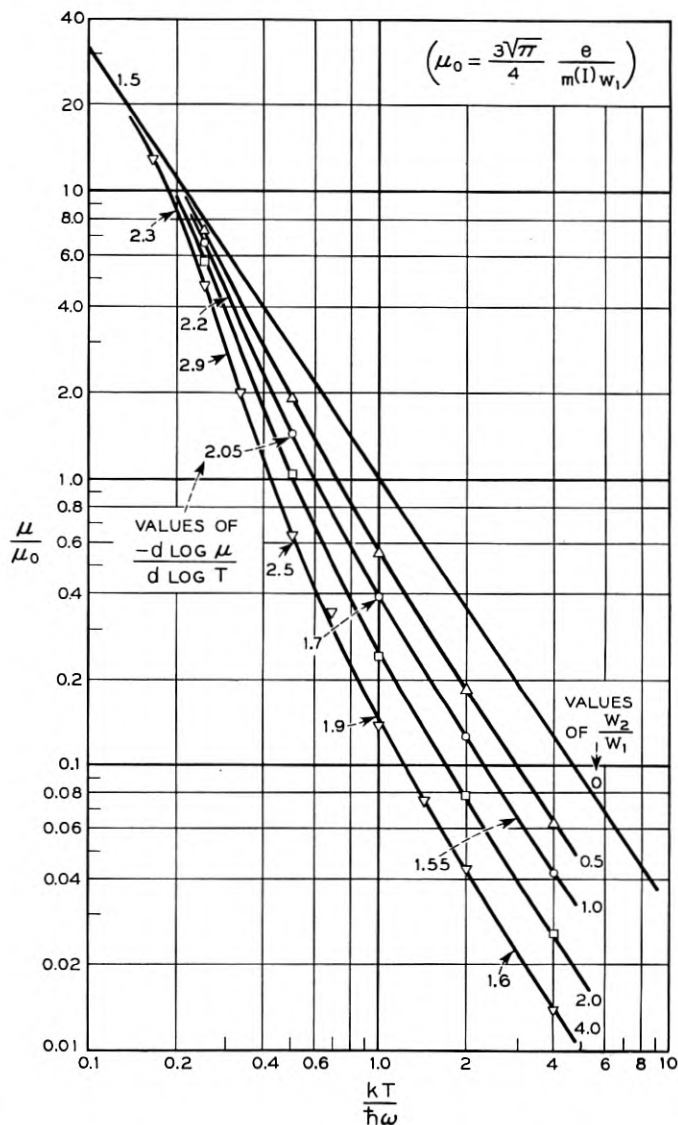
Fig. 4 — Mobility-temperature curves for pure lattice scattering, as obtained from the simplified expression (24) for the relaxation time. The quantities $w_1$ and $w_2$ measure the strength of the coupling of the carriers to intra- and inter-valley modes, respectively; $\omega$ is the frequency of the inter-valley mode. The curves have been drawn so as to smooth out irregularities in the computed points.

## 4. THERMOELECTRIC POWER

The thermoelectric power of a semiconductor is a little different from the other effects discussed in this paper, in that it involves not only the response of the distribution function of the charge carriers to a perturbing temperature gradient or electric field, but also the alteration of the distribution function of the phonon system.[9, 10] The duality manifests itself in the apperaance of two contributions to the thermoelectric power $Q$: the measured $Q$ is the sum of an electronic part $Q_e$, representing the emf necessary to counteract the tendency of charge carriers to diffuse from hot regions to cold, and a phonon part $Q_p$, representing the emf necessary to counteract the drag exerted on the carriers by the phonons which flow down the temperature gradient in thermal conduction. As the present paper is devoted to effects having to do with the response of the electronic distribution function to various influences, and as all aspects of the theory of thermoelectric power have been discussed elsewhere,[10] we shall limit the present section to a discussion only of the electronic part $Q_e$, which, fortunately, predominates greatly over $Q_p$ at high temperatures.

The expression for $Q_e$ is most simply derived by making use of the Kelvin relation $Q_e = \Pi_e/T$ between $Q_e$ and the electronic contribution $\Pi_e$ to the Peltier coefficient, which represents the energy flux, relative to the Fermi level, which accompanies the transport of unit charge in an isothermal conduction process. For an intrinsic semiconductor with low carrier concentration

$$eTQ_e = e\Pi_e = \epsilon_F - \epsilon_b - \Delta\epsilon_T \qquad (25)$$

where as before $\epsilon_F$ is the Fermi level, $\epsilon_b$ the band edge energy, and where $\Delta\epsilon_T$ is the average energy of the transported electrons relative to the band edge, a quantity $>0$ for n-type material, $<0$ for p-type, and of the order of magnitude of $kT$. Now $|\epsilon_F - \epsilon_b|$ can be expressed in terms of the carrier concentration $n$ and the effective masses. For a many-valley model the number of carriers $n^{(i)}$ in each valley is easily shown to be the same as for a simple model semiconductor with the same $|\epsilon_F - \epsilon_b|$ and with an effective mass equal to the geometric mean of the principal masses $m_1^*$, $m_2^*$, $m_3^*$ of the valley. This is because the density of states in energy is proportional to the volume of K-space inside an energy surface, a quantity which for a spherical surface goes as the cube of the radius, and for an ellipsoidal one as the product of the principal semi-

[9] H. P. R. Frederikse, Phys. Rev., **92**, p. 248, 1953.
[10] C. Herring, Phys. Rev., **96**, p. 1163, 1954.

axes. The total carrier concentration is therefore

$$n = N_V n^{(i)} = \frac{2(2\pi kT)^{3/2}}{h^3} (m_1^* m_2^* m_3^*)^{1/2} N_V \exp\left(-\frac{|\epsilon_F - \epsilon_b|}{kT}\right)$$

where $N_V$ is the number of valleys. The final expression for $Q_e$ obtained by expressing $|\epsilon_F - \epsilon_b|$ in terms of $n$ and inserting in (25) is

$$Q_e = \mp$$

$$86.2 \left[ \ell n \frac{4.70 \times 10^{15}}{n} + \tfrac{3}{2} \ell n\, N_V + \tfrac{1}{2} \ell n \left( \frac{m_1^*}{m} \cdot \frac{m_2^*}{m} \cdot \frac{m_3^*}{m} \right) \right.$$
$$\left. + \tfrac{3}{2} \ell n\, T + \frac{|\Delta\epsilon_T|}{kT} \right] \mu v/\text{deg.} \tag{26}$$

where $n$ is in cm$^{-3}$ and the upper sign is for n-type material, the lower for p-type.

If a relaxation time exists, dependent only on energy, the distribution function for isothermal conduction has the form $f^{(0)} + \mathbf{E} \cdot \mathbf{f}^{(1)}$ worked out in Section 2, and we have, for a cubic crystal,

$$|\Delta\epsilon_T| = \frac{\int \Delta\epsilon \mathbf{v} \cdot \mathbf{f}^{(1)}\, d\mathbf{P}}{\int \mathbf{v} \cdot \mathbf{f}^{(1)}\, d\mathbf{P}} = \frac{\langle \Delta\epsilon^2 \tau \rangle}{\langle \Delta\epsilon\tau \rangle} \tag{27}$$

by (5), (6) and (10), where as before the angular brackets denote Maxwellian averages as defined by (13).

It is important to know the value of (27) as accurately as possible, in the temperature range where $Q_e$ is measurable, since if (27) is known the measured $Q_e$ can be used with (26) to give information on the effective masses. For pure intra-valley scattering, (27) has the value $2kT$. Impurity scattering increases $\Delta\epsilon_T$ by causing the current to be carried more by fast carriers and less by slow; inter-valley scattering has the reverse effect. It is worth while to try to correlate the effect of inter-valley scattering on $\Delta\epsilon_T$ with its effects on two measurable properties, namely, the temperature variation of mobility (Section 3) and the ratio of Hall to drift mobility (Section 6). Accordingly, calculations of $|\Delta\epsilon_T|$ have been made using the expression (24) (model of Fig. 3) for the relaxation time. The results are shown in Fig. 5, which shows $|\Delta\epsilon_T|/kT$ as a function of $kT/\hbar\omega$.

5. PIEZORESISTANCE

As we have just seen in Section 2, the quantum states in any small region of wave number space make a contribution to the conductivity
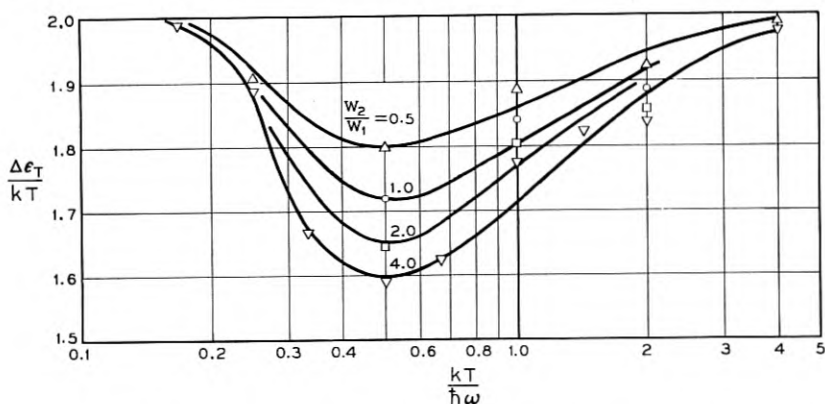
Fig. 5 — Values of the thermoelectric transport ratio $|\Delta\epsilon_T|/kT$ defined by (27), for the simplified lattice scattering law (24). The ratio $w_2/w_1$ measures the strength of the coupling of the carriers to inter-valley modes in terms of that for intra-valley scattering; $\omega$ is the frequency of the inter-valley modes. The curves have been drawn to smooth out irregularities (severe for $kT/\hbar\omega = 1$ to 2) in the calculated points.

which depends on (i) *the degree to which these states are populated* in the equilibrium distribution $f^{(0)}$, (ii) *their group velocity*, and (iii) *their relaxation time*, or more generally, the transition probabilities for scattering from these states to others. When the crystal is strained, any or all of these factors may change, and the resulting change in the sum of all the local contributions to the conductivity constitutes the piezoresistance effect recently discovered by Smith.[4] Although there are a number of processes which can contribute to the three factors (i) to (iii) just enumerated, it can be argued plausibly that for a simple many-valley model the principal effects are usually those due to a single process, namely, the strain-induced shifts of the energies $\epsilon_i$ of the band edge points $\mathbf{K}^{(i)}$. We shall consider (i) to (iii) in turn:

(i) The change in the population $f^{(0)}(K^{(i)} + \Delta K)$ depends on the shift of the energy $\epsilon(K^{(i)} + \Delta K)$, and because of the smallness of $\Delta K$ this is practically the same as the shift $\delta\epsilon^{(i)}$ of $\epsilon(\mathbf{K}^{(i)})$. In a shearing strain some of the $\delta\epsilon^{(i)}$ will be positive, some negative, and so some of the valleys will have their populations decreased, some increased, the fractional change in each case being $\delta\epsilon^{(i)}/kT$. Now it is evident from the second equation of (7) that the contribution of a single valley to the conductivity is anisotropic. If all valleys are populated equally, as we assumed in Section 1, the total conductivity will be isotropic. But if strain causes different valleys to have different populations, the overall conductivity will have

an anisotropy like that of the more populous valleys. If the ratio of $\delta\epsilon^{(i)}$ to the shear strain amplitude is of the order of magnitude of the known ratio of $\delta\epsilon_G$ to strain for isotropic compression, viz., a few volts, the fractional change of $f^{(0)}$ per unit strain will be of the order of hundreds. Since the observed fractional change in resistance per unit shear strain is of the order of $10^2$ in the more favorable orientations,[4] the change in $f^{(0)}$ is of the right order of magnitude to contribute a major part of the effect.

For an isotropic compression or dilation there exists the possibility, not present for shearing strains, that the total carrier concentration may be changed in first order. A large effect of this sort, again of the order of the $\delta\epsilon^{(i)}/kT$, will occur for a specimen in or near the intrinsic range, because of the change of energy gap $\epsilon_G$ with strain. This effect rapidly becomes negligible, however, as the specimen is made extrinsic. For example, if unit volume of an n-type specimen has an excess $n_D$ of donors over acceptors, all ionized, the hole and electron concentrations $n_h$, $n_e$, obey

$$n_e = n_D + n_h, \qquad n_e n_h = n_i^2$$

where $n_i(T)$ is the value of $n_e = n_h$ in intrinsic material. Thus if $n_D \gg n_i$

$$\frac{n_h}{n_D} = \left(\frac{n_i}{n_D}\right)^2 \tag{28}$$

and since $\partial n_e/\partial\epsilon_G = \partial n_h/\partial\epsilon_G$, the energy gap effect is negligible if $n_D$ exceeds $n_i$ by a large factor, even though the change in $n_i$ with strain may be sizable. For extrinsic specimens with incomplete ionization of impurity centers, there may of course be an effect of compression on total carrier concentration due to change in the ionization energy of the centers; however, if this ionization energy is $\ll \epsilon_G$ this effect will be of a smaller order of magnitude than the $\delta\epsilon^{(i)}/kT$.

(ii) It is easy to show that the fractional change in group velocity per unit strain must be much smaller than the $\delta\epsilon^{(i)}/kT$ just discussed, hence too small to contribute in a major way to the piezoresistance effect. For we expect the change $\delta v$ in the group velocity at $\mathbf{K}^{(i)} + \Delta\mathbf{K}$ to have an order of magnitude given by

$$\delta v \sim [\delta\epsilon(\mathbf{K}^{(i)} + \Delta\mathbf{K}) - \delta\epsilon^{(i)}]/\hbar\Delta K \sim (\Delta K/K^{(i)})^2(\delta\epsilon^{(i)}/\hbar\Delta K)$$

since the quantity in square brackets must vary as $\Delta K^2$. Since

$$v \sim [\epsilon(\mathbf{K}^{(i)} + \Delta\mathbf{K}) - \epsilon^{(i)}]/\hbar\Delta K = \Delta\epsilon/\hbar\Delta K$$

we have

$$\delta v/v \sim (\Delta K/K^{(i)})^2(\delta\epsilon^{(i)}/\Delta\epsilon) \tag{29}$$

Since a typical charge carrier has $\Delta\epsilon \sim kT$, (29) is smaller than the ratio discussed in the preceding paragraph by the factor $(\Delta K/K^{(i)})^2$. It is thus plausible to neglect strain-induced changes in group velocity, or equivalently, in the effective masses.

(iii) Consider the transition probability from a state $\mathbf{K}$ to the group of states lying in a small element of volume in $\mathbf{K}$-space, centered on a point $\mathbf{K}'$ at which the proper energy conservation law for the transition $\mathbf{K} \rightarrow \mathbf{K}'$ is satisfied. This probability, like all quantum-mechanical transition probabilities, can be expressed as the product of the square of a matrix element $M(\mathbf{K}, \mathbf{K}')$ by the number of states per unit energy in the given element of volume. We have to consider the effect of strain on each of these factors.

The matrix element $M(\mathbf{K}, \mathbf{K}')$ can be changed either by a change in the wave functions $\Psi_K$, $\Psi_{K'}$, or by a change in the physical processes determining the perturbation operator $M$, e.g., a change in the amplitudes of the thermal vibrations, or a change in the dielectric constant, which enters into scattering by charged impurities. Typical assumptions on the equation of state of a crystal suggest that the fractional change in the squared vibration amplitude, per unit strain, might be of the order of a few units, i.e., at least an order of magnitude less than the observed elastoresistance for the optimum orientations. The effect of the change in the wave functions is of similar magnitude: To effect a major change in $M(\mathbf{K}, \mathbf{K}')$ one must make a major change in the wave functions. To do this probably usually requires a strain of amplitude 0.1 to 1. Therefore it is reasonable to expect that the fractional change in $|M^2|$ per unit strain will be of the order of 10 or less, i.e., again an order of magnitude smaller than $\delta\epsilon^{(i)}/kT$, or than the observed elastoresistance.

The effect of strain on the density-of-states factor, on the other hand, can be larger. For intra-valley scattering, where initial and final states are both near the same band edge point $\mathbf{K}^{(i)}$, the effect is of course very small, since initial and final states undergo very nearly the same energy shift with strain. But for scattering from one valley $i$ to another valley $j$, the two energy shifts $\delta\epsilon^{(i)}$ and $\delta\epsilon^{(j)}$ are in general quite different, and for a given initial state application of a strain will change the set of $\mathbf{K}'$'s describing final states which conserve energy and hence will change the density of final states — e.g., the density in a given solid angle of vectors $\Delta\mathbf{K}' \equiv \mathbf{K}' - \mathbf{K}^{(j)}$. Since in a given solid angle the density of states is $\propto \Delta\epsilon'^{1/2}$, the fractional change in this density due to a strain is $\delta\epsilon^{(j)}/2\Delta\epsilon'$, which on the average is of the order of $\delta\epsilon^{(j)}/kT$, i.e., of the same order as the effect discussed under (i).

TABLE I — WAYS IN WHICH STRAIN CAN AFFECT CONDUCTIVITY

| Effect | Probable Order of Magnitude | Rank in Importance |
|---|---|---|
| (i) Population function $f^{(0)}$ . . . . . . . . . . . . | $\delta\epsilon^{(i)}/kT$ + much smaller terms | First |
| (ii) Group velocities of states . . . . . . . . . | $\ll \delta\epsilon^{(i)}/kT$ | |
| (iii) Transition probabilities | | |
| (a) Matrix elements | | |
| ($\alpha$) Wave functions . . . . . . . . . . . | $\ll \delta\epsilon^{(i)}/kT$ | |
| ($\beta$) Vibration amplitudes, etc... | Rather $< \delta\epsilon^{(i)}/kT$ | Second (?) |
| (b) Density of states | | |
| ($\alpha$) Intravalley . . . . . . . . . . . . . . . | $\ll \delta\epsilon^{(i)}/kT$ | |
| ($\beta$) Intervalley . . . . . . . . . . . . . . . | $\delta\epsilon^{(i)}/kT$ + much smaller terms | First |

Table I summarizes the foregoing discussion of the ways in which strain can affect conductivity.

Appendix C gives the mathematical treatment of the two effects which are of the order of the quantities $\delta\epsilon^{(i)}/kT$, namely, the change in $f^{(0)}$ and the change in the density-of-states factor in the transition probabilities for inter-valley scattering. This treatment, which is fairly simple and straightforward, is based on the following assumptions:

(a) Neglect of all other effects of strain on the conductivity.

(b) The assumption of the preceding sections that the scattering of the carriers is describable by a relaxation time which in each valley is a function of energy only.

(c) Carrier concentrations in the extrinsic range.

(d) Maxwell-Boltzmann statistics.

(e) Valleys lying along a threefold or fourfold symmetry axis of a cubic crystal. For such valleys the energy surfaces are ellipsoids of revolution.

The principal features of the calculation are qualtative ones which can be derived with little or no mathematics. These we shall consider here, with a little inquiry in each case as to the sensitivity of the conclusion to relaxation of the assumptions (a) and (e) above. The first such feature to be noted is that *under assumption (a) the change of mobility in an isotropic compression vanishes.* For in an isotropic compression all the band edge shifts $\delta\epsilon^{(i)}$ are equal. This means that for a given total carrier density the distribution function $f^{(0)}$ in each valley does not change, if, as we are doing, we neglect changes in the effective masses. Similarly, since all valleys are shifted together, there is no change in the density of final states corresponding to any inter-valley scattering process. The present conclusion is easily seen to be independent of

TABLE II — ISOTHERMAL ELASTORESISTANCE CONSTANTS FOR Ge AND Si (SMITH, REFERENCE 4)

| Material and Resistivity | | $m_{1212} \equiv m_{44}$ | $\dfrac{m_{1111} - m_{1122}}{2} \equiv \dfrac{m_{11} - m_{12}}{2}$ | $-\dfrac{1}{\sigma}\dfrac{d\sigma}{d \ln V} = \dfrac{m_{11} + 2m_{12}}{3}$ |
|---|---|---|---|---|
| | Ω cm | | | |
| n Ge | 1.5 | −93.0 | +0.4 | −5.3 |
| | 5.7 | −92.0 | +0.5 | −6.8 |
| | 9.9 | −92.8 | +0.1 | −9.8 |
| | 16.6 | −93.4 | +0.1 | −13.6 |
| p Ge | 1.1 | +65.1 | −2.8 | +3.9 |
| | 15.0 | +66.5 | −6.3 | +1.4 |
| n Si | 11.7 | −10.8 | −79.5 | +5.7 |
| p Si | 7.8 | +110.0 | +3.9 | +6.0 |

Here $m_{\mu\nu\alpha\beta}$, defined by (C7) of Appendix C, describes the relative change of the conductivity tensor with the strain tensor, in a coordinate system oriented along the cube axes. The abbreviation of this by $m_{rs}(r, s = 1 \text{ to } 6)$ follows the same practice as that used for elastic constants.

assumptions (b), (d), and (e). As regards assumption (a), however, it is clear that inclusion of any of the other strain effects listed in Table I will in general lead to a nonvanishing effect of compression on the mobility.

By virtue of the fact just mentioned it is possible to test the validity of assumption (a) by comparing the observed elastoresistance for isotropic compression with that for a typical shear. Table II, taken from the work of Smith,[4] shows the room temperature elastoresistance constants of Ge and Si. The entries in the last column vary with resistivity for the case of Ge, because of the energy gap effect discussed under (i) above [failure of assumption (c)]; our present interest is therefore in the values for low resistivity specimens. For these the volume coefficient (last column) is in all cases only a few percent of the larger of the shear coefficients (middle columns); this accords with the expectation that the volume variation of the squared matrix element for scattering (presumably the largest of the neglected effects) should be an order of magnitude or more smaller than the $\delta\epsilon^{(i)}/kT$. This is encouraging, but it must be remembered that the shear variation of the matrix element may well be larger than its volume variation because suitable shearing strains can usually couple a band edge state to states closer to it in energy than can isotropic dilatation.

The second important conclusion is that *under assumption (a) the change of mobility vanishes for a dilatation along a (100) direction if the valleys are on (111) axes, and for a dilatation along a (111) direction, if the valleys are on (100) axes.* In terms of the elastoresistance coefficients of Table II, $(m_{11} - m_{12})$ vanishes for (111) valleys, and $m_{44}$ vanishes for (100) valleys. This conclusion is obvious from the symmetry of the

problem: a shear compounded out of a unidirectional dilatation of the type described and an isotropic compression must shift all band edge points by the same amount. This amount must be the same as in the negative of this shear, so all $\delta\epsilon_i = 0$. This conclusion is again independent of assumptions (b) and (d), but in general breaks down if assumption (a) is relaxed to the extent of taking account of the effect of strain on the matrix element for scattering.

The third point to be made is that *the change of mobility accompanying a given strain is inversely proportional to $T$ at temperatures low enough for inter-valley scattering to be of negligible importance.* This is because the relative change of population of different valleys with strain is proportional to the $\delta\epsilon^{(i)}/kT$. The more complete treatment of Appendix C shows that, under the present assumptions, the decrease of elastoresistance with increasing $T$ should be more rapid than $1/T$ when inter-valley scattering is just becoming important, but that for very high $T$ it should again go as $1/T$. This behavior is illustrated schematically in Fig. 6. The present conclusion is not dependent on assumptions (b) or (e), but depends on the others, especially (a). The effect of strain on the matrix elements for scattering will give a contribution to the elastoresistance which is independent of $T$ in the range (if such exists) where only intra-valley lattice scattering is important; if impurity or inter-valley scattering contributes the dependence is of course more complicated.
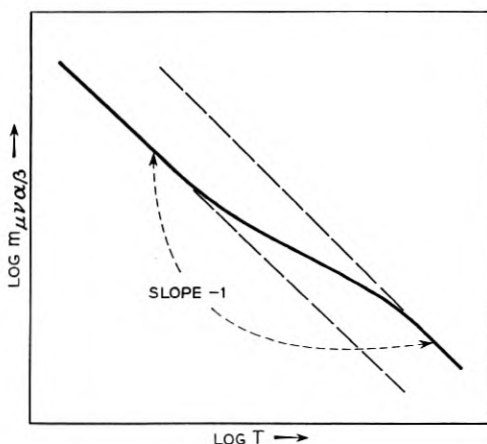


Fig. 6 — Schematic variation of any component of elastoresistance with temperature, showing the transition from the low temperature region where the only important effect of strain is to change the relative population of the valleys, to the high temperature region where the effect of strain on inter-valley scattering is of comparable importance.

In conclusion, a few words are in order regarding the extent to which the present conclusions on piezoresistance can be expected to hold for models other than the simple many-valley type to which this memorandum is restricted. First of all, we may note that for the various types of degenerate models, such as that of Fig. 1(a), different states in the neighborhood of the same band edge point can experience widely different energy shifts $\delta\epsilon$ under shear. Consequently the mobility will be affected in a major way not only by the two effects labeled "first" in Table I — change of the population function and change of the density-of-states function in inter-valley scattering — but also by changes in the group velocities. Moreover, the perturbation of the crystal Hamiltonian by a uniform shearing strain may now have sizable matrix elements between states of the same wave vector belonging to the different bands which come together at the band edge point. This can cause the form of the crystal wave functions to be much more sensitive to strain than when the perturbation only connects states a few volts apart in energy, and so the dependence on strain of the matrix element for scattering may be much larger than for the simple many-valley case. Thus at least four, rather than two, of the entries in Table I become of first magnitude.

In view of these facts, most of the conclusions reached for the simple many-valley case probably become invalid for the degenerate and degenerate many-valley cases. An exception is the conclusion concerning the smallness of the change of mobility in isotropic compression. The perturbation introduced into the crystal Hamiltonian by an isotropic compression does not mix states of a degenerate set, so the arguments previously given remain valid.

## 6. HALL EFFECT AT LOW $H$

When a magnetic field is present the $\mathbf{E}$ in the transport equation (3) must be replaced by $\mathbf{E} + \mathbf{v} \times \mathbf{H}/c$, where as before, $\mathbf{v}$ is the group velocity. Thus with the upper sign for electrons, the lower for holes, the distribution function $f$ of the carriers obeys

$$0 = \frac{\partial f}{\partial t} = \pm \left[ e\mathbf{E} \cdot \nabla_P f + e \frac{\mathbf{v} \times \mathbf{H}}{c} \cdot \nabla_P f \right] - \frac{(f - f^{(0)})}{\tau} \qquad (30)$$

We shall seek the solution of this as far as the first order in $\mathbf{E}$ and the first order in $\mathbf{H}$, i.e., we shall set

$$f = f^{(0)} + \mathbf{E} \cdot \mathbf{f}^{(10)} + \sum_{\mu\nu} E_\mu H_\nu f_{\mu\nu}^{(11)} + \ldots \qquad (31)$$

There is, of course, no term of the first order in $\mathbf{H}$ and the zeroth in $\mathbf{E}$,

since a pure magnetic field has no effect on $f$. The vector function $\mathbf{f}^{(10)}$ is of course the same as in Section 2, Equation (5), namely,

$$\mathbf{f}^{(10)} = \pm e\tau\nabla_P f^{(0)} \tag{32}$$

Putting (32) and (31) into (30) we get for $f_{\mu\nu}^{(11)}$:

$$\frac{1}{\tau}\sum_{\mu,\nu} E_\mu H_\nu f_{\mu\nu}^{(11)} = \pm e\,\frac{\mathbf{v}\times\mathbf{H}}{c}\cdot\nabla_P(\pm\,e\tau\mathbf{E}\cdot\nabla_P f^{(0)})$$

whence

$$f_{\mu\nu}^{(11)} = \frac{e^2\tau}{c}\sum_{\alpha,\beta}\delta_{\nu\alpha\beta}v_\beta\frac{\partial}{\partial P_\alpha}\left(\tau\frac{\partial f^{(0)}}{\partial P_\mu}\right) \tag{33}$$

where $\delta_{\nu\alpha\beta} = 0$ if any two of its suffixes are the same, and $\pm 1$ if the suffixes are an even (odd) permutation of $xyz$.

The physical meaning of the steps leading to (33) is just that a weak magnetic field perturbs the $f^{(1)}$ solution of Section 2 by displacing each part of the distribution in the direction of $\mathbf{v}\times\mathbf{H}$ in crystal momentum space, the displacement being proportional to $\mathbf{v}\times\mathbf{H}$ and to $\tau$.

The term (33) in the distribution function gives rise to a contribution $\mathbf{j}^{(11)}$ to the current, which is at right angles to $\mathbf{H}$ and to $\mathbf{E}$. This contribution can be described by a "Hall conductivity tensor" $\sigma_{\lambda\mu\nu}$, thus:

$$j_\lambda^{(11)} = \sum_{\mu\nu}\sigma_{\lambda\mu\nu}E_\mu H_\nu \tag{34}$$

The contribution $\sigma_{\lambda\mu\nu}^{(i)}$ of the $i$th valley to $\sigma_{\lambda\mu\nu}$ is easily obtained from (33). We shall assume Maxwellian statistics, so that $\partial f^{(0)}/\partial P_\mu = -(v_\mu/kT)f^{(0)}$. When this is inserted into (33) the last factor involves a derivative of $v_\mu f^{(0)}\tau$ with respect to $P_\alpha$. If $\tau$ depends only on energy, as we are assuming throughout this memorandum, the derivative of $f^{(0)}\tau$ with respect to $P_\alpha$ in (33) will be proportional to $v_\alpha$, and $\Sigma_{\alpha,\beta}\delta_{\nu\alpha\beta}v_\beta v_\alpha$ will vanish identically because of the anti-symmetry of $\delta_{\nu\alpha\beta}$ in $\alpha$ and $\beta$. Therefore the only term which need be retained in $\partial/\partial P_\alpha$ is that in $\partial v_\mu/\partial P_\alpha$. If the coordinate axes are chosen along the principal axes of the energy surfaces of the $i$th valley, this latter derivative is just $\delta_{\mu\alpha}/m_\mu^*$. Thus we get, with the upper sign for n-type the lower for p,

$$
\begin{aligned}
\sigma_{\lambda\mu\nu}^{(i)} &= \mp\frac{e^3}{c}\sum_{\Delta\mathbf{P}^{(i)},s}\tau\sum_{\alpha,\beta}\delta_{\nu\alpha\beta}v_\lambda v_\beta\frac{\partial}{\partial P_\alpha}\left(-\frac{\tau f^{(0)}v_\mu}{kT}\right)\\
&= \pm\frac{e^3}{c}\sum_{\Delta\mathbf{P}^{(i)},s}\frac{\tau^2 f^{(0)}}{m_\mu^*}\sum_\beta\delta_{\nu\mu\beta}\frac{v_\lambda v_\beta}{kT}
\end{aligned}
\tag{35}
$$

where as usual the first summation is over all vectors $\Delta\mathbf{P}^{(i)}$ in the $i$th

valley, per unit volume, and over both states of spin. In our present co-ordinate system the average of $v_\lambda v_\beta$ over an energy shell vanishes unless $\beta = \lambda$, while that of $v_\lambda^2$ can be evaluated from the equipartition relation (10): $v_\lambda^2 \to 2\Delta\epsilon/3m_\lambda^*$, where $\Delta\epsilon = |\epsilon - \epsilon_b|$ is the distance from the band edge. Thus with $kT = \tfrac{2}{3} \langle\Delta\epsilon\rangle$, (35) reduces to

$$\sigma_{\lambda\mu\nu}{}^{(i)} = \mp \frac{e^3 n^{(i)}}{c} \frac{\langle\Delta\epsilon\tau^2\rangle}{\langle\Delta\epsilon\rangle} \frac{\delta_{\lambda\mu\nu}}{m_\lambda^* m_\mu^*} \tag{36}$$

where $n^{(i)}$, as in (9), is the number of carriers in the $i$th valley per unit volume, and where the angular brackets are Maxwellian averages, as in (13).

The proportionality of the Hall conductivity tensor to $\langle\Delta\epsilon\tau^2\rangle$ and to the reciprocal product of two different principal masses is easy to understand physically. Without a magnetic field, an electric field in the $\mu$ direction gives a distribution, in each energy shell of the $i$th valley, which has a mean velocity in the $\mu$ direction proportional to $\Delta\epsilon\tau/m_\mu^*$ (cf. Section 2). Thus the distribution in this energy shell is acted on by a transverse magnetic force whose average value is proportional to this expression. This transverse magnetic force produces a transverse current proportional to the force and to $\tau/m_\lambda^*$, where $\lambda$ is the transverse direction.

For a cubic crystal the relation of the Hall current to $\mathbf{E}$ and $\mathbf{H}$ must be isotropic, i.e., the right of (34) must be proportional to $\mathbf{E} \times \mathbf{H}$. It is easily shown that the quantities in (34) are related to the ordinary conductivity $\sigma_0$, Hall coefficient $R$, and Hall mobility $\mu_H = R\sigma_0 c$, by

$$\mathbf{j}^{(11)} = \sigma_0^2 R \, \mathbf{E} \times \mathbf{H}$$

or

$$\sigma_{\lambda\mu\nu} = \sigma_0^2 R \delta_{\lambda\mu\nu} = \mp \frac{\sigma_0 \mu_H}{c} \delta_{\lambda\mu\nu} \tag{37}$$

where as usual the upper sign is for n-type, the lower for p. Since $\Sigma_{\lambda\mu\nu}\sigma_{\lambda\mu\nu}\delta_{\lambda\mu\nu}$ is invariant with respect to changes in the orientation of the coordinate system, we may evaluate it by evaluating each $\Sigma_{\lambda\mu\nu}\sigma_{\lambda\mu\nu}{}^{(i)}\delta_{\lambda\mu\nu}$ in the system of principal axes of the $i$th valley, and then summing on $i$. From (36) we find in this way

$$\sigma_0^2 R = \mp \frac{\sigma_0 \mu_H}{c} = \tfrac{1}{6} \sum_{\lambda\mu\nu} \sigma_{\lambda\mu\nu}\delta_{\lambda\mu\nu} = \tfrac{1}{6} \sum_i \sum_{\lambda\mu\nu} \sigma_{\lambda\mu\nu}{}^{(i)}\delta_{\lambda\mu\nu}$$

$$= \mp \frac{e^3 n}{c} \frac{\langle\Delta\epsilon\tau^2\rangle}{\langle\Delta\epsilon\rangle} \cdot \tfrac{1}{3} \left( \frac{1}{m_1^* m_2^*} + \frac{1}{m_2^* m_3^*} + \frac{1}{m_3^* m_1^*} \right) \tag{38}$$

where $n = \Sigma_i n^{(i)}$ is the total density of carriers. A neater way of pre-

TABLE III — VALUES OF THE LAST FACTOR IN (39), FOR CASES OF THE FORM

$$m_1{}^* = m_2{}^* = m_\perp{}^*, \ m_3{}^* = m_\parallel{}^*.$$

| $\dfrac{m_\parallel{}^*}{m_\perp{}^*}$ | $B = \dfrac{3\,\dfrac{m_\parallel{}^*}{m_\perp{}^*}\left(2 + \dfrac{m_\parallel{}^*}{m_\perp{}^*}\right)}{\left(1 + 2\,\dfrac{m_\parallel{}^*}{m_\perp{}^*}\right)^2}$ |
|:---:|:---:|
| 20 | 0.784 |
| 10 | 0.816 |
| 5 | 0.868 |
| 3 | 0.918 |
| 2 | 0.960 |
| 1 | 1.000 |
| 0.5 | 0.938 |
| 0.3 | 0.808 |
| 0.2 | 0.674 |
| 0.1 | 0.437 |
| 0.05 | 0.254 |

senting this result is in terms of the ratio $\mu_H/\mu$. Multiplying (38) by $c/\sigma_0\mu$ and using $\sigma_0 = ne\mu$ and (14) for $\mu$ we get

$$\frac{\mu_H}{\mu} = \frac{\langle\Delta\epsilon\tau^2\rangle\langle\Delta\epsilon\rangle}{\langle\Delta\epsilon\tau\rangle^2} \cdot \frac{3\left(\dfrac{1}{m_1{}^*m_2{}^*} + \dfrac{1}{m_2{}^*m_3{}^*} + \dfrac{1}{m_3{}^*m_1{}^*}\right)}{\left(\dfrac{1}{m_1{}^*} + \dfrac{1}{m_2{}^*} + \dfrac{1}{m_3{}^*}\right)^2} \qquad (39)$$

$$= \frac{\langle\Delta\epsilon\tau^2\rangle\langle\Delta\epsilon\rangle}{\langle\Delta\epsilon\tau\rangle^2} \cdot B, \qquad \text{say.}$$

Note that the first factor of (39) is the value of $\mu_H/\mu$ in the simple theory,[11] and that the second factor $B$, involving the anisotropy of the effective mass, is unity for zero anisotropy and $<1$ in general. Some sample values of this mass factor $B$ are given in Table III and Fig. 7. Fig. 8 gives values of the first factor in (39), for the simplified model of intra- and inter-valley scattering described by (24).

## 7. THE BENEDICT-SHOCKLEY EXPERIMENT

We turn now to the response of the assembly of carriers to an electric field which varies sinusoidally with time. As Benedict and Shockley have shown,[12] this response becomes limited at high frequencies by the inertia of the carriers, and so by measuring it one can obtain an effective mass.

[11] See, for example, Reference 7, p. 277.
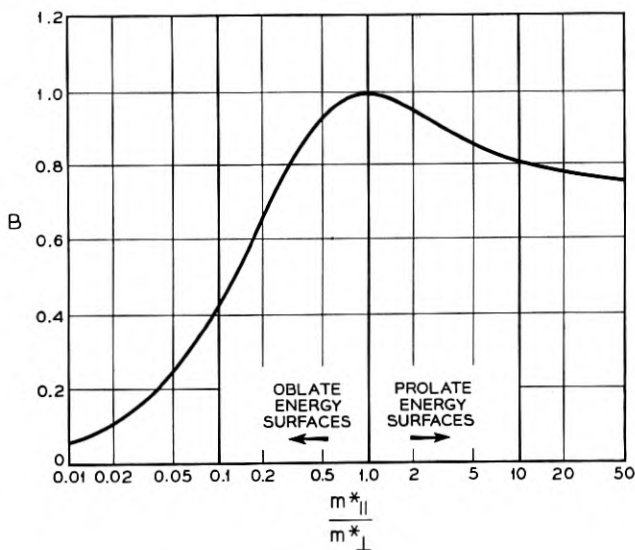[12] T. S. Benedict and W. Shockley, Phys. Rev., **89**, p. 1152, 1953.

Fig. 7 — Dependence of

$$B = \frac{3\left(\dfrac{1}{m_1{}^* m_2{}^*} + \dfrac{1}{m_2{}^* m_3{}^*} + \dfrac{1}{m_3{}^* m_1{}^*}\right)}{\left(\dfrac{1}{m_1{}^*} + \dfrac{1}{m_2{}^*} + \dfrac{1}{m_3{}^*}\right)^2}$$

on the anisotropy of the effective mass, for the case $m_1{}^* = m_2{}^* = m_\perp{}^*$, $m_3{}^* = m_\parallel{}^*$.

The solution of the transport equation for this case proceeds almost exactly as in Section 2. We assume that the scattering of the carrier is described by a relaxation time $\tau$, whose dependence on position in momentum space we shall for the moment leave unrestricted. The analysis starts as before from (3) for the distribution function $f$ of the carriers, namely, with the upper sign for electrons, the lower for holes,

$$\frac{\partial f}{\partial t} = \pm e\mathbf{E} \cdot \nabla_P f - \frac{(f - f^{(0)})}{\tau} \tag{40}$$

(We neglect the very small effect of the magnetic field generated by $\partial \mathbf{E}/\partial t$.) Instead of (4) we write, if $\mathbf{E} = \mathbf{E}_0 e^{i\omega \tau}$,

$$f(t) = f^{(0)} + E_0 \cdot \mathbf{f}^{(1)}(t) + O(E_0{}^2) \tag{41}$$

From (40) and (41) the equation for $f^{(1)}$ is

$$\frac{\partial \mathbf{f}^{(1)}}{\partial t} = \pm e\nabla_P f^{(0)} e^{i\omega t} - \frac{f^{(1)}}{\tau} \tag{42}$$
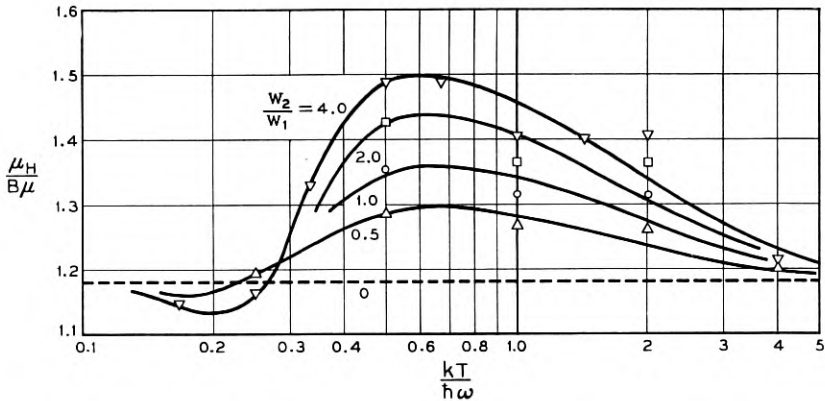
Fig. 8 — Values of the ratio

$$\frac{\mu_H}{B\mu} = \frac{\langle \Delta\epsilon\tau^2 \rangle \langle \Delta\epsilon \rangle}{\langle \Delta\epsilon\tau \rangle^2},$$

for the simplified lattice scattering law (24). The ratio $w_2/w_1$ measures the strength of the coupling of the carriers to inter-valley modes in terms of that for intra-valley scattering; $\omega$ is the frequency of the inter-valley modes. The curves have been drawn to smooth out irregularities (severe for $kT/\hbar\omega = 1$ to 2) in the calculated points.

If $\Delta\epsilon = |\epsilon - \epsilon_b|$ is the distance from the band edge, we have, for Maxwellian statistics, $f^{(0)} \propto \exp(-\Delta\epsilon/kT)$ and of course $\nabla_P \Delta\epsilon = \mathbf{v}$, the group velocity. Thus $\mathbf{f}^{(1)} = f_0^{(1)} e^{i\omega t}$ with

$$\mathbf{f}_0^{(1)} = \frac{\mp e\mathbf{v}f^{(0)}\tau}{kT(1 + i\omega\tau)} \tag{43}$$

The current density is given by the usual sum over the different valleys $i$ and the states ($\Delta\mathbf{P}^{(i)}$, spin) in each valley:

$$\mathbf{j} = \sum_i \sum_{\Delta\mathbf{P}^{(i)}, s} (\mp e)\mathbf{v}(\Delta\mathbf{P}^{(i)})f(\Delta\mathbf{P}^{(i)}) \tag{44}$$

From (41), (43) and (44),

$$\begin{aligned}
\mathbf{j} &= \sum_i \sum_{\Delta\mathbf{P}^{(i)}, s} \left( \frac{e^2\mathbf{v}\mathbf{v}\cdot\mathbf{E}_0 f^{(0)}\tau}{kT(1 + i\omega\tau)} \right) e^{i\omega\tau} \\
&= \frac{e^2\mathbf{E}_0}{3kT} \sum_i \sum_{\Delta\mathbf{P}^{(i)}, s} \left( \frac{v^2 f^{(0)}\tau}{(1 + i\omega\tau)} \right) e^{i\omega\tau}
\end{aligned} \tag{45}$$

if the crystal has cubic symmetry. Thus the semiconductor has the fre-

quency-dependent complex conductivity

$$\sigma(\omega) = \frac{ne^2}{3kT}\left\langle \frac{v^2\tau}{1 + i\omega\tau}\right\rangle \qquad (46)$$

where $n$ is the number of carriers per unit volume and the angular brackets denote a Maxwellian average. Note that in deriving (46) we have not had to assume anything about the dependence of $\epsilon$ or $\tau$ on $\Delta P$; (46) is therefore valid for all models, not merely for the many-valley case. However, (46) is still not explicit enough to be directly usable for the evaluation of experimental results, and we shall need to use the special properties of the many-valley model to express the Maxwellian average in terms of measurable or readily interpretable quantities.

Under the usual assumptions of this paper $\tau$ is a function of energy $\epsilon$ only, so the average in (46) depends only on the average of $v^2$ over an energy shell. By the equipartition principle (10) we may therefore replace $v^2$ by

$$2\left(\frac{\frac{1}{2}m_1^*v_1^2}{m_1^*} + \frac{\frac{1}{2}m_2^*v_2^2}{m_2^*} + \frac{\frac{1}{2}m_3^*v_3^2}{m_3^*}\right) = \frac{2}{3}\Delta\epsilon\left(\frac{1}{m_1^*} + \frac{1}{m_2^*} + \frac{1}{m_3^*}\right) \quad (47)$$

The average of the masses is the same $m^{(I)}$ we encountered in Section 2, namely,

$$\frac{1}{m^{(I)}} = \frac{1}{3}\left(\frac{1}{m_1^*} + \frac{1}{m_2^*} + \frac{1}{m_3^*}\right) \qquad (48)$$

Therefore

$$\left\langle\frac{v^2\tau}{1 + i\omega\tau}\right\rangle = \frac{2}{m^{(I)}}\left\langle\frac{\Delta\epsilon\tau}{1 + i\omega\tau}\right\rangle \qquad (49)$$

The expression for $\sigma(\omega)$ becomes, with its real and imaginary parts separated,

$$\sigma(\omega) = \frac{2ne^2}{3kTm^{(I)}}\left[\left\langle\frac{\Delta\epsilon\tau}{1 + \omega^2\tau^2}\right\rangle - i\omega\left\langle\frac{\Delta\epsilon\tau^2}{1 + \omega^2\tau^2}\right\rangle\right] \qquad (50)$$

The real part $\sigma_R(\omega)$ of this is what is usually called the "conductivity". The imaginary part $\sigma_I(\omega)$ is proportional to a contribution to the dielectric constant $\kappa(\omega)$, since $(4\pi)^{-1}\kappa(\omega)\partial E/\partial t$ is the sum of the displacement current $(4\pi)^{-1}\kappa_0\partial E/\partial t$ and the part of the true current $j$ which is in phase with $\partial E/\partial t$. Thus the departure of $\kappa(\omega)$ from the dielectric constant $\kappa_0$ of the crystal without its free carriers is given by

$$\kappa_0 - \kappa(\omega) = -\frac{4\pi}{\omega}\sigma_I(\omega) \qquad (51)$$

At the frequencies and temperatures which have been used for the Benedict-Shockley experiment, most of the carriers have relaxation times short compared to $\omega^{-1}$, and it is appropriate to make an expansion in powers of $\omega$:

$$\left\langle \frac{\Delta\epsilon\tau}{1 + \omega^2\tau^2} \right\rangle = \langle\Delta\epsilon\tau\rangle - \omega^2\langle\Delta\epsilon\tau^3\rangle \cdots \tag{52}$$

$$\left\langle \frac{\Delta\epsilon\tau^2}{1 + \omega^2\tau^2} \right\rangle = \langle\Delta\epsilon\tau^2\rangle - \omega^2\langle\Delta\epsilon\tau^4\rangle \cdots \tag{53}$$

(If $\tau \to \infty$ as $\Delta P \to 0$, as would be the case if the only scattering were by phonons of negligible energy, the series (52), (53) do not converge for any finite $\omega$. However, asymptotic series can be written down which differ only in order $\omega^3$ and higher from the series obtained by simply expanding the denominators). Denoting the dc conductivity by $\sigma_0$ — equal to $ne$ times the $\mu$ of (14) — we have from (50) to (53)

$$\sigma_R(\omega) = \sigma_0 \left[ 1 - \frac{\omega^2\langle\Delta\epsilon\tau^3\rangle}{\langle\Delta\epsilon\tau\rangle} + O(\omega^4) \right] \tag{54}$$

$$\kappa_0 - \kappa(\omega) = 4\pi\sigma_0 \left[ \frac{\langle\Delta\epsilon\tau^2\rangle}{\langle\Delta\epsilon\tau\rangle} - \omega^2 \frac{\langle\Delta\epsilon\tau^4\rangle}{\langle\Delta\epsilon\tau\rangle} + O(\omega^4) \right] \tag{55}$$

It is convenient to express the first term in the square bracket in (55) in terms of the Hall mobility $\mu_H$, since the same average $\langle\Delta\epsilon\tau^2\rangle$ occurs in (38) as in (55). Let us set $\sigma_0 = ne\mu$ and use the designation $B$ for the last factor in (39), a factor $\leqq 1$ dependent on the anisotropy of the effective mass in each valley and close to unity unless the anisotropy is very extreme (see Table III). Then

$$\kappa_0 - \kappa(\omega) = \frac{4\pi n\mu\mu_H m^{(I)}}{B} \left[ 1 - \omega^2 \frac{m^{(I)2}\mu\mu_H}{e^2 B} \frac{\langle\Delta\epsilon\tau^4\rangle\langle\Delta\epsilon\rangle}{\langle\Delta\epsilon\tau^2\rangle^2} + O(\omega^4) \right] \tag{56}$$

Equations (55) and (56), like all equations in this memorandum, is in Gaussian units. For rationalized MKS units, as used in the papers of Benedict and Shockley, the coefficient $4\pi$ should be replaced by $1/\epsilon_0$, where $\epsilon_0$ is the permittivity of the vacuum.

The leading term of (56) is the same as that which one would obtain by simply replacing $\mu^2$ by $\mu\mu_H/B$ in the formula used by Benedict and Shockley (simple model, $\tau$ = constant). But because the dimensionless factor $\langle\Delta\epsilon\tau^4\rangle\langle\Delta\epsilon\rangle/\langle\Delta\epsilon\tau^2\rangle^2$ is always $\geqq 1$ instead of $= 1$, the second term in the brackets in (56) is not the same as that resulting from this substitution. Thus the expression used by Benedict[13] in his later analysis of data
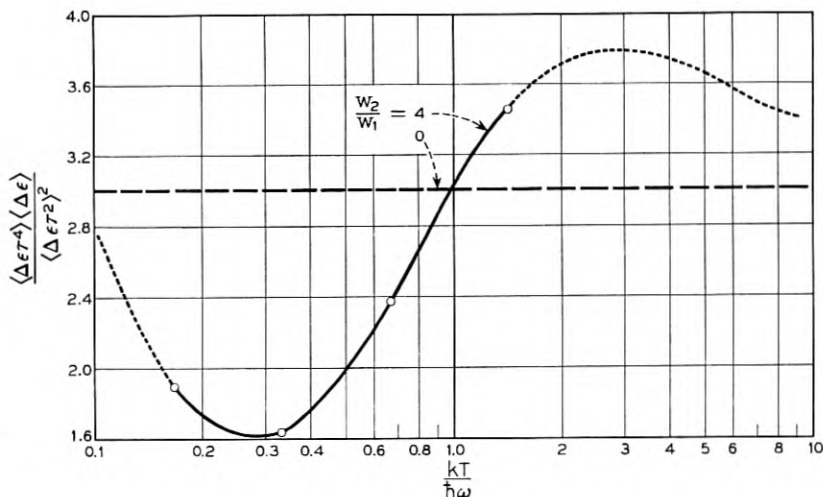
---

[13] T. S. Benedict, Phys. Rev., **91**, p. 1565, 1953.

Fig. 9 — Sample values of the quantity $\dfrac{\langle \Delta\epsilon\tau^4 \rangle \langle \Delta\epsilon \rangle}{\langle \Delta\epsilon\tau^2 \rangle^2}$ occurring in the equation (56) for the high-frequency dielectric constant, for the simplified lattice scattering law (24). The ratio $w_2/w_1$ measures the strength of the coupling of the carriers to inter-valley modes in terms of that for intra-valley scattering; $\omega$ is the frequency of the inter-valley modes. The dotted extrapolations are qualitative only.

on $p$ germanium is corrected only to the zeroth order in $\omega$ and to the approximation $B \approx 1$. Some sample values of the ratio $\langle \Delta\epsilon\tau^4 \rangle \langle \Delta\epsilon \rangle / \langle \Delta\epsilon\tau^2 \rangle^2$ have been computed for the scattering law (24), and are graphed in Fig. 9. These show that the range of possible variation of this factor is considerable.

A similar, though less useful, transformation can be made on (54), to express the second term in brackets in terms of the dimensionless coefficient $G$ defined by

$$G = \frac{\langle \Delta\epsilon\tau^3 \rangle \langle \Delta\epsilon \rangle^2}{\langle \Delta\epsilon\tau \rangle^3} \tag{57}$$

This is a coefficient which we shall encounter in the next section, in the theory of magnetoresistance, and which is graphed in Fig. 10 for the inter-valley scattering law (24). We find

$$\sigma_R(\omega) = \sigma_0 \left[ 1 - \omega^2 \frac{m^{(I)2} \mu^2}{e^2} G + O(\omega^4) \right] \tag{58}$$

The quantity $G$ is $\geqq 1$, the equality holding only if $\tau$ is a constant. Table

IV and Fig. 10 give some typical values, for scattering laws of the form $\tau \propto \Delta\epsilon^r$ or of the form (24).

## 8. LOW-FIELD MAGNETORESISTANCE

In Section 6 we set up the Boltzmann equation for the steady motion of charge carriers under the combined influence of an electric field $\mathbf{E}$ and a magnetic field $\mathbf{H}$, and solved it to the first order in $H$. We shall now undertake to solve this equation to the second and higher orders in $H$. The solution has been worked out independently for a number of cases by Abeles and Meiboom,[3] Shibuya,[3] and Shockley (unpublished). We shall not give all the details of the solution, especially at large $H$, as many of them can be found in the reference just mentioned. However, to emphasize some features not brought out in this previously published work we shall review the whole calculation briefly from the beginning.

The relation of theory and experiment in the area of magnetoresistance resembles that for piezoresistance, in that the tensor quantity which is
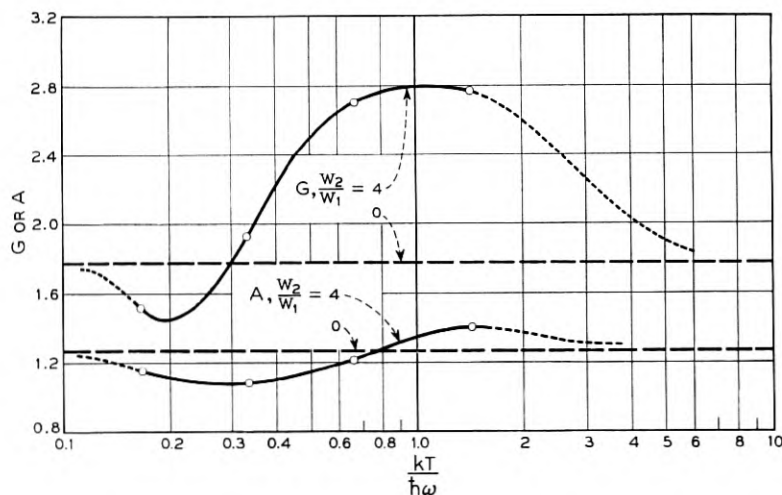


Fig. 10 — Sample values of

$$G = \frac{\langle \Delta\epsilon\tau^3 \rangle \langle \Delta\epsilon \rangle^2}{\langle \Delta\epsilon\tau \rangle^3} \quad \text{and} \quad A = \frac{\langle \Delta\epsilon\tau^3 \rangle \langle \Delta\epsilon\tau \rangle}{\langle \Delta\epsilon\tau^2 \rangle^2}$$

for the simplified lattice scattering law (24). The ratio $w_2/w_1$ measures the strength of the coupling of the carriers to inter-valley modes in terms of that for intra-valley scattering; $\omega$ is the frequency of the inter-valley modes. The $G$ curve has been drawn to be roughly consistent with the smoother $A$ curve and the curve of Fig. 8. The dotted extrapolations of the curves are intended only to show the expected qualitative behavior.

TABLE IV — SAMPLE VALUES OF THE QUANTITIES $G$ AND $A$ DEFINED BY (57) AND (68), RESPECTIVELY

| Scattering Law | $G = \dfrac{\langle\Delta\epsilon\tau^3\rangle\,\langle\Delta\epsilon\rangle^2}{\langle\Delta\epsilon\tau\rangle^3}$ | $A = \dfrac{\langle\Delta\epsilon\tau^3\rangle\,\langle\Delta\epsilon\tau\rangle}{\langle\Delta\epsilon\tau^2\rangle^2}$ |
|---|---|---|
| $\tau \propto \Delta\epsilon^{-0.6}$ | 2.58 | 4.01 |
| $\tau \propto \Delta\epsilon^{-0.5}$ | 1.77 | 1.27 |
| $\tau = $ constant | 1.00 | 1.00 |
| $\tau \propto \Delta\epsilon^{0.5}$ | 1.33 | 1.09 |
| $\tau \propto \Delta\epsilon$ | 2.52 | 1.28 |
| $\tau \propto \Delta\epsilon^{1.5}$ | 5.89 | 1.58 |
| Form (24), $w_2/w_1 = 4$, $kT/\hbar\omega = 1.43$ | 2.77 | 1.40 |
| $0.667$ | 2.70 | 1.21 |
| $0.333$ | 1.92 | 1.08 |
| $0.167$ | 1.52 | 1.15 |

simplest to calculate is reciprocal to the one which is directly measured. Thus one measures piezoresistance but calculates elastoresistance. Similarly the measured magnetoresistance is a change in the electric field $\mathbf{E}$ for given current, whereas the simplest quantity to calculate is the change of the current for given $\mathbf{E}$, i.e., the dependence of the conductivity tensor on $\mathbf{H}$. We shall see below that the neatest way of comparing theory and experiment, at least for small $H$, is to invert the observed magnetoresistivity tensor to get the magnetoconductivity tensor, and then compare the latter with theory.

The Boltzmann equation (30), from which we shall start, is

$$0 = \pm e\mathbf{E}\cdot\nabla_P f \pm e\,\frac{\mathbf{v}\times\mathbf{H}}{c}\cdot\nabla_P f - \frac{(f - f^{(0)})}{\tau} \qquad (59)$$

where as before the upper sign is for electrons, the lower for holes, $f$ is the distribution function of the carriers, $f^{(0)}$ the (Maxwellian) distribution in thermal equilibrium, $\mathbf{v} = \nabla_P \epsilon$ the group velocity. As in Section 2 we set

$$f = f^{(0)} + \mathbf{E}\cdot\mathbf{f}^{(1)} + O(E^2) \qquad (60)$$

and neglect higher order terms in $\mathbf{E}$. The resulting equation for $\mathbf{f}^{(1)}$ can be written in the condensed form[14,15]

$$0 = \pm \tau e\nabla_P f^{(0)} \pm \tau\mathbf{H}\cdot\boldsymbol{\gamma}\mathbf{f}^{(1)} - \mathbf{f}^{(1)} \qquad (61)$$

where

$$\boldsymbol{\gamma} = \frac{e}{c}\mathbf{v}\times\nabla_P \qquad (62)$$

In the notation of Davis,[14] and Seitz,[15] $\gamma = (e/\hbar^2 c)\Omega$, while in the notation of Abeles and Meiboom,[3] $\gamma = (e/c)\Omega$. The solution of (61) can be expressed formally in terms of the reciprocal of the operator $(1 \pm \tau\mathbf{H}\cdot\gamma)$:

$$\mathbf{f}^{(1)} = \pm(1 \pm \tau\mathbf{H}\cdot\gamma)^{-1}\tau e\nabla_P f^{(0)} \tag{63}$$

If we set

$$(1 \pm \tau\mathbf{H}\cdot\gamma)^{-1} = 1 \mp \tau\mathbf{H}\cdot\gamma + (\tau\mathbf{H}\cdot\gamma)(\tau\mathbf{H}\cdot\gamma)\cdots \tag{64}$$

we see that the leading term of (63) is just the solution (5) of Section 2 for $H = 0$, while the next is just the term (33) of Section 6. In other words, the series (64) corresponds to an iterative solution of (59). If we are interested only in the first few powers of $H$, this iterative solution is as simple as any; at high fields it is better to solve (61) explicitly in closed form, a procedure we shall outline in the next section.

The solution given by (63) and (64) of course applies for any dependence of the relaxation time $\tau$ and the energy $\epsilon$ on position $\mathbf{P}$ in crystal momentum space. However, we are here interested only in the case where, in each valley $i$, $\epsilon$ is a quadratic function of the components of $\Delta\mathbf{P} = \mathbf{P} - \mathbf{P}^{(i)}$ and where $\tau$ is a function of $\epsilon$ only. For this case some simplifications are possible. For one thing, $\tau$ commutes with the operator $\gamma$, since (62) acting on a function of $\epsilon$ contains the factor $\mathbf{v} \times \nabla_P\epsilon \equiv 0$. The expression for the current density $\mathbf{j}$ in powers of $H$ has the form

$$j_\mu = \sum_\nu \sigma_{\mu\nu}E_\nu + \sum_{\nu\alpha} \sigma_{\mu\nu\alpha}H_\alpha E_\nu + \sum_{\nu\alpha\beta} \sigma_{\mu\nu\alpha\beta}H_\alpha H_\beta E_\nu + \cdots \tag{65}$$

where of course $\sigma_{\mu\nu} = \sigma_0\delta_{\mu\nu}$ for a cubic substance, with $\sigma_0$ given by the equations of Section 2, and similarly $\sigma_{\mu\nu\alpha} = \sigma_0^2 R\delta_{\mu\nu\alpha}$, where $R$ is the low-field Hall constant and as in Section 6 $\delta_{\mu\nu\alpha} = \pm 1$ if $\mu\nu\alpha$ is an even (odd) permutation of 123, zero otherwise. To get the contribution of the $i$th valley to the second-order "magnetoconductivity" tensor $\sigma_{\mu\nu\alpha\beta}$, we multiply (60) by $\pm ev_\mu$, insert (63) and (64), and sum on all momentum vectors in the $i$th valley and in unit volume, and on spins. The result is

$$\sigma_{\mu\nu\alpha\beta}{}^{(i)} = \frac{e^2}{kT} \sum_{\Delta\mathbf{P}^{(i)},s} f^{(0)}\tau^3(v_\mu\gamma_\alpha\gamma_\beta v_\nu)_{\text{symm}} \tag{66}$$

where as usual we have assumed $f^{(0)}$ to be Maxwellian, and where the subscript "symm" means that the expression in parentheses is to be averaged with the expressions obtained from it by permuting $\alpha$ with $\beta$, since only the part of $\sigma_{\mu\nu\alpha\beta}$ symmetrical in $\alpha$ and $\beta$ has physical signifi-

[14] L. Davis, Phys. Rev., **56**, p. 93, 1939.
[15] F. Seitz, Phys. Rev., **79**, p. 372, 1950.

cance. (This symmetrization is not necessary, but simplifies the work by preventing the appearance of meaningless components.)

The explicit evaluation of (66) is a straightforward but tedious exercise in algebra, and will not be given in detail here. However, there are some important properties of the tensor (66) which can be established rather simply. Since the components of $\mathbf{v}$ are linear functions of the components of $\Delta\mathbf{P}$, the operator $\gamma$ defined by (62) takes any linear function of the $\Delta P_\lambda$ into another linear function. Therefore the $v_\mu\gamma_\alpha\gamma_\beta v_\nu$ in (66) is a quadratic function of the $\Delta P_\lambda$, and it is easily seen that this function contains denominators of the fourth degree in the effective masses. Now a quadratic function of the $\Delta P_\lambda$ can be written as an effective mass times the energy $\Delta\epsilon$ relative to the band edge, times a function of the direction of $\Delta\mathbf{P}$, dependent only on the ratios of the effective masses in the principal directions. Thus we may write, for example,

$$v_\mu\gamma_\alpha\gamma_\beta v_\nu = \frac{\Delta\epsilon}{m^{(I)3}} \times \text{function } (\mu\nu\alpha\beta, \text{ direction of } \Delta\mathbf{P}, \text{ mass ratios}) \quad (67)$$

where $m^{(I)}$ is the inertial average of the effective masses, defined by (15). Now let the summation on $\Delta\mathbf{P}$ be broken up into a summation over values in an energy shell $\Delta\epsilon$ to $\Delta\epsilon + d\Delta\epsilon$, and a summation over different shells. The function of direction in (67) will be the same for all the shells, and so we have the result

$$\sigma_{\mu\nu\alpha\beta}{}^{(i)} = \frac{a}{m^{(I)3}} F_{\mu\nu\alpha\beta}{}^{(i)}$$

where $F_{\mu\nu\alpha\beta}{}^{(i)}$ depends on the anistropy ratios of the effective masses in the $i$th valley, but not on the variation of $\tau$ with energy, while $a$ is proportional to the number of carriers in the valley and to the Maxwellian average of $\tau^3\Delta\epsilon$.

It is convenient to express the average of $\tau^3\Delta\epsilon$ in dimensionless form by using the quantity $G$ defined by (57), namely,

$$G = \frac{\langle\Delta\epsilon\tau^3\rangle\langle\Delta\epsilon\rangle^2}{\langle\Delta\epsilon\tau\rangle^3}$$

or else the quantity

$$A = \frac{\langle\Delta\epsilon\tau^3\rangle\langle\Delta\epsilon\tau\rangle}{\langle\Delta\epsilon\tau^2\rangle^2} \quad (68)$$

Here as usual the angular brackets denote Maxwellian averages as defined in connection with (14). We may use (14) to eliminate $m^{(I)}$ and the carrier density and if we wish we may eliminate $\mu$ in favor of $\mu_H$ by (39).

The result is

$$\sigma_{\mu\nu\alpha\beta}^{(i)} = G \frac{(\sigma_0\mu^2)}{(N_V c^2)} B^2 F_{\mu\nu\alpha\beta}^{(i)} = A \frac{(\sigma_0\mu_H^2)}{(N_V c^2)} F_{\mu\nu\alpha\beta}^{(i)} \qquad (69)$$

where $\sigma_0$, $\mu$, $\mu_H$, are the conductivity, mobility, and Hall mobility, respectively, at $H = 0$, $N_V$ is the number of valleys, and where $B \leq 1$ is the function of the effective mass ratios defined in (39) and Table III. Summing on valleys $i$ gives

$$\sigma_{\mu\nu\alpha\beta} = A \left(\sigma_0 \frac{\mu_H^2}{c^2}\right) F_{\mu\nu\alpha\beta} = G \left(\sigma_0 \frac{\mu^2}{c^2}\right) B^2 F_{\mu\nu\alpha\beta} \qquad (70)$$

$$F_{\mu\nu\alpha\beta} = \frac{1}{N_V} \sum_i F_{\mu\nu\alpha\beta}^{(i)} \qquad (71)$$

Note that $F_{\mu\nu\alpha\beta}$, as defined by (70) or (71), is dimensionless, as are $G$ and $A$; $F_{\mu\nu\alpha\beta}$ or $B^2 F_{\mu\nu\alpha\beta}$ depends on the geometry of the valleys and the ratios of the principal masses of a valley. We shall see presently how the analysis of experimental data is facilitated by this decomposition of the magnetoconductivity into the product of a scalar factor depending on the behavior of $\tau$ and a tensor factor depending on the shape of the energy surfaces.

The quantity $A$ defined by (68), like $G$, is $\geq 1$, the equality holding only if $\tau$ is a constant. Some sample graphs of $A$ and $G$ are shown in Fig. 10, for scattering laws of the form (24), and some numerical values for this case and for $\tau \propto \Delta\epsilon^r$ are given in Table IV. Note that for the ideal case of intra-valley lattice scattering only, a case approximated in very pure material at moderately low $T$, $r = -\frac{1}{2}$ and $A = 4/\pi = 1.27$, $G = 9\pi/16 = 1.77$. Table V gives values of all the nonvanishing coefficients $F_{\mu\nu\alpha\beta}^{(i)}$ relative to a coordinate system oriented along the principal axes of a valley. The middle rows of Table VI give the $F_{\mu\nu\alpha\beta}$, relative to the crystal axes, for some of the simpler possible arrangements of valleys. The entries were obtained, of course, by comparing (69) or (70) with the results of explicit evaluations of (66). For completeness, Table V also gives the directional factors involved in the contribution $\sigma_{\mu\nu}^{(i)}$ of a single valley to the conductivity tensor $\sigma_{\mu\nu}$ in the absence of a magnetic field, and to the Hall conductivity tensor $\sigma_{\mu\nu\alpha}$ defined by (34) or (65). All these table entries are similar to those given by Abeles and Meiboom.[3] However, they have given the unsymmetrized $\sigma_{\mu\nu\alpha\beta}$ etc. for the cases $r = -\frac{1}{2}$ and $+\frac{3}{2}$, in terms of mean free path, absolute values of the masses, and carrier concentration; here we have given the symmetrized $\sigma_{\mu\nu\alpha\beta}$ etc. in terms of the directly observable $\sigma_0$ and $\mu_H$, and for any $\tau(\Delta\epsilon)$.

TABLE V — ANISOTROPY FACTORS FOR CONDUCTION, HALL EFFECT, AND MAGNETORESISTANCE CONTRIBUTIONS FROM A SINGLE VALLEY

| Phenomenon and Text Reference | Type of Valley | Tensor Component | Value | Relative Value |
|---|---|---|---|---|
| Conduction, (11) | Any | $\mu_{\alpha\beta}^{(i)}$ | | $\dfrac{\delta_{\alpha\beta}}{m_\alpha^*}$ |
| Hall effect (36) | Any | $\sigma_{\lambda\mu\nu}^{(i)}$ | | $\dfrac{\delta_{\lambda\mu\nu}}{m_\lambda^* m_\mu^*}$ |
| | $\left.\begin{array}{l} m_1^* = m_2^* = m_\perp^* \\ m_3^* = m_\parallel^* \end{array}\right\}$ | $\sigma_{312}^{(i)} = \sigma_{231}^{(i)}$ | | $\dfrac{1}{m_\perp^{*2}}$ $\dfrac{1}{m_\parallel^* m_\perp^*}$ |
| Magnetoresist-ance, (69) | Any | $F_{\alpha\alpha\alpha\alpha}^{(i)}$ | 0 | 0 |
| | | $F_{\alpha\beta\alpha\beta}^{(i)} = F_{\alpha\beta\beta\alpha}^{(i)}\,(\beta \neq \alpha)$ | $-\dfrac{3(m_1^* m_2^* + m_1^* m_3^* + m_2^* m_3^*)\,m_\beta^*}{(m_1^* + m_2^* + m_3^*)^2\,m_\alpha^*}$ | $-\dfrac{m_\beta^*}{m_\alpha^*}$ |
| | | $F_{1122}^{(i)}$ | $\tfrac{3}{2}\dfrac{(m_1^* m_2^* + m_1^* m_3^* + m_2^* m_3^*)}{(m_1^* + m_2^* + m_3^*)^2}$ | $\tfrac{1}{2}$ |
| | $\left.\begin{array}{l} m_1^* = m_2^* = m_\perp^* \\ m_3^* = m_\parallel^* \end{array}\right\}$ | $F_{1133}^{(i)} = F_{2233}^{(i)}$ | $-\dfrac{3(m_\perp^* + 2m_\parallel^*)\,m_\perp^*}{(m_\parallel^* + 2m_\perp^*)^2}$ | $-1$ |
| | | $F_{3311}^{(i)} = F_{3322}^{(i)}$ | $-\dfrac{3(m_\perp^* + 2m_\parallel^*)\,m_\parallel^*}{(m_\parallel^* + 2m_\perp^*)^2}$ | $-\dfrac{m_\parallel^*}{m_\perp^*}$ |
| | | $F_{1212}^{(i)} = F_{1313}^{(i)} = F_{2323}^{(i)}$ | $-\dfrac{3(m_\perp^* + 2m_\parallel^*)\,m_\perp^{*2}}{(m_\parallel^* + 2m_\perp^*)^2\,m_\parallel^*}$ | $-\dfrac{m_\perp^*}{m_\parallel^*}$ |
| | | | $\tfrac{3}{2}\dfrac{(m_\perp^* + 2m_\parallel^*)\,m_\perp^*}{(m_\parallel^* + 2m_\perp^*)^2}$ | $\tfrac{1}{2}$ |

TABLE VI — EXPRESSIONS FOR THE $F_{\mu\nu\alpha\beta}$ OF (70) IN TERMS OF THE MASS RATIOS, AND FOR THE $\sigma_{\mu\nu\alpha\beta}$ OF (65) IN TERMS OF THE EMPIRICAL MAGNETORESISTANCE CONSTANTS $b$, $c$, $d$, OF (74)

| $\mu\nu\alpha\beta$ | $\alpha\alpha\alpha\alpha$ | $\alpha\alpha\beta\beta$ | $\alpha\beta\alpha\beta$ or $\alpha\beta\beta\alpha$ | Other |
|---|---|---|---|---|
| $F_{\mu\nu\alpha\beta}$, (100) valleys | $0$ | $-\dfrac{\left(1+2\dfrac{m_\|^*}{m_\perp^*}\right)\left[1+\dfrac{m_\|^*}{m_\perp^*}+\left(\dfrac{m_\|^*}{m_\perp^*}\right)^2\right]}{\dfrac{m_\|^*}{m_\perp^*}\left(2+\dfrac{m_\|^*}{m_\perp^*}\right)^2}$ | $\dfrac{\left(1+2\dfrac{m_\|^*}{m_\perp^*}\right)}{\left(2+\dfrac{m_\|^*}{m_\perp^*}\right)^2}$ | $0$ |
| Relative value | $0$ | $\left[\left(\dfrac{m_\|^*}{m_\perp^*}\right)^{-1}+1+\dfrac{m_\|^*}{m_\perp^*}\right]$ | $\dfrac{3}{2}$ | $0$ |
| $F_{\mu\nu\alpha\beta}$, (111) valleys | $-\dfrac{2}{3}\dfrac{\left(\dfrac{m_\|^*}{m_\perp^*}-1\right)^2\left(1+\dfrac{2m_\|^*}{m_\perp^*}\right)}{\dfrac{m_\|^*}{m_\perp^*}\left(2+\dfrac{m_\|^*}{m_\perp^*}\right)^2}$ | $-\dfrac{1}{3}\dfrac{\left(1+2\dfrac{m_\|^*}{m_\perp^*}\right)\left[2+5\dfrac{m_\|^*}{m_\perp^*}+2\left(\dfrac{m_\|^*}{m_\perp^*}\right)^2\right]}{\dfrac{m_\|^*}{m_\perp^*}\left(2+\dfrac{m_\|^*}{m_\perp^*}\right)^2}$ | $\dfrac{1}{6}\dfrac{\left(1+2\dfrac{m_\|^*}{m_\perp^*}\right)\left[2+5\dfrac{m_\|^*}{m_\perp^*}+2\left(\dfrac{m_\|^*}{m_\perp^*}\right)^2\right]}{\dfrac{m_\|^*}{m_\perp^*}\left(2+\dfrac{m_\|^*}{m_\perp^*}\right)^2}$ | $0$ |
| Relative value | $-4\left(\dfrac{m_\|^*}{m_\perp^*}-1\right)^2$ | $-2\left[2+5\dfrac{m_\|^*}{m_\perp^*}+2\left(\dfrac{m_\|^*}{m_\perp^*}\right)^2\right]$ | $\left[2+5\dfrac{m_\|^*}{m_\perp^*}+2\left(\dfrac{m_\|^*}{m_\perp^*}\right)^2\right]$ | $0$ |
| $\dfrac{\sigma_{\mu\nu\alpha\beta}}{\sigma_0}$ | $-(b+c+d)$ | $-b-\left(\dfrac{\mu_H}{c}\right)^2$ | $-\tfrac{1}{2}c+\tfrac{1}{2}\left(\dfrac{\mu_H}{c}\right)^2$ | $0$ |

The coordinate axes are assumed oriented along the cubic axes of the crystal, and $\alpha$, $\beta$ represent any two *different* ones of the three directions $x$, $y$, $z$. In the last row of the table we have adhered to established notation at the cost of a double meaning for the letter $c$: in the expressions $\mu_H/c$ it is the velocity of light; elsewhere it is the magnetoresistance constant defined by (74).

The qualitative behavior of the entries in Table V is easily understandable. The components $F_{\alpha\alpha\alpha\alpha}{}^{(i)}$ refer to the longitudinal magnetoconductivity when both electric and magnetic fields are in one of the principal directions of the valley, the $\alpha$ direction. Since a magnetic field in such an $\alpha$ direction does not change the $\alpha$ component of the velocity of the carrier, this longitudinal magnetoconductivity must vanish:

$$F_{\alpha\alpha\alpha\alpha}{}^{(i)} = 0.$$

It is easily verified that, for our model, the principal directions of a valley are the only directions in which the longitudinal magnetoconductivity contribution vanishes. Since the relative longitudinal magnetoconductivity $\Delta\sigma/\sigma$ is necessarily $\leq 0$, and for a cubic crystal is the negative of the relative longitudinal magnetoresistivity $\Delta\rho/\rho$, we can conclude that on our model the vanishing of the longitudinal magnetoresistance in any direction is possible, at least for cubic materials, only if the direction in question is a principal axis of all the valleys. It can further be shown, though we shall not give the details here, that lack of constancy of $\tau$ over an energy shell, far from upsetting this conclusion, merely makes it impossible for the longitudinal magnetoresistance to vanish in *any* direction.

The nonvanishing magnetoconductance effects can be described in terms of the current due to the force exerted by the magnetic field on the transverse Hall current. This current is proportional to the Hall current and inversely proportional to the effective mass — call it $m_J{}^*$ — in the direction normal to the Hall current and to **H**, this being the direction of the force producing the second-order current. The Hall current, as we have noted in Section 6, is proportional to the zero-order current, hence to the reciprocal of the effective mass $m_E{}^*$ in the direction of **E**, and to the reciprocal of the effective mass $m_{E\cdot H}{}^*$ in the direction normal to **E** and **H**.

To employ these ideas specifically, consider first the component $\sigma_{\alpha\alpha\beta\beta}{}^{(i)}$, which measures the change in current in the $\alpha$ direction produced by a magnetic field in the $\beta$ direction. Here $m_E{}^* = m_\alpha{}^*$, $m_{E\cdot H}{}^* = m_\gamma{}^*(\gamma \neq \alpha, \beta)$, $m_J{}^* = m_\alpha{}^*$. Thus

$$\sigma_{\alpha\alpha\beta\beta}{}^{(i)} \propto -\frac{1}{m_\alpha{}^{*2}m_\gamma{}^*}\ (\gamma \neq \alpha, \beta) \tag{72}$$

the minus sign coming in because the second-order current is in the direction opposite to **E**. When we insert the mass-dependence of $\sigma_0\mu_H{}^2$ into (69) and combine with the $F_{\alpha\alpha\beta\beta}{}^{(i)}$ of Table V, we do in fact find that $\sigma_{\alpha\alpha\beta\beta}{}^{(i)}$ contains the masses only as indicated in (72). (The fact that

$F_{\alpha\alpha\beta\beta}^{(i)}$ depends on the masses in a much more complicated way is due to our choice of the defining equation for it, (69): we chose to write this equation so that it involved only the directly measurable quantities $\sigma_0$ and $\mu_H$, and the dimensionless quantities $A$ and $F_{\mu\gamma\alpha\beta}^{(i)}$. This choice is the most convenient one for comparisons with experiment, but is less simple conceptually than a choice giving the factor (72).) The remaining independent component, $F_{\alpha\beta\alpha\beta}^{(i)}$, may be analyzed similarly. It represents the second-order current in the $\alpha$ direction due to an $\mathbf{E}$ in the $\beta$ direction and an $\mathbf{H}$ in the direction midway between $\alpha$ and $\beta$. Here $m_E^* = m_\beta^*$, $m_{E \cdot H}^* = m_\gamma^* (\gamma \neq \alpha, \beta)$, $m_J^* = m_\alpha^*$. The Hall current is weaker by $2^{1/2}$ than for the previous case, because of the 45° angle between $\mathbf{E}$ and $\mathbf{H}$, and since the force producing the second-order current is at 45° to the $\beta$ direction, we must put a second $2^{1/2}$ in the denominator. Thus

$$\sigma_{\alpha\beta\alpha\beta}^{(i)} \propto \frac{1}{2m_\alpha^* m_\beta^* m_\gamma^*} = \frac{1}{2m_1^* m_2^* m_3^*} \tag{73}$$

This, again, can be verified to follow from (69) and Table V.

To apply these results to experimental magnetoresistance data in the region of proportionality to $H^2$, it is necessary, as has been mentioned above, to derive an experimental magnetoconductivity tensor from the observed magnetoresistance. For a cubic crystal the magnetoresistivity tensor can be described by three constants $b$, $c$ (not to be confused with velocity of light), $d$, defined by

$$\frac{\Delta\rho}{\rho H^2} = bH^2 + c\frac{(\mathbf{j} \cdot \mathbf{H})^2}{j^2} + d\frac{(j_x^2 H_x^2 + j_y^2 H_y^2 + j_z^2 H_z^2)}{j^2} \tag{74}$$

where $\Delta\rho$ is the change of resistivity $\rho$ due to a small field $\mathbf{H}$, and where the axes are those of the crystal. The equations relating the constants $b, c, d$ to the corresponding constants describing the magnetoconductivity tensor have been given by Pearson and Suhl.[16] From these equations the components of $\sigma_{\mu\nu\alpha\beta}$ can be expressed in terms of the empirical constants $b, c, d$. The results are tabulated in the last row of Table VI.

If the ratios of these components are compared with the ratios of the corresponding $F_{\mu\nu\alpha\beta}$, one can check the correctness of an assumed model and determine the ratio $m_\parallel^*/m_\perp^*$. From the absolute values of the $\sigma_{\mu\nu\alpha\beta}$ one can then determine $A$. A further check is provided if data are available at more than one temperature, since the mass ratio should come out roughly independent of $T$, while the variation of $A$ with $T$ should

[16] G. L. Pearson and H. Suhl, Phys. Rev., **83**, p. 768, 1951.

accord with a reasonable picture of the effects of impurity, inter-, and intra-valley scattering on the form of $\tau(\Delta\epsilon)$.

An analysis of this sort has been carried out for n type silicon.[17] For this substance the longitudinal magnetoresistance nearly vanishes in directions of the type (100). From what has been said above, this almost requires that the band structure have valleys on the (100) axes in **K**-space, and that the relaxation time be practically a function of energy only. Fitting the remaining magnetoconductivity constants gives $m_{\parallel}^*/m_{\perp}^* \approx 5$; this ratio comes out independent of temperature as it should. It agrees with the ratio determined by cyclotron resonance.[18]

## 9. HALL EFFECT AND MAGNETORESISTANCE FOR LARGE MAGNETIC FIELDS

As the theory of Hall and magnetoresistance effects for large magnetic fields is rather complicated mathematically, it will suffice for our purposes merely to outline the approach which can be used and to quote a few results without proof. Some of the details can be found in the papers of Abeles and Meiboom[3] and of Shibuya.[3]

To treat these we solve the transport equation (61) for the distribution function $\mathbf{f}^{(1)}$ and calculate the current density from $\mathbf{f}^{(1)}$. This gives the electrical conductivity tensor $\sigma_{\mu\nu}(\mathbf{H})$, which can be inverted to give the resistivity tensor $\rho_{\mu\nu}$. The antisymmetric part of $\rho_{\mu\nu}$ determines a Hall coefficient (in general slightly orientation-dependent),' and the symmetrical part determines the magnetoresistance.

The solution of (61) can be carried out either by summing the series (64), or directly by guessing that $\mathbf{f}^{(1)}$ will be a linear function of the velocity components, with coefficients which are functions of energy. These coefficients can be determined by solving a set of three simultaneous equations.

As $H \rightarrow \infty$, the conductivity tensor $\sigma_{\mu\nu}(\mathbf{H})$ becomes singular, the contribution $\sigma_{\mu\nu}^{(i)}$ of the $i$th valley taking the form

$$\sigma_{\mu\nu}^{(i)} \rightarrow \frac{ne^2}{N_V} \frac{\langle\Delta\epsilon\tau\rangle}{\langle\Delta\epsilon\rangle} \frac{H_\mu H_\nu}{m_1^*H_1^2 + m_2^*H_2^2 + m_3^*H_3^2} \qquad \text{in general}$$

$$= \frac{\sigma_0}{N_V} \frac{m^{(I)}H_\mu H_\nu}{m_1^*H_1^2 + m_2^*H_2^2 + m_3^*H_3^2} \qquad \text{for a cubic crystal} \qquad (75)$$

where $n$ is the total density of carriers, $N_V$ the number of valleys, $\Delta\epsilon$ the distance from the band edge, $\tau(\Delta\epsilon)$ the relaxation time, $\sigma_0$ the conductivity at $H = 0$, $m^{(I)}$ the average inertial mass defined by (15), and

---

[17] G. L. Pearson and C. Herring, Physica, to appear.
[18] Dexter, Lax, Kip, and Dresselhaus, see Reference 5.

$H_1$, $H_2$, $H_3$ are the components of $\mathbf{H}$ along the principal axes of the valley. Summing (75) on $i$ gives a limiting $\sigma_{\mu\nu}$ which is still proportional to $H_\mu H_\nu$; its determinant therefore vanishes, and it cannot be inverted to give the limiting $\rho_{\mu\nu}$. It turns out that (75) suffices to give the limiting value of the longitudinal magnetoresistance as $H \to \infty$, but that to get the limiting value of the transverse magnetoresistance it is necessary to evaluate contributions to $\sigma_{\mu\nu}(\mathbf{H})$ of order $1/H^2$. The limiting Hall coefficient can be obtained from the contributions of order $1/H$ to $\sigma_{\mu\nu}(\mathbf{H})$.

The following points are noteworthy: (i) The limiting value of the Hall coefficient as $H \to \infty$ is $R = \mp(1/nec)$, independently of the arrangement and mass anisotropy of the valleys, and of the dependence of $\tau$ on energy. Even if $\tau$ is not a mere function of energy, the same limiting form obtains. (ii) For $\mathbf{E}$ parallel to $\mathbf{H}$, the ratio $\sigma(H \to \infty)/\sigma_0$ depends on the arrangement and mass anisotropy of the valleys, but not on the law of variation of $\tau$ with energy, as long as $\tau$ is a function of energy only. (iii) For $\mathbf{E}$ not parallel to $\mathbf{H}$, both $\sigma(H \to \infty)$ and $\sigma(H \to \infty)/\sigma_0$ depend on the behavior of $\tau(\Delta\epsilon)$ as well as on the band structure.

## Appendix A

### THE RELAXATION TIME ASSUMPTION

The use of the relaxation time concept represents a great simplification of the Boltzmann equation for all kinds of transport phenomena. To be strictly correct, one ought to describe the scattering processes which the carriers undergo by a transition probability $S(\mathbf{K}, \mathbf{K}')$ defined as the probability per unit time of a transition from an initial state $\mathbf{K}$ to a final state in unit volume of wave number space centered on $\mathbf{K}'$. (Conservation of energy will usually limit these transitions to a certain surface in $\mathbf{K}'$-space, so that $S$ will contain a delta function of energy; however, this complica-

tion has no bearing on the remarks of this paragraph.) The Boltzmann equation, for the case of Maxwellian statistics, therefore takes the form

$$\frac{\partial f(\mathbf{K})}{\partial t} = \left(\frac{\partial f(\mathbf{K})}{\partial t}\right)_{\text{Fields}} + \int \left[f(\mathbf{K}')S(\mathbf{K}', \mathbf{K}) - f(\mathbf{K})S(\mathbf{K}, \mathbf{K}')\right] d\mathbf{K}'. \quad \text{(A1)}$$

This is in general an integral equation for the distribution function $f$, or for the part $\mathbf{E} \cdot \mathbf{f}^{(1)}$ of $f$ which is linear in the electric field. Our task in this Appendix is to say a few words about the validity of approximating (A1) by an equation of the form

$$\frac{\partial f(\mathbf{K})}{\partial t} = \left(\frac{\partial f(\mathbf{K})}{\partial t}\right)_{\text{Fields}} - \left(\frac{f(\mathbf{K}) - f^{(0)}(\mathbf{K})}{\tau(\mathbf{K})}\right) \quad \text{(A2)}$$

and to examine the validity of the further approximation $\tau(\mathbf{K}) = \tau(\epsilon)$. We shall give only a rather brief discussion of these questions, however, as a future publication[6] will give a more thorough treatment and include a discussion of the solution of (A1) when these approximations fail.

To begin with, let us consider the special class of cases for which

$$S(\mathbf{K}', \mathbf{K}) = S(\mathbf{K}', \mathbf{K}^*) \quad \text{(A3)}$$

where $\mathbf{K}^*$ is the state in the same valley as $\mathbf{K}$ but with opposite velocity. If (4) is inserted for $f$ in the integral of (A1), the $f^{(0)}$ term contributes nothing, and if (A3) is satisfied

$$\int \mathbf{E} \cdot \mathbf{f}^{(1)}(\mathbf{K}')S(\mathbf{K}', \mathbf{K}) \, d\mathbf{K}' = 0, \quad \text{(A4)}$$

since $\mathbf{f}^{(1)}$ is an odd function of velocity while $S$ is even. Therefore, to the first order in $E$, (A1) reduces to (A2) with

$$1/\tau(\mathbf{K}) = \int S(\mathbf{K}, \mathbf{K}') \, d\mathbf{K}' \quad \text{(A5)}$$

For collision processes which do not satisfy the velocity-randomizing condition (A3) we may assess roughly the legitimacy of using a $\tau(\epsilon)$ by comparing, for different initial states on the same constant energy surface, the mean rate with which scattering destroys the component of velocity in the original direction. This rate of loss of velocity defines a $1/\tau$ which if isotropic is known to be usable in (A2) for energy-conserving scattering processes in the simple model,[19] and for the many-valley model it is a reasonable presumption, borne out by the more rigorous treatment

---

[19] W. Shockley, *Electrons and Holes in Semiconductors* (Van Nostrand 1951) p. 251 et seq.

of Reference 6, that if this rate of loss of velocity is nearly constant over an energy surface, then (A2) is legitimate with $\tau = \tau(\epsilon)$.

We shall consider five types of scattering in turn, three by phonons and two by impurities. For the phonon processes, scattering of a carrier from $\mathbf{K}$ to $\mathbf{K}'$ involves absorption of a phonon of wave vector $\mathbf{K}' - \mathbf{K}$, or emission of one with the negative wave vector. The matrix element for such a process is of the form[20]

$$M(\mathbf{K}, \mathbf{K}') = \left. \begin{matrix} N^{1/2} \\ (N+1)^{1/2} \end{matrix} \right\} \frac{C(\mathbf{K}, \mathbf{K}')}{\omega^{1/2}} \text{ for} \left\{ \begin{matrix} \text{absorption} \\ \text{emission} \end{matrix} \right. \tag{A6}$$

where $N$ is the occupation number of the lattice mode involved, $\omega$ is its frequency, and $C(\mathbf{K}, \mathbf{K}')$ is proportional to the matrix element, between states $\mathbf{K}$ and $\mathbf{K}'$, of the perturbation of the electronic Hamiltonian due to a static displacement of the nuclei of the lattice, of unit amplitude in this mode. The scattering function $S(\mathbf{K}, \mathbf{K}')$ which we have used above is given by

$$S(\mathbf{K}, \mathbf{K}') = \frac{2\pi}{\hbar} \sum | M(\mathbf{K}, \mathbf{K}') |^2 \delta[\epsilon(\mathbf{K}) - \epsilon(\mathbf{K}') \pm \hbar\omega] \tag{A7}$$

where $\delta$ is the Dirac delta function and the summation is over absorption (upper sign) and emission (lower sign), and over the various branches of the vibrational spectrum.

(i) *Inter-Valley Lattice Scattering.* Here the magnitude of $\mathbf{K}' - \mathbf{K}$ is large compared with the distance of either $\mathbf{K}$ or $\mathbf{K}'$ from the band edge point nearest it. Moreover, the total change in $\hbar\omega$ as $\mathbf{K}$ or $\mathbf{K}'$ ranges over a constant energy surface will be $\ll kT$ if the energy relative to the band edge is $\sim kT$. Therefore the percentage variation of $M(\mathbf{K}, \mathbf{K}')$ over such an energy surface will be small. And unless the energy relative to the band edge is extremely small, the surface in $\mathbf{K}'$-space which makes the argument of the delta function vanish will be nearly a constant energy surface. We conclude that $S(\mathbf{K}, \mathbf{K}')$ can be taken to be independent of $\mathbf{K}$ and $\mathbf{K}'$ when either ranges over a constant energy surface in its valley, and in particular, (A3) applies. Therefore inter-valley lattice scattering is described by a relaxation time which is given by (A5) and is a function of energy only.

(ii) *Intra-Valley Scattering by Optical Modes.* For a nonpolar crystal this case is essentially the same as the preceding, since the matrix element $C(\mathbf{K}, \mathbf{K}')$ is substantially equal to its limiting value as $\mathbf{K}' \to \mathbf{K}$. For a polar material the longitudinal polar optical modes give a $C(\mathbf{K}, \mathbf{K}') \propto$

---

[20] See, for example, Reference 19, p. 520.

$| \mathbf{K}' - \mathbf{K} |^{-1}$. In a cubic crystal this contribution is independent of the direction of $(\mathbf{K}' - \mathbf{K})$, but if the energy surfaces are anisotropic the variation of $C(\mathbf{K}, \mathbf{K}')$ with $| \mathbf{K}' - \mathbf{K} |$ will suffice to prevent $\tau(\mathbf{K})$ from being constant over an energy surface.

(iii) *Intra-Valley Scattering by Acoustical Modes.* Here (A6) and (A7) again apply, but with the simplification that for all but extremely slow carriers the $\hbar\omega$ in the argument of the delta function can be neglected, since only very low-energy phonons are involved. However, $C(\mathbf{K}, \mathbf{K}')$ need no longer be independent of direction. The reason for this is that according to the deformation potential concept,[21] $C$ depends on the strain associated with the lattice mode in question. If only the volume dilatation affected $C$, as was the case for the simple model, $C$ would be independent of the direction of the phonon wave vector $(\mathbf{K} - \mathbf{K}')$, since the dilatation amplitude in a compressional wave of unit displacement amplitude is independent of direction. But for a many-valley model both shears and dilatations can produce deformation potentials, and the nature of the shear strain in a shear wave of unit amplitude depends strongly on the direction of propagation. Therefore $C$ may be a function of the direction of $(\mathbf{K} - \mathbf{K}')$, in any particular valley.

For a valley whose $\mathbf{K}^{(i)}$ lies on a threefold or fourfold symmetry axis of a crystal the deformation potentials due to the different possible types of strains can be expressed in terms of two constants, $\Xi_d$, $\Xi_u$, which appear in the theory of piezoresistance, Eq. (C6) below, and which are defined thus: Let $u_r(r = 1$ to $6)$ be the six components of the strain tensor, relative to the principal axes of the valley, the $z$ axis being taken along the symmetry axis of the valley. Then a dilatation in the two directions normal to the symmetry axis $(u_1 = u_2 = u/2, u_3 = 0)$ produces a band edge shift $\Xi_d u$. A uniaxial shear $(u_1 = u_2 = -u/2, u_3 = u)$ produces a shift $\Xi_u u$. As will be shown in detail in Reference 6, it is not hard to evaluate the dependence of $C(\mathbf{K}, \mathbf{K}')$ on the direction of $\mathbf{K}' - \mathbf{K}$, in terms of $\Xi_d$ and $\Xi_u$ and the elastic constants of the crystal. Thus we can calculate the variation of $S(\mathbf{K}, \mathbf{K}')$ over an energy surface by (A7), in terms of the constants $\Xi_d$, $\Xi_u$, and the anisotropy of the effective mass.

We shall presently exhibit the results of some calculations, made by the procedure just outlined, of how much the rate of loss of initial velocity varies over an energy surface. However, we shall first present an argument that for given deformation potentials, the greatest variation of this quantity may be expected to occur for spherical energy surfaces,

---

[21] See, for example, Reference 19, p. 520 et seq.

and that usually very little variation will occur for extremely prolate or oblate surfaces. In other words, for almost all values of $\Xi_d$ and $\Xi_u$ , the use of a $\tau(\epsilon)$ will be justified for very prolate or very oblate surfaces. Fig. 11 shows the argument. In (a) is shown a spherical energy surface centered on a band edge point on some symmetry axis in K-space. We expect the greatest difference in relaxation time $\tau(\mathbf{K})$, defined in terms of rate of loss of initial velocity, to be that between a point $\mathbf{K}_1$ where the symmetry axis cuts the sphere and a point $\mathbf{K}_2$ 90° around the sphere from $\mathbf{K}_1$ . Backward scattering processes for $\mathbf{K}_1$ , i.e., those which almost reverse its velocity, have $\mathbf{K}' - \mathbf{K}_1$ vectors almost parallel to the symmetry axis, while forward scattering processes, which take $\mathbf{K}_1$ to a state $\mathbf{K}'$ with almost the same velocity, have $\mathbf{K}' - \mathbf{K}_1$ almost at right angles to the axis. For a carrier initially at $\mathbf{K}_2$ , on the other hand, the vector $\mathbf{K}' - \mathbf{K}_2$ is almost perpendicular to the axis for backward scattering, while for forward scattering it may range from almost parallel (plane of paper) to perpendicular (normal to paper). For this case, therefore, a dependence of the scattering function on the inclination to the symmetry axis will be quite effective in producing an anisotropy of $\tau(\mathbf{K})$. Now consider the situation for a very prolate energy surface, as shown in (b) of Fig. 11. This figure has been derived from that of (a) by a horizontal extension and a vertical contraction. Now all the dotted lines representing vectors $\mathbf{K}' - \mathbf{K}_1$ or $\mathbf{K}' - \mathbf{K}_2$ are nearly parallel to the symmetry axis, and the
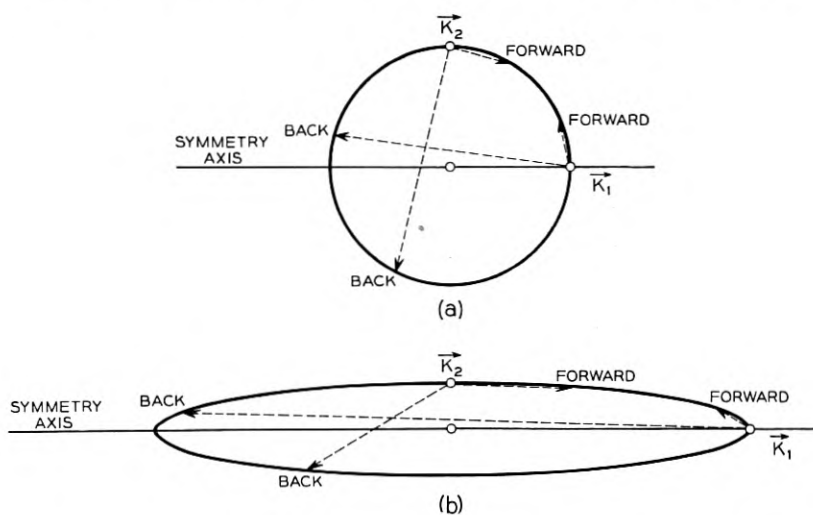


Fig. 11 — Comparison of intra-valley lattice scattering processes for a valley with spherical energy surfaces, (a), and a valley with highly prolate surfaces, (b).

corresponding scattering functions must therefore be very nearly the same. Of course, forward scattering from $\mathbf{K}_1$ through sufficiently small angles has $\mathbf{K}' - \mathbf{K}_1$ nearly normal to the axis, and the same is true for backward scattering from $\mathbf{K}_2$ in a small range close to 180°. But it is clear that as the energy surface becomes more prolate these cases form a smaller and smaller fraction of the total of possible final states $\mathbf{K}'$. Thus for extreme prolateness $\tau(\mathbf{K}_1) \rightarrow \tau(\mathbf{K}_2)$. For strongly oblate energy surfaces most of the $\mathbf{K}' - \mathbf{K}$ vectors approach normality to the symmetry axis, and a similar conclusion holds. The argument fails if, and only if, the scattering probability practically vanishes for $\mathbf{K}' - \mathbf{K}$ along the symmetry axis (prolate case) or perpendicular to it (oblate case).

Fig. 12 shows the results of some calculations of $\tau(\mathbf{K}_2)/\tau(\mathbf{K}_1)$ carried out by $E$. Vogt for the worst case, that of spherical energy surfaces, and for an actual case, that of valleys on a (111) axis with $m_\parallel{}^*/m_\perp{}^* = 1.3/.08$, the value found in cyclotron resonance experiments on $n$ germanium.[5] The anisotropy of the elastic constants has been assumed to be that for germanium. For the spherical surfaces calculations were made for valleys centered on (100) and (111) axes, but the results were found to be indistinguishable. A comparison of the full curve (spherical surfaces) with the dashed one (prolate surfaces) shows, as expected, that with a highly anisotropic effective mass the anisotropy of the relaxation time is much less than for the spherical case, except near the ratio $\Xi_d/\Xi_u = -1$. This is the ratio for which modes with wave vectors along the symmetry axis are incapable of scattering. We conclude that for intra-valley lattice scattering the assumption of a relaxation time dependent only on energy will fail over a considerable range of the deformation potential parameters if the effective mass is isotropic, but only over a moderate range near $\Xi_d/\Xi_u = -1$ if the effective mass is very anisotropic.

(iv) *Scattering by Ionized Impurities.* To date the quantum theory of this effect has been developed only on the crude basis of treating the fluctuations of potential due to a random arrangement of ions as a small perturbation on the motion of the charge carriers.[22] The result is that the effective matrix element for scattering between two states $\mathbf{K}, \mathbf{K}'$ on the same energy surface is a function of $| \mathbf{K}' - \mathbf{K} |$ which is sharply peaked at very small values of this quantity, at least when the density of impurities is not too high. This means that the principal effects come from small-angle scattering, a fact well known in the classical theory of impurity

[22] H. Brooks, Phys. Rev., **83**, p. 8 & 9, 1951 (contains typographical errors); C. Horie, Tôh. U. Sci. Rep. **34**, 27 (1950); and C. Herring, unpublished.
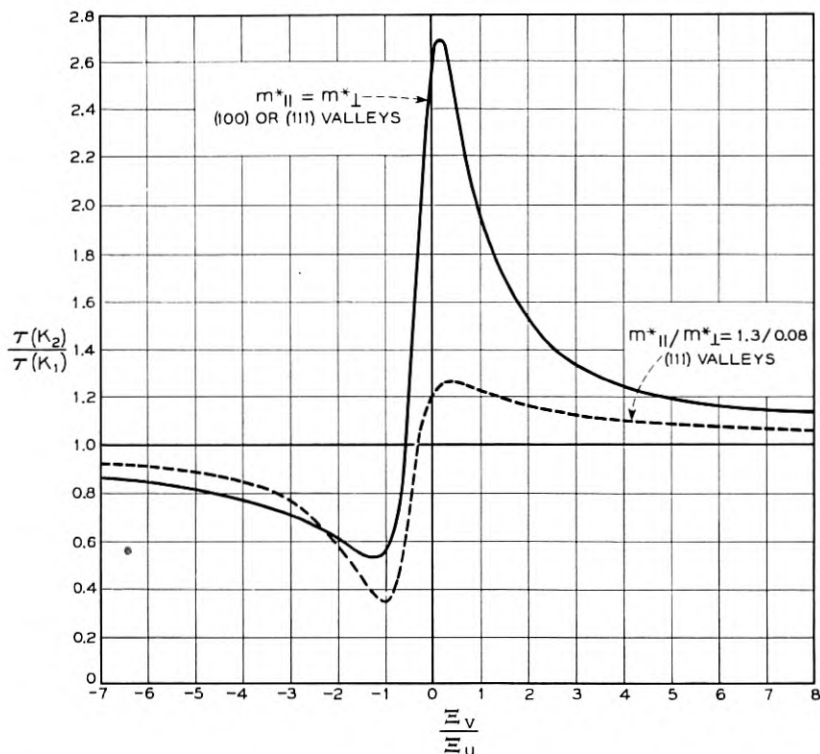
Fig. 12 — Anisotropy of the relaxation time for intra-valley lattice scattering, as a function of the ratio of the deformation potential coefficient $\Xi_d$ for two-dimensional dilatation to the coefficient $\Xi_u$ for uniaxial shear, as these quantities are defined in Appendix A or in Equation (C6). The ordinate is the ratio, for the points $K_2$ and $K_1$ of Fig. 11, of the effective relaxation time $\tau$ defined by $\tau^{-1} = -$ (initial velocity) $^{-1}$ (time rate of change of mean velocity, due to scattering). The elastic anisotropy assumed is that for germanium.

scattering.[23] Now the probability of a collision with a given range of $|\mathbf{K}' - \mathbf{K}|$ is proportional to the square of this matrix element and to the number of possible final states $\mathbf{K}'$ within the given range of distances from $\mathbf{K}$ and in a given small range of energy. A little calculation shows that this number is greater when $\mathbf{K}$ is near the $\mathbf{K}_1$ of Fig. 11(b) than when $\mathbf{K}$ is near $\mathbf{K}_2$. Moreover, the fractional loss of the velocity component in the initial direction is greater for a collision at $\mathbf{K}_1$ than at $\mathbf{K}_2$. The result

[23] E. Conwell and V. F. Weisskopf, Phys. Rev., **77**, p. 388, 1950; R. S. Cohen, L. Spitzer, Jr., and P. M. Routly, Phys. Rev., **80**, p. 230, 1950; and L. Spitzer, Jr. and R. Härm, Phys. Rev., **89,** p. 977, 1953.

is that carriers near $\mathbf{K}_1$ have a much shorter effective relaxation time for impurity scattering than do those near $\mathbf{K}_2$.

Thus if the energy surfaces are very anisotropic, the assumption $\tau = \tau(\epsilon)$ may be expected to be a poor approximation when ionized impurity scattering is important.

(v) *Neutral Impurity Scattering.* In the simple theory the scattering of charge carriers from neutral impurities in hydrogen-like states is mathematically equivalent to the scattering of electrons from hydrogen atoms.[24] At the temperatures at which such scattering is important the wavelength of the incident carrier is usually $\gg$ the diameter of the wave function of the center, so $s$ wave scattering predominates. Therefore the scattering is isotropic. For a many-valley model, however, the situation seems at first sight more complicated, since the effective mass is anisotropic, and the centers are not spherically symmetrical. However, it can be shown[6] that, at least if the energy of the carrier is low enough, the scattered wave must be describable as an $s$ wave in the space of the transformed coordinates defined by

$$\xi_\lambda = m_\lambda^{* \ 1/2} x_\lambda \tag{A8}$$

In the corresponding momentum space (the $\varphi$-space of Appendix B) the surfaces of constant energy are spheres. It follows, that, at least at low energies, neutral impurity scattering satisfies (A3) and is describable by a constant relaxation time.

## Appendix B

### EQUALITY OF THE ENERGY-SHELL AVERAGES

$$\overline{\tfrac{1}{2}m_1^* v_1^{\ 2}}, \qquad \overline{\tfrac{1}{2}m_2^* v_2^{\ 2}}, \qquad \overline{\tfrac{1}{2}m_3^* v_3^{\ 2}}$$

Choose coordinate axes along the principal axes of the energy surfaces in any given valley. For $\lambda = 1, 2, 3$, let

$$\varphi_\lambda = v_\lambda (m_\lambda^*)^{1/2} = \Delta P_\lambda / (m_\lambda^*)^{1/2} \tag{B1}$$

Then

$$\Delta \epsilon \equiv |\ \epsilon(\Delta \mathbf{P}) - \epsilon_b\ | = \tfrac{1}{2} \sum \varphi_\lambda^{\ 2} \tag{B2}$$

so that the energy surfaces are concentric spheres in $\varphi$-space. The density of $\varphi$-vectors consistent with the periodic boundary conditions is of course uniform, like that of $\Delta \mathbf{P}$. The average of $\tfrac{1}{2}m_\lambda^* v_\lambda^{\ 2}$ over an energy shell is therefore the average of $\varphi_\lambda^{\ 2}$ over a spherical shell, and thus obviously independent of $\lambda$.

---

[24] C. Erginsoy, Phys. Rev., **79**, p. 1013, 1950.

## Appendix C

### EXPLICIT CALCULATION OF ELASTORESISTANCE CONSTANTS

Let $\sigma_{\mu\nu}$ be the conductivity tensor of a crystal, defined by the relation

$$j_\mu = \sum_\nu \sigma_{\mu\nu} E_\nu \tag{C1}$$

between current density $\mathbf{j}$ and electric field $\mathbf{E}$. Let $u_{\alpha\beta}$ be the strain tensor, defined by the relation

$$\delta x_\alpha = \sum_\beta u_{\alpha\beta} x_\beta \tag{C2}$$

between displacement $\delta\mathbf{x}$ and initial position $\mathbf{x}$ of any point in the body. Then the elastoresistance of any crystal is described mathematically by the fourth-rank tensor $\partial\sigma_{\mu\nu}/\partial u_{\alpha\beta}$, or more conveniently, for a cubic crystal, by the dimensionless tensor $-\sigma^{-1}\partial\sigma_{\mu\nu}/\partial u_{\alpha\beta}$, where $\sigma$ is the scalar conductivity in the unstrained state ($\sigma_{\mu\nu} = \sigma\delta_{\mu\nu}$). The task of this appendix is to obtain this tensor by calculating the strain variation of the current contributions (7), which in terms of components take the form

$$j_\mu^{(i)} = \frac{e^2}{kT} \sum_{\Delta\mathbf{P}^{(i)},s} f^{(0)}\tau(\Delta\epsilon^{(i)}) \sum_\nu v_\mu v_\nu E_\nu \tag{C3}$$

We shall base the calculation on assumptions (a) through (e) of Section 5. Assumptions (b) and (d), regarding existence of a $\tau(\Delta\epsilon)$ and nondegenerate statistics, are already contained in (C3). According to assumption (a), we shall neglect any effect of strain on the relations between $\mathbf{v}$, $\Delta\mathbf{P}$, and $\Delta\epsilon^{(i)} = |\epsilon - \epsilon^{(i)}|$. Thus we may replace $v_\mu v_\nu$ in (C3) by its average $L_{\mu\nu}^{(i)}$ over an energy shell, and treat this average for given $\Delta\epsilon^{(i)}$ as uninfluenced by strain:

$$v_\mu v_\nu \longrightarrow \tfrac{2}{3}\Delta\epsilon^{(i)} \begin{bmatrix} \dfrac{1}{m_1{}^*} & 0 & 0 \\[2mm] 0 & \dfrac{1}{m_2{}^*} & 0 \\[2mm] 0 & 0 & \dfrac{1}{m_3{}^*} \end{bmatrix} \equiv L_{\mu\nu}^{(i)} \tag{C4}$$

When the crystal is strained, the only things that we assume to change in (C3) are the population factor $f^{(0)}$ and the relaxation time $\tau(\Delta\epsilon^{(i)})$. For some groups of electrons — i.e., values of $i$ and $\epsilon$ — the product $\tau f^{(0)}$ will be increased, and for others it will be decreased. If on the average this product is increased in the valleys whose conductivity tensors are most favorably oriented to the direction of $\mathbf{E}$, and decreased in the

others, the total current will be increased by the strain, and vice versa. When assumption (e) of Section 5 is fulfilled, so that two of the $m^*$'s, say $m_1^*$ and $m_2^*$, are equal, the tensor (C4) takes the form

$$L_{\mu\nu}{}^{(i)} = \tfrac{2}{3}\Delta\epsilon^{(i)} \left[ \frac{\delta_{\mu\nu}}{m_1^*} + \frac{K_\mu{}^{(i)}K_\nu{}^{(i)}}{K^{(i)2}} \left( \frac{1}{m_3^*} - \frac{1}{m_1^*} \right) \right] \tag{C5}$$

as is easily verified by inspection. Moreover, it is easily seen from symmetry that the only strain components which can alter $\epsilon^{(i)}$ in first order are the isotropic dilatation and the shear compounded out of an extension along $\mathbf{K}^{(i)}$ and a contraction in both directions at right angles. Mathematically expressed, we must have

$$\frac{\partial\epsilon^{(i)}}{\partial u_{\alpha\beta}} = \Xi_d\delta_{\alpha\beta} + \Xi_u K_\alpha{}^{(i)}K_\beta{}^{(i)}/K^{(i)} \tag{C6}$$

where $\Xi_d$ and $\Xi_u$ are constants independent of the valley $i$, the subscripts referring to "dilatational" and "uniaxial" effects respectively. The elastoresistance tensor is, for nondegenerate concentrations,

$$m_{\mu\nu\alpha\beta} \equiv -\frac{1}{\sigma}\frac{\partial\sigma_{\mu\nu}}{\partial u_{\alpha\beta}} = \frac{e^2}{kT\sigma}\sum_i \sum_{\Delta\mathbf{P}^{(i)},s}$$
$$\left[ \tau^{(i)}\frac{f^{(0)}}{kT}\frac{\partial\,|\,\epsilon^{(i)} - \epsilon_F\,|}{\partial u_{\alpha\beta}} - \sum_j f^{(0)}\frac{\partial\tau^{(i)}}{\partial\epsilon^{(j)}}\frac{\partial\epsilon^{(j)}}{\partial u_{\alpha\beta}} \right] L_{\mu\nu}{}^{(i)} \tag{C7}$$

where $\epsilon_F$ is the Fermi level, $\tau^{(i)}(\epsilon)$ is the relaxation time in the $i$th valley, and $L_{\mu\nu}{}^{(i)}$ is given by (C5). The second term in brackets in (C7) represents the effect of strain on the transition probability for inter-valley scattering.

We shall now combine and simplify the equations just given. The behavior of the Fermi level is simple: by assumption (a) of Section 4 it does not shift in shear, and for extrinsic concentrations it shifts with the band edge in compression. Mathematically,

$$\frac{\partial\epsilon_F}{\partial u_{\alpha\beta}} = \frac{1}{N_V}\sum_j \frac{\partial\epsilon^{(j)}}{\partial u_{\alpha\beta}} = \left( \Xi_d + \frac{\Xi_u}{3} \right)\delta_{\alpha\beta} \tag{C8}$$

where $N_V$ is the number of valleys. By virtue of the fact that $\tau_i$ is a function only of $\Delta\epsilon$ and the differences $(\epsilon^{(i)} - \epsilon^{(j)})$,

$$\sum_j \frac{\partial\tau^{(i)}}{\partial\epsilon^{(j)}} = 0 \tag{C9}$$

and it is easily seen that when (C6) and (C8) are inserted into (C7) the $\Xi_d$ terms disappear. Since by symmetry $\partial\tau^{(j)}/\partial\epsilon^{(i)} = \partial_\tau(i)/\partial\epsilon^{(j)}$, (C9)

implies

$$\sum_i \frac{\partial \tau^{(i)}}{\partial \epsilon^{(j)}} = 0 \tag{C10}$$

Using this and (C6) and (C8) in (C7), and processing further by (C5), we get

$$
\begin{aligned}
m_{\mu\nu\alpha\beta} &= \frac{e^2}{kT\sigma} \Xi_u \sum_i \sum_{\Delta \mathbf{P}^{(i)}, s} \left[ \pm \frac{\tau^{(i)}}{kT} \left( \frac{K_\alpha^{(i)} K_\beta^{(i)}}{K^{(i)2}} - \frac{\delta_{\alpha\beta}}{3} \right) \right. \\
&\qquad \left. - \sum_j \frac{\partial \tau^{(i)}}{\partial \epsilon^{(j)}} \frac{K_\alpha^{(j)} K_\beta^{(j)}}{K^{(j)2}} \right] f^{(0)} L_{\mu\nu}^{(i)} \\
&= \frac{2ne^2 \Xi_u}{3kT\sigma} \left[ \pm \frac{\langle \Delta\epsilon\tau \rangle}{kT} \left( \left\langle \frac{K_\mu^{(i)} K_\nu^{(i)}}{K^{(i)2}} \cdot \frac{K_\alpha^{(i)} K_\beta^{(i)}}{K^{(i)2}} \right\rangle_i - \frac{\delta_{\mu\nu}\delta_{\alpha\beta}}{9} \right) \right. \\
&\qquad \left. - N_V \left\langle \langle \Delta\epsilon \frac{\partial \tau^{(i)}}{\partial \epsilon^{(j)}} \rangle \cdot \frac{K_\mu^{(i)} K_\nu^{(i)}}{K^{(i)2}} \cdot \frac{K_\alpha^{(j)} K_\beta^{(j)}}{K^{(j)2}} \right\rangle_{i,j} \right] \left( \frac{1}{m_3^*} - \frac{1}{m_1^*} \right)
\end{aligned}
\tag{C11}
$$

where the upper sign is for n-type, the lower for p, $n$ is the total carrier concentration, angular brackets with a subscript $i$ or $i, j$ mean averages on valleys $i$ or $i$ and $j$, and angular brackets without subscripts mean Maxwellian averages as defined in Section 2. Substituting from (14) for the conductivity $\sigma = ne\mu$, we get finally

$$
\begin{aligned}
m_{\mu\nu\alpha\beta} &= 3\Xi_u \left[ \pm \frac{1}{kT} \left( \left\langle \frac{K_\mu^{(i)} K_\nu^{(i)}}{K^{(i)2}} \cdot \frac{K_\alpha^{(i)} K_\beta^{(i)}}{K^{(i)2}} \right\rangle_i - \frac{\delta_{\mu\nu}\delta_{\alpha\beta}}{9} \right) \right. \\
&\quad \left. - N_V \left\langle \frac{\langle \Delta\epsilon \partial\tau^{(i)}/\partial\epsilon^{(j)} \rangle}{\langle \Delta\epsilon\tau \rangle} \frac{K_\mu^{(i)} K_\nu^{(i)}}{K^{(i)2}} \cdot \frac{K_\alpha^{(j)} K_\beta^{(j)}}{K^{(j)2}} \right\rangle_{i,j} \right] \frac{(m_1^* - m_3^*)}{(m_1^* + 2m_3^*)}
\end{aligned}
\tag{C12}
$$

As in (C7), the first term in square brackets in (C12) represents the effect of the strain on the relative populations of the different valleys, and the second term represents the effect on the inter-valley scattering probabilities. We may note the following features:

(1) The trace $\sum_\alpha m_{\mu\nu\alpha\alpha}$ vanishes identically, by virtue of (C9) etc. This means that an isotropic dilatation produces no elastoresistance under the assumptions we are using.

(2) The elastoresistance is proportional to the anisotropy of the effective mass within a valley.

(3) $m_{1212} = m_{44}$ vanishes for valleys of the (100) type, while $m_{1111} = m_{1122}$ ($m_{11} = m_{12}$) for valleys of the (111) type.

(4) The elastoresistance is proportional to $1/T$ at temperatures low enough for inter-valley scattering to be frozen out. Moreover, when $kT \gg$ the inter-valley $\hbar\omega(ij)$ of Section 3, the variation of $\partial\tau_i/\partial\epsilon_j$ with energy

$\Delta\epsilon$ relative to the band edge is easily shown to be one of proportionality to $\tau/\Delta\epsilon$, if impurity scattering is negligible. Under these conditions (probably never achieved by Si and Ge in the extrinsic range) the elastoresistance is again proportional to $1/T$, but with a larger factor of proportionality than at low $T$.

NOTATIONS

| | |
|---|---|
| $A$ | dimensionless average involving relaxation time, Equation (68) |
| $B$ | dimensionless function of mass anisotropy, Equation (39) |
| $C(\mathbf{K}, \mathbf{K}')$ | factor in matrix element for a scattering process $\mathbf{K} \to \mathbf{K}'$, Equation (A6) |
| $c$ | velocity of light; magnetoresistance constant, Equation (74) |
| $D_{ij}$ | factor in matrix element for scattering from valley $i$ to valley $j$, Equation (16) |
| $\mathbf{E}$ | electric field |
| $e$ | electronic charge |
| $F_{\mu\nu\alpha\beta}$ | factor determining the anisotropy of the magnetoconductance tensor Equations (70), (71) |
| $F_{\mu\nu\alpha\beta}^{(i)}$ | contribution of the $i$th valley to above, Equation (69) |
| $f$ | distribution function for charge carriers. |
| $f^{(0)}$ | same in absence of perturbing fields, Equation (1) |
| $\mathbf{f}^{(1)}$ etc. | change of $f$ in perturbing fields, Equations (4), (31), (60) |
| $G$ | dimensionless average involving relaxation time, Equation (57) |
| $\mathbf{H}$ | magnetic field |
| $\hbar$ | Planck's constant$/2\pi$ |
| $\mathbf{j}$ | density of electric current |
| $\mathbf{K}$ | wave number vector for a charge carrier |
| $\mathbf{K}^{(i)}$ | value of $\mathbf{K}$ for the $i$th band edge point (center of the $i$th valley) |
| $k$ | Boltzmann's constant |
| $L_{\mu\nu}^{(i)}$ | average of $v_\mu v_\nu$ over an energy shell in the $i$th valley |

| $M(\mathbf{K}, \mathbf{K}'), M_{ij}$ | matrix element for lattice scattering, Equation (A6) or (16) |
|---|---|
| $m$ | normal electron mass |
| $m_\lambda{}^*$ | effective mass in $\lambda$th principal direction of a valley |
| $m^{(I)}$ | inertial average of the $m_\lambda{}^*$, defined by Equation (15) |
| $m_{\mu\nu\alpha\beta}$ | elastoresistance tensor $-\sigma^{-1}\partial\sigma_{\mu\nu}/\partial u_{\alpha\beta}$ |
| $N, N_\alpha$ | number of quanta (phonons) in a given lattice mode |
| $N_V$ | number of valleys or band edge points |
| $n$ | number of free charge carriers per unit volume |
| $\mathbf{P}$ | crystal momentum $\hbar\mathbf{K}$ |
| $Q_e$ | electronic part of thermoelectric power |
| $\mathbf{q}$ | wave number vector for a lattice mode |
| $R$ | Hall constant |
| $S(\mathbf{K}, \mathbf{K}')$ | transition probability $\mathbf{K}$ to unit $d\mathbf{K}'$ at $\mathbf{K}'$, Equation (A1) |
| $s$ | spin quantum number of a charge carrier |
| $T$ | absolute temperature |
| $t$ | time |
| $u_{\alpha\beta}, u_r$ | strain tensor components |
| $\mathbf{v}$ | group velocity of a charge carrier |
| $W_a(ij, \alpha)$ | transition probability from valley $i$ to valley $j$ with absorption of a phonon of branch $\alpha$ |
| $W_e(ij, \alpha)$ | same but with emission |
| $w_1, w_2$ | constants of dimension $(\text{time})^{-1}$, measuring coupling of carriers to intra- and intervalley modes, respectively, Equations (21)–(23) |
| $x, y, z$ | rectangular coordinates |
| $\alpha$ (when not subscript) | index labeling branches of vibrational spectrum |
| $\Upsilon$ | differential operator describing rotation of distribution by magnetic field, Equation (62) |
| $\Delta\epsilon$ | energy of a carrier relative to the band edge ($\Delta\epsilon \geqq 0$) |
| $\Delta\mathbf{P}, \Delta\mathbf{K}$ | departure of crystal momentum or wave vector from valley center |
| $\delta_{\mu\nu}$ | Kronecker delta |

| | |
|---|---|
| $\delta_{\mu\nu\sigma}$ | antisymmetric coefficient $= \pm 1$ or $0$, Equation (33) |
| $\delta\epsilon$, $\delta\mathbf{v}$, etc. | changes induced by strain |
| $\epsilon$ | energy of an electron |
| $\epsilon_b$ | band edge energy |
| $\epsilon^{(i)}$ | energy of the center of the $i$th valley (normally $= \epsilon_b$) |
| $\theta$ | Debye temperature |
| $\kappa$ | dielectric constant |
| $\mu$ | drift mobility |
| $\mu_H$ | Hall mobility |
| $\Xi_d$, $\Xi_u$ | deformation potentials, Equation (C6) |
| $\xi$ | transformed coordinate, equation (A15) |
| $\rho$ | resistivity |
| $\sigma$ | conductivity |
| $\sigma_0$ | conductivity in the absence of magnetic fields or other perturbations |
| $\sigma_{\mu\nu\alpha}$, $\sigma_{\mu\nu\alpha\beta}$ | coefficients in the expansion of conductivity in powers of $H$, Equation (65) |
| $\sigma_{\mu\nu\alpha\beta}^{(i)}$ | contribution of $i$th valley to magnetoconductivity |
| $\tau$ | relaxation time of charge carriers |
| $\tau^{(i)}$ | ditto for carriers in valley $i$ |
| $\varphi$ | transformed wave vector, Equation (B1) |
| $\omega$ | angular frequency of a lattice mode or an rf field |

# New Manufacturing Techniques for Precision Transformers for the L3 Coaxial Carrier System

By NORMAN E. EARLE

*The critical function performed by the 2504A transformer in the L3 coaxial carrier telephone system necessitated the use of completely new materials and of methods radically different from those usually associated with carrier telephone transformer production. To satisfy operational requirements, extensive use has been made of parts machined from ceramic and glass insulating materials which can now be machined to very close dimensional tolerances. A description is given of the equipment and techniques which were developed to produce these transformers on a commercial basis, and their effectiveness is evaluated.*

## INTRODUCTION

Early in 1951 the Western Electric Company set up manufacturing facilities for the production of the L3 coaxial system.* This new system transmits the frequency band from 0.3 to 8.353 mc over coaxial cables. To counteract the attenuation of these cables, line amplifiers are provided every four miles. The 2504A transformer is used both at the input and output of these amplifiers to couple them to the cable. The system is designed so that the parasitic elements of the transformer are used to help shape the transmission characteristic. Since upwards of 2,000 of these transformers are used in a long cable system, it is necessary that the constants of each transformer be held to an extraordinary degree of precision, and that the temperature variation and aging effects be virtually eliminated.

---

* L. H. Morris, G. H. Lovell, and F. R. Dickinson, The L3 Coaxial System, B.S.T.J., **32**, July, 1953. A technical description of the 2504A transformer appears on page 891.

To achieve the desired precision, a radically new type of construction was used for the 2504A transformer shown in Fig. 1. In place of the conventional windings, glass cylinders are used, having accurately machined grooves in which the conductor is embedded by a combination of firing and plating processes. Steatite forms are used to provide the housing. In turn, the manufacture of this transformer required the use of new tools and the development of new processes and techniques. The problem the



FIG. 1A — Exploded schematic of the 2504A transformer.

Fig. 1B — 2504A transformer.

engineer of manufacture had to solve, therefore, was how to set up to produce 2,000 such transformers a year on a capacity basis.

Table I gives a comparison of the structural details of the 2504A as compared to those of a carrier transformer of conventional design.

The new materials, in addition to having stability, low loss, and good insulating properties, can be held to very close mechanical tolerances, a condition which makes it possible to obtain the degree of electrical

TABLE I — COMPARISON OF STRUCTURAL DETAILS OF 2504A WITH THOSE OF A CONVENTIONAL TRANSFORMER

| Part | Conventional Transformer | 2504A Transformer |
|---|---|---|
| Core............ | Permalloy powder | Manganese-zinc ferrite |
| Winding form.... | Molded or fabricated plastic part | Threaded cylinders of special glass |
| Windings........ | Machine-applied insulated drawn copper wire | Plated copper embedded in special glass forms |
| Shielding........ | Copper or tin foil | Fired silver on special glass forms |
| Coil enclosure... | Metal or plastic container | Steatite housing and cradle |

precision required, once the techniques for performing the mechanical operations have been developed.

In the gage making art, equipment has been available for many years for grinding threads to tolerances comparable to the .0.031 ± 0.0005″ limits on the wall thickness between the inner diameter and the bottom of the threads of the glass form for the outer winding of the transformer and the 0.0005″ concentricity limits on the cylindrical inner and outer surfaces of the same part. However, considerable development was necessary before techniques could be worked out to hold such tolerances on fragile, thin-walled glass parts without excessive breakage.

Application of silver to the inner surfaces of the middle and outer forms, shown on Fig. 2, and the plated copper windings to the inner and outer forms presented similar problems. The difficulty was not the newness of the art, but the unusual surfaces on which the silver had to be applied, and the uniformity required in order to obtain the precise control of leakage fluxes and interwinding capacitances that was necessary to maintain the desired shape for the transformer's transmission characteristic.

Further problems were encountered by the manufacturing engineer in measuring some of the mechanical dimensions of the parts and some of the electrical characteristics of the completed transformer to the accuracy required and in arriving at a satisfactory correlation between changes in mechanical dimensions and the resultant changes in electrical properties that in the end were the characteristics that had to be held precisely.

MACHINING OF THE GLASS COIL FORMS AND THE CERAMIC CRADLE AND HOUSING

With the exception of a honing operation on the inside diameters of the coil forms and the final dressing operation after the parts are copper plated, all machining operations on the cradle and housing and on all three coil forms for the 2504A transformer are done on special grinding machines equipped with fine grit bonded-diamond wheels.

The configuration of the winding forms is shown diagrammatically in Fig. 2. These forms were originally made from fused quarts; but later, extensive tests showed that a cheaper 96 per cent silica glass, could be used interchangeably with the quartz. While the material is very hard, the parts themselves are thin and consequently fragile and easily distorted; so easily, in fact, that they can be distorted by more than the specified tolerance simply by holding them between one's fingers. This

factor, together with the close dimensional tolerances including the concentricity requirement of ±0.0005″ on the outer winding form, made it imperative that tubing appreciably thicker than the finished product be used as the raw material for making the part. This made possible accurate honing of the inside diameters, used as the base point for hold-
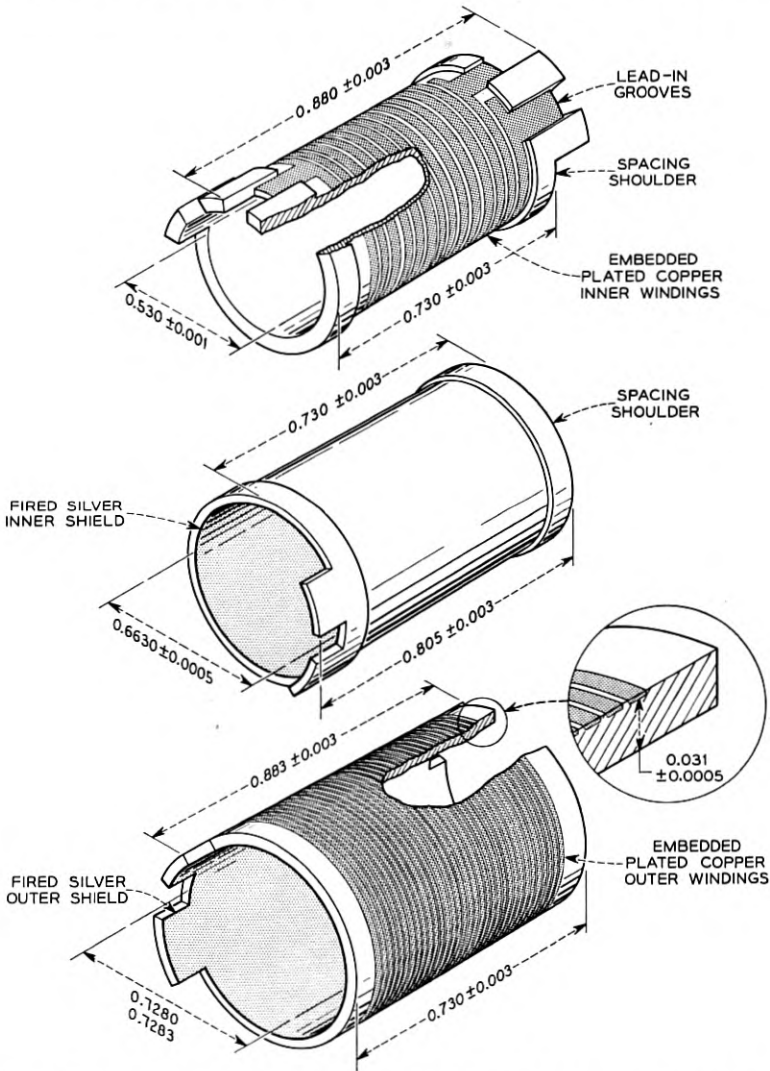


Fig. 2 — Composite view of winding forms, embedded copper windings and fired silver interwinding shields.

ing all other dimensions of the parts to the precise dimensions required. Though the specified tolerances on these dimensions were wider, in actual practice, it was found that holding the inside diameters to plus or minus 50 millionths of an inch made it much easier to meet the limits on other dimensions, particularly the highly important wall thickness dimension of the outer form.



Fig. 3 — Honing inside diameter of coil form.

In preparing the tubing for this operation, the material is pre-shrunk in a plastic state at elevated temperatures around an accurately ground mandrel with the result that the bore is reduced to within 0.001″ to 0.002″ of its final size. The keyway required in the inner form is formed by the same process. No machining of this slot is required.

Rough grinding of the outside diameters to ±0.001″ and facing the



Fig. 4 — Cutting slots in end of inner form.

ends to meet the ±0.003″ tolerance for the overall length of the parts presented no particular problems. However, the use of an Arnold gage to give a continuous indication of the outside diameter during grinding of the parts was novel to coil manufacturing.

Three special machines, developed early in 1951, make the notching, slotting and trimming, necessary in forming the lead in grooves and tabs of which they are a part, a relatively easy operation, despite the fact that the tolerance in their angular location is better than ±1°. These machines — one of which is shown pictorially in Fig. 4 — contributed materially to the uniformity of the product, a factor of utmost importance in controlling the internal parasitics of the transformer.

Undercutting of the inner and middle forms for the spacing shoulders on these parts is a straight forward cylindrical grinding operation. Tolerances on the width of the shoulders are ±0.003″, while those of the diameter of the parts are ±0.001″.

Grinding of the threads on the inner and outer forms — the final operation prior to application of the plated copper windings — as had been anticipated — turned out to be a major manufacturing problem. Production of an experimental lot of transformers on a specially equipped toolmakers' lathe indicated that the best way to obtain the uniform control desired was through the use of an automatically cycled commercial thread grinder. Other methods, including a novel approach in which the threads were formed by a high-frequency vibration process using a powdered abrasive cutting agent, were either too slow or inherently did not have sufficient controls to yield the required uniformity. Even with the automatically cycled commercial thread grinder, considerable development was required to work out the required techniques for the process and only one supplier could be found who would undertake to work out the proposed machine modifications involved.

One of the problems was to adapt the grinder to permit use of a metal bonded in place of the resinoid bonded diamond or silicon carbide wheel conventionally used for thread grinding. Factors leading to the selection of the metal bonded diamond wheel were (1) the necessity for holding as nearly square a thread as possible — in practice a 14½° Acme thread with 0.0015″ maximum radii in the corners — and (2) the extremely hard abrasive nature of the material being worked. To permit use of these wheels, specially designed silicon carbide scrubbers had to be provided, built in to automatically dress the wheel to the shape required after completion of a predetermined number of parts.

Selection of the proper grain size for this wheel was also critical. Wheels much coarser than 300 grit — a size so fine it could be used to produce a

high polish on metal parts — caused excessive chatter and breakage of parts when used to grind fine threads on glass.

The speed of approach and retraction of the wheel — controlled by a specially shaped forming cam — also had to be carefully worked out, not to minimize parts breakage alone, but to control thread over-
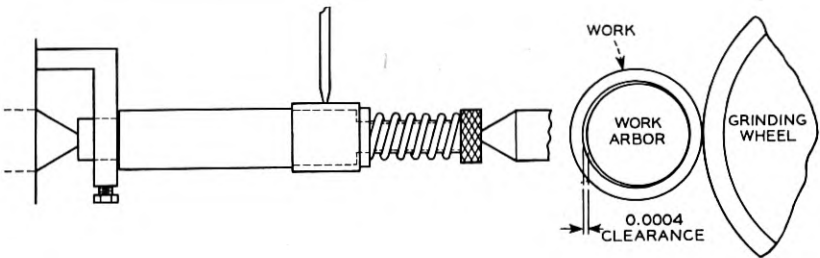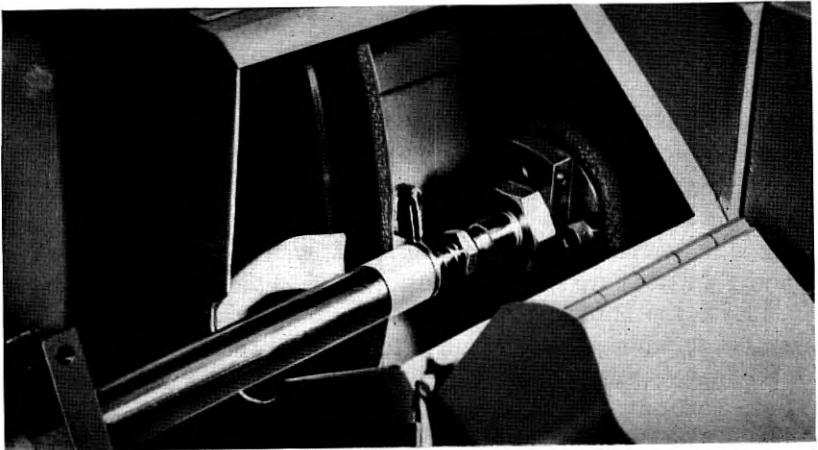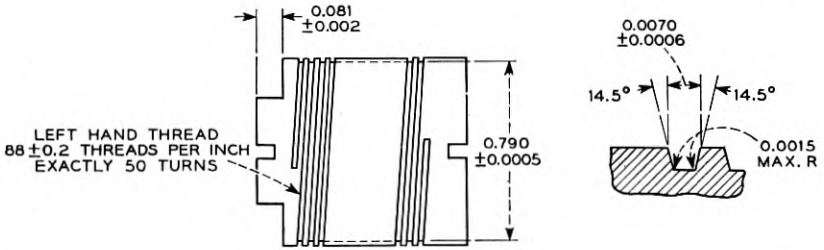


Fig. 6 — Cutting thread in outer winding form.

run, a factor affecting one of the transformer's parasitic interwinding capacitances.

Simultaneous removal of copper and glass in the final dressing of the parts after copper plating of the forms for the transformer windings required the selection of a wheel which could cut with very light pressure so as not to break the brittle glass form and at the same time not load up with copper from the windings. The most practical unit found for this operation was an "I" bond vitrified silicon carbide 220 grit wheel. Wheels with a stronger bond tended to cause excessive parts breakage, while softer ones broke down so rapidly that the required tolerances could not be held.

To provide as desirable operating conditions as possible, all of the grinding machines with the exception of one used for rough work were located in a dust free air-conditioned room in the basement of the building. In the case of the thread grinder, it was found necessary to take the further precaution of warming up the machine for a minimum of one and one-half hours prior to using it and providing more rigid guards in place of standard equipment in order to reduce vibration.

Machining of the steatite cradle and housing for the transformer presented problems of a similar nature. However, tolerances were not quite as critical and more information was available on commercial machining practices. As a result, not nearly as much development was required as on the glass winding forms. Procurement of unground parts sufficiently non-porous to avoid spalling and breakage of the housings, during firing after application of the sprayed silver shield on its inner surface, has been troublesome.

SILVERING AND PLATING OF THE GLASS AND CERAMIC PARTS

The inner surfaces of the ceramic coil housing and those of the middle and outer glass winding forms must be silvered for shielding purposes. The outer surfaces of the inner and outer coil forms are similarly processed to provide a base for the plated copper which, in combination with the silver, forms the actual coil windings.

It was imperative that these coatings be uniform in density and thickness particularly those on the inner surfaces of the coil forms. Thin non-uniform coatings result in intolerable variations in leakage inductance and interwinding capacitance, while excessively thick ones cause assembly difficulties. Other prerequisites were:

1. A strong dependable bond to the glass and ceramic surfaces.
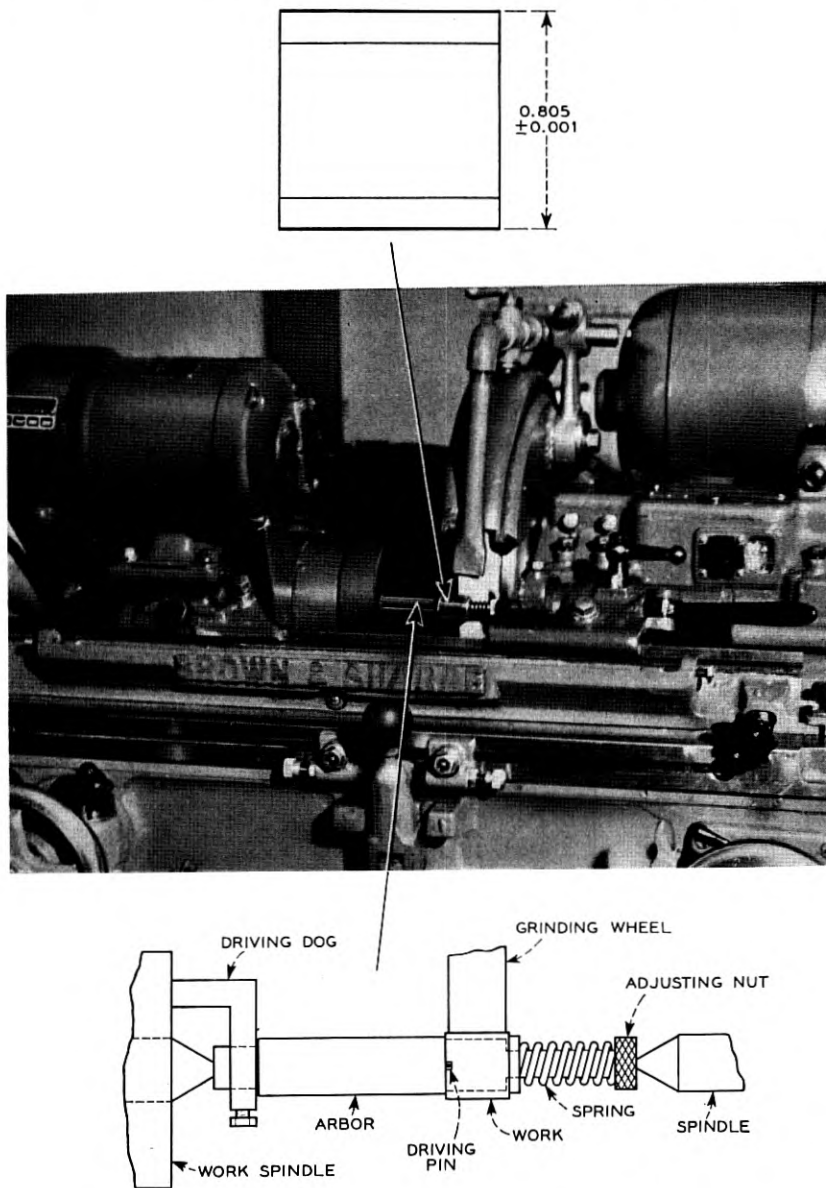2. Careful masking and process controls to insure that conductive

Fig. 7 — Grinding outside diameter of coil form.

material did not splatter onto critical areas where it might cause corona discharges while in service.

The silvering agent used was a finely ground silver powder suspended in a volatile organic solvent with small amounts of resin and glass frit. Sprayed on under controlled conditions in specially developed machines, the organic material is driven from the mixture by firing at 1250°F in a continuous belt furnace. At this temperature, the glass frit melts and wets the surfaces, bonding the silver securely to the glass and ceramic parts. Thickness of the coating is controlled by weighing control parts on a chain-o-matic laboratory balance before and after spraying and again after firing. Uniformity is improved by the application of multiple coats. However, not more than two can be used because of devitrification which results if quartz or high silica content glass parts are repeatedly fired at high temperatures.

In the early stages of production, the shield consisted of copper plating over fired silver coatings. This method required extra plating processes and prolonged the exposure of the parts to the plating solution. Since such solutions can attack the silver glass bond of the parts, plating of the shields was abandoned. To improve the continuity and conductivity of the coating, the parts are burnished in a specially developed machine which in a carefully timed cycle draws the inner surfaces of the parts spirally across the periphery of a rapidly rotating wire brush. Very good results have been obtained with the present process, the leakage impedance variation being reduced to plus or minus one-half of one per cent.

The threads which form the windings of the transformers are produced by plating copper on top of the silver fired on the outer surfaces of the inner and outer forms. The problem confronting the manufacturing engineer in this operation was that of applying as homogenous as possible a trapezoidal ribbon of copper in the narrow grooves on the outer surfaces without loosening the silver bonded to these forms.

This meant that the parts must not be permitted to remain in the plating bath any longer than absolutely necessary and that a well controlled process must be used to avoid bridging of the copper at the top before the groove is completely filled in the middle. To solve this problem, a plating process was used where initially the parts are given a flash of copper in a copper sulphate bath to protect the fired silver from attack by the copper fluoborate bath, subsequently used to permit rapid application of the heavy plating required. This flash effectively protects the silver from attack by the latter solution. To keep a fresh solution in the threads and to provide a uniform coating, the parts are rotated briskly during the plating operation. Because the quantity of parts being proc-

essed has not been sufficiently great to justify provision of a plant to provide continuously fresh plating solution, constant vigilance was required to keep impurities out of the solution. For this reason the anodes were enclosed in closely woven plastic bags and the solution was purified daily by filtering after treatment with activated carbon.

Microscopic analysis of etched cross-sections of the threads and dc resistance measurements on the overall windings indicated that exceptionally good results were being obtained by this process. Control is maintained by a continuous check of the winding for dc resistance.

PHYSICAL AND ELECTRICAL MEASUREMENTS

The performance of the transformer in its final state depends to a great extent upon the accuracy to which dimensions of the winding forms have been held. In fact, investigation has shown that half of the total variations in the overall transmission characteristic of the transformer can be attributed to variations in the dimensions of the outer form. Dimensional measurements on these parts thus becomes an extremely important consideration. Critical dimensions are the inside diameters of the three forms, wall thickness of the outer form at the root of the thread, and outside diameter of the inner form. Standard air, dial indicator and electronic gages provided have proven adequate for measuring these dimensions to the close tolerances required. Measurement of the wall thickness of the outer forms at the root diameter, however, has been very troublesome. The problem here is complicated by the fact that the measurement has to be made from a cylindrical surface in one plane to a point at the bottom of a narrow thread in a plane at the helix angle of the thread. At this point, the Acme thread is approximately 0.004″ wide including the 0.0015″ radii in each corner. The use of an optical comparator even at 100× magnification proved impractical because of the curved surfaces involved. After trial of this and other methods, some of them quite elaborate, it was found that the simple expedient of inserting a thin diamond blade in the thread and comparing the distance from the bottom of the blade to a circular rod supporting the cylindrical form with a flat ground carbide thickness gage gave the most consistent results. A diamond blade was used because it was found that the thin edge of the blade wore rapidly even with carboloy. A diamond edge lasts about six months. Results indicate that an accuracy of ±0.0002″ is attainable. While this accuracy is consistent with the ±0.0005″ limit on this dimension, greater accuracy would be useful in setting control limits. The strain gage used as the measuring instrument can be read to much greater

accuracy. However, to date, no means has been found to obtain an improved method for contacting the surfaces to be measured.

Other mechanical measurements, while made to a degree of precision unusual in the coil manufacturing field, are of more or less routine nature and no particular difficulty was experienced in providing equipment capable of making measurements to the accuracy required.

The transformer was designed to require no adjustment to meet its specified electrical requirements, even though the requirements for the transformers' transmission characteristics were closer than those on all but a very few other designs even after adjustment. The phase and transmission set* provided for making these measurements is capable of a high degree of precision. However, to obtain the needed precision, an unusual amount of maintenance is required, the set being taken out of service for one day a week for this purpose. The set proved very useful in arriving at a correlation between transformer variations and those of the amplifier stage of which it is a part.

Special sets used for measuring the leakage impedance, the capacitance across the high impedance winding and the capacitance from the high voltage end of that winding to ground were very useful in controlling variations in the transformer.

CORRELATION BETWEEN TRANSFORMER ELECTRICAL AND MECHANICAL RE-
QUIREMENTS

The purpose of electrical tests is to insure proper performance of apparatus in its end use and to provide process controls as an aid in meeting overall manufacturing objectives. Because of its unique mechanical construction and the precision to which the dimensions of its components are held, the first order causes of variation in the parasitic inductances and capacitances used to shape the transformers' transmission characteristics are under close control. Certain secondary causes of variation show up, however, which are quite different from those of the average carrier transformer. This is aggravated by the fact that the transformer is part of the amplifier feedback circuit and in consequence, certain parasitic capacitances have a much greater effect than they would have if the apparatus were used in a less complicated circuit.

Statistical control charts were very useful in meeting the test objectives on the 2504A transformer. Continuous control charts showed small

* D. A. Alsberg and D. Leed, A Precise Direct Reading Phase and Transmission Measuring System for Video Frequencies, B.S.T.J., **28,** p. 221, April, 1949, and D. A. Alsberg, A Precise Sweep-Frequency Method of Vector Impedance Measurement, Proc. I.R.E., **39,** p. 1393, Nov., 1951.

but very noticeable leakage inductance variations as measured at seven megacycles. Discovery of the reasons for the variations was not as obvious. Calculation showed that the possible contribution of the various factors in a commonly used formula for the total leakage inductance referred to one of the windings

$$L = \frac{10.6N^2l(2\ Ct\ +\ a)\ \text{henries}}{bc^2 \times 10^9}$$

where $N$ = turns in the winding
  $l$ = mean length of turn for all windings
  $c$ = number of insulation spaces
  $t$ = thickness of insulation space
  $a$ = height of window opening
  $b$ = width of winding

was too small to account for the variations experienced, small though they were by comparison with those of conventional transformers. The number of turns in the windings and the number of insulation spaces being automatically fixed and charts of other factors indicating adequate control focused attention on the assumptions used in arriving at the formula expressed above. One of these — the assumption that the leakage flux is uniform and parallel to the axis of the core — provides the clue to the variations. A certain amount of flux must intercept the shields and set up eddy currents. Differences in the effective thickness of the shields were found to be major contributors to leakage impedance variations. This led to the establishment of the elaborate controls on the spraying process described earlier in this article.

Control of parasitic capacitances presented a somewhat different problem. While the interwinding capacitances of the transformer are complexly interrelated, charts kept for control purposes showed fair correlation between the wall thickness of the outer form and the capacitance across the high impedance winding. Possibly greater correlation could have been obtained if the average wall thickness could have been arrived at through a weighted integration process. Despite the fact that the data were not entirely adequate, a number of irregularities in mechanical processes were discovered as a result of the keeping of control charts on this characteristic. Typical was the discovery of a worn condition on a carbide rod used for contacting the inner cylindrical surface of the outer winding for one of the gaging operations. Through continued use, a crescent shaped valley had been lapped out of the middle rod by the abrasive action of the glass parts in inserting them in the gage. As a

result, a condition that might have gone undetected for some time was corrected before a large quantity of expensive parts were completed.

Numerous correlations similar to the two just described have been helpful in establishing manufacturing procedures on the 2504A transformer and have contributed a great deal to the setting up of satisfactory process routines.

# Development of Reed Switches and Relays

By O. M. HOVGAARD and G. E. PERREAULT

*Improvements in the operating speed, efficiency, compactness and contact reliability of relays may be obtained from the application of sealed-in-glass reed switches. This paper is an account of the development of such switches and relay designs suitable for their use. Because of the close dependence of contact performance and switch operating characteristics upon process methods their development has been an important part of this project.*

## I — REED SWITCHES

Switches are those basic circuit elements of a telephone system which make it feasible to establish the connections desired by telephone users. To permit remote operation, switches in the telephone plant are usually electromechanical devices and they are used to select and establish needed talking paths. Since switches are used in large numbers their manufacturing and operating costs have a significant impact upon telephone economics and for this reason the development of switching apparatus receives much attention.

### HISTORICAL

In 1936 the availability of new magnetic alloys prompted a study of improved means for the operation of switching contacts. As a part of this work there was conceived an extremely simple magnetic structure giving promise of excellent electromechanical efficiency, unusually high operating speeds and lending itself to very compact designs. This paper is an account of the development of this basic concept to the point of practical application in the Bell System.

The basic concept is shown in Fig. 1. Two flat reeds of magnetic material are supported as cantilevers with their free ends overlapping and separated by a small gap. Surrounding the reeds is an operating coil so

placed that the overlap of the reeds is at its center. When the coil is energized the flux in the gap pulls the reeds together. The reeds perform the double function of a magnetic operating gap and a contact pair with which to close and open an electrical load circuit. This is a radical departure from the classical relay wherein the operating gap is external to the energizing coil and the armature operates the electrical contacts through mechanical linkages. The result is improved electromechanical efficiency through reduction in magnetic leakage and increased operating speed as the result of smaller displacements and moving masses.

In practice the reeds are supported by sealing them into the ends of a glass tube. This provides a structure sufficiently stable mechanically to assure the maintenance of the critical gap between the reeds and gives
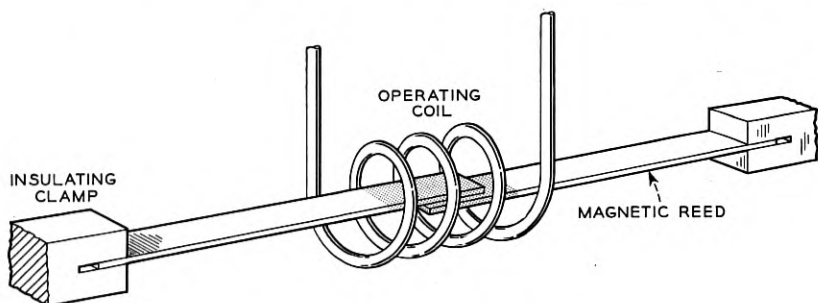


Fig. 1 — Basic reed switch concept.

excellent electrical insulation between them in their unoperated position. Furthermore, by providing a hermetic seal around the electrical contacts it frees them from the effects of environment. This is of fundamental importance since corrosive or dirt-laden atmosphere is a major cause for expensive contact maintenance.

Initial embodiments of this design concept were not as simple as indicated in Fig. 1. Materials having optimum magnetic characteristics do not in general make sound gas tight seals to available glasses, nor do they serve satisfactorily as electrical contacts if exposed to the air. Consideration of these factors led to the design shown in Fig. 2.* Reeds of the desired magnetic material are welded to supports of materials which will seal to glass and serve as electrical connections to the contacts. One of these supports is a rod, the other is a metal tube. After the switch has been assembled it is evacuated through the tube and filled with the de-

* U. S. Patent 2,289,830.

Fig. 2 — One type of early reed switch design.

sired gas and the tube is closed with a plug of solder. The proportions of this design were very suggestive and the first application, in about 1938, was in a carrier system where the switch was made a part of the central conductor in a coaxial line. The operation of the switch is controlled from a coil surrounding the line and complete isolation is obtained between the high frequency paths and the control circuits. Switches in this use have given trouble-free service for more than fifteen years.

It was natural that this radically new switching element should find important military application because of its independence of such environmental factors as corrosion, dirt and altitude. Under the impetus of World War II a dozen or more designs were evolved for specific applications and much was learned about the basic design as well as the process problems of manufacture. This experience demonstrated that with its inherent compactness, efficiency and operating speed this switch could find broader application in the telephone plant and its development for this purpose was undertaken. To fully attain this goal the design and manufacturing process must result in a cost comparable to that of a pair of nickel silver springs carrying the customary precious metal contacts and the operating life of the switch should approach one billion operations for electrical loads up to one-half an ampere.

DEVELOPMENT

The operating characteristics of a switch, i.e., the ampere turns to cause it to close, to hold it closed or let it open, depend upon the magnetic characteristics of the reeds, their dimensions, overlap, separation in the unoperated state and the amount of contact plating. These factors also influence the relationship between operating speed and power and thus affect switch efficiency. In addition to these design factors such process matters as the flatness of the reed surfaces, precision of their alignment and the accuracy with which the unoperated gap is established determine the tolerances which can be placed upon the switch operating characteristics. Switch and process design are inseparable and the degree of development of the latter becomes controlling in the determination of attainable switch operating tolerances. Therefore, studies were made to explore both of these phases of reed switch development.

The relationship between the operating characteristics and the unoperated gap of an experimental reed switch is shown in Fig. 3 for several values of reed overlap. As the gap is increased, the reed deflection and the magnetic pull required to close the switch increase. Therefore, the flux density in the gap must be increased correspondingly, and, as shown by the curves, the onset of magnetic saturation of the reeds and increased leakage across the gap cause the relationship between the ampere turns required to close the contacts and the switch gap to become increasingly nonlinear. Magnetic saturation of the reeds will therefore determine the maximum gap at which a given switch design can be closed. The restoring forces to separate the reeds, i.e., release the switch, when the operating flux is removed are obtained from the deflections of
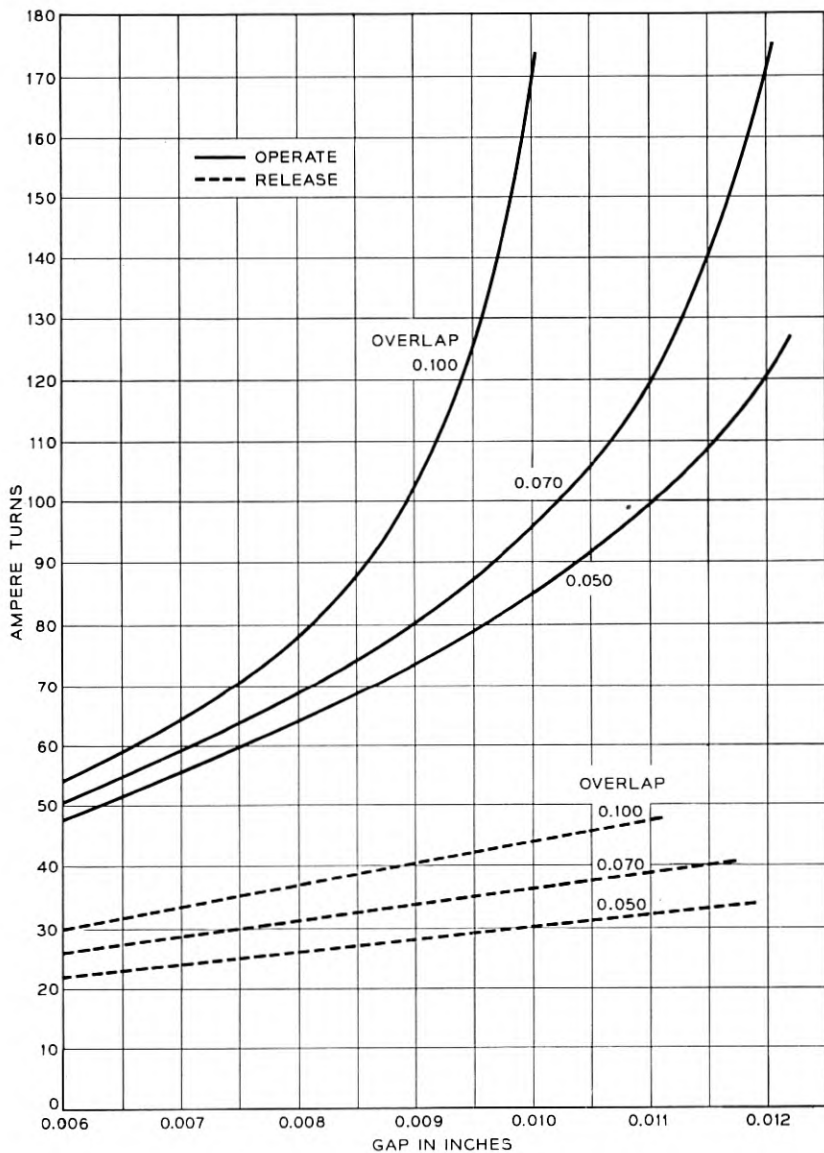
Fig. 3 — Operating characteristics of a reed switch versus the gap in the un-operated switch.

the reeds upon switch closure and therefore are proportional to the un-operated switch gap. A switch will open when the restoring forces exceed the magnetic pull between the reeds. Although under these conditions the flux density in the reeds and the gap leakage are low, there are magnetic saturation effects at the surfaces of the reeds which cause the
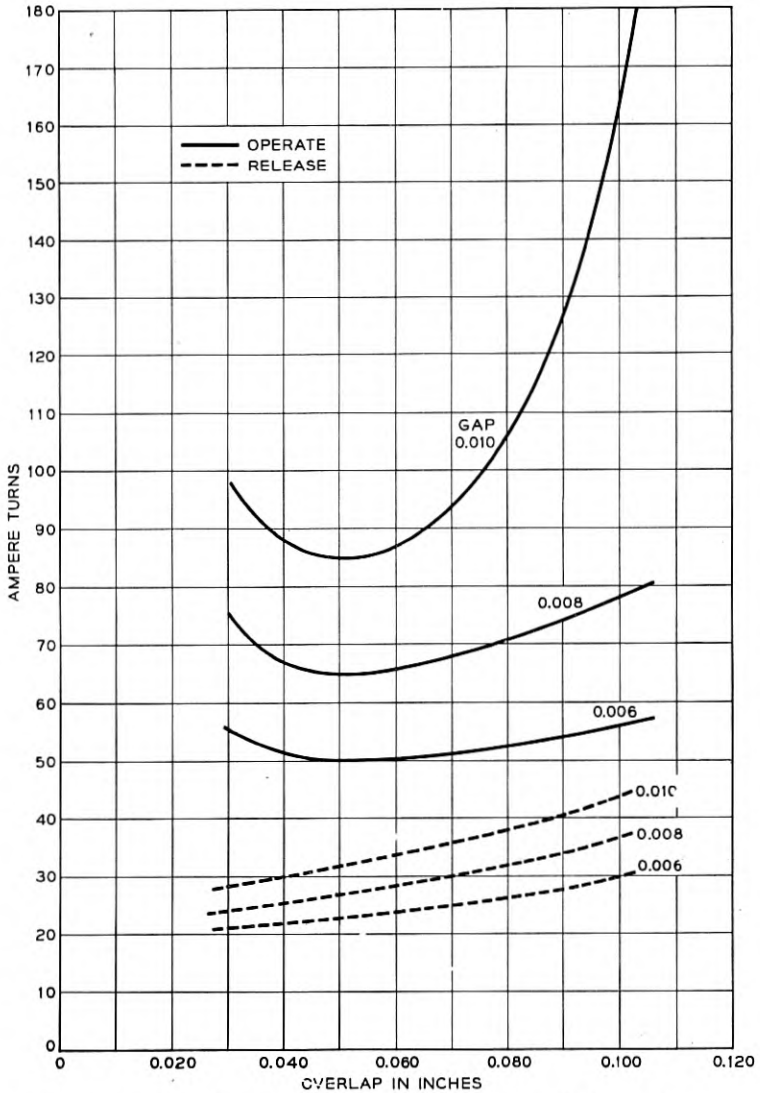


Fig. 4 — Operating characteristics of the same reed switch as in Fig. 3 versus the reed overlap.

relationship between switch release and switch gap to be approximately linear. The relationships for the same reed design as a function of reed overlap is shown in Fig. 4. The shape of the operate characteristic arises from the decrease in percentage of leakage flux around the gap and the decrease in the gap flux density as the gap area is increased. These factors have opposing effects upon the switch ampere turn operate value, the former predominating at small overlap, the latter for large overlap, and therefore the characteristic exhibits a minimum when the reluctance of the working gap equals that of the remainder of the magnetic circuit. The extreme curvature for the 0.010″ gap at the higher values of overlap is caused by magnetic saturation of the reeds. The rising character of the release curves is due to the decrease in flux density as the overlap is increased. This results in less magnetic pull and a higher switch ampere turn release value.

The contacts of reed switches are subject to electrical erosion which may take the form of a buildup on one reed associated with a pit on the mating reed. As with ordinary contacts, such buildups can give rise to mechanical locking of the reeds causing failure of the switch to open. The reeds may also become welded upon closure by excessive currents or capacitive discharges from the electrical load circuits. The restoring forces of the reeds are counted on to rupture such locks or welds as may occur within the specified limits of contact operation. Since the restoring forces are proportional to the switch gap better margins against locking or welding will be provided with larger switch gaps. However, the characteristics described above show that increases in the gap rapidly decrease the switch magnetic efficiency. Furthermore, the increase in the slope of the characteristic indicates that, for larger gaps, we should expect greater dispersion in the operate characteristic for a given precision in establishing the nominal gap in manufacture. Small gaps look attractive from the standpoint of switch efficiency but are more difficult to establish in switch assembly. In addition, the by-products of contact electrical erosion are magnetic and tend to collect in the switch gap and cause trouble if this is too small. The selection of the gap size for a given switch design is a compromise of these factors of switch application and process design. The reed overlap is selected to fall at the minimum point of the characteristic which at once provides maximum switch response and minimum sensitivity to manufacturing variations in overlap.

In the assembly of a reed switch the seals are the anchorage for the reeds and their location determines the effective cantilever lengths and the stiffness of the reeds, and the lengths of the seals influence the rigidity of the reed anchorage. Thus it is important that the sealing technique

permits good control of this process step. Since a switch is only about three inches long the close proximity of reed holding chucks and mechanisms for aligning and spacing the reeds make the open flame sealing techniques of the vacuum tube art unattractive. Furthermore, with open flame techniques it would be difficult to obtain adequate control of seal
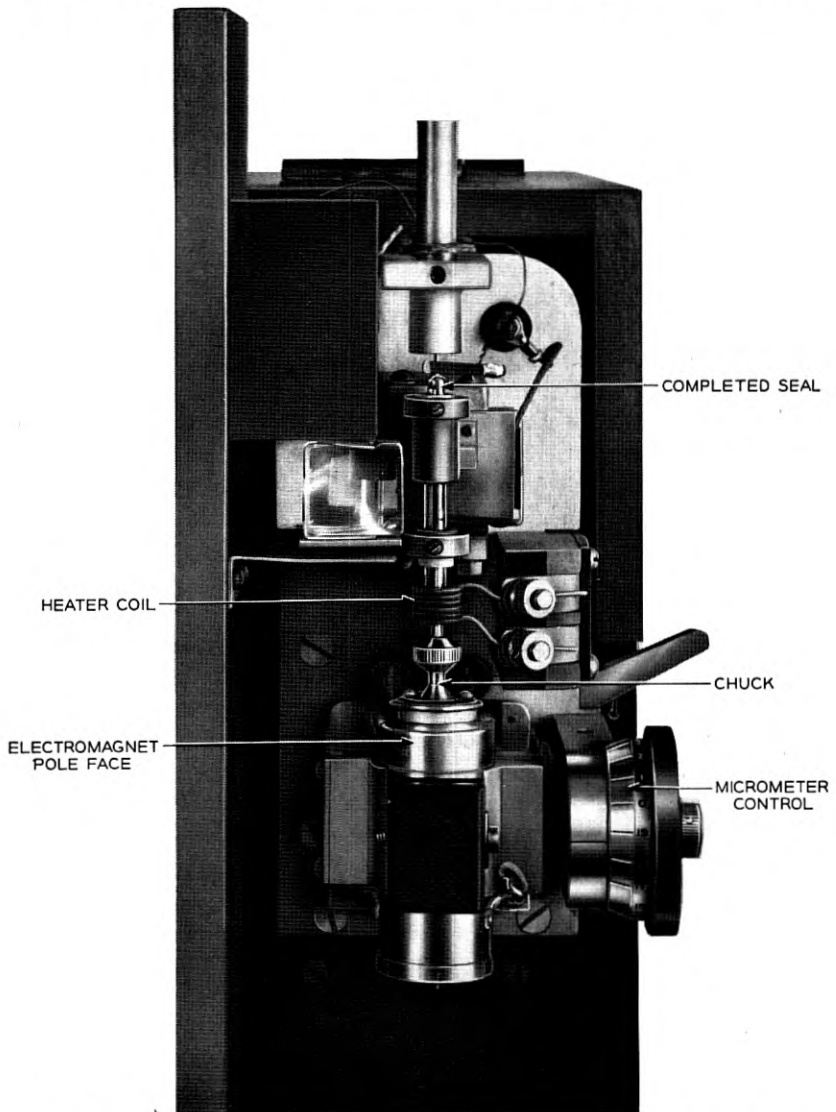


Fig. 5 — Early laboratory type machine for assembling reed switches.

location and dimensions. A new method had to be conceived and this has taken the form of an electrically heated coil of resistance wire surrounding the ends of the glass envelope. A laboratory facility for assembling switches is shown in Fig. 5. One of the reeds has already been sealed into one end of a glass envelope by a previous operation and this subassembly is shown mounted in position for receiving the second reed. The glass tube of the switch is to the right of a prism used in the process of aligning the reeds. The seal already made can be seen projecting from the mechanism just below a spring clip providing an electrical contact to the reed support. The lower end of the glass tube projects into and is closely surrounded by an electrical heater coil. Below the heater is a chuck for supporting the reed about to be sealed into the switch. This chuck is of magnetic material and can slide on a steel surface which is also the pole face of an electromagnet. The chuck may be moved in rotation and translation until the reeds are in alignment as determined by visual inspection through the prism after which it is locked in position by energizing the electromagnet. The chuck and its support are fastened to a carriage which can be moved in translation under control of a micrometer having a calibrated head, and this is the means for establishing the desired switch gap. After this operation the heater coil is energized to make the seal. The temperature of sealing is controlled by the voltage applied to the heater coil and the time of application by an electrical timer. Most of the heat is transferred to the glass by radiation and the temperature of the glass is raised sufficiently so that its surface tension can overcome its viscosity. Surface tension draws the open end of the glass body inward until the glass comes into contact with and wets the metal oxides of the reed support thereby effecting a seal. In addition to performing well the sealing function, this method also permits control of the seal cooling so as to provide a desirable strain relieving anneal of the glass.

In developing the process for making a metal-to-glass seal it is necessary to obtain detailed knowledge of the thermal characteristics of the materials to be used as well as the temperature-time phenomena involved. In the range between room temperature and the softening point of the glass the coefficient of expansion of the metal is essentially linear and its magnitude depends on the composition of the alloy. In this same range the coefficient of the glass is nonlinear but it is possible to select metal and glass combinations which will exhibit low seal strains after the combination has cooled to room temperature. At temperatures between the softening point of the glass and room temperature a seal will be subjected to varying degrees of stress depending upon the choice of

alloy composition and these stresses may become serious enough to cause the seals to crack. The stresses at intermediate temperatures may be reduced by suitably annealing the seal but they can become greatly aggravated if the thermal characteristic of the sealing machine does not permit the metal and the glass to cool at approximately the same rate. In consequence the thermal characteristics of the means for supporting the reeds during sealing, the rate of cooling of the electrical heater coils and the susceptibility to the effects of drafts are all matters which must be studied in the design of the sealing process.

To obtain switches with operating characteristics that are predictable within acceptable tolerances it is necessary that the reeds be parallel to each other. This requires that the reeds be flat and straight because their final position in the switch is determined by reference to the means for supporting them during sealing. It is obvious that if the reeds are bowed or otherwise not flat it will not be possible to obtain good alignment of the overlapping ends. As the reeds are held in some form of chuck this must be capable of seizing a succession of straight reeds in exactly the same manner and its function must not be interfered with by residual burrs or surface roughness or by reasonable manufacturing variations of the reeds. After the reeds have been placed in the chucks it is necessary to establish the operating gap between the overlapping ends. This is of the order of 0.01 of an inch and it is desired to hold this to within a range of about 0.001 of an inch. There are two basically different ways in which the desired gap can be established. One of these is to provide a means for bringing the two reeds exactly into contact and then, using this position as a reference, separate the reeds by the desired amount. This method is relatively simple to execute to the desired precision but it provides no means to allow for the effects of such process variables as the failure of the reeds to be ideally straight or minor failures of the chucks to grasp properly. Another method is to place the reeds initially so that the gap is larger than that finally desired, subjecting them to a suitable magnetic field and then advance the reeds toward each other until they are closed by the effects of the magnetic field. This has the advantage that the reed separation is determined directly in terms of one of the final switch operating characteristics. It takes full account of the consequences of all of the process variables which might affect this particular characteristic but has the disadvantage of greater operational complexity. Both of these methods of assembly have been used during the development of the switch and it appears that essentially equivalent results in terms of switch performance and yields can be obtained.

A further cause for variability in switch product arises during the

sealing of the reeds into the glass envelope. Molten glass has both surface tension and viscosity and since the seals are to be located between the chucks and the free ends of the reeds, differential cooling can cause the surface tension forces to impart displacements to the reeds with respect to their desired final positions. Observations during sealing indicate that the displacements of the reeds may be quite large and erratic and the residual errors of displacement from this cause will be determined by the viscosity of the glass and its permitted rate of cooling.

The process factors mentioned in the preceding paragraphs are important to design because they control the magnitude of the requirements which the designer can place upon the product and therefore have a marked influence on the potential value of the design in terms of circuit application. Process studies are for this reason a very important part of design. Laboratory development must usually be conducted on the basis of making a rather small number of models and it is further limited by the necessary use of hand-operated equipment differing materially in characteristics from the large scale facilities normally acquired for regular manufacture. Recourse must be had to careful design of experiments and to the extensive application of quality control and other statistical methods based on the theory of small samples. Through the use of these methods it has been possible to develop proper reed designs, chucking methods, sealing cycles and gas-filling procedures, and to lay a sound foundation for the design of large scale automatic facilities.

CONTACT INVESTIGATIONS

It is known that some airborne corrosive elements, organic material and dirt have deleterious effects on contacts but it cannot be postulated that therefore the exclusion of air will give rise to ideal contact operating conditions. The volume of gas included in a switch of acceptable size is small and the concentration of contaminants introduced in manufacture could be much higher than that likely to be encountered in large ventilated spaces. Therefore, contact investigations of sealed switches must be based upon completed units containing the materials under study and made by a process sufficiently defined to permit repetition of results. Only then can the effects of materials be differentiated from those of the process and the latter be designed for optimum results.

Early switch designs used reeds made from pure iron and from Perminvar. Although these reeds had been cleaned by heat treatment in hydrogen their brief exposure to air in the course of switch assembly permitted sufficient oxidation to adversely affect their performance as

contacts even in switches which were evacuated and filled with hydrogen. This condition was overcome by gold plating the tips of the reeds before heat treatment, and it was the favorable results with such switches that gave rise to the present development. Because they could not be sealed to inexpensive lead glasses, iron and Perminvar reeds had been welded to supports of a nickel-iron alloy specially designed for this purpose. If this material, called 52 alloy, could be used for the reeds the switch structure could be simplified and its cost reduced. Tests showed the alloy to have acceptable magnetic characteristics and a program was initiated to study its properties as a contact material. Perminvar was included in this program to provide a control against previous experience.

As pointed out above, the restoring force of a switch is proportional to the reed travel to contact closure. Since this is only about 0.005 inches, the thickness of contact metal added to the reeds must be small to avoid significant loss of restoring force. This suggests electroplating as the method for applying the contacts and studies were conducted on the basis of such a process step. Contact materials such as gold, rhodium, chromium, etc., and such gases as helium, hydrogen and nitrogen were included in the test program. Switches containing various combinations of these were placed in operating life tests at selected loads up to a half ampere. Early trends showed that there were no significant differences between Perminvar and 52 alloy as the underlying material for the contacts, and that gold plating in combination with either hydrogen or nitrogen gave better results than other combinations. Therefore subsequent interest centered on detailed studies of gold plated 52 alloy reeds in atmospheres of hydrogen and nitrogen and the effects of process upon such combinations.

Like ordinary contacts, reed switches are subject to electrical erosion and observations during life tests show that this takes three distinct forms. In one type, a buildup is formed by transfer of material from one reed to the other and eventually the buildup will lock into the pit of the mating contact and the switch will fail to open. A second type of erosion results in the formation of powdery and flaky particles and since these are magnetic they will collect in the contact gap under the influence of the operating flux. Under these conditions the erosion products appear to be brittle and give the impression that they may have been incipient buildups which were broken by the impact of the reeds as the switch closed. The area of erosion is generally large compared to that resulting in an early buildup showing a disposition of the point of contact to wander over the contact surface. The volume of particles in the contact gap is proportional to contact use and it seems likely that quite early in the

switch life, if no single predominant buildup has occurred, the electrical contact between the reeds is through these particles which are conducting and are distributed over the contact area rather than by direct contact between the reed surfaces. While some of the particles may be welded to the reed surfaces most of them are loose and under the influence of the gap flux become rearranged at each switch operation. This action may promote distribution of contact wear but can also lead to aggregations which can become welded together and initiate a predominant buildup resulting in ultimate contact locking. It has also been observed that after many operations the volume of particles may become great enough so that their rearrangement will permit them to bridge the gap of the unoperated switch and thus cause contact failure or other malfunctioning. Obviously, this type of erosion is very complex and its consequences are difficult to predict. A third type of erosion occurs in a significant number of switches. Here a buildup takes the form of a puddle-like formation covering an appreciable part of the contact surface with a corresponding shallow depression in the mating contact. Contacts with this type of erosion have much less loose material than contacts of the second type but when they fail it is due to the relatively sudden appearance of a buildup causing locking.

Switches having only the first kind of erosion frequently exhibit short life and are rather uncommon. Those of the second kind are representative of the present state of the art and account for the bulk of the failures observed. Switches having the third kind of erosion give an operating life about ten times as long as those having the second kind. The causes for these differences in contact erosion are not well understood but they appear early in the switch life. For this reason it is suspected that they may be related to such factors as the smoothness and alignment of the reed surfaces, the degree of penetration of the contact material into the underlying reed and process contaminants. These can all affect the reed surfaces and influence the nature of the initial erosion and thereby possibly determine its course through the switch life. The large number of interrelated variables bearing on switch life makes it very difficult to obtain an evaluation through small scale laboratory experiments and there can be no assurance that the same factors will exist in the process for large scale manufacture. For this reason an important part of the continuing program for contact development will be the study of the output from regular manufacture and the effects of variations introduced in the process.

Operating life tests for reed switches have been conducted mostly at one-eighth and one-half ampere resistive or protected inductive loads.

Tests of switches from the early stages of large scale manufacture show that a few of these failed as early as 200 million operations, about half had failed at one billion operations and some were still going at two billion operations. This large spread in the performance of supposedly identical switches does not appear related to the load conditions and therefore must be considered a switch attribute. It is thought that in part it must be due to the fortuitous nature of the second type of erosion described above since both it and the third type of erosion are found in the long life switches and there were no early failures having the third type. In evaluating the behavior of sealed switches it must be remembered that their contacts are not accessible for servicing and for this reason the first failure of a switch must be considered the end of its useful life. On this premise reed switch life in the load range from one eighth to one-half an ampere must be taken as about 200 million operations. For smaller loads longer life may be expected exceeding a billion operations for very small loads. The capabilities discussed are all for 48-volt load circuit conditions.

DESIGN FOR MANUFACTURE

With this background of information it became possible to adopt a switch design and a manufacturing process which would permit a sufficient degree of mechanization to predict acceptable costs. In this design shown in Fig. 6 the reeds are made from 52 alloy wire in the shape of a double-ended paddle with a round section near its middle. The round section provides optimum conditions for minimizing stress concentrations in the metal-to-glass seal. One flat end of the paddle provides the moving reed inside of the switch and the other provides the external electrical connection. Both ends are gold plated, one to act as the elec-



GLASS TUBE

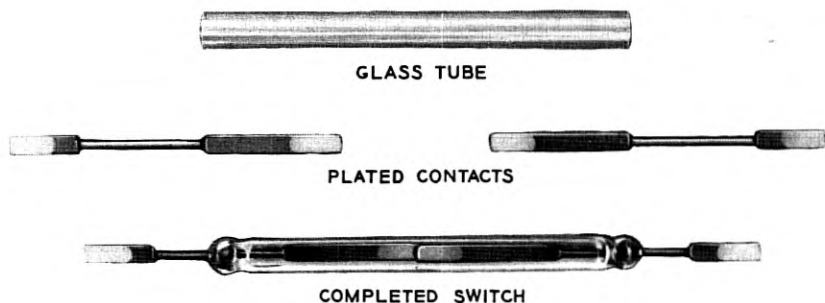PLATED CONTACTS

COMPLETED SWITCH

Fig. 6 — Reed switch design suitable for large scale manufacture.

trical contact and the other to provide for ease of soldering an external connection to the switch. The final design also had to avoid the use of metal tubulations if cost objectives were to be met, and the methods for pumping and filling used in the vacuum tube art were considered too expensive. Novel methods were explored. One of these was to enclose
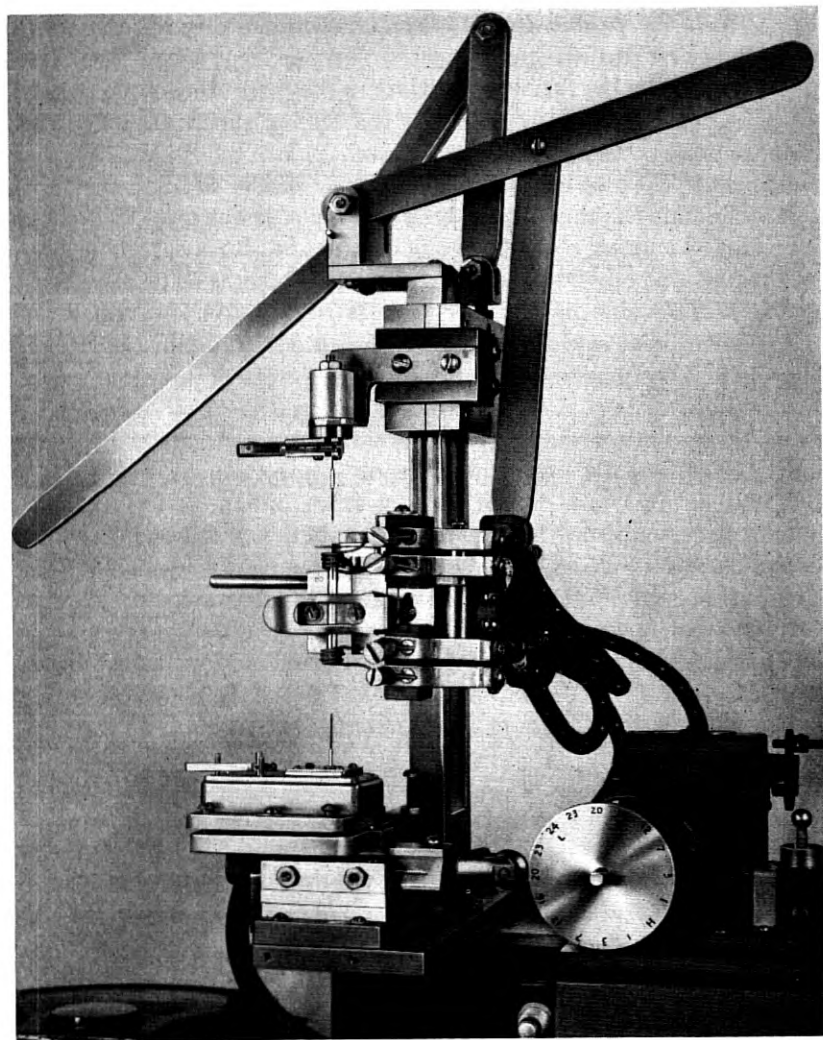


Fig. 7 — Laboratory facility for the semiautomatic assembly and flush-filling of reed switches.

the entire switch assembly machine in the atmosphere desired for switch filling. This method received considerable experimental study and it was concluded that the problems of switch cleanliness in the presence of machine lubricants and other contaminants would be difficult to solve. The means selected was to mount the reeds and the glass envelope in the proper position for sealing and then displace the air in the envelope by flushing with the desired gas. Extensive experiments were required to design means for introducing the gas so that aspirator effects would not draw in air during the flushing process and the whole technique of glass sealing had to be revised to allow for the cooling effects of the flowing gas as the seals were closing. Fig. 7 shows the laboratory setup for this development. The two heater coils are clearly visible and between them a clip for holding the glass tube. Above the upper heater is seen the chuck for holding the upper reed. At the bottom is the chuck for holding the lower reed and it provides the means for introducing the gas including a chamber for assuring its lamellar flow. It is mounted on a track and under control of a servo motor may be moved to obtain the desired switch gap. The upper chuck and the sealing heaters are mounted on vertical slides so that after the reeds and glass tube have been loaded the reeds can be inserted in the glass tube by operation of the levers shown. After this the machine operation is programmed by auxiliary facilities. The air is flushed from the glass tube while the servo establishes the desired switch gap. The upper heat coil is then energized and when this seal has been effected the velocity of gas flow is lowered and the lower seal is made. Aside from their value in supplying requirements for final machine design these automatic features were an absolute necessity in providing a sufficient degree of process stability to permit significant conclusions to be drawn from small lots of switches intended to display the effects of specific design or process changes.

## II — REED RELAYS

The combination of a reed switch and an operating coil constitutes a relay and in contrast to other relays no magnetic core, armature or contact carrying springs are required. The design of switch developed provides a simple make contact. However, as will be described below, break and transfer functions and other worth while features can be obtained by associating a permanent magnet with a make contact reed switch.

### RELAYS WITHOUT MAGNETS (NEUTRAL)

If a single switch is placed in an operating coil and the current through the coil is increased there will be a definite value of ampere turns at which
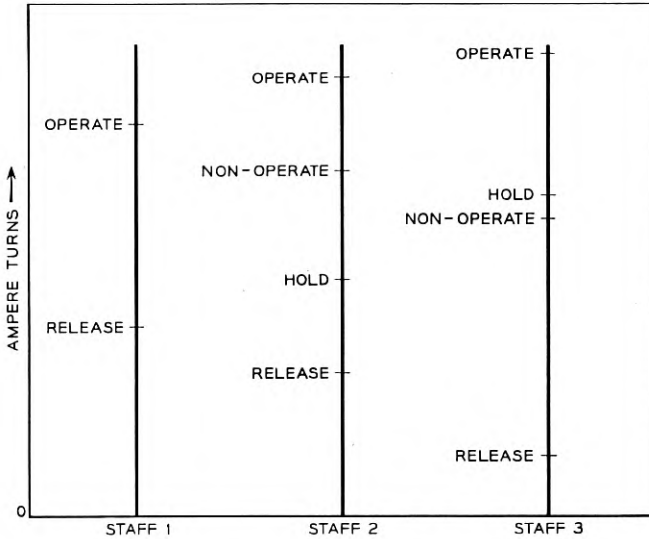
Fig. 8 — Comparison of the operating characteristics of a single reed relay with that of the relay type universe of which it is a sample.

the switch will close. When the current is then decreased, there will be another definite value at which the switch will open. These are called the operate and release values of that switch and are shown on Staff 1 of Fig. 8 where the ordinates are operating coil ampere turns. Sealed switches are not adjustable and their individual operating values are subject to manufacturing tolerances. In consequence, if a succession of switches are measured in the same operating coil the results are as indicated on Staff 2. As the coil current is increased from zero there will be an ampere turn value at which the most sensitive switch will just fail to close; this is called the non-operate value of that relay design and is the highest ampere turn value at which no relay will operate. As the ampere turn value is further increased a value will be found just sufficient to assure that all switches will close; this is called the operate value of the relay. If the coil current is then decreased, a value will be found below which one or more switches will open; this is called the hold value of the relay and the still lower point at which all switches will just have opened is called the release value of the relay. The difference between operate and non-operate and that between hold and release represent the manufacturing spread in the nominal operate and release values of a given relay design. These four characteristics are important to circuit designers and considerable development effort is justified to

obtain good margin between hold and non-operate. Staff 3 shows the consequences of a lower grade of design and process where the spread is so great that the hold value is higher than the non-operate value. Such a product would require selection of switches for relays specifying these parameters and this would not be tolerable as the basis for wide application.

If two switches are placed in the same operating coil the events are more complicated. As the current through the coil is increased, the more sensitive switch will close first and provide a magnetic shunt across the second switch which therefore will require more ampere turns to close than would otherwise be needed. Upon decreasing the ampere turns after operating both switches, it would be expected that the switch with the higher restoring force would open first. This may or may not be the case depending upon the relative closed gap reluctances of the switches. The result is an increase in the dispersion of the operating characteristics of the relay and difficulty in specifying the sequence of operation of the switches. As an example of this effect, for two similar designs we find the following results in ampere turns:

|  | Single Switch Relay | | Two Switch Relay | |
| --- | --- | --- | --- | --- |
|  | Average | Tolerance | Average | Tolerance |
| Operate.................. | 79 | ±16 | 92 | ±22 |
| Release.................. | 25 | ±12 | 33 | ±17 |

MAGNET BIASED RELAYS

By suitable association of a permanent magnet with a reed relay the operating characteristics become radically altered from those otherwise obtainable. Referring to Fig. 9 the axes are ampere turns, the ordinates representing those due to the operating coil and the abscissa those due to the permanent magnet. The polarity of the permanent magnet is assumed to be aiding the operating winding in the direction taken as positive and the curves depict events as the strength of the magnet is varied.

On the staff of ordinates for zero flux from the permanent magnet are shown the four operating parameters for the particular relay design. It should be noticed that these are repeated in mirror image fashion for negative coil ampere turns. In other words, if there is no permanent magnet the closure and opening of the switch depends only on the magnitude of the coil flux and is independent of its direction. As the strength of the permanent magnet is increased from zero, less coil flux will be

needed to close the switch and the four relay parameters will behave as shown by the curves. To illustrate, if the magnet strength were 10 ampere turns the coil operate ampere turns would be reduced from 110 to 100 and the release ampere turns from 20 to 10 ampere turns. If negative coil flux were applied for this condition it would take 120 ampere turns to close the switch compared with 110 ampere turns for the case of no permanent magnet. Therefore, increasing the magnet strength increases the sensitivity of the relay and provides additional margin against operation by unwanted negative current. As the strength is increased further a point will be reached where it is equal to the release value of the relay design. Beyond this point some of the relays will fail to release on coil open circuit. This limit, indicated by the dashed ordinate A gives the
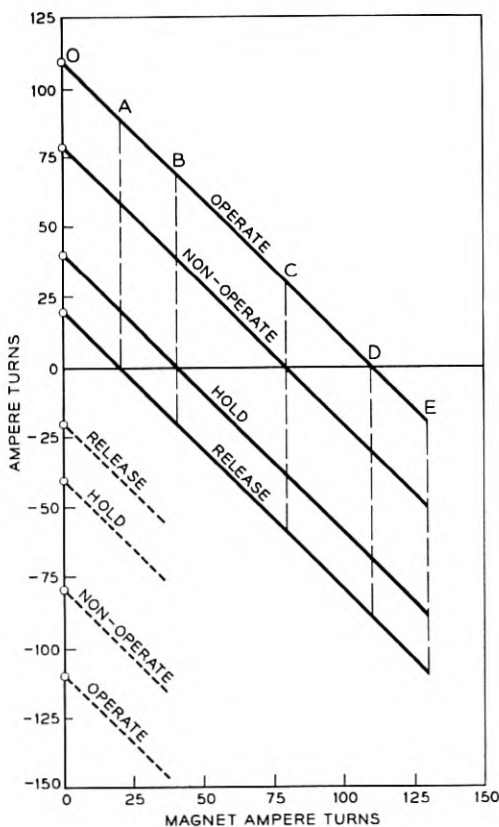


Fig. 9 — Typical operating characteristics of a magnet biased reed relay.

useful range for magnet bias as a means for increasing relay sensitivity. In the operating range from 0 to A there is the additional possibility that by adjusting the magnet strength the spreads in the operate or release characteristics of the relay can be reduced. However, both of these values cannot be controlled simultaneously.

In the operating region between A and B some of the relays will release on coil open circuit and others not. This operating region would require selection of switches and is only useful for special applications.

PULSE OPERATION

The operating region between B and C is of great interest. Any relay in this region when operated will remain operated and will require reversed coil current to be released. Due to the high speed of the switch, operation can be accomplished by pulses of a few milliseconds and no holding power will be required to keep the switch closed. Likewise the switch can be released by a negative pulse. In this region the switch can perform either the make or the break function as desired and the pulses may be symmetrical or assymmetrical in either direction. This is also the most sensitive operating region of the relay.

NORMALLY CLOSED OR BREAK OPERATION

Further increases in the magnet bias in the region from C to D yields another region of uncertainty. Operation here would be with negative coil current which should cause the relay contacts to open but some relays would fail to reclose on open circuit.

The region beyond D is again generally useful and here the operation is that of a normally closed contact which opens on negative coil current. When the magnet strength becomes equal to the sum of the operate and release ampere turns of the switch design the relay is a break relay with the same sensitivity as a neutral make relay. The additional cost of a break relay will be that of the magnet and its adjustment.

TRANSFER RELAY

The transfer function can be performed by using a make and a break relay with simultaneous excitation of their windings and suitable connections to their switches.

RELAY DESIGN

The design of a relay includes a packaging problem. Heretofore this has been based upon using the relay core as a central stucture to which
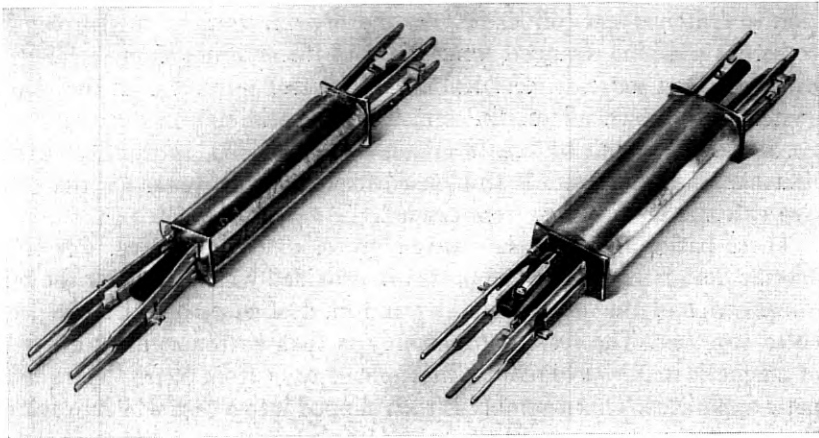
Fig. 10 — Models showing the construction of basic relay units. These may be used singly or in assemblies.

was attached the operating windings, armature, contact springs, covers and the means for mounting this assembly. Reed switch relays having no magnetic cores must use other means to provide a method for mounting.

In the case of reed switch relays a new design requirement arises from the needs of the maintenance personnel in the operating telephone plant. Over the years, maintenance personnel tracing circuit troubles have been aided by visually determining the state of operation of relays and other electromechanical devices. Additionally, by having access to the contacts they could as desired isolate portions of circuits by "toothpicking" individual contacts to cause them to remain open or closed regardless of the operating state of the relay. Because the reed switch magnetic gap and contacts are inside of the operating coil neither of these well established techniques can be used and the relay designer must provide other means for maintenance.

One design approach to the reed relay rests upon the concept of a basic building block composed of a coil containing 1, 2, 4, 6 or 10 switches. The inside dimensions of the coil and its shape conforms to the number of switches to be used and applications requring an odd number of switches can use the next higher even sized building block. Fig. 10 shows two of these basic units, one for a single switch, the other for two. Shown in the picture are the terminals which take the form of metal strips projecting beyond the switch ends. These terminals are held in position on the outside of the coil by acetate spoolheads and locked by acetate sheeting. Each terminal is shaped to provide means for making connec-

tion to switches and coil leads. One end of each terminal is also shaped to accept machine wrapped connections to the external circuit while the other end can serve as electrical access for test purposes. In the single switch unit shown all of the connections to the coil and switch were needed in the circuit so four terminals were provided. In the two switch unit one of the switches is to be used for locking purposes so that only five external connections were needed.

These basic units may be used singly or in combinations. For each specific design simple molded parts are provided which slip over the terminal strips of the desired number and kind of subassemblies and lock them together. The resulting assembly is then surrounded by a shield of magnetic material to reduce interference to or from other electromagnetic apparatus. This assembly is then slipped into a case which provides the means for mounting the whole. Here again new design principles have been applied as seen in Fig. 11. The four corners of the case are extended and provided with detents and the case is suitably slit to control the compliance of the corners. To mount this assembly it is simply pushed into suitable holes in a mounting plate. After passing through the holes the detents lock the assembly into position. This method of mounting eliminates the need for mounting screws and the space they would occupy and permits the fullest realization of the inherent compact-
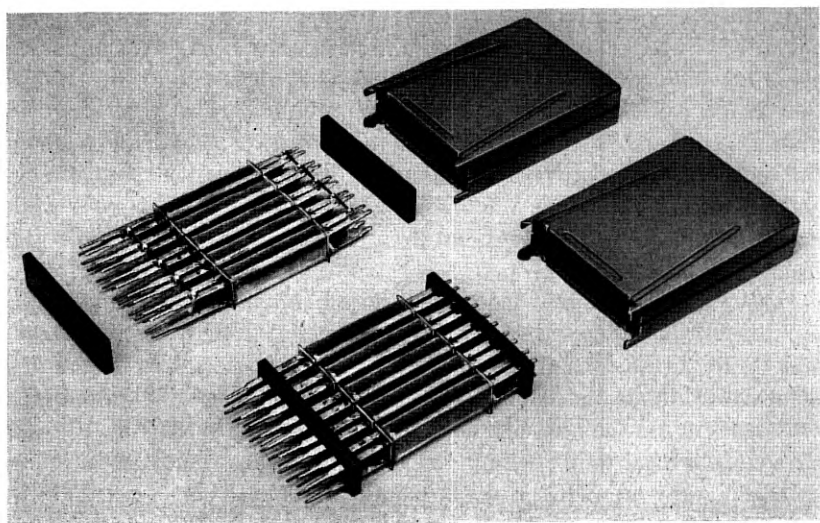


Fig. 11 — A digit register which is an assembly of five relay units each containing two reed switches. A digit is stored by operating two of the relays on the basis of a 2 out of 5 code.
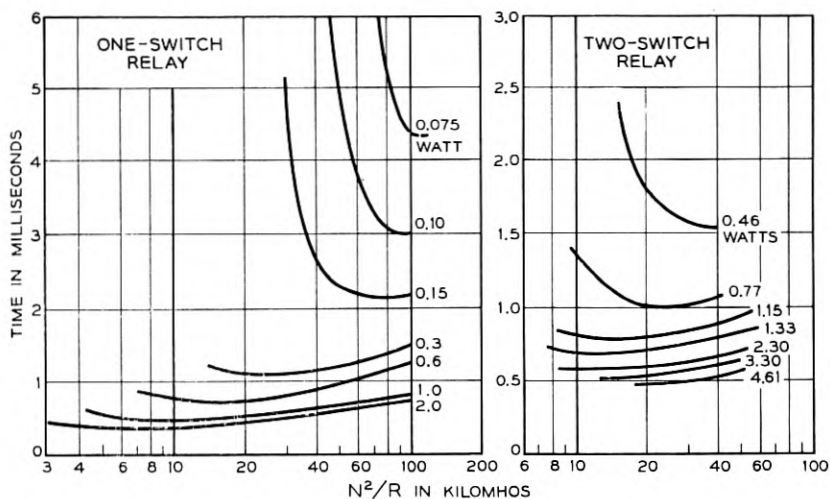
Fig. 12 — Operating speed of typical reed relays versus operating coil design constant $N^2/R$.*

ness of reed relays. Fig. 11 shows the details of an assembly of five relays each containing two switches. This is a digit register for use in storing a digit by the operation of two of the five relays. This register occupies about three square inches of panel space and performs the same functions as five standard relays requiring about ten square inches of space. It appears possible to plan the building blocks so that the height of relay assemblies will permit two rows on a standard 2″ mounting plate. By spreading the relays in two rows the resulting presentation of terminals avoids wiring congestion. The wiring terminals of the relay are at the mounting end while the test terminals are at the free end. This gives access to the relay test points from the equipment aisles and testing from the wiring sides of frames is avoided. By molding grooves in the end insulators and providing soldering lugs on the terminals means for strapping between coils and switches is made possible inside of the relay assembly at less cost than could be done in the course of equipment wiring. Wiring straps can be seen in the illustration and these further reduce wiring congestion on the equipment frames.

The operating speed versus relay coil design constant is shown for various circuit power inputs in Fig. 12 for one and two switch relays of the basic design described above. These characteristics compare very

---

* M. A. Logan, Estimation and Control of the Operate Time of Relays, B.S.T.J., **33** pp. 144–186, Jan., 1954.

favorably with those of other types of relays now used in the telephone plant.

At this writing a number of relays and relay assemblies have been designed for use with reed switches and several are already in manufacture by the Western Electric Company which is also producing the reed switches. All of these are of the neutral type without permanent magnet bias. Applications of these relays are conservative and are made in circuits where the contacts make or break only moderate loads. Further development of the contacts, coupled with field experince will permit gradual broadening of use but even so millions of these switches will be required during the next few years.

# Stability of Negative Impedance Elements in Short Transmission Lines

By J. GAMMIE and J. L. MERRILL, JR.

*Until recently, voice frequency repeaters of the two-way type have been applied almost exclusively to electrically long transmission lines. Now, negative impedance repeaters are used in quantity in the exchange telephone plant, and applications to electrically short lines arise more frequently. Because lower over-all transmission losses can be obtained by utilizing the lower phase shift in short lines, a different engineering approach to the application of E-type negative impedance repeaters is desirable.*

*This paper outlines a general method whereby transmission performance and stability can be related to the characteristics of a symmetrical repeater located in the center of a short transmission line. The theory is particularly applicable to negative impedance repeaters. A tandem arrangement of short sections of transmission line, where each section has a centrally located repeater, can be classed as a line loaded with negative impedance.*

## 1. INTRODUCTION

For a period of about 40 years voice frequency repeaters have been engineered to provide amplification for both directions of conversation in two-wire telephone lines. Some of these repeaters have been operated in lines over 50 miles long; others have been operated in lines shorter than 10 miles. Yet practically all, including negative impedance repeaters of the E-type[1], have been associated with transmission lines which can be classed as electrically long in that they have exceeded one half wavelength at the highest frequency in the pass-band.

Within the past few years the need for two-way amplification in electrically short lines in the exchange area plant has become increasingly evident. A short section of line has limited phase shift at voice frequencies, and advantage can be taken of this fact in the repeater design, to reduce the over-all attenuation below that obtainable with design

methods applicable to electrically long lines. Furthermore, the use of negative impedance devices such as E-type repeaters has made it possible to consider engineering the repeater as an integral part of an electrically short line. This method of design is a logical one because in addition to the reduction in over-all attenuation, the image impedances seen looking into the line terminals are modified by the addition of the repeater. In effect, the philosophy of the hybrid coil and the 22-type repeater is discarded along with the idea that the image impedance of the repeater must match the characteristic impedance of the line. Where the repeater is located a distance less than one quarter wavelength (at a frequency of 4,000 cps) from either line terminal, better transmission performance generally can be obtained by a mismatch between the image impedance of the repeater and the characteristic impedance of the line.

Once a change in philosophy in matching the repeater to the line impedance is made, it becomes easier to forget the repeater as a separate device and to treat it as an integral part of the line in the way a loading coil would be treated. Hence, interest is centered upon the propagation constant and image impedances of the repeatered or loaded line and the transmission characteristics of the device itself are subordinated to this end.

When a two-wire repeater, or its equivalent in the form of a network of active elements, is inserted in a transmission line, stability (freedom from oscillation) becomes a prime consideration. In electrically short lines, the image impedance as well as the loss of the over-all line is a function of the degree of stability desired, which in turn will depend upon the requirements of the system in which the repeatered line must operate.

It is the purpose of this paper to relate transmission characteristics with stability, for a repeater in the form of a symmetrical active network located in the center of an electrically short transmission line. The equations shown and the method of solution are particularly applicable to the fundamental design of E-type negative impedance repeaters in transmission lines wherein the repeater is located less than a quarter wavelength from the line terminals. The frequency at which this wavelength is determined is the highest desired in the pass-band.

The problem is attacked by taking the general case of the symmetrical two-wire repeater located in the center of a transmission line as shown in Fig. 1(a) and substituting for the repeater the equivalent lattice shown in Fig. 1(b). This lattice consists of series arms, $Z_A/2$, and shunt arms, $2Z_B$. The method described herein is general and can be applied to any
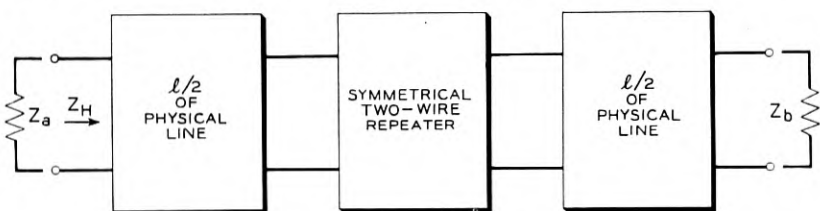
type of impedance in either the series or the shunt arm. However, be cause the specific application considered here is for the $E$-type repeater, $Z_A$ is specified as an open circuit stable negative impedance and $Z_B$ is specified as a short circuit stable negative impedance. This is designated on Fig. 1(b) where these impedances are defined as the ratio of two polynominals which are functions of the complex frequency variable $p$. It is understood that these impedances will have negative resistance components at some real frequencies because the term negative impedance is used herein to describe an impedance whose resistive component is negative within some band of frequencies.

Three conditions are considered:

(a) An E1 or E2 repeater of negative impedance $Z_A$ in series with the line (special case where $Z_B$ is infinite).

(b) An E3 repeater of negative impedance $Z_B$ shunted across the line (special case where $Z_A$ is zero).

(c) The E23 repeater in the line (case of Fig. 1(b), which is the lattice equivalent of the bridged T arrangement of the E23 repeater).



(a)   $Z_H$ = IMAGE IMPEDANCE OF LINE + REPEATER

$P\ell$ = PROPAGATION CONSTANT OF LINE + REPEATER

(b)   $Z_A = \dfrac{N(p)}{D(p)}$   WHERE $D(p)$ HAS NO ZEROS IN THE RIGHT HALF $p$–PLANE

$Z_B = \dfrac{M(p)}{B(p)}$   WHERE $M(p)$ HAS NO ZEROS IN THE RIGHT HALF $p$–PLANE

Fig. 1 — Symmetrical two-wire repeater in transmission line. (a) Schematic. (b) Equivalent circuit.

## 2. GENERAL STABILITY CRITERIA

Before the specific objective is considered, stability criteria will be reviewed briefly and a stability theorem applicable to symmetrical linear four-pole networks will be described.

### 2.1. *Basic Stability Equation*

The stability of the network of Fig. 1(a) can be determined from an examination of the roots in the complex frequency plane of the equation:

$$1 - \left[\frac{Z_H - Z_a}{Z_H + Z_a}\right] \cdot \left[\frac{Z_H - Z_b}{Z_H + Z_b}\right] e^{-2P\ell} = 0 \qquad (1)$$

where:

$Z_H$ = Image impedance of the over-all line

$P\ell$ = Propagation constant of the over-all line

$Z_a$ = Impedance of one line termination

$Z_b$ = Impedance of the other termination

The quantity on the left hand side of the equation is the reciprocal of the interaction factor[2] and its use as a measure of stability has been discussed by F. B. Llewellyn.[3]

As pointed out by Llewellyn Eq. (1) bears a striking similarity to the famous Nyquist equation for stability of feedback amplifiers, usually written

$$(1 - \mu\beta) = 0$$

In both cases the fundamental requirement for stability is that the equations should have no roots in the right half complex frequency plane.

In specific cases, the Nyquist criterion for stability can be applied by plotting on the complex plane as a function of real frequency the factors in (1) which correspond to $\mu\beta$, and seeing whether the plot encircles the point $(1, j0)$.

In general, however, the fact that the point $(1, j0)$ is outside such a plot is not in itself proof of stability. This ambiguity in the interpretation of the diagram can be resolved if the factors involved in (1) are evaluated at complex frequencies.

It should be noted that a separate plot at real frequencies would be required for each combination of terminating impedance $Z_a$ and $Z_b$. The assumption of particular values for $Z_a$ and $Z_b$ would naturally lead to specialized stability criteria and it was to avoid this that Llewellyn

made the alternative assumption that the system should be stable with any combination of passive terminating impedances. Since in a practical telephone system the network of Fig. 1(a) has to be stable when $Z_a$ and $Z_b$ are arbitrary passive impedances, Llewellyn's results are applicable to the cases considered herein. However, since his criterion is stated in terms of the image impedances and loss of the network, it is not in a form which can be readily applied in a design problem involving negative impedance loading.

For design purposes, what is required is a relationship between stability, the properties of the physical line and the negative impedance repeaters. This relationship can be found directly by means of the bisection theorem given in the following section.

## 2.2. A Stability Theorem For Active Four Poles

The symmetrical network of Fig. 1(a) is a particular case of the somewhat more general type of symmetrical structure shown on Fig. 2(a), to which the theorem to be discussed in this section applies.

Referring to Fig. 2(a), $N$ is a symmetrical network in the sense that its external characteristics are such as to make the terminal pairs $(1, 1')$ and $(2, 2')$ electrically indistinguishable. For example, $N$ may be a symmetrical T network with fairly obvious symmetry or it may be a two-way repeater with somewhat less apparent structural symmetry.

For simplicity, the theorem will be stated in terms of the network in Fig. 2(b) which has complete structural symmetry in the sense that the networks $N_1$ and $N_2$ are the mirror images of each other in the plane of symmetry AB. In this case, the open and short circuit impedances of the bisected network are the impedances looking into the terminals $(1, 1')$ or $(2, 2')$ with the terminals in the plane of symmetry AB respectively open and short circuited.

To apply the theorem in the more general case of Fig. 2(a), the open and short circuit impedances of the bisected network should be interpreted as the impedances of the series and shunt arms of the lattice network which is electrically equivalent to $N$. Methods of determining these impedances by external measurements on the network are discussed by Bode.[4]

In a particular application of the theorem in this paper, the symmetrical network consists of a transmission line of length $\ell$ with an active network at its center in the form of a lattice structure. This situation is shown in Fig. 1(b) and it is a slightly more complicated form of symmetry than the simple mirror image symmetry of Fig. 2(b). The situation at the middle of this network is indicated in Fig. 2(c) and it can be shown

that the arm impedances of the over-all equivalent lattice are given in the manner indicated in the figure. From what has been said above, these impedances will also be the $Z_{Short}$ and $Z_{Open}$ of the theorem.

With regard to the terminating impedances $Z_a$ and $Z_b$ of Fig. 2, the theorem is based on the assumption that the network must be stable when these impedances assume any arbitrary passive values. This is also the requirement in the transmission line problem considered in this paper.

*Statement Of The Theorem*

A necessary and sufficient condition for a structurally symmetrical linear four-pole to be stable with any combination of passive terminating



$$Z_{SHORT} = Z_{IN} \text{ WITH 1 AND 4 SHORTED AND 2 AND 3 OPEN}$$

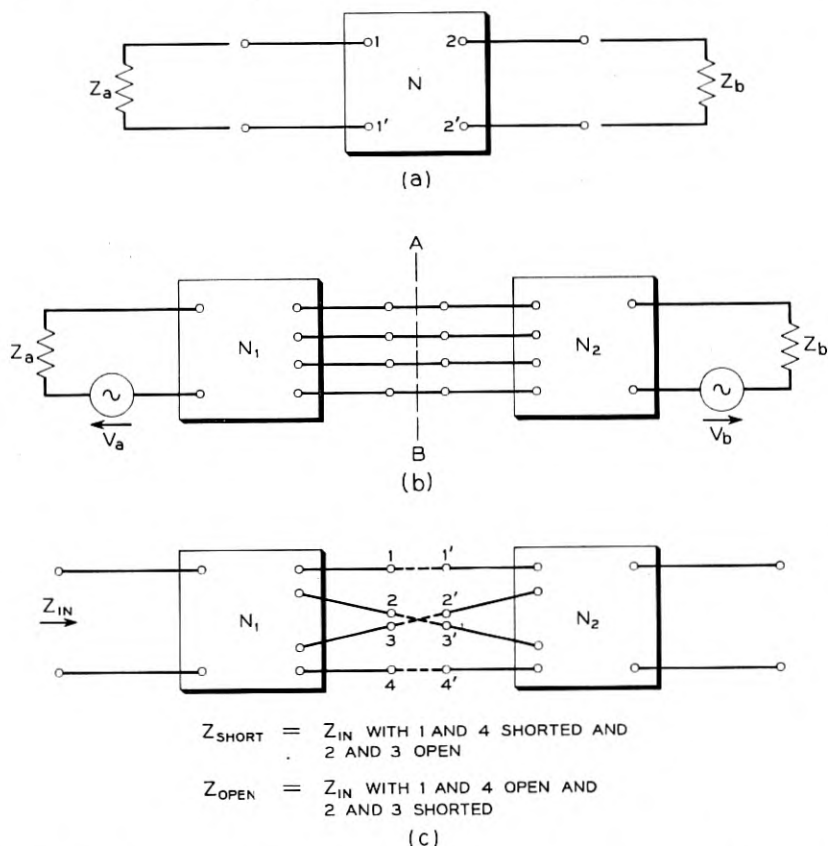$$Z_{OPEN} = Z_{IN} \text{ WITH 1 AND 4 OPEN AND 2 AND 3 SHORTED}$$

(c)

Fig. 2 — Symmetrical linear four-pole with passive terminations. (a) General symmetrical network. (b) Simple bisection. (c) Lattice bisection.

impedances is that the open and short circuit impedances of the bisected network shall be positive real impedance functions. These open and short circuit impedances are the input impedances of either half of the network when the terminals in the plane of bisection are respectively open and short circuited. The term passive impedance as used herein denotes a positive real impedance function.

The theorem is proved in Appendix A and therein also are found the requirements for an impedance function to be positive real.

### 3. SERIES NEGATIVE IMPEDANCE LOADING

This is the case of Fig. 1(b) where $Z_B$ is infinite and where $Z_A$ is a negative impedance of the open circuit stable type. It also represents the installation of an E1 or E2 repeater in the center of an electrically short transmission line.

### 3.1. *Stability*

Consider Fig. 3(a) where a negative impedance of the open circuit stable type ($Z_A$) is shown in the center of a physical transmission line. The problem is to determine the equations which relate transmission characteristics with stability for all passive terminating impedances.

According to the bisection theorem stated in 2.2 the network of Fig. 3(a) will be stable for all passive impedance terminations provided the open and short circuit impedances of the bisected network are positive real. The short circuit impedance of the bisected network is shown in Fig. 3(b) and is represented by $Z_H$ multiplied by Tanh $P\ell/2$. This must be made positive real. The open circuit impedance of the bisected network is shown in Fig. 3(c) and is expressed as $Z_H$ divided by Tanh $P\ell/2$. This open circuit impedance is positive real because it equals the open circuit impedance of one half of the physical line, $Z_{OC}$. Thus it has no direct bearing on stability but does contribute the relationship:

$$\frac{Z_H}{\text{Tanh } P\ell/2} = \frac{Z_o}{\text{Tanh } \gamma\ell/2} = Z_{OC} \tag{2}$$

where:

$Z_H$ = Image impedance of the line with $Z_A$

$P\ell$ = Propagation constant of the line with $Z_A$

$Z_o$ = Characteristic impedance of the physical line

$\gamma\ell$ = Propagation constant of length $\ell$ of physical line

Fig. 3 — Application of bisection theorem to series loading. (a) Schematic. (b) Short circuit impedance of bisected network. (c) Open circuit impedance of bisected network.

Equation (2) demonstrates an important relationship which has been known ever since the discovery of coil loading. It is worthwhile repeating here because the network in Fig. 3(a) is, in fact, a single section of line loaded with a series impedance $Z_A$. Equation (2) demonstrates that the midsection impedance and propagation constant of the loaded line bear the same relationship to each other as the corresponding parameters of the nonloaded line bear to each other. Thus the general relationship between propagation constant and midsection impedance of a loaded line is to this extent independent of the loading element.

With regard to stability, the application of the bisection theorem to Fig. 3(a) has shown that the basic criterion for stability is that $Z_H$ Tanh $P\ell/2$ must be a positive real impedance function. This impedance can be expressed in terms of physical line parameters and $Z_A$ as follows:

$$Z_H \text{ Tanh } P\ell/2 = Z_o \text{ Tanh } (M + \gamma\ell/2) \tag{3}$$

where

$$M = \text{Tanh}^{-1} \frac{Z_A}{2Z_0}$$

$Z_0$ = Characteristic impedance of the physical line

There are two requirements which must be placed upon $Z_0$ Tanh $(M + \gamma \ell/2)$ for it to be positive real. One of them is that in the following equation, $R$ shall be a positive resistance at all frequencies.

$$Z_0 \text{ Tanh } (M + \gamma \ell/2) = R \pm jX \qquad (4)$$

The other requirement will be reserved until (4) has been discussed.

The limit of stability will be approached as $R$ approaches zero. If (4) is taken to the limit of stability then:

$$Z_0 \text{ Tanh } (M + \gamma \ell/2) = \pm jX \qquad (5)$$

where as before

$$M = \text{Tanh}^{-1} \frac{Z_A}{2Z_0}$$

If at a single frequency all values of $Z_A/2$ which satisfy (5) are plotted on the $Z$-plane, their locus will trace a circle because $jX$ is a straight line and the relationship is a bilinear transformation.[5] Formulas for the centers and radii of these circles are given in Appendix B. If, as is the case with telephone cable, the characteristic impedance $Z_0$ has a negative angle, this trace will appear as shown in Fig. 4. The region inside this
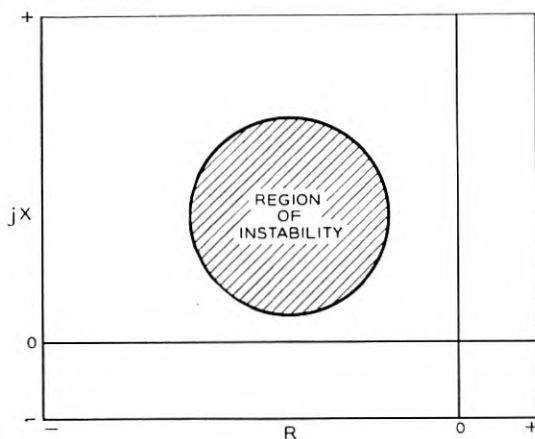


Fig. 4 — Circle for determining the stability of a loading element in a line.

circle corresponds to negative values of $R$ in (4). For stability, the negative impedance $Z_A/2$ must lie outside this circle. However, this is only a part of the stability criterion.

The other requirement on $Z_o$ Tanh $(M + \gamma\ell/2)$ for it to be positive real is discussed in Appendix A3. In order to apply this second requirement it is necessary to expand $Z_o$ Tanh $(M + \gamma\ell/2)$ in terms of the circuit parameters as follows:

$$Z_o \text{ Tanh } (M + \gamma\ell/2) = Z_{oc} \left[ \frac{\dfrac{Z_A}{2} + Z_{sc}}{\dfrac{Z_A}{2} + Z_{oc}} \right] \tag{6}$$

where

$$Z_{sc} = Z_o \text{ Tanh } \gamma\ell/2$$

$$Z_{oc} = \frac{Z_o}{\text{Tanh } \gamma\ell/2}$$

Then from Appendix A3 it can be seen that stability will be obtained providing the real part of $(Z_A/2) + Z_{sc}$ is a positive resistance and providing (4) is satisfied. Thus, the second requirement for stability is that the magnitude of the real part of $Z_A/2$ at real frequencies shall not exceed the real part of $Z_{sc}$, the short circuit impedance of $\ell/2$ of physical line.

This second requirement while sufficient is not necessary as will be discussed later. However, in many practical applications it is not unduly restrictive.

The graphical meaning of the stability requirements can be seen from Fig. 5. Here a family of stability circles has been drawn for a cable circuit at three different frequencies. As the attenuation increases with frequency the circles decrease in size. At each frequency, $-Z_{sc}$ is shown. It falls on the circumference of the corresponding circle. As the frequency increases, $-Z_{sc}$ will rotate clockwise and at some frequency will be on the left edge of the circle. The first requirement for stability means that at any given frequency $Z_A/2$ cannot lie inside the corresponding circle. The second requirement for stability, means that $Z_A/2$ must lie to the right of $-Z_{sc}$ on Fig. 5 at all frequencies. It also means that the trace of $Z_A/2$ over the frequency range cannot enclose this family of circles. The locus on the Z-plane will be similar to that shown on Fig. 5. This is difficult to show graphically on a single plane because it may

Fig. 5 — Stability circles for $Z_A/2$.

appear that the trace of $Z_A/2$ does go through the family of circles. However, where it apparently passes through the family of circles it does so at a lower frequency. At no given frequency does the trace of $Z_A/2$ lie within its stability circle.

The second requirement for stability, namely that the real part of $(Z_A/2) + Z_{SC}$ should be positive, is a sufficient condition but not a necessary one. A necessary and sufficient condition, when $Z_A/2$ is open circuit stable, is derived in Appendix A3. With reference to Fig. 5 this condition requires that the plot of $Z_A/2$ will not enclose any of the stability circles.

### 3.2. *Applications*

Two important classes of problems can be solved by applying the stability considerations of the previous section:

(a) The determination of the lowest value of attenuation possible with stability for all passive impedance terminations if the image impedance $Z_H$ is restricted by over-all system requirements;

(b) The determination of $Z_A$ and the allowable variations in it consistent with stability.

The problem of determining the lowest attenuation in a loading section like that of Fig. 3(a) consistent with stability for all passive impedance terminations when the impedance $Z_H$ is given can be solved by the following method.

It has been shown in Section 3.1 that the first requirement for stability is that the real part of the short circuit impedance (shown on Fig. 3 as $Z_H$ Tanh $P\ell/2$) shall have a positive real part at all real frequencies. By substituting for Tanh $P\ell/2$ its value from (2) this requirement may be written as

$$\frac{Z_H^{\ 2}}{Z_{oc}} = R \pm jX \tag{7}$$

where $R$ is a positive resistance at all frequencies.

As $R$ approaches zero from positive values, the limit of stability will be approached and in the limit

$$\frac{Z_H^{\ 2}}{Z_{oc}} = \pm jX \tag{8}$$

Likewise, by substituting values from (2), (8) can be expressed in the alternate form as

$$Z_{oc} \text{ Tanh}^2 P\ell/2 = \pm jX \tag{9}$$

If the phase as well as the magnitude of $Z_H$ is to be specified in the problem, it should be noted that $Z_H$ must be specified as a passive impedance. This is necessary because $Z_H$ is the square root of the product of the open and short circuit impedances of the bisected network. Since these separately must be positive real, $Z_H$ must also be positive real. This means that $Z_H^{\ 2}/Z_{oc}$, which equals the short circuit impedance of the bisected network, will have no roots in the right half, complex frequency plane. To check for stability in this case all that is required is to insure that (7) is satisfied for the length of the physical line selected. Section length enters into (7) because $Z_{oc}$ is the open circuit impedance of one half this length of line. If the system proved unstable by this check

the only recourse would be to solve for a section length which would prove stable.

If the magnitude only of $Z_H$ had been fixed then for each section length there would be a choice of phase angle for $Z_H$ within the limits prescribed by system requirements. For solution in this case (8) could be interpreted to mean that for stability

$$\text{the angle of } \frac{Z_H^2}{Z_{oc}} < 90 \text{ degrees} \tag{10}$$

Eq. (2) can be rearranged as follows:

$$\text{Tanh } P\ell/2 = \frac{Z_H}{Z_{oc}} \tag{11}$$

From (10) and (11) it may be shown that when the magnitude of $Z_H$ is fixed, the closer the loading section is brought to instability by adjusting the angle of $Z_H$, the lower will be the attenuation of the section. Thus this angle should be adjusted for minimum stability.

When the limit of stability is reached, (8) indicates that $Z_H^2/Z_{oc}$ will be a pure reactance but the question arises as to whether the angle of $Z_H$ should be adjusted to make the sign of the reactance positive or negative. The choice here is based on two observations. First, that if the angle of $Z_{oc}$ is negative, as it generally is for cable circuits, minimum overall attenuation will result when the reactance on the right hand side of (8) is positive. Second, from (6) it may be shown that if the short circuit impedance of the bisected network is a positive reactance, $Z_A/2$ will lie on the circumference of the stability circles of Fig. 5 on the arc to the right of $-Z_{sc}$. In the case of lines such as are under consideration here, where the top frequency in the pass band is less than one quarter wave length, the second stability requirement which is discussed in Section 3.1 thereby will be satisfied within the pass band of frequencies.

Therefore, when the magnitude alone of $Z_H$ has been specified the design procedure is to select an angle for $Z_H$ such that the angle of $Z_H^2/Z_{oc}$ is just less than $+90$ degrees. In (8), $Z_H^2/Z_{oc}$ then will be nearly equal to $+jX$ and the magnitude of $X$ will be given by the ratio of the magnitude of $Z_H^2$ to the magnitude of $Z_{oc}$. The magnitude and angle of $Z_H$ now is determined for any chosen frequency in the band and the value of the negative impedance $Z_A/2$ required to give the desired value of $Z_H$ can be obtained from the equation:

$$Z_A = 2Z_{oc} \frac{Z_o^2 - Z_H^2}{Z_H^2 - Z_{oc}^2} \tag{12}$$

where $Z_o$ is the characteristic impedance of the nonloaded line.

The value found for $Z_A$ may turn out to be unrealizable as such. It will have to be synthesized and a compromise made. The practical value can be checked for stability by the graphical method explained in Section 3.1. Here any correction for stability can be made which might be necessary. With the realizable value of $Z_A$, the final values for attenuation and also for $Z_H$ then will have to be found.

The value of $Z_H$ can be obtained from the following equation.

$$Z_H = Z_{oc} \sqrt{\frac{\frac{Z_A}{2} + Z_{sc}}{\frac{Z_A}{2} + Z_{oc}}} \qquad (13)$$

The propagation constant $P$ can be found from

$$\text{Tanh } P\ell/2 = \sqrt{\frac{\frac{Z_A}{2} + Z_{sc}}{\frac{Z_A}{2} + Z_{oc}}} \qquad (14)$$

The attenuation per section with $Z_A/2$ included can be determined from the following equation wherein the logarithmic rather than the hyperbolic form has been used.

$$\text{Attenuation (db per section)} = 20 \log_{10} \left| \frac{1 + \frac{Z_H}{Z_{oc}}}{1 - \frac{Z_H}{Z_{oc}}} \right| \qquad (15)$$

and the angle of $\left[ 1 + \frac{Z_H}{Z_{oc}} \middle/ 1 - \frac{Z_H}{Z_{oc}} \right]$ in degrees is the phase shift of the loaded line section.

The class of problem represented by (b) at the head of this section can be solved by the graphical method for determining the stability of $Z_A/2$ as outlined in Section 3.1.

### 3.3. *Characteristics*

This type of loading has several interesting characteristics. First, zero attenuation over any frequency band, no matter how narrow, when $Z_A/2$ is adjusted to the limit of stability as defined by (5) and (6) is unrealizable. This can be seen from (14) together with the fact that $Z_{sc}$ and $Z_{oc}$ represent passive impedances and that $Z_A/2$ is a negative impedance of the open circuit stable type. All three impedances when shown on the impedance plane with increasing frequency will rotate

in the clockwise direction. For stability, the locus of $Z_A/2$ must not enclose either $-Z_{SC}$ or $-Z_{OC}$. At the limit of stability, however, $Z_A/2$ must lie on the circles of stability which are shown on Fig. 5. This means a changing relationship between these three impedances with frequency which is incompatible with zero attenuation over any frequency band. Second, a flat response can be realized over a band of frequencies, in general, only at the expense of increasing the loss at the lower frequencies above that required for stability. Third, as the length of the physical line is increased, it becomes more and more difficult to obtain a low over-all loss and yet avoid the enclosure of the stability circles of Fig. 5 with a realizable design of $Z_A/2$. The practical limit here appears to be one quarter wave length of physical line at the highest frequency it is desired to pass.

## 4. SHUNT NEGATIVE IMPEDANCE LOADING

This case where $Z_B$ is shunted across the line conductors at the mid-point of a physical line (Fig. 1(b) where $Z_A$ is zero) can be classed as shunt type negative impedance loading. The negative impedance $Z_B$ is of the short circuit stable type such as the impedance produced by the E3 repeater. Hence, this case can represent an E3 repeater bridged across the conductors of an electrically short transmission line.



Fig. 6 — Application of bisection theorem to shunt loading. (a) Schematic. (b) Short circuit impedance of bisected network. (c) Open circuit impedance of bisected network.

4.1. *The Stability Equations*

The same method described in detail in Section 3.1 is used here to determine stability. The short and open circuit impedances of the bisected network are obtained as shown in Fig. 6.

As seen from Fig. 6(b) the short circuit impedance of the bisected network is positive real being equal to $Z_{sc}$, the short circuit impedance of the physical line of length $\ell/2$.

$$Z_H \operatorname{Tanh} P\ell/2 = Z_o \operatorname{Tanh} \gamma\ell/2 = Z_{sc} \tag{16}$$

The open circuit impedance of the bisected network [Fig. 6(c)] determines stability. Thus for stability:

$\dfrac{Z_H}{\operatorname{Tanh} P\ell/2}$ must be a positive real impedance function.

A substitution from (16) for $\operatorname{Tanh} P\ell/2$ above will yield the following requirement for stability:

$$\frac{Z_H^2}{Z_{sc}} = R \pm jX \tag{17}$$

where $R$ must be a positive resistance and $Z_H$ a positive real impedance function.

The limit of stability is reached as $R$ goes to zero. Therefore, the limit of stability for all passive impedance terminations can be expressed by:

$$\frac{Z_H^2}{Z_{sc}} = \pm jX \tag{18}$$

A similar equation can be obtained in terms of $\operatorname{Tanh} P\ell/2$ rather than $Z_H$.

$$\frac{Z_{sc}}{\operatorname{Tanh}^2 P\ell/2} = \pm jX \tag{19}$$

4.2. *Analysis*

From Eq. (18) the sole criterion for stability is that the angle of

$$\frac{Z_H^2}{Z_{sc}} < 90° \tag{20}$$

provided that $Z_H$ is a positive real impedance function.

The value of the propagation constant can be found from (16) to be:

$$\operatorname{Tanh} P\ell/2 = \frac{Z_{sc}}{Z_H} \tag{21}$$

and from the short circuit and open circuit impedances, Fig. 6(b) and Fig. 6(c), the following are found

$$\text{Tanh } P\ell/2 = \sqrt{\frac{Z_{sc}(2Z_B + Z_{oc})}{Z_{oc}(2Z_B + Z_{sc})}} \tag{22}$$

and

$$Z_H = \sqrt{Z_{sc}Z_{oc}\left[\frac{2Z_B + Z_{sc}}{2Z_B + Z_{oc}}\right]} \tag{23}$$

In order to translate stability into engineering parameters, the open circuit impedance, $Z_H/\text{Tanh } P\ell/2$, can be expressed in parameters of the physical line and $Z_B$ .

$$\frac{Z_H}{\text{Tanh } P\ell/2} = Z_{oc}\left[\frac{2Z_B + Z_{sc}}{2Z_B + Z_{oc}}\right] \tag{24}$$

By reasoning similar to that used in Appendix A and in Section 3.1 it can be shown that the requirement for stability will be met if

$$Re\left(\frac{1}{2Z_B} + \frac{1}{Z_{oc}}\right) \text{ is positive} \tag{25}$$

where $Z_B$ is a negative impedance of the short circuit stable type and providing

$$Z_{oc}\left[\frac{2Z_B + Z_{sc}}{2Z_B + Z_{oc}}\right] = R \pm jX \tag{26}$$

where $R$ is a positive resistance at all frequencies.

The limit of stability will be reached as $R$ approaches zero.

If at any given frequency all values of $2Z_B$ , which fulfill (26) when $R$ is zero, are plotted on the $Z$-plane they will trace out the same circle as found before in Section 3.1 where $Z_A$ was likewise examined (if the frequency and physical line parameters are the same in both cases). Thus $2Z_B$ should not lie inside the stability circle of Fig. 4.

However, in order to meet the restriction imposed by (25) upon $Z_B$ , the locus of $2Z_B$ when plotted over the frequency range from zero to infinity must enclose this family of circles as shown in Fig. 7. This can be established in much the same way as in Section 3.1 where it was proved that the locus of $Z_A/2$ must *not* enclose this family of circles.

Here, in the case of shunt loading, as in the case of series loading, zero attenuation is inconsistent with stability for all passive impedance terminations. Likewise, with shunt loading designed for minimum stability over the pass band the attenuation will vary with frequency.

Fig. 7 — Stability circles for $2Z_B$ .

## 5. LATTICE LOADING

The general case of Fig. 1(b) is where the combination of an E2 and an E3 repeater is located in the center of an electrically short line. Although the actual E23 repeater is connected as a bridged T arrangement[1] the equivalent lattice form is used herein for simplicity in explanation. What is said as applied to the lattice structure applies also to the bridged T.

Fig. 8(a) shows a lattice network of negative impedances connected at the midpoint of a line of length $\ell$. Negative impedance $Z_A$ is open circuit stable; $Z_B$ is short circuit stable.

The short circuit impedance $Z_H$ Tanh $P\ell/2$ is shown in Fig. 8(b). It is the same as that shown in Fig. 3(b) for the case where $Z_A$ is used alone. Thus a requirement for stability in the limit is that

$$Z_{oc}\left[\frac{\dfrac{Z_A}{2} + Z_{sc}}{\dfrac{Z_A}{2} + Z_{oc}}\right] = \pm jX \qquad (27)$$

where $Re\ [(Z_A/2) + Z_{sc})]$ is positive.

Stability can be determined exactly as explained in Section 3.1.

The open circuit impedance $Z_H/$Tanh $P\ell/2$ is shown in Fig. 8(c). It is the same as the case where $Z_B$ is used alone. Thus the limiting requirement for stability is that

$$Z_O\left[\frac{2Z_B + Z_{sc}}{2Z_B + Z_{oc}}\right] = \pm jX \qquad (28)$$

where

$$Re\left(\frac{1}{2Z_B} + \frac{1}{Z_{oc}}\right) \text{ is positive.}$$



Fig. 8 — Application of bisection theorem to lattice loading. (a) Schematic. (b) Short circuit impedance of bisected network. (c) Open circuit impedance of bisected network.

Stability can be determined as explained in Section 4.2 for (26).

Thus stability with $Z_A$ is determined independently of $Z_B$ and the converse is true also. In regard to stability each negative impedance can be designed without considering the other. The image impedance, $Z_H$, and propagation constant, $P\ell$, of the resulting line [Fig. 8(a)] will depend upon both $Z_A$ and $Z_B$, however.

Equations for the image impedance $Z_H$ and the propagation constant of the repeatered line can be expressed as follows:

$$Z_H = Z_{oc} \sqrt{\frac{\left[\dfrac{Z_A}{2} + Z_{sc}\right][2Z_B + Z_{sc}]}{\left[\dfrac{Z_A}{2} + Z_{oc}\right][2Z_B + Z_{oc}]}} \tag{29}$$

and

$$P\ell/2 = \operatorname{Tanh}^{-1} \sqrt{\frac{\left[\dfrac{Z_A}{2} + Z_{sc}\right][2Z_B + Z_{oc}]}{\left[\dfrac{Z_A}{2} + Z_{oc}\right][2Z_B + Z_{sc}]}} \tag{30}$$

From what has been said in the preceding sections it should be evident that when both $Z_A$ and $Z_B$ are designed to the limit of stability in a telephone cable section the image impedance and propagation constant will be as follows:

$$Z_H = \sqrt{(+jX_1)(-jX_2)} \tag{31}$$

and

$$\operatorname{Tanh} P\ell/2 = \sqrt{\frac{+jX_1}{-jX_2}} \tag{32}$$

where

$$+jX_1 = Z_0 \operatorname{Tanh} (M + \gamma\ell/2)$$

$$M = \operatorname{Tanh}^{-1} Z_A/2Z_0 \qquad \text{at limit of stability}$$

$$-jX_2 = Z_0 \operatorname{Tanh} (N + \gamma\ell/2)$$

$$N = \operatorname{Tanh}^{-1} 2Z_B/Z_0 \qquad \text{at limit of stability}$$

From these last two equations it is apparent that at the limit of stability $Z_H$ is a resistance in the pass band and the attenuation of the repeatered section can be zero and the system be stable for all passive impedance terminations. Furthermore, zero db attenuation can be realized theoretically over the pass band.

## 6. SUMMARY

The stability and transmission characteristics have been outlined for a single section considering three separate systems of negative impedance loading and a summary is shown in Fig. 9. A general practical restriction on the use of these systems is that the negative impedance or impedances shall be located in the center of the section and shall be less than one quarter wave length from the line terminals at the highest frequency it is desired to pass. The condition for stability has been taken that each section must be stable for all passive impedance terminations.

The important features can be outlined as follows.

### 6.1. *Series Negative Impedance Loading*

With a series loading element, $Z_A$, stability is determined solely by the short circuit impedance of the bisected section. The attenuation of the section must be finite for stability. Where the loading section is designed to the limit of stability, the transmission-frequency response will vary with frequency in the pass band; and the magnitude of the image impedance $|Z_H|$ will tend to increase with frequency within this band. A flat transmission-frequency response is possible only at the expense of greater loss at the lower frequencies than is required for stability.

### 6.2. *Shunt Negative Impedance Loading*

With a shunt loading element, $Z_B$, stability is determined solely by the open circuit impedance of the bisected section. The attenuation of the section must be finite for stability. Where the loading element is designed to the limit of stability, the attenuation will vary with frequency in the pass band and increase at frequencies outside the band. The magnitude of the image impedance $|Z_H|$ will tend to decrease in the pass band as the frequency increases.

A flat transmission-frequency response with shunt loading is possible only at the expense of greater loss at the higher frequencies than is required for stability.

### 6.3. *Loading with a Lattice or Equivalent Bridged T Network*

Loading with a lattice network having series arms of $Z_A/2$ and shunt arms of $2Z_B$ both of which are negative impedances, the former open circuit stable, the latter short circuit stable, will have the following characteristics.

| TYPE OF LOADING | SCHEMATIC | AT THE LIMIT OF STABILITY — OVERALL EQUIVALENT CIRCUIT | AT THE LIMIT OF STABILITY — STABILITY REQUIREMENT | TRANSMISSION CHARACTERISTICS (OF LOADED LINE) |
|---|---|---|---|---|
| SERIES | LINE $\frac{\ell}{2}$ — $Z_A$ — LINE $\frac{\ell}{2}$ | $+jX$, $Z_{OC}$, $Z_{OC}$, $+jX$ with $Z_H$ | $Z_{OC}\left[\dfrac{\frac{Z_A}{2}+Z_{SC}}{\frac{Z_A}{2}+Z_{OC}}\right] = +jX$ <br> WHERE: $Re\left(\dfrac{Z_A}{2}+Z_{SC}\right)$ IS POSITIVE | $\tanh P\dfrac{\ell}{2} = \sqrt{\dfrac{\frac{Z_A}{2}+Z_{SC}}{\frac{Z_A}{2}+Z_{OC}}}$ <br> $Z_H = Z_{OC}\sqrt{\dfrac{\frac{Z_A}{2}+Z_{SC}}{\frac{Z_A}{2}+Z_{OC}}}$ |
| SHUNT | LINE $\frac{\ell}{2}$ — $Z_B$ — LINE $\frac{\ell}{2}$ | $Z_{SC}$, $-jX$, $-jX$, $Z_{SC}$ with $Z_H$ | $Z_{OC}\left[\dfrac{2Z_B+Z_{SC}}{2Z_B+Z_{OC}}\right] = -jX$ <br> WHERE: $Re\left(\dfrac{1}{2Z_B}+\dfrac{1}{Z_{OC}}\right)$ IS POSITIVE | $\tanh P\dfrac{\ell}{2} = \sqrt{\dfrac{Z_{SC}}{Z_{OC}}\left[\dfrac{2Z_B+Z_{OC}}{2Z_B+Z_{SC}}\right]}$ <br> $Z_H = \sqrt{Z_{SC}Z_{OC}\left[\dfrac{2Z_B+Z_{SC}}{2Z_B+Z_{OC}}\right]}$ |
| LATTICE | LINE $\frac{\ell}{2}$ — $\frac{Z_A}{2}$, $2Z_B$, $2Z_B$, $\frac{Z_A}{2}$ — LINE $\frac{\ell}{2}$ | $+jX$, $-jX$, $-jX$, $+jX$ with $Z_H$ | REQUIREMENT SAME AS FOR SERIES PLUS THAT FOR SHUNT | $\tanh P\dfrac{\ell}{2} = \sqrt{\dfrac{\frac{Z_A}{2}+Z_{SC}}{\frac{Z_A}{2}+Z_{OC}}\left[\dfrac{2Z_B+Z_{SC}}{2Z_B+Z_{OC}}\right]}$ <br> $Z_H = Z_{OC}\sqrt{\dfrac{\frac{Z_A}{2}+Z_{SC}}{\frac{Z_A}{2}+Z_{OC}}\left[\dfrac{2Z_B+Z_{SC}}{2Z_B+Z_{OC}}\right]}$ |

Fig. 9 — Summary.

The short circuit impedance of the bisected section will determine the stability of the negative impedances in the series arms. The open circuit impedance of the bisected section will determine the stability of the negative impedances in the shunt arms.

Zero attenuation is theoretically possible as a limit in the pass band of frequencies consistent with stability for all passive impedance terminations.

The image impedance for zero attenuation will be a positive resistance

.

## Appendix A

To prove that: A necessary and sufficient condition for a structurally symmetrical linear four-pole to be stable with any combination of passive terminating impedances is that the open and short circuit impedances of the bisected network shall be positive real. These open and short circuit impedances are the input impedances of either half of the network when the terminals in the plane of bisection are respectively open and short circuited.

A.1. *Proof of Necessity*

Consider the network of Fig. 2(b), in the text, representing a linear four-pole which is structurally symmetrical in the sense that the right half of the network is the mirror image of the left half in the plane of symmetry $AB$.

Assuming the network is stable, the necessity of the condition in the theorem will be established if the open and short circuit impedances of the bisected network are shown to be positive real. Stability is used here in the sense that the response to an impressed signal will die out upon removal of the excitation.

For an impedance function $Z(p)$ of the complex frequency variable $p = \alpha + i\omega$ to be positive real it is sufficient to show that the following four conditions are satisfied.[6]

1. $Z(p)$ has no zeros in the right half $p$-plane.

2. Zeros of $Z(p)$ on the boundary of the right half $p$-plane are simple and at them $dZ/dp$ = a positive real constant.

3. The real part of $Z(i\omega) \geqq 0$ for all values of $\omega$.

4. The imaginary part of $Z(p) = 0$ whenever the imaginary part of $p = 0$.

The fourth condition is always satisfied by physical networks and will be assumed true without proof.

To show that condition three is satisfied, consider the determinant $\Delta$ of the entire network in Fig. 2(b) in terms of its open circuit impedance parameters and the arbitrary passive terminations $Z_a$ and $Z_b$

$$\Delta = \begin{vmatrix} Z_{11} + Z_a & Z_{12} \\ Z_{12} & Z_{11} + Z_b \end{vmatrix} \tag{A1}$$

Since the network is stable by hypothesis, $\Delta$ as a function of the complex variable $p$ can have no zeros in the right half of the $p$-plane. Since the definition of stability requires that a response will die out on the removal of the excitation, zeros of $\Delta$ on the boundary of the right half $p$-plane are excluded.

If $Z_a = Z_b = Z$, in Eq. (A1), $\Delta$ may be expanded into the product of two factors as follows

$$\Delta = (Z_{11} + Z - Z_{12})(Z_{11} + Z + Z_{12}) \tag{A2}$$

From what has been said above, neither of the factors in (A2) can have zeros on the imaginary $p$ axis or in other words at real frequencies.

If $Z_{11} = R_{11} + jX_{11}$ or in general if $Z_{rs} = R_{rs} + jX_{rs}$, equation (A2) may be rewritten in the following form.

$$\Delta = [R_{11} - R_{12} + R + j(X_{11} - X_{12} + X)]$$
$$[R_{11} + R_{12} + R + j(X_{11} + X_{12} + X)] \tag{A3}$$

Remembering that $Z$ is an arbitrary passive terminating impedance, $X$ can always be chosen to nullify either of the imaginary parts in the above two factors. Moreover, since neither factor has a root at real frequencies and since $R$ can be given any positive value, it is obviously necessary that their real parts shall be positive, thus:

$$R_{11} - R_{12} + R > 0 \qquad R_{11} + R_{12} + R > 0 \tag{A4}$$

Since $R$, the real part of the terminating impedance, is not negative, the limiting situation in the above conditions will occur when $R = 0$ and it follows that both $R_{11} - R_{12}$ and $R_{11} + R_{12}$ must be positive. It can also be concluded from this that $R_{11}$ is positive, though this is ob-

vious from the fact $R_{11}$ is the resistive component of the open circuit input impedance of a stable network.

The result just established which may be expressed by stating that $R_{11} > |R_{12}|$, is identical to the Gewertz condition for symmetrical linear networks mentioned by F. B. Llewellyn.[3]

From standard network theory, the open and short circuit input impedances of the bisected network of Fig. 2(b) are given by the following equations

$$Z_{\text{Open}} = Z_{11} + Z_{12} \qquad Z_{\text{Short}} = Z_{11} - Z_{12} \qquad \text{(A5)}$$

where

$Z_{\text{Open}} = $ Open circuit input impedance of the bisected network

$Z_{\text{Short}} = $ Short circuit input impedance of the bisected network

By applying the Gewertz condition to (A5) it is clear that the open and short circuit impedances must have positive real parts which establishes requirement 3 for a positive real function. To show that $Z_{\text{Open}}$ and $Z_{\text{Short}}$ satisfy conditions 1 and 2 for a positive real function, set $Z$ equal to zero in equation (A2). This reduces $\Delta$ to the product of $Z_{\text{Open}}$ and $Z_{\text{Short}}$. Since $\Delta$ has no roots inside or on the boundary of the right half $p$-plane this must also be true of $Z_{\text{Open}}$ and $Z_{\text{Short}}$. Hence, they each meet all the requirements for positive real functions which completes the proof of necessity.

## A.2. Proof of Sufficiency

In the proof of sufficiency, $Z_{\text{Open}}$ and $Z_{\text{Short}}$ are assumed to be positive real and it must be shown that the network of Fig. 2(b) is stable when terminated in arbitrary passive impedances.

The proof depends on Bartlett's Bisection Theorem.[4] According to this theorem, a lattice network with arm impedances $Z_{\text{Open}}$ and $Z_{\text{Short}}$ will have exactly the same external characteristics as the symmetrical network of Fig. 2(b). Since $Z_{\text{Open}}$ and $Z_{\text{Short}}$ are assumed to be positive real, the lattice arms can be realized with passive impedances. This means that the lattice network will be stable with any combination of passive terminating impedances and it follows that this must likewise be true of the equivalent circuit of Fig. 2(b).

## A.3. Special Applications of the Theorem

In some special applications of the stability theorem it is only necessary to ensure that the open circuit impedance and the short circuit

impedance of the bisected symmetrical network have positive real parts. This is more lenient than that these impedances shall be positive real as stated in the theorem. The other main condition for positive realness which requires that the roots be located in the left half $p$-plane (complex frequency-plane) will be guaranteed automatically by placing special requirements on some of the network elements.

As an example of such a situation consider the circuit of Fig. 3(a) consisting of a transmission line of length $\ell$ with a negative impedance $Z_A$ located at its center.

If $Z_{OC}$ and $Z_{SC}$ are respectively the open and short circuit impedances of the nonloaded line of length $\ell/2$ the open and short circuit impedances of the same length of line with loading may be written down as follows.

$$Z_{\text{Open}} = Z_{OC} \tag{A6}$$

$$Z_{\text{Short}} = Z_{OC} \frac{\left[\dfrac{Z_A}{2} + Z_{SC}\right]}{\left[\dfrac{Z_A}{2} + Z_{OC}\right]} \tag{A7}$$

In this example, the open circuit impedance of the bisected network $Z_{\text{Open}}$, is obviously positive real since it equals the open circuit impedance of a length $\ell/2$ of nonloaded line, which is passive.

It will now be shown that a sufficient condition for the short circuit impedance $Z_{\text{Short}}$ of the bisected network to be positive real is that its resistive component shall be positive at all frequencies, providing that $Z_A$ is a negative impedance of the open-circuit-stable type having a resistive component which is always less in magnitude than the real part of $2Z_{SC}$.

To show this, assume the real part of $Z_{\text{Short}}$ is positive and that all the impedances being considered are rational. To satisfy the rationality requirement in the case of transmission line impedances, the line may be considered as the limit of a lumped element network. Since $Z_A$ is open circuit stable by hypotheses, it can have no poles inside or on the boundary of the right half $p$-plane. Likewise, since $Z_{OC}$ and $Z_{SC}$ are passive impedances they have neither poles nor zeros in the right half $p$-plane. With these facts in mind, consider the expression on the right hand side of equation (A7). The only zeros which this function can have are those due to $(Z_A/2) + Z_{SC}$. As the complex variable $p$ traces a path around the boundary of the right half complex frequency plane, the impedance function $(Z_A/2) + Z_{SC}$ will trace out a closed curve in the $Z$-plane. Since it has been assumed that the magnitude of the real part of $Z_A/2$ is always less than the real part of $Z_{SC}$ the closed curve will lie entirely in the

right half $Z$-plane and cannot therefore enclose the origin. It follows from complex variable theory[7] and the rational nature of all the impedance functions involved that $(Z_A/2) + Z_{sc}$ must have an equal number of zeros and poles in the right half $p$-plane. Since it has no poles in that area it has no zeros either hence, $Z_{\text{Short}}$ has no zeros in, or on the boundary of the right half $p$-plane.

If, in addition, $Z_{\text{Short}}$ has a positive real part as assumed, requirements 1, 2 and 3 for a positive real function are met. Taking requirement 4 as being true without proof, it follows that $Z_{\text{Short}}$ is a positive real function.

Thus, if $Z_{\text{Short}}$ has a positive real part and $Z_A$ has the properties attributed to it above, the network of Fig. 3(a) will be stable with arbitrary passive terminations.

As mentioned at the beginning of this section, the situation just considered relates to a sufficient condition for stability when the impedance $Z_A$ is open circuit stable.

A necessary condition for stability when $Z_A$ is open circuit stable, may be obtained from an examination of (A7). As pointed out in considering the sufficient condition for stability, the only roots which $Z_{\text{Short}}$ can have in the right half $p$-plane are due to the factor $(Z_A/2) + Z_{sc}$. Since this factor has no poles in the right half $p$-plane it follows from complex variable theory[7] that if $Z_A/2 + Z_{sc}$ is plotted on the $Z$ plane as a function of real frequency, the number of times the plot encircles the origin gives the number of zeros which the function has in the right half $p$-plane. Since $Z_{\text{Short}}$, for stability, can have no zeros in this area, it follows that the plot of $Z_A/2 + Z_{sc}$ must not encircle the origin. This last statement is equivalent to saying that $Z_A/2$ must not encircle $-Z_{sc}$.

If the necessary condition just established, is combined with the stability requirement that $Z_A/2$ cannot enter the stability circles discussed in connection with Fig. 5, it will be seen readily that a necessary and sufficient condition for stability may be laid down as follows.

If in the circuit of Fig. 3(a), $Z_A$ is open circut stable then a necessary and sufficient condition for the transmission line with series loading to be stable is that the plot of $Z_A/2$ as a function of real frequency shall not encircle any of the stability circles associated with the line of length $\ell/2$. These stability circles are shown in Fig. 5.

## APPENDIX B

### EQUATIONS FOR THE STABILITY CIRCLE

At any given frequency the equation for the stability circle can be expressed in the following formulas in terms of the rectangular coordinates $(R, X)$ of the center and the radius $(r)$.

$$R = -\frac{(R_I^2 - X_I^2) + (R_{oc}^2 + X_{oc}^2)}{2\,R_{oc}} \tag{B1}$$

$$X = \frac{X_I R_I}{R_{oc}} \tag{B2}$$

$$r = \sqrt{R^2 + X^2 - \left[R_I^2 - X_I^2 + \frac{2 X_{oc} X_I R_I}{R_{oc}}\right]} \tag{B3}$$

where:

$R_I$ = The resistance component of the characteristic impedance, $Z_o$, of the physical line

$X_I$ = The reactance component of the characteristic impedance, $Z_o$, of the physical line

$R_{oc}$ = The resistance component of the open circuit impedance of $\ell/2$ of physical line

$X_{oc}$ = The reactance component of the open circuit impedance of $\ell/2$ of physical line

REFERENCES

1. J. L. Merrill, A. F. Rose, and J. O. Smethurst, Negative Impedance Telephone Repeaters, B.S.T.J., **33,** pp. 1055–1092, Sept., 1954.
2. K. S. Johnson, Transmission Circuits for Telephone Communication, p. 141, D. Van Nostrand Co. Inc., N. Y.
3. F. B. Llewellyn, Some Fundamental Properties of Transmission Systems, Proc I.R.E., **40,** pp. 271–273, March, 1952.
4. Hendrik W. Bode, Network Analysis and Feedback Amplifier Design, p. 268, D. Van Nostrand Co. Inc., N. Y.
5. S. A. Schelkunoff, Applied Mathematics for Engineers and Scientists, p. 25, D. Van Nostrand Co. Inc., N. Y.
6. Otto Brune, Journal of Mathematics and Physics, **10,** pp. 191–235, 1930–31.
7. E. C. Titchmarsh, The Theory of Functions (Second Edition), p. 115, Oxford University Press.

# A Telephone Switching Network and Its Electronic Controls

## By S. T. BREWER and G. HECHT

*A high-speed telephone switching network with electronic controls is described. The high speed allows one-at-a-time operation, reducing both the number and complexity of control circuits. Germanium diodes, transistors, cold-cathode tubes and fast relays perform the control functions; glass-sealed relay contacts provide paths for speech. A laboratory model of the network has performed reliably for over three years.*

## INTRODUCTION

In considering the value of speed in a telephone switching network, first let us take the subscriber's viewpoint. It is obvious that any increase in speed of establishing connections gives the subscriber better service. The actual improvement in service, however, turns out to be very modest. For example, if the subscriber saves a half second on each connection through the switching network, he saves one second over-all on the originating and terminating connections of a local call. With the present method of dailing, dialing time tends to obscure any saving in time due to a faster switching network.

Actually, the chief goal in seeking higher switching and control speeds is the reduction in the cost and complexity of the controls. This may be illustrated by citing an example from the present switching art. Consider a large No. 5 crossbar switching network. To handle the desired traffic, six markers may be required. Each marker is capable of controlling the switching network, but due to its low operating speed (and the slow speed of the network's crossbar switches), a single marker cannot establish all of the originating and terminating connections required during a busy hour. In addition to the duplication of markers, elaborate connectors are required so that various markers can be associated with different portions of the switching network. The markers and marker connectors furthermore require lockout protection so that two markers cannot at-

tempt to control the same portion of the switching network at the same time.

Now consider a switching network and controls much faster than those presently used. Such a system could handle a 10,000-line central office on a "one-at-a-time basis" without requiring duplication of the controls. Thus, the connector problem is almost eliminated and there is no necessity for lockout among competing controls. The switching network and controls which we describe in this paper constitute a model of this type of high-speed system. In addition to reducing the system to one set of control circuits we have devised circuits which perform terminal selection and channel control functions in a greatly simplified fashion.

GENERAL OBJECTIVES

In any project such as this, certain general objectives are established at the outset. These objectives guide the designer in his choice of tools or devices and act as a framework within which the designer works. In our particular case, these general objectives were:

(a) Speed of network and controls should permit operation with one set of control circuits.

(b) The design should be based on a large central office carrying heavy traffic, specifically, an office with 10,000 subscribers and 2,000 trunks, and handling a total of 50,000 originating and terminating calls per busy hour.

(c) The control circuits should be active during 70 per cent or less of a busy hour. (This limit was set to avoid excessive delay in gaining access to the control circuits.)

(d) A minimum of maintenance should be required.

(e) Circuits should function reliably throughout a 40-year life.

(f) The switching network should operate with the present type of subscriber instrument and existing methods of ringing, dialing, and supervision.

(g) The switching network and controls should be incorporated in a functionally complete but skeletanized high speed telephone system to provide a practical working test of the system and all of its components.

DEVICES USED IN DESIGN

General objectives (a), (b), and (c) dictated that we use high-speed devices as the backbone of our design. From these objectives and the assumption that the controls are used twice on a particular connection (once to establish and once to release the connection), we calculate that
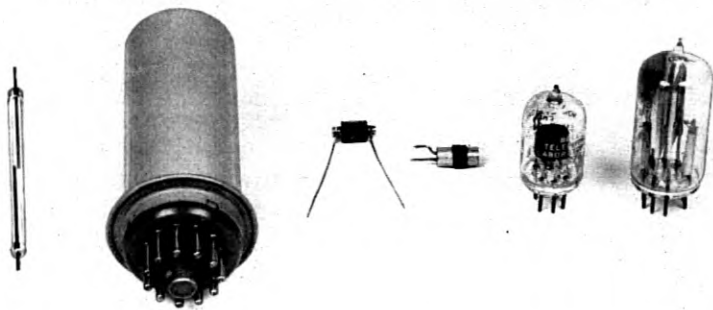
FIG. 1 — Typical high-speed devices used in design. From left to right: glass-sealed reed switch; mercury contact (3-transfer) relay, point contact germanium diode, point contact germanium transistor, cold-cathode diode, grid-controlled cold-cathode tetrode.

the maximum overall time for the connect and release operations is 50 milliseconds. Thus, crossbar switches and relays of conventional design are too slow to meet our minimum objectives. However, objective (f) indicated the use of a metallic path through the switching network. These objectives and objectives (d) and (e) regarding reliability and low maintenance led to the selection of the dry reed switch for use in the switching network. In control circuits subject to heavy use and requiring a metallic path, mercury contact relays were selected. Where circuits required electron tubes, cold cathode diodes, triodes, and tetrodes were used. These are simple and rugged and have the big advantage that they do not consume power or deteriorate during stand-by periods.

Germanium diodes were selected for the circuits performing code translation because they are simple and fast. Transistors were chosen for the channel selecting circuit. This choice was made to gain design and operating experience with a relatively new tool. The basic high speed devices used in our design are shown in Fig. 1.

GENERAL DESIGN CONSIDERATIONS

A high-speed switching circuit can be designed by one-to-one substitution of new devices in existing designs. In general, if the new devices are very much faster or different than the old, this method fails to take full advantage of the capabilities of the new tools. Hence, we have used the method of designing "from scratch." It is beyond the scope of this article to detail all of the advantages obtained by the second and more complete approach to design. However, a result which may be considered typical

is the evolution of a new method of channel control which distinguishes the system.

## THE SERIES MARKING METHOD OF CHANNEL CONTROL

The general philosophy behind the series marking method of channel control may be understood by comparing it with a present method. As an example of a present method, let us take that used in a four-stage No. 5 crossbar switching network. For a channel between a particular line and trunk to be useful its A, B, and C links must be idle. Leads associated with these links are carried through individual contacts on connector relays to each of the markers in the central office. In one marker, a channel testing relay performs a busy test on the A, B, and C links comprising each channel. Those channels which are found to be completely idle compete in a relay lockout circuit which selects one of the channels for use. This is what we call a parallel method of switching network control. It can be seen that every marker in the office requires individual access to each A, B, and C link of the switching network.

In the series marking method of control which we use, marking voltages are applied to the line and trunk. Enabling voltages are applied to the A and C links through high resistance. If a particular A or C link is associated with a busy path, the holding circuit, being low in impedance, will determine the link voltage and block the enabling signal. In this simple manner, busy paths are automatically isolated and excluded from the path marking operation. The voltage marks then extend fan-wise over idle paths from the selected line and trunk toward the center of the switching network. In the center of the switching network, these marks are passed through a connector relay to the mactors.* There, a simple coincidence circuit matches the voltage received from the line with that received from the trunk side of the switching network. If a match is obtained, it indicates that the complete line to trunk path is idle. In this simple manner, which requires no individual access to the A and C links, each mactor determines whether its associated line-to-trunk channel is idle.

## FUNCTIONAL RELATIONSHIP OF SWITCHING NETWORK AND CONTROLS

With this preamble of our objective and philosophy, let us now consider the switching network and controls which we evolved to meet our requirements. To demonstrate the feasibility of the electronic controls

---

* The mactor is a new type of control circuit which takes its name through abbreviation of the functionally descriptive term, MAtcher-seleCTor-connectOR.

and the speed and reliability of the circuits, the switching network and controls were skeletonized and built as an integral part of an automatic telephone system known as the DIAD system.* The other parts of the system have been described earlier; furthermore, one can follow the operation of the switching network and its controls as separate entities. For these reasons we will consider the system operation only from the time the switching network receives an order to establish or release a connection until the time the controls answer, "Have taken the desired action," or "Could not perform the desired action."

A general idea of the physical relation of the various circuits and the type of construction used may be obtained from Fig. 2, which is a front view of the two switching bays. In addition to power supplies, the left bay contains test and monitor panels and the three line frames. The right bay contains three power supplies, the two trunk frames, the mactors, the switching control or sequence circuit and the switching number groups. A rear view of these two bays is shown in Fig. 3. This view gives a general picture of the wiring and cabling. Fig. 4 shows a close-up of line frame 11. On the left of the panel are three 3 × 3 primary switches, toward the right of the panel are four mactor connector relays. Fig. 5 is a rear close-up of line frame 11. This gives a better conception of the division of the crosspoint array into switches than the previous view. It also illustrates the very simple method of wiring the horizontal and vertical multiples within the switches. Having examined the apparatus itself through photographs, let us now see how these various pieces of equipment function together.

The functional arrangement of the switching network and associated controls is illustrated in block schematic form on Fig. 6. Although only 4 lines, 4 trunks, and 2 line frames are shown, the functional relations are unchanged in going to larger networks. The general operation of the various blocks may be summarized as follows:

(1) The DIAD, through the SWITCHING NUMBER GROUP CONNECTOR, tells the NUMBER GROUPS and the SWITCHING SEQUENCE CIRCUIT what terminal(s) (line and trunk or trunk alone) of the SWITCHING NETWORK to manipulate, and whether to establish a connection or release a connection.

(2) The SWITCHING SEQUENCE CIRCUIT applies enabling voltages to the NUMBER GROUP and MACTOR circuits where appropriate, in response to connect or relase signals from the DIAD.

---

* The word "DIAD" derives from "Drum Information Assembler and Dispatcher." The magnetic drum memory and other aspects of the system were described by W. A. Malthaner and H. E. Vaughan, An Automatic Telephone System Employing Magnetic Drum Memory, Proc. I.R.E. Oct., 1953.

Fig. 2 — Switching network and controls — front view.

(3) The LINE NUMBER GROUP CIRCUIT accepts the line information furnished by the DIAD and converts it into codes, voltages, and impedances which can be applied to control the desired line of the switching network on connect.

(4) The TRUNK NUMBER GROUP CIRCUIT accepts the trunk information furnished by the DIAD and converts it into codes, voltages, and impedances which can be applied to the proper trunk of the SWITCHING NETWORK on connect or release.

Fig. 3 — Switching network and controls — rear view.

(5) The SWITCHING NETWORK passes on voltage marks from the marked line and trunk to trace idle paths on connect, provides a metallic talking path between line and trunk during the conversation, and releases the connection upon receipt of a release signal from the trunk number group.*

---

* It is interesting to note that after the connection is established, it is sufficient to furnish only trunk information to perform a release. This is true because the series holding path from trunk to line constitutes a memory of what line and what crosspoints are associated with a particular trunk.

Fig. 4 — Front view of line frame 11.

(6) A pair of MACTOR CONNECTORS, on connect, receive line frame and trunk frame marks and operate to associate the MACTORS with the paths leading between the indicated LINE FRAME and the TRUNK FRAME.

(7) Each mactor, on connect, (through the operated connector) "matches" the two idle-path marking signals received from LINE FRAME and TRUNK FRAME to determine whether the mactor's associated line to trunk path is idle; if the path is found to be idle, the mactor competes in lockout with other mactors which found idle paths, and the mactor which is selected by the lockout operates its associated talking path.

With this description of the principal functions of the switching network and associated control circuits as background, we shall proceed to examine the various components in greater detail.

Fig. 5 — Rear view of line frame 11.

THE SWITCHING NETWORK

*The Switch Crosspoint*

The basic building block of the switching network is the crosspoint. One crosspoint appears at each intersection of horizontal and vertical multiples. The packaging design of the crosspoint is such as to achieve maximum flexibility in the laboratory operation and testing of the network. The crosspoint has two parts, (1) a switch assembly using reed contacts, and (2) a small cold-cathode diode.

The crosspoint is shown on Fig. 7. It consists of four reed contacts or switches in an actuating coil with appropiate mounting and connecting parts. Each reed contact consists of two 52 alloy* rods sealed in opposite ends of a piece of glass tubing containing helium gas. The contacting ends of the 52 alloy* rods are rhodium plated. The soft

---

* This is an iron-nickel alloy containing approximately 52 per cent nickel.

steel cylinder housing the switch assembly (1) serves as a mechanical housing and a support for the plug and socket; (2) acts as a magnetic shield between adjacent crosspoints; and (3) acts as a partial return path for the magnetic circuit of the reed contacts. The plug on one end of the switch assembly serves to connect the crosspoint assembly to the external circuit, while the socket on the other end receives the cold-cathode diode. (Where the switch assembly is used as a simple relay, a short-circuiting plug replaces the cold-cathode diode.)

The operating margins of the switching network are to a large extent a function of the variation in diode breakdown and sustain voltages with



Fig. 6 — Relation of switching circuits.

Fig. 7 — Crosspoint assembly.

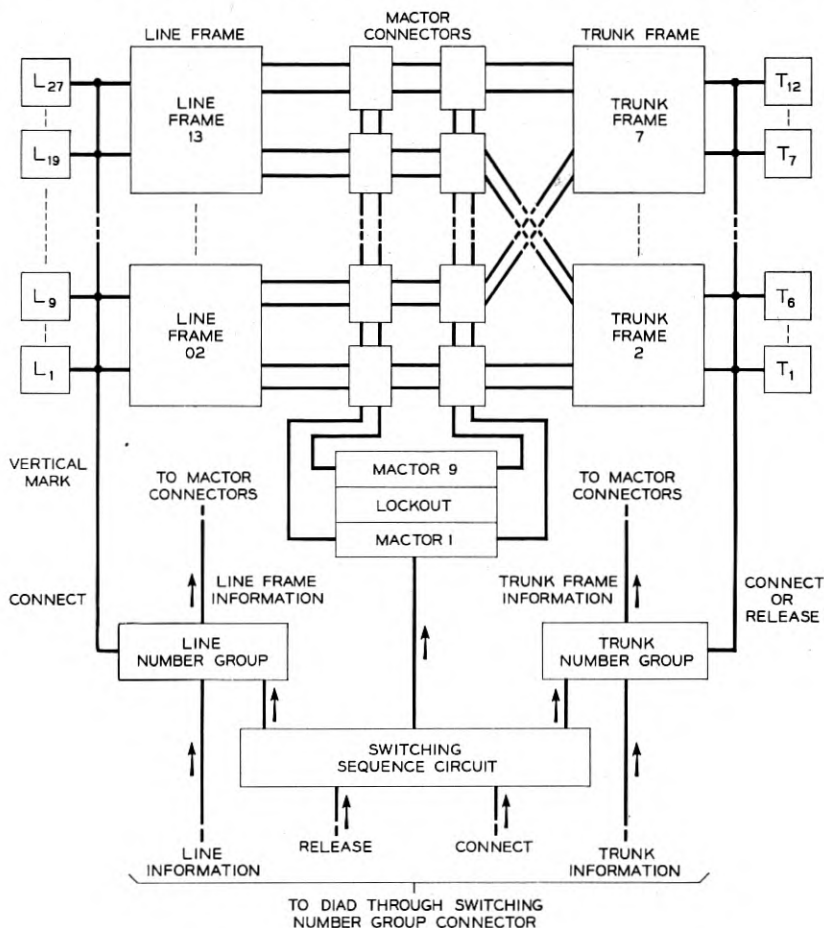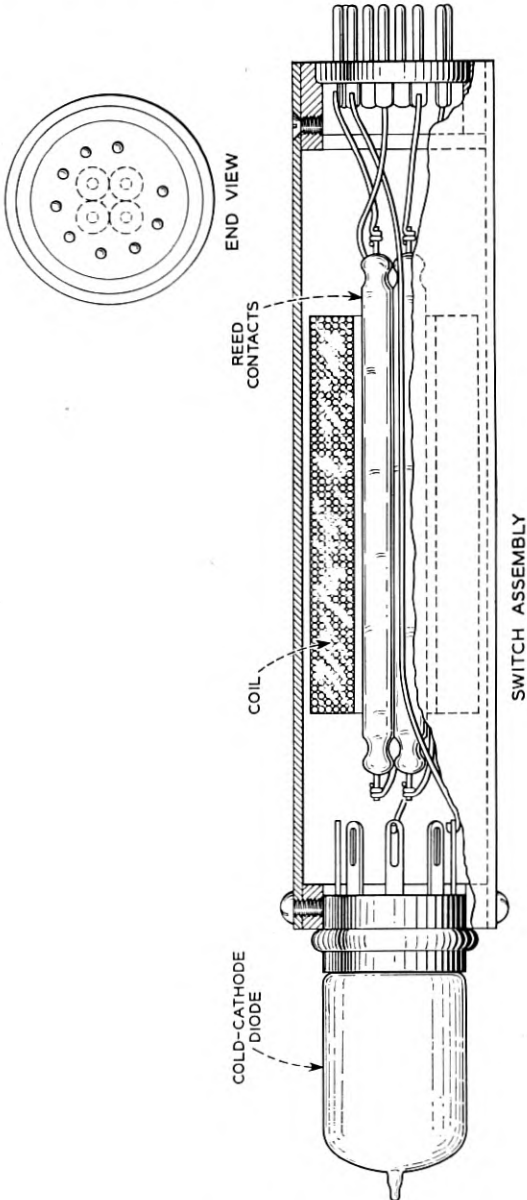age, temperature and also the variation from unit to unit. To assure satisfactory operation, electrical and life requirements were specified on the basis of the switching network needs and the A-1676 diode was designed to meet these requirements. In designing the switching network, arrangements which would lead to overly severe requirements on the diode were avoided. The A-1676 diode is a simple, rugged structure which meets or exceeds the requirements in every regard.

*Single Path Circuit*

A single path through the switching network is shown on Fig. 8. In each of the four crosspoints, the general operation is the same. The cold-cathode diode in series with the coil in the control lead is used (1) to pass tracing voltages through idle paths while remaining an open circuit for busy paths and (2) to break up parasitic paths in parallel with operated crosspoints, and thus to avoid marginal release requirements on the crosspoint relays. The reed contact in the sleeve lead, s, is used to establish a series path through the switching network for the purpose of holding the crosspoint relays of the selected path operated during the conversation. The reed contacts in the tip, T, and ring, R, paths are used to establish a series metallic path for signaling and talking between subscriber and trunk. Within the switching network the route of the tip and ring leads parallels that of the sleeve lead, and hence, tip and ring leads will not be shown or described in detail in the sections to follow. The fourth reed contact is used for display purposes only and has not been shown on Fig. 8.

In the following description of Fig. 8, a complete switching network is assumed, although for simplicity, branching paths have been indicated by "multiple whiskers". Starting from the idle condition, a path through the switching network is established as follows:

(1) Line frame and trunk frame identifying signals (marks of +130 v and −85 v, respectively) are applied to operate the proper connector relays to associate the mactors with the paths which go between the desired LINE FRAME and the desired TRUNK FRAME.

(2) A switch enabling mark (−85v) is applied to the idle A links leading out of all line primary switches having the same switch number as the desired line.

(3) A vertical identifying mark (+130v) is applied to all lines having the same vertical number as the desired line.

(4) All diodes appearing between verticals marked in (3) and A links enabled in (2) break down and, in doing so, produce voltage shifts on the associated A links leading into the line secondary switches.
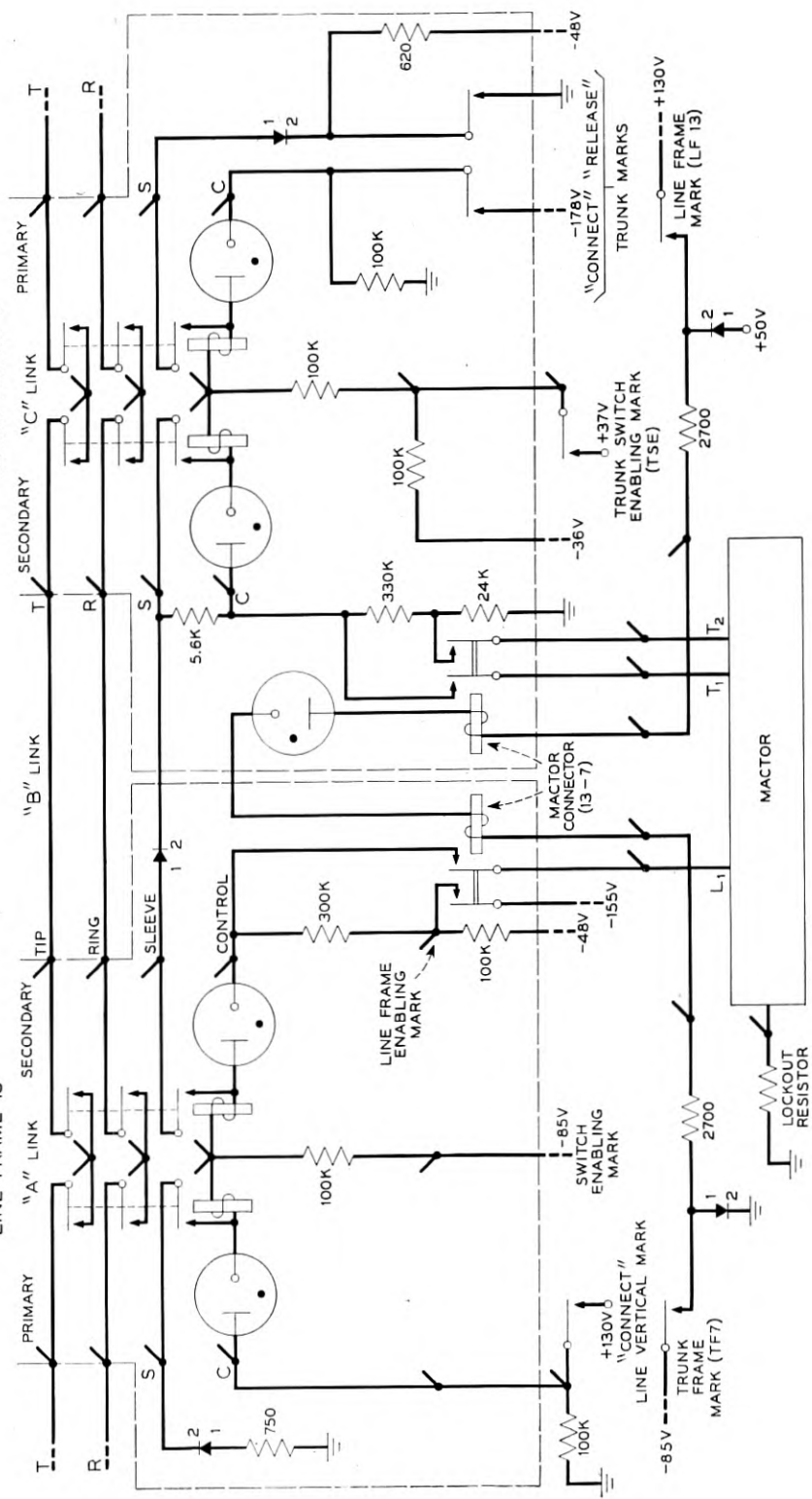
Fig. 8 — Single path through switching network.

(5) The left contact on the line side MACTOR CONNECTOR RELAY during step (1) has applied an enabling bias of – 155v to the control leads which lead out of the desired line frame to the mactors.

(6) The voltage shifts on the marked A links, in conjunction with the enabling biases of step (5), cause breakdown of idle line secondary diodes connecting between marked A links and the enabled B link control leads.

(7) Breakdown of the secondary diodes produces voltage shifts or marks which are fed on the control leads into the $L_1$ terminals of the mactors, each mark representing an idle path between the desired line and the mactor in question.

(8) On the trunk side, an enabling mark (+37v) is applied to all idle C links.

(9) A trunk mark (–178v) is applied to the individual trunk to which connection is desired.

(10) As a result of (8) and (9), the trunk primary diodes representing connections between idle c links and the desired trunk are broken down, producing voltage shifts on these links.

(11) Each mactor, through an idle-test circuit, samples the voltage on its $T_1$ lead, and, if this voltage indicates an idle B link, the mactor connects an enabling bias of +107v to its $T_1$ and $T_2$ leads.

(12) The voltage shifts on the marked c links, combined with the enabling biases applied by the mactors' idle-test circuits, cause breakdown of all of the trunk secondary diodes representing idle paths usable for the trunk half of the desired connection.

(13) The breakdown of the trunk secondary diodes puts voltage shifts or marks on the $T_1$ leads going to the mactors, each mark representing an idle path between the desired trunk and the mactor in question.

(14) Mactors which receive marks from both line and trunk sides indicate through a matching or coincidence circuit that their associated paths through the switching network between the desired line and the desired trunk are idle and available for use.

(15) The mactors which indicate a match energize their respective transistors, which compete in a common lockout circuit.

(16) The transistor which wins the lockout competition (i.e., achieves the fully conducting state) thereby selects the corresponding switching path for use.

(17) In response to the lockout operation, the selected mactor raises the current in both line and trunk halves of the selected path and operates the crosspoint relays associated with this path.

(18) During the conversation the selected path is held operated or "locked up" over the series sleeve lead path between ground on the line

end of the connection and the –48v battery on the trunk end of the connection.

(19) The mactors are retired and the paths which were marked but not used are restored to normal by the removal of the line and trunk marks and the removal by the switching sequence circuit of various enabling voltages.

(20) At the completion of the conversation the connection is released by the temporary application of ground to the sleeve lead at the trunk, thereby bringing both ends of the holding circuit to the same potential and causing release of the crosspoint relays involved in the connection.

Regarding the holding circuit, it will be noticed that there is a 750-ohm resistor in series with the sleeve lead at the line end of the connection. This resistor simulates the winding of a cut-off relay. Hence, if a cut-off relay were required, it would be a simple matter to substitute the operate winding of a relay for the 750-ohm resistor.

A more detailed description of the control of the switching network is provided in Appendix I.

CHANNEL SELECTORS OR MACTORS

*Functional Description of Mactors*

As the description of the switching network has indicated, the mactors perform the functions of (1) matching of line and trunk idle-path tracing marks, and (2) selecting and connecting for use one of the idle paths revealed by the matching operation. To these general functions, in the present system, we have added the function of performing an idle-test on the B links of the switching network. The idle-test precedes the other portion of the mactor operation; hence, a mactor is enabled to continue its functions only if its associated B link is idle. Each mactor contains, as its principal elements, 4 cold-cathode gas tubes and a reed relay associated with each, and a transistor. One gas tube and relay combination performs the idle-test; another combination receives the path tracing voltage from the line side of the switching network; the third such combination receives the path tracing voltage from the trunk side of the network; the transistor (in conjunction with those in other mactors) functions to select one of the idle paths for use, while locking out all others; and finally, the fourth gas tube and relay combination receives a triggering signal from the transistor and functions to establish a connection over the selected path. For a detailed description of how these functions are performed, we refer the reader to Appendix II.

*Mactor Design Features*

In general, the design of the mactor in the present system was based on maximizing the operating margins of switching network, rather than minimizing the amount of apparatus in the mactor. However, it should be noted that there are only as many mactors as there are alternate paths or channels in the switching network (3 in our scaled-down model, perhaps 10 in a full-sized office) and this means that the amount of apparatus in the mactors is a very small proportion of the total office equipment. The use of a separate idle-test, and the detection of line and trunk path tracing marks on different gas tubes lead to increased margins. The idle-test circuit used in the mactor is of particular interest because it uses a new grid-controlled cold-cathode tube, type M-1652. This tube, vi of Fig. 15, has control characteristics similar to those of a hot-cathode gas tube, while possessing all of the long life advantages of a cold-cathode tube in applications involving a low duty ratio. The main cathode, MC of vi also functions as a keep-alive anode, while KC is a keep-alive cathode. Discharge is initiated in this keep-alive gap by the application of −178 volts to terminal CE 1 a few milliseconds before the tube is required to respond to the control signal. At the same time the keep-alive gap is energized, the main anode, MA, is enabled by the application of +130 volts to terminal AE 1. The tube is now ready to respond to a signal on its control grid. A negative potential of 20 volts or more on the control grid prevents transfer of discharge to the main gap while ground potential assures that transfer to the main gap will take place. In addition to the advantages of this tube from the standpoint of life, the fact that it will trigger on relatively small voltage shifts in high impedance circuits is utilized in the present circuit design.

The advantage from a margin standpoint of receiving line and trunk path-tracing marks on separate tubes may be understood if we consider what happens when these marks are received across the start gap of a single tube. It is evident in such a case that the variations of line and trunk path tracing marks must be summed in determining the least favorable circuit operating conditions. However, when the two marks are received on individual tubes, the marks are quantized before combining or matching takes place, so that variations in the two marks are never added together.

Perhaps the most interesting design feature of the mactor is the negative resistance lockout circuit using the switching transistor, type A-1698. In the usual operating routine, the lockout in this circuit is decided on a time basis; i.e., due to statistical time variations one transistor is operated and fully conducting before other transistors are energized. In the oc-

casional "dead heat", a plurality of transistors can reach the high current condition together. In this case, it is significant to note that the transistor in each mactor has a high (29,800-ohm) base load resistance which gives it a negative emitter-to-ground impedance throughout the operating region. The emitter-to-ground impedances of the various mactors' transistors are connected in parallel, and this parallel combination is connected in series with a common lockout resistor (29,800-ohms). It can be shown that this circuit is unstable if more than one transistor is conducting in the operating region. Thus, even though two or more transistors should reach the operating region, one of them quickly assumes most of the circuit current and drives the other(s) into the quiescent region of positive emitter-to-ground impedance and virtually zero emitter current. The maximum time required for one transistor to reach the high current state and drive the other(s) to the quiescent state is the "severance time". For the lockout circuit employed in the mactor, there is a severance time of approximately 0.5 microsecond. It is essential, of course, to avoid producing an output before the lockout has been resolved. This requires that the circuit responding to the transistor have a minimum response time considerably exceeding the severance time. In the mactor circuit this is assured by the fact that the minimum triggering time of the cold-cathode tube, v3, is approximately ten microseconds. In addition, it will be noticed that there is a delay in response caused by feeding the transistor output. which occurs across the base load impedance, through an rc filter before placing it on the start cathode, sc, of v3. This rc filter was designed to have a delay great enough to take care, not only of the severance time, but also of possible temporary ambiguities caused by contact chatter in the l4 relays. These provisions result in a lockout circuit which is "a priori failure-proof".

FUNCTIONAL DESCRIPTION OF SWITCHING SEQUENCE CIRCUIT

The switching sequence circuit, see Fig. 16, is the common control circuit which functions whenever a connection through the switching network is established or released.

In the first case, when the switching sequence circuit receives an (es) establish signal from the switching information dispatcher, (1) it prepares the line number group circuit and trunk number group circuit to handle the coded signals on the line and trunk information leads (signals which correspond to the line and trunk to be connected together) (2) it applies operating voltage to the mactors, (3) it applies the connect trunk marking voltage to all the trunk frames of the switching network, and (4) it applies enabling bias to the c links of the switching network.

After one of the idle paths between the line and trunk in question has been selected and the path established, the switching sequence circuit restores the number group circuits, mactors and mactor connectors to normal and removes the connect voltage mark and c link enabling bias from the switching network. The switching sequence circuit then transmits an OK signal to indicate to DIAD that the required connection has been made. Finally, the switching sequence circuit restores itself to normal.

In the second case, that in which the switching sequence circuit receives a (DS) disconnect signal from the switching information dispatcher, the sequence circuit enables the trunk number group (to handle the signals on the trunk information leads) and applies a release ground to the trunk number group elements. (As was noted earlier, a connection is released from the trunk side only of the switching network.) Sufficient time is allowed for the crosspoint relays in the path to release (about 10 milliseconds). Subsequently, the release ground is removed, and the trunk number group circuit is restored to normal. Finally, a pulse is sent out on the OK lead which signifies to DIAD that the path in question has been released and at the same time the switching sequence circuit restores itself to normal.

Connect and release cycles for the switching sequence circuit are described in detail in Appendix III.

SWITCHING NUMBER GROUP CIRCUIT

The switching number group circuit comprises the following parts: line vertical number group, line switch number group, trunk number group, and mactor connectors. All of these circuits have been skeletonized to conform with the skeletonized switching network. One of them, the line vertical number group, is shown in detail in Fig. 17, while other portions of the number group appear in block diagram.

The switching number group circuit performs two functions: (a) it decodes the signals on the line and trunk information leads from "2 out of 5" to "1 out of 10" and (b) it provides the means for selecting and connecting to the proper points in the reed-diode switching network and for applying the several connect and release voltage marks. The number group circuits were designed on the basis that the signals on the line and trunk information leads are in the form of the equipment designations of the line and trunk which are being served; that is, frame, switch and vertical. Each digit is indicated by the presence of +150v on two out of a group of five information leads, the other three leads being at +100v.

(In the quiescent state, all the information leads would be at +100v, which means that no information is being supplied to the switching network.) The switching number group circuit is enabled and released by the switching sequence circuit. In setting up a connection, the switching sequence circuit applies +200v to the main anodes of the line number group digit tubes via the LNG lead and +200v to the main anodes of the trunk number group digit tubes via the TNG lead. On the release of a connection, only the trunk number group is enabled.

The decoding and selection operations of the number group circuit are described in Appendix IV.

## COMPONENT AND CIRCUIT TESTS

Of the five types of cold-cathode gas tubes used in the design, three are essentially commercial products. These are the WE 376-B, the BTL A-1703, and the BTL A-1704. The latter two tubes are versions of the WE 425-A and WE 426-A without the plastic mounting flange and external resistors.

The other two tubes types, the BTL A-1676 and the BTL M-1652 are not commercial products. The A-1676 diode which appears in each crosspoint was tested as part of the tube design and development effort. Life tests on this gas diode gave a life expectancy of 3,000 active hours. Since the tube is used only momentarily each time a connection is established, this corresponds to an in-circuit life of 240,000 years. Such figures, of course, merely mean that a high degree of reliability may be expected in using this tube. The other non-commercial tube, the M-1652, is of interest because it is a grid-controlled cold-cathode tube which can be triggered reliably from low voltage, low power signals. The control characteristics of eleven of these tubes were measured in detail. Each tube was found to be suitable for use in the mactor circuit, where it is used as the idle-test tube. A typical control characteristic of one of these tubes is shown on Fig. 9.

Each reed-diode crosspoint was tested individually in four respects: (1) dc coil resistance, (2) insulation resistance, (3) operate and release current, and (4) operate and release time. The last two tests were made with positive and negative polarities of operate current and with contacts strapped in series and in parallel. The reversed polarity test established that any residual effect due to hysteresis in the magnetic material was very small. The operate and release tests with contacts in series and in parallel gave a measure of the total spread from the most sensitive contact to the least sensitive contact. Control limits were calculated for

Fig. 9 — Trigger characteristic of grid-controlled cold-cathode tube.

the data on dc resistance, operate and release currents, and operate and release times in order to reveal any crosspoints which were out of control in any of these parameters.

The switching number group, switching sequence circuit, mactors, line frames, and trunk frames were tested individually before being interconnected. The tests consisted generally of making breadboard circuits which would apply the proper operating and information signals to the circuit under test. The test input information was then put through a variety of conditions calculated to exercise every element in the circuit.

Following the testing of individual circuits, they were assembled on the two switching bays, interconnected, and tests of the combined apparatus were made. Principally, these combined tests consisted of: (1) measurement of the timing of all of the important steps in establishing and releasing a connection, and (2) checking the satisfactory operation and release of paths involving all possible combinations of line, trunk, and interconnecting path.

*Timing Measurements*

The timing measurements were made by putting the switching network and control circuits through connect-release cycles repetitively at about 15 cycles per second, and observing operation on an oscilloscope. For the measurements which led to the sequence chart of the normal

connect cycle, Fig. 10, the sweep of the oscilloscope was triggered from the ES pulse which starts the connect cycle.

Taking the start of the connect cycle as zero on the time scale, and the times in milliseconds, the following relations are apparent: (1) at $t = 4$, the major control circuits; (line number group, trunk number group, and mactors) receive their enabling voltage from the switching sequence circuit; (2) at $t = 6$, the line and trunk number groups operate; (3) at $t = 8$, the mactor connector relays operate; (4) at $t = 10$, the mactors complete their idle-tests, receive voltage marks from line and trunk sides of the path, and one or more of them returns to the switching sequence circuit a "match" signal (MS); (5) at $t = 12$, the mactor lockout circuit operates to perform the path selection, and the current is increased in the selected path; (6) at $t = 13$, the crosspoints in the selected path operate; (7) at $t = 15$, the main anode voltage is removed from V1 of the switching sequence circuit; (8) at $t = 19$, the enabling voltages are removed from the number groups and mactors; (9) at $t = 23$, the main anode voltage is re-applied to V1 of the switching sequence circuit; (10) at $t = 23$, the switching sequence circuit sends an OK signal



Fig. 10 — Connect cycle when call is completed.

Fig. 11 — Connect cycle when call is blocked.

through the switching number group connector to the switching information dispatcher, indicating that the switching network and controls are ready for a new connect or disconnect operation. The above cycle description assumed that an idle path was found.

The time relations which were measured with all usable channels busy are shown on Fig. 11. In this case, no idle path is found and the cycle differs from the above beyond (3), as follows: (4) at $t = 10$, the mactors complete their operation without finding an idle path, and therefore without sending a match signal to Switching Sequence Circuit; (5) at $t = 28$, the time-out circuit removes main anode voltage from v1 of the switching sequence circuit; (7) at $t = 34$, number groups, mactors, and mactor connectors release; (8) at $t = 36.5$, main anode voltage is re-applied to v1 and the switching sequence circuit sends a BLOCK signal through the switching number group connector to the switching information dispatcher, indicating that the desired connection could not be obtained, but that the switching network and controls are ready for the next operation.

The disconnect cycle has only one possible sequence, since there is no question of blocking forcing the system into an alternate cycle. The disconnect cycle is simpler, too, since the connection is released from the trunk end alone. Thus, the line number group and mactors are not required to function on a release operation. Time readings on the release cycle, which are displayed in the form of a sequence chart on Fig. 12,

were taken with the start of the oscilloscope sweep synchronized with the
DS pulse which is used to initiate the release cycle. Taking the start of
the DS pulse as zero on the time scale, and the time in milliseconds, the
following relations are apparent: (1) at $t = 3.5$, the trunk number group
receives enabling voltage from the switching sequence circuit; (2) at $t
= 6$, the trunk number group operates, (3) at $t = 9$, the trunk release
(TR) voltage is applied to the trunk, and main anode voltage is removed
from v2 of the switching sequence circuit; (4) at $t = 12$, the crosspoints
comprising the connection release; (5) at $t = 18$, the trunk release signal
is removed from the trunk, enabling voltage is removed from the trunk
number group, main anode voltage is re-applied to v2, and an OK signal
is sent through the switching number group connector to the switching
information dispatcher; (6) at $t = 19$, the trunk number group releases
and all circuits are back to normal.

OPERATIONAL TESTS

In addition to the tests described earlier, it was decided to check both
conductors of every path which could be set up, using the switching
network's controls to connect and disconnect. Since there are 27 lines,
12 trunks, and three possible paths between each line and trunk, and
two conductors for each path, a total of 1944 connections were tested.
In testing all possible combinations as described here, any single portion
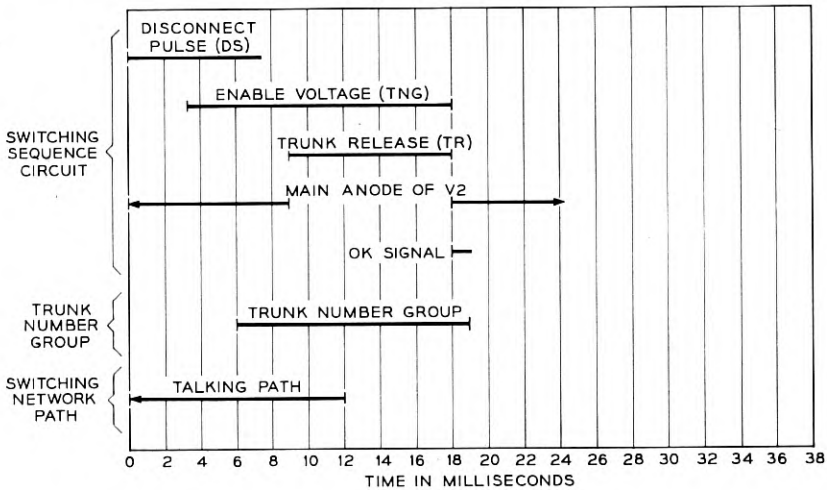of the circuit is involved in a number of tests, giving a certain degree of



Fig. 12 — Release cycle.

redundancy. Nevertheless, a full testing schedule was considered desirable in the present instance for two reasons: (1) the operate and hold characteristics of crosspoints in a marginal case can be a function of how they are combined (i.e., crosspoints $x$ and $y$ could operate satisfactorily in different paths, but fail to operate in a path which involved both of them); and (2) a full test increased the likelihood of discovering any incipient troubles. In testing a larger switching system by this method, the total number of operations required and the degree of redundancy associated with a full test would be excessive. In such a case the ideal test would operate on a fully automatic basis and would follow an abbreviated schedule which would test each element of the switching system the desired number of operations only.

PERFORMANCE OF NETWORK AND CONTROLS

For forty-two months the switching network and controls were operated to gain information on reliability. The results of this operation are summarized briefly on Fig. 13. The left column lists the major circuits which required adjustment or replacement of elements. The second and third columns, give, respectively, element replacements made as a result of preventive maintenance and those made as a result of equipment failures. Most of the germanium diodes which were replaced were used as relay contact protection. Diodes used as logical combining elements have functioned reliably with no failures. Three means of improving the situation regarding diodes used as contact protection are: (1) install all diodes on plug-in mounts so preventive maintenance can catch incipient troubles, (2) use a diode with a higher inverse voltage rating, or (3) use resistance-capacitance contact protection. The germanium diode failure rate diminished with time, indicating that we were weeding out the weaker diodes. Early in the tests the diodes were transferred to plug-in mounts; therefore, most later diode replacements were made as a result of routine maintenance tests without incurring circuit failures.

Twenty-eight vacuum tubes were replaced during the tests. Ten were replaced as a result of circuit failure while eighteen were replaced as a result of routine tests. All of these vacuum tubes were used in the regulated power supplies. To catch weak tubes before they produced failures, the following monthly power supply regulation check was applied. With 80 per cent of the maximum rated load applied, the power supply output voltage should not sag more than 3 per cent. If a power supply passed this test, no further check was made. If it failed, tubes were

| CIRCUIT SERVICED | ELEMENTS REPLACED ON ROUTINE BASIS | ELEMENTS REPLACED BECAUSE OF CIRCUIT FAILURE | TYPE OF ELEMENTS RESPONSIBLE FOR SERVICING | REMARKS |
|---|---|---|---|---|
| POWER SUPPLIES | 18 | 10 | VACUUM TUBES | TYPE REQUIRING MOST REPLACEMENTS: SERIES REGULATOR TUBES, 6Y6 |
| NUMBER GROUPS | 9 | 24 | POINT CONTACT GERMANIUM DIODES | 18 DIODES FAILED SIMULTANEOUSLY DUE TO ABNORMAL POWER VOLTAGE. DIODES TYPE 400E |
| SWITCHING SEQUENCE CIRCUIT | — | 2 | POINT CONTACT GERMANIUM DIODES | 2 DIODES TYPE 400E HAD LOW BACK RESISTANCE; REPLACED WITH GERMANIUM AREA JUNCTION DIODES |
| TRUNK FRAMES | — | 5 | POINT CONTACT GERMANIUM DIODES | DIODES IN TRUNK HOLDING CIRCUIT (TYPE 400E) OPEN CIRCUITED |
| GENERAL | — | — | COLD CATHODE GAS TUBES | BREAKDOWN TIMES ERRATIC; LIGHTS INSTALLED TO PROVIDE INITIAL PHOTO—IONIZATION |

Fig. 13 — Maintenance summary.

checked individually and replacements made as needed. Vacuum tube failures were the chief source of trouble requiring further design effort. The higher degree of power reliability which would be necessary in a commercial system could be achieved by (1) duplicate power supplies with automatic switching to the spares, and (2) by designing individual supplies to a higher order of reliability.

Early in the tests, we discovered that gas tube breakdown times were longer and more erratic than desired. We know that all cold-cathode tubes for their operation depend on a small residual ionization. In the tubes we used, this ionization was provided by a small smear of radium bromide in each tube. The ionization from this source alone proved insufficient to assure stable sensitivities and breakdown times. Therefore, we installed three line-type incandescent light sources of 60 watts each behind the dividing strips of the bays. The resulting photo-ionization produced very consistent tube performance.

Despite the troubles mentioned above, the overall reliability of the switching network and controls has been very satisfactory. Examination of our maintenance log reveals only 23 circuit failures in 42 months, most of the failures occurring early in the test. In particular, the reliability of the cold cathode gas tubes, mercury relays and reed relays has been outstanding. Where germanium diodes function in logic circuits,

we have encountered no failures. The transistor lockout circuit which performs the channel selection has operated throughout the period of test without any transistor replacements.

CONCLUSION

The construction and testing of a 27-line, 12-trunk model of a telephone switching network and its associated electronic control circuits and their operation in the laboratory for over three years have demonstrated the technical feasibility and the reliability of the basic design. On the basis of 50,000 calls per busy hour and 57 per cent active time of the control circuits, this network and its controls are fast enough to permit one-at-a-time operation in a busy 10,000 line telephone exchange. In the design of this skeletonized model, one of the objectives achieved was the use of relatively simple components having practical tolerances. Although the network and its controls were designed specifically for the DIAD telephone system, these circuits are sufficiently flexible that they could be adapted to other information handling systems such as digital computers.

ACKNOWLEDGEMENTS

APPENDIX I

THE SWITCHING NETWORK AND ITS CONTROL

The switching network and mactors are shown on Fig. 14. The network is laid out in standard 4-stage fashion except that the entire circuit is skeletonized. Thus, a network which might contain 10,000 lines and 2,000 trunks in a large office is scaled down to contain 27 lines and 12 trunks. In Fig. 14 only one line frame and one trunk frame have been shown. The other line frames and trunk frames would be similar to those shown. The numbering of lines, switches and frames is such as to

Fig. 14—Switching network and mactors

indicate the full-sized office on which the skeletonized model is based. In the body of this article we described the operation in terms of a single path through the switching network. Now we will see how this description applies when we consider the complete switching network.

As a specific example, let us assume a connection is desired between the line appearing at VERTICAL 56 on SWITCH NO. 8 of LINE FRAME 13 and the trunk appearing at VERTICAL 4 on SWITCH 15 of TRUNK FRAME 7. This line and trunk appear near the top of Fig. 14. Leads involved in establishing and releasing a connection between this line and trunk have been shown by heavy lines on Fig. 14.

At the bottom center of Fig. 14 there are leads which carry line frame and trunk frame information on the basis of which the proper mactor connector relays may be operated. In the example we are considering, the LF13 and TF7 leads would be energized, thus operating the two mactor connector relays which are shown on Fig. 14. Their operation associates the mactors with the B links extending between LINE FRAME 13 and TRUNK FRAME 7. Now let us examine the line frames and see how the proper line terminal is marked. This is done by first applying enabling voltage to the idle A links of the desired switch and all other switches having the same switch number, and then applying marking voltage to the proper vertical number on all line primary switches. Since both of these voltages are required to produce breakdown in the line primary diodes, only diodes associated with the desired line or with similar numbered lines on other line frames are broken down. Referring again to Fig. 14, at the bottom center of the line frame are line switch leads designated LS 0, LS 5, LS 8. In this case, the LS 8 lead is energized, resulting in the application of an enabling bias to the idle A links leading out of line primary switches numbered 8 on all line frames. At the bottom left of Fig. 14 are terminals designated LV 56, LV 09, and LV 00. In our present example, LV 56 has a marking voltage which is applied to VERTICAL 56 on every line primary switch. However, since only the line switches numbered 8 have enabling marks on their links, the vertical mark is effective in breaking down only the crosspoint diodes associated with VERTICAL 56 on LINE PRIMARY SWITCHES 8 of each line frame. At this point, it should be noted that if any of the A links had been busy, the voltage on the links would be determined primarily by the holding circuit of the busy path (which has a low impedance) rather than the link enabling voltage, which passes through a 100K link resistor. This link voltage for busy paths is arranged to act as a bias which prevents the breakdown of primary diodes. This same method of marking links busy applies to the c links.

MACTOR NO. 9

* TO SWITCHING SEQUENCE CIRCUIT

The foregoing paragraphs have described the process of marking the idle line links. In the ensuing description, let us assume that all of these links were found idle, and therefore were marked. The marks or voltage shifts are passed on to the appropriate inlets of the various line secondary switches of the line frames, here represented by LS 1, LS 4, and LS 9 on LINE FRAME 13.

It will be noted that the line side relay of the mactor connector pair, when it operated, applied an enabling mark of –155v to the control leads of the B links leading out of LINE FRAME 13 toward TRUNK LEAD 7. Inter-frame links which go to other trunk frames or which come from other line frames receive no enabling marks and hence cannot be marked in the next step of the path tracing operation.

As a result of the marks appearing at the inlets of the line secondary switches, plus the enabling marks applied by the mactor connector relay, the secondary diodes which could be used in establishing the desired line-to-trunk path are broken down. (Note that secondary diodes on frames other than LINE FRAME 13 are not broken down because of the absence of the –155v B link enabling signal.) The breakdown of secondary diodes on LINE FRAME 13 places path tracing marks on the corresponding control leads which pass through the contacts of the left mactor connector relay. These marks then proceed to the $L_1$ terminals of the various mactors. In this case, the marks from LS 1, LS 4 and LS 9 go, respectively, to MACTORS 1, 4 and 9. This completes the path tracing and marking operation for the line half of the connection.

The marking of the trunk half of the connection (again see Fig. 14) is similar to the marking of the line half of the connection in general principle but varies somewhat as to specific detail. First, all of the idle c links are enabled by the application of +37v to the terminals marked TSE on the TRUNK FRAMES. (Those c links which are busy remain at a voltage which is determined by the division of the –48v holding battery in the busy paths.) Then a –85v trunk frame and vertical identifying signal is applied to the terminal TFV 74 while a trunk switch identifying signal of +130v is applied to terminal TS 15. As a result of these two signals, a trunk number group element shown at the top right of Fig. 14 is selected and operated. The other elements have either no signals applied, or a signal on one side only, and hence do not operate. The SWITCHING SEQUENCE CIRCUIT then applies a "trunk connect" voltage of –178v to terminal TC. This voltage is passed up the number group multiple, through the connect contact of the operated number group relay, and placed on the control lead of VERTICAL 4 of TRUNK PRIMARY SWITCH 15 on TRUNK FRAME 7. The trunk connect voltage, in connection with the

enabling voltage on idle c links, results in the breakdown of all primary switch diodes connecting between the marked trunk and the enabled links. The breakdown of these diodes places a voltage shift, or mark, on the idle links connecting TRUNK PRIMARY SWITCH 15 to the various trunk secondary switches of TRUNK FRAME 7. For purposes of further description, we will assume that all of the c links extending from TRUNK PRIMARY SWITCH 15 to trunk secondary switches of TRUNK FRAME 7 were idle and hence received path tracing marks.

We will leave the path tracing marks for the time being and consider what is happening on the trunk side of the mactors and mactor connectors. It will be recalled that a pair of mactor connector relays has operated. The contacts on the trunk side mactor connector relay extend control leads from the trunk secondary switches of TRUNK FRAME 7 to the mactors. The control leads which are thus extended are those associated with the B links going from TRUNK FRAME 7 to LINE FRAME 13, the frames on which the desired trunk and line appear. The control leads are carried to the $T_1$ and $T_2$ terminals of the various mactors. In this particular case, the control leads from TRUNK SECONDARY SWITCHES 1, 4 and 9 are associated with MACTORS 1, 4 and 9 respectively. By means of an idle-test circuit which is a part of each mactor, the mactors examine the voltages on their $T_1$ terminals. Voltages of approximately −24 volts at these points represent busy B links and prevent the idle-test circuits from functioning, while ground voltages represent idle paths and trigger the idle-test circuits. Assuming that each of the B links under test was found to be idle, each mactor idle-test circuit applies an enabling mark of +107v to its $T_1$ and $T_2$ terminals. These enabling marks travel back over the mactor connector multiple and through the contacts on the mactor connector relay to reach the control terminals of the appropriate links leading from TRUNK SECONDARY SWITCHES 1, 4 and 9. These enabling marks, in connection with the path tracing marks which earlier were placed on the c links, cause breakdown of the trunk secondary diodes of the crosspoints which receive both of these signals. The breakdown of these trunk secondary diodes extends the idle-path tracing marks to the $T_1$ terminals of the mactors. In the case we are considering, the marks from TS 1, TS 4, and TS 9 are received on the $T_1$ terminals of MACTORS 1, 4 and 9, respectively.

The mactors at this point have idle-path tracing marks on their $L_1$ terminals and also on their $T_1$ terminals. Each mactor performs a match of these two marking voltages and thereby indicates the idle or busy status of the corresponding path between the desired line and trunk. Assuming that the match indicates the path to be idle, each mactor

energizes its portion of the common transistor lockout circuit. The lockout circuit is designed so that only one transistor can remain in the conducting state. In this case, we will assume that it is the transistor in MACTOR 9 which conducts and selects the path, and that MACTORS 4 and 1 are locked out. MACTOR 9 then shorts out resistance in series with its $L_1$ and $T_1$ terminals, thereby operating the 4 reed-diode crosspoints in its associated line to trunk path. The operated path is the one running across the top of Fig. 14 whose central link extends between LINE SECONDARY SWITCH 9 and TRUNK SECONDARY SWITCH 9. The operation of the various crosspoint relays comprising this path establishes a metallic holding path between ground on the line end of the connection and –48v on the trunk end of the connection. The various enabling and marking voltages which were applied to the switching network by the switching sequence and number group circuits are now removed. Removal of these voltages causes the various gas tubes shown on Fig. 14 to de-ionize and restore to normal, while the newly operated path remains locked up to the –48v battery.

When it is desired to release the connection after the conversation is completed, the trunk number group is operated. Referring to Fig. 14, the trunk to be released is selected by the application of –85v to terminal TFV 74 and +130v to terminal TS 15. These two voltages operate a trunk number group element shown at the top right of Fig. 14. Then the switching sequence circuit grounds the TR lead and this ground is passed up the number group multiple, through a contact on the operated number group element, and onto the sleeve lead of the trunk. Ground appearing on the sleeve lead reduces the holding voltage to zero, and the crosspoint relays comprising the path release, freeing the crosspoints and links involved for subsequent use. Ground is now removed from the TR terminal, and the –85v and +130v signals are removed from terminals TFV 74 and TS 15, respectively. The previously operated trunk number group element releases, and at this point the entire switching network and controls are back in the original state.

<div align="center">APPENDIX II</div>

*Mactor Operating Cycle*

The mactor operating cycle may be followed by referring to Fig. 15. The sequence is as follows:

(1) Anode and cathode enabling voltages of +130v, ground, and –178v are applied by the switching sequence circuit to terminals AE 1, AE 2, and CE 1, respectively.

(2) Through the operation of the appropriate mactor connector the $L_1$ terminal, and the $T_1$ and $T_2$ terminals are associated with control leads from the line frame and trunk frame, respectively, of one of the B links which could be used for the desired connection.

(3) If a particular B link is busy, a potential of about −24v appears on $T_1$, and hence on the control grid (CG) of $v_1$, preventing transfer of the discharge to the main cathode-anode gap (MC-MA); however, assuming the B link to be idle, ground potential appears on CG via $T_1$, causing the transfer of discharge in $v_1$ to the MC-MA gap and causing the operation of relay $L_1$.

(4) Operation of relay $L_1$ applies +107v to terminals $T_1$ and $T_2$ to enable the control leads associated with the trunk side of the B link and also extends the connection from terminal $L_1$ down to the start anode, SA, of $v_2$.

(5) Assuming the line half of the path through the switching network to be idle, a path tracing mark of 0 volts (nominal) appears at the normally negative start anode, SA, of $v_2$, initiating a discharge from this electrode to the main cathode, MC.

(6) Discharge in $v_2$ shifts to the main gap (MC-MA), thus operating relay $L_2$.

(7) Operation of relay $L_2$ applies −36v to the main cathode of $v_4$ and also extends the control terminal, $T_1$, to the start cathode, SC, of $v_4$.

(8) Assuming the trunk half of the path through the SWITCHING NETWORK to be idle, a path tracing mark of −48v (nominal) appears at terminal, $T_1$, whence it is connected to the start cathode, SC, of $v_4$.

(9) The appearance of −48v on the start cathode, SC, in conjunction with the +130v on the start anode, SA, causes the start gap of $v_4$ to break down.

(10) With the start gap broken down, the discharge in $v_4$ transfers to the main cathode-to-anode gap (MC-MA), operating relay $L_4$,

(11) Operation of relay $L_4$ sends the switching sequence circuit a "match" signal of +107v via the terminal labeled MS, to indicate that the mactor has found its complete line-to-trunk path idle; operation of $L_4$ also places the mactor's transistor in lockout competition with those in other Mactors by energizing the collector of the transistor with −67v.

(12) Assuming that this transistor wins the lockout competition, it will go to the high current condition, and, in doing so, will produce in the common lockout resistor a voltage drop which will prevent transistors in other mactors from reaching the high current condition.

(13) The high current condition produces a negative voltage at the base, b, of the transistor, which negative voltage is passed to the start cathode, SC, of $v_3$, initiating discharge in its start gap (SC-SA).

(14) Discharge in v3 transfers to the main gap (MC-MA), operating relay L3.

(15) Operation of L3 increases the current and causes the crosspoint relays to operate in the line and trunk halves of the selected path through the switching network.

(16) Previously applied enabling voltages appearing at terminals AE 1, AE 2, and CE 1 are removed from the mactor by the switching sequence circuit, interrupting discharges in tubes v1, v2, and v3 and releasing relays L1, L2 and L3.

(17) Release of L2 interrupts discharge in v4 by removing voltage from its main cathode, MC, while simultaneously causing relay L4 to release.

(18) During the release of the various relays, the NE-2 cold-cathode diodes absorb some of the inductively stored energy, thus limiting the magnitude of voltage transients.

(19) After a deionization time of a few milliseconds, the entire mactor is again in its normal state and ready to function on a new call.

(20) Because the connection which was just established remains held up over the sleeve lead's series metallic path, it requires no further servicing from the mactor.

## APPENDIX III

### DETAILED DESCRIPTION OF SWITCHING SEQUENCE CIRCUIT

*Connect Cycle*

Refer to Fig. 16, which shows the SWITCHING SEQUENCE CIRCUIT. On "connect" the voltage on the ES lead rises abruptly from approximately +100v to +150v. This causes the control anode of tube v1 to rise momentarily from ground to +50v. Thereupon, v1 fires, conduction transfers to its main gap, and the two mercury contact relays L1 and L2 operate. However, before L1 and L2 operate, the switching number group connector operates and applies the appropriate signals to the information leads of the number group. The operation of L1: (a) enables all the trunk primary switches by applying +37v to the idle c links over the TSE lead, and (b) enables the line and trunk number group decoding tubes of the switching number group circuit by applying +200v to their main anodes over the LNG and TNG leads. Relay L2 applies –178v, +130v and ground to the mactors, +130v to the time-out tube (v3) and +130v to the match tube (v4) and –178v connect mark (TC) to the trunk frames.

If there is at least one idle path available, the mactor or mactors so indicate by applying +107v to the MS lead of the switching sequence circuit. This voltage fires the match tube, v4, which operates relay L4.
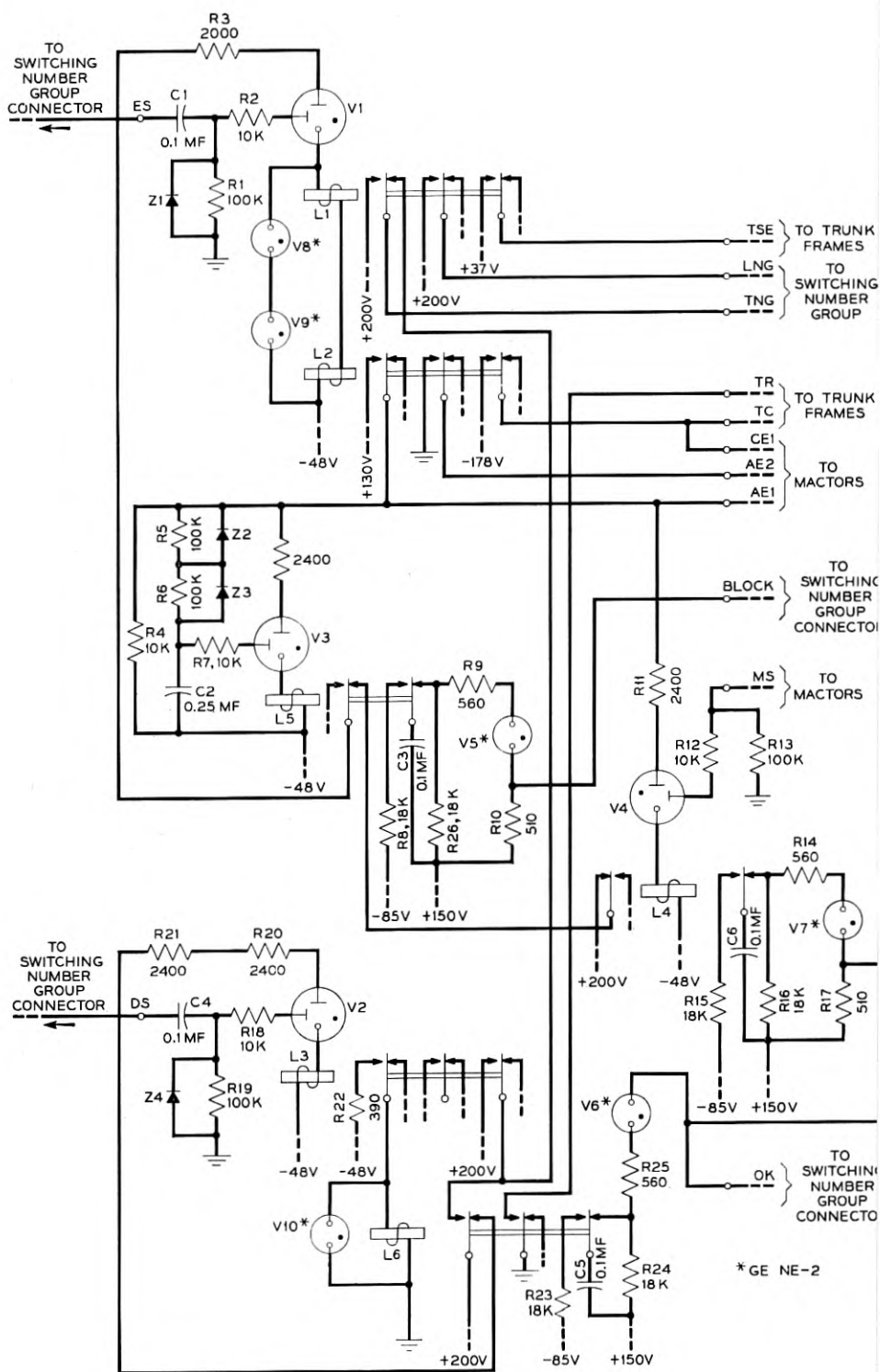
Fig. 16 — Switching sequence circuit.

394

When relay L4 operates, it removes the +200v from the main anode of V1 and at the same time starts the charging of condenser C6 through resistance R15. Tube V1 stops conducting and relays L1 and L2 release. In releasing, these relays remove the +37v and –178v connect marks from the trunk frames, remove +200v from the switching number group circuit, remove –178v, ground and +130v from the mactors and remove +130v from tubes V3 and V4. During the approximately 7 milliseconds required to operate relay L4 and release relays L1 and L2, the crosspoint relays of the talking path are being operated.

After the switching number group and mactors are disabled and all the "connect" voltages are removed from the switching network, the connection is held up between –48v and ground via the sleeve lead. After relay L2 releases, relay L4 releases, thereby restoring +200v to the main anode of V1. At the same time condenser C6, which had become charged between +150v and –85v (to a total EMF of 235 volts), is discharged through gas diode V7. This discharge current develops a negative voltage pulse across resistance R17 about 60 volts in peak amplitude and approximately 0.3 millisecond in duration. Thus, the voltage on the OK lead drops momentarily from +150v to +90v, and this voltage shift is transmitted via the switching number group connector to the DIAD circuits, indicating that the connection called for has been made and that the switching sequence circuit is back to normal and ready to handle another connect or release.

It is recognized that the match (MS) signal does not furnish a positive check that the path has been closed through from line to trunk. However, since this can be verified later as a part of supervision, it was felt the extra time and increased circuit complexity required to make a positive determination that tip and ring are connected through the switching network was not justified.

In connection with the circuit for generating the match signal, it should be noted that one of the purposes of resistance R16 is to prevent tube V7 from firing when the contacts of mercury relay L4 are "bunched" on operate. (The voltage, being divided between resistances R15 and R16, is not high enough across R16 to break down V7, which has a nominal breakdown of 180v.) Also, R16 completes the discharge of C6 after its potential has dropped below the sustain voltage of V7.

If no complete idle path is available between the line and trunk in question, none of the mactors will be fully operated; that is, none of the L4 relays in the mactors will be operated. Thus, no match signal will be sent out by the mactors and condenser C2 in the switching sequence circuit will continue to charge toward +130v through the series-parallel

combination of R5, R6, Z2, and Z3 until the potential across C2 reaches the breakdown voltage of the control gap of V3. This occurs about 25 milliseconds after the appearance of the ES signal. V3 fires, conducts in its main gap and relay L5 operates to perform in the same way as did relay L4 to restore the number group circuit, mactors, and switching network to normal and release relays L1 and L2. Relay L5 also initiates the charging of condenser C3 through resistance R8. Relay L2 in releasing removes voltage from V3, allowing it to de-ionize and relay L5 to release. When L5 releases, condenser C3 discharges through gas diode V5 and a 60v (negative) pulse similar to the OK signal appears on the BLOCK lead to indicate that the connection called for was not made. At the same time, voltage is restored to the main anode of tube V1 and the switching sequence circuit is thus ready to accept another connect or release signal. While relay L5 is releasing, condenser C2 is discharged through resistance R4 and the series-parallel combination of R5, R6, Z2 and Z3.

When the ES signal is removed, that is, the voltage on the ES lead drops abruptly from +150v to +100v, condenser C1 is discharged rapidly due to the low forward conducting resistance of the germanium varistor Z1 in parallel with resistance R1.

The two small gas diodes V8 and V9 which are connected in series across relays L1 and L2 have a nominal breakdown of 90v and are employed to limit the inductive voltage developed by the windings of these relays on release, without substantially increasing the relay operate or release time.

### Release Cycle

The indication that a connection through the switching network is to be released is an abrupt rise in voltage on the DS lead from its normal potential of +100v to +150v. This abrupt rise fires gas tube V2 which operates relay L3. Relay L3 (a) applies +200v to the trunk number group and (b) operates relay L6. Before the contacts of relay L6 close, however, the trunk number group functions and the number group element corresponding to the trunk in question is operated. Then relay L6 applies release ground (via a contact on the trunk number group element) to the trunk sleeve terminal. This reduces the holding voltage on the control path of the connection through the switching network to zero, causing the crosspoint relays to release.

The operation of relay L6 also starts the charging of condenser C5 through resistance R23, and interrupts operating current to tube V2 and relay L3. When L3 releases it de-energizes L6 and removes a connection of +200v from the trunk number group. (However, the trunk number

group is still receiving +200v through L6.) Relay L6 in releasing: (a) removes release ground from the trunk control terminal, (b) removes +200v from the trunk number group, (c) restores main anode voltage to tube v2 (which meanwhile has become de-ionized) and (d) allows condenser c5 to discharge through gas diode v6 and produce a 60v negative pulse on the OK lead in the same way that the "connect" OK pulse is generated. About a millisecond following the OK pulse, the trunk number group releases. However, the switching sequence circuit is ready to handle another connect or release just as soon as the OK pulse appears, since the trunk number group is restored to normal before any of the relays L1, L2 or L3 could be re-operated on a subsequent ES or DS signal.

The sequence of operation in applying and removing the release ground is designed so that the relatively heavy release current (100 MA) is established and interrupted with mercury contacts instead of reed contacts of the trunk number group element. The release ground is held on the trunk for sufficient time to insure the release of the switching path crosspoint relays.

As in the case of a connect, no positive check of the release of the path is made and for much the same reasons.

The gas diode v10 (nominal breakdown 90v) limits the inductive voltage surge from the winding of relay L6 without substantially increasing either its operate or release time.

## APPENDIX IV

### DETAILED DESCRIPTION OF SWITCHING NUMBER GROUP CIRCUIT

*Decoding*

The decoding is accomplished by means of varistor-resistor networks and sensitive cold-cathode gas triodes. For an example, let us look at the Fig. 17, where the line vertical units digit of our skeletonized version is represented by the LV2, LV4 and LV7 leads (5 would be required for a full digit LV0, LV1, LV2, LV4, and LV7). Let us assume that the LV2 and LV4 leads are at +150v, while the LV7 lead is at +100v. Furthermore, let us note the following circuit conditions: (a) the backward resistance of the varistors z3, z4, etc. is greater than 100,000 ohms, (b) the forward resistance of the varistors is less than 1,000 ohms, (c) the resistance of the resistors R3, R4, etc., is 0.5 megohm ± 5 per cent, (d) the breakdown voltage of the control gap of the gas triodes LV9, LV6, and LV0 varies from 62 to 89 volts, (e) the impedance of each source supplying voltage to the

information leads is approximately 5,000 ohms, and (f) the bias voltage on the gas triodes is +50v ± 3 per cent. Under the foregoing conditions, the control anode of tube LV6 will rise to a potential sufficient to fire the control gap of LV6, while the control anodes of LV0 and LV9 will remain below the potential required to fire their control gaps. While we have supplied only three input leads and three output digits in our skeletonized number group, the circuit operating margins are very satisfactory with a fully equipped digit, i.e., one with five input leads and ten outputs.

It will be noted in Fig. 17 that not all of the information was coded; as for example, the line vertical tens 0 and 5 (LVT0 and LVT5). Since only two numbers out of the possible ten of a digit were actually provided, it was more convenient to fire the number group tens digit tubes (such as LVT0 and LVT5) directly from the +150v on one of the two information leads and omit the decoding network. This arrangement does not detract from the test value of the overall circuit, since the decoding network is being tested adequately in places where three numbers out of a possible ten are used.

*Selection*

As we have stated before, one of the functions of the number group circuits is to provide a means for selecting and marking the various points in the switching network. These selection circuits: line vertical, line switch, trunk, and mactor connector, are, in effect, digit combining circuits. The line vertical number group, for example, combines tens and units digit information to select one out of a possible 70 line primary verticals which might be provided in a full-sized office. Referring again to Fig. 17, let us assume that the signals on the information leads correspond to LV09. This means that leads LVT0, LV2, and LV7 will carry +150v, and leads LVT5 and LV4 will carry +100v. These voltages fire the start gaps of tubes LVT0 and LV9. When +200v is applied to the LNG lead by the switching sequence circuit, these tubes will transfer, conduct in their main gaps, and operate the reed-contact relays in their respective cathode circuits. Observe that the anode transfer current for digit tube LV9 is conducted through resistors R3 and R4 and varistors Z3 and Z4 effectively in parallel, with the varistors conducting in the backward direction. The circuit has been designed so that each start gap, when fired, receives sufficient current to insure anode transfer, even if the reverse resistance of the varistors is infinite. Relay LV9 applies +130v to the positive side of the LV09 vertical number group element which consists of a diode and reed relay in series. This diode and relay are identical with those used as crosspoint elements in the switching network. Relay LVT0 applies –85v to the nega-

tive side of both the v00 and v09 vertical number group elements. The breakdown potential of the diodes is 180v ± 20v; hence, the only tube which receives sufficient voltage to fire is LV09. When LV09's associated reed relay operates, it applies the +130v connect mark on all the line primary switches in the switching network.

The line switch number group operates in the same manner as the line vertical number group. It selects a line primary switch and applies a −85v mark to the control leads of all the outlets of that switch. Instead of marking only one switch on one frame, however, it marks the same numbered switch on all of the line frames. For instance, if the information leads say that the switch in question is LP8 of LINE FRAME 13, the mark is applied to the idle A links of LP8 switches on all of the line frames.

The trunk number group is a four digit selection circuit in which the two digits, trunk-switch tens and trunk-switch units, are combined to select one out of a group of 20 tubes (one out of three in our skeleton version), each of which corresponds to a particular trunk switch. Also the two digits, trunk frame and trunk vertical, are combined to select one out of a group of 100 tubes (one out of 4 in the skeletonized circuit), each of which corresponds to a particular vertical on a particular trunk frame. Then these two primary selection groups of 20 and 100 are combined to select one out of 2,000 trunk number group elements (one of 12 in the laboratory model). Each trunk number group element corresponds to a particular vertical on a particular switch on a particular trunk frame. On a connect, this trunk number group element puts a −178v mark on the control lead of the particular trunk appearance. On a release, a ground release mark is applied to the control lead of the trunk appearance, as described in the operation of the switching network.

In the setting up of a connection, a mactor connector is used to connect the mactors to the line and trunk control leads of the unique group of channels or alternate paths between the line frame and the trunk frame on which the line and trunk in question appear. There is a mactor connector for the group of paths from each line frame to each trunk frame. For instance, in a large system having 20 line frames and 10 trunk frames, there would be 200 such mactor connectors, while in the laboratory version, there are 3 × 2 or 6. The relays and the diode comprising one connector are in series. One relay and the diode are mounted on the trunk frame, while the other relay is mounted on the line frame.

As indicated in Fig. 17, the proper mactor connector is selected by combining line frame and trunk frame information in a manner similar to that employed in the other selecting circuits in the number group. Suppose we desire to connect a line on LINE FRAME 13 to a trunk on TRUNK FRAME

2. We see that relay LF13 applies +130v to a mactor connector appearance on TRUNK FRAME 2 and to a connector appearance on TRUNK FRAME 7, while relay TF2 applies –85v to a mactor connector on each of the three line frames. The connectors and their diodes function as a combining



Fig.

matrix; only one connector receives both voltages which are required to fire its gas diode and operate its relays.

When the +200v is removed from the LNG and TNG terminals by the switching sequence circuit, the main gaps of the conducting tubes in the number group circuit are extinguished and the relays which had been operated are released. The only portions of the circuit remaining "on" are the control gaps of those digital tubes which had been originally ionized directly from the signals on the information leads. When the information signals are removed (the potential on all the information leads is returned to +100v), the control gaps of the digital tubes are extinguished and the switching number group circuit is restored to normal.

It will be noted that all of the selections, whether they are 2, 3 or 4



:hing number group circuit.

digit, are made with nothing more complicated than gas diodes and triodes and reed-contact relays. In general, a selection of one out of any number can be made in this same fashion with the same simple elements, by breaking up the number into groups of digits and combining them in one or more stages as illustrated in the foregoing paragraphs. The gas diodes, triodes, and reed relays which perform the selection can have wide tolerances on their characteristics which means that inexpensive elements can be employed.

## APPENDIX V

### POWER SUPPLIES

The power voltages used in the switching network, number group and switching sequence circuits are: $-178$, $-155$, $-85$, $-67$, $-48$, $-36$, $-24$, $-12$, $+37$, $+50$, $+107$, $+130$, $+150$, and $+200$. Of these, the $-67$ volts is obtained by a tap on the $-85$-volt supply, the $-36$, $-24$, and $-12$ by taps on the $-48$-volt supply, and the $+37$ by a tap on the $+50$-volt supply. These taps consist simply of resistance divider circuits with an adequate capacitive by-pass to ground. The $+150$-volt supply is associated only with the generation of an OK or BLOCK pulse for transmission from the switching sequence circuit to the switching information section of the DIAD. Since the switching information section of the DIAD uses $+150$ volts also, no separate source is provided on the switching bays. This leaves eight separate power supplies to consider. One of these, the $-48$-volt supply, is a heavy-duty, regulated rectifier in our laboratory model; however, all of the circuit design has been based on the use of standard $-48$-volt office battery with limits of 44 to 50 volts. The other seven supplies are small-capacity regulated rectifiers. These supplies are assumed to have a variation in output voltage due to all causes of $\pm 3$ per cent. Because of the one-at-a-time nature of the control, the peak load on the power supplies does not increase as the size of the office increases, so that the present supplies have sufficient capacity to handle the largest presently used switching networks. The one exception to this statement occurs in the case of the $-48$-volt supply which is used for holding the path operated and to supply talking battery. The loads due to these two functions have peaks which are directly proportional to the number of simultaneous conversations. It should be mentioned, however, that the holding power for the reed-diode switching network is only one-half watt per path for the four stages.

# Signal-Detection Studies, with Applications

## By E. L. KAPLAN

*Curves are given relating the probability β of detection of a signal in noise to the signal-to-noise power-ratio $S/N$, to the proportion α of false detections that can be tolerated, and to the time available (more specifically, to the number m of independent squared samples of the envelope of the filter output that are averaged in making a single attempt at detection). For the curves, three types of sinusoidal signals are assumed, according as the amplitude is constant, varies at random very slowly (fading), or varies as rapidly as the filtered noise itself. The second case is of course very unfavorable for reliable detection, and the third is also if one is limited to one or two sample points. The curves are applied to problems of optimizing radar parameters such as pulse energy, scan rate, averaging time, etc.*

*The Appendix gives the mathematical background for the foregoing and then proceeds to consider additional types of signal (dc and arbitrary Gaussian) and methods of detection (failure to rectify or take the envelope, counting of samples above a threshold, and averaging by continuous integration).*

## 1. INTRODUCTION

The word detection suggests a decision between just two alternatives: either a signal was transmitted, or it was not. Typically, however, such decisions are made repeatedly at short intervals of time, with the effect that a multiplicity of possibilities are involved, namely the time intervals in which signals are received, and perhaps also the carrier frequencies at which they are transmitted. In radar and sonar systems, rotating directional receivers map one, two, or three dimensions of space upon the time axis, so that the observation of time is translatable into an observation of position. In these applications, of course, the received signal is a reflection of one transmitted by the observer, and the aircraft or other object to be detected is not attempting to communicate with the ob-

403

server. Telegraphy and PCM are applications in which the signals do represent attempts to communicate. In all cases, random noise is present, and will hide signals whose intensity is sufficiently low.

If the bandwidth of the noise accompanying the signal is considerably larger than that of the signal, the best practical procedure is well known to be to pass the signal-plus-noise through a bandpass filter whose bandwidth approximates the larger of the following: (a) The bandwidth of the signal. (b) The reciprocal of the time available for detecting the signal.[1] If the frequency of the signal is not known in advance, one of course moves the pass-band of the filter up and down the frequency scale and thus searches for the signal in frequency as well as in time; the present discussion applies equally well to this case.

In this way much of the noise may be eliminated, but that is not the end of the problem. Some noise still gets through. However, the presence of a signal may be expected to produce an increase in the magnitude of the rectified output. The purpose of this discussion is to indicate how great this increase should be required to be before one decides a signal is present, to draw conclusions therefrom regarding optimum procedures, and to illustrate the great importance of the nature of the signal and the method of detection. For simplicity, square-law rectification is assumed; it is known that linear rectification does not give very different results.[2] The noise, and in some cases the signal also, is assumed to be Gaussian.

We assume the following detection procedure, which may apply literally or else be approximately equivalent to that used in a physical system (e.g., averaging by continuous integration, as in Section 15; or by a human observer, or a phosphor). The rectified output is sampled at $m$ different instants of time,[3] which are far enough apart so that the values of the noise are effectively independent (uncorrelated). If the average output at these $m$ instants exceeds the average noise level $N$ by an amount $kN$ or more, we decide that a signal is present; if not, we decide no signal is present. For simplicity, the value of $N$ is assumed to be known as the result of past observation; since this is only approximately true in practice, the results will somewhat exaggerate the effectiveness of the detection procedure, especially for large values of $m$ and low signal strengths. Another assumption made in calculating the curves is that if

---

[1] Alternative (b) is superfluous if, in determining the bandwidth, the signal is defined to be zero outside the time interval within which it is available for analysis.

[2] M. Schwartz in his Harvard (1951) dissertation, "A Statistical Approach to the Automatic Search Problem," finds agreement to within 0.2 db.

[3] This assumes a narrow-band output whose envelope is sampled; if no envelope is obtained, the number of samples is denoted by $2m$. The point will be discussed below.

a signal is present at all, it is present at each of the $m$ instants sampled; the methods for removing this assumption are discussed later.

It is well known (and shown in Section 5) that the averaging after rectification described above does not fully compensate for a filter whose passband is too wide; but when the filter has been made as selective as it can or should be, the post-rectification averaging, when it is possible, gives a further increase in sensitivity.

In practice many of these averages must be formed, and many decisions made, corresponding to the search for the signal in time and possibly also in frequency and in geometrical position (e.g., radar detection of an airplane). The probability (for any one average or any one decision) of deciding that a signal is present is denoted by $\beta$ or $\alpha$ according as the signal is actually present or not. Thus $\alpha$ is the proportion of false alarms among the total number of averages containing no signal, while $\beta$ is the proportion of valid detections among all those averages that do contain the signal. A detection is considered valid only if it is associated with the average, or one of the averages, in which the signal actually occurs. The average time between decisions divided by $\alpha$ gives the average time between false alarms (in case the signal is usually absent).

The problem could now be formulated thus: to calculate the values of $k$ and $\beta$ corresponding to given values of $m$, $\alpha$, and the signal-to-noise power-ratio $S/N$. Actually it has been convenient to regard $m$ and $S/N$ as the principal variables, and so to plot their relationship for a few different values of $\beta$ and $\alpha$. The number $k$ depends only on $m$ and $\alpha$, and is numerically equal to the $S/N$ values (*not* measured in db) given by the curves for $\beta = 0.50$ in Fig. 1. (These curves should be rigorously correct for $k$, whereas their use for $S/N$ involves an approximation as discussed below.)

## 2. THE THREE TYPES OF SIGNAL AND EXPLANATION OF THE GRAPHS

Three kinds of signal are considered: (1) Steady sinusoid. (2) Fading sinusoid. (3) Noise-like signal. All of these are in a sense special cases, but other signals will generally be of some type intermediate to these.

Cases 1 and 2 are alike in that for both it is assumed that the $m$ signal amplitudes occurring in an average are identical, whereas in Case 3 these $m$ amplitudes are assumed to be independent random variables, just as the corresponding noise values are.

Cases 1 and 2 differ in that in Case 1 the signal-to-noise ratio $S/N$ refers to the instantaneous signal power received, while in Case 2, $S/N$

Fig. 1 — Signal-to-noise ratio required for detection of steady sinusoidal signal.

refers to a long-term average of the signal power, which may be either greater or less than the instantaneous power. Thus the signal amplitude in Case 2 is a random variable, but one that does not change much during the period of time that is averaged over. The Rayleigh or circular Gaussian distribution has been assumed for the signal amplitude, corresponding to a Gaussian signal. For the signal power, this becomes the exponential distribution.

The pulses of carrier used in communication, and possibly the radar return from a fixed object, can be considered as steady sinusoids (Case 1). The radar return from an airplane belongs to Cases 2 or 3 or intermediate cases, since the phase relations of the returns from different parts of the airplane change by different amounts as the aspect of the airplane changes. Case 2 may arise in at least three ways:

(a) The signal as generated may be Gaussian and have a very narrow but non-zero bandwidth; i.e., its Nyquist interval may be as large as, or larger than, the averaging time employed.

(b) The signal may consist of several sinusoids of slightly different

frequencies beating against one another. If there are reflecting surfaces involved and the source of the signal, the observer, and the reflecting surfaces have any relative motion, the same phenomenon of beats occurs even with a single sinusoid, whenever multiple paths for its transmission exist. Interference between radar reflections from an airplane and from its image in the sea surface is a familiar example.

(c) Even if the signal in each instance is a steady sinusoid, one may still have to regard the collection of amplitudes encountered on different occasions as values of a random variable. Thus, since various airplanes have different speeds and radar cross-sections and are detected at different ranges, the strengths of their radar returns will differ accordingly. Another example is provided by the interference between direct and reflected signals in the absence of any relative motion among the signal



Fig. 2 — Signal-to-noise ratio required for detection of a fading sinusoidal signal (amplitude steady during each detection, but variable from one detection to another).

source, the observer, and the reflecting surface. Then the signal strength may be constant with time, but it will still be a random variable depending on the positions that the source and the observer happen to occupy.

The exponential distribution of signal power assumed in the present treatment of Case 2 is theoretically exact for situation (a) above, and possibly as good as one can do for situation (b). In situation (c) the proper distribution may be quite different, though its effect should be qualitatively similar.

In all cases the $S/N$ value is defined as if the signal were continuously present; i.e., it does not depend on the duration of the signal. In Case 1, $S$ is the peak signal power, while in Case 2 it is the average of the various possible peak signal powers. In Case 3, $S$ might be described as the peak value of the statistical expectation of the power.

Curves of $S/N$ versus $m$ are given for the three cases in Figs. 1 to 3



Fig. 3 — Signal-to-noise ratio required for detection of a noise-like signal.

Fig. 4 — Comparison of the signal-to-noise ratios required for detecting various kinds of signals.

respectively, for $\alpha = 10^{-3}$ (dotted lines) and $10^{-6}$ (solid lines) and three values of $\beta$. In Fig. 4 one curve from each of the other figures is reproduced for comparison, and also one new curve (the broken line) to be discussed in Section 6. Curves for $\beta = 90$ per cent may be interpolated in Figure 1 (and in Fig. 3 when $m \geqq 10$) at about 0.58 of the distance from the curve for $\beta = 50$ per cent toward the curve for $\beta = 99$ per cent. This may be done by eye, or numerically with $S/N$ expressed in decibels. The curves for $\beta = 1$ per cent on Figure 1 are probably of low accuracy.

These same three cases are also considered (for somewhat different values of the parameters) by M. Schwartz in his dissertation cited previously, and by J. I. Marcum and P. Swerling in their work at the Rand Corporation.

The curve in the lower left corner of Fig. 1 shows the relation

$$S/N = 2^{1/m} - 1$$

that would exist between $S/N$ and $m$ if the maximum information-carrying capacity of the system were utilized. The maximum capacity is $(1/2m)$ bit per Nyquist interval $(1/2W$; see next Section), achieved by having the signal present with an independent probability of $1/2$ in each average of $m$ samples. According to C. E. Shannon's theorem,[4] the maximum capacity in terms of $S/N$ is $\log_2(1 + S/N)^{1/2}$ bits per Nyquist interval. When $m$ is large (or equivalently, $S/N$ small), the relation is $S/N \doteq (\log_e 2)/m$.

## 3. DETAILS: DC SIGNALS, ENVELOPES, ETC.

This section discusses various other cases of less importance than the three defined in the preceding section; namely, dc signals, low-pass and broadband filters, the relations between instantaneous values and the envelope, and the distinction between narrow-band noise and a noise-modulated carrier. The first three topics are intimately connected; just as sinusoidal (ac) signals, narrow-band filters, and envelopes naturally go together (Sections 1 and 2), so are dc signals, low-pass filters, and instantaneous samples usually associated.

The optimum detection of a dc signal, given independent samples, is trivially easy. No filter is needed, and rectification is undesirable, as shown by the two lowest curves on Fig. 4. The lowest curve, labelled "steady dc," is obtained by the optimum procedure of taking the simple average of the unrectified sample values, some of which may be negative. If the squared values are averaged, the curve will coincide with the next higher curve, labelled "steady sinusoid," as shown later in connection with equations (13b) and (13c). The "fading dc" signal is one which is constant during each detection but has a Gaussian distribution of mean zero from detection to detection. Although it is the worst curve on the figure, quadratic averaging would make it still poorer. [On the other hand, (12c) shows that rectification is as indispensable for the detection of a lowpass noise-like signal as for any other rapidly oscillating signal.] Both of the dc curves are straight lines of 45° inclination, and for them the number of samples is not $m$ but $2m$. This adjustment keeps the total time the same, and in the case of quadratic averaging produces complete equivalence between the dc and the sinusoidal signal, as shown in Section 13. The simple formulas applying to the averaging of unrectified dc signals are given later as equations (12a) to (12c).

If a filter is used in detecting a dc signal, it must be of the low-pass

---

[4] B.S.T.J., **27**, pp. 379–423, 623–656, 1948. Proc. Inst. Radio Eng. **37**, pp. 10–21, 1949. Mathematical Theory of Communication (with Warren Weaver). Univ. Illinois Press, 1949.

type,[5] and its effect is not essentially different from that of the assumed averaging of $m$ sample values. In detecting an ac signal, a narrow-band filter is decidedly preferable to a low pass because it eliminates more noise; but it also has the effect of permitting the (squared) envelope of the filter output to be obtained, provided the frequencies contained in this output are confined to a range of less than 3 to 1. In this case the output of the square-law rectifier will contain two separate bands of frequencies of equal widths which (given a reasonable separation between them) can be separated by another filtration. The lower frequency-band is usually the only one that is wanted, and it gives half the square of the envelope of the output of the first filter. If the frequency range of the first filter output is 3 to 1 or more, the envelope is still susceptible of a theoretical definition but it cannot be isolated so easily and is then a somewhat artificial concept.

The general effect of a second filtration after rectification of an ac signal, whether or not the envelope is actually obtained, is to increase the reliability of a single sample, without necessarily having a marked effect on an average over a long (fixed) period of time. The elimination of the sinusoidal component resulting from the signal, and having double its original frequency, is perhaps the main object. If this is not done, it will be necessary to do considerable averaging to eliminate the possibility of missing the signal because of sampling it at its troughs. In this respect instantaneous sampling of a sinusoidal signal is qualitatively similar to the various cases of the sampling of the noise-like signal.

The properties of the envelope will now be considered. Let $n_1(t)$ and $n_2(t)$ be Gaussian noises containing no radian frequency above $\omega/2$. Then the spectrum of

$$n(t) = n_1(t) \cos \Omega t + n_2(t) \sin \Omega t \tag{3a}$$

is limited to the interval $\Omega - \omega/2$, $\Omega + \omega/2$. Conversely, consideration of the Fourier transform of $n(t)$ shows[6] that any (stationary) noise whose spectrum is limited to $\Omega \pm \omega/2$ can be expressed in the form (3a). It is found furthermore that

$$En_1(t)n_1(u) = En_2(t)n_2(u) = \int_0^\infty P(\lambda) \cos (\lambda - \Omega)(t - u) \, d\lambda \tag{3b}$$

$$En_1(t)n_2(u) = -En_1(u)n_2(t) = \int_0^\infty P(\lambda) \sin (\lambda - \Omega)(t - u) \, d\lambda$$

[5] This may not be true in a rigorous sense. If indefinitely prolonged dc signals are not encountered, a cut-off at a sufficiently low frequency is permissible.

[6] For example, see S. O. Rice, B.S.T.J., **24**, pp. 46–156, 1945, Section 3.7.

where $P(\lambda)$ is the power per radian/sec. for frequencies in $n(t)$ whose absolute values are in the neighborhood of $\lambda$ radians/sec. Let $\psi_0(t - u)$ and $\psi_1(t - u)$ represent the two distinct covariances of (3b), the subscripts 0 and 1 reflecting the even and the odd characters of the two functions. Then

$$En(t)n(u) = \psi_0(t - u) \cos \Omega(t - u) - \psi_1(t - u) \sin \Omega(t - u) \quad (3c)$$

Replacing $\Omega$ in (3a) by $\Omega_0 + \Omega_1$ gives

$$n(t) = [n_1(t) \cos \Omega_1 t + n_2(t) \sin \Omega_1 t] \cos \Omega_0 t$$

$$+ [-n_1(t) \sin \Omega_1 t + n_2(t) \cos \Omega_1 t] \sin \Omega_0 t$$

so that the same envelope

$$\sqrt{n_1^2(t) + n_2^2(t)}$$

is obtained regardless of the value adopted for the "carrier" frequency. However, the highest frequency occurring in $n_1(t)$ and $n_2(t)$ will be minimized by taking $\Omega$ in the center of band; and (3b) shows that if the power spectrum of $n(t)$ has a center of symmetry which is taken as $\Omega$, then $En_1(t)n_2(u) = 0$.

The noise-like signal used in this study, like the filtered noise itself, has the form (3a). The curves do not apply quantitatively to a single term $n_1(t) \cos \Omega t$ obtained by modulating the carrier $\cos \Omega t$ with a low-frequency Gaussian noise-like signal $n_1(t)$, though the necessary tools are given in equations (13g) and (13h) and the remarks following thereafter.

The most obvious difference between the two signals resides in their envelopes, which are

$$\sqrt{n_1^2(t) + n_2^2(t)}$$

(with a Rayleigh distribution) and $| n_1(t) |$ (the absolute value of a Gaussian variable) respectively, but other differences may be discovered. The former can be, and is assumed to be, stationary, Gaussian, and ergodic. The latter is stationary only if the phase $\varphi$ of the carrier $\cos (\Omega t + \varphi)$ is taken to be random, but the relative phases of the sinusoidal components of $n_1(t) \cos (\Omega t + \varphi)$ are still not completely random, and the signal is neither Gaussian nor ergodic. [The Gaussian variable $n_1(t)$ is multipled by $\cos (\Omega t + \varphi)$, which has the distribution $dz/\pi\sqrt{1 - z^2}$.] Where detection is concerned, the two cases are qualitatively similar, but the single term $n_1(t) \cos \Omega t$ has a greater amount of fluctuation for the same amount of power, and is therefore a little harder to detect than the sum of the two terms.

The results of this study have been phrased in terms of the narrow-band two-term form $n_1(t) \cos \Omega t + n_2(t) \sin \Omega t$ for the noise and the noise-like signal. However, the curves apply also to the less important case of the detection of a dc signal by rectification (which is not the optimal procedure, as shown above). The word "sinusoid" is to be replaced by "dc signal," and *two* independent samples of the rectifier output are to be taken for every one that is prescribed in the narrow-band case. From a noise or signal that (at the input to the square-law rectifier) has a flat spectrum of width $W$ cps, independent samples are obtainable at interval $1/W$ seconds in the narrow-band case, and at interval $1/2W$ in the lowpass (dc) case,[7] so that for large $m$ the averaging time is $m/W$ seconds in both cases. The amplitude of the "fading" signal of Case 2 with its Rayleigh distribution is equivalent in the lowpass case to the $RMS$ value formed from two independent samples of a Gaussian noise-like signal. The case of only one sample could be calculated but would give even less satisfactory detection.

4. NUMERICAL EXAMPLES (RADAR)

The following hypothetical example illustrates the concepts involved. Suppose one has a radar that searches in range and azimuth with a beam 3° wide, a pulse rate of 5,000 per second, a pulse length of 1 microsecond, and a scan rate of 10 per minute. On the average one false detection in 15 minutes can be tolerated. Pulses returned by a target during a single scan are averaged (after rectification). What signal-to-noise ratio is required to give a probability of 90 per cent of detection in one scan?

This is the narrow-band case. The pulse length is usually about equal to the reciprocal of the bandwidth and so can be taken as the sampling interval. 200 samples of the radar return could then be taken in the 0.0002 sec. elapsing between consecutive transmitted pulses. Suppose that 150 of these are relevant, giving 150 values of range that can be distinguished. Similarly the azimuth scan of 360° is covered by 120 beam-widths (this does not assume that the azimuth accuracy is no better than 3°). If a factor of 2 is allowed for overlap (in range or azimuth or both) among the averages, the number of decisions per scan is $2 \times 120 \times 150 = 36,000$, or 5,400,000 in a period of 15 minutes. So $\alpha = 1/5,400,000 = 2 \times 10^{-7}$. In each scan the beam is on a target for 0.05 sec., the time required to turn 3°. In this time 250 pulses are transmitted; this is the value of $m$,

---

[7] The latter interval is well-known; the former is twice as long because the smoothing accompanying rectification in the narrow-band case makes the sampling equivalent to the sampling of $n_1(t)$ and $n_2(t)$, and their bandwidth is not $W$ but only $W/2$.

since returns from different ranges and/or azimuths are averaged separately (e.g., by a PPI scope). The 250 samples in an average are thus not consecutive in time in this case.

The values of $S/N$ can be read just above the curves for $\alpha = 10^{-6}$ (which are the solid curves). For a non-fading return, Fig. 3 then gives $S/N$ equal to $-4.2$ db for $\beta = 50$ per cent and $-2.2$ db for $\beta = 99$ per cent; by interpolation, $S/N$ is about $-3$ db for $\beta = 90$ per cent. It should be remembered that $S$ is the peak (not the average) signal power. For a fading return, Fig. 2 gives $S/N$ equal to 5.5 db. Use of the latter figure would imply that the returned pulses had virtually the same amplitude during the time of 0.05 sec. during each scan when the beam was on the target (but still had a random amplitude when viewed for a sufficiently longer period). On the other hand, use of $S/N$ equal to $-3$ db would imply that the values of the signal at interval 0.0002 sec. were virtually independent samples. These conditions may also be expressed in terms of the width of the lines or narrow bands, if any, in the spectrum of the signal (the spacing between the lines is the pulse repetition rate of 5000 cps, which is irrelevant here). If the width is 5000 cps or more (i.e., the spectrum continuous), one has $S/N$ equal to $-3$ db; if the width is 20 cps or less, one has $S/N$ equal to 5.5 db. The intermediate region is rather extensive, corresponding to a factor of $m = 250$. It might be explored in an approximate manner by the formulas (13g) and (13k).

If the line-width is not much more than 20 cps, but not less than say 1 cps, then conclusion (A) of Section 8 shows that the probability of detecting the airplane within 6 seconds ($=1$ scan, as first stated) is increased by increasing the scan rate to as much as 1 per second, so that a number of scans are completed within the 6 seconds.

Suppose now that one has a radar receiver with an antenna pattern 3° wide scanning in azimuth at 10 rps as before, whose object is to detect a distant radar C-W transmitter whose frequency is only known to lie within a band 150-mc wide. The detection is accomplished by passing the signal through a filtering device that passes 150 different frequency bands of 1 mc each in succession (one at a time). The problem is then numerically the same as before, the search in frequency having replaced the search in range. There is the difference that here the samples in any average are presumably all consecutive in time, while previously they were taken at intervals of 0.0002 sec., the time between transmitted pulses. The search in frequency also occurs in the determination of the velocity of an airplane by its Doppler frequency, but the bandwidth of the frequency-analyzer is then much less than 1 megacycle, and not nearly so many independent samples can be obtained.

If the signal is applied simultaneously to a bank of filters (e.g., vibrating reeds), the number of decisions is the same as before, but the saving of time would permit a corresponding increase in the number $m$ of samples in an average.

B. McMillan has pointed out that with long-range search radars whose resolution is much better in range than in the other coordinates, some of the range resolution might profitably be sacrificed in order to increase the pulse length. Although the longer pulse length could be used to increase the value of $m$, it would be preferable to decrease the receiver bandwidth instead, and thus decrease the noise power $N$ in the same ratio.

## CONCLUSIONS

### 5. THREE MAJOR EFFECTS

The curves show that the steady sinusoid is the easiest signal to detect, as one would expect, while the fading sinusoid is the most difficult to detect reliably. It is therefore very desirable to avoid the latter situation if reliable detection is wanted, either by reducing the severity of the fading and so moving toward Case 1, or by sampling the output over a period of time long enough to average out some of the fluctuation in the signal amplitude, and so moving toward Case 3 (noise-like signal). The latter is of intermediate difficulty of detection, and coincides with Case 2 when $m = 1$ and approaches coincidence with Case 1 as $m \rightarrow \infty$ . The difficulty with the fading sinusoid is this: If the average signal power is equal to that which gives good detection in Case 1, little reliability is left to be gained in those detections where additional signal power happens to be available, but much may be lost when the signal power happens to be lower than the average.

A second major conclusion results from the steepness of the left-hand portion of the curves for $\beta = 99$ per cent in Fig. 3: If reliable detection of a noise-like signal is required, it is highly desirable that at least 4 or 5 *independent* samples of the signal be available and made use of. In fact, if a 99 per cent chance of detection is desired, 4 samples require only $\frac{1}{50}$ of the signal power, or $\frac{1}{12}$ of the energy required by one sample. The principle is well-known in connection with search radars, which are designed to return about four pulses from a target.

When $m$ is significantly greater than unity, and the samples being averaged are adjacent to one another in time, one has the option of using a more selective filter and thus doing the averaging linearly (with preservation of signs) before rectification rather than after. The former is well-

known to be the more effective in Cases 1 and 2. (Case 3, on the contrary, represents the situation that results when the selectivity of the filter has already been made as great as it can profitably be.) Let the time available for the detection be fixed. Then $m$, the number of independent noise samples available, is proportional to the filter bandwidth; the latter is inversely proportional to $S/N$, provided the noise spectrum can be regarded as flat in the region considered, and the filter does not reject any of the signal frequencies. Thus $m$ and $S/N$ are inversely proportional, and one is operating along a 45° line in Figures 1 and 2 (both $m$ and $S/N$ being plotted logarithmically). Evidently $\beta$ is maximized by keeping $m$ (and hence the bandwidth) small.

## 6. MATCHING THE DURATIONS OF SAMPLE AND SIGNAL

It is fairly obvious intuitively that the chance of detection would be greatest if some one average coincided exactly with the period during which the signal was present. (The signal may be intermittent, as with the first radar of Section 4; if so, the sampling also should be intermittent in the same pattern.) There is no way to insure this if the duration of the signal is not known in advance, but if it is known, the desired coincidence could be approximated to some extent by having the averages overlap; e.g., run from 0 to 1 time unit, $\frac{1}{2}$ to $1\frac{1}{2}$, 1 to 2, $1\frac{1}{2}$ to $2\frac{1}{2}$, etc. This increase in $n$ necessitates use of a smaller value of $\alpha$ and hence requires an increase in signal strength, but the latter increase is insignificant: When $\beta \geqq 0.50$, the curves show that it is 0.3 db or less for a factor of 2 in $\alpha$. Similarly, the passbands may overlap when one searches in frequency, and the radar antenna patterns may overlap when one searches in direction.

If the averages include fewer samples (though all the samples in some average contain the signal), the effect in any of the Figs. 1 to 4 is to move to the left along a horizontal line. The detection probability $\beta$ is decreased because there is no increase in the concentration of the signal energy, while the fluctuation of the noise is averaged out less effectively.

If the averages include too many samples, then some of these samples will always consist of pure noise, even when a signal is present, and the concentration of the signal energy is reduced relative to that of the noise. If a number $m'$ of samples containing the signal are increased to $m$ by adding $m - m'$ samples of pure noise, the effect in Figs. 1 and 2 is to move down and to the right along a 45° line, since the signal-to-noise ratio is effectively reduced by the same factor $m/m'$ by which the number of samples is increased. Evidently this decreases $\beta$.

For the noise-like signal of Fig. 3, new calculations are required when $m' < m$, and the broken line in Fig. 4 shows some results obtained by the chi-square approximation of equation (13g), with $m'' = m'$. This curve is plotted against $m$ (assumed for this example to be twice $m'$) and so shows the effect of halving the duration of the signal. To show the effect of adding an equal number of pure noise samples (the duration of the signal being kept fixed), the curve should be plotted against $m'$, which is equivalent to moving the curve to the left a distance corresponding to a factor of 2. The resulting curve still lies above the solid curve (for which $m = m'$) for a noise-like signal, but only by about 1 db, which is about the same loss that one finds in Figs. 1 and 2.

It is natural to ask which one should guard against the more, taking too many samples or too few? The penalty for the former increases rather slowly as we have just seen. The latter is a more serious mistake (especially for the noise-like signal of Fig. 3) if one remains limited to a single opportunity to detect the signal. However, if the signal samples missed by one average are included in another and so give another opportunity for detection, the loss is reduced but not eliminated, as shown below in connection with repeated searches.

## 7. DISTRIBUTION OF A FIXED AMOUNT OF SIGNAL ENERGY

The preceding discussion is closely related to the question of the optimum utilization of a fixed amount of signal energy; assuming that $m = m'$, is it better to have a big pulse of signal occurring in only one sample point, or a lower signal power occurring in a proportionately greater number of samples? Here the product of $S/N$ and $m$ is held constant, and one moves along a 45° line in *all* cases, including Case 3. Reference to the curves shows that with a sinusoidal signal one should take $m = 1$, or at least concentrate the signal energy enough to make the power ratio $S/N$ very much above unity. With a noise-like signal, the same conclusion ($m = 1$) holds under conditions such that $\beta \leqq 50$ per cent, but when better detection is possible a larger value of $m$ may be preferable. For example, take the solid curves ($\alpha = 10^{-6}$) in Fig. 3 and suppose that a one-pulse signal makes $S/N = 18$ db. The corresponding 45° line is tangent to the curve $\beta = 99$ per cent at about $m = 15$. Thus the maximum detection probability is 99 per cent and is attained by taking $S/N = 6.5$ db, and so distributing the signal energy among 15 independent samples. Detection with 99 per cent probability by means of a single sample would require 20 times as much energy.

## 8. REPEATED SEARCHES

Repeated searches have two important characteristics: (a) There is no carry-over of data from one search to the next. This is not desirable, as we shall see, and would expect on theoretical grounds; but it may be unavoidable. (b) They involve repeated opportunities to detect the signal (or its source).

Repeated searches arise non-trivially when the signal recurs periodically, or can be made to do so by the detecting agency. This assumes that it is sufficient to detect the signal at one of its appearances, so that it is not necessary to regard each appearance as a separate signal. Examples are a signal of long duration but unknown frequency, and the radar return to a search radar from an airplane whose position in space is unknown. Then if one search of frequency or of space does not detect the signal, it may be possible to repeat the search one or more times before the source of the signal disappears.

An equivalent of repeated searches arises when the averaging time is short enough so that two or more averages fall within the duration of the signal. The results below apply to this case and show that a longer average (i.e., carry-over of data from one search to the next) is somewhat preferable. These results do not apply to the case where over-lapping averages give several opportunities to detect the signal, because such averages have some data in common; but it has already been shown that such overlapping is desirable, other things being equal.

Suppose now that the signal power is constant and the time consumed is proportional to the product of $m$ and the number $\lambda$ of searches. How can one find the signal most quickly? If the probability $\beta$ of detection in one search is $1/2$, the probability that exactly $\lambda$ searches are required for detection is $1/2^\lambda$, and this also happens to be the probability that the signal remains undetected after $\lambda$ searches. Thus one would need to make nearly 7 searches before one could conclude with 99 per cent assurance that no signal was present. However, if the signal is in fact present, the *average* number of searches required to detect the signal with 99 per cent assurance is

$$1/2 + 2/2^2 + 3/2^3 + 4/2^4 + 5/2^5 + 6/2^6 + 7/2^6 \doteq 2$$

searches if the superfluous part of the successful search is included, or

$$(1/2)\cdot(1/2) + (3/2)\cdot(1/2^2) + \cdots$$
$$+ (11/2)(1/2^6) + (13/2)(1/2^6) \doteq 3/2$$

searches if the search is terminated the moment the signal is detected.

The latter result is the less favorable for repeated searching, since it is 3 times the result ($\frac{1}{2}$) obtained with 99 per cent detection in a single search (which on the average could perhaps be terminated, due to detection of the signal, when it was half completed).

The alternative way of increasing $\beta$ from 50 per cent to 99 per cent is to increase the value of $m$ by a factor lying between 2.3 and 5 for Figs. 1 and 3, while the fading sinusoid of Figure 2 requires a factor of the order of 1,000. Comparing these with the factors 2 (or 3) and 7 obtained previously gives the following conclusions, which have been confirmed by considering some other values of $\beta$:

A. With a fading sinusoidal signal, repeated searching is extremely advantageous, *provided* it is true as assumed that independent samples of the signal are obtainable from search to search but not within a search.

B. Repeated searching may have a small advantage in Cases 1 and 3 (the steady sinusoid and the noise-like signal) provided (1) each search has $\beta \geq 50$ per cent, and (2) the criterion is the average time required to detect the signal when a signal is present, and one is not concerned with the (increased) time required to conclude that no signal is present (which of course is also the time required for the detection of some of the signals).

C. Repeated searching is somewhat disadvantageous in other cases, as one would expect; e.g., in Cases 1 and 3, if a fixed amount of time is available, the greatest probability of detection is achieved by using all of the time for a single search, rather than dividing it among several less sensitive searches.

## Mathematical Appendix

### 9. Pure Noise (and Noise-Like Signal = Case 3)

If narrow-band Gaussian noise is applied to a square-law rectifier, the value of the output at any instant has a Rayleigh or exponential distribution, as is well known. The average of $m$ such values, taken far enough apart to be virtually independent has a chi-square ($\chi^2$) distribution with $2m$ degrees of freedom. The standard form of the distribution used in tables chooses the units so that the mean of the distribution is $2m$, whereas we want the mean to equal the noise power $N$. The exact value of $k$ can therefore be found by means of the relation

$$Pr(N\chi^2_{2m}/2m > (k + 1)N) = \alpha \tag{9a}$$

or

$$Pr(\chi^2_{2m} > 2m(k + 1)) = \alpha$$

where $\alpha$ is the expected ratio of false detections to total number of decisions, when no signal is actually present.

The usual short tables of $\chi^2$ (e.g., that in Fisher and Yates' Statistical Tables or that of Hartley and Pearson in Biometrika **37**, p. 313) are useful but not entirely adequate for the needs of the problem. The most extensive table is that of K. Pearson,[8] which goes to $m = 50$. In terms of the function I tabulated by Pearson,

$$Pr(\chi^2_{2m} > A) = 1 - I(A/2\sqrt{m}, m-1) \qquad (9b)$$

For $m > 50$ it is sufficient to use the expansion[9]

$$\chi^2_{2m} = 2m + 2um^{1/2} + \tfrac{2}{3}(u^2 - 1) + \tfrac{1}{18}(u^3 - 7u)m^{-1/2} - \cdots \qquad (9c)$$

Here $u$ is the standard normal[10] deviate. Thus if $u_\alpha$ is defined by

$$\alpha = Pr(u > u_\alpha) = \int_{u_\alpha}^{\infty} e^{-u^2/2} \, du/\sqrt{2\pi}$$

one has

$$k = u_\alpha/\sqrt{m} + (u_\alpha^2 - 1)/3m + \cdots \qquad (9d)$$

If a noise-like signal is present, this is mathematically the same as an increase in the average noise power from $N$ to $N + S$, the critical power level remaining at $(k + 1) N$. The probability of exceeding the critical level is now the probability $\beta$ of a true detection. Thus

$$\beta = Pr[(N + S)\chi^2_{2m}/2m > (k + 1)N]$$
$$= Pr[\chi^2_{2m} > 2m(k + 1)/(1 + S/N)] \qquad (9e)$$

Writing $\chi^2_{2m}(\beta)$ for the number that is exceeded by the variable $\chi^2_{2m}$ with probability $\beta$, one gets from (9a) and (9e)

$$1 + S/N = \chi^2_{2m}(\alpha)/\chi^2_{2m}(\beta) \qquad (9f)$$

For large $m$ (and only in that case), (9c) then gives

$$S/N \doteq (u_\alpha - u_\beta)\left[\frac{1}{\sqrt{m}} + \frac{u_\alpha - 2u_\beta}{3m}\right] \qquad (9g)$$

---

[8] K. Pearson, Tables of the Incomplete Gamma-Function, Biometrika Office, London, 1934.

[9] The author first discovered this formula in T. Lewis, Biometrika **40**, p. 424, 1953; but S. O. Rice has pointed out that G. A. Campbell published it as early as 1923 in the B.S.T.J., **2**, page 95, in connection with Poisson distribution. Other references and inversion formulas are given by John Riordan. Ann. Math. Stat., **20**, p. 417, 1949.

[10] "Normal" is a synonym for "Gaussian." However a new application of the distribution is involved here, and so there is no disadvantage in making the conventional change in terminology.

where the standard normal deviate $u_\beta$ is defined by $\beta = Pr(u > u_\beta)$. It is negative when $\beta > 50$ per cent.

Wherever it occurs in equations, $S/N$ naturally means the ratio of $S$ to $N$, and not the value in db, which is $10 \log_{10} S/N$.

## 10. NOISE PLUS STEADY SINUSOID (CASE 1)

When a pure sinusoidal signal of constant amplitude is added to the noise, the distribution of the rectifier output has what is called a non-central chi-square distribution, which has been little tabulated. Fortunately the normal distribution gives a fair approximation in this case (with the probable exception of small values of $\beta$, which are of little interest anyway). More accuracy could be obtained at the cost of additional labor.[11]

Before rectification, the signal plus noise can be represented, as remarked in Section 3, by the expression

$$\sqrt{2S} \cos (\Omega t + \varphi) + n_1(t) \cos (\Omega t + \varphi) + n_2(t) \sin (\Omega t + \varphi) \quad (10a)$$

where the two Gaussian noise variables $n_1(t)$ and $n_2(t)$ are independent and have zero means and a common variance $N$. After rectification and smoothing, half the square of the envelope is obtained, namely

$$[(\sqrt{2S} + n_1(t))^2 + n_2(t)^2]/2$$

Its mean value is $N + S$, and its variance is the sum of the variances of the terms in the expanded form, namely $0 + 2NS + N^2/2 + N^2/2$ or $2NS + N^2$. The average of $m$ such independent variables has the same mean $N + S$ but the variance $N(N + 2S)/m$.

The critical value $(k + 1)N$ of the rectifier output is reduced to standard units by subtracting the mean output $N + S$ and dividing by the square root of the variance, giving

$$u_\beta = (k - S/N) \sqrt{\frac{m}{1 + 2S/N}}$$

After $u_\beta$ is found from

$$\beta = Pr(u > u_\beta) = \int_{u_\beta}^{\infty} e^{-u^2/2} \, du / \sqrt{2\pi}$$

by using a table of the normal distribution, one can solve the preceding

---

[11] P. B. Patnaik, Biometrika **36**, p. 202, 1949.

equation for $S/N$, giving

$$\frac{S}{N} \doteq k + \frac{u_\beta^2}{m} - \frac{u_\beta}{\sqrt{m}} \left[ 1 + 2k + \frac{u_\beta^2}{m} \right]^{1/2} \tag{10b}$$

This is then the signal-to-noise ratio giving the probability $\beta$ of detection.

By using (9d) and taking $1 + u_\alpha/\sqrt{m}$ as an approximate value of the radical, one obtains from (10b)

$$\frac{S}{N} \doteq \frac{u_\alpha - u_\beta}{\sqrt{m}} + \frac{u_\beta(u_\beta - u_\alpha) + (u_\alpha^2 - 1)/3}{m} \tag{10c}$$

This result differs most (a little over 1 db) from the curves of Figure 1 when $m = 1$ and $\beta = 0.01$, but (10b) itself is not very accurate in that case.

## 11. NOISE PLUS FADING SINUSOID (CASE 2)

In this case the signal plus noise has the form

$$[s + n(t)] \cos \Omega t + [s' + n'(t)] \sin \Omega t \tag{11a}$$

where $n'(t)$ is *not* a derivative, and the components $s$ and $s'$ of the signal amplitude are essentially constant during the period that is averaged over, although like $n$ and $n'$ they are Gaussian variables of mean zero. Detection is accomplished by means of the expression

$$R = \frac{1}{2m} \left[ \sum_1^m (s + n_i)^2 + \sum_1^m (s' + n_i')^2 \right] \tag{11b}$$

where the $n_i$ and $n_i'$ are (independent) values of $n(t)$ and $n'(t)$ at successive sampling points.

Now (11b) can be written

$$\begin{aligned} 2R = {} & \left( s + \frac{1}{m} \sum_1^m n_i \right)^2 + \left( s' + \frac{1}{m} \sum_1^m n_i' \right)^2 \\ & + \frac{1}{m} \sum n_i^2 - \left( \frac{1}{m} \sum n_i \right)^2 + \frac{1}{m} \sum n_i'^2 - \left( \frac{1}{m} \sum n_i' \right)^2 \end{aligned} \tag{11c}$$

as in the analysis of variance in statistics. The sum of two squares in the first line of (11c) is then distributed as $(S + N/m)\chi_2^2$; the second line is distributed as $\chi_{2m-2}^2 \cdot N/m$; and these portions are independent of one another, because they correspond to the mean and the variance respec-

tively of the sample of $n_i$ (or $n_i'$), and the sample mean and variance are known to be independent in the Gaussian case. Here the signal power $S$ is the variance of each of $s$ and $s'$, and $N$ is the variance of each of the $n_i$ and $n_i'$.

For detection $R$ must exceed $N(1 + k)$. Hence if $S/N$ is denoted by $r$, one has

$$\beta = Pr[(1 + mr)\chi_2^2 + \chi_{2m-2}^2 > 2m(1 + k)] \tag{11d}$$

Since $(1 + mr)\chi_2^2$ has the c.d.f. $1 - e^{-t/2(1+mr)}$ and $\chi_{2m-2}^2$ has the density function

$$e^{-t/2}t^{m-2}/2^{m-1}\Gamma(m - 1),$$

$\beta$ is obtained by convolution as

$$\beta = 1 - \int_0^z [1 - e^{-(z-t)/2(1+mr)}]e^{-t/2}t^{m-2}\,dt/2^{m-1}\Gamma(m - 1)$$

$$= Pr[\chi_{2m-2}^2 > z] + e^{-z/2(1+mr)}\left(\frac{mr}{1 + mr}\right)^{1-m} Pr\left[\chi_{2m-2}^2 < \frac{mrz}{1 + mr}\right] \tag{11e}$$

where $z = 2m(1 + k)$.

An excellent approximation is obtained by replacing the variable $\chi_{2m-2}^2$ by its mean value $2m - 2$; this is reasonable, since the ratio of the variances of $\chi_{2m-2}^2$ and $(1 + mr)\chi_2^2$ is $(m - 1)/(1 + m\ S/N)^2$, which is very small throughout the interesting area ($\beta \geqq 50$ per cent). (11d) then gives

$$\beta \doteq \exp\left[-(1 + km)/(1 + m\ S/N)\right] \tag{11f}$$

This result will be generalized in (13k) below. Using (9d) then gives

$$\frac{S}{N}\ \ell n\ (1/\beta) \doteq \frac{u_\alpha}{\sqrt{m}} + \frac{(u_\alpha^2 + 2)/3 - \ell n\ (1/\beta)}{m} \tag{11g}$$

to a little lower order of accuracy when $m$ is small, and still with $\beta \geqq 50$ per cent.

Before the foregoing results were obtained, the curves of Fig. 2 had already been calculated[12] on IBM equipment by the rather tedious process of integrating the results for Case 1 with respect to the random signal

---

[12] By Mrs. L. R. Lee, at the request of the author.

strength. This gives

$$\beta \doteq \int_0^\infty \Phi[(xS/N - k)\sqrt{m/(1 + 2xS/N)}]p(x)\,dx \qquad (11\text{h})$$

for Case 2, where

$$p(x) = e^{-x}$$

and

$$\Phi(u) = \int_{-\infty}^u e^{-t^2/2}\,dt/\sqrt{2\pi}$$

If the fading is due to interference between two sinusoidal signals of powers $S_1$ and $S_2 (S_1 + S_2 = S)$ and very slightly different frequencies, it can be shown that the appropriate form for $p(x)$ is $1/\pi\sqrt{A^2 - (1 - x)^2}$ for $1 - A < x < 1 + A$, and zero elsewhere, where $A = 2\sqrt{S_1 S_2}/(S_1 + S_2)$. Here $x$ can approach zero only in the special case $S_1 = S_2$. Since one or both of $S_1$ and $S_2$, and hence $A$, is likely to be a random variable, one apparently cannot say a priori that any particular distribution $p(x)$ is the appropriate one in this case.

## 12. COMPARISON OF 12 CASES

The three principal cases defined in Section 2 and analyzed in the last three sections may be expanded to 12 by considering dc signals as well as ac, and considering also a second method of detection. The original 3 cases will be labeled "ac envelope." (More precisely, half of its square is the quantity averaged.) If the envelope is not isolated but instantaneous squared values are sampled, the situation is "ac instantaneous." A precise analysis of these latter cases would apparently be difficult. As indicated, quadratic rectification is assumed for all 6 of the ac cases. On the other hand, for the 6 dc cases, distinguished as rectified and unrectified, it is the instantaneous sampling that is assumed throughout.

The reduction of the rectified dc cases to the standard ac envelope cases has been indicated at the end of Section 3, and will be demonstrated at the beginning of the following section. The unrectified dc cases, which represent the better way to detect a steady or fading dc signal, involve nothing but the simple Gaussian distributions specified in Table I. If the number of samples averaged is denoted by $2m$, one easily finds the theoretically exact relations

$$S/N = (u_\alpha - u_\beta)^2/2m \qquad (12\text{a})$$

$$S/N = (u_{\alpha/2}^2/u_{\beta/2}^2 - 1)/2m \qquad (12\text{b})$$

$$S/N = u_{\alpha/2}^2/u_{\beta/2}^2 - 1 \qquad (12\text{c})$$

## TABLE I — VARIABLES TO BE AVERAGED

|  | DC Unrectified | DC Rectified |
|---|---|---|
| Steady............. | $\sqrt{S} + n$ | $S + 2n\sqrt{S} + n^2$ |
| Fading }<br>Noise-like } ............. | $s + n$ | $s^2 + 2ns + n^2$ |

|  | AC Envelope |
|---|---|
| Steady.............<br>Fading }<br>Noise-like } ......... | $S + n\sqrt{2S} + (n^2 + n'^2)/2$<br>$[(s + n)^2 + (s' + n')^2]/2$ |

|  | AC Instantaneous |
|---|---|
| Steady............. | $[(\sqrt{2S} + n)\cos\varphi + n'\sin\varphi]^2$ |
| Fading............. | $[(s + n)\cos\varphi + (s' + n')\sin\varphi]^2$ |
| Noise-like............. | Ditto or $s^2 + 2ns + n^2$ |

for the steady, fading, and noise-like unrectified dc (or low-pass) signals. As usual, $u_p$ is defined by

$$\int_{u_p}^{\infty} e^{-t^2/2}\, dt / \sqrt{2\pi} = p$$

In (12a) it is assumed that the steady dc signal of magnitude $\sqrt{S}$ is of known sign. Otherwise it would be necessary, as with fading and noise-like signals, to be prepared to detect both positive and negative dc signals, and quantities like $u_{\alpha/2}$ enter in place of $u_\alpha$ via the relation

$$\int_{u_{\alpha/2}}^{\infty} + \int_{-\infty}^{-u_{\alpha/2}} = \alpha.$$

As $\beta$ approaches unity, $u_{\beta/2}$ approaches zero and $S/N$ rapidly approaches infinity for the fading and the noise-like signal. In the latter case (12c), averaging has no effect ($m$ does not appear), because the unrectified signal averages to zero as fast as the noise does. In fact, with instantaneous sampling of a noise-like signal, there is no essential difference between the lowpass ("dc") and narrow-band (ac) cases.

If $n$, $n'$, $s$, $s'$ are independent Gaussian variables with zero means and variances $N$, $N$, $S$, $S$ respectively, and $\varphi$ is uniformly distributed between 0 and $2\pi$, the variables that are averaged in the detection process have the forms given in Table I in the various cases, when a signal is present (if not, put $s$, $s'$, and $S$ equal to zero).

Table II gives the variances of the variables of Table I for the cases

TABLE II — VARIANCES

| | DC Unrectified | DC Rectified |
|---|---|---|
| Steady........................ | $N$ | $2N^2 + 4NS$ |
| Noise-like.................... | $N + S$ | $2N^2 + 4NS + 2S^2$ |
| | AC Envelope | AC Instantaneous |
| Steady........................ | $N^2 + 2NS$ | $2N^2 + 4NS + \frac{1}{2}S^2$ |
| Noise-like.................... | $N^2 + 2NS + S^2$ | $2N^2 + 4NS + 2S^2$ |

of a steady or a noise-like signal. The means are $\sqrt{S}$ and zero, respectively, for the dc unrectified variables, while all the others have the power $N + S$ as their mean.

Table III gives the data for the fading signal. The reducible variance and the mean are calculated under the condition that $s$ and $s'$ are fixed; the irreducible variance is then the variance of the resulting mean when $s$ and $s'$ do vary. The distinction between the two variances is that by averaging, the first is reduced in the usual way (divided by the number of samples), while the second retains its full value. The variance already given for the steady signal can be derived from the reducible variances above by replacing $s^2$ and $s'^2$ by $S$. The variances already given for the noise-like signal are equal to the irreducible variance plus the expected value of the reducible variance.

The steady ac instantaneous case will serve to illustrate the calculation of the variance. If the rectifier output is written in the form

$$S + n\sqrt{2S} + (n^2 + n'^2)/2 + (n'\sqrt{2S} + nn') \sin 2\varphi$$
$$+ [S + n\sqrt{2S} + (n^2 - n'^2)/2] \cos 2\varphi$$

it may be verified that the various terms are uncorrelated with one another, so that their individual variances may be added. Also $E \sin 2\varphi = E \cos 2\varphi = 0$, $E \sin^2 2\varphi = E \cos^2 2\varphi = \frac{1}{2}$, $En^4 = En'^4 = 3N^2$,

TABLE III — FADING SIGNALS

| | Reducible Variance | Mean | Irreducible Variance |
|---|---|---|---|
| DC unrectified...... | $N$ | $s$ | $S$ |
| DC rectified........ | $2N^2 + 4Ns^2$ | $N + s^2$ | $2S^2$ |
| AC envelope........ | $N^2 + N(s^2 + s'^2)$ | $N + (s^2 + s'^2)/2$ | $S^2$ |
| AC instantaneous... | $2N^2 + 2N(s^2 + s'^2) + (s^2 + s'^2)^2/8$ | $N + (s^2 + s'^2)/2$ | $S^2$ |

var $n^2 = $ var $n'^2 = 2N^2$. These last relations are not obvious, but are a special case ($t = u$) of the relation cited below just prior to (15e). The variance is now obtained as

$$0 + 2NS + 4N^2/4 + \tfrac{1}{2}E(2Sn'^2 + n^2n'^2)$$
$$+ \tfrac{1}{2}E[S^2 + 2Sn^2 + (n^4 - 2n^2n'^2 + n'^4)/4]$$
$$= 2N^2 + 4NS + S^2/2$$

The variances tabulated indicate the relative merit of the various cases. If $m$ is significantly greater than unity, the anomalous case of an unrectified noise-like signal represented by (12c) gives the poorest detection; next come the fading signals with their irreducible variances. The steady dc unrectified signal with $S/N$ proportional to $m^{-1}$ (by 12a) is the easiest to detect. The steady signals are in all cases more easily detected than the noise-like at the higher signal strengths. However, instantaneous sampling of a steady ac signal loses some of this advantage; a term $S^2/2$ appears, in addition to the expected doubling of the variance of a single sample.

13. DETECTION OF AN ARBITRARY GAUSSIAN SIGNAL

General formulas will now be derived which include Cases 1, 2, and 3 and also give approximations to intermediate cases. It will be more convenient to deal with the Gaussian sample values $s_i + n_i$ of the signal plus noise in the dc or lowpass case, and $2m$ of them will be taken so that the results obtained will have the same form as those already given for the ac or narrow-band case. Detection is then accomplished by means of the variable

$$R = \sum_1^{2m} (s_i + n_i)^2/2m = \sum_1^{2m} (s_i^2 + 2s_in_i + n_i^2)/2m \qquad (13a)$$

It is assumed that the $n_i$ are independent of one another and of the $s_i$, while the $s_i$ may be constant or mutually correlated. Also $En_i = 0$ and $En_i^2 = N$, a constant.

An orthogonal transformation (rotation in $2m$-dimensional space) can be used to show that (13a) has the same distribution as

$$\sum_1^{2m} (h + n_i')^2/2m \qquad (13b)$$

or

$$\left[\sum_1^m (h\sqrt{2} + n_i'')^2 + \sum_{m+1}^{2m} n_i''^2\right]\Big/ 2m \qquad (13c)$$

where $h^2 = \sum_1^{2m} s_i^2/2m$, and the $n_i'$ or $n_i''$ are new noise variables with exactly the same properties as the original $n_i$. Evidently $h^2$ is so defined as to leave unchanged the terms independent of the noise variables, while the orthogonal transformation by definition makes $\sum n_i^2 = \sum n_i'^2 = \sum n_i''^2$. Thus it is only necessary to choose the transformation so that the original linear terms $\sum_1^{2m} s_i n_i/m$ go into $h \sum_1^{2m} n_i'/m$ or $\sqrt{2}h \sum_1^m n_i''/m$. This is always possible since all three expressions have the same norm (square root of sum of squares of coefficients), namely $h\sqrt{2/m}$.

The forms (13b) and (13c) are appropriate to the instantaneous sampling of the square of a dc signal plus noise, and the sampling of the square of the envelope of an ac signal plus narrow-band noise, respectively, the latter being the standard case. Thus the equivalence of the two cases is established. It is only necessary to observe that the number of samples taken is $2m$ in (13b) but only $m$ in (13c). The noise powers are $En_i^2$ and

$$E(n_i \cos \varphi + n_{m+i} \sin \varphi)^2 = E(n_i^2 + n_{m+i}^2)/2 = En_i^2.$$

The signal powers are $h^2$ and $E \cdot 2h^2 \cos^2\varphi = h^2$. Another consequence of (13b) or (13c) is that when the distribution of $h$ has a known form, the distribution of $R$ could be obtained from the results for Case 1 by integrating over $h$, as stated at the end of Section 11.

The mean and the variance of the general form (13a) will now be calculated. It is assumed that $Es_i = 0$ and $Es_i^2 = S$ for the first $2m'$ values of $i(m' \leq m)$, while $s_i \equiv 0$ for the other $2(m - m')$ values of $i$. Then one has $ER = N + m'S/m$. If $i$ and $j$ are any two distinct integers, the five variables $s_i^2$, $s_i n_i$, $s_j n_j$, $n_i^2$, $n_j^2$ are all uncorrelated (though not all independent), and so one has

$$4m^2 \operatorname{var} R = \operatorname{var}\left(\sum s_i^2\right) + 4 \sum \operatorname{var}(s_i n_i) + 2m \operatorname{var} n_i^2$$

Since $\operatorname{var}(s_i n_i) = NS$ or 0 and $\operatorname{var} n_i^2 = 2N^2$, this gives

$$\operatorname{var} R = \operatorname{var} h^2 + 2m'NS/m^2 + N^2/m \tag{13d}$$

with $h^2 = \sum_1^{2m'} s_i^2/2m$. In Cases 1, 2, and 3, $\operatorname{var} h^2 = 0$, $(m'S/m)^2$, and $m'S^2/m^2$ respectively. (In Case 2, $m'$ of the $s_i$ have one identical value and $m'$ have another identical value, independent of the first.) In general

$$\operatorname{var} h^2 = \sum_{-2m'+1}^{2m'-1} (2m' - |i|)\psi^2(i)/2m^2 \tag{13e}$$

where $\psi(i) = Es_j s_{i+j}$ for all $j$ for which $s_j s_{i+j} \neq 0$. This is the discrete analogue of the integral appearing in (15k) below. It is convenient to define a number $m''$ (which need not be an integer) such that $\operatorname{var} h^2 =$

$(m'S/m)^2/m''$. Then putting $S' = m'S/m$ gives

$$ER \doteq N + S', \qquad \mathrm{var}\, R = S'^2/m'' + (2NS' + N^2)/m. \quad (13\mathrm{f})$$

The mean and variance of $R$ are identical with those of $[S'^2/m'' + (2NS' + N^2)/m]/2(N + S')$ times a chi-square variable with

$$2\bar{m} = 2(N + S')^2/[S'^2/m'' + (2NS' + N^2)/m]$$

degrees of freedom. So one has approximately

$$\beta = Pr(R > (k + 1)N)$$
$$\doteq Pr\{\chi^2_{2\bar{m}} > 2(k + 1)N(N + S')/[S'^2/m'' + (2NS' + N^2)/m]\} \quad (13\mathrm{g})$$
$$\doteq Pr\{u > (kN - S')/\sqrt{S'^2/m'' + (2NS' + N^2)/m}\}$$

The last line can be used to verify that for fixed $m$, the optimum value of $m'$ is $m$, and for fixed $m'$, the optimum value of $m$ is $m'$. Defining $u_\beta$ by $Pr(u > u_\beta) = \beta$ and solving for $S/N = mS'/m'N$ gives

$S/N$
$$= \frac{mk + u_\beta^2 - u_\beta\sqrt{(2k + 1)m + k^2m^2/m'' - u_\beta^2(m/m'' - 1)}}{m'(1 - u_\beta^2/m'')} \quad (13\mathrm{h})$$

The approximation breaks down if $u_\beta^2 \geqq m''$.

It can be shown that by replacing the symbols $m'$, $m''$, and $S$ in the above formulas by $m'/2$, $m''/2$, and $2S$, one obtains the results for the narrow-band case in which the signal is a noise-modulated carrier $s(t)$ $\cos \Omega t$, while the noise as usual has the form $n(t) \cos \Omega t + n'(t) \sin \Omega t$.

In Cases 1, 2, and 3, $m'' = \infty$, 1, and $m'$ respectively; also, $m'$ has been assumed equal to $m$. In the dc or lowpass case, $m''$ could be as small as $\frac{1}{2}$. One has $m \geqq m'$ always, and $m' \geqq m''$ when the signal is Gaussian (i.e., contains no steady sinusoidal or steady dc component).

When $m''$ is very small, (13g) and (13h) are of low accuracy, and it is much better to use (11e) or (11f), or the following generalization of them. The derivation proceeds as if $m'/m''$ were an integer, although this is probably not a necessary condition for the usefulness of the results. The averaged rectifier output is represented by

$$2mR = \sum_{i=1}^{2m''} \sum_{j=1}^{m'/m''} (s_i + n_{ij})^2 + \sum_{1}^{2m-2m'} n_i^2 \quad (13\mathrm{i})$$

where the $n_{ij}$ are $2m'$ independent Gaussian noise variables accompanying the signal variables $s_i$ (which have only $2m''$ independent values), and the $n_i$ are $2m - 2m'$ additional noise variables that are not accom-

panied by a signal. Of course, the signal variables would not in reality fall into independent sets of identical values, so that an approximation enters here even if $m'/m''$ is an integer.

(13$i$) may be written as

$$2mR = (m'/m'') \sum_{i=1}^{2m''} (s_i + \bar{n}_{i+})^2 + \sum_{i=1}^{2m''} \sum_{j=1}^{m'/m''} (n_{ij} - \bar{n}_{i+})^2$$

$$+ \sum_{1}^{2m-2m'} n_i^2 = (N + m'S/m'')\chi^2_{2m''} + N\chi^2_{2m-2m''} \qquad (13j)$$

where

$$\bar{n}_{i+} = (m''/m') \sum_{j=1}^{m'/m''} n_{ij}$$

$$S = \text{var } s_i$$

$$N = \text{var } n_{ij} = \text{var } n_i$$

and the two $\chi^2$ variables are independent. If the ratio $m''(m - m'')/(m'' + m'S/N)^2$ of the variances of the latter (including their multipliers) is small, one may replace $\chi^2_{2m-2m''}$ by its mean value $2m - 2m''$ and obtain

$$\beta = Pr[R > (1 + k)N]$$

$$\doteq Pr[\chi^2_{2m''} > 2(m'' + mk)/(1 + m'S/m''N)] \qquad (13k)$$

## 14. COUNTING SAMPLES ABOVE A THRESHOLD

It is sometimes suggested that instead of averaging the $m$ samples, the number of such samples exceeding some threshold might be counted and used as the detection criterion. This is equivalent to replacing the average of the $m$ samples by one of their order statistics, such as the median.

It is of interest to ask which order statistic is best to use, and how it compares in efficiency with the average. M. Schwartz (in the dissertation cited previously) made numerical calculations for the case of a steady sinusoid and $m \leq 49$, and concluded that the method of coincidences, as he called it, required $S/N$ to be about 1.4 db above that which sufficed for equal performance using averages. For the larger values of $m$ there appeared to be a small advantage in requiring less than half the samples to exceed the (suitably chosen) threshold. These results are confirmed by the following asymptotic analysis.

In the absence of a signal, a single sample of half the square of the

envelope has an exponential distribution $e^{-x/N} \, dx/N$. The same (with increased power $N$) is true in the presence of a noise-like signal, and approximately so for a sinusoidal signal of low intensity. Since $m$ is assumed to be fairly large, low-intensity signals are the interesting ones. In both of these cases detection can be based on an estimate of $N$, denoted by $\hat{N}$.

In the present method $\hat{N}$ is determined from $\hat{p} = e^{-K/\hat{N}}$, where $\hat{p}$ is the proportion of samples observed to exceed the threshold $K$. The standard deviation $\sigma_p$ of $\hat{p}$ is well-known to be $\sqrt{p(1-p)/m}$, where $p = e^{-K/N}$ = the "true" or expected value of $\hat{p}$. For large $m$ the sampling fluctuations are small and the standard deviation $\sigma'_N$ of $\hat{N}$ is approximately equal to $\sigma_p$ times $dN/dp$ or $\sigma_p \div e^{-K/N} \cdot K/N^2$ or

$$\sigma'_N = N(1-p)^{1/2}/(mp)^{1/2} \, \ell n \, (1/p) \tag{14a}$$

for counting samples as compared with

$$\sigma_N = N/\sqrt{m}$$

when the average is used. So the efficiency of the counting procedure (expressed in terms of the (inverse) number of samples required for equivalent reliability) is

$$\sigma_N^2/\sigma_N'^2 = (\ell n \, 1/p)^2 p/(1-p) \tag{14b}$$

This expression has its maximum value of 64.7 % when $p = e^{-2(1-p)}$ $\doteq 0.203$. The median $(p = \frac{1}{2})$ has an efficiency of only 48.0 %. For large $m$, the required $S/N$ varies as $m^{-1/2}$ and so the minimum loss due to counting is $10 \log_{10} 0.647^{-1/2} \doteq 1.0$ db.

In the case of the steady unrectified dc signal, the detection problem is equivalent to estimating the mean $\mu$ of a Gaussian distribution

$$\varphi[(x-\mu)/\sigma] \, dx = e^{-(x-\mu)^2/2\sigma^2} \, dx/\sigma\sqrt{2\pi}.$$

The estimate $\hat{\mu}$ of $\mu$ is determined from

$$\hat{p} = 1 - \Phi\left(\frac{K-\hat{\mu}}{\sigma}\right) = \int_K^\infty e^{-(x-\hat{\mu})^2/2\sigma^2} \, dx/\sigma\sqrt{2\pi}$$

and the same relation holds between the true values $p$ and $\mu$. Then $dp/d\mu = \varphi[(K-\mu)/\sigma]/\sigma$ and

$$\sigma'_\mu = \sigma\sqrt{\Phi(\xi)[1-\Phi(\xi)]}/\varphi(\xi)\sqrt{2m} \tag{14c}$$

where $\xi = (K-\mu)/\sigma$ and $2m$ is the number of samples. The variance $\sigma_\mu$ of the mean is $\sigma/\sqrt{2m}$, so the efficiency of the counting procedure is

$$\sigma_\mu^2/\sigma_\mu'^2 = \varphi^2(\xi)/\Phi(\xi)[1-\Phi(\xi)] \tag{14d}$$

The following are some values of this expression:

| $\xi$ | 0 | 0.5 | 1 | 2 |
|---|---|---|---|---|
| Efficiency | 63.7% | 58.0% | 43.9% | 13.1% |

Thus the highest efficiency of 63.7% is attained when $\xi = 0$ or the threshold $K$ equals the mean $\mu$. This corresponds to using the median as the statistic on which detection is based. This maximum efficiency is very nearly the same as that obtained in the preceding case. However, in the present case $S/N$ varies as $m^{-1}$ by (12a), so that the equivalent loss in signal strength is not 1.0 but $10 \log_{10} 1/0.637 \doteq 2.0$ db.

### 15. AVERAGING BY CONTINUOUS INTEGRATION

In practice, in place of the discrete sums of squares such as (11b) and (13a) one may have integrals such as

$$R_U = \int_0^T Z(t) \, dt/T \qquad \text{and} \tag{15a}$$

$$R_E = \int_{-\infty}^0 Z(t) e^{t/T} \, dt/T \tag{15b}$$

where $Z(t)$ is the variable to be averaged (usually a rectifier output), and the subscripts $U$ and $E$ refer to uniform and exponential weighting respectively. The purpose of this section is to calculate the mean $\mu$ and variance $\sigma^2$ of each of these expressions for various cases. It will then be shown that for a flat spectrum and $T \to \infty$, $\mu$ and $\sigma^2$ are the same as those already given for the discrete averages of samples. The results are closely related to those given by Rice.[13]

It will suffice to consider steady signals. If the signal is absent, put $S = 0$. If the signal is fading, $S$ is a random variable. If the signal is noise-like (with the same spectrum as the noise), put $S = 0$ and replace $N$ by $N + S$. The four combinations of ac and dc signals with two methods of detection (see Section 12) will be considered separately.

The exceptional case of an unrectified dc signal may be disposed of first. Then $Z(t) = \sqrt{S} + n(t)$ and $R_U$ and $R_E$ both have Gaussian distributions with mean $\sqrt{S}$. The variance of $R_U$ is

$$E \left( \int_0^T n(t) \, dt/T \right)^2$$

Writing the square as a double integral and taking the expectation under

[13] S. O. Rice, B.S.T.J., **24**, pp. 46–156, 1945, Eqns. 3.9–8 and 3.9–28. Also J. Acous. Soc. of Amer., **14**, pp. 216–227, 1943.

the integral sign gives

$$\text{var } R_U = E \int_0^T \int_0^T n(t)n(u) \, dt \, du/T^2$$

$$= \int_0^T \int_0^T \psi(t - u) \, dt \, du/T^2 = 2 \int_0^T (T - v)\psi(v) \, dv/T^2$$

(15c)

where $\psi(v) = En(t)n(t + v)$, and $\psi(0) = N$. Similarly

$$\text{var } R_E = \int_0^\infty \int_0^\infty e^{-(t+u)/T}\psi(t - u) \, dt \, du/T^2$$

$$= \int_0^\infty e^{-v/T}\psi(v) \, dv/T$$

(15d)

In the remaining three cases, $R_U$ and $R_E$ always have the power $N + S$ as their mean. Their distributions are not known exactly but they might be assumed to be distributed approximately like $(\sigma^2/2\mu)\chi^2_{2\mu^2/\sigma^2}$ where $\mu$ is the mean $N + S$ and $\sigma^2$ is the appropriate variance given below. A probably more convenient procedure is to use the $\mu = N$ and $\sigma^2 = \sigma_0^2$ (see 15k and $\ell$) for the noise alone to determine an equivalent value $N^2/\sigma_0^2$ for $m$, and then use Figs. 1–3.

Consider next the rectified dc signal, so that

$$Z(t) = S + 2\sqrt{S}n(t) + n^2(t)$$

The variance of $R_U$ is $E[R_U - S - N]^2$ or

$$E\left(\int_0^T [n^2(t) - N] \, dt/T + 2\sqrt{S} \int_0^T n(t) \, dt/T\right)^2.$$

The cross-product has zero expectation. Expressing the squares as double integrals gives

$$E\int_0^T \int_0^T ([n^2(t) - N][n^2(u) - N] + 4Sn(t)n(u)) \, dt \, du/T^2.$$

We now take the expectation under the integral sign. For

$$En^2(t)n^2(u) = \psi^2(0) + 2\psi^2(t - u)$$

see M. G. Kendall, The Advanced Theory of Statistics, Volume I, Section 3.28, equation

$$\mu_{22} = (1 + 2\rho^2)\sigma_1^2\sigma_2^2$$

This gives

$$\operatorname{var} R_U = \int_0^T \int_0^T [2\psi^2(t - u) + 4S\psi(t - u)]\, dt\, du/T^2$$

$$= 4 \int_0^T (T - v)[\psi^2(v) + 2S\psi(v)]\, dv/T^2 \tag{15e}$$

A similar treatment of $R_E$ gives

$$\operatorname{var} R_E = 2 \int_0^\infty e^{-v/T}[\psi^2(v) + 2S\psi(v)]\, dv/T. \tag{15f}$$

The third case is that of a sinusoidal signal with instantaneous sampling, for which

$$Z(t) = 2S \cos^2 (\Omega t + \varphi) + 2\sqrt{2S} n(t) \cos (\Omega t + \varphi) + n^2(t)$$

where $\varphi$ is uniformly distributed in $(0, 2\pi)$. The variance of $R_U$ is now given by

$$E \left( \int_0^T [n^2(t) - N]\, dt/T + 2 \int_0^T \sqrt{2S} \cos (\Omega t + \varphi) n(t)\, dt/T \right.$$

$$\left. + \int_0^T S[2 \cos^2 (\Omega t + \varphi) - 1]\, dt/T \right)^2$$

$$= E \int_0^T \int_0^T [n^2(t) - N][n^2(u) - N]\, dt\, du/T^2$$

$$= E \cdot 8S \int_0^T \int_0^T \cos (\Omega t + \varphi) \cos (\Omega u + \varphi) n(t) n(u)\, dt\, du/T^2$$

$$+ E \cdot S^2 \int_0^T \int_0^T \cos (2\Omega t + 2\varphi) \cos (2\Omega u + 2\varphi)\, dt\, du/T^2$$

To find $E \cos (\Omega t + \varphi) \cos (\Omega u + \varphi)$, expand the cosines, note that $E \cos^2 \varphi = E \sin^2 \varphi = \frac{1}{2}$ and $E \cos \varphi \sin \varphi = 0$, and combine the resulting terms. This gives

$$\operatorname{var} R_U = \int_0^T \int_0^T (2\psi^2(t - u) + 4S\psi(t - u) \cos \Omega(t - u)$$

$$+ \tfrac{1}{2}S^2 \cos 2\Omega(t - u))\, dt\, du/T^2 \tag{15g}$$

$$= 4 \int_0^T (T - v)[\psi^2(v) + 2S\psi(v) \cos \Omega v]\, dv/T^2 + (S^2/2T^2\Omega^2) \sin^2 \Omega T.$$

A similar treatment of $R_E$ gives

var $R_E =$

$$2 \int_0^\infty e^{-v/T}[\psi^2(v) + 2S\psi(v) \cos \Omega v] \, dv/T + S^2/2(1 + 4T^2\Omega^2) \tag{15h}$$

of which the last term is the evaluation of

$$\frac{S^2}{2} \int_0^\infty \int_0^\infty e^{-(t+u)/T} \cos 2\Omega(t - u) \, dt \, du/T^2$$

It remains to consider the narrow-band case, with filtering after rectification to give half the square of the envelope. As in Sections 3 and 10, the final output has the form

$$Z(t) = S + \sqrt{2S}n_1(t) + [n_1^2(t) + n_2^2(t)]/2$$

The variance of $R_U$ is then

$$E \left( \frac{1}{2} \int_0^T [n_1^2(t) + n_2^2(t) - 2N] \, dt/T + \sqrt{2S} \int_0^T n_1(t) \, dt/T \right)^2$$

The usual method of evaluation gives

$$\text{var } R_U = 2 \int_0^T (T - v)[\psi_0^2(v) + \psi_1^2(v) + 2S\psi_0(v)] \, dv/T^2 \tag{15i}$$

$$\text{var } R_E = \int_0^\infty e^{-v/T}[\psi_0^2(v) + \psi_1^2(v) + 2S\psi_0(v)] \, dv/T \tag{15j}$$

where

$$\psi_0(v) = En_1(t)n_1(t + v) = En_2(t)n_2(t + v)$$

and

$$\psi_1(v) = En_1(t)n_2(t + v) = - En_1(t + v)n_2(t).$$

Then $T \to \infty$, the ratio of the variances of $R_E$ and $R_U$ approaches one-half in all cases, so it will suffice to consider the latter. The spectrum of the noise will be assumed to be flat, of width $\omega$ radians per second. Then in (15c) and (15e) one has $\psi(v) = N \sin \omega v / \omega v$, and so the variance of $R_U$ is asymptotically equal to

$$\frac{2}{T} \int_0^\infty \frac{N \sin \omega v}{\omega v} \, dv$$

and

$$\frac{4}{T} \int_0^\infty \left[ \frac{N^2 \sin^2 \omega v}{(\omega v)^2} + \frac{2NS \sin \omega v}{\omega v} \right] dv$$

respectively.

Putting $\psi(v) = (2N/\omega v) \sin (\omega v/2) \cos \Omega v$ in (15g) gives

$$\frac{4}{T} \int_0^\infty \left[ \frac{N^2 \sin^2 \omega v/2}{(\omega v/2)^2} + \frac{2NS \sin \omega v/2}{\omega v/2} \right] \cos^2 \Omega v \, dv$$

Putting

$$\psi_1(v) = 0, \qquad \psi_0(v) = (2N/\omega v) \sin (\omega v/2)$$

in (15i) gives

$$\frac{2}{T} \int_0^\infty \left[ \frac{N^2 \sin^2 \omega v/2}{(\omega v/2)^2} + \frac{2NS \sin \omega v/2}{\omega v/2} \right] dv$$

The first integral has the value $\pi N/\omega T$ and the last three all have the value $(2\pi/\omega T)(N^2 + 2NS)$. Comparing these with the variances of single samples given for steady signals in Table II (Section 12), and writing $\omega/2\pi = W$ cps we see that $R_U$ is asymptotically equivalent to $2WT$ independent samples in the two dc cases, and $WT$ in the ac envelope case. The number is something above $2WT$ for instantaneous samples of ac. These results agree with those arrived at in Section 3 and 13. With the exception of the unrectified dc signal, the differences are seen to lie in the efficacy of an isolated sample, rather than in the long-term rate of transport of information.

When the signal is absent, (15e) to (15h) reduce to the results of Rice cited above:

$$\text{var } R_U = 4 \int_0^T (T - v)\psi^2(v) \, dv/T^2 \tag{15k}$$

and

$$\text{var } R_E = 2 \int_0^\infty e^{-v/T}\psi^2(v) \, dv/T \tag{15\ell}$$

The same is nearly true for (15i) and (15j) also, since by (3c), $\psi_0^2(v) + \psi_1^2(v)$ is the square of the envelope of $\psi^2(v)$ in the narrow-band (ac) case.

16. OPTIMUM PROCEDURES AND THE BASIC ASSUMPTIONS

The rigorous determination of the optimum detection procedure is a deep problem with more or less complicated answers depending on the spectra and probably involving both discrete sampling and continuous integration, or even differentiation.[14-16] However, the solution of the prob-

[14] U. Grenander, Stochastic Processes and Statistical Inference, Arkiv för Matematik, **1**, p. 195, 1950, Sections 4.11 and 5.4.
[15] E. Reich and P. Swerling, The Detection of a Sine Wave in Gaussian Noise. J. Appl. Phys., **24**, p. 289, 1953.
[16] D. Slepian, Estimation of Signal Parameters in the Presence of Noise. Trans. I.R.E., Professional Group on Information Theory, p. 68, March, 1954.

lem is known as soon as one restricts oneself to independent samples of the signal-plus-noise. Then square-law detection and averaging is the rigorously optimal procedure for detecting a noise-like signal, and it is virtually optimal for a steady sinusoid. The rigorous optimum in the latter case is well-known and depends on the assumed signal power $S$; if $v$ is the amplitude of the envelope (output of a linear rectifier), then a non-linear rectifier or other device is to be used to convert $v$ to $\log I_0(v\sqrt{2S}/N)$, and the values of the latter are averaged. $I_0(x)$ is the Bessel function $J_0(x\sqrt{-1})$. This corresponds very nearly to square-law rectification when $S/N$ is small, and linear rectification when $S/N$ is large. Since square-law is little different from linear rectification, as noted in Section 1, it is also little different from the optimal. For a steady dc signal, the optimal procedure uses the average of the algebraic values of the samples (without any rectification); the performance is given by (12a).

# Motion of Individual Domain Walls In a Nickel-Iron Ferrite

## Erratum

### By J. K. GALT

In footnote (20) to a paper[1] of the above title it was asserted that the theory of magnetic after-effect* given by L. Néel[2] leads to zero loss for large motions of a 180° domain wall. It has since become clear that this assertion is based on a misinterpretation of Section 10 of Néel's paper, and is therefore not correct. In fact, Professor Néel has pointed out in a private communication that if the general analysis in his paper is used to calculate the viscous drag on a 180° domain wall, a result substantially the same as that given by equation (38) in Reference 1 is obtained. Néel's theory therefore does not lead to a result inconsistent with the domain wall data given in Reference 1; it appears to be possible to account for this set of data with either his theory or the analysis presented in Reference 1.

The following is a derivation, which is due to Néel, of the result which is to be compared with equation (38) in Reference 1. We start from equation (26) of Reference 2, and use Néel's notation:

$$P = -W_0 \int_0^t f(U) g(t - \tau) \, d\tau \tag{1}$$

Here $P$ is the pressure due to magnetic drag, $W_0$ is a constant which determines the magnitude of the energy to be gained by rearranging carbon atoms (valence electrons in the case of the ferrite) and $g(t - \tau)$ is a weighting factor which takes the form

$$\frac{1}{\theta} \exp \frac{(\tau - t)}{\theta}$$

if only one relaxation time, $\theta$, is involved. $U$ is the distance between the

---

* Traînage.

[1] J. K. Galt, B.S.T.J., **33**, p. 1023, Sept., 1954.
[2] L. Néel, J. Phys. et Radium, **13**, p. 249, 1952.

positions of the wall at time $\tau$ and at time $t$; $f(U)$ is defined thus:

$$f(U) = \frac{\partial F(U)}{\partial U} \tag{2}$$

where

$$F(U) = -\frac{1}{W_0} \int_{-\infty}^{\infty} E_d(\tau) \, dx \tag{3}$$

The function $F(U)$ is an integral, over a cylinder of unit cross-section normal to the wall, of the angular dependence of the energy $E_d$ to be gained by rearranging carbon atoms (valence electrons in the case of the ferrite). Further details will be found in Reference 2.

In the case of a domain wall moving with constant velocity $v$,

$$U - v\tau - vt,$$

and we assume that only one relaxation time is important in the loss mechanism. In this case Equation (1) becomes

$$P = -\frac{W_0}{v\theta} \int_{-vt}^{0} f(U) e^{U/v\theta} \, dU \tag{4}$$

If we note that $f(U)$ is an odd function (Section 8 in Reference 2) and rearrange the limits of integration, this becomes:

$$P = \frac{W_0}{v\theta} \int_{0}^{\infty} f(U) e^{-U/v\theta} \, dU. \tag{5}$$

Now because of the factor $e^{-U/v\theta}$, we only get contributions to the integral in Equation (5) from the region where $U$ is comparable to or less than $v\theta$. If $d$ is the thickness of the domain wall, and if the velocity of the domain wall is slow enough so that $d \gg v\theta$, $U \ll d$ in the region of importance. We may therefore use the first term of the series for $f(U)$ discussed in Section 8 of Reference 2. From this we find

$$P = -\frac{4W_0}{3v\theta d} \int_{0}^{\infty} U e^{-U/v\theta} \, dU$$

$$= -\frac{4W_0 v\theta}{3d} \tag{6}$$

If we set this pressure, due to viscous drag, equal to the pressure from the applied steady field on the domain wall, $2 M_s H_0$, we find

$$v = \frac{1}{\theta} \frac{3M_s d}{2W_0} H_0 \tag{7}$$

This relation is to be compared with Equation 38 in Reference 1. They are of the same form, and in particular they both lead to the same dependence of $v$ on applied field and temperature (note that the relaxation time $\theta$ depends exponentially on the temperature). They both can therefore be used to fit the experimental data in Reference 1, and comparisons with other data will be necessary to distinguish between the two approaches.

# Bell System Technical Papers Not Published in this Journal

ANDERSON, O. L.,[1] and STUART, D. A.[9]

**The Calculation of the Activation Energy of Ionic Conductivity in Silica Glasses by Classical Methods,** J. Am. Ceramic Society, **37,** pp. 573–580, Dec., 1954.

BENÉS, V. E.[1]

**Partial Model for Quine's "New Foundations",** J. Symbolic Logic, **19,** pp. 187–200, Sept., 1954.

BOGAN, L. B.,[2] and YOUNG, K. D.[2]

**Simplified Transmission Engineering in Exchange Cable Plant Design,** A.I.E.E. Commun. and Electronics, No. 15, pp. 498–502, Nov., 1954.

BOWER, FRANK H.[3]

**Manufacturing Grown Junction Transistors,** Electronics, **27,** pp. 130–134, Dec., 1954.

BRERETON, D. S.,[8] and DONNELLY, H. J.[7]

**480-Wye/277-Volt Power System in Telephone Building at Menands, N. Y.,** Elec. Eng., **73,** pp. 1100–1105, Dec., 1954.

BURNS, R. M.[1]

**Crisis in Science Teaching,** J. Electrochem. Soc., **101,** pp. 261C–262C, Nov., 1954.

CALBICK, C. J.[1]

**Ice,** Letter to the Editor, Physics Today, **7,** p. 27, Dec., 1954.

---

[1] Bell Telephone Laboratories, Inc.
[2] American Telephone and Telegraph Company.
[3] Western Electric Company.
[7] New York Telephone Company.
[8] General Electric Company.
[9] Cornell University.

CAMPBELL, MARY E., see Luke, C. L.

CLARKE, K. B.,[3] and COURAGE, J. W.[3]

**Making Small Parts,** Electronics, **27,** pp. M15–M22, Oct., 1954.

COURAGE, J. W., see Clarke, K. B.

CRUSER, V. I.[1]

**Equipment and Mechanical Features of the AN/TRC-24 Radio Set,** A.I.E.E. Commun. and Electronics, No. 15, pp. 544–547, Nov., 1954.

DANIELSON, W. E., see Pierce, J. R.

DANSER, J. W., see Hazen, D. F.

DITZENBERGER, J. A., see Fuller, C. S.

DONNELLY, H. J., see Brereton, D. S.

EBERS, J. J.,[1] and MOLL, J. L.[1]

**Large-Signal Behavior of Junction Transistors,** I.R.E., Proc., **42,** pp. 1761–1772, Dec., 1954.

FULLER, C. S.,[1] and DITZENBERGER, J. A.[1]

**Diffusion of Boron and Phosphorus into Silicon,** Letter to the Editor, J. Appl. Phys., **25,** pp. 1439–1440, Nov., 1954.

HAZEN, D. F.,[4] DANSER, J. W.,[5] and ZILIS, G. S.[5]

**A Private Microwave Radio System for Power Company Use,** A.I.E.E. Commun. and Electronics, No. 15, pp. 492–498, Nov., 1954.

HELM, H. A.[1]

**Frequency Response Approach to the Design of a Mechanical Servo,** Trans. A.S.M.E., **76,** pp. 1195–1214, Nov., 1954.

HOFFMANN, J. P.[1]

**New Military Carrier Telephone Systems Equipment Features,** A.I.E.E. Commun. and Electronics, No. 15, pp. 509–515, Nov., 1954.

HUBER, G. H.,[1] MILLER, W. F.,[1] and SCHRAMM, C. W.[1]

**New Military Carrier Telephone Systems,** A.I.E.E. Commun. and Electronics, No. 15, pp. 515–525, Nov., 1954.

JENSEN, A. G.[1]

**Evolution of Modern Television,** J.S.M.P.T.E., **63,** pp. 174–187, Nov., 1954.

KELLY, H. P.[1]

**Differential Phase and Gain Measurements in Color Television Systems,** A.I.E.E. Commun. and Electronics, No. 15, pp. 565–569, Nov., 1954.

KOHN, W.,[1] see Luttinger, J. M.

LUKE, C. L.,[1] and CAMPBELL, MARY E.[1]

**Photometric Determination of Magnesium in Electronic Nickel,** Analyt. Chem., **26,** pp. 1778–1780, Nov., 1954.

LUTTINGER, J. M.,[1] and KOHN, W.[1]

**Hyperfine Splitting of Donor Status in Silicon,** Letter to the Editor, Phys. Rev., **96,** p. 802, Nov. 1, 1954.

MASON, D. R.[1]

**Design Method for the Calculation of Stagewise Reaction Systems,** (in French), Chimie and Industrie, **72,** pp. 241–251, Aug., 1954.

MASON, W. P.,[1] and WICK, R. F.[1]

**Ferroelectrics and the Dielectric Amplifier,** Proc. I.R.E., **42,** pp. 1606–1620, Nov., 1954.

McLEAN, D. A.,[1] and WEHE, H. G.[1]

**Miniature Lacquer Film Capacitors,** I.R.E., Proc., **42,** pp. 1799–1805 Dec., 1954.

McMILLAN, B.[1]

**Absolutely Monotone Functions,** Annals of Math., **60,** pp. 467–501, Nov., 1954.

---

[1] Bell Telephone Laboratories, Inc.

MERRILL, J. L., JR.[1]

**Theory of E-Type Repeaters,** A.I.E.E. Commun. and Electronics, No. 15, pp. 443–447, Nov., 1954.

MILLER, W. F., see Huber, G. H.

MINER, E. S.[2]

**Is Your Safety Program Adequate?,** Telephony, **147,** pp. 17–19, Nov., 1954.

MOLL, J. L.[1]

**Large-Signal Transient Response of Junction Transistor,** I.R.E. Proc., **42,** pp. 1773–1784, Dec., 1954.

MOLL, J. L., see Ebers, J. J.

NORDAHL, J. G.[1]

**A New Ultra High Frequency Multichannel Military Radio Relay System,** A.I.E.E. Commun. and Electronics, No. 15, pp. 526–531, Nov., 1954.

PFANN, W. G.[1] and VAN ROOSBROECK, W.[1]

**Radioactive and Photoelectric p-n Junction Power Sources,** J. Appl. Phys., **25,** pp. 1422–1434, Nov., 1954.

PIERCE, J. R.,[1] and DANIELSON, W. E.[1]

**Minimum Noise Figure of Traveling-Wave Tubes with Uniform Helices,** J. Appl. Phys., **25,** pp. 1163–1165, Sept., 1954.

PIERCE, J. R.[1]

**Some Recent Advances in Microwave Tubes,** I.R.E. Proc., **42,** pp. 1735–1747, Dec., 1954.

QUINLAN, A. L.[3]

**Automatic Percussion Welding,** A.I.E.E. Commun. and Electronics, No. 15, pp. 561–565, Nov., 1954.

READ, W. T., JR.[1]

**Statistics of the Occupation of Dislocation Acceptor Centres,** Phil. Mag., **45,** pp. 1119–1128, Nov., 1954.

---

[1] Bell Telephone Laboratories, Inc.
[2] American Telephone and Telegraph Company.
[3] Western Electric Company.

Rose, A. F.[2]

**Negative Impedance Telephone Repeaters — Application in the Bell System,** A.I.E.E. Commun. and Electronics, No. 15, pp. 430–435, Nov., 1954.

Rowen, J. H., and Vonaulock, W.[1]

**Measurement of the Complex Tensor Permeability of Ferrites,** Letter to the Editor, Phys. Rev., **96,** p. 1151, Nov. 15, 1954.

Schramm, C. W., see Huber, G. H.

Sernelius, W. C., see Van Uitert, L. G.

Smethurst, J. O.[1]

**E-Type Telephone Repeaters — Description, Equipment, and Testing,** A.I.E.E. Commun. and Electronics, No. 15, pp. 435–443, Nov. 1954.

Stuart, D. A., see Anderson, O. L.

Sullivan, J. W.[1]

**A Wide Band Voltage Tunable Oscillator,** Proc. I.R.E., **42,** pp. 1658–1665, Nov., 1954.

Thomas, D. E.[1]

**Stability Considerations in VHF Point-Contact Transistor Parameter Measurements,** Proc. I.R.E., **42,** pp. 1636–1644, Nov., 1954.

Van Roosbroeck, W., see Pfann, W. G.

Van Uitert, L. G.,[1] and Sernelius, W. C.[6]

**Determination of Thermodynamic Equilibrium Constants in Mixed Solvents,** J. Am. Chem. Soc., **76,** pp. 5887–5888, Nov. 20, 1954.

Vonaulock, W.,[1] see Rowen, J. H.

Wehe, H. G.,[1] see McLean, D. A.

Wick, R. F., see Mason, W. P.

Young, K. D., see Bogan, L. B.

Zilis, G. S., see Hazen, D. F.

# Recent Monographs of Bell System Technical Papers Not Published in This Journal

CHRISTENSEN, H.

**Surface Conduction Channel Phenomena in Germanium,** Monograph 2304.

CORENZWIT, E., see Matthias, B. T.

DANIELSON, W. E., see Pierce, J. R.

FAGEN, M. D.

**Bibliography on Ultrasonic Delay Lines,** Monograph 2317.

FULLER, C. S., and SEVERIENS, J. C.

**Mobility of Impurity Ions in Germanium and Silicon,** Monograph 2309.

GEBALLE, T. H., see Matthias, B. T.

GELLER, S., see Matthias, B. T.

HAGSTRUM, H. D.

**Auger Ejection of Electrons from Tungsten,** Monograph 2313.

HAGSTRUM, H. D.

**Auger Ejection of Electrons from Metals by Ions,** Monograph 2351.

KARP, A.

**Traveling-Wave Tube Experiments at Millimeter Wavelengths,** Monograph 2369.

LAW, J. T.

Effect of Water on Grown Germanium N-P Junctions, Monograph 2315.

LOGAN, R. A., and SCHWARTZ, M.

**Thermal Effects on Lifetime of Minority Carriers in Germanium,** Monograph 2312.

MAHONEY, J. J., see Perkins, E. H.

MATTHIAS, B. T., GEBALLE, T. H., GELLER, S., and CORENZWIT, E.
**Superconductivity of Nb₃Sn,** Monograph 2355.

MERZ, W. J.

**Domain Formation and Domain Wall Motions in Ferroelectric BaTiO₃ Single Crystals,** Monograph 2286.

PERKINS, E. H., and MAHONEY, J. J.

**Type-N Carrier Telephone Deviation Regulator,** Monograph 2318.

PIERCE, J. R., and DANIELSON, W. E.

**Minimum Noise Figure of Traveling-Wave Tubes with Uniform Helices,** Monograph 2237.

PIERCE, J. R., and TIEN, P. K.

**Coupling of Modes in Helixes,** Monograph 2303.

SCHWARTZ, M., see Logan, R. A.

SEVERIENS, J. C., see Fuller, C. S.

TIEN, P. K., see Pierce, J. R.

TIEN, P. K.

**Focusing of a Long Cylindrical Electron Stream,** Monograph 2293.

# Recent Monographs of Bell System Technical Papers Not Published in This Journal

CHRISTENSEN, H.

**Surface Conduction Channel Phenomena in Germanium,** Monograph 2304.

CORENZWIT, E., see Matthias, B. T.

DANIELSON, W. E., see Pierce, J. R.

FAGEN, M. D.

**Bibliography on Ultrasonic Delay Lines,** Monograph 2317.

FULLER, C. S., and SEVERIENS, J. C.

**Mobility of Impurity Ions in Germanium and Silicon,** Monograph 2309.

GEBALLE, T. H., see Matthias, B. T.

GELLER, S., see Matthias, B. T.

HAGSTRUM, H. D.

**Auger Ejection of Electrons from Tungsten,** Monograph 2313.

HAGSTRUM, H. D.

**Auger Ejection of Electrons from Metals by Ions,** Monograph 2351.

KARP, A.

**Traveling-Wave Tube Experiments at Millimeter Wavelengths,** Monograph 2369.

LAW, J. T.

Effect of Water on Grown Germanium N-P Junctions, Monograph 2315.

LOGAN, R. A., and SCHWARTZ, M.

**Thermal Effects on Lifetime of Minority Carriers in Germanium,** Monograph 2312.

MAHONEY, J. J., see Perkins, E. H.

MATTHIAS, B. T., GEBALLE, T. H., GELLER, S., and CORENZWIT, E.

**Superconductivity of $Nb_3Sn$,** Monograph 2355.

MERZ, W. J.

**Domain Formation and Domain Wall Motions in Ferroelectric $BaTiO_3$ Single Crystals,** Monograph 2286.

PERKINS, E. H., and MAHONEY, J. J.

**Type-N Carrier Telephone Deviation Regulator,** Monograph 2318.

PIERCE, J. R., and DANIELSON, W. E.

**Minimum Noise Figure of Traveling-Wave Tubes with Uniform Helices,** Monograph 2237.

PIERCE, J. R., and TIEN, P. K.

**Coupling of Modes in Helixes,** Monograph 2303.

SCHWARTZ, M., see Logan, R. A.

SEVERIENS, J. C., see Fuller, C. S.

TIEN, P. K., see Pierce, J. R.

TIEN, P. K.

**Focusing of a Long Cylindrical Electron Stream,** Monograph 2293.

# Contributors to this Issue

S. Theodore Brewer, B.S. in E.E., M.S. in E.E., Purdue University, 1937 and 1938; Bell Telephone Laboratories, 1937–. His early work involved broad-band carrier systems and transformer coupled video amplifiers, and the design of measuring equipment associated with these systems. Later, he was concerned with electronically controlled automatic switching. He holds patents on control and feedback systems and switching networks. Mr. Brewer is presently working on the application of transistors to high-speed switching networks. During World War II, he served as radar staff officer with the 62nd Fighter Wing. Member of I.R.E., Eta Kappa Nu, Tau Beta Pi and Sigma Xi.

Norman E. Earle, B.S. and M.S., M.I.T., 1929 and 1930. Bell Telephone Laboratories 1930–31, Timken Silent Automatic 1932–1933, Postal Telegraph 1934–1935, Western Electric Company 1936–. From 1936 to 1947 he engaged in product engineering on various types of transmission coils and transformers. Since 1947 Mr. Earle has been Engineer on Engineering Development and Test Set and Machine Design at the Haverhill Coil Shops. He is a member of the I.R.E.

James Gammie, B.Sc. in Electrical Engineering, University of Aberdeen, 1944; B.Sc. in Mathematics, Birkbeck College, University of London, 1951. Mr. Gammie was employed by the Standard Telephones and Cables Ltd., North Woolwich, London, from 1944 to 1951. He joined Bell Telephone Laboratories in 1952. For the British concern he was engaged on the design of test and frequency translating equipment associated with its coaxial and other carrier systems. At Bell Laboratories he has worked on the L3 coaxial system. Recently, Mr. Gammie has been doing development work on a short haul radio system.

George Hecht, B.S., 1930; E.E., 1951, Cooper Union. Mr. Hecht joined the Western Electric Company in 1924 and transferred to the Bell Telephone Laboratories the following year. His early work entailed fundamental studies of quartz crystals and tuning forks as resonators, and their applications as precision oscillators and frequency standards. In 1938 he began four years of work on circuit research in the field of ac

key pulsing receivers for toll systems. From 1942 to 1945 he worked on military projects, particularly radar range indicators. Since 1945 he has been connected with research in switching, studying means of applying electronics to telephone systems. Senior Member, Institute of Radio Engineers.

CONYERS HERRING, A.B., University of Kansas, 1933; Ph.D., Princeton University, 1937; National Research Fellow, Massachusetts Institute of Technology, 1937–39; Research Associate, Princeton University, 1939–40; Instructor of Physics, University of Missouri, 1940–41; Division of War Research, Columbia, 1941–45; Professor of Applied Mathematics, University of Texas, 1946; Bell Telephone Laboratories, 1945–. Dr. Herring has been engaged principally in research in physical electronics and solid state physics. He has given a number of lectures on solid state physics at the Institute for Advanced Study, Princeton, N. J. Fellow of the American Physical Society and member of the executive committee of the Society's Division of Solid State Physics. Member, A.A.A.S.; Board of Editors of the Physical Review, 1952–54. Awarded the Army-Navy Certificate of Appreciation.

OLE M. HOVGAARD, B.S. in E.E., Massachusetts Institute of Technology, 1926; Tropical Radio Telegraph Company, 1921–23; Briggs and Stratton Corp., 1927–28; Bell Telephone Laboratories, 1928–. Mr. Hovgaard's early work at Bell Laboratories was concerned with the design of broadcasting transmitter antennas. Later, he was engaged in the development of quartz crystals, precious metal contacts engineering, and military developments. He is presently concerned with the development of sealed switches. Mr. Hovgaard is the author of several technical papers. Senior member of I.R.E.

EDWARD L. KAPLAN, B.S., Carnegie Institute of Technology, 1941; Ph.D., Princeton University, 1951. Dr. Kaplan worked at the U. S. Naval Ordnance Laboratory in Silver Spring, Md. from 1941 to 1948 and was a research assistant at Princeton University from 1948 to 1950. Since 1950 he has been concerned with military and Bell System applications of probability theory and statistical methods at Bell Telephone Laboratories. Member of American Mathematical Society, Institute of Mathematical Statistics, American Statistical Association, Tau Beta Pi and Sigma Xi.

J. L. MERRILL, JR., B.S. and M.S., Pennsylvania State University,

1928 and 1930; Elliott Research Fellow, 1928–1930; American Telephone and Telegraph Company, 1930–1934; Bell Telephone Laboratories, 1934–. Mr. Merrill spent his first years with the Laboratories on transmissions features of such projects as the time and weather announcement systems and operator training programs. During World War II, he engaged in planning system operation of air raid warnings as well as work on tactical wire and radio networks for the armed forces. Later, he was concerned with the design and application of negative impedance repeaters for the improvement of exchange transmission. He is presently engaged in long range studies of transmission systems. He holds several patents and is the author of a number of technical articles.

GEORGE E. PERREAULT, B.S. in M.E., Worcester Polytechnic Institute, 1930; Bell Telephone Laboratories, 1930–. His early work was concerned with machine design of sound picture machines, picture transmission apparatus, vibration pickups, and hearing-test machines for the Bell System exhibit in the New York World's Fair. Later, he was engaged in the mechanical design of plotting boards and servo-mechanisms for various electrical gun directors. More recently he has been engaged in the apparatus development for vibrating reed selectors and crossbar switches. Mr. Perreault is presently concerned with development work on dry reed and mercury contact relays. He has been awarded several patents.