# The Bell System Technical Journal

## Reactance Tube Modulation of Phase Shift Oscillators

### By F. R. DENNIS and E. P. FELCH

This paper describes a basic circuit for reactance tube modulation of phase shift oscillators. The design of suitable phase shift oscillators for frequencies from audio through the ultra-high frequencies is discussed. Experimental performance data derived from several types of frequency modulated phase shift oscillators are presented.

### INTRODUCTION

FREQUENCY modulation of oscillators is finding wide-spread use in such diverse fields as FM broadcasting, telemetering systems for guided missiles and measuring apparatus for observing transmission frequency characteristics on cathode ray tubes. Design objectives for such oscillators may be listed briefly as:

1. A wide range of frequency modulation or, alternatively, high modulation sensitivity.
2. A linear relationship between instantaneous values of modulation input voltage and frequency deviation.
3. Freedom from accompanying amplitude modulation.
4. Inherent center frequency stability.
5. Ease and stability of adjustment.
6. A minimum number of components, none of which should be critical.
7. Modulation by dc, audio, or video inputs.
8. Operation anywhere in the frequency spectrum from low audio frequencies through the ultra-high frequency region.

The circuits described in this paper were developed in the course of an investigation of various frequency modulation circuits for use in visual transmission measuring systems. The oscillators had to be capable of linear modulation at 60 cycles over a $\pm 3$ megacycle band at 25 megacycles and over a $\pm 9$ megacycle band at 80 megacycles. Existing designs fell short of meeting the requirements with respect to several of the characteristics listed above. The reactance tube modulated phase shift oscillator circuit was found to perform satisfactorily in the transmission set and proved superior in many respects to all the other circuits tried. Tests of the circuit at other frequencies disclosed that the advantages were not peculiar to the frequency range and the following description is presented with the expectation that it may prove useful to others.
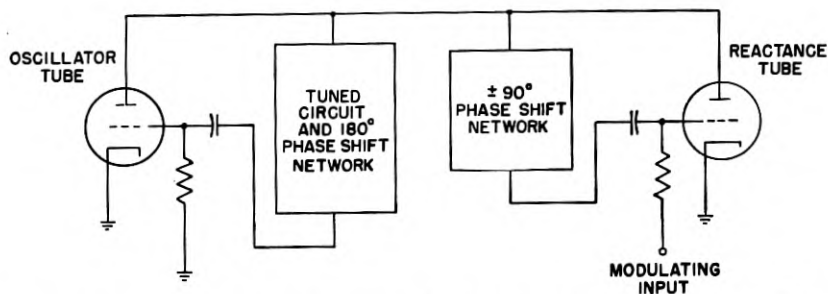
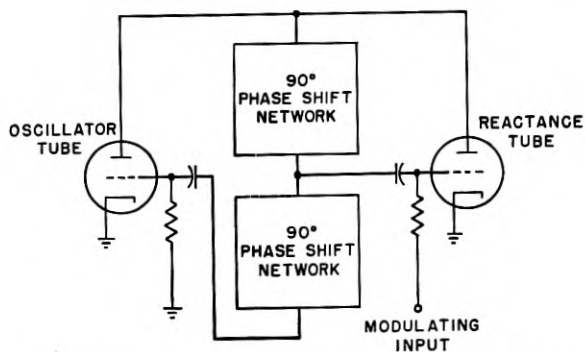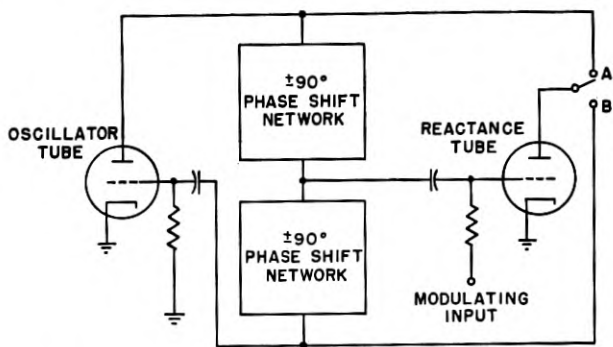Fig. 1—Simplified schematic of conventional reactance tube modulated oscillator.



Fig. 2—Simplified schematic of phase shift reactance tube modulated oscillator.



| TUBE CONNECTION | A | | B | |
|---|---|---|---|---|
| NETWORKS | LEADING (+90°) | LAGGING (-90°) | LEADING (+90°) | LAGGING (-90°) |
| OSCILLATION FREQ. | DECREASES | INCREASES | DECREASES | INCREASES |

Fig. 3—Direction of frequency deviation for increasing $G_M$ of reactance tube.
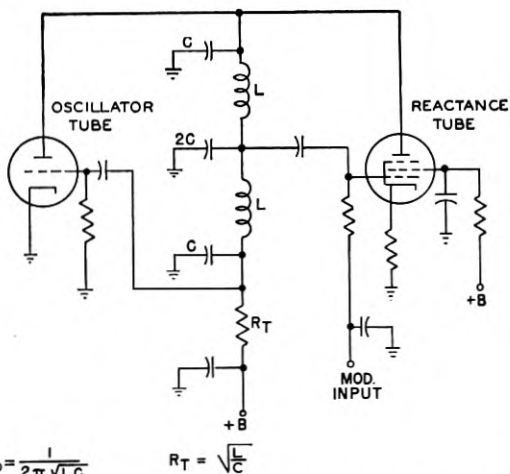
$$f_o = \frac{1}{2\pi\sqrt{LC}} \qquad R_T = \sqrt{\frac{L}{C}}$$

Fig. 4—LC reactance tube modulated phase shift oscillator.
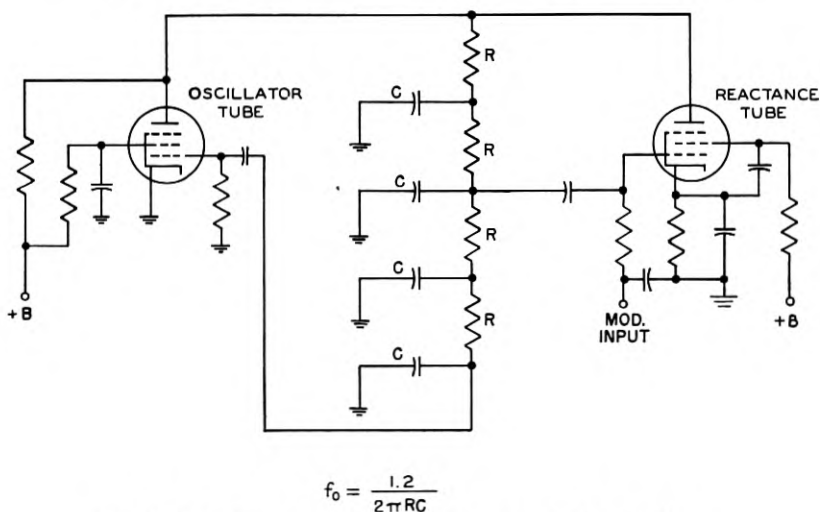


$$f_o = \frac{1.2}{2\pi RC}$$

Fig. 5—RC reactance tube modulated phase shift oscillator.

## CIRCUIT DESCRIPTION

The theory and design of conventional reactance tube modulated oscillators has been discussed adequately in the literature[1,2,3]. A schematic in

[1] "Frequency Modulation" (book) by August Hund—McGraw-Hill, New York, 1942. Page 155.
[2] "Automatic Tuning, Simplified Circuits and Design Practice," D. E. Foster, and S. W. Seeley. *Proc. I. R. E.*, Vol. 25, 1937, page 289.
[3] ATC Systems—*Wireless World*, February 19, 1937, page 177.

Fig. 6—Transmission line reactance tube modulated oscillator.

$$F_0 = \frac{150}{\sqrt{K}\, L} \text{ (MEGACYCLES)}$$

WHERE
K = DIELECTRIC CONSTANT.
L = LENGTH OF LINE IN METERS.



Fig. 7—Performance curves of typical LC reactance tube modulated phase shift oscillator.

simplified form is shown in Fig. 1. The input and output of a vacuum tube amplifier are connected together by a tuned circuit and feedback network which introduces 180° phase shift at the undeviated frequency $F_0$.
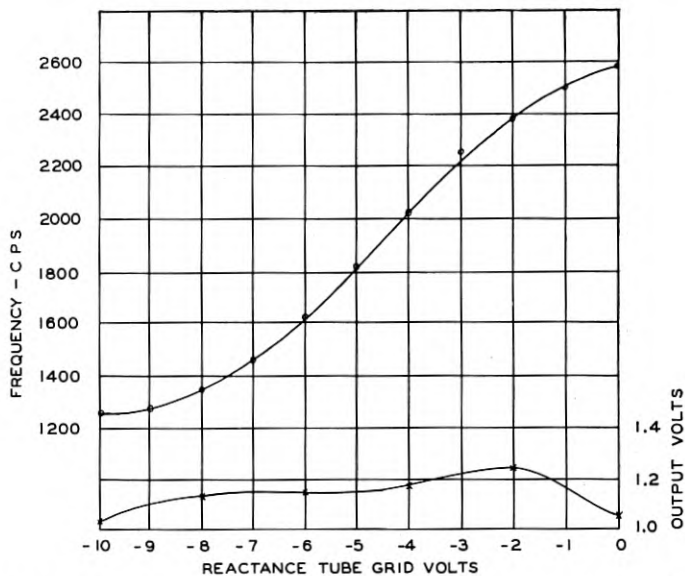
Fig. 8—Performance curves of typical RC reactance tube modulated phase shift oscillator.



Fig. 9—Performance curves of typical transmission line reactance tube modulated oscillator.

An auxiliary path contains the reactance tube fed from a 90° phase shift network connected as shown. The direction of frequency deviation is determined by the sign of the 90° phase shift. The amount of the deviation is



Fig. 10—Construction of transmission line reactance tube modulated oscillator. (a) Tube side. (b) Line side.

determined by the transconductance variation of the reactance tube, by the impedance across which the reactance tube is connected and by the loss in the 90° phase shift network. The linearity is a function of all of these factors. In general the frequency deviation may be increased by increasing

the $L/C$ ratio in the oscillator tuned circuit, but only at the expense of frequency stability.

A simplified schematic of the reactance tube modulated phase shift oscillator is shown in Fig. 2. The mathematical theory of operation is an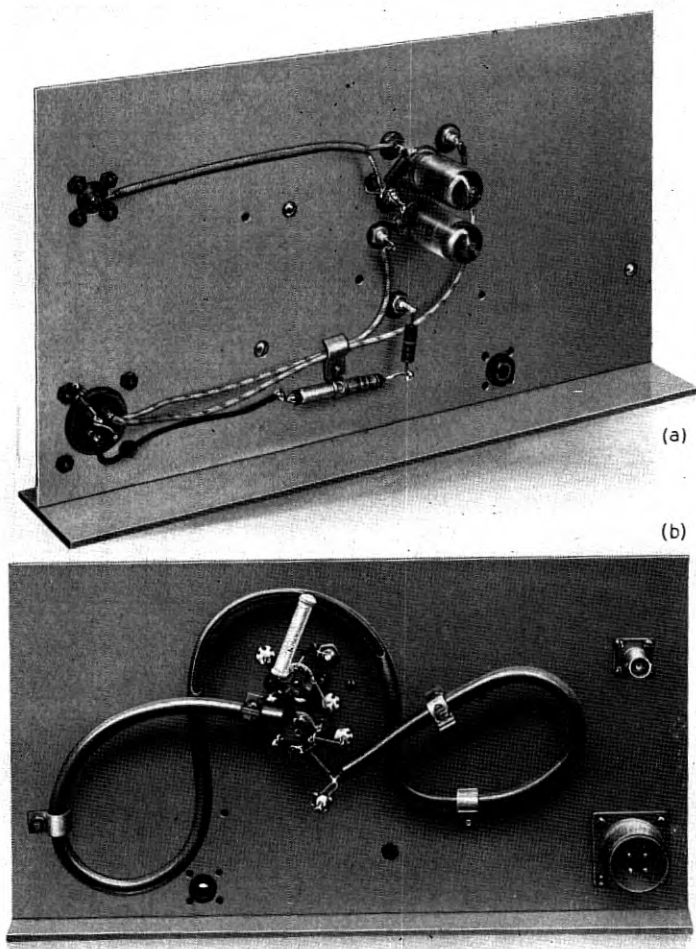alogous to that of the conventional reactance tube modulated oscillator, and the same methods of analysis may be applied. The 90° phase shift network required in the reactance tube grid circuit is a portion of the feedback network and provides half of the 180° phase shift required for oscillation. In this circuit the reactance tube is tightly coupled into the oscillating circuit with minimum loss in the 90° phase shift network. Hence small values of $L/C$ ratio may be employed with a consequent increase in the inherent frequency stability. In practice, oscillators comparable in stability to good nonmodulated oscillators may be realized. The direction of deviation is determined by whether the phase of the reactance tube grid voltage leads or lags the reactance tube plate current. The permutations of connections and signs of the 90° phase shift networks are shown on Fig. 3 with the corresponding directions of frequency deviation.

The phase shift networks need not be of the LC lumped constant variety. For example, RC networks or sections of transmission line may be employed to particular advantage at the lower and higher frequencies respectively. A few of the many possible circuit configurations are shown in Figs. 4, 5, 6.

## EXPERIMENTAL DATA

Frequency deviation and output variation curves for some typical oscillators are shown in Figs. 7, 8, and 9.

The oscillator of Fig. 9 which was built by Mr. D. Leed, is shown in Fig. 10. The transmission line is a section of RG59U cable with the shield removed, encased in a copper tube with a slot for bringing out the center tap of the line to the reactance tube grid. The tubes are 6J6's with both sections connected in parallel.

## CONCLUSION

Frequency modulated phase shift oscillators of several types have been described. These offer interesting possibilities for applications over a wide range of frequencies wherever stable, simple frequency modulated oscillators are required. With respect to range, linearity, and freedom from amplitude modulation their performance, as shown, is superior to that of conventional circuits and is at least equal to that of the complex circuits employed in the most critical applications.

# A Broad-Band Microwave Noise Source

## By W. W. MUMFORD

Measurements of the microwave noise power available from gaseous discharges, such as in an ordinary fluorescent lamp, show remarkable uniformity and stability. Such tubes are therefore suitable for a new type of standard noise source.

### INTRODUCTION

A STANDARD noise source, such as a hot resistance or a temperature limited diode, has been used advantageously for making measurements of the noise figure of radio receivers in the short-wave and the ultra-short wave region. The use of such a tool eliminates the possible errors which are practically inescapable when using the large amounts of attenuation which are needed for the determination of the ratio of power levels encountered in measuring noise figures with a standard signal generator. For example, the power from a standard signal generator might be measurable and known accurately at a level of 40 db below a watt, whereas the noise power available from a resistance might be 141 db below one watt.[1] It is difficult to ascertain accurately power ratios of this magnitude, $10^{10}$.

Another advantage of using a standard noise source arises from the fact that ordinarily the bandwidth of the receiver need not be considered, thereby eliminating a time consuming measurement. This assumes, of course, that the bandwidth of the noise source is much greater than that of the amplifier under test.

In the microwave region it is possible to match a resistive element to the waveguide over a wide enough band, but ordinary resistive materials will not stand the high temperatures (5000 degrees or more) needed to measure the noise figures encountered in practice. The noise diode is capable of furnishing adequate noise power, but one with wide bandwidth has yet to be developed. A good, stable, broadband microwave noise generator is needed.

Another possible source of noise power consists of a gaseous discharge.[2] Before we examine the data which have led us to conclude that the gaseous discharge is a good, broad-band, stable microwave noise generator and possibly a calculable noise standard, we review our definitions of noise figure

---

[1] This figure, 141 db below one watt, assumes that the effective bandwidth is 2 mc. The resistance noise power available from a generator at 290° Kelvin is 204 db below one watt per cycle.

[2] G. C. Southworth, *Journal of the Franklin Institute*, Vol. 239, ⅋14, pp. 285-298, April 1945.

and gain,[3] and discuss the factors involved in making noise figure measurements by means of a noise source.

### NOTES ON NOISE FIGURE

*Definition:* The NOISE FIGURE of a network, with a generator connected to its input terminals, is the ratio of the available signal-to-noise power ratio at the signal generator terminals (weighted by the network bandwidth) to the available signal-to-noise power ratio at its output terminals.

*Definition:* The GAIN of a network is the ratio of the available signal power at the output terminals of the network to the available signal power at the output terminals of the signal generator.



Fig. 1—Schematic diagram of generator, network and output power meter.

These definitions apply to a circuit consisting of a generator, a network and an output power meter as shown schematically in Fig. 1. The signal power available from the generator, having an open circuit voltage $e$ and an internal resistance $R_1$, is:

$$P_{SA} = \frac{e^2}{4R_1} \tag{1}$$

The noise power available from the signal generator resistance, $R_1$, at absolute temperature $T_1$, is

$$P_{NA} = \frac{4KT_1R_1B}{4R_1} = kT_1B \tag{2}$$

where $B$ is the effective bandwidth of the network, by which the generator noise is weighted in this case.

[3] H. T. Friis, *Proc. I. R. E.*, Vol. 32, # 17, pp. 419–422, July, 1944.

The weighted available signal-to-noise ratio at the generator terminals is:

$$\frac{P_{SA}}{P_{NA}} = \frac{\dfrac{e^2}{4R_1}}{KT_1B} \tag{3}$$

The network amplifies (or attenuates) the generator's signal power by the factor $G$, the gain of the network, so that the available signal power at the output terminals of the network is:

$$P_{SO} = G \frac{e^2}{4R_1} \tag{4}$$

The network amplifies (or attenuates) the generator noise power by the same factor $G$, and also delivers noise power which originates within itself, $N_N$, so that the total available noise power at the output terminals of the network is:

$$P_{NO} = GkT_1B + N_N \tag{5}$$

The available signal-to-noise ratio at the output terminals of the network is then:

$$\frac{P_{SO}}{P_{NO}} = \frac{G\dfrac{e^2}{4R_1}}{GkT_1B + N_N} \tag{6}$$

We now express the noise figure of the network, $F$, which by definition is the ratio of equation (3) to equation (6), thus,

$$F = \frac{GkT_1B + N_N}{GkT_1B} \tag{7}$$

We should pause at this point to consider this equation further, for it leads us to a simpler definition of noise figure.

*Definition:* The noise figure of a network is the ratio of the noise power output of that network to the noise power output which would exist if the network were noiseless. The temperature of the signal generator resistance is 290 degrees Kelvin.

The choice of generator temperature of 290 degrees is an arbitrary one, which makes $kT_1 = 4(10)^{-21}$ watts per cycle bandwidth; $-10 \log kT_1 = 204$ db below one watt per cycle. Putting $T_1 = 290$ in equation (7) gives:

$$F = \frac{Gk \, 290 \, B + N_N}{Gk \, 290 \, B} \tag{8}$$

Rearranging (8) we have:

$$N_N = (F - 1)Gk \, 290 \, B \tag{9}$$

Equation (9) will now be used to illustrate one method of measuring noise figures. In this method, the network output noise power is measured for two known values of the temperature of the generator resistance, $T_2$ and $T_1$. When the generator is hot, the output noise power is, by equation (5):

$$P_{NOH} = GkT_2B + N_N \qquad (10)$$

When the generator is cool, the output noise power is:

$$P_{NOC} = GkT_1B + N_N \qquad (11)$$

Calling the ratio of these two noise powers $Y$:

$$Y = \frac{P_{NOH}}{P_{NOC}} = \frac{GkT_2B + N_N}{GkT_1B + N_N} \qquad (12)$$

Substituting for $N_N$ the value given in equation (9), we have for the noise figure:

$$F = \frac{\left(\dfrac{T_2}{290} - 1\right) - Y\left(\dfrac{T_1}{290} - 1\right)}{Y - 1} \qquad (13)$$

In practice $T_1$ is often near enough to 290 degrees so that the second term in the numerator of equation (13) is negligible. Setting $T_1$ equal to 290 degrees, equation (13) becomes:

$$F = \frac{\dfrac{T_2}{290} - 1}{Y - 1} \qquad (14)$$

The determination of noise figure by this method is independent of the gain of the network, the degree of mismatch and the bandwidth, provided that the band of the noise source is broad compared with the overall RF band of the network and the output power meter.

## THE NOISE SOURCE

The limitations at microwaves of a noise source such as a heated wire will now be discussed. In particular we are interested in measuring amplifiers which have noise figures between 10 and 100 (10 db to 20 db) and bandwidths up to 200 mc. If a hot wire could be matched to the impedance of a waveguide over a wide enough band, and raised to a temperature of $10 \times 290$ degrees our $Y$ factor would be (rearranging eq. 14):

$$Y = \frac{\dfrac{T_2}{290} - 1}{F} + 1 \qquad (15)$$

and setting $T_2 = 2900$ degrees Kelvin

$$Y = 1.9 \text{ for } F = 10$$

$$Y = 1.09 \text{ for } F = 100$$

Assuming that $Y$ can be read to within $\pm 1\%$ our accuracy in determining $F$ would be within about $\pm 1\%$ for $F = 10$ but only within about $\pm 10\%$ for $F = 100$. If the noise source had a temperature of $40 \times 290$ degrees, our experimental errors would be reduced accordingly to about $\pm 1/4\%$ for $F = 10$ and $\pm 2.5\%$ for $F = 100$. Since metal wires will not stand such temperatures, we must look to something different for the noise source if these accuracies are to be achieved.

In view of the foregoing considerations, the nonoscillating reflex klystron presented one possibility of a suitable microwave noise source. This, however, was not exploited because the bandwidth was not wide enough.

Another possibility was found to be an electrical gas discharge. This type of source was determined to generate noise at microwave lengths when the open end of the input-waveguide of a sensitive microwave receiver was directed toward various gaseous discharge tubes, including a 721A TR tube containing water vapor and hydrogen, a neon light in a stroboscope, a mercury vapor rectifier and an ordinary fluorescent desk lamp. Of these, the commercial fluorescent lamp appeared to lend itself most readily to mounting in a waveguide without the complication of the effects of the internal metal electrodes, so further tests were performed on it.

## Microwave Measurements

A T-5, 6-watt, daylight fluorescent lamp,[4] lighted from a d-c. source, was mounted with its axis parallel to the magnetic vector in a waveguide as illustrated in Fig. 2. The lamp itself was 9″ long, with cathodes at each end. These could be isolated from the field in the 1″ x 2″ waveguide by enclosing the portion of the lamp which extended beyond the walls of the waveguide in cylindrical metal shields which formed waveguides beyond cutoff. Thus, energy was kept from reaching the cathodes, and the noise source was effectively confined to that part of the discharge which appeared inside the main waveguide. A piston in back of the gaseous discharge tube served to tune out the susceptance and a trimming screw provided an additional adjustment. The conductance could be adjusted by varying the direct current.

The admittance of the combination could be adjusted for an impedance

---

[4] A commercial fluorescent lamp contains about two mm. of argon and six to ten microns of mercury gas. The argon merely facilitates the initiation of the discharge; the mercury furnishes the radiation which excites the fluorescent material.

match at any operating frequency from 3700 mc to 4500 mc. The admittance diagram when the circuit was adjusted for match at 3960 mc is shown in Fig. 3; the standing wave ratio was less than 2.9 db from 3700 to 4240 mc.

At 3960 mc the conductance of the gaseous discharge varied directly with the direct current, while the negative susceptance had a broad maximum of $-j.62\ Y_0$ mhos at a current of 65 to 100 milliamperes, as shown in Fig. 4. These values are for the gaseous discharge; the susceptances of the enclosing glass tubing, the back piston and the holes in the sidewalls have been subtracted from the measured results. It is interesting to note that the discharge appears to be inductive.

The waveguide circuit containing the gaseous discharge tube was connected to the input waveguide of a sensitive microwave receiver which was used as a relative noise power meter. The noise power available from the
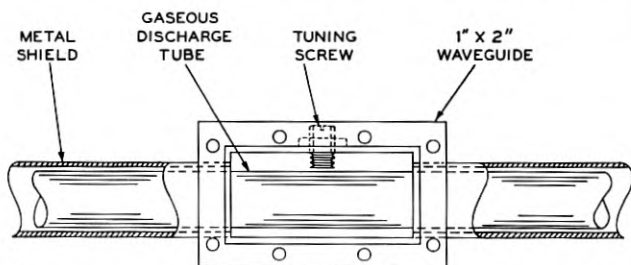


Fig. 2—Waveguide circuit for microwave noise generator using a gaseous discharge tube.

gaseous discharge was substantially independent of the direct current from 40 ma to 140 ma. These data are plotted in Fig. 5, which gives $10\log\left(\dfrac{T}{290}-1\right)$ versus direct current in milliamperes. The ordinate has been chosen so as to conform with absolute measurements made subsequently. The r.m.s. deviation from the straight line which represents a probable coefficient of only $-.003$ db per milliampere was about $\pm.05$ db. We do not claim to be able to achieve even this degree of accuracy with our present measuring equipment and hence do not place much confidence in the numerical value of this coefficient. Actually the decrease in noise with increasing current may have been associated with a change in the ambient temperature rather than with the increased current density. At least it is in the right direction for this to be the case.

The temperature coefficient of the noise from the discharge was found to be negative; when a piece of dry ice was held on the tubular shield of the circuit for a few minutes (long enough for frost to form on the brass) the output noise power of the discharge increased 0.6 db. The circuit was heated

on a hot plate and allowed to return to room temperature gradually, then cooled with an air stream and allowed to warm up gradually while the output noise and the temperature of the waveguide were being recorded. This revealed the temperature coefficient of $-.055$ db per degree centigrade. The data (plotted in Fig. 6) show an r.m.s. deviation of $\pm.114$ db from this coefficient.



Fig. 3—Admittance diagram of microwave noise generator.

The ambient temperature of the waveguide circuit had very little effect on the admittance of the gaseous discharge.

As a check on variability with respect to time, two of these noise sources were compared, one against the other, at five-minute intervals for 65 minutes. During this time the waveguide temperature of source #1 rose from 34° to 35.2° C and that of source #2 rose from 33.7° to 36.1°. Each comparison was corrected, according to the coefficient of $-.055$ db per degree C

and the observed temperature, to a common temperature of 34° C. Assuming that the noise figure of the microwave receiver was constant, source ⚹ 1



Fig. 4—Admittance of the gaseous discharge at 3960 mc as a function of the direct current in the discharge.



Fig. 5—The microwave noise power is practically independent of the discharge current.

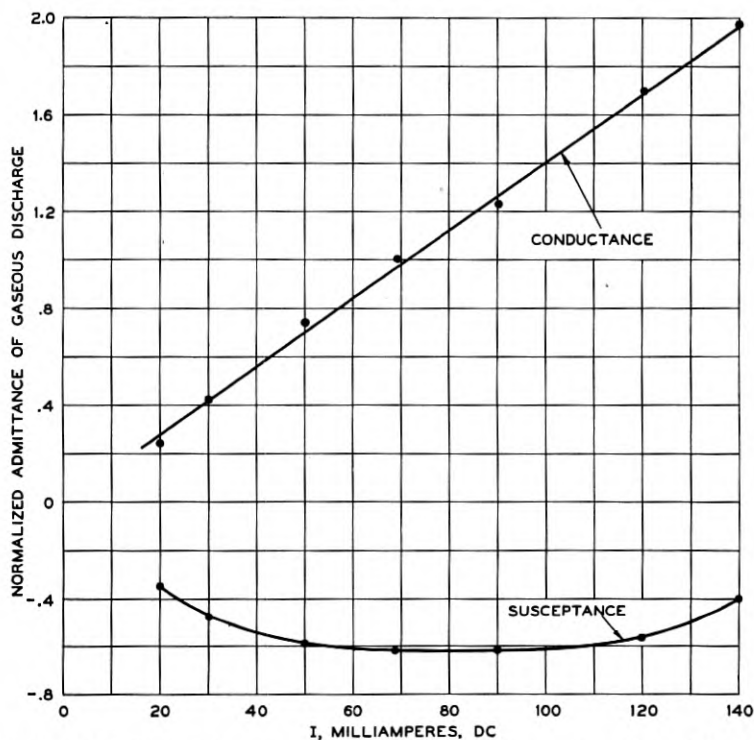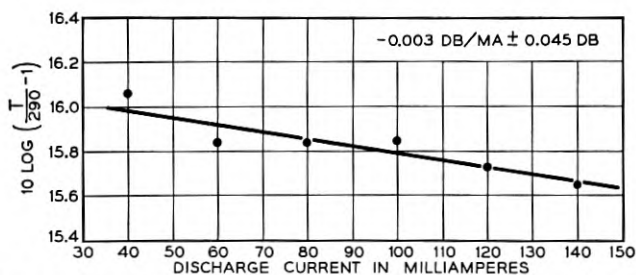showed variations whose r.m.s. deviation was ± 0.11 db, while source ⚹ 2 had similar deviations of ±.092 db. Assuming on the other hand that source ⚹ 1 held constant and that the microwave measuring set varied with time,

source #2 displayed r.m.s. deviations of ±.088 db. These variations are in fact comparable with the probable experimental error, and the proof that they actually exist at all still remains to be demonstrated.

Of thirty-two different lamps, including 10 different types of fluorescent coatings such as used in the pink, red, gold, soft white, daylight, green, white, 4500° white, black light and blue, thirty-one[5] were all within ±0.25 db of each other as was also a germicidal lamp with no fluorescent coating. Thus it appears that the source of the microwave noise energy lies chiefly in the gaseous discharge rather than in the fluorescent coating.



Fig. 6—The microwave noise power depends slightly upon the temperature of the waveguide circuit.

If this noise is tied up with the electron temperature of the discharge, we should expect the noise to be flat, or "white" noise. Corroborative evidence of this was observed when the spectrum of the noise was examined over the band from 3700 to 4500 mc at points 20 mc apart and no irregularities were found. The nature of the experiment was such that frequency bands of excessive noise power would have been observed had they been present. Further tests should indicate whether or not a gradual change in noise with frequency exists. It appears, however, unlikely that such a slope exists at 4000 mc.

Furthermore, since the level of the noise energy is so constant with respect to time, reproducible from tube to tube, practically independent of the current and only slightly affected by the ambient temperature, we might expect that it is being controlled or limited by some invariant physical property of the atoms and ions within the gaseous discharge. If this is the case, an absolute measurement of the noise power might lead us to some

---

[5] One of the 32 lamps flickered erratically. At times its noise was $\frac{1}{2}$ db higher than the average.

theoretical explanation which, when applied to the case in hand, would explain the observed results qualitatively and quantitatively, thereby establishing a new absolute standard noise source for microwave measurements.

The microwave noise power from such a discharge tube was measured at 3950 mc in cooperation with Mr. C. F. Edwards on his calibrated measuring set on two different occasions, 16 days apart. The values obtained were 15.86 db and 15.80 db respectively for $10 \log \left(\frac{T}{290} - 1\right)$.[6] This places the temperature, $T$, in the neighborhood of 11,400 degrees Kelvin. It is believed that the absolute measurements are correct to within $\pm.25$ db or better.

Having determined the temperature of this noise source, we might ask, "If we should terminate our waveguide in a black body at 11,400 degrees, how much microwave noise power would we get from it?" The black body radiates with three polarizations, only one of which is propagated along the waveguide, and this available power is given by Nyquist:[7]

$$P_{NA} = \frac{hfB}{e^{hf/kT} - 1} \tag{16}$$

where $h = 6.61 \ (10)^{-34}$ joule sec.

$\quad\quad k = 1.381 \ (10)^{-23}$ joule/deg.

$\quad\quad f = $ frequency in cycles per sec.

$\quad\quad B = $ bandwidth in cycles per sec.

At 4000 mc, $\frac{hf}{kT}$ is, for $T = 290$ degrees, $6.6 \ (10)^{-4}$ which is so small that the denominator of (16) can be replaced by $\frac{hf}{kT}$. This gives us the familiar expression for thermal noise:

$$P_{NA} = kTB \text{ watts} \tag{17}$$

In other words, thermal noise is black body radiation with but one polarization.

Going one step further we might also ask the question, "If we should examine the radiation from this black body with an optical spectroscope, at what wavelength would we find its maximum radiated energy?" The spectroscope detects radiation having three polarizations, and Planck's radiation law applies. From Wien's displacement law, the wavelength of maximum radiation is given by the relation:

$$\lambda_m T = 0.289 \text{ cm deg.} \tag{18}$$

---

[6] The temperature of the waveguide was 32°C when these values were measured.
[7] H. Nyquist, *Phys. Rev.*, Second Series, Vol. 32, pp. 110–113, July 1928.

Substituting $T = 11{,}400$ degrees,

$$\lambda_m = 2535 \; (10)^{-8} \; \text{cm} \tag{19}$$

This is indeed an interesting result, since the mercury vapor discharge in the fluorescent lamp radiates most of its energy at $\lambda = 2536.52 \; (10)^{-8}$ cm. The design of the lamp was guided by the effort to accentuate the radiation at this wavelength, and the manufacturers state that this has been achieved so that no other spectral line is excited to radiate more than two percent of the input power.[8] The conversion loss from dc to $2536 \; (10)^{-8}$ cm is only 2 or 3 db.

The striking similarity between the black body and the mercury vapor discharge at these two wavelengths, 7.6 cm and $2536 \; (10)^{-8}$ cm, suggests the following hypothesis:

*Hypothesis:* In a gaseous discharge which is radiating light energy substantially monochromatically at a particular wavelength, $\lambda_m$, the microwave noise energy is the same as that available from a black body which radiates its maximum energy at that wavelength.

Applying this hypothesis to the case in hand, where $\lambda_m$ is $2536.52 \; (10)^{-8}$ cm, and using Wien's displacement law (eq. 18) we calculate the temperature to be

$$T = \frac{0.289}{2536.52} = 11{,}394° \tag{21}$$

$$\frac{T}{290} = 39.29$$

$$\left(\frac{T}{290} - 1\right) = 38.29 \tag{22}$$

$$10 \log \left(\frac{T}{290-1}\right) = 15.84 \; \text{db} \tag{23}$$

Since this calculated value is so close to the measured values of 15.8 db and 15.86 db, it will be assumed to be correct until proved otherwise.

### CONCLUSIONS

A commercial fluorescent lamp is a reliable source of microwave noise energy. At 4000 mc its effective temperature is 11,394 degrees Kelvin which is convenient for measuring noise figures of 20 db or less. The noise power is practically independent of the fluorescent coating, the current density and only slightly affected by the room temperature. The lamp lends itself readily to a broad-band impedance match in the waveguide.

[8] G. E. Inman and R. N. Thayer, *A. I. E. E. Transactions*, Vol. 57, pp. 723–726, Dec. 1938.

# Electronic Admittances of Parallel-Plane Electron Tubes at 4000 Megacycles

By SLOAN D. ROBERTSON

This paper reports the results of some measurements of the electronic admittances of close-spaced parallel-plane diodes and "1553" triodes at a frequency of 4060 megacycles. These results reveal that the diode admittance and the input short-circuit admittance of the triode depart considerably from the values predicted by single-velocity theory. The triode transadmittance, however, is only slightly lower in magnitude than the low-frequency value.

THE high-frequency admittances of electron streams flowing between parallel-plane electrodes have stimulated considerable theoretical interest. Llewellyn[1,2,3,5] has given an analysis of the particular case in which all electrons in any plane perpendicular to the direction of flow are assumed to have identical velocities. In practice, this approximation gives a reasonably accurate expression for electron stream admittances if the electrode spacing is relatively large, and if the frequency is not so high that the actual spread in electron velocities represents an appreciable fraction of the transit time. Others have treated various aspects of the general problem[4,5,6,7,8,9,10]. Theoretical consideration has also been given to the problem of electron flow in which the electrons possess a Maxwellian velocity distribution[11,12,13,14]. There has been, however, no complete analysis of the microwave-frequency case which takes account of the Maxwellian velocities.

In order to orient the present work properly with previous work let us consider briefly the parallel plane diode shown in Fig. 1, which shows three representative potential distribution curves. If only a relatively few electrons are available at the cathode, the potential distribution between electrodes will be approximately equal to the space-charge-free distribution indicated by curve $a$. If an ample supply of electrons is provided by the cathode and if all electrons leave the cathode with zero velocity, then the space charge is complete in accordance with Child's law, and the potential distribution follows curve $b$. If, on the other hand, the cathode is capable of supplying an ample supply of electrons, the electrons being emitted with a Maxwellian velocity distribution, the potential distribution will be represented by a curve of the type shown by $c$. The cases shown by curves $a$ and $b$ can be treated by the Llewellyn analysis. With wide spacings and at lower frequencies the admittances obtained with distributions of the $c$ type may be approximated by the results obtained by analysis of distributions of the $b$ type. With the very close spacings encountered in the Bell

Laboratories 1553 triode[15] the theoretical analysis no longer represents a valid approximation.

Let us consider curve $c$ in greater detail. The fact that electrons are emitted with a Maxwellian velocity distribution, instead of being emitted at zero velocity as in the Child's law or complete space charge case, means that more electrons are introduced in the space between the electrodes than can flow to the anode in accordance with Child's law. The surplus electrons depress the potential in front of the cathode to a value below that of the cathode. This potential minimum is indicated by $Vm$ in the figure. Electrons which have insufficient energy to cross this barrier return to the cathode.

In the space between the cathode and the potential minimum, electrons are found traveling with various velocities in both directions. Between the potential minimum and the anode, electrons travel in one direction only,



Fig 1—Potential distributions in a diode

toward the anode, but with multiple velocities. With close spacings and higher frequencies the distance between the cathode and the potential minimum may be an appreciable part of the total cathode-anode spacing, with the result that the electrons returning to the cathode may absorb a substantial amount of power from the high-frequency field.

This argument also applies to the cathode-grid region of a microwave triode such as the 1553. In order to increase the transconductance of the triode, it is desirable to locate the grid as close to the cathode as possible. The close spacing, however, leads to a greater loss of power to the returning electrons, which prevents a realization of the full benefits expected from the reduced spacing. All of these difficulties are a result of the Maxwellian velocity distribution of the emitted electrons.

In view of the importance of electron stream admittances in the design of microwave amplifiers and of the need for a better understanding of the performance of the 1553, a program was initiated to investigate some of

these effects experimentally. It seemed best to start this work with a study of the electron stream admittances of simple diodes, with the object of extending the measurements to the triode as the work progressed.

## Diodes

The diodes used in this work were identical in construction with the 1553 triode, but for the substitution of a solid copper anode in place of the grid. In all cases the cathode-anode spacing was approximately 0.65 mil, and the area of the cathode was 0.164 square centimeters. With this spacing one would expect the potential minimum to be relatively close to the anode such that a considerable portion of the cathode-anode region would contain electrons moving in both directions. The potential distribution then would be something like that shown in Fig. 2.



Fig. 2—Electron motion in a close-spaced diode.

The method used in measuring the microwave-frequency input admittances of diodes was based largely on a technique used by Mr. J. A. Morton, and will be described in some detail.

In a typical amplifier, radio-frequency power is fed from a waveguide source to the cathode-grid input region of a 1553 triode through a waveguide-cavity transformer. A similar circuit can be used for measuring diode admittances. The fundamental problem is to learn how to relate admittances measured with a standing wave detector located in the waveguide supply line to the equivalent two-terminal admittances located at the cathode-anode gap of the diode itself. In other words, we have to know the trans-formation-ratio between an admittance across the cathode-anode gap of the diode and the corresponding admittance which will be measured in the waveguide.

Let us refer to the circuit in Fig. 3. The circuit shows an input trans-

mission line which, for example, may be a waveguide having a characteristic impedance $Z_{oy}$, connected through an ideal transformer to an output line having a characteristic impedance $Z_{ox}$. The output line is connected to the transformer at the point $x_o$, where $x_o$ represents the gap terminals of the diode. Suppose for the moment that provision has been made for connecting the output line at the point in the circuit normally occupied by the cathode-anode planes of the diode. This can be done by means of the special testers shown in Fig. 4. In these testers the anode has been omitted and provision has been made for attaching a coaxial line across the gap between the cathode and anode planes. The diodes used in later tests were identical with the device of Fig. 4, except that the coaxial output fitting was replaced by a sheet copper anode.

Referring again to Fig. 3, assume that the output line is shorted at point $x_0$. If power is introduced in the input line at the left, a standing wave pattern in the input line will pass through a minimum at some point $y_0$.



Fig. 3—Equivalent circuit of diode measuring equipment.

If the short circuit is now moved to the right by an increment $\Delta x$, the standing wave minimum will move by an increment $\Delta y$. The relation between $\Delta x$ and $\Delta y$ is given by the following equation:

$$\frac{1}{Z_{0y}} \cot \frac{2\pi\Delta y}{\lambda_y} = \frac{\phi}{Z_{0x}} \cot \frac{2\pi\Delta x}{\lambda_x} - \phi B_0 \tag{1}$$

where $\lambda_y$ and $\lambda_x$ are the respective wavelengths in the two lines (which may not be equal if, for example, one is a coaxial and the other is a waveguide). $\phi$ is the transformation ratio of the ideal transformer, and $B_0$ is the effective leakage susceptance of the tube and transformer referred to the terminals at $x_0$. If $\frac{2\pi\Delta y}{\lambda_y}$ is plotted as a function of $\frac{2\pi\Delta x}{\lambda_x}$ on cot-cot coordinate paper, a straight line is obtained whose slope $m$ is

$$m = \phi \frac{Z_{0y}}{Z_{0x}} \tag{2}$$

and whose ordinate intercept $\rho$ is

$$\rho = -\phi B_0 Z_{0y} \tag{3}$$

A typical cot–cot plot is shown in Fig. 5.

Now, assume that the right-hand transmission line is removed and that the diode gap is connected at the transformer terminals $x_0$. The normalized admittance referred to the point $y_0$ on the input line can be measured by a simple standing wave measurement. Represent this admittance by $\bar{Y}_{wg}$.



Fig. 4—Coaxial tester.

Let the unknown diode admittance be represented by $Y_x$. $Y_x$ is then given by the following relation:

$$Y_x = \frac{1}{Z_{0x}m} [\bar{Y}_{wg} + j\rho] \qquad (4)$$

Hence, having determined $y_0$, it is only necessary to measure the slope $m$ and the intercept $\rho$ on the cot–cot curve in order to relate $Y_x$ to $\bar{Y}_{wg}$. The characteristic impedance of the output line $Z_{0x}$ used in obtaining the cot–cot plot must also be known. Since a coaxial is used for this line, its characteristic impedance is easily calculated.

If no losses were associated with the transformer or the parts of the diode external to the actual cathode-anode region, such as the metal vacuum

envelope and certain ceramic details of the tube, the above measurements would give complete information regarding the circuit. Certain losses have been found, however. These are measured as follows: At the time when terminals $x_0$ are shorted a standing wave measurement is made in the wave-



Fig. 5—Typical cotangent–cotangent plot.

guide line at the left. From this measurement and the cot–cot data it is possible to compute an equivalent resistance in series with the gap caused by losses present in the circuit. This equivalent series resistance is given by

$$R_s = \frac{Z_{0x} m}{SWR} \qquad (5)$$

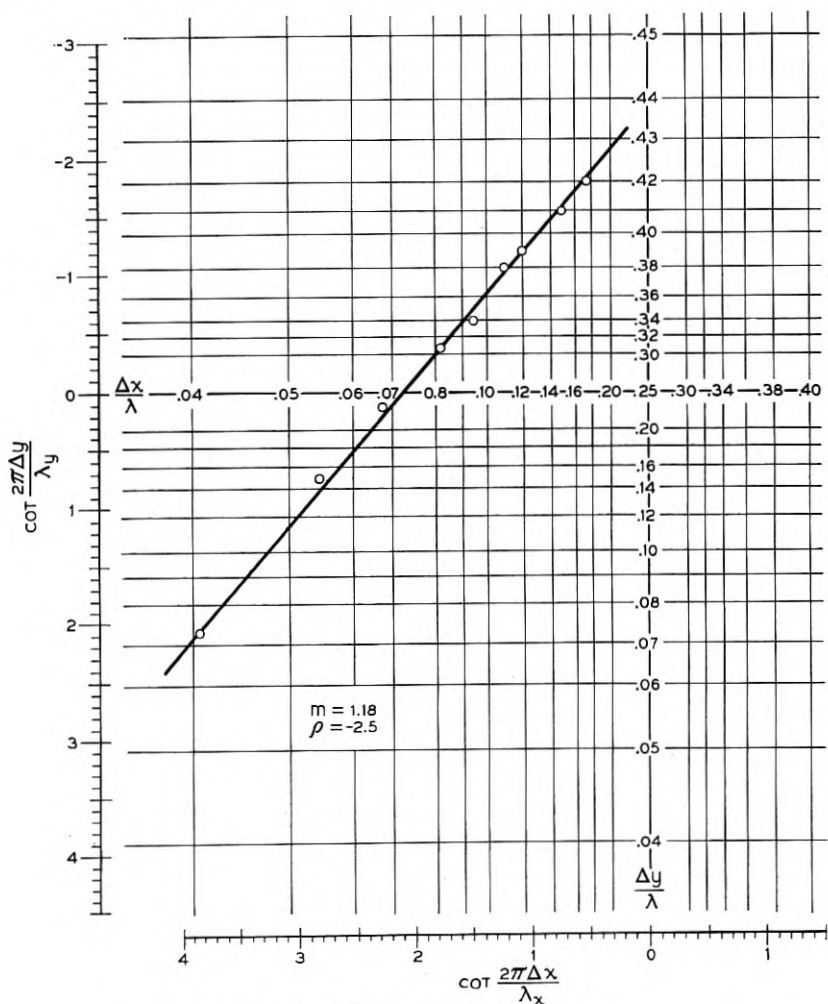where $SWR$ is the voltage standing wave ratio mentioned above. The determination of a series loss resistance in this manner is quite analogous to the short-circuit test used in determining the losses in a power transformer.

There is one other factor in the cot–cot technique which is worthy of mention. If, at the very beginning, the output line is terminated in $Z_{0x}$ and if the transformer is adjusted so that the input line is matched, then the value of $m$ will be unity and $\rho$ will equal zero. It is then unnecessary to take a cot–cot curve. It is, however, still necessary to locate $y_0$ by shorting the terminals at $x_0$.

### DIODE ADMITTANCE AT 4060 MEGACYCLES

Electron stream admittance measurements with diodes were made in the following way: A coaxial tester was installed and the circuit was adjusted for a slope $m$ of about one. This coaxial tester was then removed and replaced by another in order to learn whether the slope obtained with one tester would be the same with another, supposedly identical, tester. This process was repeated several times, and the slope was found to vary no more than about 10% from one tester to the other.

The procedure was then to replace the coaxial tester with a diode and make admittance measurements with the assumption that the slope would be the same for the diode as for the tester. This assumption was believed to be reasonable since the structure of the diode was identical with that of the tester except that an anode was substituted for the coaxial output connector. In either case all elements that were located inside the waveguide cavity were presumably identical.

Electron stream measurements were made at a frequency of 4060 megacycles with a number of diodes over a wide range of anode and heater voltages. In making these measurements, the radio-frequency power was kept at a relatively low level (0.2 milliwatt) in order that the measured admittances would be independent of the radio frequency voltage.

Results for several diodes are shown in Figs. 6 through 13. The various symbols used in the figures are defined as follows:

$V_H$ = heater voltage

$I_H$ = heater current

$V_0$ = anode voltage (neglecting contact potentials)

$I_0$ = anode current in ma

$J_0$ = anode current density in ma/cm²

$g_0$ = low-frequency diode conductance measured with an audio frequency bridge

$g$ = high-frequency diode conductance measured as described above

$b$ = high-frequency diode susceptance

$R_s$ = equivalent resistance in series with diode

In computing the admittance of the electron stream it was necessary to

Fig. 6—Admittance of a diode.

allow for the circuit and tube losses previously discussed. The equivalent series resistance $R_s$ of the diode circuit was determined by biasing the tube negatively to the point where a further increase in bias failed to produce a

perceptible change in the waveguide standing wave ratio. Under such conditions the electrons experienced a large retarding field at the cathode and did not emerge an appreciable distance into the cathode-anode region. Any resistance measured at this time was due to the series loss and was not



Fig. 7—Effect of heater voltage upon diode conductance.

produced electronically. The diode series resistances varied from about 1.3 to 5.0 ohms with an average value around 3 ohms.

Figure 6 shows the results of admittance measurements of a diode. As expected, the high-frequency conductance is considerably greater than the low-frequency value $g_0$. In fact $g$ is seen to have a value of several thousand micromhos when the negative bias of the tube is such that no perceptible anode current flows. The susceptance $b$ for large negative anode potentials

has a value of 150,000 micromhos, which agrees fairly well with the value computed from the geometrical capacitance. As anode current is drawn and a space charge condition prevails, $b$ drops to a value of 125,000 micromhos. Theoretical considerations would predict a drop of about 40% in the case of a single-velocity electron stream. This is somewhat greater than the drop exhibited in Fig. 6.



Fig. 8—Effect of heater voltage upon $g/g_0$.

Figures 7 and 8 show the effect of cathode temperature on $g_0$ and the ratio $g/g_0$. The parameter used to represent the cathode temperature is the heater voltage $V_H$. As the heater voltage is raised the total conductance $g$ increases. The ratio $g/g_0$, however, decreases, particularly for low or negative anode voltages. This means that, with a given anode voltage, as the cathode temperature is raised, $g_0$ increases more rapidly than $g$. If the curves of Fig. 8 are replotted in terms of $J_0$ rather than $V_0$, the ratio $g/g_0$ is relatively independent of $V_H$. This is shown in Fig. 9.

The results of measurements on another diode are shown in Fig. 10.

These are very similar in all respects to those of the preceding figure. It is probable that the cathode-anode spacings of the two diodes of Figs. 6 and 10 were somewhat greater than the 0.65 mil for which they were designed. In both cases the capacitances measured at low frequency were somewhat low.



Fig. 9—Variation of $g/g_0$ with current density and heater voltage.

In Fig. 11, results are shown for a third diode. In this case the susceptance at a large negative bias is in almost exact agreement with the value to be expected with the intended diode spacing of 0.65 mil. It is interesting to observe that, with this tube, $b$ drops a greater amount as the current increases. Moreover, the ratio $g/g_0$ is greater than that found with earlier diodes.

In Fig. 12 data are shown for a diode having a very high value of $g_0$. From the standpoint of cathode activity this was the best tube that was

Fig. 10—Admittance of a diode.

tried. At maximum current the susceptance $b$ dropped to 50% of the initial value. The data of Fig. 12 have been replotted in Fig. 13 in terms of the variable $126x^{\frac{1}{3}}/\lambda J_0^{\frac{1}{3}}$, where $x$ is the cathode-anode spacing. In the Llewellyn

Fig. 11—Admittance of a diode.

theory this variable is equal to the transit time. The solid curves in the figure are the theoretical results of the Llewellyn theory, whereas the broken curves present the corresponding experimental values. In the latter it should

Fig. 12—Admittance of a diode.

be understood that the abscissa do not represent transit time. The curves do serve, however, to compare the theoretical diode resulting from a single-valued electron velocity assumption with the actual diode in which a Max-

wellian velocity distribution prevails. In the experimental case it is probable that, for values of the abscissa greater than 6 or 7, the actual transit time is considerably greater than in the theoretical case. In fact, at a value of 11.4 the anode voltage was zero, the anode current being maintained by the thermal energy of the electrons.



Fig. 13—Comparison of theoretical and experimental values of diode conductance and susceptance.

Other diodes were tested, but they exhibited results substantially equivalent to those already disclosed. In a few cases anomalous results were obtained. With some diodes the capacitance with no electron flow did not approach the low-frequency value. These were rejected on the assumption that there was some mechanical imperfection in the tube which changed the calibration of the measuring equipment.

With the realization that sufficient data are not available to define the phenomena in all detail, it is believed that certain general conclusions can be drawn. From the present work and that of Lavoo[16] and others[17,18,19], it is apparent that the microwave conductance of a close-spaced diode is substantially greater than the low-frequency value. The ratio $g/g_0$ appears to increase as the spacing decreases. This increase will probably continue until the position of the potential minimum approaches the anode plane. The susceptance decreases with increasing current and appears to level off at high-current densities. The final value at a current density of 240 ma/cm² varied between 0.5 and 0.9 of the initial value.

For a given current density, the ratio $g/g_0$ does not appear to vary appreciably as the cathode temperature is changed.

An attempt was made to study the available diodes at 10,000 megacycles. It was found, however, that the value of $R_s$ was so high at this frequency and that variations in tube conductance were so small in comparison with $R_s$ that accurate results could not be obtained.



Fig. 14—Equivalent circuit of a triode.

## FOUR-POLE ADMITTANCES OF A TRIODE

A triode may be considered as an active linear four-pole transducer, and may be defined by the network of Fig. 14. It is apparent that

$y_{11}$ is the input admittance with the output shorted,

$y_{22}$ is the output admittance with the input shorted,

$y_{12}$ is the feedback admittance with the input shorted,

$y_{21}$ is the transadmittance with the output shorted.

The values of the parameters $y_{11}$, $y_{22}$, $y_{12}$, and $y_{21}$ to be measured at the grid, cathode, and anode terminals differ from the values of the $y$ admittance coefficients given by Llewellyn and Peterson[5] who define $y_{11}$ as the admittance of the diode coinciding with the cathode and the fictitious equivalent grid plane, and $y_{22}$ as the admittance between the equivalent grid plane and the anode, and finally $y_{21}$ as the transadmittance between the two. The relations between the $y$ admittance coefficients of Llewellyn and Peterson and the coefficients measured by the author are given by Peterson.[6] It turns out that, with a high-mu tube, such as the 1553 triode, the two sets of

coefficients differ in the order of 10–20% over the useful operating range of current densities; so, for practical considerations, the measure coefficients may be regarded as substantially equivalent to the coefficients referred to the fictitious grid plane. Not that they will be equal to the theoretical values, but they may be regarded as being associated with the same geometry and will serve at least as a qualitative test of the validity of the theoretical values for the physical tube.

In order to measure the four-pole parameters, the 1553 triode was mounted in a coaxial circuit of the type shown in Fig. 15. The grid-anode output circuit of the tube is seen to connect directly with the coaxial output line. The input circuit required a more careful design. Due to the size of the base of the tube it was necessary to taper the input coaxial as shown. In the early stages of this work, difficulty was experienced with higher order modes in the large diameter section of the input coaxial. It was believed that these modes were generated by the action of the parallel wire grid which lacked the



Fig. 15—Detail of coaxial mount for measuring four-pole admittances of a triode.

radial symmetry appropriate to coaxial transmission. The difficulty was overcome by constricting the outer diameter of the coaxial line in the immediate vicinity of the grid of the tube, thus inhibiting generation of the higher order mode.

Before measurements could be made it was necessary first to calibrate both the input and the output circuits in a manner similar to that used and described in connection with the diode measurements. The coaxial tester used for calibrating the input circuit was identical with that used for the diode work. For the output circuit a similar tester was use. As one might expect, the value of the cot–cot slope of the output circuit was close to unity. The value actually turned out to be 0.9. In the input circuit the slope was so great that it was difficult to measure, so that it was necessary to introduce a transformer in the coaxial input circuit to permit tuning.

The complete apparatus necessary to measure $y_{11}$ and $y_{22}$ is shown in Fig. 16. This equipment, save for the details already discussed, is quite conventional in every respect.

In order to measure $y_{11}$, the output coaxial line was short-circuited at a point an integral number of half-wave-lengths from the grid-anode terminals of the tube. The admittance measured in the input line could then be used in computing $y_{11}$. To measure $y_{22}$, the procedure was reversed, the input line being shorted, and the corresponding admittance being measured in the output line. In either case the normalized line admittances were measured by the standard procedure of determining the standing wave ratio in the line and locating the position of the standing wave minimum with respect to the equivalent terminals of the tube.

The transfer admittances were measured with the equipment shown in Fig. 17. The equipment shown here has been fully described in a recent



Fig. 16—Circuit connected for measuring input short-circuit admittance of a triode.

paper[20] and will be described only briefly here. The output of a signal oscillator is divided into two portions. One portion is applied to a balanced modulator where it is modulated by an audio-frequency signal. The suppressed-carrier, double-sideband signal from the modulator is applied to the input circuit of the triode. Probes are provided for sampling the voltages $V_1''$ and $V_2''$ at points an integral number of half wavelengths from the input and output gaps of the tube respectively. The other portion of the oscillator power is fed through a calibrated phase shifter and is applied to a crystal detector in the manner of a local oscillator of a double-detection receiver. The signal samples at $V_1''$ and $V_2''$ are then alternately applied to the crystal detector where they are demodulated by the action of the homodyne carrier. In each case the phase shifter is adjusted so that the audio signal disappears in the detector output. This occurs when the phase of the homodyne carrier

is in quadrature with the signal sidebands. The difference in phase between the two adjustments of the phase shifter is equal to the phase between $V''_1$ and $V''_2$. In measuring the transfer phase from $V''_1$ to $V''_2$ the output coaxial line is terminated in its characteristic impedance. By reversing this procedure it is possible, of course, to measure the ratio of $V''_2$ to $V''_1$ with the input circuit terminated in $Z_0$. The ratio of the magnitudes of $V''_1$ and $V''_2$ may be measured either with the equipment shown in Fig. 17 by adjusting the phase of the homodyne carrier to maximize the signals in each case and



Fig. 17—Circuit for measuring transfer phase of a triode.

comparing the levels, or by using the equipment in Fig. 16 in the conventional way.

Figure 18 is a photograph of the triode circuit which shows the input and output coaxial standing-wave detectors with the triode mounted in the enlarged section at the center.

As in the case of the diode it was found that, with the tube biased negatively such that no electrons could leave the immediate vicinity of the cathode, the input circuit exhibited an equivalent series resistance $R_s$. The latter had to be allowed for in reducing the experimental data.

Fig. 18—Coaxial mount for measuring triode admittances.

The experimental data obtained as described above were sufficient for computing the four-pole parameters. The calculations necessary for the reduction of the data can best be understood by referring to Fig. 19. The various symbols used in connection with the figure are defined as follows:

$\bar{Y}_1$ = Normalized admittance measured at 1–1 with 2–2 shorted

$\bar{Y}_2$ = Normalized admittance measured at 2–2 with 1–1 shorted

$\dfrac{V_1''}{V_2''}$ = $\gamma_{21}$ (measured with output line terminated in $Z_0$)

The above parameters represent those obtained by the measurements described above.



Fig. 19—Equivalent circuit of triode and associated measuring equipment.

In calibrating the circuit the following parameters were obtained:

$\rho_1$ = ordinate intercept of input cot–cot curve

$\rho_2$ = ordinate intercept of output cot–cot curve

$m_1$ = slope of input cot–cot curve

$m_2$ = slope of output cot–cot curve

$$B_{01} = -\frac{\rho_1}{66m_1} \qquad B_{02} = -\frac{\rho_2}{66m_2}$$

$Z_0$ = characteristic impedance of input and output coaxial lines.

$R_s$ was measured by shorting the output line, placing a large negative bias on the tube, and measuring the admittance of the input line. Then

$$R_s = 66m_1 Re(\bar{Y}_1) \tag{6}$$

where the number 66 represents the characteristic impedance of the coaxial line used in calibrating the input circuit, corresponding to $Z_{0x}$: in Equation 4.

Fortunately for simplicity, the series resistance in the output circuit was negligible.

The computations are then as follows:

$$y_{11}' = \frac{\bar{Y}_1}{66m_1} \qquad y_{22}' = \frac{\bar{Y}_2}{66m_2} \tag{7}$$

$$y_{11} = \frac{\bar{Y}_1}{66m_1 - \bar{Y}_1 R_s} + \frac{j\rho_1}{66m_1} \tag{8}$$

$$y_{22} = \frac{1}{66m_2} [\bar{Y}_2 + j\rho_2] \tag{9}$$

In order to compute $y_{21}$, the following four-pole equations are used:

$$V_1' = \frac{I_1'}{y_{11}'} + \frac{V_2 y_{12}'}{y_{11}'} \tag{10}$$

$$V_1 = \frac{I_1}{y_{11}} + \frac{V_2 y_{12}}{y_{11}} \tag{11}$$

$$V_2 = \frac{I_2}{y_{22}} + \frac{V_1 y_{21}}{y_{22}} \approx \frac{I_2}{y_{22}} + \frac{V_1' y_{21}'}{y_{22}} \tag{12}$$

It follows that

$$V_1 y_{21} \approx V_1' y_{21}' \tag{13}$$

$$y_{21} \approx y_{21}' \frac{V_1'}{V_1} \tag{14}$$

Referring to Fig. 19, one may write

$$V_1 = V_1' - (I_1' + V_2 y_{12}')R_s \tag{15}$$

Combining (10) and (15)

$$\frac{V_1'}{V_1} = \frac{1}{1 - y_{11}' R_s} \tag{16}$$

$y_{21}'$ can be evaluated by making use of the relation

$$V_2 \approx \frac{I_2'}{y_{22}'} + \frac{V_1' y_{21}'}{y_{22}'} \tag{17}$$

Dividing (17) by $V_2$ and rearranging terms

$$y_{21}' \approx \frac{y_{22}'}{\left(\dfrac{V_1'}{V_2}\right)} \left[1 - \frac{I_2'}{V_2 y_{22}'}\right] \tag{18}$$

where $I_2'/V_2$ can be expressed as

$$\frac{I_2'}{V_2} = \frac{1}{66m_2 \bar{Z}_0} = \frac{1}{66m_2} \tag{19}$$

where $\bar{Z}_0 = 1$.

$V_1'/V_2$ can be expressed in terms of $\gamma_{21}$, $m_1$, and $m_2$ by using the relations:

$$\frac{V_1'}{V_1''} = \sqrt{\frac{66m_1}{Z_0}} \qquad \frac{V_2}{V_2''} = \sqrt{\frac{66m_2}{Z_0}} \tag{20}$$

Solving (20) for $V_1'/V_2$ and remembering that $V_1''/V_2'' = \gamma_{21}$,

$$\frac{V_1'}{V_2} = \gamma_{21} \sqrt{\frac{m_1}{m_2}} \tag{21}$$

If (19) and (21) are substituted in (18), one finds

$$y_{21}' \approx \frac{y_{22}'}{\gamma_{21}} \sqrt{\frac{m_2}{m_1}} \left[ 1 + \frac{1}{66 m_2 y_{22}'} \right] \tag{22}$$

By using (14) and (16), $y_{21}$ can then be written as

$$y_{21} \approx \frac{y_{22}'}{\gamma_{21}} \sqrt{\frac{m_2}{m_1}} \left[ 1 + \frac{1}{66 m_2 y_{22}'} \right] \left[ \frac{1}{1 - y_{11}' R_s} \right] \tag{23}$$

Several 1553 triodes were available for study. Typical experimental results obtained with two of them are shown in Figs. 20, 21, and 22. The triode used in obtaining the data of Fig. 20 had input and output spacings of 0.65 and 12 mils, respectively. The cathode and anode diameters were 180 mils. The grid opening was 250 mils and was wound with 0.3 mil tungsten wire at 1000 strands per inch. In the figures, $V_g$ and $V_p$ represent the d-c. grid and plate potentials, respectively.

There are a number of interesting things to observe in Fig. 20. As with the diode, $b_{11}$ for a large negative bias approaches the "cold" value computed from the capacitance. However, as anode current is drawn, $b_{11}$ drops rapidly to a much lower value than was the case for the diodes. The conductance $g_{11}$ behaves somewhat like $g$ for the diode. $b_{22}$ is equal to the value computed from the grid-anode capacitance and is not appreciably influenced by the electron stream. $g_{22}$ was very low with a magnitude of slightly less than 1000 micromhos at maximum anode current. It is not shown in the figure. The transadmittance $y_{21}$ is worth considering. When the bias is several volts negative, $y_{21}$ has a value of about 9000 micromhos. This is about 50 times as high as one would expect from a consideration of the electrostatic capacitance between the cathode and anode of the tube. This effect has been investigated more fully and is discussed in another paper.[21] As the tube starts to draw plate current, $y_{21}$ rises and reaches a maximum of about 40,000 micromhos. The low-frequency transconductance was measured and is plotted in the figure. It will be observed that the high-frequency transadmittance is only slightly lower than $g_m$. This is in agreement with the theories of Llewellyn.[5] The agreement appears reasonable when one remembers that, in the theoretical analysis, the magnitude of the ratio $y_{21}/g_0$ is relatively independent of the transit time in the input space.

Figure 21 shows the results of measurements on a triode identical with that of Fig. 20 except that the grid consists of a mesh of 0.3 mil tungsten wires wound at 550 strands per inch in both directions. It will be noted

Fig. 20—Four-pole admittances of a triode having a parallel-wire grid.

that $y_{21}$ is much lower when this tube is biased beyond cutoff than in the previous case. The electromagnetic coupling is therefore much less for the mesh grid. This has also been treated in the above reference.[21] With high negative bias the feedback admittance $y_{12}$ was substantially equal to $y_{21}$

Fig. 21—Four-pole admittances of a triode having a cross-lateral grid.

but, as the current density increased, $y_{12}$ tended to decrease. The feedback admittance was always lower for the mesh grid than for the parallel-wire grid.

The remaining parameters for the triode of Fig. 21 are very similar to those of Fig. 20.

Figure 22 shows the variation of the phase of the transadmittances $y_{21}$ for the two triodes. The figure also shows the theoretical curve of the Llewellyn analysis for purposes of comparison. As in the case of Fig. 13 the abscissa do not represent transit time for the experimental values. The quantity $x$ is equal to the cathode-grid spacing.

It is of interest to compare the triode measurements with those of the diode. It was expected that $g_{11}$ for the triode should correspond with $g$ for the diode. Within the limits of reasonable experimental accuracy this appears to be the case. For the triode at low frequencies $g_0 \approx g_m$. The triode



Fig. 22—Phase of triode transadmittance.

results indicate that the ratio $g_{11}/g_m$ is quite comparable in magnitude with the corresponding ratio $g/g_0$ for the diode. This was expected. The behavior of $b_{11}$ for the triode was unexpected. It was thought that, as the grid voltage was varied so that the input space changed from a condition of zero space charge to one of maximum space charge, $b_{11}$ would vary from its initial "cold" value to a value approaching 60% of the latter. This was not so. In the figures one observes that it drops to a much lower value. This effect has not been explained from a theoretical standpoint. There are several qualitative interpretations, but as yet no way of determining which of them is correct in a quantitative sense has been found. The observed phenomenon could, for example, be explained by an increase in the effective series resistance of the tube caused perhaps by an increase in the resistance of the

cathode coating.[14] Since the effect was not observed to such a marked degree in the case of the diodes, it seems probable that this is not the correct explanation.

It is probable that the observed variation in $b_{11}$ is a space charge effect. It is evident in examining the diode curves that tubes which possessed the higher values for $g_0$ exhibited a greater variation in $b$. If maximum $g_0$ can be taken as a measure of the cathode activity, we can then perhaps relate the variation in susceptance with cathode activity and hence with the location of the potential minimum. A shift in the position of the potential minimum, however, may produce two effects. It varies the transit time of the electrons and changes the degree of space charge in the input space. Either effect might account for the variation of $b_{11}$. A clue to this effect might be discovered by making measurements on structures with different cathode-grid spacings.

The following experiments were performed to determine the effect of plate voltage on the input admittance of the triode of Fig. 20. The plate and grid voltages were varied simultaneously in such a way that the sum of the direct currents to the grid and plate remained constant at 30 milliamperes corresponding to a current density of 184 ma/cm². The input admittance did not vary from the value shown for this same current density in Fig. 20 even though the plate voltage was varied from 250 volts to 40 volts. In a second experiment the plate potential was maintained at −90 volts with respect to the cathode and the grid potential was varied such that the direct grid current varied over a range of 0 to 10 milliamperes. Again the admittances were found to be equal to those of Fig. 20 for the corresponding total currents. These two experiments suggest that, for a given geometry, the value of $b_{11}$ is primarily a function of the total current density in the input circuit.

## REFERENCES

1. "Electron Inertia Effects," F. B. Llewellyn, Cambridge University Press.
2. "Equivalent Networks of Negative Grid Vacuum Tubes at Ultra-High Frequencies," F. B. Llewellyn, *B. S. T. J.*, Vol. 15, pp. 565–586, October 1936.
3. "Operation of Ultra-High Frequency Vacuum Tubes," F. B. Llewellyn, *B. S. T. J.*, Vol. 14, pp. 632–665, October 1935.
4. "Theory of the Internal Action of Thermionic Systems at Moderately High Frequencies," W. E. Benham, *Phil. Mag.*, Vol. 5, pp. 641–662, March 1928; and Vol. 11, pp. 457–517, Feb. 1931.

5. "Vacuum-Tube Networks," F. B. Llewellyn and L. C. Peterson, *Proc. I. R. E.*, Vol. 32, no. 3, pp. 144–166, March 1944.
6. "Equivalent Circuits of Linear Active Four-Terminal Networks," L. C. Peterson, *B. S. T. J.*, Vol. XXVII, No. 4, pp. 593–622, October 1948.
7. "Impedance Properties of Electron Streams," L. C. Peterson, *B. S. T. J.*, Vol. 18, pp. 465–481, July 1939.
8. "Klystron and Microwave Triodes," Hamilton, Knipp and Kuper, pp. 97–169, *Radiation Laboratory Series*, Vol. 7, McGraw-Hill, 1948.
9. "High-Frequency Total Emission Loading in Diodes," Nicholas A. Begovich, *Journal Applied Physics*, Vol. 20, No. 5, pp. 457–461, May 1949.
10. "On the Velocity-Dependent Characteristics of High-Frequency Tubes," Julian K. Knipp, *Journal Applied Physics*, Vol. 20, No. 5, pp. 425–431, May 1949.
11. I. Langmuir, *Phys. Rev.*, 21, pp. 419–435, 1923.
12. "Extension and Application of Langmuirs' Calculations on a Phase Diode with Maxwellian Velocity Distribution of the Electrons," A. Van Der Ziel, *Philips Research Reports*, Vol. 1, No. 2, pp. 97–118, January 1946.
13. "Extension of Langmuir's ($\xi$, $\eta$) Tables for a Plane Diode with a Maxwellian Distribution of the Electrons," P. H. J. A. Klegmen, *Philips Research Reports*, Vol. 1, No. 2, January 1946, pp. 81–96.
14. "Some Characteristics of Diodes with Oxide-Coated Cathodes," W. R. Ferris, *R. C. A. Review*, Vol. X, No. 1, pp. 134–149, March 1949.
15. "A Microwave Triode for Radio Relay," J. A. Morton, *Bell Laboratories Record*, Vol. XXVII, No. 5, May 1949.
16. "Transadmittance and Input Conductance of a Lighthouse Triode at 3000 Megacycles," Norman T. Lavoo, *Proc. I. R. E.*, Vol. 35, No. 11, pp. 1248–1251, November 1947.
17. "Total Emission Damping in Diodes," A. Van Der Ziel, *Nature*, Vol. 159, No. 4046, May 17, 1947, pp. 675–676, (52 mc).
18. "Total Emission Damping," J. Thomson, *Nature*, Vol. 161, No. 4100, pp. 847, May 29, 1948.
19. "Total Emission Damping with Space-Charge Limited Cathodes," C. N. Smyth, *Nature*, 157, 841, June 22, 1946.
20. "A Method of Measuring Phase at Microwave Frequencies," S. D. Robertson, *B. S. T. J.*, Vol. XXVIII, No. 1, pp. 99–103, January 1949.
21. "Passive Four-Pole Admittance of Microwave Triodes," S. D. Robertson, this issue of the *B. S. T. J.*
22. "Total Emission Noise in Diodes," A. Van Der Ziel and A. Versnel, *Nature*, Vol. 159, No. 4045, pp. 640–641, May 10, 1947.
23. "Ultra-High-Frequency Oscillations by Means of Diodes," F. B. Llewellyn and A. E. Bowen, *B. S. T. J.*, Vol. 18, pp. 280–291, April 1939.

# Passive Four-Pole Admittances of Microwave Triodes

## By SLOAN D. ROBERTSON

Measurements have been made of the passive, four-pole admittances of parallel-plane triodes over a wide range of cathode-to-grid and grid-to-plate spacings at a frequency of 4060 megacycles. Results are given for a parallel wire grid and a cross-lateral grid. The microwave transadmittances are found to be much higher than the values measured at low frequencies.

DURING the course of an experimental study of the active four-pole admittances[1] of the 1553 close-spaced triode,[2] a question arose as to whether the grid wires were introducing any appreciable inductance or resistance in the circuit used for measurement. It appeared necessary, therefore, to learn something of the passive four-pole parameters of the triode in order to separate the electronic from the passive admittances. It was generally believed that the electrostatic analyses of the passive admittances which have been successfully applied at the lower frequencies would no longer be valid with close-spaced structures at microwave frequencies. For example, it was considered possible that the grid wires themselves might possess an effective inductive reactance, so that the admittances between the grid and cathode or between the grid and anode might not be equal to the values computed from the electrostatic capacitances. Moreover, it was thought likely that energy might be transmitted from the cathode-grid region to the cathode-plate region or vice versa, not only by the medium of the electrostatic coupling, but also by means of an electromagnetic coupling through the grid. The measurements to be reported below indicate that the first of these conjectures was false, but that the second was true.

In view of the lack of available information on these questions in general, it seemed highly desirable to employ the available measuring equipment, not only to determine the passive parameters of a triode having electrode spacings corresponding with those of the 1553, but to extend the scope of the measurements to include a wide range of electrode spacings in order that the results would be of more general interest.

Although these measurements were in principle very simple, in practice the mechanical problem of achieving the desired degree of accuracy proved rather difficult. It was required that the cathode, grid, and anode planes be almost perfectly parallel and that the spacings between them be adjustable

[1] S. D. Robertson, "Electronic Admittances of Parallel-Plane Electron Tubes at 4000 Megacycles," this issue of the B. S. T. J.

[2] J. A. Morton, "A Microwave Triode for Radio Relay," *Bell Laboratories Record*, Vol. XXVII, No. 5, pp. 166–170, May 1949.

to specific values with a high degree of precision. In order to equal the dimensional tolerances of the 1553 it was necessary that parallelism and spacing be accurate to 0.1 mil.

A schematic diagram of the apparatus is shown in Fig. 1. A flat, circular disc having a 250-mil diameter aperture, across which the grid was stretched, was mounted upon the face of the hollow micrometer screw ⚡1. The latter was mounted so that its face was accurately parallel with the end face of the central conductor of the input coaxial line in the upper part of the figure. By means of the micrometer ⚡1 the input spacing $S_1$, which we shall consider as representing the cathode-grid spacing, could be adjusted. The central conductor of the coaxial line was insulated at d.c. from the outer conductor; hence it was possible to use an ohmmeter to indicate when the grid was just touching the coaxial face. The micrometer could then be backed away from the grid by any desired amount. The input coaxial was fitted with a standing wave detector in the form of a probe which could be moved along the line and placed at any arbitrary distance $h$ from the grid.

On the output side of the circuit, in the lower part of the figure, there was another coaxial line arranged so that its center conductor could be driven by micrometer ⚡2. The latter was insulated from the outer conductor of the coaxial by means of a condenser in order that an ohmmeter could be used to determine the position of the micrometer which caused the central conductor to just touch the grid. Spacing $S_2$ could then be adjusted. The output coaxial line was terminated in its characteristic impedance of 62 ohms. At a distance of $\lambda/2$ from the grid a probe was located for sampling the power in the output line.

The diameter of the input coaxial conductor was 180 mils at the end. In the figure it will be noted that at a short distance from the end the diameter increased to a larger diameter (250 mils). Because of the required length of the central conductor, it was necessary to increase its size in this way for mechanical rigidity. The effect of this change in cross-section was computed and allowed for in the final results. The output coaxial conductor was relatively short, so that it was possible to assign a diameter of 180 mils for its entire length. The 180-mil diameter was selected to correspond with the diameters of the cathode and anode in the 1553 triode.

The procedure for making the measurements was as follows: With a particular set of spacings $S_1$ and $S_2$ the standing wave ratio in the input line was measured. This ratio, together with the measurement of the position of a standing wave minimum, permitted the calculation of an input admittance $Y$ to be made. Then with the standing wave detector probe placed at a distance $h = \lambda/2$ from the grid, the ratio of the voltage at the input terminals of the triode to the voltage appearing at the output probe was measured both as to magnitude and phase as described in a recent

Fig. 1—Apparatus for measuring passive admittances of triode.

paper.[3] This quantity will be called $\gamma$. These measurements were sufficient for an evaluation of the four-pole parameters of the structure. All measurements were made at a frequency of 4060 megacycles.

The equivalent circuit of the passive triode structure is shown in Fig. 2. The desired parameters are $y_{11}$, $y_{12}$, and $y_{22}$. The following equations indi-



Fig. 2—Equivalent passive circuit of a triode.



PARALLEL–WIRE GRID

CROSS–LATERAL GRID

Fig. 3—Types of grids used in the measurements.

cate the relation between these parameters and the measured quantities $Y$ and $\gamma$:

$$y_{11} = Y + \frac{62y_{12}^2}{62y_{22} + 1} \approx Y \tag{1}$$

$$y_{12} = \frac{y_{22}}{\gamma}\left[1 + \frac{1}{62y_{22}}\right] \tag{2}$$

where the number 62 represents the output terminating impedance. For all cases to be described here the second term on the right side of Equation 1 is small in comparison with $Y$. This is a result of the small values encountered for $y_{12}$. To a good approximation $y_{11}$ is equal to the measured input admittance $Y$. This was verified by observing the variation in input admittance as the output spacing was varied while keeping the input spacing fixed. Only a slight variation in admittance was observed, which indicated that the fractional term in Equation 1 was small in comparison with $Y$.

Suppose, then, that for a given input and output spacing $S_1$ and $S_2$,

[3] "A Method of Measuring Phase at Microwave Frequencies," S. D. Robertson, *Bell System Technical Journal*, Vol. XXVIII, No. 1, pp. 99–103, January 1949.

$Y$ and $\gamma$ are known. $y_{22}$ can readily be determined by readjusting the input spacing to equal the output spacing and measuring a second admittance $Y'$. $y_{22}$ will be approximately equal to this value. There is, then, sufficient information to compute $y_{12}$.



Fig. 4—Variation of passive input and output admittances with spacing.

Two grids were used in this work. The first was a parallel wire grid of 0.3 mil tungsten wire wound at 1000 turns per inch. The second was also of 0.3 mil tungsten wound in a crisscross fashion at 550 turns per inch. Both grids are shown in Fig. 3. It will be noted that the cross-lateral grid has an aperture 220 mils in diameter.

The values of $y_{11}$ and $y_{22}$ were found to be almost purely capacitive and were the same for both types of grid. These values are shown in Fig. 4. $y_{11}$ and $y_{22}$ correspond to capacitances $C_{11}$ and $C_{22}$, which agree surprisingly well with the calculated capacitances between the grid and cathode, and grid and plate planes, respectively. Figure 5 shows the experimentally

determined values of $C_{11}$ and $C_{22}$ plotted as a dashed curve. The theoretical values (neglecting fringing capacitance) are shown by the solid curve. Since fringing was neglected, it is not surprising that the measured capacitances should exceed the calculated values by the amount shown.



Fig. 5—Comparison of theoretical and experimental values of input and output capacitances.

The magnitudes of $y_{12}$ for each grid over a range of values of $S_1$ and $S_2$ are shown in Figs. 6 and 7. It will be noted that, for a given set of spacings $S_1$ and $S_2$, $y_{12}$ is much greater for the parallel wire grid than for the cross-lateral. This is the sort of result one would expect if $y_{12}$ resulted from electromagnetic coupling through the grid, since the parallel wire grid would be expected to offer a better transmission path than the cross-lateral grid. It was not practicable with the equipment used in these experiments to measure the values of $y_{12}$ at low frequencies where $y_{12}$ would be determined by the cathode-plate capacitance. Data were available, however, for the low-frequency, cathode-plate capacitance of the standard, parallel-wire

Fig. 6—Passive transadmittances of a triode having a parallel wire grid of 0.3 mil wire wound at 1000 turns per inch.

Fig. 7—Passive transadmittances of a triode having a cross-lateral grid of 0.3 mil wire wound at 550 turns per inch.



Fig. 8—Phase of the transadmittance of the parallel-wire grid.

grid, 1553 triode having input and output spacings of 0.5 and 12 mils respectively. The capacitances averaged about 0.008 $\mu\mu$f, which would correspond to a value of $y_{12}$ of 0.0002 mho at 4060 megacycles. The latter is about 50 times lower than the measured 4060 megacycle value. Evidently, therefore, electromagnetic coupling plays a dominant role.

Reciprocity should give a reasonable idea of the accuracy of these measurements. Thus, for $S_1 = 0.001''$ and $S_2 = 0.012''$, one would expect the same $y_{12}$ as for the case where $S_1 = 0.012''$ and $S_2 = 0.001''$. An examination of the data will indicate that the reciprocal differences are of the order of 10% in some cases. These differences may be partly the result of the change in line cross section encountered in going from the input to the output. That is to say, the two cases being compared are not quite reciprocal in geometry.

Figure 8 shows the phase of $y_{12}$ for the parallel wire grid. Because of the low transmission through the grids there was not sufficient energy to determine the transfer phases with any very great accuracy, particularly for wide spacings in the case of the parallel wire grid and for all spacings in the case of the cross-lateral. Consequently, Fig. 8 shows only those results which are believed to be reasonably accurate.

The author wishes to acknowledge the contribution of Mr. F. A. Braun who ably assisted in this work.

# Communication Theory of Secrecy Systems*

## By C. E. SHANNON

### 1. Introduction and Summary

THE problems of cryptography and secrecy systems furnish an interesting application of communication theory.[1] In this paper a theory of secrecy systems is developed. The approach is on a theoretical level and is intended to complement the treatment found in standard works on cryptography.[2] There, a detailed study is made of the many standard types of codes and ciphers, and of the ways of breaking them. We will be more concerned with the general mathematical structure and properties of secrecy systems.

The treatment is limited in certain ways. First, there are three general types of secrecy system: (1) concealment systems, including such methods as invisible ink, concealing a message in an innocent text, or in a fake covering cryptogram, or other methods in which the existence of the message is concealed from the enemy; (2) privacy systems, for example speech inversion, in which special equipment is required to recover the message; (3) "true" secrecy systems where the meaning of the message is concealed by cipher, code, etc., although its existence is not hidden, and the enemy is assumed to have any special equipment necessary to intercept and record the transmitted signal. We consider only the third type—concealment systems are primarily a psychological problem, and privacy systems a technological one.

Secondly, the treatment is limited to the case of discrete information, where the message to be enciphered consists of a sequence of discrete symbols, each chosen from a finite set. These symbols may be letters in a language, words of a language, amplitude levels of a "quantized" speech or video signal, etc., but the main emphasis and thinking has been concerned with the case of letters.

The paper is divided into three parts. The main results will now be briefly summarized. The first part deals with the basic mathematical structure of secrecy systems. As in communication theory a language is considered to

be represented by a stochastic process which produces a discrete sequence of symbols in accordance with some system of probabilities. Associated with a language there is a certain parameter $D$ which we call the redundancy of the language. $D$ measures, in a sense, how much a text in the language can be reduced in length without losing any information. As a simple example, since $u$ always follows $q$ in English words, the $u$ may be omitted without loss. Considerable reductions are possible in English due to the statistical structure of the language, the high frequencies of certain letters or words, etc. Redundancy is of central importance in the study of secrecy systems.

A secrecy system is defined abstractly as a set of transformations of one space (the set of possible messages) into a second space (the set of possible cryptograms). Each particular transformation of the set corresponds to enciphering with a particular key. The transformations are supposed reversible (non-singular) so that unique deciphering is possible when the key is known.

Each key and therefore each transformation is assumed to have an *a priori* probability associated with it—the probability of choosing that key. Similarly each possible message is assumed to have an associated *a priori* probability, determined by the underlying stochastic process. These probabilities for the various keys and messages are actually the enemy cryptanalyst's *a priori* probabilities for the choices in question, and represent his *a priori* knowledge of the situation.

To use the system a key is first selected and sent to the receiving point. The choice of a key determines a particular transformation in the set forming the system. Then a message is selected and the particular transformation corresponding to the selected key applied to this message to produce a cryptogram. This cryptogram is transmitted to the receiving point by a channel and may be intercepted by the "enemy*." At the receiving end the inverse of the particular transformation is applied to the cryptogram to recover the original message.

If the enemy intercepts the cryptogram he can calculate from it the *a posteriori* probabilities of the various possible messages and keys which might have produced this cryptogram. This set of *a posteriori* probabilities constitutes his knowledge of the key and message after the interception. "Knowledge" is thus identified with a set of propositions having associated probabilities. The calculation of the *a posteriori* probabilities is the generalized problem of cryptanalysis.

As an example of these notions, in a simple substitution cipher with random key there are 26! transformations, corresponding to the 26! ways we

---

*The word "enemy," stemming from military applications, is commonly used in cryptographic work to denote anyone who may intercept a cryptogram.

can substitute for 26 different letters. These are all equally likely and each therefore has an *a priori* probability $1/26!$. If this is applied to "normal English" the cryptanalyst being assumed to have no knowledge of the message source other than that it is producing English text, the *a priori* probabilities of various messages of $N$ letters are merely their relative frequencies in normal English text.

If the enemy intercepts $N$ letters of cryptogram in this system his probabilities change. If $N$ is large enough (say 50 letters) there is usually a single message of *a posteriori* probability nearly unity, while all others have a total probability nearly zero. Thus there is an essentially unique "solution" to the cryptogram. For $N$ smaller (say $N = 15$) there will usually be many messages and keys of comparable probability, with no single one nearly unity. In this case there are multiple "solutions" to the cryptogram.

Considering a secrecy system to be represented in this way, as a set of transformations of one set of elements into another, there are two natural combining operations which produce a third system from two given systems. The first combining operation is called the product operation and corresponds to enciphering the message with the first secrecy system $R$ and enciphering the resulting cryptogram with the second system $S$, the keys for $R$ and $S$ being chosen independently. This total operation is a secrecy system whose transformations consist of all the products (in the usual sense of products of transformations) of transformations in $S$ with transformations in $R$. The probabilities are the products of the probabilities for the two transformations.

The second combining operation is "weighted addition."

$$T = pR + qS \qquad p + q = 1$$

It corresponds to making a preliminary choice as to whether system $R$ or $S$ is to be used with probabilities $p$ and $q$, respectively. When this is done $R$ or $S$ is used as originally defined.

It is shown that secrecy systems with these two combining operations form essentially a "linear associative algebra" with a unit element, an algebraic variety that has been extensively studied by mathematicians.

Among the many possible secrecy systems there is one type with many special properties. This type we call a "pure" system. A system is pure if all keys are equally likely and if for any three transformations $T_i$, $T_j$, $T_k$ in the set the product

$$T_i T_j^{-1} T_k$$

is also a transformation in the set. That is enciphering, deciphering, and enciphering with any three keys must be equivalent to enciphering with some key.

With a pure cipher it is shown that all keys are essentially equivalent—they all lead to the same set of *a posteriori* probabilities. Furthermore, when a given cryptogram is intercepted there is a set of messages that might have produced this cryptogram (a "residue class") and the *a posteriori* probabilities of messages in this class are proportional to the *a priori* probabilities. All the information the enemy has obtained by intercepting the cryptogram is a specification of the residue class. Many of the common ciphers are pure systems, including simple substitution with random key. In this case the residue class consists of all messages with the same pattern of letter repetitions as the intercepted cryptogram.

Two systems $R$ and $S$ are defined to be "similar" if there exists a fixed transformation $A$ with an inverse, $A^{-1}$, such that

$$R = AS.$$

If $R$ and $S$ are similar, a one-to-one correspondence between the resulting cryptograms can be set up leading to the same *a posteriori* probabilities. The two systems are crypt analytically the same.

The second part of the paper deals with the problem of "theoretical secrecy." How secure is a system against cryptanalysis when the enemy has unlimited time and manpower available for the analysis of intercepted cryptograms? The problem is closely related to questions of communication in the presence of noise, and the concepts of entropy and equivocation developed for the communication problem find a direct application in this part of cryptography.

"Perfect Secrecy" is defined by requiring of a system that after a cryptogram is intercepted by the enemy the *a posteriori* probabilities of this cryptogram representing various messages be identically the same as the *a priori* probabilities of the same messages before the interception. It is shown that perfect secrecy is possible but requires, if the number of messages is finite, the same number of possible keys. If the message is thought of as being constantly generated at a given "rate" $R$ (to be defined later), key must be generated at the same or a greater rate.

If a secrecy system with a finite key is used, and $N$ letters of cryptogram intercepted, there will be, for the enemy, a certain set of messages with certain probabilities, that this cryptogram could represent. As $N$ increases the field usually narrows down until eventually there is a unique "solution" to the cryptogram; one message with probability essentially unity while all others are practically zero. A quantity $H(N)$ is defined, called the equivocation, which measures in a statistical way how near the average cryptogram of $N$ letters is to a unique solution; that is, how uncertain the enemy is of the original message after intercepting a cryptogram of $N$ letters. Various properties of the equivocation are deduced—for example, the equivocation

of the key never increases with increasing $N$. This equivocation is a theoretical secrecy index—theoretical in that it allows the enemy unlimited time to analyse the cryptogram.

The function $H(N)$ for a certain idealized type of cipher called the random cipher is determined. With certain modifications this function can be applied to many cases of practical interest. This gives a way of calculating approximately how much intercepted material is required to obtain a solution to a secrecy system. It appears from this analysis that with ordinary languages and the usual types of ciphers (not codes) this "unicity distance" is approximately $H(K)/D$. Here $H(K)$ is a number measuring the "size" of the key space. If all keys are *a priori* equally likely $H(K)$ is the logarithm of the number of possible keys. $D$ is the redundancy of the language and measures the amount of "statistical constraint" imposed by the language. In simple substitution with random key $H(K)$ is $\log_{10} 26!$ or about 20 and $D$ (in decimal digits per letter) is about .7 for English. Thus unicity occurs at about 30 letters.

It is possible to construct secrecy systems with a finite key for certain "languages" in which the equivocation does not approach zero as $N \to \infty$. In this case, no matter how much material is intercepted, the enemy still does not obtain a unique solution to the cipher but is left with many alternatives, all of reasonable probability. Such systems we call *ideal* systems. It is possible in any language to approximate such behavior—i.e., to make the approach to zero of $H(N)$ recede out to arbitrarily large $N$. However, such systems have a number of drawbacks, such as complexity and sensitivity to errors in transmission of the cryptogram.

The third part of the paper is concerned with "practical secrecy." Two systems with the same key size may both be uniquely solvable when $N$ letters have been intercepted, but differ greatly in the amount of labor required to effect this solution. An analysis of the basic weaknesses of secrecy systems is made. This leads to methods for constructing systems which will require a large amount of work to solve. Finally, a certain incompatibility among the various desirable qualities of secrecy systems is discussed.

## PART I

## MATHEMATICAL STRUCTURE OF SECRECY SYSTEMS

### 2. Secrecy Systems

As a first step in the mathematical analysis of cryptography, it is necessary to idealize the situation suitably, and to define in a mathematically acceptable way what we shall mean by a secrecy system. A "schematic" diagram of a general secrecy system is shown in Fig. 1. At the transmitting

end there are two information sources—a message source and a key source. The key source produces a particular key from among those which are possible in the system. This key is transmitted by some means, supposedly not interceptible, for example by messenger, to the receiving end. The message source produces a message (the "clear") which is enciphered and the resulting cryptogram sent to the receiving end by a possibly interceptible means, for example radio. At the receiving end the cryptogram and key are combined in the decipherer to recover the message.



Fig. 1—Schematic of a general secrecy system.

Evidently the encipherer performs a functional operation. If $M$ is the message, $K$ the key, and $E$ the enciphered message, or cryptogram, we have

$$E = f(M, K)$$

that is $E$ is a function of $M$ and $K$. It is preferable to think of this, however, not as a function of two variables but as a (one parameter) family of operations or transformations, and to write it

$$E = T_i M.$$

The transformation $T_i$ applied to message $M$ produces cryptogram $E$. The index $i$ corresponds to the particular key being used.

We will assume, in general, that there are only a finite number of possible keys, and that each has an associated probability $p_i$. Thus the key source is represented by a statistical process or device which chooses one from the set of transformations $T_1, T_2, \cdots, T_m$ with the respective probabilities $p_1, p_2, \cdots, p_m$. Similarly we will generally assume a finite number of possible messages $M_1, M_2, \cdots, M_n$ with associated *a priori* probabilities $q_1, q_2, \cdots, q_n$. The possible messages, for example, might be the possible sequences of English letters all of length $N$, and the associated probabilities are then

the relative frequencies of occurrence of these sequences in normal English text.

At the receiving end it must be possible to recover $M$, knowing $E$ and $K$. Thus the transformations $T_i$ in the family must have unique inverses $T_i^{-1}$ such that $T_i T_i^{-1} = I$, the identity transformation. Thus:

$$M = T_i^{-1} E.$$

At any rate this inverse must exist uniquely for every $E$ which can be obtained from an $M$ with key $i$. Hence we arrive at the definition: A secrecy system is a family of uniquely reversible transformations $T_i$ of a set of possible mssages into a set of cryptograms, the transformation $T_i$ having an associated probability $p_i$. Conversely any set of entities of this type will be called a "secrecy system." The set of possible messages will be called, for convenience, the "message space" and the set of possible cryptograms the "cryptogram space."

Two secrecy systems will be the same if they consist of the same set of transformations $T_i$, with the same message and cryptogram space (range and domain) and the same probabilities for the keys.

A secrecy system can be visualized mechanically as a machine with one or more controls on it. A sequence of letters, the message, is fed into the input of the machine and a second series emerges at the output. The particular setting of the controls corresponds to the particular key being used. Some statistical method must be prescribed for choosing the key from all the possible ones.

To make the problem mathematically tractable we shall assume that *the enemy knows the system being used*. That is, he knows the family of transformations $T_i$, and the probabilities of choosing various keys. It might be objected that this assumption is unrealistic, in that the cryptanalyst often does not know what system was used or the probabilities in question. There are two answers to this objection:

1. The restriction is much weaker than appears at first, due to our broad definition of what constitutes a secrecy system. Suppose a cryptographer intercepts a message and does not know whether a substitution, transposition, or Vigenère type cipher was used. He can consider the message as being enciphered by a system in which part of the key is the specification of which of these types was used, the next part being the particular key for that type. These three different possibilities are assigned probabilities according to his best estimates of the *a priori* probabilities of the encipherer using the respective types of cipher.

2. The assumption is actually the one ordinarily used in cryptographic studies. It is pessimistic and hence safe, but in the long run realistic, since one must expect his system to be found out eventually. Thus,

even when an entirely new system is devised, so that the enemy cannot assign any *a priori* probability to it without discovering it himself, one must still live with the expectation of his eventual knowledge.

The situation is similar to that occurring in the theory of games[3] where it is assumed that the opponent "finds out" the strategy of play being used. In both cases the assumption serves to delineate sharply the opponent's knowledge.

A second possible objection to our definition of secrecy systems is that no account is taken of the common practice of inserting nulls in a message and the use of multiple substitutes. In such cases there is not a unique cryptogram for a given message and key, but the encipherer can choose at will from among a number of different cryptograms. This situation could be handled, but would only add complexity at the present stage, without substantially altering any of the basic results.

If the messages are produced by a Markoff process of the type described in ([1]) to represent an information source, the probabilities of various messages are determined by the structure of the Markoff process. For the present, however, we wish to take a more general view of the situation and regard the messages as merely an abstract set of entities with associated probabilities, not necessarily composed of a sequence of letters and not necessarily produced by a Markoff process.

It should be emphasized that throughout the paper a secrecy system means not one, but a set of many transformations. After the key is chosen only one of these transformations is used and one might be led from this to define a secrecy system as a single transformation on a language. The enemy, however, does not know what key was chosen and the "might have been" keys are as important for him as the actual one. Indeed it is only the existence of these other possibilities that gives the system any secrecy. Since the secrecy is our primary interest, we are forced to the rather elaborate concept of a secrecy system defined above. This type of situation, where possibilities are as important as actualities, occurs frequently in games of strategy. The course of a chess game is largely controlled by threats which are *not* carried out. Somewhat similar is the "virtual existence" of unrealized imputations in the theory of games.

It may be noted that a single operation on a language forms a degenerate type of secrecy system under our definition—a system with only one key of unit probability. Such a system has no secrecy—the cryptanalyst finds the message by applying the inverse of this transformation, the only one in the system, to the intercepted cryptogram. The decipherer and cryptanalyst in this case possess the same information. In general, the only difference between the decipherer's knowledge and the enemy cryptanalyst's knowledge

[3] See von Neumann and Morgenstern "The Theory of Games," Princeton 1947.

is that the decipherer knows the particular key being used, while the crypt-analyst knows only the *a priori* probabilities of the various keys in the set. The process of deciphering is that of applying the inverse of the particular transformation used in enciphering to the cryptogram. The process of crypt-analysis is that of attempting to determine the message (or the particular key) given only the cryptogram and the *a priori* probabilities of various keys and messages.

There are a number of difficult epistemological questions connected with the theory of secrecy, or in fact with any theory which involves questions of probability (particularly *a priori* probabilities, Bayes' theorem, etc.) when applied to a physical situation. Treated abstractly, probability theory can be put on a rigorous logical basis with the modern measure theory ap-proach.[4,5] As applied to a physical situation, however, especially when "subjective" probabilities and unrepeatable experiments are concerned, there are many questions of logical validity. For example, in the approach to secrecy made here, *a priori* probabilities of various keys and messages are assumed known by the enemy cryptographer—how can one determine operationally if his estimates are correct, on the basis of his knowledge of the situation?

One can construct artificial cryptographic situations of the "urn and die" type in which the *a priori* probabilities have a definite unambiguous meaning and the idealization used here is certainly appropriate. In other situations that one can imagine, for example an intercepted communication between Martian invaders, the *a priori* probabilities would probably be so uncertain as to be devoid of significance. Most practical cryptographic situations lie somewhere between these limits. A cryptanalyst might be willing to classify the possible messages into the categories "reasonable," "possible but un-likely" and "unreasonable," but feel that finer subdivision was meaningless.

Fortunately, in practical situations, only extreme errors in *a priori* prob-abilities of keys and messages cause significant errors in the important parameters. This is because of the exponential behavior of the number of messages and cryptograms, and the logarithmic measures employed.

## 3. Representation of Systems

A secrecy system as defined above can be represented in various ways. One which is convenient for illustrative purposes is a line diagram, as in Figs. 2 and 4. The possible messages are represented by points at the left and the possible cryptograms by points at the right. If a certain key, say key 1, transforms message $M_2$ into cryptogram $E_4$ then $M_2$ and $E_4$ are connected

---

[4] See J. L. Doob, "Probability as Measure," *Annals of Math. Stat.*, v. 12, 1941, pp. 206–214.
[5] A. Kolmogoroff, "Grundbegriffe der Wahrscheinlichkeits rechnung," *Ergebnisse der Mathematic*, v. 2, No. 3 (Berlin 1933).

by a line labeled 1, etc. From each possible message there must be exactly one line emerging for each different key. If the same is true for each cryptogram, we will say that the system is *closed*.

A more common way of describing a system is by stating the operation one performs on the message for an arbitrary key to obtain the cryptogram. Similarly, one defines implicitly the probabilities for various keys by describing how a key is chosen or what we know of the enemy's habits of key choice. The probabilities for messages are implicitly determined by stating our *a priori* knowledge of the enemy's language habits, the tactical situation (which will influence the probable content of the message) and any special information we may have regarding the cryptogram.



CLOSED SYSTEM          NOT CLOSED
Fig. 2—Line drawings for simple systems.

### 4. SOME EXAMPLES OF SECRECY SYSTEMS

In this section a number of examples of ciphers will be given. These will often be referred to in the remainder of the paper for illustrative purposes.

1. *Simple Substitution Cipher.*

In this cipher each letter of the message is replaced by a fixed substitute, usually also a letter. Thus the message,

$$M = m_1 m_2 m_3 m_4 \cdots$$

where $m_1$, $m_2$, $\cdots$ are the successive letters becomes:

$$E = e_1 e_2 e_3 e_4 \cdots$$
$$= f(m_1)f(m_2)f(m_3)f(m_4) \cdots$$

where the function $f(m)$ is a function with an inverse. The key is a permutation of the alphabet (when the substitutes are letters) e.g. $X\ G\ U\ A\ C\ D$ $T\ B\ F\ H\ R\ S\ L\ M\ Q\ V\ V\ Y\ Z\ W\ I\ E\ J\ O\ K\ N\ P$. The first letter $X$ is the substitute for $A$, $G$ is the substitute for $B$, etc.

## 2. *Transposition (Fixed Period d).*

The message is divided into groups of length $d$ and a permutation applied to the first group, the same permutation to the second group, etc. The permutation is the key and can be represented by a permutation of the first $d$ integers. Thus, for $d = 5$, we might have 2 3 1 5 4 as the permutation. This means that:

$$m_1\ m_2\ m_3\ m_4\ m_5\ m_6\ m_7\ m_8\ m_9\ m_{10}\ \cdots\ \text{becomes}$$
$$m_2\ m_3\ m_1\ m_5\ m_4\ m_7\ m_8\ m_6\ m_{10}\ m_9\ \cdots\ .$$

Sequential application of two or more transpositions will be called compound transposition. If the periods are $d_1$, $d_2$, $\cdots$, $d_s$ it is clear that the result is a transposition of period $d$, where $d$ is the least common multiple of $d_1$, $d_2$, $\cdots$, $d_s$.

## 3. *Vigenère, and Variations.*

In the Vigenère cipher the key consists of a series of $d$ letters. These are written repeatedly below the message and the two added modulo 26 (considering the alphabet numbered from $A = 0$ to $Z = 25$. Thus

$$e_i = m_i + k_i \,(\text{mod } 26)$$

where $k_i$ is of period $d$ in the index $i$. For example, with the key $G\,A\,H$**,** we obtain

| | |
|---|---|
| message | $N\ O\ W\ I\ S\ T\ H\ E \cdots$ |
| repeated key | $G\ A\ H\ G\ A\ H\ G\ A \cdots$ |
| cryptogram | $T\ O\ D\ O\ S\ A\ N\ E \cdots$ |

The Vigenère of period 1 is called the Caesar cipher. It is a simple substitution in which each letter of $M$ is advanced a fixed amount in the alphabet. This amount is the key, which may be any number from 0 to 25. The so-called Beaufort and Variant Beaufort are similar to the Vigenère, and encipher by the equations

$$e_i = k_i - m_i \,(\text{mod } 26)$$

and

$$e_i = m_i - k_i \,(\text{mod } 26)$$

respectively. The Beaufort of period one is called the reversed Caesar cipher.

The application of two or more Vigenères in sequence will be called the compound Vigenère. It has the equation

$$e_i = m_i + k_i + l_i + \cdots + s_i \,(\text{mod } 26)$$

where $k_i$, $l_i$, $\cdots$, $s_i$ in general have different periods. The period of their sum,

$$k_i + l_i + \cdots + s_i$$

as in compound transposition, is the least common multiple of the individual periods.

When the Vigenère is used with an unlimited key, never repeating, we have the Vernam system,[6] with

$$e_i = m_i + k_i \ (\text{mod } 26)$$

the $k_i$ being chosen at random and independently among 0, 1, $\cdots$, 25. If the key is a meaningful text we have the "running key" cipher.

## 4. *Digram, Trigram, and N-gram substitution.*

Rather than substitute for letters one can substitute for digrams, trigrams, etc. General digram substitution requires a key consisting of a permutation of the $26^2$ digrams. It can be represented by a table in which the row corresponds to the first letter of the digram and the column to the second letter, entries in the table being the substitutes (usually also digrams).

## 5. *Single Mixed Alphabet Vigenère.*

This is a simple substitution followed by a Vigenère.

$$e_i = f(m_i) + k_i$$
$$m_i = f^{-1}(e_i - k_i)$$

The "inverse" of this system is a Vigenère followed by simple substitution

$$e_i = g(m_i + k_i)$$
$$m_i = g^{-1}(e_i) - k_i$$

## 6. *Matrix System.*[7]

One method of $n$-gram substitution is to operate on successive $n$-grams with a matrix having an inverse. The letters are assumed numbered from 0 to 25, making them elements of an algebraic ring. From the $n$-gram $m_1 \, m_2 \cdots m_n$ of message, the matrix $a_{ij}$ gives an $n$-gram of cryptogram

$$e_i = \sum_{j=1}^{n} a_{ij} m_j \qquad i = 1, \cdots, n$$

---

[6] G. S. Vernam, "Cipher Printing Telegraph Systems for Secret Wire and Radio Telegraphic Communications," *Journal American Institute of Electrical Engineers*, v. XLV, pp. 109–115, 1926.

[7] See L. S. Hill, "Cryptography in an Algebraic Alphabet," *American Math. Monthly*, v. 36, No. 6, 1, 1929, pp. 306–312; also "Concerning Certain Linear Transformation Apparatus of Cryptography," v. 38, No. 3, 1931, pp. 135–154.

The matrix $a_{ij}$ is the key, and deciphering is performed with the inverse matrix. The inverse matrix will exist if and only if the determinant $|\ a_{ij}\ |$ has an inverse element in the ring.

### 7. *The Playfair Cipher.*

This is a particular type of digram substitution governed by a mixed 25 letter alphabet written in a 5 x 5 square. (The letter $J$ is often dropped in cryptographic work—it is very infrequent, and when it occurs can be replaced by $I$.) Suppose the key square is as shown below:

$$L \quad Z \quad Q \quad C \quad P$$
$$A \quad G \quad N \quad O \quad U$$
$$R \quad D \quad M \quad I \quad F$$
$$K \quad Y \quad H \quad V \quad S$$
$$X \quad B \quad T \quad E \quad W$$

The substitute for a digram $AC$, for example, is the pair of letters at the other corners of the rectangle defined by $A$ and $C$, i.e., $LO$, the $L$ taken first since it is above $A$. If the digram letters are on a horizontal line as $RI$, one uses the letters to their right $DF$; $RF$ becomes $DR$. If the letters are on a vertical line, the letters below them are used. Thus $PS$ becomes $UW$. If the letters are the same nulls may be used to separate them or one may be omitted, etc.

### 8. *Multiple Mixed Alphabet Substitution.*

In this cipher there are a set of $d$ simple substitutions which are used in sequence. If the period $d$ is four

$$m_1 \ m_2 \ m_3 \ m_4 \ m_5 \ m_6 \ \cdots$$

becomes

$$f_1(m_1) \ f_2(m_2) \ f_3(m_3) \ f_4(m_4) \ f_1(m_5) \ f_2(m_6) \ \cdots$$

### 9. *Autokey Cipher.*

A Vigenère type system in which either the message itself or the resulting cryptogram is used for the "key" is called an autokey cipher. The encipherment is started with a "priming key" (which is the entire key in our sense) and continued with the message or cryptogram displaced by the length of the priming key as indicated below, where the priming key is COMET. The message used as "key":

| Message | S E N D S U P P L I E S $\cdots$ |
|---|---|
| Key | C O M E T S E N D S U P $\cdots$ |
| Cryptogram | U S Z H L M T C O A Y H |

The cryptogram used as "key":[8]

| Message | S E N D S U P P L I E S ··· |
|---|---|
| Key | C O M E T U S Z H L O H ··· |
| Cryptogram | U S Z H L O H O S T S ··· |

## 10. *Fractional Ciphers.*

In these, each letter is first enciphered into two or more letters or numbers and these symbols are somehow mixed (e.g. by transposition). The result may then be retranslated into the original alphabet. Thus, using a mixed 25-letter alphabet for the key, we may translate letters into two-digit quinary numbers by the table:

$$
\begin{array}{c|ccccc}
 & 0 & 1 & 2 & 3 & 4 \\
0 & L & Z & Q & C & P \\
1 & A & G & N & O & U \\
2 & R & D & M & I & F \\
3 & K & Y & H & V & S \\
4 & X & B & T & E & W \\
\end{array}
$$

Thus $B$ becomes 41. After the resulting series of numbers is transposed in some way they are taken in pairs and translated back into letters.

## 11. *Codes.*

In codes words (or sometimes syllables) are replaced by substitute letter groups. Sometimes a cipher of one kind or another is applied to the result.

## 5. VALUATIONS OF SECRECY SYSTEMS

There are a number of different criteria that should be applied in estimating the value of a proposed secrecy system. The most important of these are:

## 1. *Amount of Secrecy.*

There are some systems that are perfect—the enemy is no better off after intercepting any amount of material than before. Other systems, although giving him some information, do not yield a unique "solution" to intercepted cryptograms. Among the uniquely solvable systems, there are wide variations in the amount of labor required to effect this solution and in the amount of material that must be intercepted to make the solution unique.

---

[8] This system is trivial from the secrecy standpoint since, with the exception of the first $d$ letters, the enemy is in possession of the entire "key."

## 2. *Size of Key.*

The key must be transmitted by non-interceptible means from transmitting to receiving points. Sometimes it must be memorized. It is therefore desirable to have the key as small as possible.

## 3. *Complexity of Enciphering and Deciphering Operations.*

Enciphering and deciphering should, of course, be as simple as possible. If they are done manually, complexity leads to loss of time, errors, etc. If done mechanically, complexity leads to large expensive machines.

## 4. *Propagation of Errors.*

In certain types of ciphers an error of one letter in enciphering or transmission leads to a large number of errors in the deciphered text. The errors are spread out by the deciphering operation, causing the loss of much information and frequent need for repetition of the cryptogram. It is naturally desirable to minimize this error expansion.

## 5. *Expansion of Message.*

In some types of secrecy systems the size of the message is increased by the enciphering process. This undesirable effect may be seen in systems where one attempts to swamp out message statistics by the addition of many nulls, or where multiple substitutes are used. It also occurs in many "concealment" types of systems (which are not usually secrecy systems in the sense of our definition).

## 6. THE ALGEBRA OF SECRECY SYSTEMS

If we have two secrecy systems $T$ and $R$ we can often combine them in various ways to form a new secrecy system $S$. If $T$ and $R$ have the same domain (message space) we may form a kind of "weighted sum,"

$$S = pT + qR$$

where $p + q = 1$. This operation consists of first making a preliminary choice with probabilities $p$ and $q$ determining which of $T$ and $R$ is used. This choice is part of the key of $S$. After this is determined $T$ or $R$ is used as originally defined. The total key of $S$ must specify which of $T$ and $R$ is used and which key of $T$ (or $R$) is used.

If $T$ consists of the transformations $T_1$, $\cdots$, $T_m$ with probabilities $p_1$, $\cdots$, $p_m$ and $R$ consists of $R_1$, $\cdots$, $R_k$ with probabilities $q_1$, $\cdots$, $q_k$ then $S = pT + qR$ consists of the transformations $T_1$, $T_2$, $\cdots$, $T_m$, $R_1$, $\cdots$, $R_k$ with probabilities $pp_1$, $pp_2$, $\cdots$, $pp_m$, $qq_1$, $qq_2$, $\cdots$, $qq_k$ respectively.

More generally we can form the sum of a number of systems.

$$S = p_1T + p_2R + \cdots + p_mU \qquad \sum p_i = 1$$

We note that any system $T$ can be written as a sum of fixed operations

$$T = p_1T_1 + p_2T_2 + \cdots + p_mT_m$$

$T_i$ being a definite enciphering operation of $T$ corresponding to key choice $i$, which has probability $p_i$.

A second way of combining two secrecy systems is by taking the "product," shown schematically in Fig. 3. Suppose $T$ and $R$ are two systems and the domain (language space) of $R$ can be identified with the range (cryptogram space) of $T$. Then we can apply first $T$ to our language and then $R$



Fig. 3—Product of two systems $S = RT$.

to the result of this enciphering process. This gives a resultant operation $S$ which we write as a product

$$S = RT$$

The key for $S$ consists of both keys of $T$ and $R$ which are assumed chosen according to their original probabilities and independently. Thus, if the $m$ keys of $T$ are chosen with probabilities

$$p_1 \ p_2 \ \cdots \ p_m$$

and the $n$ keys of $R$ have probabilities

$$p_1' \ p_2' \ \cdots \ p_n' \, ,$$

then $S$ has at most $mn$ keys with probabilities $p_ip_j'$. In many cases some of the product transformaions $R_iT_j$ will be the same and can be grouped together, adding their probabilities.

Product encipherment is often used; for example, one follows a substitution by a transposition or a transposition by a Vigenère, or applies a code to the text and enciphers the result by substitution, transposition, fractionation, etc.

It may be noted that multiplication is not in general commutative, (we do not always have $RS = SR$), although in special cases, such as substitution and transposition, it is. Since it represents an operation it is definitionally associative. That is $R(ST) = (RS)T = RST$. Furthermore we have the laws

$$p(p'T + q'R) + qS = pp'T + pq'R + qS$$

(weighted associative law for addition)

$$T(pR + qS) = pTR + qTS$$
$$(pR + qS)T = pRT + qST$$

(right and left hand distributive laws)
and

$$p_1T + p_2T + p_3R = (p_1 + p_2)T + p_3R$$

It should be emphasized that these combining operations of addition and multiplication apply to secrecy systems as a whole. The product of two systems $TR$ should not be confused with the product of the transformations in the systems $T_iR_j$, which also appears often in this work. The former $TR$ is a secrecy system, i.e., a set of transformations with associated probabilities; the latter is a particular transformation. Further the sum of two systems $pR + qT$ is a system—the sum of two transformations is not defined. The systems $T$ and $R$ may commute without the individual $T_i$ and $R_j$ commuting, e.g., if $R$ is a Beaufort system of a given period, all keys equally likely,

$$R_iR_j \neq R_jR_i$$

in general, but of course $RR$ does not depend on its order; actually

$$RR = V$$

the Vigenère of the same period with random key. On the other hand, if the individual $T_i$ and $R_j$ of two systems $T$ and $R$ commute, then the systems commute.

A system whose $M$ and $E$ spaces can be identified, a very common case as when letter sequences are transformed into letter sequences, may be termed *endomorphic*. An endomorphic system $T$ may be raised to a power $T^n$.

A secrecy system $T$ whose product with itself is equal to $T$, i.e., for which

$$TT = T,$$

will be called idempotent. For example, simple substitution, transposition of period $p$, Vigenère of period $p$ (all with each key equally likely) are idempotent.

The set of all endomorphic secrecy systems defined in a fixed message space constitutes an "algebraic variety," that is, a kind of algebra, using the operations of addition and multiplication. In fact, the properties of addition and multiplication which we have discussed may be summarized as follows:

*The set of endomorphic ciphers with the same message space and the two combining operations of weighted addition and multiplication form a linear associative algebra with a unit element, apart from the fact that the coefficients in a weighted addition must be non-negative and sum to unity.*

The combining operations give us ways of constructing many new types of secrecy systems from certain ones, such as the examples given. We may also use them to describe the situation facing a cryptanalyst when attempting to solve a cryptogram of unknown type. He is, in fact, solving a secrecy system of the type

$$T = p_1 A + p_2 B + \cdots + p_r S + p'X \qquad \sum p = 1$$

where the $A, B, \cdots, S$ are known types of ciphers, with the $p_i$ their *a priori* probabilities in this situation, and $p'X$ corresponds to the possibility of a completely new unknown type of cipher.

## 7. PURE AND MIXED CIPHERS

Certain types of ciphers, such as the simple substitution, the transposition of a given period, the Vigenère of a given period, the mixed alphabet Vigenère, etc. (all with each key equally likely) have a certain homogeneity with respect to key. Whatever the key, the enciphering, deciphering and decrypting processes are essentially the same. This may be contrasted with the cipher

$$pS + qT$$

where $S$ is a simple substitution and $T$ a transposition of a given period. In this case the entire system changes for enciphering, deciphering and decryptment, depending on whether the substitution or transposition is used.

The cause of the homogeneity in these systems stems from the group property—we notice that, in the above examples of homogeneous ciphers, the product $T_i T_j$ of any two transformations in the set is equal to a third transformation $T_k$ in the set. On the other hand $T_i S_j$ does not equal any transformation in the cipher

$$pS + qT$$

which contains only substitutions and transpositions, no products.

We might define a "pure" cipher, then, as one whose $T_i$ form a group. This, however, would be too restrictive since it requires that the $E$ space

be the same as the $M$ space, i.e. that the system be endomorphic. The fractional transposition is as homogeneous as the ordinary transposition without being endomorphic. The proper definition is the following: A cipher $T$ is *pure* if for every $T_i$, $T_j$, $T_k$ there is a $T_s$ such that

$$T_i T_j^{-1} T_k = T_s$$

and every key is equally likely. Otherwise the cipher is mixed. The systems of Fig. 2 are mixed. Fig. 4 is pure if all keys are equally likely.

*Theorem 1: In a pure cipher the operations $T_i^{-1} T_j$ which transform the message space into itself form a group whose order is $m$, the number of different keys.*

For

$$T_j^{-1} T_k T_k^{-1} T_j = I$$

so that each element has an inverse. The associative law is true since these are operations, and the group property follows from

$$T_i^{-1} T_j T_k^{-1} T_l = T_s^{-1} T_k T_k^{-1} T_l = T_s^{-1} T_l$$

using our assumption that $T_i^{-1} T_j = T_s^{-1} T_k$ for some $s$.

The operation $T_i^{-1} T_j$ means, of course, enciphering the message with key $j$ and then deciphering with key $i$ which brings us back to the message space. If $T$ is endomorphic, i.e. the $T_i$ themselves transform the space $\Omega_M$ into itself (as is the case with most ciphers, where both the message space and the cryptogram space consist of sequences of letters), and the $T_i$ are a group and equally likely, then $T$ is pure, since

$$T_i T_j^{-1} T_k = T_i T_r = T_s .$$

*Theorem 2: The product of two pure ciphers which commute is pure.*

For if $T$ and $R$ commute $T_i R_j = R_l T_m$ for every $i, j$ with suitable $l, m$, and

$$\begin{aligned} T_i R_j (T_k R_l)^{-1} T_m R_n &= T_i R_j R_l^{-1} T_k^{-1} T_m R_n \\ &= R_u R_v^{-1} R_w T_r T_s^{-1} T_t \\ &= R_h T_g . \end{aligned}$$

The commutation condition is not necessary, however, for the product to be a pure cipher.

A system with only one key, i.e., a single definite operation $T_1$, is pure since the only choice of indices is

$$T_1 T_1^{-1} T_1 = T_1 .$$

Thus the expansion of a general cipher into a sum of such simple transformations also exhibits it as a sum of pure ciphers.

An examination of the example of a pure cipher shown in Fig. 4 discloses

certain properties. The messages fall into certain subsets which we will call *residue classes*, and the possible cryptograms are divided into corresponding residue classes. There is at least one line from each message in a class to each cryptogram in the corresponding class, and no line between classes which do not correspond. The number of messages in a class is a divisor of the total number of keys. The number of lines "in parallel" from a message $M$ to a cryptogram in the corresponding class is equal to the number of keys divided by the number of messages in the class containing the message (or cryptogram). It is shown in the appendix that these hold in general for pure ciphers. Summarized formally, we have:



**PURE SYSTEM**
Fig. 4—Pure system.

*Theorem 3: In a pure system the messages can be divided into a set of "residue classes" $C_1, C_2, \cdots, C_s$ and the cryptograms into a corresponding set of residue classes $C'_1, C'_2, \cdots, C'_s$ with the following properties:*

(1) *The message residue classes are mutually exclusive and collectively contain all possible messages. Similarly for the cryptogram residue classes.*

(2) *Enciphering any message in $C_i$ with any key produces a cryptogram in $C'_i$. Deciphering any cryptogram in $C'_i$ with any key leads to a message in $C_i$.*

(3) *The number of messages in $C_i$, say $\varphi_i$, is equal to the number of cryptograms in $C'_i$ and is a divisor of $k$ the number of keys.*

(4) *Each message in $C_i$ can be enciphered into each cryptogram in $C'_i$ by exactly $k/\varphi_i$ different keys. Similarly for decipherment.*

The importance of the concept of a pure cipher (and the reason for the name) lies in the fact that in a pure cipher all keys are essentially the same. Whatever key is used for a particular message, the *a posteriori* probabilities of all messages are identical. To see this, note that two different keys applied to the same message lead to two cryptograms in the same residue class, say $C'_i$. The two cryptograms therefore could each be deciphered by $\dfrac{k}{\varphi_i}$ keys into each message in $C_i$ and into no other possible messages. All keys being equally likely the *a posteriori* probabilities of various messages are thus

$$P_E(M) = \frac{P(M)P_M(E)}{P(E)} = \frac{P(M)P_M(E)}{\Sigma_M P(M)P_M(E)} = \frac{P(M)}{P(C_i)}$$

where $M$ is in $C_i$, $E$ is in $C'_i$ and the sum is over all messages in $C_i$. If $E$ and $M$ are not in corresponding residue classes, $P_E(M) = 0$. Similarly it can be shown that the *a posteriori* probabilities of the different keys are the same in value but these values are associated with different keys when a different key is used. The same set of values of $P_E(K)$ have undergone a permutation among the keys. Thus we have the result

*Theorem 4: In a pure system the* a posteriori *probabilities of various messages $P_E(M)$ are independent of the key that is chosen. The* a posteriori *probabilities of the keys $P_E(K)$ are the same in value but undergo a permutation with a different key choice.*

Roughly we may say that any key choice leads to the same cryptanalytic problem in a pure cipher. Since the different keys all result in cryptograms in the same residue class this means that all cryptograms in the same residue class are cryptanalytically equivalent—they lead to the same *a posteriori* probabilities of messages and, apart from a permutation, the same probabilities of keys.

As an example of this, simple substitution with all keys equally likely is a pure cipher. The residue class corresponding to a given cryptogram $E$ is the set of all cryptograms that may be obtained from $E$ by operations $T_j T_k^{-1} E$. In this case $T_j T_k^{-1}$ is itself a substitution and hence any substitution on $E$ gives another member of the same residue class. Thus, if the cryptogram is

$$E = X\ C\ P\ P\ G\ C\ F\ Q,$$

then

$$E_1 = R\ D\ H\ H\ G\ D\ S\ N$$
$$E_2 = A\ B\ C\ C\ D\ B\ E\ F$$

etc. are in the same residue class. It is obvious in this case that these crypto-grams are essentially equivalent. All that is of importance in a simple sub-stitution with random key is the *pattern* of letter repetitions, the actual letters being dummy variables. Indeed we might dispense with them en-tirely, indicating the pattern of repetitions in $E$ as follows:

This notation describes the residue class but eliminates all information as to the specific member of the class. Thus it leaves precisely that information which is cryptanalytically pertinent. This is related to one method of attack-ing simple substitution ciphers—the method of pattern words.

In the Caesar type cipher only the first differences mod 26 of the crypto-gram are significant. Two cryptograms with the same $\Delta e_i$ are in the same residue class. One breaks this cipher by the simple process of writing down the 26 members of the message residue class and picking out the one which makes sense.

The Vigenère of period $d$ with random key is another example of a pure cipher. Here the message residue class consists of all sequences with the same first differences as the cryptogram, for letters separated by distance $d$. For $d = 3$ the residue class is defined by

$$m_1 - m_4 = e_1 - e_4$$
$$m_2 - m_5 = e_2 - e_5$$
$$m_3 - m_6 = e_3 - e_6$$
$$m_4 - m_7 = e_4 - e_7$$
$$\cdot$$
$$\cdot$$
$$\cdot$$

where $E = e_1, e_2, \cdots$ is the cryptogram and $m_1, m_2, \cdots$ is any $M$ in the corresponding residue class.

In the transposition cipher of period $d$ with random key, the residue class consists of all arrangements of the $e_i$ in which no $e_i$ is moved out of its block of length $d$, and any two $e_i$ at a distance $d$ remain at this distance. This is used in breaking these ciphers as follows: The cryptogram is written in successive blocks of length $d$, one under another as below ($d = 5$):

$$
\begin{array}{ccccc}
e_1 & e_2 & e_3 & e_4 & e_5 \\
e_6 & e_7 & e_8 & e_9 & e_{10} \\
e_{11} & e_{12} & \cdot & \cdot & \cdot \\
\cdot & \cdot & \cdot & \cdot & \cdot
\end{array}
$$

The columns are then cut apart and rearranged to make meaningful text. When the columns are cut apart, the only information remaining is the residue class of the cryptogram.

*Theorem 5: If $T$ is pure then $T_i T_j^{-1} T = T$ where $T_i T_j$ are any two transformations of $T$. Conversely if this is true for any $T_i T_j$ in a system $T$ then $T$ is pure.*

The first part of this theorem is obvious from the definition of a pure system. To prove the second part we note first that, if $T_i T_j^{-1} T = T$, then $T_i T_j^{-1} T_s$ is a transformation of $T$. It remains to show that all keys are equiprobable. We have $T = \sum_s p_s T_s$ and

$$\sum_s p_s T_i T_j^{-1} T_s = \sum_s p_s T_s .$$

The term in the left hand sum with $s = j$ yields $p_j T_i$. The only term in $T_i$ on the right is $p_i T_i$. Since all coefficients are nonnegative it follows that

$$p_j \leq p_i .$$

The same argument holds with $i$ and $j$ interchanged and consequently

$$p_j = p_i$$

and $T$ is pure. Thus the condition that $T_i T_j^{-1} T = T$ might be used as an alternative definition of a pure system.

## 8. SIMILAR SYSTEMS

Two secrecy systems $R$ and $S$ will be said to be *similar* if there exists a transformation $A$ having an inverse $A^{-1}$ such that

$$R = AS$$

This means that enciphering with $R$ is the same as enciphering with $S$ and then operating on the result with the transformation $A$. If we write $R \approx S$ to mean $R$ is similar to $S$ then it is clear that $R \approx S$ implies $S \approx R$. Also $R \approx S$ and $S \approx T$ imply $R \approx T$ and finally $R \approx R$. These are summarized by saying that similarity is an equivalence relation.

The cryptographic significance of similarity is that if $R \approx S$ then $R$ and $S$ are equivalent from the cryptanalytic point of view. Indeed if a cryptanalyst intercepts a cryptogram in system $S$ he can transform it to one in system $R$ by merely applying the transformation $A$ to it. A cryptogram in system $R$ is transformed to one in $S$ by applying $A^{-1}$. If $R$ and $S$ are applied to the same language or message space, there is a one-to-one correspondence between the resulting cryptograms. Corresponding cryptograms give the same distribution of *a posteriori* probabilities for all messages.

If one has a method of breaking the system $R$ then any system $S$ similar

to $R$ can be broken by reducing to $R$ through application of the operation $A$. This is a device that is frequently used in practical cryptanalysis.

As a trivial example, simple substitution where the substitutes are not letters but arbitrary symbols is similar to simple substitution using letter substitutes. A second example is the Caesar and the reversed Caesar type ciphers. The latter is sometimes broken by first transforming into a Caesar type. This can be done by reversing the alphabet in the cryptogram. The Vigenère, Beaufort and Variant Beaufort are all similar, when the key is random. The "autokey" cipher (with the message used as "key") primed with the key $K_1 K_2 \cdots K_d$ is similar to a Vigenère type with the key alternately added and subtracted Mod 26. The transformation $A$ in this case is that of "deciphering" the autokey with a series of $d$ $A$'s for the priming key.

## PART II

## THEORETICAL SECRECY

### 9. Introduction

We now consider problems connected with the "theoretical secrecy" of a system. How immune is a system to cryptanalysis when the cryptanalyst has unlimited time and manpower available for the analysis of cryptograms? Does a cryptogram *have* a unique solution (even though it may require an impractical amount of work to find it) and if not how many reasonable solutions does it have? How much text in a given system must be intercepted before the solution becomes unique? Are there systems which never become unique in solution no matter how much enciphered text is intercepted? Are there systems for which no information whatever is given to the enemy no matter how much text is intercepted? In the analysis of these problems the concepts of entropy, redundancy and the like developed in "A Mathematical Theory of Communication" (hereafter referred to as MTC) will find a wide application.

### 10. Perfect Secrecy

Let us suppose the possible messages are finite in number $M_1, \cdots, M_n$ and have *a priori* probabilities $P(M_1), \cdots, P(M_n)$, and that these are enciphered into the possible cryptograms $E_1, \cdots, E_m$ by

$$E = T_i M.$$

The cryptanalyst intercepts a particular $E$ and can then calculate, in principle at least, the *a posteriori* probabilities for the various messages, $P_E(M)$. It is natural to define *perfect secrecy* by the condition that, for all $E$ the *a posteriori* probabilities are equal to the *a priori* probabilities independently of the values of these. In this case, intercepting the message has

given the cryptanalyst no information.[9] Any action of his which depends on the information contained in the cryptogram cannot be altered, for all of his probabilities as to what the cryptogram contains remain unchanged. On the other hand, if the condition is *not* satisfied there will exist situations in which the enemy has certain *a priori* probabilities, and certain key and message choices may occur for which the enemy's probabilities do change. This in turn may affect his actions and thus perfect secrecy has not been obtained. Hence the definition given is necessarily required by our intuitive ideas of what perfect secrecy should mean.

A necessary and sufficient condition for perfect secrecy can be found as follows: We have by Bayes' theorem

$$P_E(M) = \frac{P(M)P_M(E)}{P(E)}$$

in which:

$P(M)$ = *a priori* probability of message $M$.

$P_M(E)$ = conditional probability of cryptogram $E$ if message $M$ is chosen, i.e. the sum of the probabilities of all keys which produce cryptogram $E$ from message $M$.

$P(E)$ = probability of obtaining cryptogram $E$ from any cause.

$P_E(M)$ = *a posteriori* probability of message $M$ if cryptogram $E$ is intercepted.

For perfect secrecy $P_E(M)$ must equal $P(M)$ for all $E$ and all $M$. Hence either $P(M) = 0$, a solution that must be excluded since we demand the equality independent of the values of $P(M)$, or

$$P_M(E) = P(E)$$

for every $M$ and $E$. Conversely if $P_M(E) = P(E)$ then

$$P_E(M) = P(M)$$

and we have perfect secrecy. Thus we have the result:

*Theorem 6: A necessary and sufficient condition for perfect secrecy is that*

$$P_M(E) = P(E)$$

*for all M and E. That is, $P_M(E)$ must be independent of M.*

Stated another way, the total probability of all keys that transform $M_i$

---

[9] A purist might object that the enemy has obtained some information in that he knows a message was sent. This may be answered by having among the messages a "blank" corresponding to "no message." If no message is originated the blank is enciphered and sent as a cryptogram. Then even this modicum of remaining information is eliminated.

into a given cryptogram $E$ is equal to that of all keys transforming $M_j$ into the same $E$, for all $M_i$, $M_j$ and $E$.

Now there must be as many $E$'s as there are $M$'s since, for a fixed $i$, $T_i$ gives a one-to-one correspondence between all the $M$'s and some of the $E$'s. For perfect secrecy $P_M(E) = P(E) \neq 0$ for any of these $E$'s and any $M$. Hence there is at least one key transforming any $M$ into any of these $E$'s. But all the keys from a fixed $M$ to different $E$'s must be different, and therefore *the number of different keys is at least as great as the number of $M$'s.* It is possible to obtain perfect secrecy with only this number of keys, as



Fig. 5—Perfect system.

one shows by the following example: Let the $M_i$ be numbered 1 to $n$ and the $E_i$ the same, and using $n$ keys let

$$T_i M_j = E_s$$

where $s = i + j$ (Mod $n$). In this case we see that $P_E(M) = \dfrac{1}{n} = P(E)$ and we have perfect secrecy. An example is shown in Fig. 5 with $s = i + j - 1$ (Mod 5).

Perfect systems in which the number of cryptograms, the number of messages, and the number of keys are all equal are characterized by the properties that (1) each $M$ is connected to each $E$ by exactly one line, (2) all keys are equally likely. Thus the matrix representation of the system is a "Latin square."

In MTC it was shown that information may be conveniently measured by means of entropy. If we have a set of possibilities with probabilities $p_1$, $p_2$, $\cdots$, $p_n$, the entropy $H$ is given by:

$$H = -\sum p_i \log p_i.$$

In a secrecy system there are two statistical choices involved, that of the message and of the key. We may measure the amount of information produced when a message is chosen by $H(M)$:

$$H(M) = -\sum P(M) \log P(M),$$

the summation being over all possible messages. Similarly, there is an uncertainty associated with the choice of key given by:

$$H(K) = -\sum P(K) \log P(K).$$

In perfect systems of the type described above, the amount of information in the message is at most $\log n$ (occurring when all messages are equiprobable). This information can be concealed completely only if the key uncertainty is at least $\log n$. This is the first example of a general principle which will appear frequently: that there is a limit to what we can obtain with a given uncertainty in key—the amount of uncertainty we can introduce into the solution cannot be greater than the key uncertainty.

The situation is somewhat more complicated if the number of messages is infinite. Suppose, for example, that they are generated as infinite sequences of letters by a suitable Markoff process. It is clear that no finite key will give perfect secrecy. We suppose, then, that the key source generates key in the same manner, that is, as an infinite sequence of symbols. Suppose further that only a certain length of key $L_K$ is needed to encipher and decipher a length $L_M$ of message. Let the logarithm of the number of letters in the message alphabet be $R_M$ and that for the key alphabet be $R_K$. Then, from the finite case, it is evident that perfect secrecy requires

$$R_M L_M \leq R_K L_K.$$

This type of perfect secrecy is realized by the Vernam system.

These results have been deduced on the basis of unknown or arbitrary *a priori* probabilities for the messages. The key required for perfect secrecy depends then on the total number of possible messages.

One would expect that, if the message space has fixed known statistics, so that it has a definite mean rate $R$ of generating information, in the sense of MTC, then the amount of key needed could be reduced on the average in just this ratio $\dfrac{R}{R_M}$, and this is indeed true. In fact the message can be passed through a transducer which eliminates the redundancy and reduces the expected length in just this ratio, and then a Vernam system may be applied to the result. Evidently the amount of key used per letter of message is statistically reduced by a factor $\dfrac{R}{R_M}$ and in this case the key source and information source are just matched—a bit of key completely conceals a

bit of message information. It is easily shown also, by the methods used in MTC, that this is the best that can be done.

Perfect secrecy systems have a place in the practical picture—they may be used either where the greatest importance is attached to complete secrecy— e.g., correspondence between the highest levels of command, or in cases where the number of possible messages is small. Thus, to take an extreme example, if only two messages "yes" or "no" were anticipated, a perfect system would be in order, with perhaps the transformation table:

| M    K | A | B |
|--------|---|---|
| yes    | 0 | 1 |
| no     | 1 | 0 |

The disadvantage of perfect systems for large correspondence systems is, of course, the equivalent amount of key that must be sent. In succeeding sections we consider what can be achieved with smaller key size, in particular with finite keys.

## 11. EQUIVOCATION

Let us suppose that a simple substitution cipher has been used on English text and that we intercept a certain amount, $N$ letters, of the enciphered text. For $N$ fairly large, more than say 50 letters, there is nearly always a unique solution to the cipher; i.e., a single good English sequence which transforms into the intercepted material by a simple substitution. With a smaller $N$, however, the chance of more than one solution is greater; with $N = 15$ there will generally be quite a number of possible fragments of text that would fit, while with $N = 8$ a good fraction (of the order of $1/8$) of all reasonable English sequences of that length are possible, since there is seldom more than one repeated letter in the 8. With $N = 1$ any letter is clearly possible and has the same *a posteriori* probability as its *a priori* probability. For one letter the system is perfect.

This happens generally with solvable ciphers. Before any material is intercepted we can imagine the *a priori* probabilities attached to the various possible messages, and also to the various keys. As material is intercepted, the cryptanalyst calculates the *a posteriori* probabilities; and as $N$ increases the probabilities of certain messages increase, and, of most, decrease, until finally only one is left, which has a probability nearly one, while the total probability of all others is nearly zero.

This calculation can actually be carried out for very simple systems. Table I shows the *a posteriori* probabilities for a Caesar type cipher applied to English text, with the key chosen at random from the 26 possibilities. To enable the use of standard letter, digram and trigram frequency tables, the

text has been started at a random point (by opening a book and putting a pencil down at random on the page). The message selected in this way begins "creases to ..." starting inside the word increases. If the message were known to start a sentence a different set of probabilities must be used, corresponding to the frequencies of letters, digrams, etc., at the beginning of sentences.

TABLE I

A *Posteriori* Probabilities for a Caesar Type Cryptogram

| Decipherments | | | | | $N = 1$ | $N = 2$ | $N = 3$ | $N = 4$ | $N = 5$ |
|---|---|---|---|---|---|---|---|---|---|
| C | R | E | A | S | .028 | .0377 | .1111 | .3673 | 1 |
| D | S | F | B | T | .038 | .0314 | | | |
| E | T | G | C | U | .131 | .0881 | | | |
| F | U | H | D | V | .029 | .0189 | | | |
| G | V | I | E | W | .020 | | | | |
| H | W | J | F | X | .053 | .0063 | | | |
| I | X | K | G | Y | .063 | .0126 | | | |
| J | Y | L | H | Z | .001 | | | | |
| K | Z | M | I | A | .004 | | | | |
| L | A | N | J | B | .034 | .1321 | .2500 | | |
| M | B | O | K | C | .025 | | .0222 | | |
| N | C | P | L | D | .071 | .1195 | | | |
| O | D | Q | M | E | .080 | .0377 | | | |
| P | E | R | N | F | .020 | .0818 | .4389 | .6327 | |
| Q | F | S | O | G | .001 | | | | |
| R | G | T | P | H | .068 | .0126 | | | |
| S | H | U | Q | I | .061 | .0881 | .0056 | | |
| T | I | V | R | J | .105 | .2830 | .1667 | | |
| U | J | W | S | K | .025 | | | | |
| V | K | X | T | L | .009 | | | | |
| W | L | Y | U | M | .015 | | .0056 | | |
| X | M | Z | V | N | .002 | | | | |
| Y | N | A | W | O | .020 | | | | |
| Z | O | B | X | P | .001 | | | | |
| A | P | C | Y | Q | .082 | .0503 | | | |
| B | Q | D | Z | R | .014 | | | | |
| H (decimal digits) | | | | | 1.2425 | .9686 | .6034 | .285 | 0 |

The Caesar with random key is a pure cipher and the particular key chosen does not affect the *a posteriori* probabilities. To determine these we need merely list the possible decipherments by all keys and calculate their *a priori* probabilities. The *a posteriori* probabilities are these divided by their sum. These possible decipherments are found by the standard process of "running down the alphabet" from the message and are listed at the left. These form the residue class for the message. For one intercepted letter the *a posteriori* probabilities are equal to the *a priori* probabilities for letters[10] and are shown in the column headed $N = 1$. For two intercepted letters the probabilities are those for digrams adjusted to sum to unity and these are shown in the column $N = 2$.

[10] The probabilities for this table were taken from frequency tables given by Fletcher Pratt in a book "Secret and Urgent" published by Blue Ribbon Books, New York, 1939. Although not complete, they are sufficient for present purposes.

Trigram frequencies have also been tabulated and these are shown in the column $N = 3$. For four- and five-letter sequences probabilities were obtained by multiplication from trigram frequencies since, roughly,

$$p(ijkl) = p(ijk)p_{jk}(l).$$

Note that at three letters the field has narrowed down to four messages of fairly high probability, the others being small in comparison. At four there are two possibilities and at five just one, the correct decipherment.

In principle this could be carried out with any system but, unless the key is very small, the number of possibilities is so large that the work involved prohibits the actual calculation.

This set of *a posteriori* probabilities describes how the cryptanalyst's knowledge of the message and key gradually becomes more precise as enciphered material is obtained. This description, however, is much too involved and difficult to obtain for our purposes. What is desired is a simplified description of this approach to uniqueness of the possible solutions.

A similar situation arises in communication theory when a transmitted signal is perturbed by noise. It is necessary to set up a suitable measure of the uncertainty of what was actually transmitted knowing only the perturbed version given by the received signal. In MTC it was shown that a natural mathematical measure of this uncertainty is the conditional entropy of the transmitted signal when the received signal is known. This conditional entropy was called, for convenience, the equivocation.

From the point of view of the cryptanalyst, a secrecy system is almost identical with a noisy communication system. The message (transmitted signal) is operated on by a statistical element, the enciphering system, with its statistically chosen key. The result of this operation is the cryptogram (analogous to the perturbed signal) which is available for analysis. The chief differences in the two cases are: first, that the operation of the enciphering transformation is generally of a more complex nature than the perturbing noise in a channel; and, second, the key for a secrecy system is usually chosen from a finite set of possibilities while the noise in a channel is more often continually introduced, in effect chosen from an infinite set.

With these considerations in mind it is natural to use the equivocation as a theoretical secrecy index. It may be noted that there are two significant equivocations, that of the key and that of the message. These will be denoted by $H_E(K)$ and $H_E(M)$ respectively. They are given by:

$$H_E(K) = \sum_{E,K} P(E, K) \log P_E(K)$$

$$H_E(M) = \sum_{E,M} P(E, M) \log P_E(K)$$

in which $E$, $M$ and $K$ are the cryptogram, message and key and

$P(E, K)$   is the probability of key $K$ and cryptogram $E$

$P_E(K)$   is the *a posteriori* probability of key $K$ if cryptogram $E$ is intercepted

$P(E, M)$ and $P_E(M)$ are the similar probabilities for message instead of key.

The summation in $H_E(K)$ is over all possible cryptograms of a certain length (say $N$ letters) and over all keys. For $H_E(M)$ the summation is over all messages and cryptograms of length $N$. Thus $H_E(K)$ and $H_E(M)$ are both functions of $N$, the number of intercepted letters. This will sometimes be indicated explicitly by writing $H_E(K, N)$ and $H_E(M, N)$. Note that these are "total" equivocations; i.e., we do not divide by $N$ to obtain the equivocation rate which was used in *MTC*.

The same general arguments used to justify the equivocation as a measure of uncertainty in communication theory apply here as well. We note that zero equivocation requires that one message (or key) have unit probability, all others zero, corresponding to complete knowledge. Considered as a function of $N$, the gradual decrease of equivocation corresponds to increasing knowledge of the original key or message. The two equivocation curves, plotted as functions of $N$, will be called the equivocation characteristics of the secrecy system in question.

The values of $H_E(K, N)$ and $H_E(M, N)$ for the Caesar type cryptogram considered above have been calculated and are given in the last row of Table I. $H_E(K, N)$ and $H_E(M, N)$ are equal in this case and are given in decimal digits (i.e. the logarithmic base 10 is used in the calculation). It should be noted that the equivocation here is for a particular cryptogram, the summation being only over $M$ (or $K$), not over $E$. In general the summation would be over all possible intercepted cryptograms of length $N$ and would give the average uncertainty. The computational difficulties are prohibitive for this general calculation.

## 12. Properties of Equivocation

Equivocation may be shown to have a number of interesting properties, most of which fit into our intuitive picture of how such a quantity should behave. We will first show that the equivocation of key or of a fixed part of a message decreases when more enciphered material is intercepted.

*Theorem 7: The equivocation of key $H_E(K, N)$ is a non-increasing function of $N$. The equivocation of the first $A$ letters of the message is a non-increasing function of the number $N$ which have been intercepted. If $N$ letters have been intercepted, the equivocation of the first $N$ letters of message is less than or equal to that of the key. These may be written:*

$$H_E(K, S) \leq H_E(K, N) \qquad S \geq N,$$
$$H_E(M, S) \leq H_E(M, N) \qquad S \geq N \; (H \text{ for first } A \text{ letters of text})$$
$$H_E(M, N) \leq H_E(K, N)$$

The qualification regarding $A$ letters in the second result of the theorem is so that the equivocation will not be calculated with respect to the amount of message that has been intercepted. If it is, the message equivocation may (and usually does) increase for a time, due merely to the fact that more letters stand for a larger possible range of messages. The results of the theorem are what we might hope from a good secrecy index, since we would hardly expect to be worse off on the average after intercepting additional material than before. The fact that they can be proved gives further justification to our use of the equivocation measure.

The results of this theorem are a consequence of certain properties of conditional entropy proved in MTC. Thus, to show the first or second statements of Theorem 7, we have for any chance events $A$ and $B$

$$H(B) \geq H_A(B).$$

If we identify $B$ with the key (knowing the first $S$ letters of cryptogram) and $A$ with the remaining $N - S$ letters we obtain the first result. Similarly identifying $B$ with the message gives the second result. The last result follows from

$$H_E(M) \leq H_E(K, M) = H_E(K) + H_{E,K}(M)$$

and the fact that $H_{E,K}(M) = 0$ since $K$ and $E$ uniquely determine $M$.

Since the message and key are chosen independently we have:

$$H(M, K) = H(M) + H(K).$$

Furthermore,

$$H(M, K) = H(E, K) = H(E) + H_E(K),$$

the first equality resulting from the fact that knowledge of $M$ and $K$ or of $E$ and $K$ is equivalent to knowledge of all three. Combining these two we obtain a formula for the equivocation of key:

$$H_E(K) = H(M) + H(K) - H(E).$$

In particular, if $H(M) = H(E)$ then the equivocation of key, $H_E(K)$, is equal to the *a priori* uncertainty of key, $H(K)$. This occurs in the perfect systems described above.

A formula for the equivocation of message can be found by similar means. We have:

$$H(M, E) = H(E) + H_E(M) = H(M) + H_M(E)$$
$$H_E(M) = H(M) + H_M(E) - H(E).$$

If we have a product system $S = TR$, it is to be expected that the second enciphering process will not decrease the equivocation of message. That this is actually true can be shown as follows: Let $M$, $E_1$, $E_2$ be the message and the first and second encipherments, respectively. Then

$$P_{E_1 E_2}(M) = P_{E_1}(M).$$

Consequently

$$H_{E_1 E_2}(M) = H_{E_1}(M).$$

Since, for any chance variables, $x$, $y$, $z$, $H_{xy}(z) \leq H_y(z)$, we have the desired result, $H_{E_2}(M) \geq H_{E_1}(M)$.

*Theorem 8: The equivocation in message of a product system $S = TR$ is not less than that when only $R$ is used.*

Suppose now we have a system $T$ which can be written as a weighted sum of several systems $R, S, \cdots, U$

$$T = p_1 R + p_2 S + \cdots + p_m U \qquad \sum p_i = 1$$

and that systems $R, S, \cdots, U$ have equivocations $H_1, H_2, H_3, \cdots, H_m$.

*Theorem 9: The equivocation $H$ of a weighted sum of systems is bounded by the inequalities*

$$\sum p_i H_i \leq H \leq \sum p_i H_i - \sum p_i \log p_i.$$

*These are best limits possible. The $H$'s may be equivocations either of key or message.*

The upper limit is achieved, for example, in strongly ideal systems (to be described later) where the decomposition is into the simple transformations of the system. The lower limit is achieved if all the systems $R, S, \cdots, U$ go to completely different cryptogram spaces. This theorem is also proved by the general inequalities governing equivocation,

$$H_A(B) \leq H(B) \leq H(A) + H_A(B).$$

We identify $A$ with the particular system being used and $B$ with the key or message.

There is a similar theorem for weighted sums of languages. For this we identify $A$ with the particular language.

*Theorem 10: Suppose a system can be applied to languages $L_1, L_2, \cdots, L_m$ and has equivocation characteristics $H_1, H_2, \cdots, H_m$. When applied to the weighted sum $\sum p_i L_i$, the equivocation $H$ is bounded by*

$$\sum p_i H_i \leq H \leq \sum p_i H_i - \sum p_i \log p_i.$$

*These limits are the best possible and the equivocations in question can be either for key or message.*

The total redundancy $D_N$ for $N$ letters of message is defined by

$$D_N = \log G - H(M)$$

where $G$ is the total number of messages of length $N$ and $H(M)$ is the uncertainty in choosing one of these. In a secrecy system where the total number of possible cryptograms is equal to the number of possible messages of length $N$, $H(E) \leq \log G$. Consequently

$$H_E(K) = H(K) + H(M) - H(E)$$
$$\geq H(K) - [\log G - H(M)].$$

Hence

$$H(K) - H_E(K) \leq D_N.$$

This shows that, in a closed system, for example, the decrease in equivocation of key after $N$ letters have been intercepted is not greater than the redundancy of $N$ letters of the language. In such systems, which comprise the majority of ciphers, it is only the existence of redundancy in the original messages that makes a solution possible.

Now suppose we have a pure system. Let the different residue classes of messages be $C_1, C_2, C_3, \cdots, C_r$, and the corresponding set of residue classes of cryptograms be $C_1', C_2', \cdots, C_r'$. The probability of each $E$ in $C_1'$ is the same:

$$P(E) = \frac{P(C_i)}{\varphi_i} \qquad E \text{ a member of } C_i$$

where $\varphi_i$ is the number of different messages in $C_i$. Thus we have

$$H(E) = -\sum_i \varphi_i \frac{P(C_i)}{\varphi_i} \log \frac{P(C_i)}{\varphi_i}$$
$$= -\sum P(C_i) \log \frac{P(C_i)}{\varphi_i}$$

Substituting in our equation for $H_E(K)$ we obtain:

*Theorem 11: For a pure cipher*

$$H_E(K) = H(K) + H(M) + \sum_i P(C_i) \log \frac{P(C_i)}{\varphi_i}.$$

This result can be used to compute $H_E(K)$ in certain cases of interest.

### 13. Equivocation for Simple Substitution on a Two Letter Language

We will now calculate the equivocation in key or message when simple substitution is applied to a two letter language, with probabilities $p$ and $q$ for 0 and 1, and successive letters chosen independently. We have

$$H_E(M) \;=\; H_E(K) \;=\; -\sum \; P(E)P_E(K) \log P_E(K)$$

The probability that $E$ contains exactly $s$ 0's in a particular permutation is:

$$\tfrac{1}{2}(p^s q^{N-s} + q^s p^{N-s})$$



Fig. 6—Equivocation for simple substitution on two-letter language.

and the *a posteriori* probabilities of the identity and inverting substitutions (the only two in the system) are respectively:

$$P_E(0) \;=\; \frac{p^s q^{N-s}}{(p^s q^{N-s} + q^s p^{N-s})} \quad P_E(1) \;=\; \frac{p^{N-s} q^s}{(p^s q^{N-s} + q^s p^{N-s})} \cdot$$

There are $\binom{N}{s}$ terms for each $s$ and hence

$$H_E(K, N) \;=\; -\sum_{s} \binom{N}{s} p^s q^{N-s} \log \frac{p^s q^{N-s}}{(p^s q^{N-s} + q^s p^{N-s})} \cdot$$

For $p = \frac{1}{3}$, $q = \frac{2}{3}$, and for $p = \frac{1}{8}$, $q = \frac{7}{8}$, $H_E(K, N)$ has been calculated and is shown in Fig. 6.

14. THE EQUIVOCATION CHARACTERISTIC FOR A "RANDOM" CIPHER

In the preceding section we have calculated the equivocation characteristic for a simple substitution applied to a two-letter language. This is about the simplest type of cipher and the simplest language structure possible, yet already the formulas are so involved as to be nearly useless. What are we to do with cases of practical interest, say the involved transformations of a fractional transposition system applied to English with its extremely complex statistical structure? This complexity itself suggests a method of approach. Sufficiently complicated problems can frequently be solved statistically. To facilitate this we define the notion of a "random" cipher.

We make the following assumptions:

1. The number of possible messages of length $N$ is $T = 2^{R_0 N}$, thus $R_0 = \log_2 G$, where $G$ is the number of letters in the alphabet. The number of possible cryptograms of length $N$ is also assumed to be $T$.

2. The possible messages of length $N$ can be divided into two groups: one group of high and fairly uniform a priori probability, the second group of negligibly small total probability. The high probability group will contain $S = 2^{RN}$ messages, where $R = H(M)/N$, that is, $R$ is the entropy of the message source per letter.

3. The deciphering operation can be thought of as a series of lines, as in Figs. 2 and 4, leading back from each $E$ to various $M$'s. We assume $k$ different equiprobable keys so there will be $k$ lines leading back from each $E$. For the random cipher we suppose that the lines from each $E$ go back to a random selection of the possible messages. Actually, then, a random cipher is a whole ensemble of ciphers and the equivocation is the average equivocation for this ensemble.

The equivocation of key is defined by

$$H_E(K) = \sum P(E) P_E(K) \log P_E(K).$$

The probability that exactly $m$ lines go back from a particular $E$ to the high probability group of messages is

$$\binom{k}{m} \left(\frac{S}{T}\right)^m \left(1 - \frac{S}{T}\right)^{k-m}$$

If a cryptogram with $m$ such lines is intercepted the equivocation is $\log m$. The probability of such a cryptogram is $\dfrac{mT}{SK}$, since it can be produced by

$m$ keys from high probability messages each with probability $\dfrac{T}{S}$. Hence the equivocation is:

$$H_E(K) = \frac{T}{Sk} \sum_{m=1}^{k} \binom{k}{m} \left(\frac{S}{T}\right)^m \left(1 - \frac{S}{T}\right)^{k-m} m \log \dot{m}$$

We wish to find a simple approximation to this when $k$ is large. If the expected value of $m$, namely $\overline{m} = Sk/T$, is $\gg 1$, the variation of $\log m$ over the range where the binomial distribution assumes large values will be small, and we can replace $\log m$ by $\log \overline{m}$. This can now be factored out of the summation, which then reduces to $\overline{m}$. Hence, in this condition,

$$H_E(K) \doteq \log \frac{Sk}{T} = \log S - \log T + \log k$$

$$H_E(K) \doteq H(K) - DN,$$

where $D$ is the redundancy per letter of the original language ($D = D_N/N$).

If $\overline{m}$ is small compared to the large $k$, the binomial distribution can be approximated by a Poisson distribution:

$$\binom{k}{m} p^m q^{k-m} \doteq \frac{e^{-\lambda}\lambda^m}{m!}$$

where $\lambda = \dfrac{Sk}{T}$. Hence

$$H_E(K) \doteq \frac{1}{\lambda} e^{-\lambda} \sum_{2}^{\infty} \frac{\lambda^m}{m!} m \log m.$$

If we replace $m$ by $m + 1$, we obtain:

$$H_E(K) \doteq e^{-\lambda} \sum_{1}^{\infty} \frac{\lambda^m}{m!} \log (m + 1).$$

This may be used in the region where $\lambda$ is near unity. For $\lambda \ll 1$, the only important term in the series is that for $m = 1$; omitting the others we have:

$$\begin{aligned} H_E(K) &\doteq e^{-\lambda}\lambda \log 2 \\ &\doteq \lambda \log 2 \\ &\doteq 2^{-ND}k \log 2 \,. \end{aligned}$$

To summarize: $H_E(K)$, considered as a function of $N$, the number of intercepted letters, starts off at $H(K)$ when $N = 0$. It decreases linearly with a slope $-D$ out to the neighborhood of $N = \dfrac{H(K)}{D}$. After a short transition region, $H_E(K)$ follows an exponential with "half life" distance

$\frac{1}{D}$ if $D$ is measured in bits per letter. This behavior is shown in Fig. 7, together with the approximating curves.

By a similar argument the equivocation of message can be calculated. It is

$$H_E(M) = R_0N \text{ for } R_0N \ll H_E(K)$$
$$H_E(M) = H_E(K) \text{ for } R_0N \gg H_E(K)$$
$$H_E(M) = H_E(K) - \varphi(N) \text{ for } R_0N \sim H_E(K)$$

where $\varphi(N)$ is the function shown in Fig. 7 with $N$ scale reduced by factor of $\frac{D}{R_0}$. Thus, $H_E(M)$ rises linearly with slope $R_0$, until it nearly intersects



Fig. 7—Equivocation for random cipher.

the $H_E(K)$ line. After a rounded transition it follows the $H_E(K)$ curve down.

It will be seen from Fig. 7 that the equivocation curves approach zero rather sharply. Thus we may, with but little ambiguity, speak of a point at which the solution becomes unique. This number of letters will be called the unicity distance. For the random cipher it is approximately $H(K)/D$.

## 15. Application to Standard Ciphers

Most of the standard ciphers involve rather complicated enciphering and deciphering operations. Furthermore, the statistical structure of natural languages is extremely involved. It is therefore reasonable to assume that the formulas derived for the random cipher may be applied in such cases. It is necessary, however, to apply certain corrections in some cases. The main points to be observed are the following:

1. We assumed for the random cipher that the possible decipherments of a cryptogram are a random selection from the possible messages. While not strictly true in ordinary systems, this becomes more nearly the case as the complexity of the enciphering operations and of the language structure increases. With a transposition cipher it is clear that letter frequencies are preserved under decipherment operations. This means that the possible decipherments are chosen from a more limited group, not the entire message space, and the formula should be changed. In place of $R_0$ one uses $R_1$ the entropy rate for a language with independent letters but with the regular letter frequencies. In some other cases a definite tendency toward returning the decipherments to high probability messages can be seen. If there is no clear tendency of this sort, and the system is fairly complicated, then it is reasonable to use the random cipher analysis.

2. In many cases the complete key is not used in enciphering short messages. For example, in a simple substitution, only fairly long messages will contain all letters of the alphabet and thus involve the complete key. Obviously the random assumption does not hold for small $N$ in such a case, since all the keys which differ only in the letters not yet appearing in the cryptogram lead back to the same message and are not randomly distributed. This error is easily corrected to a good approximation by the use of a "key appearance characteristic." One uses, at a particular $N$, the effective amount of key that may be expected with that length of cryptogram. For most ciphers, this is easily estimated.

3. There are certain "end effects" due to the definite starting of the message which produce a discrepancy from the random characteristics. If we take a random starting point in English text, the first letter (when we do not observe the preceding letters) has a possibility of being any letter with the ordinary letter probabilities. The next letter is more completely specified since we then have digram frequencies. This decrease in choice value continues for some time. The effect of this on the curve is that the straight line part is displaced, and approached by a curve depending on how much the statistical structure of the language is spread out over adjacent letters. As a first approximation the curve can be corrected by shifting the line over to the half redundancy point—i.e., the number of letters where the language redundancy is half its final value.

If account is taken of these three effects, reasonable estimates of the equivocation characteristic and unicity point can be made. The calculation can be done graphically as indicated in Fig. 8. One draws the key appearance characteristic and the total redundancy curve $D_N$ (which is usually sufficiently well represented by the line $ND_\infty$). The difference between these out to the neighborhood of their intersection is $H_E(M)$. With a simple substitution cipher applied to English, this calculation gave the

curves shown in Fig. 9. The key appearance characteristic in this case was estimated by counting the number of different letters appearing in typical English passages of $N$ letters. In so far as experimental data on the simple substitution could be found, they agree very well with the curves of Fig. 9, considering the various idealizations and approximations which have been made. For example, the unicity point, at about 27 letters, can be shown experimentally to lie between the limits 20 and 30. With 30 letters there is



Fig. 8—Graphical calculation of equivocation.

nearly always a unique solution to a cryptogram of this type and with 20 it is usually easy to find a number of solutions.

With transposition of period $d$ (random key), $H(K) = \log d!$, or about $d \log d/e$ (using a Stirling approximation for $d!$). If we take .6 decimal digits per letter as the appropriate redundancy, remembering the preservation of letter frequencies, we obtain about $1.7d \log d/e$ as the unicity distance. This also checks fairly well experimentally. Note that in this case $H_E(M)$ is defined only for integral multiples of $d$.

With the Vigenère the unicity point will occur at about $2d$ letters, and this too is about right. The Vigenère characteristic with the same key size

Fig. 9—Equivocation for simple substitution on English.

Fig. 10—Equivocation for Vigenère on English.

698 BELL SYSTEM TECHNICAL JOURNAL

as simple substitution will be approximately as shown in Fig. 10. The Vigenère, Playfair and Fractional cases are more likely to follow the theoretical formulas for random ciphers than simple substitution and transposition. The reason for this is that they are more complex and give better mixing characteristics to the messages on which they operate.
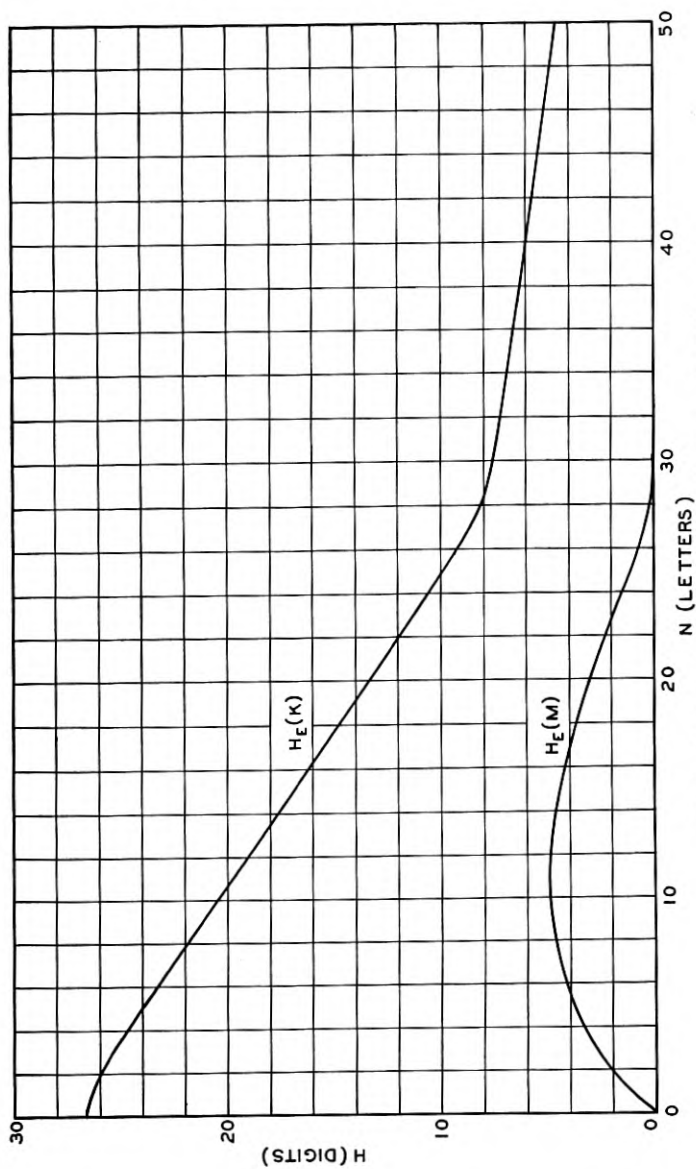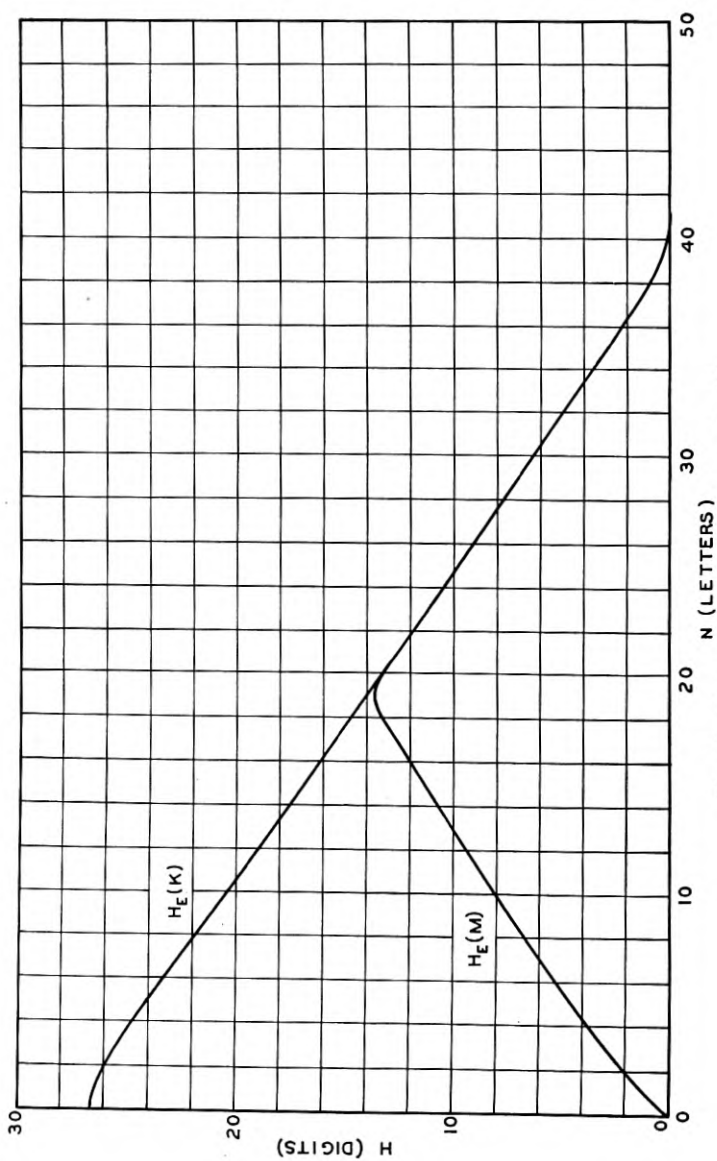
The mixed alphabet Vigenère (each of $d$ alphabets mixed independently and used sequentially) has a key size,

$$H(K) = d \log 26! = 26.3d$$

and its unicity point should be at about $53d$ letters.

These conclusions can also be put to a rough experimental test with the Caesar type cipher. In the particular cryptogram analyzed in Table I, section 11, the function $(H_E(K, N)$ has been calculated and is given below, together with the values for a random cipher.

| $N$ | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| $H$ (observed) | 1.41 | 1.24 | .97 | .60 | .28 | 0 |
| $H$ (calculated) | 1.41 | 1.25 | .98 | .54 | .15 | .03 |

The agreement is seen to be quite good, especially when we remember that the observed $H$ should actually be the average of many different cryptograms, and that $D$ for the larger values of $N$ is only roughly estimated.

It appears then that the random cipher analysis can be used to estimate equivocation characteristics and the unicity distance for the ordinary types of ciphers.

## 16. VALIDITY OF A CRYPTOGRAM SOLUTION

The equivocation formulas are relevant to questions which sometimes arise in cryptographic work regarding the validity of an alleged solution to a cryptogram. In the history of cryptography there have been many cryptograms, or possible cryptograms, where clever analysts have found a "solution." It involved, however, such a complex process, or the material was so meager that the question arose as to whether the cryptanalyst had "read a solution" into the cryptogram. See, for example, the Bacon-Shakespeare ciphers and the "Roger Bacon" manuscript.[10]

In general we may say that if a proposed system and key solves a cryptogram for a length of material considerably greater than the unicity distance the solution is trustworthy. If the material is of the same order or shorter than the unicity distance the solution is highly suspicious.

This effect of redundancy in gradually producing a unique solution to a cipher can be thought of in another way which is helpful. The redundancy is essentially a series of conditions on the letters of the message, which

[10] See Fletcher Pratt, *loc. cit.*

insure that it be statistically reasonable. These consistency conditions pro-
duce corresponding consistency conditions in the cryptogram. The key gives
a certain amount of freedom to the cryptogram but, as more and more
letters are intercepted, the consistency conditions use up the freedom al-
lowed by the key. Eventually there is only one message and key which
satisfies all the conditions and we have a unique solution. In the random
cipher the consistency conditions are, in a sense "orthogonal" to the "grain
of the key" and have their full effect in eliminating messages and keys as
rapidly as possible. This is the usual case. However, by proper design it
is possible to "line up" the redundancy of the language with the "grain of
the key" in such a way that the consistency conditions are automatically
satisfied and $H_E(K)$ does not approach zero. These "ideal" systems, which
will be considered in the next section, are of such a nature that the trans-
formations $T_i$ all induce the same probabilities in the $E$ space.

## 17. Ideal Secrecy Systems.

We have seen that perfect secrecy requires an infinite amount of key if
we allow messages of unlimited length. With a finite key size, the equivoca-
tion of key and message generally approaches zero, but not necessarily so.
In fact it is possible for $H_E(K)$ to remain constant at its initial value $H(K)$.
Then, no matter how much material is intercepted, there is not a unique
solution but many of comparable probability. We will define an "ideal"
system as one in which $H_E(K)$ and $H_E(M)$ do not approach zero as $N \to \infty$.
A "strongly ideal" system is one in which $H_E(K)$ remains constant
at $H(K)$.

An example is a simple substitution on an artificial language in which
all letters are equiprobable and successive letters independently chosen.
It is easily seen that $H_E(K) = H(K)$ and $H_E(M)$ rises linearly along a line
of slope $\log G$ (where $G$ is the number of letters in the alphabet) until it
strikes the line $H(K)$, after which it remains constant at this value.

With natural languages it is in general possible to approximate the ideal
characteristic—the unicity point can be made to occur for as large $N$ as is
desired. The complexity of the system needed usually goes up rapidly when
we attempt to do this, however. It is not always possible to attain actually
the ideal characteristic with any system of finite complexity.

To approximate the ideal equivocation, one may first operate on the
message with a transducer which removes all redundancies. After this almost
any simple ciphering system—substitution, transposition, Vigenère, etc.,
is satisfactory. The more elaborate the transducer and the nearer the
output is to the desired form, the more closely will the secrecy system ap-
proximate the ideal characteristic.

*Theorem 12:* *A necessary and sufficient condition that $T$ be strongly ideal is that, for any two keys, $T_i^{-1}T_j$ is a measure preserving transformation of the message space into itself.*

This is true since the *a posteriori* probability of each key is equal to its *a priori* probability if and only if this condition is satisfied.

## 18. Examples of Ideal Secrecy Systems

Suppose our language consists of a sequence of letters all chosen independently and with equal probabilities. Then the redundancy is zero, and from a result of section 12, $H_E(K) = H(K)$. We obtain the result

*Theorem 13:* *If all letters are equally likely and independent any closed cipher is strongly ideal.*

The equivocation of message will rise along the key appearance characteristic which will usually approach $H(K)$, although in some cases it does not. In the cases of $n$-gram substitution, transposition, Vigenère, and variations, fractional, etc., we have strongly ideal systems for this simple language with $H_E(M) \rightarrow H(K)$ as $N \rightarrow \infty$.

Ideal secrecy systems suffer from a number of disadvantages.

1. The system must be closely matched to the language. This requires an extensive study of the structure of the language by the designer. Also a change in statistical structure or a selection from the set of possible messages, as in the case of probable words (words expected in this particular cryptogram), renders the system vulnerable to analysis.

2. The structure of natural languages is extremely complicated, and this implies a complexity of the transformations required to eliminate redundancy. Thus any machine to perform this operation must necessarily be quite involved, at least in the direction of information storage, since a "dictionary" of magnitude greater than that of an ordinary dictionary is to be expected.

3. In general, the transformations required introduce a bad propagation of error characteristic. Error in transmission of a single letter produces a region of changes near it of size comparable to the length of statistical effects in the original language.

## 19. Further Remarks on Equivocation and Redundancy

We have taken the redundancy of "normal English" to be about .7 decimal digits per letter or a redundancy of 50%. This is on the assumption that word divisions were omitted. It is an approximate figure based on statistical structure extending over about 8 letters, and assumes the text to be of an ordinary type, such as newspaper writing, literary work, etc. We may note here a method of roughly estimating this number that is of some cryptographic interest.

A running key cipher is a Vernam type system where, in place of a random sequence of letters, the key is a meaningful text. Now it is known that running key ciphers can usually be solved uniquely. This shows that English can be reduced by a factor of two to one and implies a redundancy of at least 50%. This figure cannot be increased very much, however, for a number of reasons, unless long range "meaning" structure of English is considered.

The running key cipher can be easily improved to lead to ciphering systems which could not be solved without the key. If one uses in place of one English text, about 4 different texts as key, adding them all to the message, a sufficient amount of key has been introduced to produce a high positive equivocation. Another method would be to use, say, every 10th letter of the text as key. The intermediate letters are omitted and cannot be used at any other point of the message. This has much the same effect, since these spaced letters are nearly independent.

The fact that the vowels in a passage can be omitted without essential loss suggests a simple way of greatly improving almost any ciphering system. First delete all vowels, or as much of the message as possible without running the risk of multiple reconstructions, and then encipher the residue. Since this reduces the redundancy by a factor of perhaps 3 or 4 to 1, the unicity point will be moved out by this factor. This is one way of approaching ideal systems—using the decipherer's knowledge of English as part of the deciphering system.

## 20. Distribution of Equivocation

A more complete description of a secrecy system applied to a language than is afforded by the equivocation characteristics can be found by giving the *distribution of equivocation*. For $N$ intercepted letters we consider the fraction of cryptograms for which the equivocation (for these particular $E$'s, not the mean $H_E(M)$) lies between certain limits. This gives a density distribution function

$$P(H_E(M), N) \ dH_E(M)$$

for the probability that for $N$ letters $H$ lies between the limits $H$ and $H + dH$. The mean equivocation we have previously studied is the mean of this distribution. The function $P(H_E(M), N)$ can be thought of as plotted along a third dimension, normal to the paper, on the $H_E(M)$, $N$ plane. If the language is pure, with a small influence range, and the cipher is pure, the function will usually be a ridge in this plane whose highest point follows approximately the mean $H_E(M)$, at least until near the unicity point. In this case, or when the conditions are nearly verified, the mean curve gives a reasonably complete picture of the system.

On the other hand, if the language is not pure, but made up of a set of pure components

$$L = \sum p_i L_i$$

having different equivocation curves with the system, then the total distribution will usually be made up of a series of ridges. There will be one for each $L_i$ weighted in accordance with its $p_i$. The mean equivocation characteristic will be a line somewhere in the midst of these ridges and may not give a very complete picture of the situation. This is shown in Fig. 11. A similar effect occurs if the system is not pure but made up of several systems with different $H$ curves.

The effect of mixing pure languages which are near to one another in statistical structure is to increase the width of the ridge. Near the unicity



Fig. 11—Distribution of equivocation with a mixed language $L = \frac{1}{2}L_1 + \frac{1}{2}L_2$.

point this tends to raise the mean equivocation, since equivocation cannot become negative and the spreading is chiefly in the positive direction. We expect, therefore, that in this region the calculations based on the random cipher should be somewhat low.

## PART III

## PRACTICAL SECRECY

### 21. THE WORK CHARACTERISTIC

After the unicity point has been passed in intercepted material there will usually be a unique solution to the cryptogram. The problem of isolating this single solution of high probability is the problem of cryptanalysis. In the region before the unicity point we may say that the problem of cryptanalysis is that of isolating all the possible solutions of high probability (compared to the remainder) and determining their various probabilities.

Although it is always possible in principle to determine these solutions (by trial of each possible key for example), different enciphering systems show a wide variation in the amount of work required. The average amount of work to determine the key for a cryptogram of $N$ letters, $W(N)$, measured say in man hours, may be called the work characteristic of the system. This average is taken over all messages and all keys with their appropriate probabilities. The function $W(N)$ is a measure of the amount of "practical secrecy" afforded by the system.

For a simple substitution on English the work and equivocation characteristics would be somewhat as shown in Fig. 12. The dotted portion of



Fig. 12—Typical work and equivocation characteristics.

the curve is in the range where there are numerous possible solutions and these must all be determined. In the solid portion after the unicity point only one solution exists in general, but if only the minimum necessary data are given a great deal of work must be done to isolate it. As more material is available the work rapidly decreases toward some asymptotic value— where the additional data no longer reduces the labor.

Essentially the behavior shown in Fig. 12 can be expected with any type of secrecy system where the equivocation approaches zero. The scale of man hours required, however, will differ greatly with different types of ciphers, even when the $H_E(M)$ curves are about the same. A Vigenère or compound Vigenère, for example, with the same key size would have a

much better (i.e., much higher) work characteristic. A good practical secrecy system is one in which the $W(N)$ curve remains sufficiently high, out to the number of letters one expects to transmit with the key, to prevent the enemy from actually carrying out the solution, or to delay it to such an extent that the information is then obsolete.

We will consider in the following sections ways of keeping the function $W(N)$ large, even though $H_E(K)$ may be practically zero. This is essentially a "max min" type of problem as is always the case when we have a battle of wits.[11] In designing a good cipher we must maximize the minimum amount of work the enemy must do to break it. It is not enough merely to be sure none of the standard methods of cryptanalysis work—we must be sure that no method whatever will break the system easily. This, in fact, has been the weakness of many systems; designed to resist all the known methods of solution, they later gave rise to new cryptanalytic techniques which rendered them vulnerable to analysis.

The problem of good cipher design is essentially one of finding difficult problems, subject to certain other conditions. This is a rather unusual situation, since one is ordinarily seeking the simple and easily soluble problems in a field.

How can we ever be sure that a system which is not ideal and therefore *has* a unique solution for sufficiently large $N$ will require a large amount of work to break with *every* method of analysis? There are two approaches to this problem; (1) We can study the possible methods of solution available to the cryptanalyst and attempt to describe them in sufficiently general terms to cover any methods he might use. We then construct our system to resist this "general" method of solution. (2) We may construct our cipher in such a way that breaking it is equivalent to (or requires at some point in the process) the solution of some problem known to be laborious. Thus, if we could show that solving a certain system requires at least as much work as solving a system of simultaneous equations in a large number of unknowns, of a complex type, then we would have a lower bound of sorts for the work characteristic.

The next three sections are aimed at these general problems. It is difficult to define the pertinent ideas involved with sufficient precision to obtain results in the form of mathematical theorems, but it is believed that the conclusions, in the form of general principles, are correct.

---

[11] See von Neumann and Morgenstern, *loc. cit.* The situation between the cipher designer and cryptanalyst can be thought of as a "game" of a very simple structure; a zero-sum two-person game with complete information, and just two "moves." The cipher designer chooses a system for his "move." Then the cryptanalyst is informed of this choice and chooses a method of analysis. The "value" of the play is the average work required to break a cryptogram in the system by the method chosen.

## 22. Generalities on the Solution of Cryptograms

After the unicity distance has been exceeded in intercepted material, any system can be solved in principle by merely trying each possible key until the unique solution is obtained—i.e., a deciphered message which "makes sense" in the original language. A simple calculation shows that this method of solution (which we may call *complete trial and error*) is totally impractical except when the key is absurdly small.

Suppose, for example, we have a key of 26! possibilities or about 26.3 decimal digits, the same size as in simple substitution on English. This is, by any significant measure, a small key. It can be written on a small slip of paper, or memorized in a few minutes. It could be registered on 27 switches, each having ten positions, or on 88 two-position switches.

Suppose further, to give the cryptanalyst every possible advantage, that he constructs an electronic device to try keys at the rate of one each microsecond (perhaps automatically selecting from the results by a $\chi^2$ test for statistical significance). He may expect to reach the right key about half way through, and after an elapsed time of about $2 \times 10^{26}/2 \times 60^2 \times 24 \times 365 \times 10^6$ or $3 \times 10^{12}$ years.

In other words, even with a small key complete trial and error will never be used in solving cryptograms, except in the trivial case where the key is extremely small, e.g., the Caesar with only 26 possibilities, or 1.4 digits. The trial and error which is used so commonly in cryptography is of a different sort, or is augmented by other means. If one had a secrecy system which required complete trial and error it would be extremely safe. Such a system would result, it appears, if the meaningful original messages, all say of 1000 letters, were a random selection from the set of all sequences of 1000 letters. If any of the simple ciphers were applied to this type of language it seems that little improvement over complete trial and error would be possible.

The methods of cryptanalysis actually used often involve a great deal of trial and error, but in a different way. First, the trials progress from more probable to less probable hypotheses, and, second, each trial disposes of a large group of keys, not a single one. Thus the key space may be divided into say 10 subsets, each containing about the same number of keys. By at most 10 trials one determines which subset is the correct one. This subset is then divided into several secondary subsets and the process repeated. With the same key size ($26! \doteq 2 \times 10^{26}$) we would expect about $26 \times 5$ or 130 trials as compared to $10^{26}$ by complete trial and error. The possibility of choosing the most likely of the subsets first for test would improve this result even more. If the divisions were into two compartments (the best way to

minimize the number of trials) only 88 trials would be required. Whereas complete trial and error requires trials to the order of the number of keys, this subdividing trial and error requires only trials to the order of the key size in bits.

This remains true even when the different keys have different probabilities. The proper procedure, then, to minimize the expected number of trials is to divide the key space into subsets of equiprobability. When the proper subset is determined, this is again subdivided into equiprobability subsets. If this process can be continued the number of trials expected when each division is into two subsets will be

$$h = \frac{H(K)}{\log 2}$$

If each test has $S$ possible results and each of these corresponds to the key being in one of $S$ equiprobability subsets, then

$$h = \frac{H(K)}{\log S}$$

trials will be expected. The intuitive significance of these results should be noted. In the two-compartment test with equiprobability, each test yields one bit of information as to the key. If the subsets have very different probabilities, as in testing a single key in complete trial and error, only a small amount of information is obtained from the test. Thus with 26! equiprobable keys, a test of one yields only

$$-\left[\frac{26! - 1}{26!} \log \frac{26! - 1}{26!} + \frac{1}{26!} \log \frac{1}{26!}\right]$$

or about $10^{-25}$ bits of information. Dividing into $S$ equiprobability subsets maximizes the information obtained from each trial at $\log S$, and the expected number of trials is the total information to be obtained, that is $H(K)$, divided by this amount.

The question here is similar to various coin weighing problems that have been circulated recently. A typical example is the following: It is known that one coin in 27 is counterfeit, and slightly lighter than the rest. A chemist's balance is available and the counterfeit coin is to be isolated by a series of weighings. What is the least number of weighings required to do this? The correct answer is 3, obtained by first dividing the coins into three groups of 9 each. Two of these are compared on the balance. The three possible results determine the set of 9 containing the counterfeit. This set is then divided into 3 subsets of 3 each and the process continued. The set of coins corresponds to the set of keys, the counterfeit coin to the correct key, and the weighing procedure to a trial or test. The original uncertainty is $\log_2 27$

bits, and each trial yields $\log_2 3$ bits of information; thus, when there is no "diophantine trouble," $\log_2 27/\log_2 3$ or 3 trials are sufficient.

This method of solution is feasible only if the key space can be divided into a small number of subsets, with a simple method of determining the subset to which the correct key belongs. One does not need to assume a complete key in order to apply a consistency test and determine if the assumption is justified—an assumption on a part of the key (or as to whether the key is in some large section of the key space) can be tested. In other words it is possible to solve for the key bit by bit.

The possibility of this method of analysis is the crucial weakness of most ciphering systems. For example, in simple substitution, an assumption on a single letter can be checked against its frequency, variety of contact, doubles or reversals, etc. In determining a single letter the key space is reduced by 1.4 decimal digits from the original 26. The same effect is seen in all the elementary types of ciphers. In the Vigenère, the assumption of two or three letters of the key is easily checked by deciphering at other points with this fragment and noting whether clear emerges. The compound Vigenère is much better from this point of view, if we assume a fairly large number of component periods, producing a repetition rate larger than will be intercepted. In this case as many key letters are used in enciphering each letter as there are periods. Although this is only a fraction of the entire key, at least a fair number of letters must be assumed before a consistency check can be applied.

Our first conclusion then, regarding practical small key cipher design, is that a considerable amount of key should be used in enciphering each small element of the message.

## 23. STATISTICAL METHODS

It is possible to solve many kinds of ciphers by statistical analysis. Consider again simple substitution. The first thing a cryptanalyst does with an intercepted cryptogram is to make a frequency count. If the cryptogram contains, say, 200 letters it is safe to assume that few, if any, of the letters are out of their frequency groups, this being a division into 4 sets of well defined frequency limits. The logarithm of the number of keys within this limitation may be calculated as

$$\log 2! \, 9! \, 9! \, 6! = 14.28$$

and the simple frequency count thus reduces the key uncertainty by 12 decimal digits, a tremendous gain.

In general, a statistical attack proceeds as follows: A certain statistic is measured on the intercepted cryptogram $E$. This statistic is such that for all reasonable messages $M$ it assumes about the same value, $S_K$, the value

depending only on the particular key $K$ that was used. The value thus obtained serves to limit the possible keys to those which would give values of $S$ in the neighborhood of that observed. A statistic which does not depend on $K$ or which varies as much with $M$ as with $K$ is not of value in limiting $K$. Thus, in transposition ciphers, the frequency count of letters gives no information about $K$—every $K$ leaves this statistic the same. Hence one can make no use of a frequency count in breaking transposition ciphers.

More precisely one can ascribe a *"solving power"* to a given statistic $S$. For each value of $S$ there will be a conditional equivocation of the key $H_S(K)$, the equivocation when $S$ has its particular value, and that is all that is known concerning the key. The weighted mean of these values

$$\sum P(S) \quad H_S(K)$$

gives the mean equivocation of the key when $S$ is known, $P(S)$ being the *a priori* probability of the particular value $S$. The key size $H(K)$, less this mean equivocation, measures the "solving power" of the statistic $S$.

In a strongly ideal cipher *all* statistics of the cryptogram are independent of the particular key used. This is the measure preserving property of $T_j T_k^{-1}$ on the $E$ space or $T_j^{-1} T_k$ on the $M$ space mentioned above.

There are good and poor statistics, just as there are good and poor methods of trial and error. Indeed the trial and error testing of an hypothesis *is* is a type of statistic, and what was said above regarding the best types of trials holds generally. A good statistic for solving a system must have the following properties:

1. It must be simple to measure.
2. It must depend more on the key than on the message if it is meant to solve for the key. The variation with $M$ should not mask its variation with $K$.
3. The values of the statistic that can be "resolved" in spite of the "fuzziness" produced by variation in $M$ should divide the key space into a number of subsets of comparable probability, with the statistic specifying the one in which the correct key lies. The statistic should give us sizeable information about the key, not a tiny fraction of a bit.
4. The information it gives must be simple and usable. Thus the subsets in which the statistic locates the key must be of a simple nature in the key space.

Frequency count for simple substitution is an example of a very good statistic.

Two methods (other than recourse to ideal systems) suggest themselves for frustrating a statistical analysis. These we may call the methods of *diffusion* and *confusion*. In the method of diffusion the statistical structure of $M$ which leads to its redundancy is "dissipated" into long range sta-

tistics—i.e., into statistical structure involving long combinations of letters in the cryptogram. The effect here is that the enemy must intercept a tremendous amount of material to tie down this structure, since the structure is evident only in blocks of very small individual probability. Furthermore, even when he has sufficient material, the analytical work required is much greater since the redundancy has been diffused over a large number of individual statistics. An example of diffusion of statistics is operating on a message $M = m_1, m_2, m_3, \cdots$ with an "averaging" operation, e.g.

$$y_n = \sum_{i=1}^{s} m_{n+i} \ (\text{mod } 26),$$

adding $s$ successive letters of the message to get a letter $y_n$. One can show that the redundacy of the $y$ sequence is the same as that of the $m$ sequence, but the structure has been dissipated. Thus the letter frequencies in $y$ will be more nearly equal than in $m$, the digram frequencies also more nearly equal, etc. Indeed any reversible operation which produces one letter out for each letter in and does not have an infinite "memory" has an output with the same redundancy as the input. The statistics can never be eliminated without compression, but they can be spread out.

The method of *confusion* is to make the relation between the simple statistics of $E$ and the simple description of $K$ a very complex and involved one. In the case of simple substitution, it is easy to describe the limitation of $K$ imposed by the letter frequencies of $E$. If the connection is very involved and confused the enemy may still be able to evaluate a statistic $S_1$, say, which limits the key to a region of the key space. This limitation, however, is to some complex region $R$ in the space, perhaps "folded over" many times, and he has a difficult time making use of it. A second statistic $S_2$ limits $K$ still further to $R_2$, hence it lies in the intersection region; but this does not help much because it is so difficult to determine just what the intersection is.

To be more precise let us suppose the key space has certain "natural co-ordinates" $k_1, k_2, \cdots, k_p$ which he wishes to determine. He measures, let us say, a set of statistics $s_1, s_2, \cdots, s_n$ and these are sufficient to determine the $k_i$. However, in the method of confusion, the equations connecting these sets of variables are involved and complex. We have, say,

$$f_1(k_1, k_2, \cdots, k_p) = s_1$$
$$f_2(k_1, k_2, \cdots, k_p) = s_2$$
$$\vdots$$
$$f_n(k_1, k_2, \cdots, k_p) = s_n,$$

and all the $f_i$ involve all the $k_i$. The cryptographer must solve this system simultaneously—a difficult job. In the simple (not confused) cases the functions involve only a small number of the $k_i$—or at least some of these do. One first solves the simpler equations, evaluating some of the $k_i$ and substitutes these in the more complicated equations.

The conclusion here is that for a good ciphering system steps should be taken either to diffuse or confuse the redundancy (or both).

## 24. THE PROBABLE WORD METHOD

One of the most powerful tools for breaking ciphers is the use of probable words. The probable words may be words or phrases expected in the particular message due to its source, or they may merely be common words or syllables which occur in any text in the language, such as *the, and, tion, that,* and the like in English.

In general, the probable word method is used as follows: Assuming a probable word to be at some point in the clear, the key or a part of the key is determined. This is used to decipher other parts of the cryptogram and provide a consistency test. If the other parts come out in the clear, the assumption is justified.

There are few of the classical type ciphers that use a small key and can resist long under a probable word analysis. From a consideration of this method we can frame a test of ciphers which might be called the acid test. It applies only to ciphers with a small key (less than, say, 50 decimal digits), applied to natural languages, and not using the ideal method of gaining secrecy. The acid test is this: How difficult is it to determine the key or a part of the key knowing a small sample of message and corresponding cryptogram? Any system in which this is easy cannot be very resistant, for the cryptanalyst can always make use of probable words, combined with trial and error, until a consistent solution is obtained.

The conditions on the size of the key make the amount of trial and error small, and the condition about ideal systems is necessary, since these automatically give consistency checks. The existence of probable words and phrases is implied by the assumption of natural languages.

Note that the requirement of difficult solution under these conditions is not, by itself, contradictory to the requirements that enciphering and deciphering be simple processes. Using functional notation we have for enciphering

$$E = f(K, M)$$

and for deciphering

$$M = g(K, E).$$

Both of these may be simple operations on their arguments without the third equation

$$K = h(M, E)$$

being simple.

We may also point out that in investigating a new type of ciphering system one of the best methods of attack is to consider how the key could be determined if a sufficient amount of $M$ and $E$ were given.

The principle of confusion can be (and must be) used to create difficulties for the cryptanalyst using probable word techniques. Given (or assuming) $M = m_1, m_2, \cdots, m_s$ and $E = e_1, e_2, \cdots, e_s$ the cryptanalyst can set up equations for the different key elements $k_1, k_2, \cdots, k_r$ (namely the enciphering equations).

$$e_1 = f_1(m_1, m_2, \cdots, m_s ; k_1, \cdots, k_r)$$
$$e_2 = f_2(m_1, m_2, \cdots, m_s ; k_1, \cdots, k_r)$$

$$\cdot$$
$$\cdot$$
$$\cdot$$

$$e_s = f_s(m_1, m_2, \cdots, m_s ; k_1, \cdots, k_r)$$

All is known, we assume, except the $k_i$. Each of these equations should therefore be complex in the $k_i$, and involve many of them. Otherwise the enemy can solve the simple ones and then the more complex ones by substitution.

From the point of view of increasing confusion, it is desirable to have the $f_i$ involve several $m_i$, especially if these are not adjacent and hence less correlated. This introduces the undesirable feature of error propagation, however, for then each $e_i$ will generally affect several $m_i$ in deciphering, and an error will spread to all these.

We conclude that much of the key should be used in an involved manner in obtaining any cryptogram letter from the message to keep the work characteristic high. Further a dependence on several uncorrelated $m_i$ is desirable, if some propagation of error can be tolerated. We are led by all three of the arguments of these sections to consider "mixing transformations."

## 25. MIXING TRANSFORMATIONS

A notion that has proved valuable in certain branches of probability theory is the concept of a *mixing transformation*. Suppose we have a probability or measure space $\Omega$ and a measure preserving transformation $F$ of the space into itself, that is, a transformation such that the measure of a

transformed region $FR$ is equal to the measure of the initial region $R$. The transformation is called mixing if for any function defined over the space and any region $R$ the integral of the function over the region $F^n R$ approaches, as $n \to \infty$, the integral of the function over the entire space $\Omega$ multiplied by the volume of $R$. This means that any initial region $R$ is mixed with uniform density throughout the entire space if $F$ is applied a large number of times. In general, $F^n R$ becomes a region consisting of a large number of thin filaments spread throughout $\Omega$. As $n$ increases the filaments become finer and their density more constant.

A mixing transformation in this precise sense can occur only in a space with an infinite number of points, for in a finite point space the transformation must be periodic. Speaking loosely, however, we can think of a mixing transformation as one which distributes any reasonably cohesive region in the space fairly uniformly over the entire space. If the first region could be described in simple terms, the second would require very complex ones.

In cryptography we can think of all the possible messages of length $N$ as the space $\Omega$ and the high probability messages as the region $R$. This latter group has a certain fairly simple statistical structure. If a mixing transformation were applied, the high probability messages would be scattered evenly throughout the space.

Good mixing transformations are often formed by repeated products of two simple non-commuting operations. Hopf[12] has shown, for example, that pastry dough can be mixed by such a sequence of operations. The dough is first rolled out into a thin slab, then folded over, then rolled, and then folded again, etc.

In a good mixing transformation of a space with natural coordinates $X_1$, $X_2$, $\cdots$, $X_S$ the point $X_i$ is carried by the transformation into a point $X_i'$, with

$$X_i' = f_1(X_1, X_2, \cdots, X_S) \quad i = 1, 2, \cdots, S$$

and the functions $f_i$ are complicated, involving all the variables in a "sensitive" way. A small variation of any one, $X_3$, say, changes all the $X_i'$ considerably. If $X_3$ passes through its range of possible variation the point $X_i'$ traces a long winding path around the space.

Various methods of mixing applicable to statistical sequences of the type found in natural languages can be devised. One which looks fairly good is to follow a preliminary transposition by a sequence of alternating substitutions and simple linear operations, adding adjacent letters mod 26 for example. Thus we might take

12 E. Hopf, "On Causality, Statistics and Probability," *Journal of Math. and Physics*, v. 13, pp. 51–102, 1934.

$$F = LSLSLT$$

where $T$ is a transposition, $L$ is a linear operation, and $S$ is a substitution.

## 26. Ciphers of the Type $T_k F S_j$

Suppose that $F$ is a good mixing transformation that can be applied to sequences of letters, and that $T_k$ and $S_j$ are any two simple families of transformations, i.e., two simple ciphers, which may be the same. For concreteness we may think of them as both simple substitutions.

It appears that the cipher $TFS$ will be a very good secrecy system from the standpoint of its work characteristic. In the first place it is clear on reviewing our arguments about statistical methods that no simple statistics will give information about the key—any significant statistics derived from $E$ must be of a highly involved and very sensitive type—the redundancy has been both diffused and confused by the mixing transformation $F$. Also probable words lead to a complex system of equations involving all parts of the key (when the mix is good), which must be solved simultaneously.

It is interesting to note that if the cipher $T$ is omitted the remaining system is similar to $S$ and thus no stronger. The enemy merely "unmixes" the cryptogram by application of $F^{-1}$ and then solves. If $S$ is omitted the remaining system is much stronger than $T$ alone when the mix is good, but still not comparable to $TFS$.

The basic principle here of simple ciphers separated by a mixing transformation can of course be extended. For example one could use

$$T_k F_1 S_j F_2 R_i$$

with two mixes and three simple ciphers. One can also simplify by using the same ciphers, and even the same keys as well as the same mixing transformations. This might well simplify the mechanization of such systems.

The mixing transformation which separates the two (or more) appearances of the key acts as a kind of barrier for the enemy—it is easy to carry a known element over this barrier but an unknown (the key) does not go easily.

By supplying two sets of unknowns, the key for $S$ and the key for $T$, and separating them by the mixing transformation $F$ we have "entangled" the unknowns together in a way that makes solution very difficult.

Although systems constructed on this principle would be extremely safe they possess one grave disadvantage. If the mix is good then the propagation of errors is bad. A transmission error of one letter will affect several letters on deciphering.

## 27. Incompatibility of the Criteria for Good Systems

The five criteria for good secrecy systems given in section 5 appear to have a certain incompatibility when applied to a natural language with its complicated statistical structure. With artificial languages having a simple statistical structure it is possible to satisfy all requirements simultaneously, by means of the ideal type ciphers. In natural languages a compromise must be made and the valuations balanced against one another with a view toward the particular application.

If any one of the five criteria is dropped, the other four can be satisfied fairly well, as the following examples show:

1. If we omit the first requirement (amount of secrecy) any simple cipher such as simple substitution will do. In the extreme case of omitting this condition completely, no cipher at all is required and one sends the clear!
2. If the size of the key is not limited the Vernam system can be used.
3. If complexity of operation is not limited, various extremely complicated types of enciphering process can be used.
4. If we omit the propagation of error condition, systems of the type *TFS* would be very good, although somewhat complicated.
5. If we allow large expansion of message, various systems are easily devised where the "correct" message is mixed with many "incorrect" ones (misinformation). The key determines which of these is correct.

A very rough argument for the incompatibility of the five conditions may be given as follows: From condition 5, secrecy systems essentially as studied in this paper must be used; i.e., no great use of nulls, etc. Perfect and ideal systems are excluded by condition 2 and by 3 and 4, respectively. The high secrecy required by 1 must then come from a high work characteristic, not from a high equivocation characteristic. If the key is small, the system simple, and the errors do not propagate, probable word methods will generally solve the system fairly easily, since we then have a fairly simple system of equations for the key.

This reasoning is too vague to be conclusive, but the general idea seems quite reasonable. Perhaps if the various criteria could be given quantitative significance, some sort of an exchange equation could be found involving them and giving the best physically compatible sets of values. The two most difficult to measure numerically are the complexity of operations, and the complexity of statistical structure of the language.

## APPENDIX

*Proof of Theorem 3*

Select any message $M_1$ and group together all cryptograms that can be obtained from $M_1$ by any enciphering operation $T_i$. Let this class of crypto-

grams be $C_1'$. Group with $M_1$ all messages that can be obtained from $M_1$ by $T_i^{-1}T_jM_1$, and call this class $C_1$. The same $C_1'$ would be obtained if we started with any other $M$ in $C_1$ since

$$T_sT_j^{-1}T_iM_1 = T_lM_1.$$

Similarly the same $C_1$ would be obtained.

Choosing an $M$ not in $C_1$ (if any such exist) we construct $C_2$ and $C_2'$ in the same way. Continuing in this manner we obtain the residue classes with properties (1) and (2). Let $M_1$ and $M_2$ be in $C_1$ and suppose

$$M_2 = T_1T_2^{-1}M_1.$$

If $E_1$ is in $C_1'$ and can be obtained from $M_1$ by

$$E_1 = T_\alpha M_1 = T_\beta M_1 = \cdots = T_\eta M_1,$$

then

$$E_1 = T_\alpha T_2^{-1}T_1M_2 = T_\beta T_2^{-1}T_1M_2 = \cdots$$
$$= T_\lambda M_2 = T_\mu M_2 \cdots$$

Thus each $M_i$ in $C_1$ transforms into $E_1$ by the same number of keys. Similarly each $E_i$ in $C_1'$ is obtained from any $M$ in $C_1$ by the same number of keys. It follows that this number of keys is a divisor of the total number of keys and hence we have properties (3) and (4).

# The Design of Reactive Equalizers*

## By A. P. BROGLE, Jr.

This paper describes a systematic method of approximating with a finite number of network elements a transfer characteristic which is a prescribed function of frequency, rather than a constant, over the useful frequency band. Although applied here only to input and output coupling networks as reactive equalizers and where loss equalization to an extremely high degree of precision over a wide frequency band is desired, the mathematical expressions which form the basis for the design are applicable to any 4-terminal network whose transfer characteristic is specified in a similar manner over the real frequency range.

The selection of the appropriate form of the transfer function for equalization purposes is the fundamental consideration. A squared Tchebycheff polynomial is found to be particularly suitable to produce a desired cut-off characteristic without impairing the precision of equalization in the useful band.

A method of polynomial approximation based on the transformation $\omega = \tan \varphi/2$ is used to obtain the coefficients of the in-band approximating function. Predistorting the transfer specification and minimizing the mean-square error, the coefficients become the Fourier cosine coefficients for an infinite frequency range; and are the solutions of a linear set for a finite range, $o \leq \varphi \leq \pi/2$.

## 1. INTRODUCTION

IN MOST broad-band communication systems, the problems of loss equalization and distortion correction are fundamental. Of the various types of electrical networks which are found useful as equalizers and compensators, the most frequently employed are the so-called constant resistance networks. In particular, they are of three usual types, as indicated in Fig. 1.

In all cases, the relationship $Z_1 Z_2 = R^2$, which is always possible to fulfill if $Z_1$ and $Z_2$ are built up of resistive and reactive components in the well-known manner, provides the means of altering the transmission properties of the circuit without affecting its impedance.[1] Methods are also available which extend the problem to more complicated configurations having these constant resistance properties. However, in some applications, where signal-to-noise ratio considerations are of importance, the resistive elements included as components of $Z_1$ and $Z_2$ in these circuits place a limitation on the final performance of the system. Hence, the satisfactory transmission and *impedance matching* properties of these circuits are purchased at the expense of a substantially increased noise level. As a consequence of this limitation on the performance of standard constant resistance equalizers, recent work

---

* The work presented in this paper is part of a thesis, "Design of Reactive Equalizers with Prescribed Parasitic Capacitance," submitted by the author in partial fulfillment of the requirements for the degree of Master of Science at the Massachusetts Institute of Technology (Feb. 1949).

[1] Ref. 5, pp. 1–2.

has indicated the advantage of adapting reactive input and output coupling networks, ordinarily employed solely as impedance matching devices, to the additional role of partial distortion equalization.[2]

As a reactive equalizer, a lossless input or output coupling network partially equalizes the loss characteristic of a transmission line or cable by providing an insertion gain characteristic to compensate for the line loss characteristic. However, before the rigorous formulation of the problem is undertaken in the following section, it is necessary to discuss briefly the role of input and output coupling networks as equalizers in communications
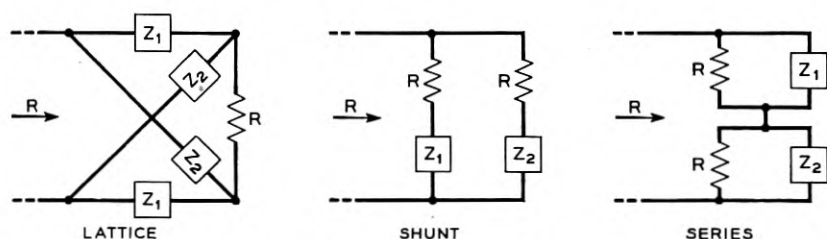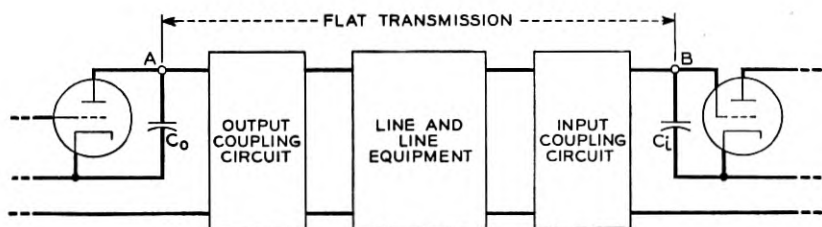
Fig. 1—Constant resistance networks.

Fig. 2—Simplified section of a broad-band transmission system.

systems, and to outline the external requirements and limitations imposed by the system itself on these networks.

The characteristics of input and output coupling networks which are of engineering interest are:

(1) The contribution of the coupling circuits to the transmission performance of the system as a whole.

(2) The impedance matching requirements between the coupling networks and the transmission line.

(3) The limitation on the maximum performance of a coupling network imposed by the parasitic capacitance usually present in the termination.

These characteristics are perhaps best illustrated by a somewhat idealized section of a broad-band transmission system. Figure 2 represents the output

[2] Ref. 1, pp. 383–392.

stage of a repeater, a section of the associated transmission line, and the first stage of the succeeding repeater of a simplified system.

The specification of a flat transmission characteristic over the useful frequency band between A and B in the figure indicates that equalization for the line loss of the section must occur in either or both coupling circuits, in the line equipment, or in all three of these circuits. For feedback amplifiers, the most desirable type, a flat characteristic between A and B can be specified only if the feedback circuits, or $\beta$ circuits, of the amplifiers are designed to have no transmission variation with frequency. In general, it is possible to suppose the feedback factor, $\beta$, of the amplifiers to be the appropriately varying function of frequency to equalize a part of the line loss, thus altering the transmission specification from A to B. However, the $\beta$ circuits must include regulation of other types in most cases. Hence, it is impractical to include much loss equalization in these circuits.

Since satisfactory performance of the section is dependent also on the maintenance of a large signal-to-noise ratio, it is important that the line contain no sources of additional loss. It is clear, then, that the best transmission performance is obtained (1) without the use of equalization in the line[3] and (2) when the reactive input and output coupling circuits equalize as large a percentage as possible of the total line loss.

Physically, the coupling circuits will be transformers, plus any number of tuning and shaping elements. In addition to the primary function of metallically separating the line from the repeater amplifiers, it will be seen later that the transformers provide the means of adjusting, independent of the value of the prescribed line impedance, the final impedance level of the network to conform with the value of the parasitic capacitance present.

Besides the contribution of the various networks in the system to the overall transmission performance, there is the problem of matching the coupling circuits to the line. For constant-resistance equalization, this problem is immediately solved by the relationship $Z_1Z_2 = R^2$. Well-established techniques make it a relatively simple matter to design for a specified attenuation variation with frequency at the same time that the impedance of the equalizer is matched to the line. This same procedure, with certain modifications, can be carried over to the design of reactive equalizers. In Fig. 2, the transformers of the input and output coupling circuits are unterminated. That is, the input of the output circuit and the output of the input circuit are terminated in substantially open circuits. In order to prevent the reflection of power at the junctions of the coupling circuits and the line, the impedances of the input and output circuits as viewed from the line must be made equal to the impedance of the line. This impedance re-

---

[3] In practice, the $\beta$ circuits and constant resistance networks associated with the line actually equalize a certain percentage of the total line loss characteristic.

quirement is fulfilled by providing both coupling circuits with a balancing network connected as shown in Fig. 3. By accepting a small constant transmission loss,[4] the relationship $Z_1Z_2 = R^2$ is satisfied if the impedance $Z_2$ of the balancing network is made the inverse of the transmission circuit impedance $Z_1$. Because of the relative ease of designing an inverse impedance $Z_2$, once $Z_1$ is known in the final stages of a particular design, it is appropriate to omit from further discussion the presence of the balancing networks.

The fundamental theoretical limitation in the maximum transmission performance of these coupling networks is due directly to the presence of the parasitic tube capacitances $C_0$ and $C_i$. If the parasitic capacitances were not present, the turns ratios of the transformers in the coupling circuits could quite evidently be made extremely high in order to produce over any specified frequency band as large a transmission response as desired. However, even though these capacitances are usually small, they always tend to short circuit the coupling networks whenever the impedance ratios of the
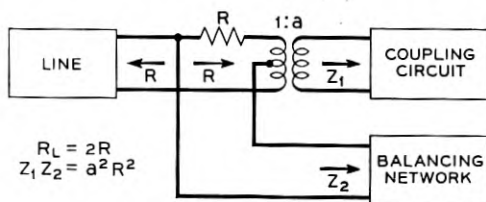


Fig. 3—Balancing network arrangement.

transformers are made too high. The determination of the maximum response of these networks over a prescribed frequency range is thus a basic problem in the design of reactive equalizers.

The fundamental limitation on the response of these networks is expressed in terms of the total area available under the transfer characteristic.[5] When this characteristic is a desired function over a finite frequency band, the maximum utilization of the area available is obviously attained when all the area is included in the useful band. This condition is described as a *resistance efficiency* of 100 per cent. A smaller resistance efficiency, 75 per cent for example, means that three-fourths of the total area under the characteristic is available in the useful frequency region, while the remainder of the area may be utilized to decrease the rate at which the characteristic is *cut-off*. Hence, the realization of a prescribed resistance efficiency in the

---

[4] The effective impedance of the line as viewed from the coupling circuit is equal to twice the actual line impedance. Thus, a penalty of $10 \log \dfrac{R_L}{R} = 3db$ is imposed by the presence of the balancing network.

[5] See eq. (4) and discussion in the following section.

design of a reactive equalizer places a definite requirement on the behavior of the transfer characteristic outside the useful frequency band.

Although the precision of equalization as a design requirement actually is inclusive in the term *transmission performance* as used previously, it is included here as a separate requirement to emphasize its importance in this problem. The specification of a flat transmission from A to B in Fig. 2 provides the means of assigning to the tolerance of equalization a quantitative meaning. Hence, the tolerance per repeater section of the system may be expressed as the maximum allowable db deviation from the flat transmission characteristic, A to B, over the useful frequency band. For extremely broad-band systems, such as a coaxial system for simultaneous long-distance telephone and television transmission, many repeater sections appear in tandem between terminals. Thus, the deviations in each of these sections contribute to the system as a whole. In addition to the distances usually involved, repeater spacing becomes closer as the effective transmission band of these systems is increased. In order to design new systems with increasingly better overall tolerances, at the same time that the broad-banding requirements call for a greatly increased number of repeater sections per system, the tolerances imposed on the individual sections become exceedingly small. As a consequence, the maximum tolerance for an individual section must be specified as perhaps less than $\pm 0.05$ db deviation.

## 2. The Problem of Reactive Equalization

In this section the problem of reactive equalization will be formulated in terms of the special problems of input and output coupling circuit design. Broadly speaking, the general characteristics of input and output coupling networks, as outlined in the introduction to establish the practical basis for reactive equalization, will be further developed in order to give them a quantitative meaning. Because of the complexity of some derivations and their extensive treatment elsewhere, detailed proofs in general will be merely outlined. The method of analysis follows Bode's treatment of the problem while the principal results taken from network theory are Guillemin's.

As previously stated, the unterminated case for input and output coupling circuits arises whenever the terminating resistance is infinite in comparison with the other impedances of the network.[6] Figures 4 and 5 represent, respectively, an output and an input coupling network of the type illustrated in Fig. 2 with infinite terminations. In each figure, $R_L$ represents the line, $N$ is the lossless coupling network, and $C_n$ is the parasitic shunt capacitance

---

[6] The so-called *terminated case* exists when the parasitic capacitance $C_0$ or $C_i$ in Fig. 2 is shunted by a finite resistance. Since no essential differences exist between the two cases with respect to the approximation problem, an analysis for the unterminated case alone is sufficient to clarify the more important design considerations.

which limits the response over any specified frequency band. For purposes of analysis and design, it is convenient to represent the coupling transformers in the manner indicated. By adopting this equivalent representation of a physical transformer, the so-called high-side equivalent circuit of the transformer, which includes the leakage reactance, the magnetizing inductance, and the input and output winding capacitances, is incorporated as part of the coupling network itself.

By excluding the ideal transformer portion of the equivalent representation of the physical transformer from the network itself, a simplification is possible. As shown in Figs. 6 and 7, the combination of the resistance $R_L$
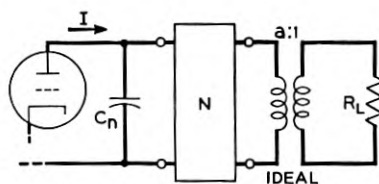


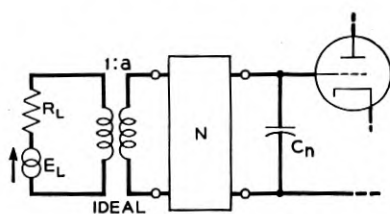Fig. 4—Output coupling circuit.



Fig. 5—Input coupling circuit.

and the ideal transformer may, in each case, be replaced by a resistance $R_0 = a^2 R_L$, where "$a$" is the step-up turns ratio of the ideal transformer. $R_L$ is the specified resistance, and $R_0$ and "$a$" are determined in the design procedure from the maximum response obtainable with the prescribed capacitance $C_n$ in the termination.

The starting point for the study of these circuits is a consideration of the limitation on the amplitude response of these networks with frequency due to the presence of $C_n$ in the terminations. Since the current ratio $\dfrac{I_L}{I}$ in Fig. 6 and the voltage ratio $\dfrac{E}{E_L}$ in Fig. 7 might be as large as desired if it were not for the presence of $C_n$, the immediate problem is that of relating the magnitude of these ratios, as functions of the real frequency, to the capacitance $C_n$. This relationship is dependent on a necessary condition for the physical

realizability of a driving-point impedance function. If this function is chosen as the $Z = R + jX$ in the figures, the necessary condition of interest is that $Z$, as an analytic function, have no poles in the right half of the complex frequency plane and that $Z$ approach $\dfrac{1}{\omega C_n}$ as $\omega$ approaches infinity. By integrating this function over the appropriate path in the right half of the $\lambda$ (complex frequency) plane and setting the result equal to zero, the desired expression becomes

$$\int_0^\infty R \, d\omega = \frac{\pi}{2C_n}.^7 \tag{1}$$

To show that the resistance $R$ is related to the ratios $\left|\dfrac{I_L}{I}\right|$ and $\left|\dfrac{E}{E_L}\right|$ it is
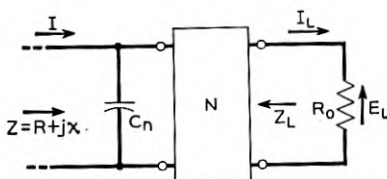


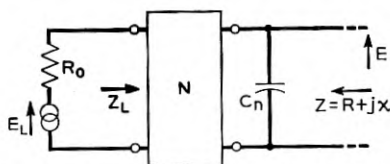Fig. 6—Modified output coupling circuit of Fig. 4.



Fig. 7—Modified input coupling circuit of Fig. 5.

necessary to examine the transfer of power through the output circuit of Fig. 6. The power driven into this circuit is $|I|^2 R$. Since the network $N$ is lossless, this is the same power, $|I_L|^2 R_0$, which reaches the line. In addition, if the transfer impedance of the circuit is defined as $Z_{12}(j\omega) = \dfrac{E_L}{I} = R_0 \dfrac{I_L}{I}$, the relationship sought is

$$\left|\frac{I_L}{I}\right|^2 = \left|\frac{Z_{12}(j\omega)}{R_0}\right|^2 = \frac{R}{R_0}. \tag{2}$$

For the input coupling circuit, the ratio $\left|\dfrac{E}{E_L}\right|$ is related to the transfer impedance and $R$ in a similar manner.

⁷ Ref. 1, pp. 278–281.

$$\left| \frac{E}{E_L} \right|^2 = \left| \frac{Z_{12}(j\omega)}{R_0} \right|^2 = \frac{R}{R_0}. \quad ^8 \qquad (3)$$

Finally, the transmission gain $\alpha$ (in nepers) is related to the current ratio $\left| \frac{I_L}{I} \right|$, or the voltage ratio $\left| \frac{E}{E_L} \right|$, by $e^{\alpha}$. Hence, the quantitative statement for the limitation on the response of these coupling circuits becomes

$$\int_0^\infty e^{2\alpha}\, d\omega = \int_0^\infty \left| \frac{Z_{12}(j\omega)}{R_0} \right|^2 d\omega = \frac{\pi}{2C_n R_0}. \qquad (4)$$

Equation (4) is the general formula which relates the response character-istic over the complete frequency range to the prescribed capacitance $C_n$ and the resistance $R_0$. This formula is especially helpful in attaching an analytical meaning to the term partial reactive equalization. If $\alpha' = f(\omega)$ is used to describe the attenuation characteristic of a line or cable over a specified finite frequency band, $\alpha = kf(\omega)$ will be the transmission response, in nepers, which is required to equalize a stated fraction of this loss at every frequency in the specified range. $k$ is then the constant ($k \leq 1$) which numerically expresses the degree of equalization.[9]

Thus, the $\alpha = kf(\omega)$ in eq. (4) is the desired insertion gain characteristic to compensate partially for the line loss characteristic, and is directly related to this loss over a specified frequency range by a constant $k$. The limitation on the response expressed by eq. (4) will be clear if the transmission $\alpha$ is now defined as $\alpha = \alpha_0 + kf(\omega)$, where $\alpha_0$ represents the general response level. Before this expression is substituted in eq. (4), however, it is necessary to change the limits of integration. Thus, the specification of a maximum re-sponse over a finite frequency band requires that the limits become $\omega_1$ and $\omega_2$, the extreme frequencies of the useful band. Since $R$ must be positive, this condition requires that $e^{2\alpha}$ be zero everywhere outside the useful range. Carrying out the integration, the result becomes

$$\alpha_0 \leq \tfrac{1}{2} \ln \left[ \frac{\pi}{2C_n R_0 \displaystyle\int_{\omega_1}^{\omega_2} e^{2kf(\omega)}\, d\omega} \right]. \qquad (5)$$

Since $kf(\omega)$ is always prescribed, $\alpha_0$ is readily computed.

So far, the equations have considered only the ideal case when the transfer characteristic $e^{2\alpha}$ is zero outside the useful band. As previously stated, this condition specifies a resistance efficiency of 100 per cent. In practical appli-cations, where a finite number of network elements are employed to approxi-

---

[8] By (1) substituting the equivalent current source for $E$, (2) applying the principle of reciprocity to the input circuit, and (3) writing the relations for the transfer of power through the circuit, eq. (3) is readily derived.

[9] In practice, this constant is called the "slope" of equalization.

mate a transfer characteristic to a specified degree of precision over the useful band, it is not possible for the transfer function chosen to represent the transfer characteristic to approximate zero outside the useful band in a manner to produce a resistance efficiency of 100 per cent. This limitation is then the prerequisite for modifying the performance which the coupling networks are required to achieve. The usual range of resistance efficiencies specified for input and output coupling network applications is approximately 45 to 80 per cent.

This modification of the final performance of the coupling networks may be examined quantitatively by referring to eqs. (1), (4), and (5). In the first two of these equations the integral may be taken only over the useful frequency range, $\omega_1$ to $\omega_2$, provided that the right-hand side of each of these equations is multiplied by the specified resistance efficiency expressed as a fraction.[10] In eq. (5) the equal sign holds only in the limiting case when the resistance efficiency is 100 per cent. If these equations are modified in the manner indicated, the variation of the transfer characteristic outside the useful frequency range may be chosen in any way which satisfies the total area requirements in eqs. (1) and (4) as they stand.

Following the choice of a satisfactory transfer characteristic, the next general problem is the realization of a physical network which will approximate this specified characteristic to the required degree of precision over the complete frequency spectrum. The solution of this problem is the main purpose of this paper.

As is well-known in network theory, the general form of the squared magnitude of the transfer impedance of any physical two-terminal-pair reactive network terminated in resistance may be expressed as the quotient of two polynomials in $\omega^2$.

$$\left| \frac{Z_{12}(j\omega)}{R_0} \right|^2 = \frac{A_0 + A_1 \omega^2 + A_2 \omega^4 + \cdots + A_n \omega^{2n}}{B_0 + B_1 \omega^2 + B_2 \omega^4 + \cdots + B_n \omega^{2n}}. \tag{6}$$

Before the necessary and sufficient conditions that the $\dfrac{Z_{12}(\lambda)}{R_0}$ derived from eq. (6) be the transfer impedance of a lossless network terminated in resistance are stated, it is appropriate to develop the modifications which must be made in eq. (6) if $\left| \dfrac{Z_{12}(j\omega)}{R_0} \right|^2$ is to approximate the transfer characteristic, $e^{2\alpha}$, in this problem. This requires that a closer examination be made of the physical limitation that the coupling networks correspond, in part, in structure to the equivalent circuit of the coupling transformer to be used. Figure 8 shows the high-side equivalent circuit of either coupling transformer of Figs. 4 and 5.

---

[10] $\omega_1$ is usually chosen as zero.

In the figure, $L_m$ represents the magnetizing inductance, $L_2$ represents the leakage reactance, and $C_1$ and $C_3$ represent, respectively, the low-side and high-side parasitic winding capacitances. The magnetizing inductance $L_m$, since it is usually large so that its impedance is substantially infinite compared with the other impedances of the circuit at high frequencies, affects the response of the transformer at low frequencies only. Since the useful band ordinarily specified does not include the range of frequencies where the effects of $L_m$ are noticeable, its presence may be omitted from further consideration. In addition, it is never practical to retain $C_3$ as the final element of the reactive coupling network $N$. In this case, the parallel combination of $C_3$ and $C_n$ would, of course, seriously limit the final response of the network. Thus, the least number of shaping elements is a series inductance $L_4$ which *splits* the high-side winding capacitance $C_3$ from the prescribed terminating capacitance $C_n$. Hence, in general, the reactive coupling network $N$ is an $(n - 1)$ element unbalanced ladder structure of alternating series inductances and shunt capacitances beginning with a shunt capacitance
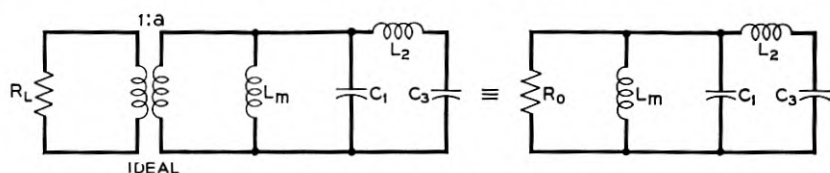


Fig. 8—High-side equivalent circuit of either coupling transformer of Figs. 4 and 5.

and ending with a series inductance. Figure 9, then, indicates the general form of the coupling network to be realized by the function chosen to approximate $e^{2\alpha}$ in this problem.

Without loss of generality, it is convenient at this point to modify Figs. 6 and 7 in the manner indicated in Figs. 10 and 11. By including $C_n$ as part of $N'$ the problem has not been altered. However, it is necessary to recognize that the final adjustment of the impedance level, i.e., the choice of $R_0$, must be made in such a manner that the total area requirement, as specified in eq. (4), is still met. In each figure $z'_{11}$, $z'_{22}$, and $z'_{12}$ are the open-circuit driving-point and transfer impedances of the network $N'$.

With the element configuration specified and the reactive coupling network $N'$ defined, it is now appropriate to carry out the modification in the form of $\left| \dfrac{Z_{12}(j\omega)}{R_0} \right|^2$ indicated previously. Thus, the fact that $\dfrac{R}{R_0} = 1$ at $\omega = 0$, and that an $n$ element unbalanced ladder structure of alternating series inductances and shunt capacitances terminated in a resistance has only an nth order zero of the transfer impedance, $\dfrac{Z_{12}(\lambda)}{R_0}$, at infinity, allows the

squared magnitude of the transfer impedance in this problem to be written as

$$\left|\frac{Z_{12}(j\omega)}{R_0}\right|^2 = \frac{1}{1 + B_1\omega^2 + B_2\omega^4 + \cdots + B_n\omega^{2n}}, \qquad (7)$$

where the $n$ constants $B_1 \cdots B_n$ are related to the $n$ elements of the network by the relation

$$\frac{Z_{12}(j\omega)}{R_0} = \frac{z_{12}'/R_0}{1 + z_{22}'/R_0}. \qquad (8)$$

Since the desired transfer characteristic $e^{2\alpha}$ determines the variation of the polynomial $B(\omega^2) = 1 + B_1\omega^2 + \cdots + B_n\omega^{2n}$, a major factor in the design
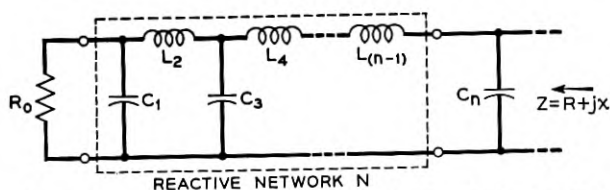


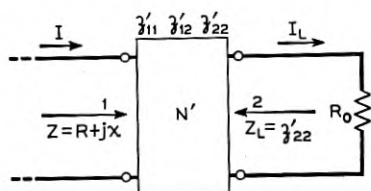Fig. 9—General form of the coupling networks of Figs. 6 and 7.



Fig. 10—Output circuit of Fig. 6 with $C_n$ included as part of $N'$.



Fig. 11—Input circuit of Fig. 7 with $C_n$ included as part of $N'$.

is the choice of the real coefficients, $B_1 \cdots B_n$, by a suitable method of polynomial approximation.

The necessary and sufficient conditions for physical realizability place a restriction on the $B$'s of eq. (7). The sufficient condition that $\left|\dfrac{Z_{12}(j\omega)}{R_0}\right|^2$ represent the squared magnitude of the transfer impedance of a physical

network of the type described is that $\left| \dfrac{Z_{12}(j\omega)}{R_0} \right|^2 \geq 0$ for $\omega \geq 0$. This condition will be insured if the polynomial, $1 + B_1\omega^2 + \cdots + B_n\omega^{2n}$, has no negative real $\lambda^2$ roots of odd multiplicity.[11] In addition to the sufficiency of eq. (7), if the $\dfrac{Z_{12}(\lambda)}{R_0} = \dfrac{g(\lambda)}{h(\lambda)}$ derived from $\left| \dfrac{Z_{12}(\lambda)}{R_0} \right|^2$ in the usual manner is to be the transfer impedance of a lossless network terminated in resistance, it is necessary that $g(\lambda)$ be either even or odd and that $h(\lambda)$ be a Hurwitz polynomial.[12] In this problem $g(\lambda) = 1$ is surely even since all zeros of $\left| \dfrac{Z_{12}(\lambda)}{R_0} \right|^2$ occur at infinity; and the method of forming $\dfrac{Z_{12}(\lambda)}{R_0}$ always insures that $h(\lambda) = m + n$, where $m$ is the even part and $n$ is the odd part of $h(\lambda)$, is a Hurwitz polynomial. Thus, the fulfillment of the sufficient condition that there be no negative real $\lambda^2$ roots of odd multiplicity of $B(\omega^2)$ is the assurance that the $B$'s of eq. (7) will always produce a physical network of the configuration of Fig. 9.

Once the conditions for physical realizability have been fulfilled, and a $\dfrac{Z_{12}(\lambda)}{R_0}$ has been found in the final stages of a particular design, the network elements are easily calculated from a partial fraction expansion of $z'_{22} = \dfrac{m}{n}$ according to the following relation:

$$\frac{Z_{12}(\lambda)}{R_0} = \frac{z'_{12}(\lambda)/R_0}{1 + z'_{22}(\lambda)/R_0} = \frac{g(\lambda)}{m + n} = \frac{g(\lambda)/n}{1 + m/n}, \qquad (9)$$

where $z'_{12}(\lambda) = \dfrac{g(\lambda)}{n}$ and $z'_{22}(\lambda) = \dfrac{m}{n}$.

The previous discussion of the special problems of input and output coupling circuit design has been based, broadly, on (1) a consideration of the terminating or load impedance, (2) a consideration of the shape of the transfer characteristic, and (3) a consideration of the conditions for physical realizability. A major problem in the design is the choice of an approximating function which satisfactorily matches the stated transfer characteristic over the useful frequency band and, at the same time, sharply changes slope near the cut-off frequency so that it approximates zero outside the useful band in a prescribed manner. When the transfer characteristic is a constant over the useful frequency band, e.g., the impedance matching and low-pass filter cases, techniques which employ Tchebycheff polynomials as the ap-

---

[11] Ref. 4.
[12] A Hurwitz polynomial is defined as a polynomial in $\lambda$ which has the property that the quotient of its even and odd parts, $\varphi(\lambda) = \dfrac{m}{n}$, yields a reactance function.

proximating functions are available which make it a relatively simple matter to design physically realizable networks exhibiting this property of a sharp cut-off to zero outside the useful band.[13] However, a similar method of applying Tchebycheff polynomials to transfer characteristics which vary with frequency in a prescribed manner over a finite band has not been evolved. In order to illustrate the preceding statements, Figs. 12 and 13 have been included as representative of typical transfer characteristics.
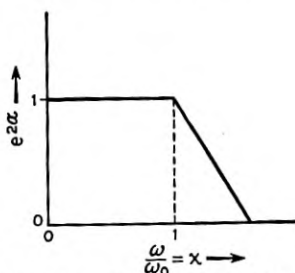


Fig. 12—Transfer characteristic for impedance matching or low-pass filter case.
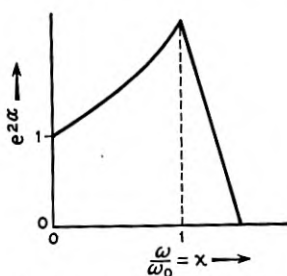


Fig. 13—Transfer characteristic for reactive equalizer case.

## 3. Derivation of Special Transfer Function

In accordance with the brief discussion at the conclusion of the previous chapter, it is now appropriate to state that it is the purpose of this paper (1) to derive a transfer function which is especially suited to the problem of reactive equalization, and (2) to develop a systematic method which utilizes this special transfer function to approximate satisfactorily, with a finite number of network elements, a specified transfer characteristic over the entire frequency spectrum. This section will consider in detail the first of these two main tasks in the formulation of a design method for reactive equalizers.

With reference to Fig. 13, it is convenient to divide the complete transfer

[13] Ref. 4. Also Ref. 2, pp. 53–79.

characteristic into two separate regions. The specification over the useful band, $0 \leq \omega \leq \omega_0$, may be called the in-band region while the specification outside the useful band, $\omega_0 < \omega \leq \infty$, may be called the out-band region. Thus, it is seen that the transfer characteristic over the in-band region depends exclusively on the $\alpha = kf(\omega)$ which is required to equalize a stated fraction of the power loss between repeaters while the transfer characteristic in the out-band region depends only on the specified resistance efficiency.

The first step in the derivation of the special transfer function for equalization purposes is a normalization of the transfer characteristic of Fig. 13 in terms of eq. (7). As indicated in Fig. 14, a constant, $K$, is chosen so that $Ke^{2\alpha}(K < 1)$ is equal to unity at $\dfrac{\omega}{\omega_0} = x = 1$. This choice of the transfer characteristic is convenient since the transfer characteristic is now expressed in a form similar to the familiar form of the transfer characteristic of a low-
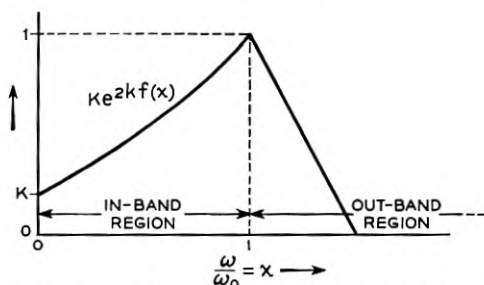


Fig. 14—Normalized transfer characteristic of Fig. 13.

pass filter and, hence, suitable for the addition of a Tchebycheff polynomial.[14]

With the transfer characteristic appropriately specified, the next step is to show the manner in which the denominator $B(x^2)$ of eq. (7), where this equation is multiplied by the factor $K$, can be broken up into two functions of $x^2$ so that one of these functions approximates the reciprocal of the in-band region of the transfer characteristic while the other produces the desired cut-off characteristic.

The derivation of the desired denominator, $B(x^2)$, begins by writing the transfer characteristic of Fig. 14 for the in-band region as

$$\frac{1}{B(x^2)} = Ke^{2kf(x)} \ .^{[15]} \tag{10}$$

[14] In order to make the following derivation clear, it is suggested that the discussion of Tchebycheff polynomials, pp. 733–734, be examined at this time.

[15] The transmission $\alpha = \alpha_0 + kf(x)$ will be written as $kf(x)$ for the remainder of this analysis. The general transmission level $\alpha_0$ may be found in the final stages of a particular design when the impedance level is adjusted to conform with the prescribed $C_n$.

In terms of $B(x^2)$ directly and a desired transmission $\alpha_0'$ at the angular cut-off frequency $\omega_0$, equation (10) becomes

$$B(x^2) = e^{2\alpha_0'} e^{-2kf(x)}, \tag{11}$$

where $K = e^{-2\alpha_0'}$. Equation (11) now represents the characteristic that is to be approximated over the useful frequency band while Fig. 15 shows a plot of this function.



Fig. 15—Specification for $B(x^2)$ over useful frequency band.



Fig. 16—Combined approximating function for $B(x^2)$ over entire frequency band.

Now, if $B(x^2)$ is broken up into two parts and represented as

$$B(x^2) = f(x^2) + \epsilon^2 V_n^2(x),^{16} \tag{12}$$

[16] It is important to note that eq. (12) now represents the approximating function over the entire frequency range as compared to eq. (11) which represents the function to be approximated only over the useful range.

where $f(x^2)$ is the rational function which approximates $e^{2\alpha'_0}e^{-2kf(x)}$ over the useful band, $V_n(x)$ is a Tchebycheff polynomial of order $n$ (odd), and $\epsilon$ is the coefficient of the Tchebycheff polynomial, $B(x^2)$ in Fig. 15 will be modified as shown in Fig. 16. In this figure it 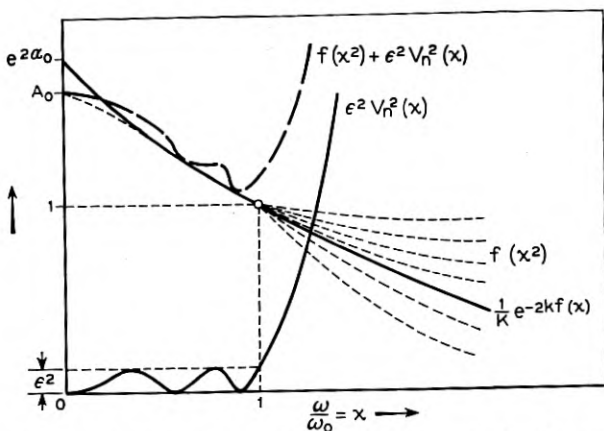is to be noted that $f(x^2)$, the in-band approximating function, is represented as having a variety of variations outside the useful band. The function has been indicated in this manner to emphasize that a fairly wide latitude in the choice of the behavior of $f(x^2)$ outside the useful is permitted since $\epsilon^2 V_n^2(x)$, the out-band approximating function, is the predominant function in this region. In addition, the variations of $\epsilon^2 V_n^2(x)$ in the in-band region have been exaggerated in order to demonstrate their effect on the combined approximating function, $f(x^2) + \epsilon^2 V_n^2(x)$, over the useful frequency band.



Fig. 17—Resultant transfer function for equalization purposes.

Finally, when the relation expressed by eq. (12) is reciprocated and replotted in terms of $K \left| \dfrac{Z_{12}(jx)}{R_0} \right|^2$, the result shown in eq. (13) and Fig. 17 is obtained.

$$K \left| \frac{Z_{12}(jx)}{R_0} \right|^2 = \frac{1}{f(x^2) + \epsilon^2 V_n^2(x)}. \tag{13}$$

Comparing the resultant special transfer function shown in Fig. 17 with the transfer characteristic shown in Fig. 14, and assuming that $f(x^2)$ and the coefficient of the Tchebycheff polynomial have been suitably chosen, it is established contingently that the combination of functions chosen to represent $B(x^2)$ produces the desired result.

This brief derivation serves as a guide to the main problem of choosing a particular $f(x^2)$ and a particular $\epsilon^2 V_n^2(x)$ which, when added together and

reciprocated, approximate the transfer characteristic to the specified degree of precision.

The choice of these approximating functions begins by finding a polynomial

$$f(x^2) = A_0 + A_1 x^2 + A_2 x^4 + \cdots + A_n x^{2n} \tag{14}$$

which approximates $e^{2\alpha_0'} e^{-2kf(x)}$ to the required degree of precision throughout the useful band and has an out-band variation subject to the initial requirements that $f(x^2)$ be positive and that the slope of $f(x^2)$ not vary rapidly in the immediate out-band region (approximately $1 \leq x \leq 1.5$). For values of $x$ greater than about 1.5, the Tchebycheff polynomial is the determining function, and variations in $f(x^2)$ are no longer of importance. A precise statement of these conditions and the exact frequency range in which they are valid depend on the degree of equalization and the desired resistance efficiency in a particular design. However, a more critical examination of Figs. 16 and 17 indicates that the generalized conditions stated above are a reasonable guide in the choice of $f(x^2)$ for most applications.

The main criteria for judging the acceptability of a particular out-band variation which accompanies the choice of in-band variation of $f(x^2)$ to produce optimum precision are physical realizability and the attainment of a desired resistance efficiency. Considering first the condition for physical realizability, $\dfrac{1}{f(x^2) + \epsilon^2 V_n^2(x)} \geq 0$ for $0 \leq x \leq \infty$, and referring to Fig. 16, a negative value of $f(x^2)$ in the immediate out-band region might be of sufficient magnitude to cancel the positive effect of $\epsilon^2 V_n^2(x)$ and, hence, produce a negative value of $f(x^2) + \epsilon^2 V_n^2(x)$. However, at higher frequencies, the squared Tchebycheff polynomial takes on very large positive values. Thus, negative values and variations in $f(x^2)$ are effectively reduced in the magnitude of their effect on

$$K \left| \frac{Z_{12}(jx)}{R_0} \right|^2 = \frac{1}{f(x^2) + \epsilon^2 V_n^2(x)}$$

in direct relation to the increase in the magnitude of $\epsilon^2 V_n^2(x)$.

In order that an accurate prediction of the resistance efficiency may be made, it is necessary that the slope of $f(x^2) + \epsilon^2 V_n^2(x)$ increase in a uniform manner in the immediate out-band region. Since variations in the slope of $f(x^2)$ have their largest effect in the region just outside the useful band, it is, of course, best to prevent rapid variations in this region.

The remaining condition on the form of $f(x^2)$ is that $A_0$ should be adjusted so that $A_0 < e^{2\alpha_0'}$. By providing the transfer specification with a less steep slope requirement at low frequencies it is possible to obtain over the valuable

portion of the useful band an increased precision of equalization.[17] This adjustment represents an increased transmission at low frequencies. Thus, it is sometimes necessary to employ an equalizer of the constant resistance type when additional equalization is desired at low frequencies. Figures 16 and 17 have been drawn to reflect this condition on $A_0$.

After an $f(x^2)$ which conforms with the requirements outlined above has been found, it is necessary to find a

$$\epsilon^2 V_n^2(x) = A_1' x^2 + A_2' x^4 + \cdots + A_n' x^{2n} \qquad (15)$$

which, when added to $f(x^2)$, produces the desired $B(x^2)$. This procedure is greatly facilitated by the known properties of Tchebycheff polynomials:

A Tchebycheff polynomial of order $n$ is defined by

$$V_n(x) = \cos (n \cos^{-1} x). \qquad (16)$$

This function oscillates between plus one and minus one for $|x| < 1$ and approaches $\pm \infty$ for $|x| > 1$. Tabulated below are the expanded analytical expressions for the polynomials for $n = 1$ through $n = 8$.

| | |
|---|---|
| $V_1(x) = x$ | $V_5(x) = 16x^5 - 20x^3 + 5x$ |
| $V_2(x) = 2x^2 - 1$ | $V_6(x) = 32x^6 - 48x^4 + 18x^2 - 1$ |
| $V_3(x) = 4x^3 - 3x$ | $V_7(x) = 64x^7 - 112x^5 + 56x^3 - 7x$ |
| $V_4(x) = 8x^4 - 8x^2 + 1$ | $V_8(x) = 128x^8 - 256x^6 + 160x^4 - 32x^2 + 1$ |

With the help of the recursion formula,

$$x V_n(x) = \tfrac{1}{2}[V_{n+1}(x) + V_{n-1}(x)], \qquad (17)$$

the corresponding expressions for $n > 8$ may be systematically calculated. Figure 18 shows a plot of the Tchebycheff polynomial for $n = 5$.

In the case of low-pass filters[18] and impedance matching networks,[19] Tchebycheff polynomials are often used for the solution of the approximation problem. The function $|Z_{12}(jx)|^2$ in these cases has an oscillatory behavior which approximates unity in the useful band, and has all its zeros at infinity so that the network consists of $n$ elements of an unbalanced ladder structure of alternating series inductances and shunt capacitances. The appropriate function for $|Z_{12}(jx)|^2$ is

$$|Z_{12}(jx)|^2 = \frac{1}{1 + \epsilon^2 V_n^2(x)}, \qquad (18)$$

[17] There is a practical limit to the reduction of $A_0$ below $\epsilon^2 \alpha_0'$. Referring to Figs. 13 and 14, it is apparent that $K = \dfrac{1}{A_0}$. Thus, $A_0$ is a direct measure of the impedance level over the useful band, and must not be made too small if the highest practical level of response is to be attained.

[18] Ref. 2, pp. 53–79.
[19] Ref. 3, pp. 26–34.

where $\epsilon$ is an arbitrary constant. Figure 19 shows the plot of the squared Tchebycheff polynomial, $\epsilon^2 V_n^2(x)$, for the values of $n = 5$, and $\epsilon = 0.5$ and $\epsilon = 0.1$, while Fig. 20 shows a plot of the transfer function expressed in eq. (18).

It is to be noted that the oscillatory behavior with equal maxima and minima of squared Tchebycheff polynomials for values of $x < 1$ and the rapid approach to $+\infty$ for values of $x > 1$ make their use particularly suitable as the solution of the approximation problem for low-pass filters and impedance matching networks. It is now apparent that these same



Fig. 18—Tchebycheff polynomial, $V_n(x)$, for $n = 5$.

properties validate their use as the out-band approximating function for reactive equalizers.[20]

Another useful property of squared Tchebycheff polynomials as approximating functions for low-pass filters and impedance matching networks is the inclusion of the specification of the tolerance as a factor in the transfer function. The allowable db deviation over the useful band is related to $\epsilon$ by

$$\epsilon^2 = e^{2\alpha_p} - 1,$$

where $\alpha_p$ is the maximum pass-band loss in nepers. Thus, the appropriate choice of $\epsilon$ always realizes the specified tolerance over the useful band.

[20] When better tolerances are required and when the network configuration is not rigidly specified, Jacobian elliptic functions, rather than Tchebycheff polynomials, might be employed.

However, it is important to observe that a given value of $\epsilon$ automatically determines both the pass-band tolerance and the rate of cut-off in the out-band region. Hence, if a specified tolerance is to be realized in the useful band, no control exists over the determination of the resistance efficiency. Also, it is apparent from Figs. 19 and 20 that small in-band deviations are always obtained at the expense of lower resistance efficiencies, and vice versa.
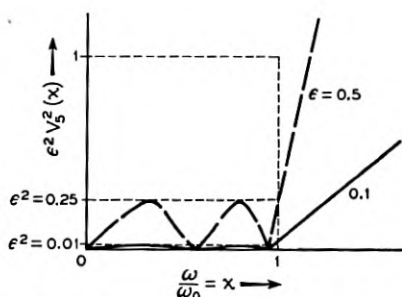


Fig. 19—Squared Tchebycheff polynomials, $\epsilon^2 V_n^2(x)$, for $n = 5$, and $\epsilon = 0.5$ and $\epsilon = 0.1$.



Fig. 20—Transfer function expressed in eq. (18) for the values of $n$ and $\epsilon$ shown in Fig. 19.

Returning to the problem of reactive equalization, for $n$ odd, $\epsilon^2 V_n^2(x)$ may be expressed as

$$\epsilon^2 V_n^2(x) = \epsilon^2(C_1 x^2 + C_2 x^4 + \cdots + C_n x^{2n}). \tag{19}$$

Thus, any $A_\nu'$ of eq. (15) is given by $A_\nu' = \epsilon^2 C_\nu$. By using the expressions for $V_1(x)$ through $V_8(x)$ tabulated previously, or eq. (17), it is a very simple task to find the $C_\nu$ for any desired $n$. Thus, $V_n^2(x) = C_1 x^2 + C_2 x^4 + \cdots + C_n x^{2n}$ is readily ascertained, and the only real problem is the choice of $\epsilon^2$. If $f(x^2)$ has already been chosen, this is accomplished by an addition of $f(x^2)$ and $\epsilon^2 V_n^2(x)$ for several values of $\epsilon^2$. When a $\epsilon^2$ is found such that the combination, when reciprocated, very closely approximates the specified resistance efficiency, $B(x^2)$ is completely defined.

The final expression for $B(x^2)$ may now be written as

$$B(x^2) = f(x^2) + \epsilon^2 V_n^2(x) = (A_0 + A_1 x^2 + \cdots + A_n x^{2n}) +$$
$$(A_1' x^2 + \cdots + A_n' x^{2n}). \qquad (20)$$

In terms of eq. (20), the corresponding expression for the special transfer function for equalization purposes becomes

$$K \left| \frac{Z_{12}(jx)}{R_0} \right|^2$$

$$= \frac{1}{A_0 + (A_1 + A_1') x^2 + (A_2 + A_2') x^4 + \cdots + (A_n + A_n') x^{2n}}. \qquad (21)$$

When all the $A_\nu$ and $A_\nu'$ are known in a particular design, the coefficients $B_1 \cdots B_n$ of eq. (7) may be readily determined. Hence, the elements of the network may be found by using the appropriate equations of Section 2.

## 4. Approximation Method

This section will consider the second of the two main tasks in the formulation of the design method. Broadly speaking, the special transfer function derived in the previous section, eq. (13), provides the approximating functions to be used in this problem while this section develops the systematic method of determining the coefficients of these functions for a finite number of network elements. The function of most interest in the approximation problem is the in-band approximating function $f(x^2)$. Thus, the development of the approximation method for reactive equalizers is concerned specifically with the determination, consistent with the previous limitations and requirements, of the coefficients, $A_0 \cdots A_n$, of the polynomial $f(x^2)$.

The Fourier method of polynomial approximation, first introduced by Wiener,[21] is characterized by a transformation of the independent variable to make the approximating function in the new frequency domain a periodic function. Thus, the well-known method of Fourier analysis is available as a general polynomial approximation method. This method has not been applied extensively in practical applications. However, the uniform nature of $B(x^2)$ over the useful frequency range makes its application to the design of reactive equalizers of the type described here seem feasible.

By the transformation $x = \tan \varphi/2$ the frequency domain, $0 \leq x \leq \infty$, is transformed to a corresponding $\varphi$ domain, $0 \leq \varphi \leq \pi$. Since the range of interest is 0 to $\pi$ in the $\varphi$ domain, all functions may be assumed to be either even or odd with a period $2\pi$. Thus, any amplitude approximating function

[21] Ref. 4.

may be written in the $\varphi$ domain as a Fourier cosine series,

$$f_1(\varphi) = a_0 + a_1 \cos \varphi + a_2 \cos 2\varphi + \cdots + a_n \cos n\varphi = \sum_{k=0}^{n} a_k \cos k\varphi. \quad (22)$$

In particular, the correspondence of the $x$ domain and $\varphi$ domain may be conveniently illustrated as in Fig. 21. It is to be noted that the comparatively limited region of the useful band, $0 \le x \le 1$, in the $x$ domain goes into half of the available range, $0 \le \varphi \le \frac{\pi}{2}$, in the $\varphi$ domain. It is apparent, then, that some advantage has already been gained by this transformation.

Before attention can be confined to the evaluation of the coefficients, $a_k$, it is necessary to establish the form of the approximating function in the $\varphi$ domain which corresponds to $f(x^2)$ in the frequency domain, and to relate



Fig. 21—Graphical representation of the transformation $x = \tan \frac{\varphi}{2}$.

the $A_k$ in eq. (14) to the $a_k$ in eq. (22). This is accomplished by means of the following relationships:

$$x = \tan \frac{\varphi}{2} = \sqrt{\frac{1 - \cos \varphi}{1 + \cos \varphi}}$$

$$\cos \varphi = \frac{1 - x^2}{1 + x^2}$$

$$\cos n\varphi = V_n(\cos \varphi).$$

Thus, the corresponding expression for eq. (22) in the frequency domain becomes

$$f_1(\varphi) = a_0 + a_1 V_1(\cos \varphi) + a_2 V_2(\cos \varphi)$$

$$+ a_3 V_3(\cos \varphi) + \cdots + a_n V_n(\cos \varphi)$$

$$f_1(\cos\varphi) = b_0 + b_1 \cos\varphi + b_2 \cos^2\varphi + b_3 \cos^3\varphi + \cdots + b_n \cos^n\varphi$$

$$f_1(x^2) = b_0 + b_1 \left(\frac{1-x^2}{1+x^2}\right) + b_2 \left(\frac{1-x^2}{1+x^2}\right)^2$$

$$+ b_3 \left(\frac{1-x^2}{1+x^2}\right)^3 + \cdots + b_n \left(\frac{1-x^2}{1+x^2}\right)^n$$

$$f_1(x^2) = \frac{A_0 + A_1 x^2 + A_2 x^4 + A_3 x^6 + \cdots + A_n x^{2n}}{(1+x^2)^n} = f(x^2)f_2(x^2),$$

where $f_2(x^2) = \dfrac{1}{(1+x^2)^n}$.

Therefore, it is necessary to predistort the approximated function $B(x^2)$ by redefining the $f(\varphi)$ corresponding to $f(x^2)$ as

$$f(\varphi) = \frac{f_1(\varphi)}{f_2(\varphi)} \rightarrow \sum_{k=0}^{n} A_k x^{2k} = f(x^2), \qquad (22)'$$

where

$$f_1(\varphi) = \sum_{k=0}^{n} a_k \cos k\varphi \rightarrow \frac{\sum_{k=0}^{n} A_k x^{2k}}{(1+x^2)^n} = f_1(x^2),$$

and

$$f_2(\varphi) = \cos^{2n}\frac{\varphi}{2} \rightarrow \frac{1}{(1+x^2)^n} = f_2(x^2).$$

Hence, $f_1(\varphi)$, which corresponds to the approximating function $f(x^2)$ multiplied by $\dfrac{1}{(1+x^2)^n}$ in the frequency domain, is the approximating function in the $\varphi$ domain. In practice, the indicated predistortion of $B(x^2)$ may be carried out either before or after the specification has been transformed to the $\varphi$ domain. Table I shows the relation of the $A_k$ to the $a_k$ for $n = 3$ and $n = 5$.

TABLE I

RELATION OF THE $A_k$ OF $f(x^2)$ TO THE $a_k$ OF $f_1(\varphi)$ FOR $n = 3$ AND $n = 5$

| $n = 3$ | $n = 5$ |
|---|---|
| $A_0 = a_0 + a_1 + a_2 + a_3$ | $A_0 = a_0 + a_1 + a_2 + a_3 + a_4 + a_5$ |
| | $A_1 = 5a_0 + 3a_1 - 3a_2 - 13a_3 - 27a_4 - 45a_5$ |
| $A_1 = 3a_0 + a_1 - 5a_2 - 15a_3$ | $A_2 = 10a_0 + 2a_1 - 14a_2 - 14a_3 + 42a_4 + 210a_5$ |
| | $A_3 = 10a_0 - 2a_1 - 14a_2 + 14a_3 + 42a_4 - 210a_5$ |
| $A_2 = 3a_0 - a_1 - 5a_2 + 15a_3$ | $A_4 = 5a_0 - 3a_1 - 3a_2 + 13a_3 - 27a_4 + 45a_5$ |
| $A_3 = a_0 - a_1 + a_2 - a_3$ | $A_5 = a_0 - a_1 + a_2 - a_3 + a_4 - a_5$ |

It is to be recognized in the following derivation and procedure that $f_1(\varphi)$ represents the actual response of the network while $B(\varphi) \cos^{2n} \frac{\varphi}{2}$, the predistorted specification for $B(x^2)$ in the $\varphi$ domain, represents the desired response. For convenience, $B(\varphi) \cos^{2n} \frac{\varphi}{2}$ may be called the amplitude function $a(\varphi)$. In addition, it is important to note that $a(\varphi)$ is specified only over the range $0 \leq \varphi \leq \frac{\pi}{2}$, and the restrictions on the behavior of the approximating function $f_1(\varphi)$ outside this range are related to the restrictions on $f(x^2)$ in the out-band region of the $x$ domain. The general problem is thus one of approximating the amplitude function $a(\varphi)$ by a Fourier cosine series, $\sum_{k=0}^{n} a_k \cos k\varphi$.

The first step towards a systematic method of obtaining the Fourier cosine coefficients, $a_0 \cdots a_n$, is the specification of the manner in which the tolerance of match is to be minimized. In this case, the approximation is always specified in the mean-square sense, i.e., the optimum coefficients are obtained by solving the set of linear equations which are determined when the integral of the error squared,

$$I = \int \left[ a(\varphi) - \sum_{k=0}^{n} a_k \cos k\varphi \right]^2 d\varphi, \tag{23}$$

is minimized.

The set of linear equations which relates the $a_k$ of the approximating function $f_1(\varphi)$ to the approximated function $a(\varphi)$ is derived for a range 0 to $s$ in the $\varphi$ domain with $s \leq \pi$ by minimizing eq. (23).[22] The minimum condition is specified when the derivative with respect to each coefficient $a_j$ is zero. Thus,

$$\frac{\partial I}{\partial a_j} = \int_0^s 2 \left[ a(\varphi) - \sum_{k=0}^{n} a_k \cos k\varphi \right] [-\cos j\varphi] \, d\varphi = 0 \tag{24}$$

is the analytical expression for this condition. Collecting terms,

$$\frac{\partial I}{\partial a_j} = -2 \int_0^s [a(\varphi) \cos j\varphi] \, d\varphi + 2 \int_0^s \left[ \sum_{k=0}^{n} a_k \cos k\varphi \right] [\cos j\varphi] \, d\varphi$$

$$= -2 \int_0^s [a(\varphi) \cos j\varphi] \, d\varphi + 2a_j \int_0^s \cos j\varphi \cos k\varphi \, d\varphi = 0,$$

and letting $P_{jk} = \int_0^s \cos j\varphi \cos k\varphi \, d\varphi$ and $C_k = \int_0^s [a(\varphi) \cos j\varphi] d\varphi$, the set of

---

[22] This derivation is similar to one given by R. M. Redheffer in Ref. 6, pp. 8–10.

linear equations becomes

$$\sum_{j=0}^{n} P_{jk} a_j = C_k. \qquad (j = 0, 1, 2, \cdots, n) \qquad (25)$$

Therefore, the procedure for determining the optimum coefficients for the range 0 to $s$ in the $\varphi$ domain is as follows: First, compute the $C_k$ which depend on the approximated function $a(\varphi)$.

$$C_k = \int_0^s [a(\varphi) \cos k\varphi] \, d\varphi. \qquad (26)$$

Next, compute the elements of $P_{jk}$ given by

$$P_{jk} = \frac{\sin (j - k)s}{2(j - k)} + \frac{\sin (j + k)s}{2(j + k)} \quad (k \neq j);$$

$$P_{jj} = \frac{s}{2}; \qquad P_{00} = s. \qquad (27)$$

These elements depend only on the range $s$ and terminate with the desired $n$ in any design. For convenience, these numbers may be arranged in the form of a symmetrical matrix $[P_{jk}]$. Hence, the optimum coefficients are found by solving the matrix equation,

$$[P_{jk}] \times [a_j] = [C_k]. \qquad (j, k = 0, 1, 2, \cdots, n) \qquad (28)$$

In this problem of approximating $B(x^2)$ to a high degree of precision over the useful frequency range, the range in the $\varphi$ domain of most interest is 0 to $\frac{\pi}{2}$. However, before the approximation over only part of the frequency range is considered, it is helpful to set down the relations which apply when $a(\varphi)$ is approximated over the whole frequency range, $s = \pi$. In this case, the matrix $[P_{jk}]$ takes on a form in which all non-diagonal entries are zero. Thus,

$$[P_{jk}] = \begin{bmatrix} P_{00} & P_{01} & \cdot & \cdot & P_{0n} \\ P_{10} & P_{11} & \cdot & \cdot & P_{1n} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ P_{n0} & \cdot & \cdot & \cdot & P_{nn} \end{bmatrix} = \begin{bmatrix} \pi & 0 & 0 & 0 & \cdot & \cdot & \cdot \\ 0 & \frac{\pi}{2} & 0 & 0 & \cdot & \cdot & \cdot \\ 0 & 0 & \frac{\pi}{2} & 0 & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & \frac{\pi}{2} & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \frac{\pi}{2} & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix}.$$

The solution in this case is particularly simple, and gives the well-known Fourier coefficients,

$$a_0 = \frac{1}{\pi} \int_0^\pi a(\varphi) \, d\varphi \qquad\qquad (j = 0),$$

$$a_j = \frac{2}{\pi} \int_0^\pi a(\varphi) \cos j\varphi \, d\varphi \qquad\qquad (j \neq 0).$$

Hence, each coefficient $a_j$ is dependent only on the area under the corresponding function $a(\varphi) \cos j\varphi$.

This result, even though it simplifies the procedure of calculating the $a_j$ in eq. (28), has only limited usefulness in this problem. As mentioned above, the range of direct interest extends only to $s = \frac{\pi}{2}$. Thus, an approximation over the whole range requires that an $f(x^2)$ be arbitrarily specified in the out-band region. Such a procedure, in this case, is an unnecessary restriction on the form of $f(x^2)$ outside the useful frequency range. Thus, an approximation over a finite range 0 to $\frac{\pi}{2}$ is the procedure to be considered in detail.

Starting as before, the system of equations in matrix notation which corresponds to eq. (28) is

$$
\begin{bmatrix}
\frac{\pi}{2} & 1 & 0 & -\frac{1}{3} & 0 & \frac{1}{5} & \cdot & \cdot \\
1 & \frac{\pi}{4} & \frac{1}{3} & 0 & -\frac{1}{15} & 0 & \cdot & \cdot \\
0 & \frac{1}{3} & \frac{\pi}{4} & \frac{3}{5} & 0 & -\frac{5}{21} & \cdot & \cdot \\
-\frac{1}{3} & 0 & \frac{3}{5} & \frac{\pi}{4} & \frac{3}{7} & 0 & \cdot & \cdot \\
0 & -\frac{1}{15} & 0 & \frac{3}{7} & \frac{\pi}{4} & \frac{5}{9} & \cdot & \cdot \\
\frac{1}{5} & 0 & -\frac{5}{21} & 0 & \frac{5}{9} & \frac{\pi}{4} & \cdot & \cdot \\
\cdot & \cdot & \cdot & \cdot & \cdot & \cdot & & \\
\cdot & \cdot & \cdot & \cdot & \cdot & \cdot & &
\end{bmatrix}
\times
\begin{bmatrix}
a_0 \\ a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ \cdot \\ \cdot
\end{bmatrix}
=
\begin{bmatrix}
C_0 \\ C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ \cdot \\ \cdot
\end{bmatrix},
$$

where the elements of $[P_{jk}]$ up to and including $P_{55}$ have been evaluated. Hence, the problem is the solution of the first $(n + 1)$ of these equations for the coefficients $a_0 \cdots a_n$. In practice, this solution may be simplified for a desired $n$ by computing once and for all the elements of the inverse matrix $[P_{jk}]^{-1}$. This matrix is formed by replacing each element of the determinant $\| P_{jk} \|$ by its minor, dividing each minor by this determinant,

and interchanging rows and columns. Thus, the solution of the $a_j$ is expressed directly in terms of the $C_k$ and becomes

$$[a_j] = [P_{jk}]^{-1} \times [C_k] \text{ or } a_j = \sum_{j=0}^{n} P_{jk}^{-1} C_k. \tag{29}$$

The sufficiency of this procedure is established when it is proved that the determinant $\| P_{jk} \|$ is different from zero for the particular value of $s$ considered. Since $s$ is a rational multiple of $\pi$ in this case and all non-diagonal entries are algebraic numbers, $\pi$ cannot satisfy an equation with algebraic coefficients to make $\| P_{jk} \| = 0$. Thus, the system of eq. (29) is a unique solution, and this solution gives the absolute minimum in the sense that no other set of $a_j$ will produce a smaller mean-square error over the range

$0$ to $\dfrac{\pi}{2}$.

However, for some values of $n$ the determinant of coefficients becomes extremely small. This condition produces very large numerical values of the elements of $[P_{jk}]^{-1}$. Since the $a_j$ and $C_k$ are usually small compared with these elements, the accuracy of the solution is impaired. Hence, the system of eq. (29) in some cases represents a set of nearly dependent equations with a fairly wide range of solution. This practical limitation on the uniqueness of these equations may be overcome quite readily by arbitrarily changing one of these equations to produce, for calculation purposes, a dependent set of equations. It turns out that the most expedient choice of this change is to replace the $P_{00} = \dfrac{\pi}{2}$ of $[P_{jk}]$ by $P_{00} = \dfrac{\pi}{4}$. This, in effect, modifies the weighting of $a_0$ in these equations and does not, in general, limit the usefulness of the result. Hence, the system of eq. (28) with $\dfrac{\pi}{2}$ replaced by $\dfrac{\pi}{4}$ determines a set of coefficients, $a_0 \cdots a_n$, which are reasonably close to the optimum for $s = \dfrac{\pi}{2}$.

It is appropriate at this point to indicate a practical modification in the approximation method which serves, incidentally, to clarify the reasons for accepting as suitable a set of coefficients that are not the optimum $a_j$ over the useful band in the $\varphi$ domain.

This modification arises since the foregoing method has considered only the average error over the range $0$ to $\dfrac{\pi}{2}$. However, an analysis of the percentage error in $f(x^2)$, and of the corresponding deviation in $\alpha$ over this range, shows that the approximation to $a(\varphi)$ is most critical at high frequencies and becomes decreasingly critical as lower frequencies are reached. Thus, in any design, it is necessary to make a slight adjustment of the

coefficients $a_0 \cdots a_n$ after they have been obtained from eq. (29) in order to compensate for this decreased tolerance of $\sum_{j=0}^{n} a_j \cos j\varphi$ at high frequencies in the useful band. The exact method of accomplishing this modification depends on the particular design and the ingenuity of the designer. Nevertheless, no more than a few trials are necessary, in general, to produce the desired precision at all frequencies in the useful band.

In practice, then, it is not appropriate that the Fourier cosine coefficients finally chosen represent the optimum coefficients in the mean-square sense. However, the important result established is that a systematic method which realizes a satisfactory set of coefficients $A_0 \cdots A_n$ of $f(x^2)$ has been developed.
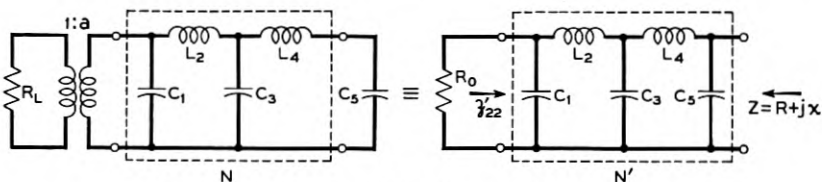


Fig. 22—Input coupling network configuration.

## 5. Illustrative Design

The numerical example which will be considered is the design of an input coupling network to equalize partially the loss characteristic of a coaxial line. On the basis of the previous discussion of the design method it is advantageous to break down the procedure into four general operations:

(1) Network Specifications
(2) Transfer Specifications
(3) Solution of Approximation Problem
(4) Realization of Non-dissipative Network

The first two of these operations are the choice of the appropriate form of the design requirements while the last two represent the major divisions in the procedure for designing the network to meet these requirements.

In this design, a set of network requirements which are consistent with the requirements indicated in Section 2 may be chosen as indicated in Fig. 22. Thus, in order that the network $N'$ correspond to the high-side equivalent circuit of the coupling transformer and, at the same time, have a final capacitance $C_n$, the least number of elements which may be chosen in a practical design is $n = 5$. The specified elements of Fig. 22 are the parasitic terminating capacitance $C_5$ and the effective impedance of the line, $R_L$.[23]

[23] See footnote 4.

Practical values for these elements may be chosen as $C_5 = 20\ \mu\mu f$ and $R_L = 150$ ohms.

Next, the transfer specifications for this illustration may be summarized as

  (a) Degree of equalization—$k = 0.25$
  (b) Useful band—2.5 to 8.0 mc
  (c) Useful band distortion—$< \pm 0.10$ db
  (d) Resistance efficiency—65%

The computation of the desired transfer characteristic $Ke^{2kf(x)}$ begins with the consideration of the degree of equalization. In order to equalize one-quarter of the power loss between coaxial repeaters, the transfer characteristic over the useful band must vary as $Ke^{\alpha'/2}$ where $\alpha'$ represents the complete line loss between repeaters. If it is assumed that $\alpha'$ is 4 nepers (34.7 $db$)[24] at 8.0 mc ($x = 1$) and varies as $\alpha' = f(x) = 4\sqrt{x}$, the transfer characteristic over the range, $0 \le x \le 1$, according to eq. (10), becomes

$$Ke^{2kf(x)} = e^{-2\alpha'_0(1-\sqrt{x})} = e^{-2(1-\sqrt{x})},$$

where $\alpha = kf(x) = \sqrt{x}$ and $\alpha'_0 = kf(1) = 1$.

The specification of a useful band from 2.5 to 8.0 mc (or $x = 0.3$ to $x = 1.0$) in this example is chosen to illustrate the practical limitation on the precision of equalization at low frequencies. The dashed curve of Fig. 23 indicates a low-frequency response which seems realistic for this illustration.

The computation of the desired transfer characteristic is completed when the out-band portion of the characteristic is chosen to satisfy the specified resistance efficiency. The assumption of a linear cut-off characteristic is suitable as an initial requirement. Hence, the transfer characteristic may be summarized as shown in Fig. 23. The solid curve of this figure represents the transfer characteristic which would be required for equalization over the range, $0 \le x \le 1$, while the dashed curve indicates the modification in this curve resulting from the choice of a conservative low-frequency response and the specification of a useful band of $0.3 \le x \le 1$.

The solution of the approximation problem consists of three main operations. First, is the determination of the amplitude function $a(\varphi)$ from the transfer characteristic specified in Fig. 23. Second, is the determination of the Fourier cosine coefficients, $a_0 \cdots a_n$, of the approximating function $f_1(\varphi)$ and the calculation of the coefficients, $A_0 \cdots A_n$, of $f(x^2)$. Third, is the choice of the coefficient $\epsilon^2$ of the squared Tchebycheff polynomial.

The amplitude function $a(\varphi)$ is calculated from the specified transfer characteristic by using the relations expressed by eq. $(22)'$. According to eq. (11) of Section 3, the specification for $B(x^2)$ over the useful band,

---

[24] This discrimination is correct for 4 or 5 miles of coaxial cable. The attenuation on a coaxial line varies as the square root of frequency.
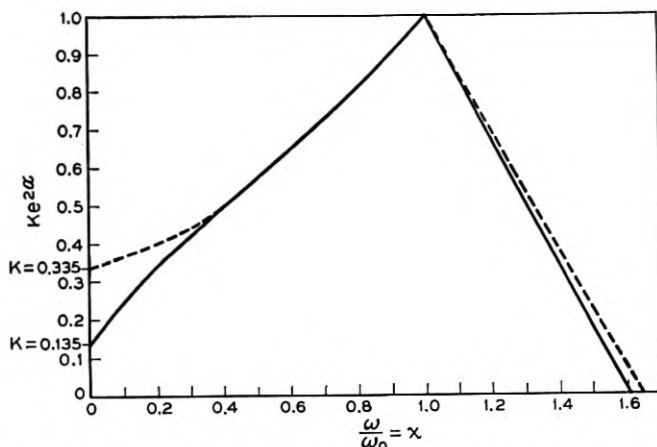
Fig. 23—Transfer characteristic for the network of Fig. 22. The dashed curve indicates the modification which results from the choice of a conservative low-frequency response.

TABLE II

RESULTS OF CALCULATIONS IN THE $x$ DOMAIN AND IN THE $\varphi$ DOMAIN

| $x$ | $B(x^2)$ | $f(x^2)$ | $\varphi$ | $B(\varphi)$ | $B(\varphi)\cos^6\frac{\varphi}{2}$ | $f_1(\varphi)$ | $f(\varphi)$ |
|-----|----------|----------|-----------|--------------|--------------------------------------|----------------|--------------|
| 0   | 3.00 | 2.98 | 0° | 3.00 | 3.00 | 2.98 | 2.98 |
| 0.1 | 2.87 | 2.91 | 10° | 2.88 | 2.80 | 2.77 | 2.87 |
| 0.2 | 2.69 | 2.74 | 20° | 2.74 | 2.49 | 2.48 | 2.73 |
| 0.3 | 2.49 | 2.48 | 30° | 2.56 | 2.09 | 2.09 | 2.57 |
| 0.4 | 2.09 | 2.17 | 40° | 2.21 | 1.54 | 1.58 | 2.28 |
| 0.5 | 1.80 | 1.85 | 50° | 1.87 | 1.05 | 1.07 | 1.95 |
| 0.6 | 1.57 | 1.57 | 60° | 1.60 | 0.68 | 0.70 | 1.65 |
| 0.7 | 1.37 | 1.39 | 70° | 1.37 | 0.42 | 0.43 | 1.39 |
| 0.8 | 1.22 | 1.23 | 80° | 1.17 | 0.24 | 0.24 | 1.17 |
| 0.9 | 1.11 | 1.13 | 90° | 1.00 | 0.13 | 0.13 | 1.00 |
| 1.0 | 1.00 | 1.00 | | | | | |
| 1.1 | — | 0.56 | | | | | |
| 1.2 | — | −0.32 | | | | | |
| 1.3 | — | −2.12 | | | | | |
| 1.5 | — | −11.4 | | | | | |
| 2.0 | — | −115.0 | | | | | |

$0.3 \leq x \leq 1$, becomes

$$B(x^2) = e^{2\alpha_0'} e^{-2kf(x)} = e^{2(1-\sqrt{x})}.$$

In addition, the specification for $B(x^2)$ may be extended to zero frequency by reciprocating the dashed portion of the curve of Fig. 23 in the range $0 \leq x < 0.3$.

In this illustration a simplified $f(x^2) = A_0 + A_1x^2 + A_2x^4 + A_3x^6$ of order $(n - 2)$ may be chosen such that the transfer characteristic is matched within the specified tolerance over the useful band.[25] The specification $a'(\varphi)$ is determined from $B(x^2)$ by (1) calculating the $B(\varphi)$ which corresponds to $B(x^2)$ in the $\varphi$ domain, and (2) multiplying $B(\varphi)$ by $\cos^{2n}\frac{\varphi}{2}$ to obtain $a(\varphi) = B(\varphi) \cos^{2n}\frac{\varphi}{2}$. The results of these calculations in the $\varphi$ domain are indicated by the fifth and sixth columns of Table II.

The Fourier cosine coefficients, $a_0 \cdots a_n$, are found by solving the set of linear equations expressed by eq. (25) for $n = 3$ and $s = \frac{\pi}{2}$. The $C_k$ which depend on the approximated function $a(\varphi)$ are computed from eq. (26). After the indicated graphical integration is carried out, these constants have the following values in this illustration:

$$C_0 = 2.323$$
$$C_1 = 1.964$$
$$C_2 = 1.148$$
$$C_3 = 0.452$$

The matrix $[P_{jk}]$ for $n = 3$ according to eq. (27) is

$$[P_{jk}] = \begin{bmatrix} \frac{\pi}{2} & 1 & 0 & -\frac{1}{3} \\ 1 & \frac{\pi}{4} & \frac{1}{3} & 0 \\ 0 & \frac{1}{3} & \frac{\pi}{4} & \frac{3}{5} \\ -\frac{1}{3} & 0 & \frac{3}{5} & \frac{\pi}{4} \end{bmatrix}.$$

The existence of a solution of eq. (28) depends on $\| P_{jk} \| \neq 0$. In this case this determinant becomes

$$\| P_{jk} \| \cong 0.00009.$$

Thus, for all practical purposes, the linear equations for $n = 3$ represent a dependent set. However, when $P_{00} = \frac{\pi}{4}$ is substituted for $\frac{\pi}{2}$ above,[26] the

---

[25] For the value of the tolerance specified in this illustration, an $f(x^2)$ of order 3 turns out to be satisfactory. In the general case, where a higher degree of precision is desired, it is, of course, expedient to choose an $f(x^2)$ of order $n$.

[26] See discussion on p. 742.

solution for the $a_j$ according to eq. (29) is

$$
\begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix} =
\begin{bmatrix}
-1.273 & 2.117 & -1.166 & 0.350 \\
2.117 & -1.273 & -0.350 & 1.166 \\
-1.166 & -0.350 & 4.320 & -3.798 \\
0.350 & -1.166 & -3.798 & 4.320
\end{bmatrix}
\times
\begin{bmatrix} 2.323 \\ 1.964 \\ 1.148 \\ 0.452 \end{bmatrix}
=
\begin{bmatrix} 0.016 \\ 2.527 \\ -0.150 \\ 0.698 \end{bmatrix}
$$

As previously stated, these coefficients represent the practical minimum of the average error in the mean-square sense over the range 0 to $\frac{\pi}{2}$ in the $\varphi$ domain. However, they do not represent the best match over the useful band for this illustration. The adjustment of these coefficients to produce a more satisfactory match at high frequencies in the useful band begins by changing the value of $a_0$ to make $f_1\left(\frac{\pi}{2}\right) = a_0 - a_2 = 0.125$. This condition is satisfied when the general level of response is lowered so that $a_0 = -0.025$. The only further adjustment that is necessary in order to compensate for the decreased tolerance of $f_1(\varphi) = \sum_{j=0}^{3} a_j \cos j\varphi$ at high frequencies in the useful band is a change in the value of $a_3$. When $a_3$ is adjusted to $a_3 = 0.623$ a suitable approximating function for $a(\varphi)$ in this illustration is

$$
f_1(\varphi) = \sum_{j=0}^{n} a_j \cos j\varphi = -0.025 + 2.527 \cos \varphi
$$
$$
- 0.150 \cos 2\varphi + 0.623 \cos 3\varphi.
$$

Hence, the approximating function for $B(\varphi)$ is

$$
f(\varphi) = \frac{f_1(\varphi)}{f_2(\varphi)} = \frac{-0.025 + 2.527 \cos \varphi - 0.150 \cos 2\varphi + 0.623 \cos 3\varphi}{\cos^6 \frac{\varphi}{2}}.
$$

These functions are tabulated in the last two columns of Table II.

The coefficients $A_0 \cdots A_3$ of $f(x^2)$ are easily calculated from the $f_1(\varphi)$ and $f(\varphi)$ above by the relation of the $A_k$ to the $a_j$ expressed in Table I. Thus,

$$
f(x^2) = 2.975 - 6.143x^2 + 7.493x^4 - 3.325x^6.
$$

The final operation in the solution of the approximation problem is the choice of the squared Tchebycheff polynomial, $\epsilon^2 V_n^2(x)$, which satisfies a resistance efficiency of 65 per cent. The Tchebycheff polynomial for $n = 5$ is

$$
V_5(x) = 5x - 20x^3 + 16x^5.
$$

Thus, $V_5^2(x)$ becomes

$$V_5^2(x) = 25x^2 - 200x^4 + 560x^6 - 640x^8 + 256x^{10}.$$

A $\epsilon^2 = 0.01$ is easily found such that the resistance efficiency calculated from a graphical integration of $\dfrac{1}{f(x^2) + \epsilon^2 V_5^2(x)}$ equals 65 per cent. Hence, the analytical expression for $K \left| \dfrac{Z_{12}(jx)}{R_0} \right|^2$ becomes

$$\frac{1}{f(x^2) + \epsilon^2 V_n^2(x)} = \frac{1}{\begin{array}{c}(2.975 - 6.143x^2 + 7.493x^4 - 3.325x^6) \\ + (0.25x^2 - 2.00x^4 + 5.60x^6 - 6.40x^8 + 2.56x^{10})\end{array}}.$$
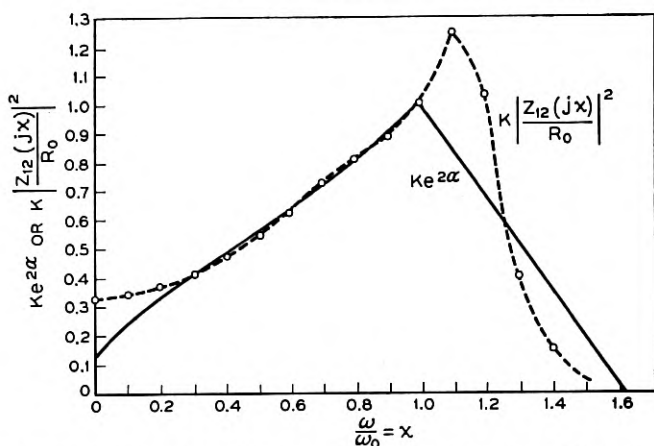


Fig. 24—Comparison of the resultant special transfer function with the transfer characteristic of Fig. 23.

This expression is the resultant special transfer function which satisfactorily approximates the transfer characteristic of Fig. 23. Fig. 24 shows a plot of these functions for comparison purposes.

The squared magnitude of the transfer impedance of the network $N'$ is found from the analytical expression for the special transfer function by adjusting the value of $K$ so that $KA_0 = 1$. Therefore,

$$\left| \frac{Z_{12}(jx)}{R_0} \right|^2 = \frac{1}{1 - 1.981x^2 + 1.846x^4 + 0.765x^6 - 2.157x^8 + 0.861x^{10}}.$$

The elements of the network $N'$ are found from the squared magnitude of the transfer impedance by methods standard in circuit theory.[27] The network elements of Fig. 22 in terms of unit impedance and unit radian

[27] Ref. 2, pp. 25–53.

frequency turn out to be

$$C_1 = 0.470 \text{ farads} \qquad L_2 = 1.250 \text{ henrys}$$
$$C_3 = 1.201 \text{ farads} \qquad L_3 = 2.220 \text{ henrys.}$$
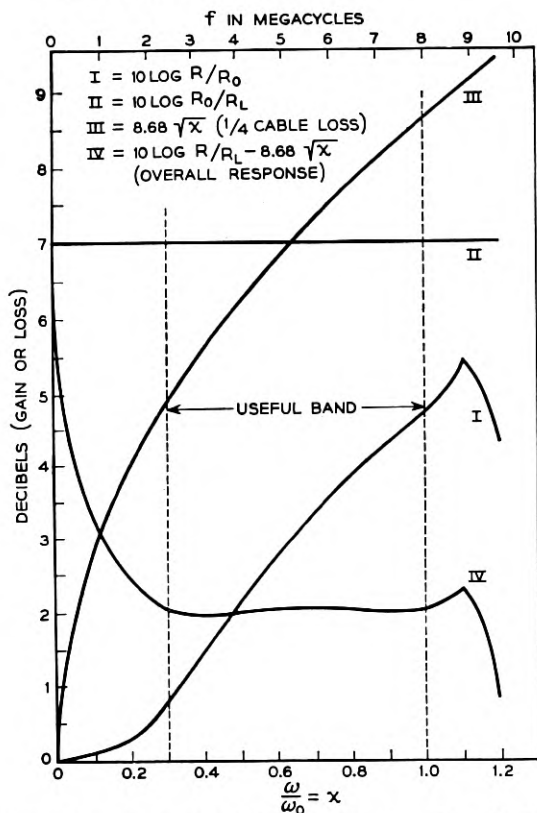$$C_5 = 0.594 \text{ farads}$$



Fig. 25—Computed gain characteristic of the input coupling circuit of Fig. 22.

$R_0$ is calculated from the equation which relates to normalized value of $C_5$ above to $\omega_0$ and the actual value of $C_5 = 20 \times 10^{-12}$ farads. Thus

$$\frac{0.594}{R_0 \, \omega_0} = 20 \times 10^{-12} \text{ farads,}$$

and $R_0 = 591$ ohms.

The actual values of the network elements of Fig. 22 are found as

$$C_1 = 15.8 \ \mu\mu f \qquad L_2 = 14.7 \text{ mh}$$
$$C_3 = 40.5 \ \mu\mu f \qquad L_4 = 26.2 \text{ mh,}$$
$$C_5 = 20.0 \ \mu\mu f$$

and the step-up turns ratio, $a$, of the ideal transformer is

$$a = \sqrt{\frac{R_0}{R_L}} = 1.98.$$

These values then represent the input coupling network which theoretically equalizes to the specified degree of precision one-quarter of the power loss between coaxial repeaters over a frequency band from 2.5 to 8.0 mc. The computed gain characteristic of this network is plotted in Fig. 25, Curve I. The presence of the ideal transformer represents an added constant gain, Curve II, given by db $= 10 \log \dfrac{R_0}{R_L} = 5.96$. The total gain inserted by the network, the sum of Curves I and II, is db $= 10 \log \dfrac{R}{R_L} = 10 \log \dfrac{R}{R_0} + 5.96$.

Since Curve III represents one-quarter of the power loss between repeaters, Curve IV is the overall transmission gain of the line and equalizer.[28] The deviation of Curve IV from a constant transmission over the useful band is less than $\pm 0.08$ db. It may be concluded, then, that a satisfactory non-dissipative design has been obtained.

References

1. H. W. Bode, "Network Analysis and Feedback Amplifier Design," Chap. 16, Van Nostrand, New York, 1945.
2. S. Darlington, "Synthesis of Reactance 4-Poles," *Journal of Mathematics and Physics*, Vol. XVIII, pp. 275–353, September, 1939, Also *Bell System Monograph* B-1186.
3. R. M. Fano, "Theoretical Limitations on the Broadband Matching of Arbitrary Impedances," *Report No.* 41, *M. I. T. Research Laboratory of Electronics*, January, 1948.
4. E. A. Guillemin, "Classroom Notes on Network Synthesis" (Notes dictated at M. I. T., course 6.561 and 6.562—as yet unpublished).
5. E. L. Norton, "Constant Resistance Networks with Applications to Filter Groups," *B. S. T. J.*, Vol. XVI, pp. 178–193, April, 1937. Also *Bell System Monograph* B-991.
6. R. M. Redheffer, "Design of a Circuit to Approximate a Prescribed Amplitude and Phase," *Report No.* 54, *M. I. T. Research Laboratory of Electronics*, November, 1947.

---

[28] Criticism may well be directed at the gain peak above the useful band. However, this condition is somewhat exceptional and probably would not occur with an in-band approximating function of order $n$ rather than $(n - 2)$.

# Abstracts of Technical Articles by Bell System Authors

*Testing Cathode Materials in Factory Production.*[1] J. T. ACKER. The paper deals with the methods of testing radio-tube cathode materials in factory production, and especially with a comparison of several specific lots of materials of variable content. It is believed that this is the first time the electron-tube industry has made mass tests on a well-controlled engineering basis of cathode materials which vary in single component elements.

*Advances in the Theory of Ferromagnetism.*[2] R. M. BOZORTH. This article presents the results of the most recent investigations in the field of ferromagnetism. There have been a number of new ideas brought forth through research along these lines, of which three of the most outstanding ones are explained and illustrated.

*On Magnetic Remanence.*[3] R. M. BOZORTH. The magnetic retentivity of many materials is about half of the magnetization at saturation, a fact accounted for by simple domain theory. In some materials, however, the retentivity is only a small fraction of saturation, sometimes less than 10 per cent. The explanation of this fact is discussed. It is suggested that in materials with almost zero magnetic anisotropy the Bloch walls between domains increase in thickness until they envelop the whole specimen and the domain structure disappears.

*Multifrequency Pulsing in Switching.*[4] C. A. DAHLBOM, A. W. HORTON, JR., and D. L. MOODY. Applications of multifrequency pulsing in switching are described in this article. Today, many installations of this type are being made in cities throughout the nation. This system permits operators or senders to complete calls to crossbar offices without the aid of other operators.

*Circuits for Cold Cathode Glow Tubes.*[5] W. A. DEPP and W. H. T. HOLDEN. This paper discusses fundamental operating characteristics and typical circuits using cold cathode glow tubes for relays, impulse generators, pulse counting and interlocking functions.

*The Substitution Method of Measuring the Open Circuit Voltage Generated by a Microphone.*[6] M. S. HAWLEY. An analysis of the substitution method of measuring the open circuit voltage generated by a microphone is given

[1] *Proc. I.R.E.—Waves and Electrons Section*, v. 37, pp. 688–690, June 1949.
[2] *Elec. Engg.*, v. 68, pp. 471–476, June 1949.
[3] *Zeits. f. Physik*, v. 124, 7/12, pp. 519–527, 1948.
[4] *Elec. Engg.*, v. 68, pp. 505–510, June 1949.
[5] *Elec. Mfg.*, v. 44, pp. 92–97, July 1949.
[6] *Jour. Acous. Soc. Amer.*, v. 21, pp. 183–189, May 1949.

which shows that the "normal" substitution voltage equals the open circuit voltage for all types of acoustic measurements and for any value of electric impedance loading the microphone. It is shown that the method recently proposed by some authors of removing the acoustic load from the microphone when applying the substitution voltage results in a substitution voltage which does not equal the open circuit voltage. It is also shown that a formula for the response of a transducer derived for a system in which the microphones are open-circuited may be used when the microphones are terminated by finite electrical impedances, by replacing the generated open circuit voltages in the formula by the corresponding "normal" substitution voltages.

Consideration is given to the restriction in the definition of the pressure response of a transducer made necessary by the fact that the pressure on a microphone diaphragm is a function of the electrical impedance terminating the microphone.

An experiment is described which involves a microphone coupled to a chamber, the acoustical impedance of which is high relative to that of the microphone. The results of this experiment agree with the conclusions of the analysis.

*A Note on Filter-Type Traveling-Wave Amplifiers.*[7] J. R. PIERCE[*] and NELSON WAX. A small-signal analysis of systems in which an electron beam interacts with a circuit composed of discrete filter elements is given here. The effects of a line beam interacting with a series of gaps, which are capacitive elements of a filter structure, are calculated, and it is shown that an admittance can be introduced which arises from the presence of the electrons. This admittance is in parallel with the gap capacitance, and thus will alter the propagation factor of the filter circuit. It is shown that traveling-wave solutions exist for the combination of electron beam and filter circuit, and that there is a solution which has a positive real part, indicating that gain will be exhibited.

[7] *Proc. I.R.E.*, v. 37, pp. 622–625, June 1949.
[*] Of Bell Tel. Labs.

# Contributors to This Issue

A. P. Brogle, Jr., S.B. and S.M. in Electrical Engineering, Massachusetts Institute of Technology, 1949; Signal Corps, U. S. Army, 1942–46; Bell Telephone Laboratories, 1948. While at the Bell Telephone Laboratories, Mr. Brogle worked on the development of high-frequency networks as part of the M.I.T. Cooperative Course in Electrical Engineering. He is now engaged in Facsimile development at the Signal Corps Engineering Laboratories, Red Bank, New Jersey.

F. R. Dennis, State College of Washington, B.S. in E.E., 1929. Bell Telephone Laboratories, 1929–. Mr. Dennis has been engaged in the design of electronic measuring apparatus, particularly of the type providing visual display of transmission characteristics. During the war he was engaged in Application Engineering of Airborne Magnetometers in this country and overseas.

E. P. Felch, Dartmouth College, A.B. in Physics, 1929. Bell Telephone Laboratories, 1929–. Mr. Felch has been concerned with the development of electronic measuring apparatus. During the war he was engaged in the development of Airborne Magnetometers and other magnetic detectors.

William W. Mumford, B.A., Willamette University, 1930. Bell Telephone Laboratories, 1930–. Mr. Mumford has been engaged in work that is chiefly concerned with ultra-short-wave and microwave radio communication.

Sloan D. Robertson, B.E.E., University of Dayton, 1936; M.Sc., Ohio State University, 1938, Ph.D., 1941; Instructor of Electrical Engineering, University of Dayton, 1940. Bell Telephone Laboratories, 1940–. Dr. Robertson was engaged in microwave radar work in the Radio Research Department during the war. He is now engaged in fundamental microwave radio research.

Claude E. Shannon, B.S. in Electrical Engineering, University of Michigan, 1936; S.M. in Electrical Engineering and Ph.D. in Mathematics, M.I.T., 1940. National Research Fellow, 1940. Bell Telephone Laboratories, 1941–. Dr. Shannon has been engaged in mathematical research principally in the use of Boolean Algebra in switching, the theory of communication, and cryptography.