# THE BELL SYSTEM

# Technical Journal

# THE BELL SYSTEM TECHNICAL JOURNAL

# Use of Frequencies Above 10 GHz for Common Carrier Applications

### By LeROY C. TILLOTSON

*We discuss problems and opportunities inherent in the use of frequencies above 10 GHz for common carrier applications. We also argue that if short hops are used to ameliorate the severe attenuation caused by excessive rainfall and if integrated solid-state electronics are used to achieve cost and reliability goals, a viable system can result. Interference-tolerant modulation methods and high performance antennas are used in order to permit many repeaters to operate co-channel in a restricted geographical area. The paper also includes several references to other papers in this issue and elsewhere which describe studies and experiments conducted to demonstrate feasibility of the system concept and components and subassemblies designed especially for broadband short-hop system application.*

I. PROLOGUE

One of the distinguishing features of the history of radio is an inexorable trend toward higher frequencies. We have witnessed in succession the use of long waves, the "broadcast" band, short waves, ultrashort waves, and microwaves; we are now on the threshold of large scale use of frequencies above 10 GHz. In each case the use of

the new frequency band was accompanied by a significant change in technology. Indeed, such changes were necessary to enable us to exploit the new portion of the radio spectrum. Radio wave propagation at higher and higher frequencies has in each case been significantly different, at each step, from that at the older more familiar wavelengths. Thus the straightforward application of the prevailing art to a new frequency band was fraught with many obvious and very severe difficulties which tended to make the new frequencies appear useless; but fundamental advances in the associated technology, new system concepts, and an ever increasing need for communications have made it possible to exploit in turn each of these new regions of the radio frequency spectrum.

The several papers accompanying this introduction are a description, in part, of work being carried on at Bell Telephone Laboratories to exploit the potential of frequencies above 10 GHz for common carrier applications. Greatest emphasis is placed on possible terrestrial systems since these are judged to be nearest in time. The potential for domestic communication by satellite repeaters operating at frequencies above 10 GHz has also been studied briefly and is reported elsewhere.[1]

The thesis of the present series of papers is to show that progress in the radio and related arts, particularly solid state millimeter wave devices and the availability of quantitative statistical data on attenuation caused by rain, has now provided a portion of the base required for such an undertaking. In addition, the vast expansion of conventional common carrier microwave systems has begun to exhaust the frequencies presently allocated and thus make it imperative to open up new regions of the radio frequency spectrum for these services.

## II. GENERAL CONSIDERATIONS

Most present day common carrier radio relay systems operating at frequencies below 10 GHz are characterized by 20- to 30-mile spacings between repeaters. They use low index frequency modulation. That these systems have proven very useful is indicated by the increasing density of such routes in the United States. In fact, in some urban areas this density is approaching saturation in the sense that, because of radio interference, it is becoming more and more difficult and will eventually be impossible to add another system. If the use of radio for such applications is to grow and to prosper, new space in

the radio frequency spectrum must be found; this will be easier to accomplish if frequencies not now used for other services can be used. This paper and the companion papers describe current work which is basic to the design of systems for use at frequencies above 10 GHz where liquid water (normally rain) severely attenuates electromagnetic waves and hence profoundly influences system design.

A measure of the problem one faces in designing a multilink radio system for use at frequencies above 10 GHz is shown on Fig. 1. These repeater spacings were determined by assuming arbitrary but reasonable system parameters (40 dB fade margin) and a very arbitrary rain rate of 100 millimeters per hour, falling uniformly over the entire path. At 4 and 6 GHz, attenuation caused by rain plays no part in determining repeater spacing; spacing is determined by terrain features and tower heights and is usually in the range 20 to 30 miles as shown. Above 10 GHz the situation is very different; repeater spacing is almost entirely determined by attenuation caused by rain and becomes only a few miles at 18 and 30 GHz as is also shown. This analysis is oversimplified but it is an adequate basis for an important conclusion: at 10 GHz and above, radio repeater spacing must be decreased, which at the highest frequencies results in very short hops. The consequences of this fact are many. If radio systems with ten times the number of repeaters needed at the lower frequencies are to be competitive, a new low-cost system concept must be evolved. Since an order of magnitude increase in the num-



Fig. 1 — Repeater spacing as determined by a rain rate of 100 mm per hour (4 inches per hour). This assumes uniform rain over the entire path and a 40 dB fade margin.

ber of intermediate repeaters produces a like increase in the amount of transmission impairment incurred, either better repeaters or more tolerant modulation methods must be used. Further, when such a vital parameter as repeater spacing is controlled by the characteristics of rainstorms, it becomes important to know much more about this natural phenomenon.

The present state of our knowledge concerning attenuation of radio waves by rainfall has been summarized in a paper by Hogg from which Figs. 2 and 3 were selected.[2] Figure 2 shows, for various parts of the United States and a locality in England, the number of minutes per year that a given level of attenuation at 30 GHz is exceeded. Large differences exist between Corvallis, Oregon, and Miami, Florida; hence allowable repeater spacings vary widely in different parts of the United States and appear to be prohibitively short in some regions. A possible improvement has been suggested by Hogg which is based on the intuitively reasonable assumption that the most intense rainstorms are also the most limited in area.[2] Figure 3, which



Fig. 2 — The number of minutes per year that a given level of attenuation at 30 GHz is exceeded for parts of the United States and England.

Fig. 3 — Improvement made possible by providing a parallel path spaced by 2 km and by providing for switching the traffic to the better path.

is based on a four-year accumulation of data, illustrates the improvement made possible by provision of a parallel path spaced by 2 km and provision for switching of the traffic to the better path. These and other related data lead us to hope that successful system designs can be devised for even the most hostile regions if only we know enough about the rainfall environment.[3-6]

One of the earliest suggestions for use of 10 to 30 GHz for terrestrial radio repeater applications was made in 1949 by J. R. Pierce and W. D. Lewis.[7] Many of the system features described below, such as short hops to ameliorate the rain problem, use of binary pulse-code modulation, and pole-mounted repeaters were discussed. But the technology available in 1949 was not adequate to support such a system. Another paper which has influenced the present studies was written by C. B. Feldman and W. R. Bennett.[8] One of the main features of this paper is an argument that band spread modulation techniques are not wasteful of the radio frequency spectrum if all factors, especially interference, are adequately taken into account. While these arguments were valid in 1949, just as they are today, no system of the suggested type resulted. In this second instance, the reasons are more complex, but a major factor once again was the lack of an adequate technology to support the system proposals.

Another important feature of radio systems having many closely spaced repeaters is tolerance to interference. This problem becomes

especially acute where, assuming a viable system concept, one must eventually reckon with many parallel and crossing routes which results in many repeaters in a given area, all of which can be within line-of-sight of each other. These considerations lead to emphasis on two aspects of system design: (*i*) well-shielded antennas having sharp beams and little off-angle radiation and (*ii*) a modulation method which provides the maximum possible communication through a given area. More antenna discrimination always makes things better; but more rugged modulation methods, while more tolerant of interference, always require more bandwidth per unit of information. Thus total communication capacity increases with bandwidth expansion as long as the capacity of the added routes, made possible by the increased tolerance to interference, exceeds the loss in capacity per radio channel. Hence an optimum bandwidth expansion exisits for a given environment, as is demonstrated in a companion paper.[9]

These problems have been studied; some of the results are reported in the accompanying papers, some in more detail than others. In particular, the repeater concept which was devised to meet the challenge of many closely spaced repeaters and the experimental system designed and built to verify the ideas are discussed in considerable detail, whereas the study of interference is not since the problem of maximizing communication through an area is a large and complex one and is as yet incomplete.[9-15] As noted above, a knowledge of the amount of attenuation which will be caused by rainstorms is vital to the entire program; this high priority problem has received considerable attention. Some of the results are reported in this issue.[3-6]

While frequencies above 10 GHz suffer the disadvantage of attenuation caused by rain and this becomes a basic limitation in the system design, there are also some advantages. Antennas with quite sharp beams and high gain are feasible at these short wavelengths, that is, antennas one meter or less in aperture provide 2° beams and gains of about 40 dB. Thus a repeater package of modest size can provide high gain and sharp beams; if the antenna is designed to control wide-angle radiation, good interference protection is also provided. These are valuable attributes for any radio system. High gain antennas together with short hops and low noise receivers which are feasible with present technology result in a requirement of only a few tens of milliwatts of transmitter power for bandwidths of the order of 100 MHz. Hence an all solid-state repeater which will handle at least 100 millon bits per second (adequate for one color television channel in digital form or about 1400 voice circuits) is feasible. This is im-

portant since at the present time our best hope of building a repeater
with adequate reliability at a competitive cost is based on solid-state
devices and integrated microwave circuits. In fact, a useful system
will result only if the repeaters require no maintenance and very
little electrical power. A good model for the design of a repeater of
the type envisaged here is one intended for submarine cable or satel-
lite application; in each of these cases, frequent repair or replacement
of a faulty repeater is so expensive as to be prohibitive. Realization
of such repeaters at a low cost will be a challenge to the skill and
ingenuity of both designers and production people, but several years'
experience with solid-state microwave systems leaves little doubt
that this can be accomplished.

We mentioned that the amount of electrical power required was
an important factor in a system design. This is not because power
itself is expensive but because reliable power delivered at the re-
peater site is expensive. Such power usually takes the form of an on-
line power supply plus a standby which is typically a battery. Bat-
teries have life, maintenance, and temperature problems, must be
kept charged, and require accessible housing. While it may appear
at the outset to be a detail, providing reliable power without main-
tenance or a bulky and unsightly housing at the base of the tower
is a crucial element in the design of the repeater. The solution
proposed here is to keep the power required to an absolute minimum
by using efficient solid-state devices to provide a small amount of
transmitter power. High transmitter power is not a panacea in any
case since a 10 dB increase under the assumptions accompanying
Fig. 1 would make possible an increase in path length of only 0.33
miles at 30 GHz. The situation is much like that on wire lines or
waveguides where a twofold increase in repeater spacing results in
a twofold increase in attenuation measured in dB instead of a 6
dB increase as in free-space propagation. The use of short hops and
lower power can keep the total power requirement for a multichan-
nel repeater to a value small enough to make feasible use of a reliable
on-site source of primary power without backup as discussed in a
companion paper.[10]

Another factor of importance is appearance. While it is undoubtedly
true that for many locations where one might presently wish to
locate a repeater, appearance would not matter, there are certainly
other locations where it would be a controlling factor. It also seems
likely that short hop systems are apt to be most needed in urban and
suburban areas where appearance is becoming increasingly import-

ant. Subjective matters are difficult for engineers to treat, and final judgment must be left to those skilled in the art of product design. The approach used here was to avoid cumbersome towers, guy wires, equipment huts, and so forth, and to try for a clean appearance modeled after the familiar street lighting standard and fixture. Evaluation of the degree to which we have succeeded must necessarily be left to the judgment of the reader. In the very long run, appearance and general acceptability may place more important constraints on the use of microwaves than purely economic or technical problems.

III. REPEATER CONCEPT

It is a consequence of the required close repeater spacing that a successful short-hop system requires repeaters which are economical to install and which need no maintenance. In this application suitably designed solid-state repeaters can open up new possibilities in contrast to their use at lower frequencies where present solutions to the repeater problem are adequate and only marginal improvement is possible. This comes about in part because at the longer wavelengths where hops of 20 to 30 miles are common, microwave electronics is only a small fraction of the total repeater station cost. For the short-hop system, the opposite is true; studies have shown that one half to two thirds (depending on number of channels) of the total repeater station cost will be for repeater electronics.

In the past, towers for microwave radio repeaters have been regarded as platforms on which antennas, waveguides, and associated electronics are mounted in somewhat the same manner as switching equipment is installed in a central office building; design and construction of the tower has been kept separate from the design of the repeater electronics. Where high towers for large and expensive repeater stations are involved, each of which tends to be a project in its own right, this is a reasonable approach.

When we are considering closely spaced repeaters where the tower cost can be a crucial part of the total system cost, the entire problem must be considered at one time. At one extreme we can design the radio repeater to use antennas which are as large as seem feasible and then insist that the tower be stable enough to keep the antennas aimed properly under any weather conditions; this is more or less what has been done in the past. At frequencies above 10 GHz, where antenna sizes which are quite feasible result in sharp beams requiring very stable towers, this approach can result in a tower so costly that the entire system concept becomes unattractive.

In the present study, we have examined another alternative. The repeater package which includes the antenna and all of the required electronics is made compact and lightweight in order to minimize wind loading on the tower which is designed to be just adequate to prevent its destruction under any likely weather conditions. The radio repeater must then be designed to give satisfactory transmission performance in face of maximum movement of the tower. This can be accomplished by proper choice of antenna size, transmitter power output, receiver noise figure, repeater spacing, modulation method, and possibly in the future some form of self-steering array antenna. Thus it will be possible to achieve a better balance between the cost of the tower and the cost of repeater amplifying and signal processing equipment.

In selecting a tower for the present system, we assume that the enclosure for the repeater electronics and associated antenna is completely symmetrical and of small diameter so that any twisting moment is small enough that changes in azimuth can be neglected. Hence we can use an antenna whose beamwidth in the azimuth plane is governed only by the precision of initial line-up, provided that its beamwidth in the vertical plane is broad enough to allow the full expected movement of the tower. These requirements can be fulfilled by using an antenna with a smaller effective aperture in the vertical plane than in the horizontal plane.*

Figure 4 shows a repeater concept which involves most of the features required for short-hop service. In this design, all of the electronic circuitry serving a specific direction is contained in a tower-mounted package which also houses the antenna and power conversion equipment. A stack of two or more such assemblies mounted at the top of a tapered cylindrical self-supporting tower makes up the complete repeater. Power and other leads are contained within the tower. Since individual repeater assemblies are free to turn relative to each other preliminary line-up to within a few degrees can easily be accomplished during installation; final line-up is done electrically with the repeater in place on the tower using vernier adjustments built into the repeater housing. This arrangement eliminates long waveguides which are required when radio apparatus is mounted at the tower base and also maintains a good antenna pattern which is vital in an interference environment.

---

* The initial repeater design as described in companion papers does not use this feature, but the idea has importance for the future when interference will become a paramount problem.[10,11]

Fig. 4 — Alternative repeater concepts.

Two problems are created: (i) once it has been installed the repeater apparatus is not readily accessible, and (ii) it must withstand the entire range of ambient temperatures. Our solution to both of these problems in a solid-state repeater designed accordingly. If the repeater is in fact reliable, the lack of accessibility is not a real handicap; initial installation and changes in repeater apparatus to accommodate changes in traffic can be planned ahead of time. In addition, route diversity, as discussed above, offers a possible solution to the problem of rain outage and also affords a worthwhile degree of protection against equipment failure. The very wide range of ambient temperatures which must be accommodated by a system which is expected to operate anywhere in the United States makes the repeater design problem considerably more difficult, but it does not introduce any new problems; every repeater must operate over at least a limited temperature range. If we understand the circuits and their interactions well enough and provide adequate margins for changes in component characteristics, a reliable repeater can be built; a first step toward this goal has been taken.[10]

IV. SHORT-HOP REPEATER DESIGN

A wide variety of possible repeater designs for short-hop systems exist, but best economy will usually be achieved when broadband

radio channels are used up to the limits set by traffic needs and by available technology. This comes about because a broadband repeater will not cost much more to build than a narrowband design; in fact, some design problems such as heat losses in filter and the stability with time and temperature of filters and certain other components are made easier by this approach. Of course, when bandwidth is doubled, transmitter power must also be doubled if margins are to remain intact; but this does not cost much as long as the power is about 100 mW where reliable solid-state devices are available. On those routes where a large volume of traffic must be handled, considerable economy can be achieved by use of broadband channels; a multiplicity of radio channels may also be used at each repeater to further increase the total system capacity. Of course this is possible only if an adequate frequency band has been allocated to such service.

Transmission through a large number of tandem connected repeaters, many of which are exposed to interference, can be tolerated only when a suitable bandwidth-expanding modulation technique is used. A basic decision has to be made whether to expand bandwidth using digital or analog techniques. A digital modulation method such as PCM is clearly superior when a very high signal-to-noise ratio must be obtained using a noisy transmission medium. On the other hand, some analog techniques, such as large index FM, are equally effective in overcoming noise and interference at S/N values of interest for common carrier applications. Analog terminal equipment is simpler, but large numbers of tandem connected repeaters which are exposed to severe interference cannot be accommodated.

In practice, however, all of these considerations are likely to be overruled by another factor—the nature of the traffic to be carried. Frequent analog-to-digital conversion is to be avoided because of signal degradation and cost; hence a radio system is most useful if it can be designed to accommodate both analog and digital signals.* For short distances, analog signals can be accommodated by using large index FM; signals which are supplied in a suitable PCM format can be carried over either short or long distances.

Table I gives a set of parameters for a possible repeater and shows the communication capability of such a repeater equipped to handle either digital or analog traffic. This is intended only as a hypothetical example; actual designs need to be tailored to a specific application

---

* As discussed in the companion paper interchangeable plug-in repeaters suited either for analog or digital signal can be used, as required.[10]

TABLE I—ILLUSTRATIVE REPEATER PARAMETERS

| | |
|---|---|
| Radio frequency | 18.5 GHz |
| Antenna effective area | 0.375 m² |
| Receiver noise figure | 7 dB |
| Receiver noise bandwidth | 150 MHz |

Transmitted power required to provide 40 dB for attenuation (by rain) and a C/N = 14 dB at the receiver is:

| Distance between repeaters (Km) | Transmitted power (mW) |
|---|---|
| 1 | 1.3 |
| 2 | 5.2 |
| 4 | 20.8 |
| 8 | 83.2 |

COMMUNICATION CAPABILITY

| *Digital*—two-phase PCM | | *Analog*—Large index phase modulation | |
|---|---|---|---|
| Pulse rate | $130 \times 10^6$/s | Baseband width | 7 MHz |
| Bit rate | 130 Mb/s | Modulation index | 8 |
| Approximate voice | | Approximate voice | |
| circuit capacity | 1900 | circuit capacity | 1600 |
| For two-phase coherent PSK with one interferer and a C/I = 15 dB, a C/N = 14 dB will result in $P_\epsilon < 10^{-9}$ in one link. | | S/N at baseband (one hop and no rain) | >70 dB |
| | | S/N, 40 dB attenuation | 36 dB |
| | | About 30 dB reduction in cochannel interference results from the large index. | |

since, as noted above, per circuit costs computed for broadband channels are always less than those computed for narrowband channels; this calculation is valid only when a reasonable "fill" is achieved.

Binary PCM impressed on the radio carrier by angle modulation and coherently detected at the following receiver is used for digital transmission. While this makes complete regeneration at each repeater possible, a hybrid scheme which uses mostly linear repeaters interspersed with a few of the digital variety, as described by Chang and Freeny, may prove desirable in situations where interference does not affect many repeaters since intense rain occurring simultaneously on several hops is also expected to be rare.[16] For the repeater parameters shown in Table I even 40 dB excess attenuation by rain does not cause many errors per hop, and hence a route with many links in tandem could be used. As is also indicated, a single co-channel interferer only 15 dB below the desired signal could be tolerated with

an error less than 1 in $10^9$ pulses.* This is a calculated value based on a perfect receiver;[17] in practice a few dB margin would be required.

In the analog case, interfering co-channel signals will be reduced about 30 dB as a result of the large index, but such interference accumulates along the system and hence is more of a problem than in the digital case where frequent regeneration is possible. In addition, the analog repeater will handle fewer voice channels as Table I indicates. This difference will be even larger if in the future improved technology and understanding make possible a doubling of the bit rate with only a slight increase in bandwidth by use of four-phase angle modulation instead of two-phase angle modulation. Before this is undertaken, however, the effect of such a change on tolerance to interference and other system parameters needs more study. In summary, a system which uses quantized modulation permits a large number of repeaters connected in tandem, provides greater tolerance to interference, and makes more efficient use of the RF channel.† Analog repeaters are simpler, but fewer can be used in tandem; they seem particularly well suited to short-haul service where the traffic to be handled is presented in analog form.

V. SUMMARY AND CONCLUSIONS

A program of research has been described which has as its objective the establishment of an adequate base for the design of common carrier radio systems operating at radio frequencies above 10 GHz. Because these short radio waves are highly attenuated by intense rainfall, closely spaced repeaters are mandatory and new lower cost repeaters are required to make such a system concept competitive. A possible solution to this problem is described which uses low-power solid-state pole-mounted repeaters. Work is also being done to learn more about intense rainfall (in New Jersey) for which a new fast-acting rain gauge has been devised and used to instrument a dense rain-gauge network. Since a successful system concept could result in many repeaters within radio line-of-sight of each other, interference is a vital aspect of system design; it has been suggested that interference tolerant modulation methods and high-quality antennas be used as a means for obtaining a large amount of communication through a given area.

---

* Such an interfering signal might be produced, for example, during heavy rain by a cross-polarized co-channel signal used to derive a second channel.
† This statement is true when the repeater must operate in the presence of severe interference; if thermal noise is the only limitation, a repeater using analog modulation methods can be designed to be more efficient.

VI. ACKNOWLEDGMENTS

In addition to people identified in the text and the references to other papers in this issue of the B.S.T.J., R. Kompfner has contributed to many fruitful discussions and has made specific suggestions, particularly concerning over-all repeater design and appearance.

REFERENCES

1. Tillotson, L. C., "A Model of a Domestic Satellite Communication System," B.S.T.J., 47, No. 10 (December 1968), pp. 2111–2136.
2. Hogg, D. C., "Millimeter-Wave Communication Through the Atmosphere," Science, 159, No. 3810 (January 5, 1968), pp. 39–46.
3. Semplak, R. A., and Keller, H. E., "A Dense Network for Rapid Measurement of Rainfall Rate," B.S.T.J., this issue, pp. 1745–1756.
4. Freeny, Mrs. A. E., "Statistical Treatment of Rain Gauge Calibration Data," B.S.T.J., this issue, pp. 1757–1766.
5. Semplak, R. A., and Turrin, R. H., "Some Measurements of Attenuation By Rainfall at 18.5 GHz," B.S.T.J., this issue, pp. 1767–1787.
6. Freeny, Mrs. A. E., and Gabbe, J. D., "A Statistical Description of Intense Rainfall," B.S.T.J., this issue, pp. 1789–1851.
7. Pierce, J. R., and Lewis, W. D., unpublished work.
8. Feldman, C. B., and Bennett, W. R., "Bandwidth and Transmission Performance," B.S.T.J., 28, No. 3 (July 1949), pp. 490–595.
9. Ruthroff, C. L., and Tillotson, L. C., "Interference in a Dense Radio Network," B.S.T.J., this issue, pp. 1727–1743.
10. Ruthroff, C. L., Osborne, T. L., and Bodtmann, W. F., "Short Hop Radio System Experiment," B.S.T.J., this issue, pp. 1577–1604.
11. Crawford, A. B., and Turrin, R. H., "A Packaged Antenna for Short Hop Microwave Radio Systems," B.S.T.J., this issue, pp. 1605–1622.
12. Osborne, T. L., "Design of Efficient Broadband Varactor Upconverters," B.S.T.J., this issue, pp. 1623–1649.
13. Osborne, T. L., Kibler, L. U., and Snell, W. W., "Low Noise Receiving Downconverters," B.S.T.J., this issue, pp. 1651–1663.
14. Bodtmann, W. F., and Guilfoyle, F. E., "Broadband 300 MHz IF Amplifier Design," B.S.T.J., this issue, pp. 1665–1686.
15. Bodtmann, W. F., "Design of Efficient Broadband Variolossers," B.S.T.J., this issue, pp. 1687–1702.
16. Chang, R. W., and Freeny, S. L., "Hybrid Digital Transmission Systems— Part 1: Joint Optimization of Analog and Digital Repeaters. Part 2: Information Rate of Hybrid Coaxial Cable Systems," B.S.T.J., 47, No. 8 (October 1968), pp. 1687–1711.
17. Prabhu, V. K., "Error Rate Considerations for Coherent Phase-Shift Keyed Systems with Co-Channel Interference," B.S.T.J., 48, No. 3 (March 1969), pp. 743–767.

# Short Hop Radio System Experiment

By C. L. RUTHROFF, T. L. OSBORNE, and
W. F. BODTMANN

(Manuscript received December 31, 1968)

*This paper reports a two-hop system experiment, designed to demonstrate the viability of the concept described in previous papers. A single RF channel, instrumented in the 11 GHz common carrier band, is transmitted from a terminal to a repeater one and a half miles away. After passing through the repeater, the signal is transmitted back to the terminal. We discuss the system parameters in relation to the requirements of operation at frequencies above 10 GHz; we also describe design, construction, and performance, as well as the lessons learned from one year of operation.*

## I. INTRODUCTION

Because the ever-increasing need for more communications circuits is crowding the frequency spectrum below 10 GHz, it is imperative to study the problems associated with microwave radio systems at frequencies above 10 GHz. The problems involved in the use of these frequencies and the resulting system concepts and constraints are discussed in a companion paper.[1] To bring these system concepts closer to reality and to find out whether the associated hardware is practical in the light of present technology, an experimental radio line has been built which embodies many of the essential system concepts. This paper describes the experimental system and its relation to the fundamental system concepts.

The experimental system is a research experiment on a system level and not a developmental radio system. Consequently many circuits and measurements necessary for a developmental system, such as order wire, alarm circuits, detailed cross modulation, and thermal noise measurements, were not made. Rather, efforts were confined to the more basic and unproven aspects of the radio line part of the system concept.

1577

II. GENERAL

2.1 *Rain Attenuation and System Design*

The main consideration in designing radio systems operating at frequencies above 10 GHz is the attenuation caused by rainfall. During a uniform rainfall of four inches per hour the rain attenuation in the 11 GHz common carrier band is 6 dB per mile.[1] At higher frequencies the attenuation increases sharply; at 18 and 30 GHz the attenuations are 15 and 30 dB per mile, respectively. Such heavy rain occurs infrequently, but to achieve the reliability required of common carrier service, radio systems must be designed so that such losses do not cause excessive system outages.

An important system constraint follows immediately from the magnitude and nature of rain attenuation: *The repeater path length must be short compared with the path lengths used at lower frequencies.* Since the rain attenuation increases in decibels per mile, the path length cannot be extended by the brute force method of increasing the transmitter power. For example, with a 30 GHz system designed to work through a uniform rainfall rate of four inches per hour, an increase in path length of one mile would require a thousand-fold increase in transmitter power.

Many fundamental system characteristics follow directly from the requirements of closely spaced repeaters. These are shown in Fig. 1 which illustrates the necessary and desirable characteristics of a viable radio system operating at frequencies above 10 GHz.

The concept of an attractive low cost tower with no building, and located on a public right-of-way, reflects a fundamental change in the philosophy of radio system design whereby tower and site performance is relaxed in favor of more adaptable and higher performance electronics. One approach is to house the entire repeater in a small enclosure on top of a metal pole which is just strong enough to keep the receiving antenna within the beam width in the strongest wind. With the repeater at the top of the pole, the need for a building and for long waveguide runs is eliminated, keeping right-of-way requirements to a minimum. The electronics must be extremely reliable, maintenance free, and able to operate in the outdoor environment. Since a system has many repeaters, there will be many interference exposures, and an interference resistant modulation method is needed requiring a broadband channel.[2] Figure 2 is a block diagram of a one-way repeater; a repeater layout, using a two frequency plan and two polarizations, is illustrated in Fig. 3.

FREQUENCIES ABOVE 10 GHz

RAIN ATTENUATION

SHORT HOPS

MANY REPEATERS PER SYSTEM

MANY INTERFERENCE EXPOSURES

MANY DISTORTION EXPOSURES

LOW COST RELIABLE REPEATER

GOOD APPEARANCE

NARROW BEAM ANTENNA WITH LOW SIDELOBES

LOW POWER CONSUMPTION

SOLID-STATE ELECTRONICS

SITES ON RIGHT-OF-WAY

LOW COST TOWER

NON-LINEAR POWER AMPLIFIER

NO BUILDING

SHIELDED ANTENNA

DISTORTION AND INTERFERENCE RESISTANT MODULATION

ANGLE MODULATION

DIGITAL: PSK-PCM
ANALOG: LARGE INDEX FM OR PM

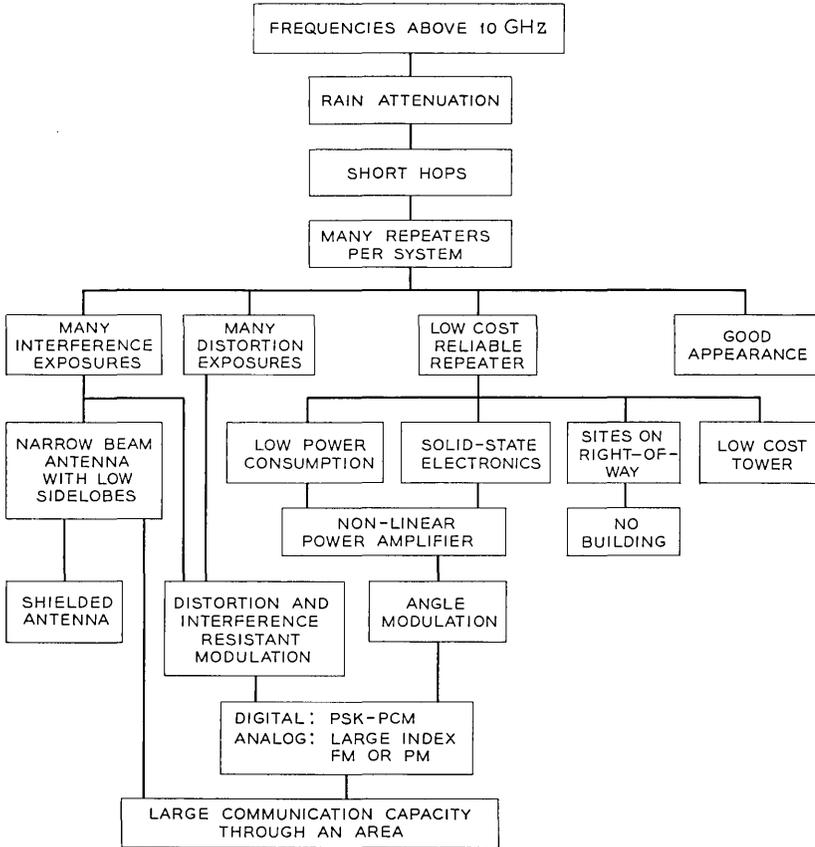LARGE COMMUNICATION CAPACITY THROUGH AN AREA

Fig. 1 — Characteristics of radio systems operating at frequencies above 10 GHz.

Illustrative parameters of an 11 GHz system of this type are summarized in Table I. The most significant parameters are: the short path length because of rain attenuation; the relatively low transmitted power, a consequence of the solid state and lower power consumption requirements; the receiver noise figure, required because of the fade margin and low transmitted power; the low dc power consumption; the bandwidth, which follows from the requirement of interference resistant modulation; and the intermediate frequency.

The gain of the varactor upconverter power amplifier is one of the main considerations in the selection of the system intermediate frequency. A low intermediate frequency favors low power consumption and increases the gain of the upconverter power amplifier. The
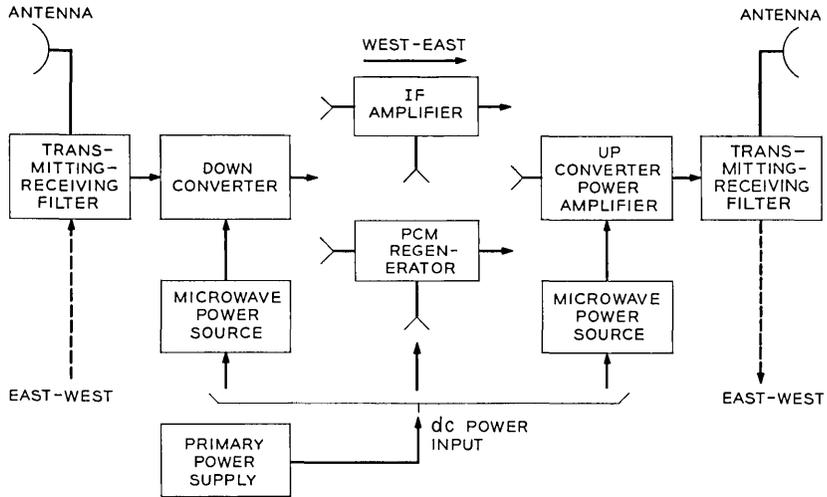
Fig. 2 — Block diagram of a one-way repeater.

300 MHz intermediate frequency chosen for this application is low enough to achieve substantial power gain and is still sufficiently high to obtain the desired 120 MHz IF bandwidth. To obtain an RF channel bandwidth of 120 MHz, a 1 dB bandwidth of at least 120 MHz is required for the upconverter, downconverter, and IF amplifier.

2.2 *Purpose of System Experiment*

Because of the new concepts and the device requirements involved, it is necessary to perform a system experiment to demonstrate that such a system can be realized. Accordingly, the system experiment described in this paper was performed to demonstrate that:

(*i*) The necessary devices, circuits, and hardware can be built,

(*ii*) Satisfactory operation in outside weather conditions is possible,

(*iii*) It is possible to operate a repeater on a few watts of dc power from a thermoelectric generator,

(*iv*) A pole mounted solid state repeater can be reliable and have long life,

(*v*) The appearance of a pole mounted repeater is acceptable,

(*vi*) And to reveal unforeseen problems.

2.3 *General Description of the Experimental System*

The experimental system has as many of the characteristics of a complete system as practical. It was not necessary to duplicate all

the transmitters and receivers and we decided that the channel com-
bining networks need not be designed. Consequently, for the experi-
ment, the repeater configuration shown in Fig. 3 was simplified to that
shown in Fig. 4. However, the repeater housing was designed to ac-
commodate a complete repeater consisting of four transmitters, four
receivers, channel combining networks, and two power supplies; the
units used were designed in the same plug-in form and size as would
be required in a complete repeater. The tower was designed to meet
the wind load and weight requirements of a complete repeater.

The radio system experiment is illustrated by Fig. 5. It consists of a
terminal transmitter and receiver on Crawford Hill and a single
channel repeater near the Holmdel laboratory 1.5 miles away. At each
location a receiver, transmitter, dc-dc converter-regulator, and an
antenna are housed inside a cylinder at the top of a 60-foot aluminum
pole. A second, empty housing is mounted on the Holmdel pole to
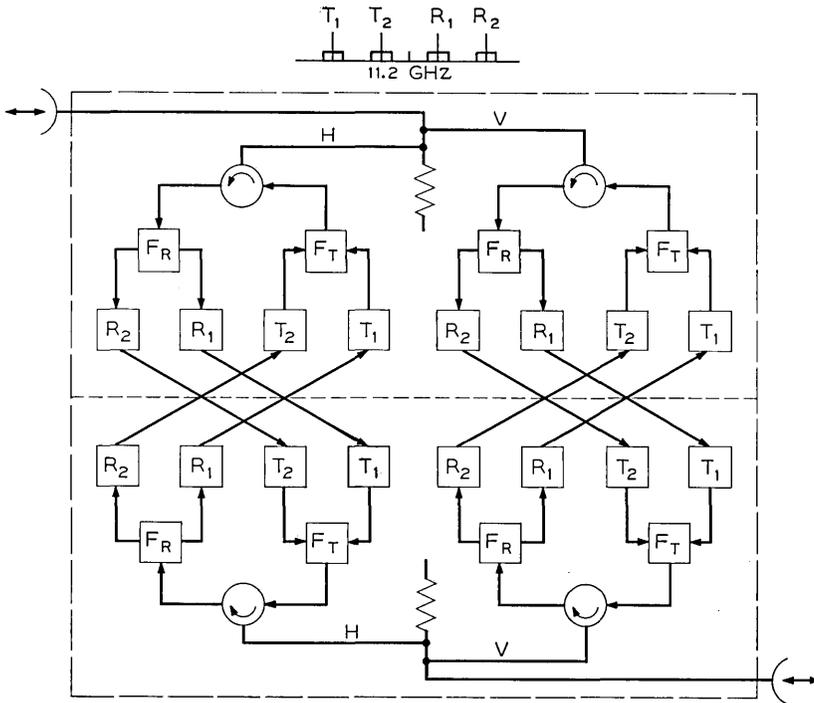simulate a complete repeater. Propane gas for the thermoelectric



Fig. 3 — Block diagram of a repeater with four two-way RF channels, a two
frequency plan, and two polarizations.

TABLE I — ILLUSTRATIVE REPEATER PARAMETERS FOR
AN 11 GHz SHORT HOP RADIO SYSTEM

| | |
|---|---|
| Radio frequency | 10.7 to 11.7 GHz |
| Transmitted power | +14 dBm |
| Antenna diameter | 2.5 feet |
| Antenna beamwidth | 2.5 degrees |
| Antenna gain (antenna efficiency = 0.5) | 36 dB |
| Path length | 3 miles |
| Section loss | 54 dB |
| Received signal power | −40 dBm |
| Fading margin (6 inches per hour uniform rate) | 30 dB |
| Receiver noise figure | 7 dB |
| Intermediate frequency | 300 MHz |
| IF bandwidth | 120 MHz |
| Width of base band | 4.5 MHz |
| Approximate voice circuit capacity per RF channel | 1000 channels |
| dc power consumption per 2-way RF channel | 10 W |

generator, mounted near the bottom of the pole, is kept in two underground tanks. At Crawford Hill, a building near the pole is used as a terminal building for test purposes. The Holmdel site has no building or commercial power; it is an isolated repeater. At the terminal the IF signals are connected to the test equipment through 100-foot coaxial cables which run down the inside of the pole and underground into the
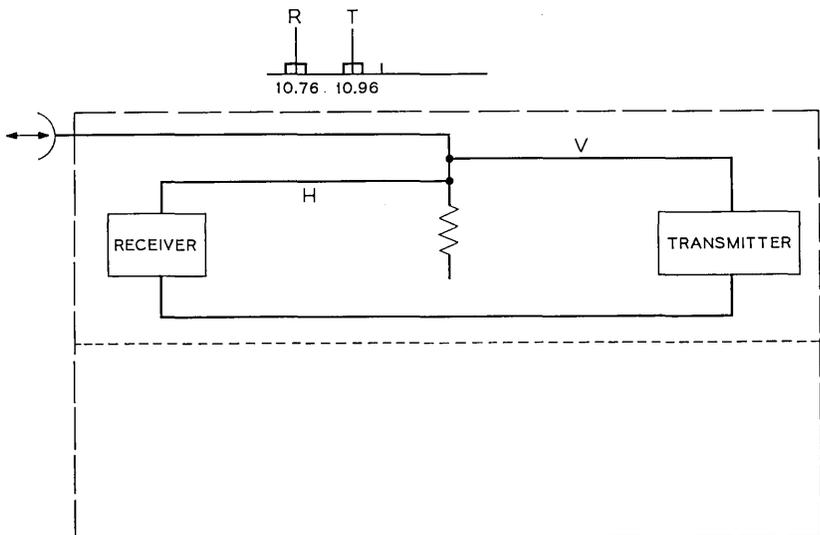


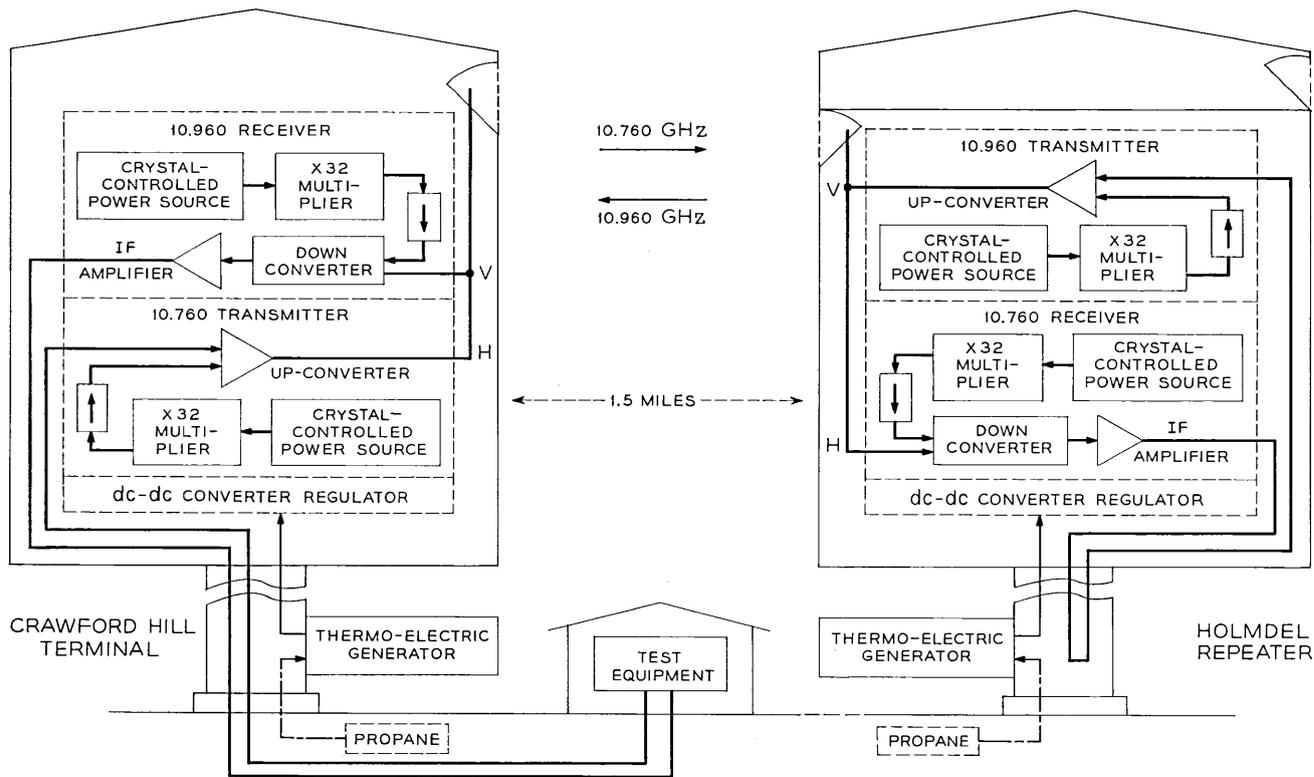Fig. 4 — Block diagram of the experimental system repeater.

Fig. 5 — Pictorial block diagram of the system experiment.

building. At the Holmdel repeater, the IF connection is made at the bottom of the pole for test purposes; normally the connection would be made at the top of the pole running only from one housing to the other. The Crawford Hill to Holmdel transmission is horizontally polarized and the return transmission is vertically polarized.

To check the operation of the system, the transmitter power, receiver AGC voltage, dc supply voltages, and thermoelectric generator power, can be monitored at both terminal and repeater. At the terminal, temperature sensors monitor the hottest points in the power supply, transmitter, and receiver, and the inside and outside ambient temperatures.

Satisfactory operation from −40°F to +135°F was set as a system objective. This was chosen as being approximately the range of temperatures common within the contiguous United States. In field use, no one repeater would normally encounter such a wide range. Because of a stability problem at the interface between the driver and X32 multiplier chain in the microwave power sources, the range actually achieved was −30°F to +125°F although all other circuits met or exceeded the original objective. This problem has subsequently been solved and microwave power sources suitable for use in the system have operated over the −40°F to +135°F temperature range.

III. REPEATER ELECTRONICS

3.1 *Basic Circuits*

3.1.1 *The dc Power Circuits*

The dc power circuits generate the required voltages from a propane gas thermoelectric generator, and provide both source and load regulation. Figure 6 shows the circuits and power levels required for a one-way repeater.

Electrical power is generated by a thermoelectric generator rated at 28 watts minimum at 75°F. Because of the increase in cold junction temperature with the ambient temperature, the output power decreases as the outside temperature increases. The generator power rating was chosen so that system requirements would just be met at the highest temperature of +140°F. At the fuel consumption rate of one gallon per day, operation for two thirds of a year without refueling is possible with two 120 gallon tanks.

The converter-limiter is a component of the thermoelectric generator which converts the 4 volt generator voltage to approximately 26.5

| TRANSMITTER | dc POWER | 3.0 WATTS |
| RECEIVER | dc POWER | 2.2 WATTS |

CONVERTER – REGULATOR

VOLTAGES −24 −6.5
POWER OUT 4 WATTS 1 WATT
REGULATION ±2% ±2%
EFFICIENCY 25%

CONVERTER LIMITER

EFFICIENCY 80%
POWER 20 WATTS
VOLTAGE 26.5 VOLTS

THERMOELECTRIC GENERATOR

POWER, 75° F 28 WATTS MIN
GENERATOR 3M − 515 OR GI − U3P
FUEL 1 GAL / DAY, $ 60 / YR
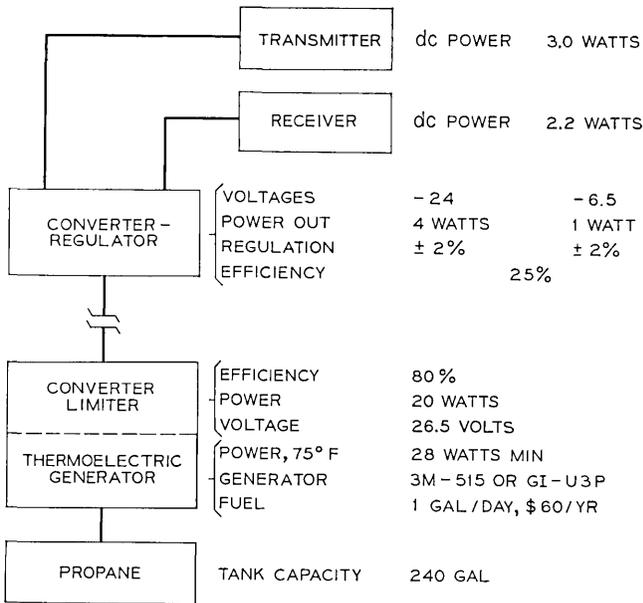
PROPANE   TANK CAPACITY   240 GAL

Fig. 6 — The dc power circuits and power levels.

volts, a more useful value for transistors and distribution circuits, and limits the converter voltage to that value if the load is reduced. The limiter also provides some regulation of converter and generator power changes.

The converter-regulator, also a commercial unit, is located in the repeater housing and supplies −24 volts regulated against load and supply variations. The −6.5 volts is derived from the −24 volts by use of a regulator diode. Both voltages are regulated to less than ±2 percent for combined load changes, input power changes, and temperature variations.

A total power of 5.2 watts is used for one transmitter and one receiver. Thus dc power requirements are approximately 10 watts per two-way RF channel. The overall dc conversion efficiency is only 20 percent because of the use of two commercially available converters and the method of deriving the −6.5 volts. A power converter has been reported which uses a variable pulse width technique to regulate against input voltage changes and provides series regulators to regulate against load changes; efficiencies as high as 69 percent were obtained.[3] This technique of combining series regulators with

preregulation appears well suited for this application where the converter can transform the thermoelectric generator voltage to the several output voltages required for the system and regulate against generator power variations. The voltages can be transmitted to series regulators located at the radio package to regulate against load changes. Overall efficiencies of 60 to 70 percent should be obtained since the power levels are about 25 watts. Assuming 10 watts per two-way channel and 50 percent converter efficiency, an 80 watt thermo-electric generator would be required for the complete repeater in Fig. 3.

### 3.1.2 *Microwave Power Sources*

The microwave power sources generate the microwave power required to pump the varactor upconverters and the Schottky-barrier diode downconverters. Figure 7 is a block diagram of a power source and typical results for transmitter and receiver sources.

The microwave power source consists of a 327 MHz crystal-controlled power source followed by a X32 multiplier chain. The frequency stability is set by the crystal controlled oscillator and is well within the ±0.005 percent requirements over the temperature range. The power amplifier is operated class C for maximum efficiency since



|  | TRANSMITTER | RECEIVER |
|---|---|---|
| OUTPUT FREQUENCY | 10.460 GHz | 10.660 GHz |
| OUTPUT POWER | 19.7 dBm | 9.8 dBm |
| RF EFFICIENCY | 8.0 dB | 10.3 dB |
| OUTPUT FREQUENCY | 326.875 | 333.125 MHz |
| OUTPUT POWER | 27.7 dBm | 20.0 dBm |
| dc POWER | 2.36 WATTS | 1.15 WATTS |
| EFFICIENCY | 25% | 9% |
| OVERALL | | |
| EFFICIENCY | 4% | 0.8% |
| TEMPERATURE RANGE | −30°F TO +125°F | |
| STABILITY | 0.005% | |

Block diagram labels: X4 MULTIPLIER, X8 MULTIPLIER, POWER AMPLIFIER, X4 MULTIPLIER, BUFFER AMPLIFIER, CRYSTAL–CONTROLLED OSCILLATOR, CRYSTAL–CONTROLLED POWER SOURCE
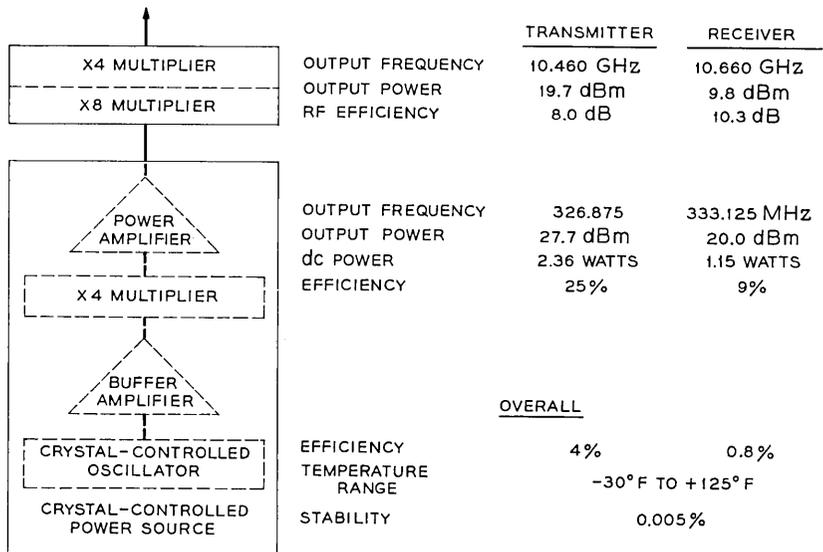
Fig. 7 — Microwave power source.

it is at the highest power level in the chain. A X8 varactor multiplier connected directly to a X4 multiplier converts the 327 MHz to the required 10.460 GHz. The overall efficiency of the high power microwave source is 4 percent.

Since the transmitter microwave power source uses almost half of the dc power required for a one-way channel, it was designed for maximum efficiency subject to the constraints of environmental stability and reliable starting. For this experiment the receiver source consists of a high power unit adjusted for the reduced power required by the receiver. A separate design of the low power source would result in efficiencies at least as good as obtained with the high power source.

In accordance with previously established system concepts, the microwave power source must be efficient, reasonably simple, and compact—attributes which tend toward low cost. Consequently, very efficient circuits with high order multiplication and no isolators were used. However, isolators are commonly used to solve the stability and starting problems inherent in nonlinear circuits; both of which are aggravated by efficient circuits and temperature and voltage variations. Therefore the design of a chain without isolators and their obvious disadvantages required theoretical and practical solutions to the problems of stability and starting.

Stability criteria for high order varactor multipliers have been established by Dragone and Prabhu.[4,5] With the aid of these criteria the X32 chain was successfully designed and built without the use of isolators. In the connection of the crystal controlled power source to the X32 chain the same instability mechanism was encountered because of the nonlinear class C power amplifier and X4 multiplier in the crystal controlled power source. It can easily be demonstrated that phase modulation is propagated in both directions through the class C power amplifier and thus through the first X4 multiplier. Therefore, the total instability feedback path includes a phase modulation gain of 128 and the round-trip phase gain of the power amplifier which can be large.

In the microwave sources actually used in the system these stability problems limited the temperature range to −30°F to +125°F. Subsequently, procedures for power amplifier design and adjustment have been found which insure that the amplifier round trip phase gain is less than unity and that the amplifier supplies sufficient power into the varactor multiplier at all power levels to insure reliable starting. Microwave power sources, similar to those used in the system, have

operated over the temperature range of −40° to +135°F with a small penalty in efficiency.

### 3.1.3 Receiving Downconverter

The receiving downconverter converts the low level received microwave signal to the intermediate frequency with a minimum of noise degradation and sufficient gain to reduce the effect of the main IF amplifier noise. The downconverter is an image rejection balanced mixer using Schottky-barrier diodes, followed by an integrally connected low noise preamplifier.[6]

The combination of intermediate frequency, bandwidth, and noise figure, with the other system constraints, require extending downconverter performance beyond existing technology. An average noise figure of 5.6 dB was obtained; Table II summarizes the results.

### 3.1.4 IF Amplifier

The IF amplifier provides most of the repeater gain and the automatic gain control necessary to maintain a constant transmitted power. This is accomplished with three low level variable-gain stages and three high level fixed-gain stages.[7,8]

The amplifier has a 1 dB bandwidth of 120 MHz. With a normal input signal level, the gain is 33 dB with an output power of +7.5 dBm into 50 ohms with 1 dB gain compression. The automatic gain control keeps the output power constant over a 43 dB range in the input level.

Most of the gain and all of the gain control is accomplished in the low level section to minimize the total power consumption and noise figure. The total dc power consumption is 0.65 watt. With 33

TABLE II — SUMMARY OF DOWNCONVERTER PERFORMANCE

| | |
|---|---|
| Radio signal frequency | 10.760 GHz |
| RF signal power (normal) | −40 dBm |
| Pump frequency | 10.460 GHz |
| Pump power (for 2 diodes) | +9.8 dBm |
| Intermediate frequency | 300 MHz |
| IF power | −21.4 dBm |
| RF to IF gain | 18.6 dB |
| Frequency response | flat ± 0.1 dB, 240 to 350 MHz |
| | 0.3 dB down at 230, 360 MHz |
| Noise figure, 75°F | 5.3 dB at 300 MHz |
| | 5.6 dB average, 240 to 360 MHz |
| Preamplifier noise figure | 2.0 dB |
| Temperature range | −40 to +140°F |

dB gain the noise figure is 13.6 dB, contributing only 0.4 dB to the receiver noise figure.

### 3.1.5 *Varactor Upconverter Power Amplifier*

The varactor upconverter power amplifier converts the 300 MHz IF signal from the main IF amplifier to the 11 GHz RF band with sufficient power for transmission over the radio path. It consists of a varactor diode mounted across the waveguide, an integrally attached IF amplifier, and two waveguide filters which provide the proper reactive terminations.[9]

Basic considerations in the design of the upconverter are the bandwidth and the pump efficiency. The bandwidth was obtained by mismatching the interface between the amplifier and the varactor diode in an optimum way. The upconverter was designed for maximum pump efficiency in order to get the maximum output power possible with the available pump power.

The bandwidth obtained was slightly less than the desired 120 MHz between 1 dB points because of a slight asymmetry. The output power is typically 16.2 dBm with a pump power of 20 dBm; the IF to RF gain is 13.2 dB, including the amplifier, when the amplifier input power is +3 dBm.

### 3.2 *Receiver*

The receiver consists of a low power microwave power source, a downconverter, and a main IF amplifier as shown in Fig. 5. An isolator is required between the output of the microwave power source and the narrowband pump input filter of the downconveter to satisfy the stability criteria for multiplier chains.[4,5] The output of the preamplifier is connected to the input of the main IF amplifier through a very short connection, easing the output return loss requirement on the preamplifier.

Since the available transmission path is approximately one-half of the nominal design length of 3 miles, a 6 dB pad was added at the receiver input to simulate the additional 1.5 miles. The actual downconverter gain is 4 dB more than the design value; so rather than operate the variolossers at their limit, a 4 dB pad was added to the main IF amplifier input. With these adjustments, typical receiver power levels are: received power at the down converter, −40.2 dBm; IF power out of the down converter, −21.6 dBm; IF power out of the main IF amplifier, +7.5 dBm. The power out of the microwave source,

less 0.4 dB isolator loss, is 9.4 dBm at the local oscillator port of the down converter. Using the 300 MHz, 75°F, spot noise figures of 5.3 dB for the down converter and 13.6 dB for the main IF amplifier, and the 18.6 dB gain of the down converter, the overall receiver spot noise figure is 5.7 dB. The addition of the 4 dB pad degrades the noise figure to 6.2 dB.

The transmission frequency response of the complete receiver is shown in Fig. 8 for temperatures of +140°F, +75°F, and −40°F and at each temperature for input power levels in the range of −40 to
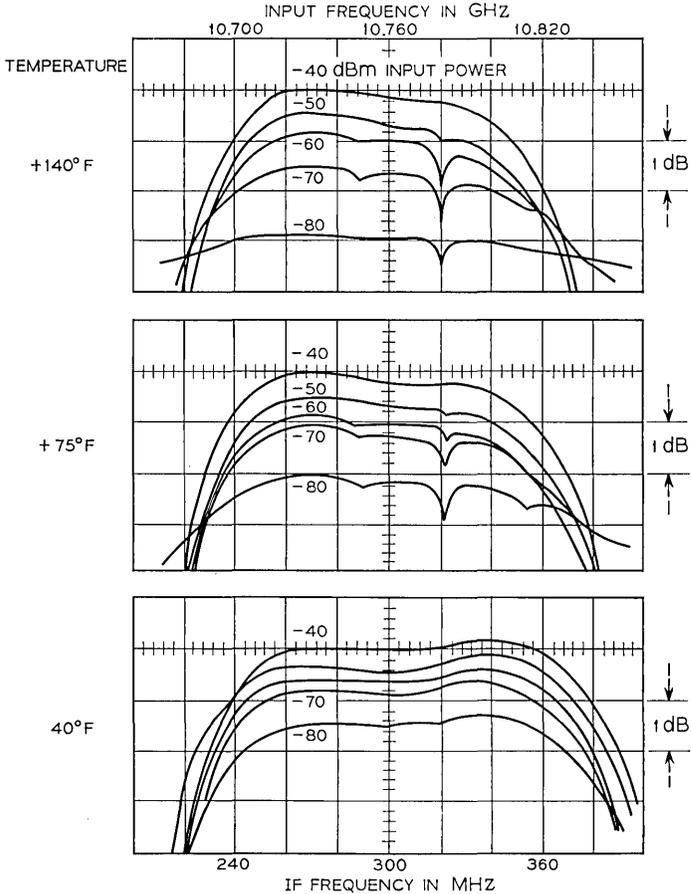
Fig. 8 — Frequency response of a complete receiver at +140°F, +75°F, and −40°F, and −40 to −80 dBm input power levels.

−80 dBm. At 75°F and −40 dBm input level, the response is down 1 dB at 240 and 360 MHz, but has a small slope across the band. At +140°F the slope increases causing the response to be down 1.8 dB at 360 MHz; at −40°F the slope is reversed causing the response to be 0 dB down at 360 MHz. These response variations are basically those of the main IF amplifier with some change resulting from changes in pump power into the down converter. The dip in the response at about 320 MHz is caused by an interfering tone which is the 33rd harmonic of the 327 MHz power source present at the output of the X32 multiplier chain. Its equivalent input level is −87 dBm and is a function of the selectivity of the down converter pump input filter.

Figure 9 shows the 10.960 GHz receiver. The entire receiver is mounted on a 12 × 12 × 2.5-inch aluminum frame and weighs 14.3 pounds. This layout is designed to allow installation or replacement by simply plugging in the unit. Waveguide alignment is assured by using a short piece of flexible waveguide with alignment pins. Choke joints with metalized gaskets are used to prevent leakage.

The crystal-controlled power source is insulated to conserve the heat generated within the box and reduce the effects of high relative humidity. For example, if the outside air is saturated at 80°F, an 8°F increase in temperature reduces the relative humidity inside the box to about 70 percent.

### 3.3 Transmitter

The transmitter contains a high power microwave power source and a varactor upconverter power amplifier connected through an isolator. The output of the IF amplifier in a repeater or the terminal equipment is applied to the IF input of the upconverter at a +3 dBm level. The 19.7 dBm out of the microwave power source, less 0.4 dB isolator loss, gives 19.3 dBm of pump power at the upconverter, for which the output power is about 15.5 dBm. The frequency response of the transmitter is the same as that of the upconverter.

Figure 10 shows the 10.760 GHz transmitter. The transmitter weighs 12 pounds and is mounted on the same type of aluminum frame as the receiver. The flexible waveguide used for waveguide alignment and the dc power plug are on the right side of the unit.

### 3.4 Antenna

The antenna consists of a small aperture feed, a focussing paraboloidal reflector, and a plane reflector in an inverted periscope ar-
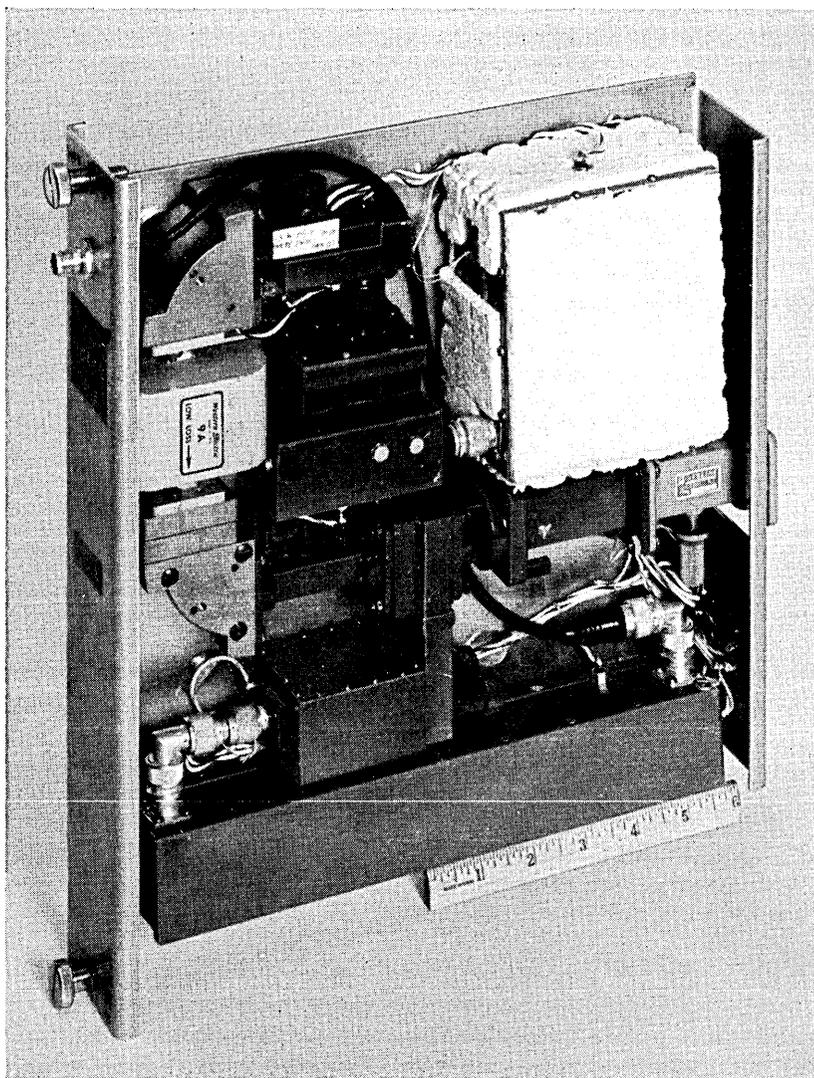
Fig. 9 — The 10.960 GHz receiver.

rangement. The axis of the feed and parabola coincide with the axis of the housing which supports and shields the radiating elements. Thus the packaged antenna has the highly desirable low side lobe and rear lobe radiation characteristics necessary for radio interference suppression.

The important characteristics of the antenna are listed in Table III

and extensive measurements are reported in a companion paper.[10] Of particular importance are the return loss and cross polarization coupling loss because they determine the amount of filtering required to isolate adjacent and cross polarized channels and thereby affect system channel capacity.
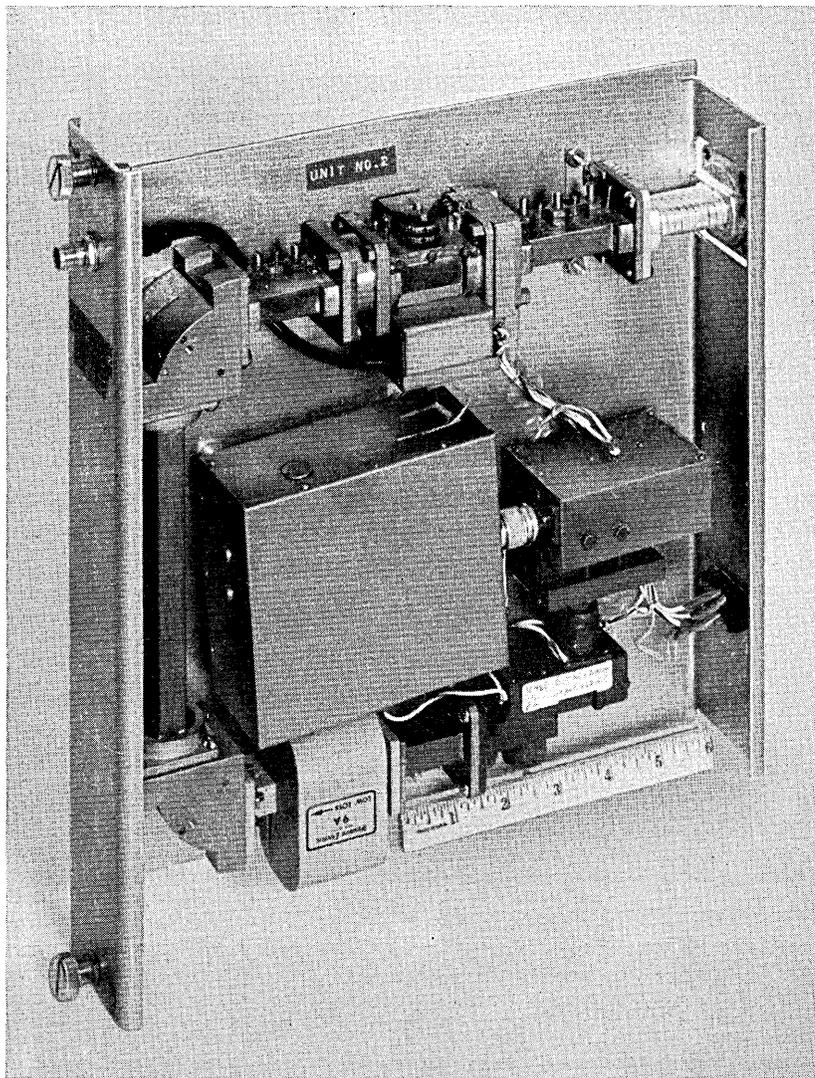


Fig. 10 — The 10.760 GHz transmitter.

TABLE III — MAJOR CHARACTERISTICS OF THE ANTENNA
WITH WEATHER COVER

| | |
|---|---|
| Aperture diameter | 2.5 feet |
| Beamwidth, 3 dB | 2.4 degrees |
| Gain (57% aperture efficiency) | 36.6 dB |
| Return loss (10.7 − 11.7 GHz) | >18 dB |
| Cross polarization coupling | >35 dB |
| Cross coupling between packages | >75 dB |
| Side and rear radiation | >60 dB for > 120° |

3.5 *System Measurements and Performance*

Before final field installation, the receiver, transmitter, and converter-regulator were tested in the laboratory as an IF repeater. Figure 11 shows the RF-to-RF transmission frequency response of the Crawford Hill repeater at +140°F, +76°F, and −40°F. The band shape, especially the fine grain structure, is caused by the varactor upconverter. In-band ripples are less than 0.4 dB; the response is down less than 1.5 dB at the band edges. The usable bandwidth is much larger than 120 MHz; for example, the 3 and 6 dB bandwidths at 75°F are 185 and 220 MHz, respectively.

The distortion resulting from the amplitude response of the repeater transfer function has been computed for large index analog frequency modulation.[11] For a noise modulation bandwidth of 5 MHz, and an rms frequency deviation of 15 MHz, the computed signal-to-distortion ratio at the highest baseband frequency is approximately 56 dB. The room temperature response of Fig. 11 was used for the computation and a linear phase response was assumed.

Field installation was completed in mid-November 1967; operation was intermittent at first. The trouble, which was traced to humidity effects in the crystal-controlled power source, is discussed in Section V. The trouble was corrected and the system was placed in operation January 2, 1968. The electronics have operated continuously since then except for a few days when the system was off because of thermoelectric generator failures and in July when a lightning stroke on Crawford Hill damaged the main IF amplifier output transistors. These problems are also discussed in Section V and VI. During this period the temperature ranged from 0°F to 100°F, wind gusts as high as 65 miles per hour were recorded, and the humidity averaged 99 percent for as long as five days. When the electronic components were removed for repair in late July, there was no visible evidence of deterioration or corrosion and there was no deterioration of performance.

IV. REPEATER MECHANICS

4.1 *Repeater Housing*

The repeater housing is a weatherproof but well ventilated shelter for the repeater electronics and the antenna. It is a cylindrical aluminum canister 34 inches in diameter and 42 inches high. About half of the volume is used for the inverted periscope antenna which was designed for such applications. Figure 12 shows the antenna side of the housing. The volume under the plane reflector is sufficiently large
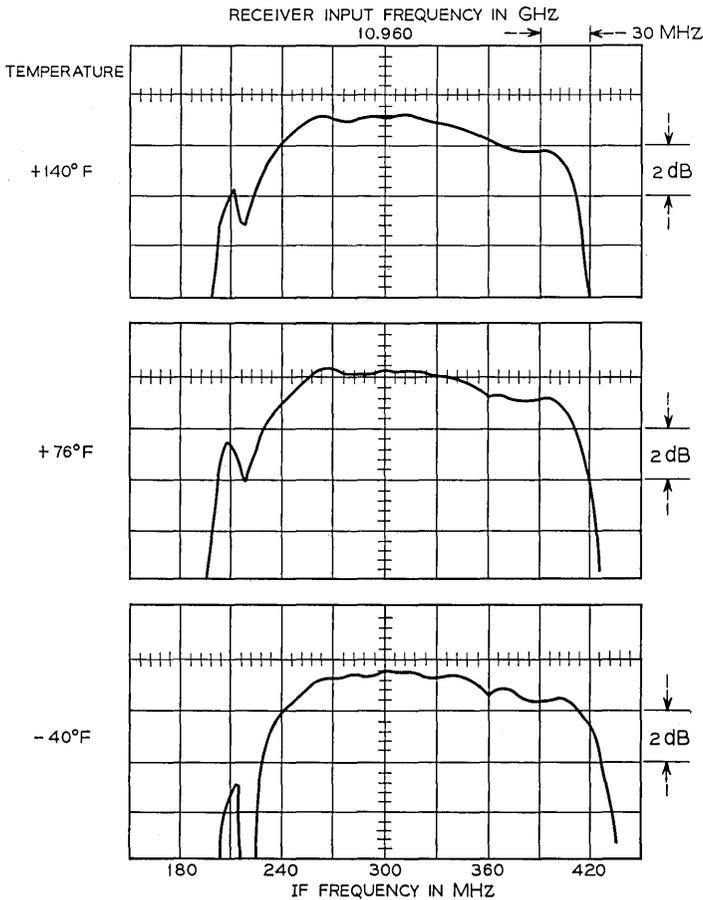


Fig. 11 — RF-to-RF transmission frequency response of the Crawford Hill repeater measured in the laboratory.
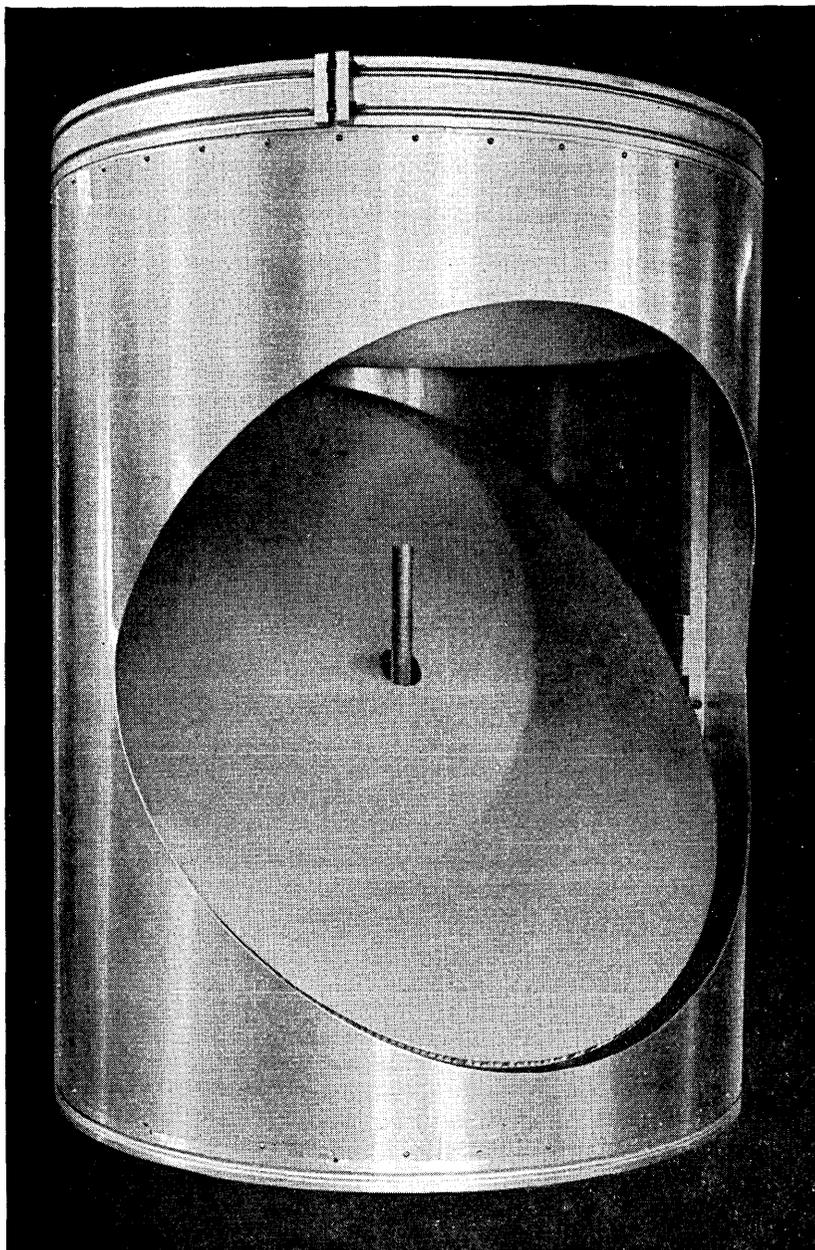
Fig. 12 — Repeater housing and antenna with weather cover removed.

for the four transmitters, four receivers, channel combining networks, and two power supplies of the system configuration of Fig. 3. Figure 13 shows the electronics side of the housing with access covers removed and four model plug-in units for the horizontally polarized channels installed. The channel combining networks are between the rear of the plug-in units and the back panel. The two power supply chassis fit in the shorter slots toward the outside of the housing.

Also visible in Figure 13 are the vent holes in the bottom of the housing; the larger hole in the center is used for IF interconnection cables between canisters. The canisters are attached to each other and to the mast by a grooved band which fits matching grooves on the top and bottom of each housing and on the mast flange. This permits azimuth alignment by rotating the entire housing. Elevation beam alignment is accomplished by tilting the antenna plane reflector; an elevation position indicator and vernier adjustment is located on the outside of the housing above the access covers. This method of alignment does not require movement of the electronics, antenna feed, or parabola and is one of the attractive features of the inverted periscope antenna.

Based on a wind loading area equal to two thirds of the projected area of the housing and a 100 miles per hour wind pressure of 40 pounds per square inch, the total wind load for each housing is approximately 300 pounds, or 600 pounds for a complete repeater. The weight of the housing and antenna without electronics is about 85 pounds which, using the previously mentioned weights for the plug-in units, gives a weight of about 200 pounds per housing or 400 pounds for a completely equipped repeater.

4.2 *Repeater Tower and Foundation*

The short repeater spacing allows the use of towers which need only be tall enough to clear the trees. Spun aluminum poles are ideally suited for this application for the following reasons: (*i*) they have an attractive and commonly accepted appearance, (*ii*) aluminum remains clean and attractive with no maintenance in contrast with non-stainless steels, an important consideration since repeaters will be located on public rights-of-way, (*iii*) spun aluminum poles are readily available because of their many other uses, and (*iv*) aluminum is cheaper than stainless steel for the same stiffness. A tapered pole is used because it has a more pleasing appearance and because truncated hollow cones generally are more efficient load carrying structures than cylindrical tubing in column and end-loading applications.[12]
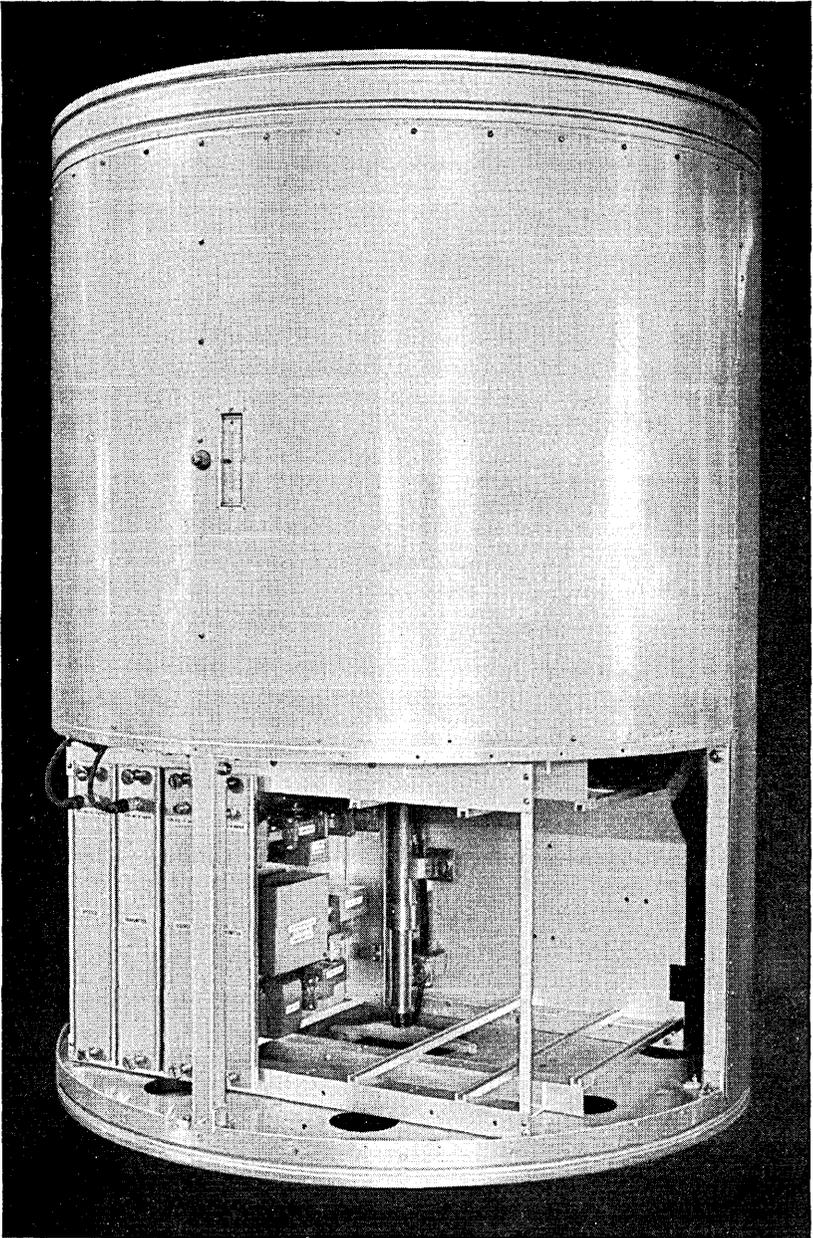
Fig. 13 — Repeater housing and repeater electronics with access covers removed.

The cost of the tower is proportional to its weight. For a given weight, the strength of the tower is a function of the diameter and the wall thickness. The parameter of importance is the end slope which is equal to the antenna beam deflection. These relationships are illustrated in Fig. 14 for tower diameters of 16, 20, and 24 inches. For a 1,650 pound tower, increasing the diameter from 16 to 24 inches decreases the beam deflection from 3 to 1.5 degrees in a 100 miles per hour wind.

The foundation must support the weight of the repeater and tower and the overturning moment caused by the wind load. For all but special situations, a foundation formed of a sufficiently deep steel reinforced concrete cylinder is preferred because it is simple, easily
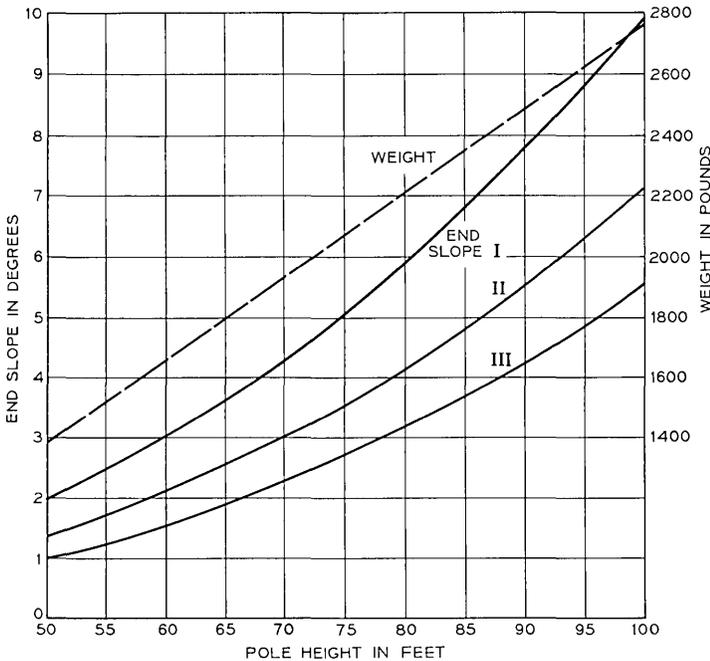


Fig. 14 — Tower end slope and weight as a function of height and diameter assuming a 100 miles per hour wind load of 600 pounds.

Material: Aluminum
Wind pressure: 40 psf

| Tower | $D_1$ | $D_2$ | $t$ |
|-------|-------|-------|--------|
| I | 16″ | 8″ | 0.625″ |
| II | 20″ | 10″ | 0.500″ |
| III | 24″ | 12″ | 0.417″ |

installed with a drill rig, and requires a minimum of ground area. A foundation of this type, 15 feet deep and 18 inches in diameter, was used for the Crawford Hill terminal. For soils with low bearing load, a concrete pad is required. A foundation of this type, 9 feet square and 2 feet thick, was used for the Holmdel repeater.

### 4.3 *Installation and Alignment*

Because of the number of repeaters in a system, it is obvious that field installation and alignment time should be minimal. The procedures used in the experimental system are: (*i*) complete adjustment and testing of the electronics in the shop, (*ii*) antenna alignment and elevation adjustment in the shop, (*iii*) field installation of the packaged antenna and repeater on the tower, (*iv*) erection of tower with complete repeater, (*v*) final vertical and azimuth adjustment of tower, and (*vi*) installation of thermoelectric generator and wiring. In the experimental system, physical alignment of the antenna elements in the shop was sufficient to assure electrical boresight to within 0.15 degree.[10] The towers of the experimental system, with electronics installed, were erected with a crane in about 1.5 hours per tower. Vertical and azimuth alignments were done with a theodolite using alignment marks on the top and bottom tower flanges. From a later measurement it was found that this mechanical alignment was accurate to within 0.3 degree in azimuth and 0.05 degree in elevation.

The heights of the towers in short hop systems are such that cherry-pickers can be used in antenna alignment and for site selection and access to the electronics.

### 4.4 *Appearance*

Figure 15 shows the Crawford Hill terminal. The antenna weather cover can be seen on the front of the single repeater housing on top of the aluminum pole. The thermoelectric generator is mounted on the side of the pole about 10 feet above the ground and the two propane fuel tank covers can be seen near the foundation. The control building on the left contains the test equipment.

The Holmdel repeater, shown in Fig. 16, has two packages to demonstrate the appearance of a complete repeater. Since it is powered by a thermoelectric generator with no standby source, there is no building or commercial power line associated with the repeater.

### V. DISCUSSION

The stability and starting problems of efficient chains of nonlinear amplifiers and varactor multipliers have already been discussed and

solutions are given by Dragone.[4] In all harmonic generators, harmonics adjacent to the desired output frequency are present to some degree in the output of the chain. This is the source of the interfering tone apparent in the amplitude response of Fig. 8. Furthermore, harmonics of the 87 MHz oscillator frequency generated by the first X4 multiplier appear as modulation of the fourth harmonic and are
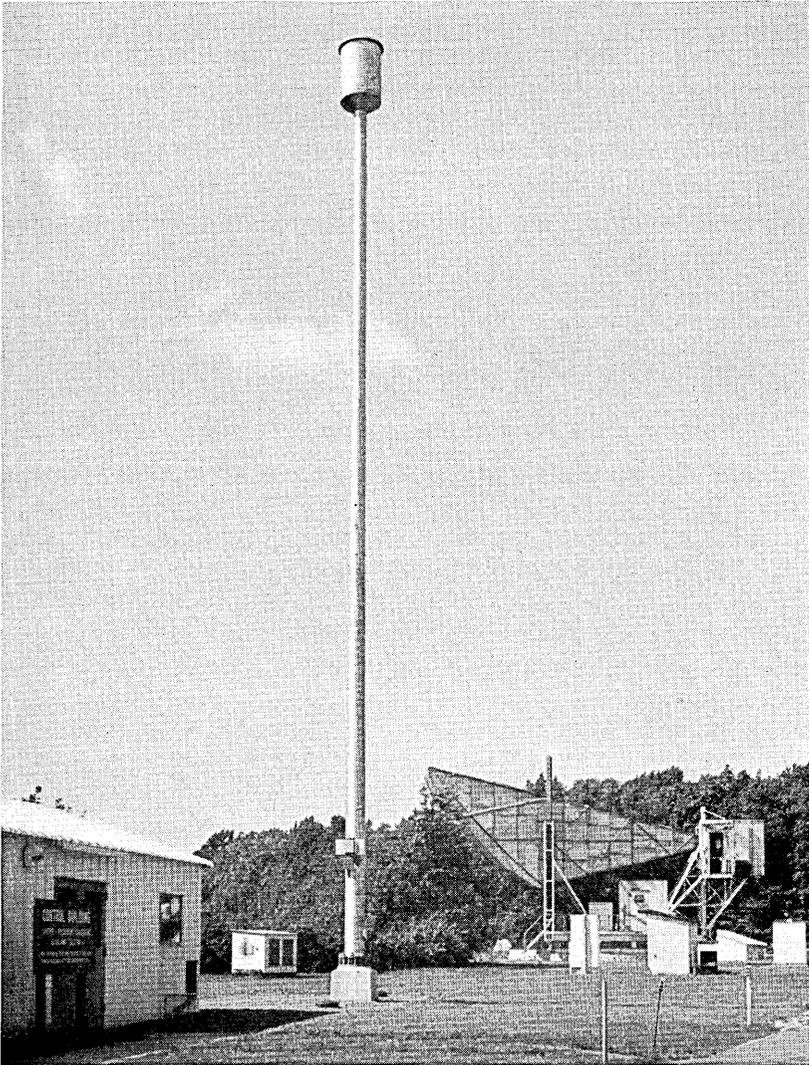
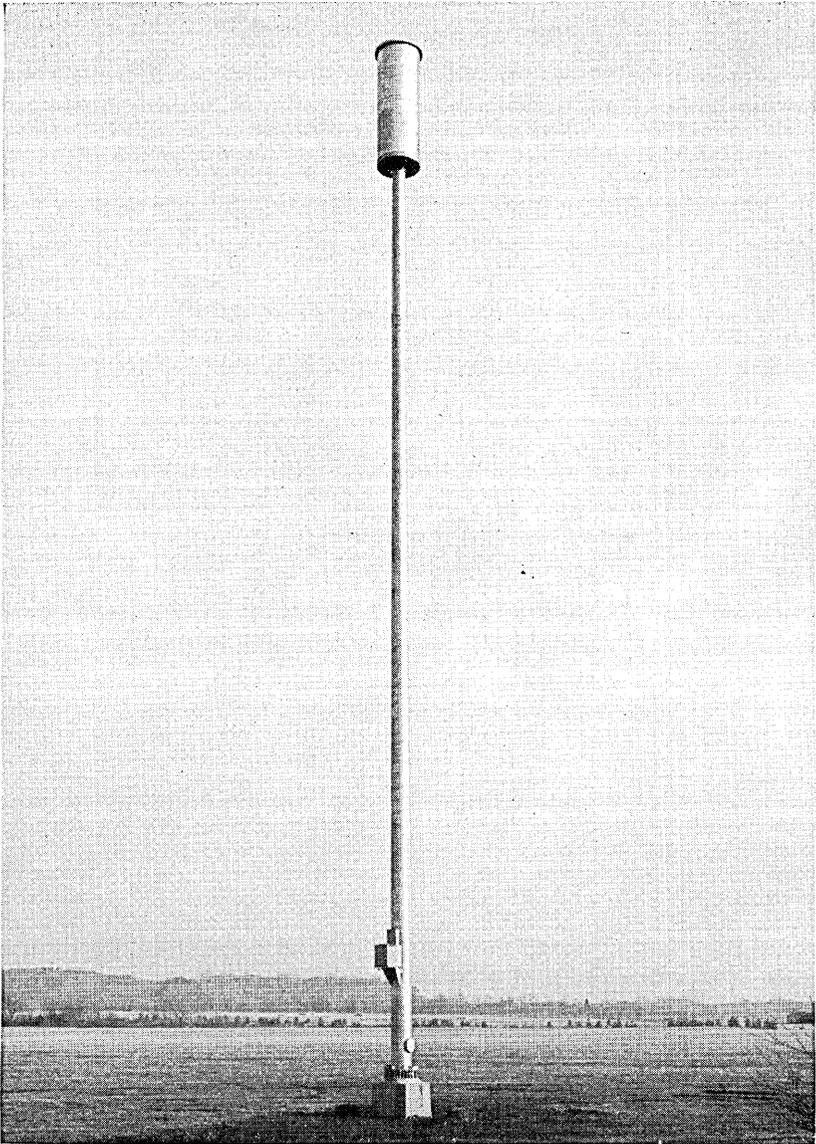Fig. 15 — The Crawford Hill terminal.

Fig. 16 — The Holmdel repeater.

transmitted through the X32 chain giving rise to additional tones at ±87 MHz around the X32 output harmonics. The best solution for these tone problems is better filtering at the output of the crystal-controlled source and at the pump input to the upconverter and downconverter. This would cost at most a 1 dB decrease in available pump power.

Immediately after field installation it was found that the crystal-controlled power source was sensitive to humidity. The sensitive elements were found to be certain ceramic trimmer capacitors which were particularly susceptible to contamination from handling and soldering. They were replaced by air capacitors which have glazed ceramic for support. As further protection, the low power unit was insulated to obtain several degrees rise in temperature which limits the maximum relative humidity inside the box to about 75 per cent. The system has since operated without degradation through several periods when the relative humidity averaged 99 per cent for several consecutive days.

Failures in thermoelectric generators have caused several interruptions of the system experiment. There are basically two types of failures: burner failures and thermopile failures. Most of the burner failures have been caused by an accumulation of sulfur from the propane in the burner nozzle orifice which changes the fuel-air mixture. Thermopile failures are caused by oxygen contamination and changes in junction material properties. Thermopiles with the required reliability have been made for military and space applications indicating that generators with the required reliability are within the realm of present technology.

VI. SUMMARY AND CONCLUSIONS

A two-hop radio system experiment, operating in the 11 GHz common carrier band, is described in detail. Careful attention has been paid to the fundamental characteristics of systems at frequencies above 10 GHz which must operate during periods of large rain attenuation; these characteristics are illustrated in Fig. 1.

The most important objectives of a system experiment are to demonstrate successful operation in the presence of difficult interface problems and in a real environment. The solution of interface and environmental problems by brute force, such as sprinkling isolators throughout the harmonic generator chain or enclosing sensitive components in constant temperature ovens, is straightforward; unfortunately such an approach could easily be fatal to the system concept

in cost and power consumption. The approach taken here has been to determine the performance of components with respect to basic circuit limitations and to approach this performance as closely as possible consistent with the overall system objectives. This long way around, requiring solutions to problems such as stability and starting of high order harmonic generator chains, optimizing varactor up-converters with respect to efficiency, bandwidth and gain, realizing broadband variolossers with low minimum loss, IF amplifiers with low power consumption, and low noise downconverters, is the surest way to good system design.

The system experiment has been in operation since January 2, 1968. Several interruptions caused by faulty operation of the thermoelectric generators occurred and a single interruption was caused by a lightning stroke on Crawford Hill which damaged the output transistors of the main IF amplifier. The inadequate lightning protection for the cables between the tower and the terminal building has been corrected. With these exceptions there has been no degradation in performance of the system.

Most of the objectives of the experiment, discussed in Section II, have been met. Long life and reliability cannot be demonstrated in one year but are certainly indicated by the stable performance of the repeater electronics. As to appearance, who will say?

REFERENCES

1. Tillotson, L. C., "Use of Frequencies Above 10 GHz for Common Carrier Applications," B.S.T.J., this issue, pp. 1563–1576.
2. Ruthroff, C. L., and Tillotson, L. C., "Interference in a Dense Radio Network," B.S.T.J., this issue, pp. 1727–1743.
3. Cassady, C. D., Gayer, D. G., and Leshe, C. M., unpublished work.
4. Dragone, C., "Phase and Amplitude Modulation in High-Efficiency Varactor Frequency Multipliers of Order $N = 2^n$ — Stability and Noise," B.S.T.J., 46, No. 4 (April 1967), pp. 797–834.
5. Dragone, C., and Prabhu, V., "Some Considerations of Stability in Lossy Varactor Harmonic Generators," B.S.T.J., 47, No. 67 (July–August 1968), pp. 887–896.
6. Osborne, T. L., Kibler, L. U., Snell, W. W., "Low Noise Receiving Downconverters," B.S.T.J., this issue, pp. 1651–1663.
7. Bodtmann, W. F., and Guilfoyle, F. E., "Broadband 300 MHz IF Amplifier Design," B.S.T.J., this issue, pp. 1665–1686.
8. Bodtmann, W. F., "Design of Efficient Broadband Variolossers," B.S.T.J., this issue, pp. 1687–1702.
9. Osborne, T. L., "Design of Efficient Broadband Varactor Upconverters," B.S.T.J., this issue, pp. 1623–1649.
10. Crawford, A. B. and Turrin, R. H., "A Packaged Antenna for Short-Hop Radio Systems," B.S.T.J., this issue, pp. 1605–1622.
11. Ruthroff, C. L., "Computation of FM Distortion in Linear Networks for Bandlimited Periodic Signals," B.S.T.J., 47, No. 6 (July–August 1968), pp. 1043–1063.
12. Schick, W. F., "How to Calculate End Slope and Deflection of Conical Tubular Beams," Machine Design, 36, No. 24 (October 8, 1964), pp. 193–194.

# A Packaged Antenna for Short-Hop Microwave Radio Systems

By A. B. CRAWFORD and R. H. TURRIN

*This paper describes a packaged antenna specially designed for mounting on a slender tapered aluminum mast and gives typical measured electrical characteristics. The antenna is used in a 1.5-mile experimental repeater installation operating at 11 GHz. The masts and foundations for the experimental system are described briefly.*

*The antenna package is an upright cylinder, a shape chosen to minimize mast twisting caused by wind, and to present a pleasing appearance in combination with a slender mast. The radiating elements consist of a waveguide aperture feed, a 30-inch parabolic reflector mounted with its axis vertical at the top of the package, and a 45-degree flat reflecting plate similar to an inverted periscope. The space below the 45-degree reflector houses all the repeater electronics. Because of the shielding effect of the cylindrical housing, this antenna like the horn-reflector antenna, has very low radiation in the far side and back lobe regions.*

## I. INTRODUCTION

When considering antenna designs suitable for short-hop microwave radio relay application, the more important factors are: low side and rear lobe radiation for suppressing radio interference, reasonably good aperture efficiency, and a structural design which permits inexpensive and simple fabrication.[1] In addition, the antenna should have symmetry permitting the use of orthogonal polarizations for increased channel capacity. Other considerations include a suitable enclosure for the antenna and electronic equipment and provisions for initial radio beam alignment. A supporting structure is required which is high enough to permit radio beam clearance over natural and man-made obstacles, is pleasing in appearance and low in wind loading, and has sufficient structural stiffness to prevent excessive radio beam tilting in heavy wind.

This paper describes such an antenna system which is used in a 1.5-mile experimental repeater installation in New Jersey operating at 11 GHz with transmitting and receiving terminals located on Crawford Hill and a repeater at Bell Laboratories property in Holmdel.[2]

## II. THE ANTENNA AND ENCLOSURE

Electrically, the antenna consists of the basic components shown in Fig. 1: a small aperture feed, a paraboloidal reflector, and a 45-degree plane reflector. This arrangement permits the antenna package to have the shape of an upright cylinder which is desirable for minimum wind loading, ease of azimuthal positioning, and vertically stacking two or more antennas. In addition, the pivoted plane reflector provides a convenient means for initial beam elevation adjustment.

The paraboloidal reflector is a spun aluminum dish with a focal length of 14.5 inches and a diameter of 30 inches. The plane reflector is an elliptically shaped aluminum honey-comb-core sandwich, 0.5 inch thick, supported by two free pivots at either end of its minor
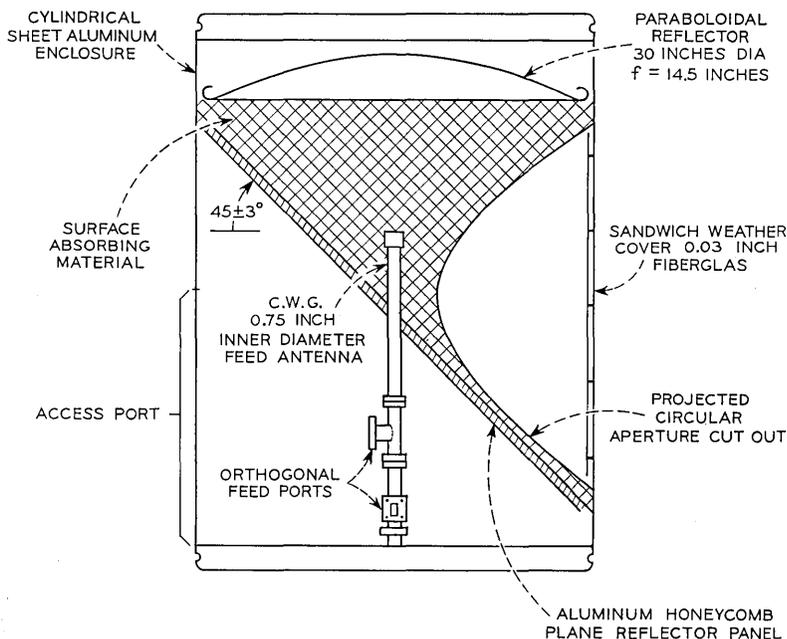


Fig. 1 — Short-hop system antenna package showing components essential to the electrical characteristics.

axis; a third support, a linkage for the elevation adjustment, is located at the upper end of the reflector. The space below the plane reflector is available for the repeater electronics.

The antenna feed is a circular dominant mode waveguide aperture, one inch in diameter, with a two-step transformer for matching to free space. The measured radiation pattern of the feed is almost circularly symmetric over the paraboloidal reflector and provides a field illumination taper of −11 dB at the edge of the paraboloid with respect to its center. Dual polarization is achieved by two orthogonal sidewall couplers.[3] Linear polarizations, vertical and horizontal with respect to the earth, are used for the experimental installation. Orthogonal circular polarizations could also have been used, and would have the advantage over linear polarizations that no cross-polarized coupling would result from the component of mast tilt orthogonal to the radio beam. The disadvantages of orthogonal circular polarizations are that a low loss, broadband 90° phase shifter is required and the reflections from the paraboloid appear as cross talk between the orthogonal input ports rather than as return loss in each port.

The antenna system and the repeater electronics are enclosed in a cylindrical housing 34 inches in diameter and 42 inches high (Fig. 2a). Fabrication is entirely of aluminum with stainless steel fasteners. The cylindrical shell is 0.050-inch sheet aluminum. The inside of the shell is covered with an electrically absorbent material (not shown in the photograph) which eliminates spurious side lobes caused by multiple internal reflection of spill-over radiation from the feed.

Since the large aperture severely weakens the cylindrical shell, a structural framework is provided as shown in Fig. 2b. The top and bottom sections of the structure each consist of two spun aluminum pans, $\frac{1}{16}$ inch thick, joined by radial ribs to form a stiff sandwich. The large perforations in the pans are for interconnecting cables and for ventilation; the external holes are covered with wire mesh. The uprights and bracing struts are square aluminum tubing. The antenna and housing, without electronics, weigh about 85 pounds.

A portion of the cylindrical shell at the rear of the housing can be removed to provide access to the feed assembly and the electronic packages as seen in the photograph. A weather cover (not shown) is provided for the aperture. It consists of a curved sandwich of two $\frac{1}{32}$-inch Fiberglas sheets, appropriately spaced to minimize reflections, and attached to the enclosure by an escutheon plate, self tapping screws, and a waterproof nonsetting cement.
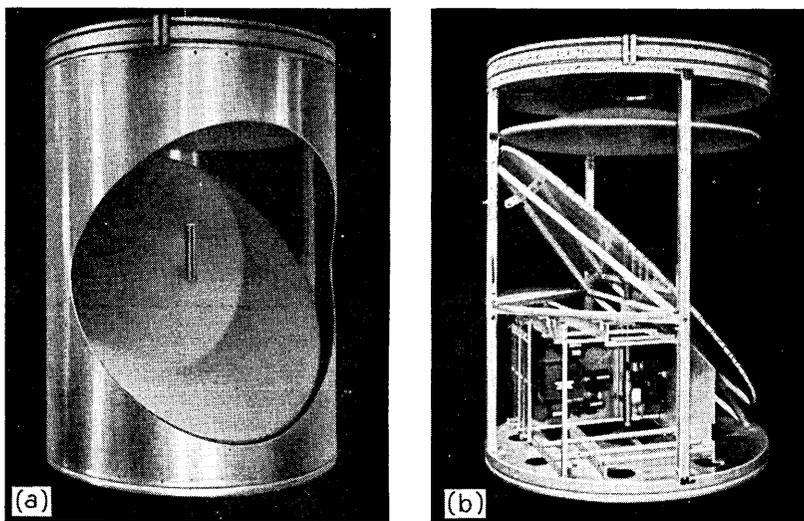
Fig. 2 — Two views of the antenna structure: (a) A front view without weather cover or inner absorbing liner. The embossed girth band and clamp are also shown at the top of the package. (b) A side view without outer skin showing the internal structure.

For installation in the field, the package is secured to a cast aluminum fixture at the top of the supporting mast, or to another package in the stacked configuration, by the girth band and clamp seen on the top of the canister in Fig. 2. This band is made from $\frac{1}{32}$-inch sheet aluminum and has embossed ridges which mate with circumferential grooves, $\frac{3}{16}$-inch deep, in the spun aluminum pans forming the top and bottom sections of the antenna enclosure. Thin Teflon spacers between the canisters reduce friction thus making it possible, after loosening the girth band clamp, to rotate the canister for azimuthal beam positioning, smoothly and accurately, while maintaining a secure hold.

An important attribute of the packaged antenna design is that the housing can be assembled and the antenna components physically aligned in the shop with sufficient precision, using square, level, and plumb bob, that electrical adjustments for bore sight and focal distance are not required. This was demonstrated with the two models assembled for the experimental system; the measured electrical bore sight agreed with the shop-established mechanical bore sight within 0.15 degrees. In making the shop alignments, the focal distance is ad-

justed by moving the paraboloidal reflector by means of slotted holes in mounting brackets attached to the rim; since the radio packages are semirigidly attached to the feed line ports, it is necessary to move the reflector rather than the feed. Screw adjustments are provided for centering the feed aperture on the axis of the paraboloid. A third shop adjustment consists of placing the plane reflector at 45 degrees to the axis of the paraboloid and zeroing the pointer on a calibrated elevation scale on the outside of the canister; a slotted screw, also accessible from the outside, drives a linkage which adjusts the plane reflector for setting the beam in elevation.

III. ANTENNA ELECTRICAL CHARACTERISTICS

3.1 *Radiation Patterns and Gain. Weather Covers*

The radiation characteristics of the antenna system and enclosure are extremely good, particularly without the weather cover. The aluminum cylindrical enclosure, with its electrically absorbent lining, provides excellent shielding which results in very low radiation in the far side lobe region, similar to the performance of the horn-reflector antenna.[4]

Typical measured radiation patterns of the antenna without a weather cover are shown in Figs. 3 through 7. These measurements were made in the principal planes at 11.2 GHz with linear polarizations. Patterns measured at the extremes of the 10.7 to 11.7 GHz band were similar. The line labeled "Isotropic Level" represents the relative power that a hypothetical non-directional antenna would receive. In the azimuthal plane, Figs. 3 and 4, the patterns for both polarizations demonstrate the shielding effect of the enclosure; the response falls abruptly at about ±70 degrees. Figure 5, an expanded version of Fig. 3, shows the main beam and near side lobes in greater detail and illustrates the symmetry of the radiation pattern. This symmetry further attests to the accurate mechanical shop alignment of the antenna components. In the elevation plane (Figs. 6 and 7) a small spill-over lobe is evident at about +90 degrees. This lobe is not of great concern in the pole-line application since in this case it is directed toward the earth.

Initially, the weather cover used was a single sheet of Fiberglas, $\frac{1}{32}$ inch thick. This material has a relatively high dielectric constant (about four). Because the weather cover is curved in the horizontal plane and the angle of incidence approaches grazing at the edges of the aperture, prominent side lobes appeared at about ±90° in azimuth. These lobes, caused by reflections at the cover, were 3 dB above the
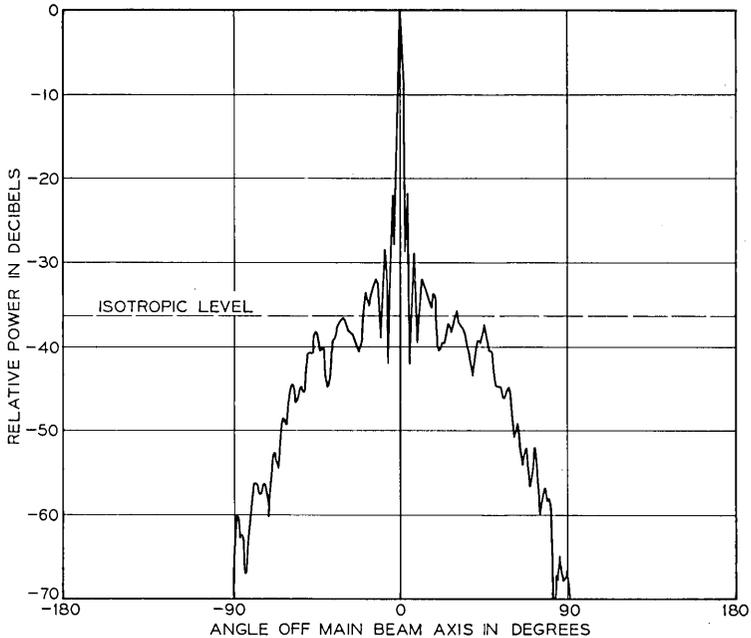
Fig. 3 — Far field radiation pattern in the horizontal plane with horizontal polarization at 11.2 GHz and without weather cover.

isotropic level for vertical polarization and 5 dB below the isotropic level for horizontal polarization.

To reduce these wide angle side lobes, a double-layer cover was designed and installed. Since the cover is curved only in one plane, it was readily constructed from flat sheets using appropriate spacers. The spacing required to cancel the reflections was computed; it varies from 0.15 inches at the center to 0.35 at the edges of the aperture.[5] Radiation patterns taken in the azimuthal plane with the double-layer cover in place are shown in Figs. 8 and 9 for vertical and horizontal polarizations. That the reflections were not cancelled completely may be seen by comparing the far side lobes in Figs. 8 and 9 with those in Figs. 3 and 4. The cancellation could probably be improved by careful adjustment of the spacing. A further advantage of the double-layer cover is that the loss resulting from reflection was reduced to 0.1 dB from 0.5 dB for the single layer cover.

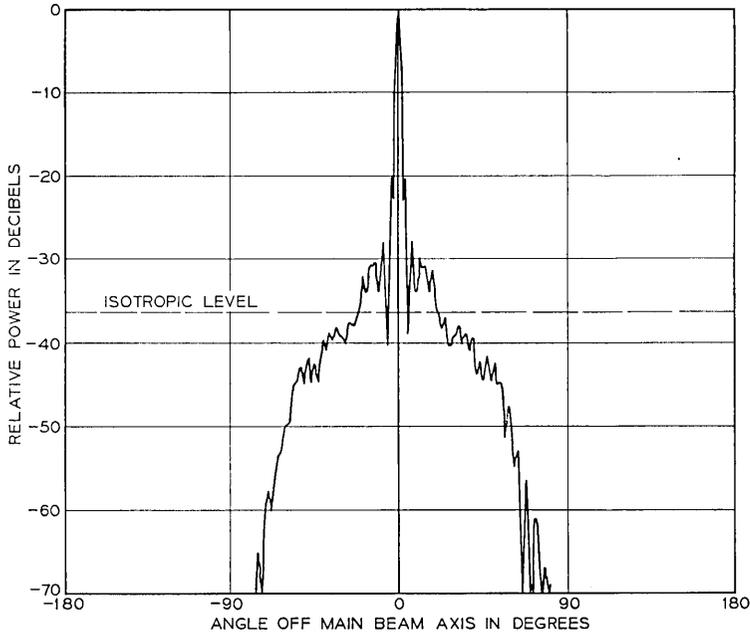The gain of the antenna was measured on a 1500-foot range using a

Fig. 4 — Far field radiation pattern in the horizontal plane with vertical polarization at 11.2 GHz and without weather cover.
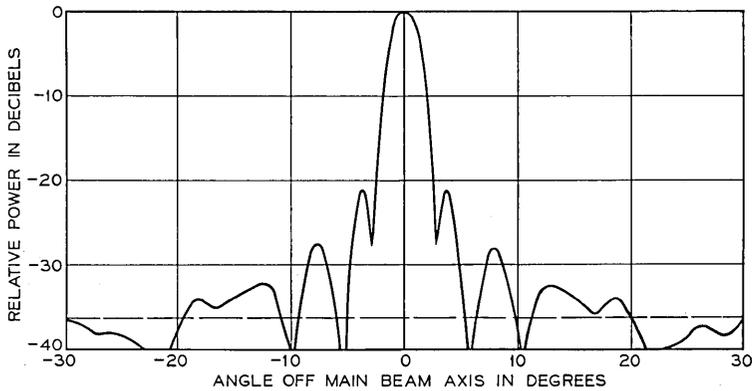


Fig. 5 — An expanded pattern of Fig. 3 centered on the main beam to show more detail and evidence of symmetry.
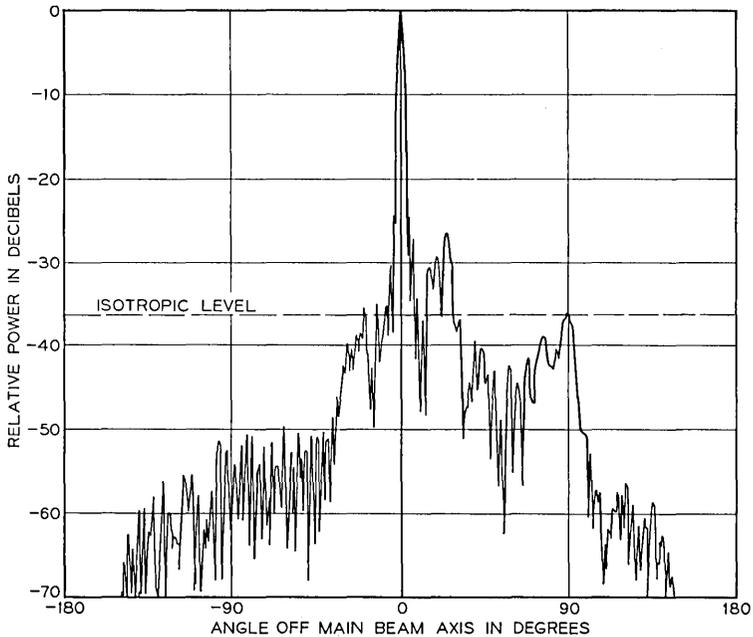
Fig. 6 — Far field radiation pattern in elevation with vertical polarization at 11.2 GHz and without weather cover. The high side lobe at +90 degrees elevation is directed toward the ground in normal applications.

standard gain horn as a reference. Two identical antenna packages, constructed for the experimental short-hop radio system, were measured. Table I summarizes the results of these measurements. The half-power beamwidth of the main lobe is 2.4 degrees at 11.2 GHz, and the gain is 36.6 dB corresponding to an aperture efficiency of about 57 percent.

## 3.2 Cross Polarization. Cross Coupling. Return Loss

Since the short-hop radio systems are expected to use both vertically and horizontally polarized radiations as a means of increasing channel capacity, the cross-polarized characteristics of the antenna are of importance. A parabolic antenna, using a feed such as a dipole or dominant mode circular waveguide, converts some energy from one polarization to its orthogonal counterpart. In theory, this cross-polarized energy appears in planes at 45 degrees to the principal planes, the principal planes being null planes of cross-polarized radiation. However, in practice it is difficult to maintain the principle planes as null

planes and it is also difficult to measure cross-polarization in the 45 degree planes. Some measurements have been made which indicate that, with careful initial alignment of two antennas, it is possible to obtain −35 to −40 dB cross-polarization suppression on the axis. The principal difficulty in maintaining this cross-polarized level will be tilting of the antenna resulting from bending of the mast during high wind conditions. It can be shown that cross-polarized radiation from a parabolic reflector may, in principle, be eliminated by using a feed antenna whose E- and H-plane radiation patterns are identical over the area of the reflector. Such feeds with circularly symmetric radiation patterns can be achieved by the use of dual waveguide modes.[6,7]

Another parameter which may be of importance in some applications is the cross coupling between two antennas stacked on the same mast. This cross coupling was measured to be less than −80 dB for most orientations and polarization combinations. As expected, the worst case
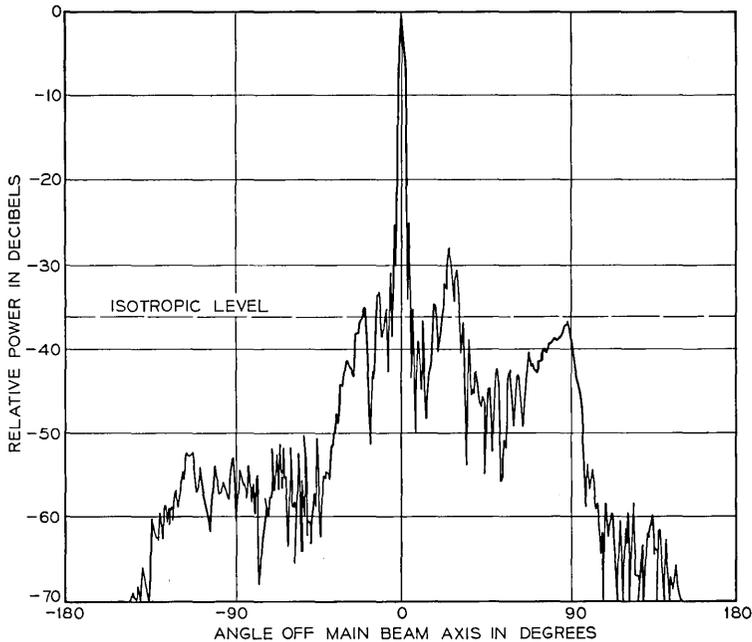


Fig. 7 — Far field radiation pattern in elevation with horizontal polarization at 11.2 GHz and without weather cover. The high side lobe at +90 degrees elevation is directed toward the ground in normal applications.
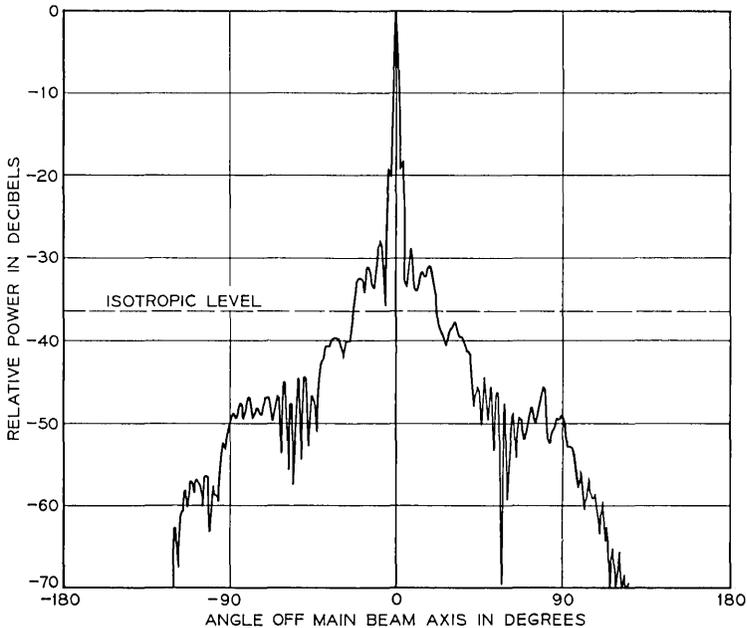
Fig. 8 — Far field radiation pattern in the horizontal plane with horizontal polarization and with the sandwich weather cover installed.

is for vertical polarizations with both beams pointing in the same general direction; for this case, the coupling increases to a maximum of −75 dB. Figures 10a, b, and c show the measured cross coupling as a function of the angle between the beam axes for the three combinations of polarization. The cross coupling between orthogonal feed ports on a single antenna is less than −40 dB.

The feed port return loss measured over the 10.7 to 11.7 GHz band is greater than 20 dB except for a few spot frequencies where it is about 18 dB. No special matching techniques are used to minimize reflections from the paraboloidal surface. The return loss, in almost equal parts, results from reflections at the aperture of the feed and at the surface of the parabolic reflector.

IV. SUPPORT MAST AND FOUNDATION

The choice of a mast to support the antenna packages for the experimental short-hop radio relay installation was based on a number of criteria, some of which were subjective. It was felt that the mast

should be 60 feet high to simulate typical installations where trees might be present. It should be strong enough to withstand the most adverse weather conditions and stiff enough to maintain the radio beam alignment within adequate limits during winds of gale force. The mast should present a pleasing appearance in combination with the antenna packages and, preferably, should be of aluminum for low maintenance. Finally, for the experimental system, it should be readily available.

A mast fulfilling the above requirements was easily obtained since it was being manufactured for use as a lighting support. It is a tapered spun aluminum tube, 60 feet in length, fabricated in two 30-foot sections joined at the time of the installation by a field joint. The outside diameter of the mast is 16 inches at the base, tapering to 10 inches at the top. The lower section has a wall thickness of 0.625 inches; the wall thickness of the upper section is 0.50 inches. The mast alone weighs about 1700 pounds. The end slope of the composite mast with two antenna packages in place, as in a repeater configuration, is computed to be about 0.1 degree for a steady wind force of one pound per square
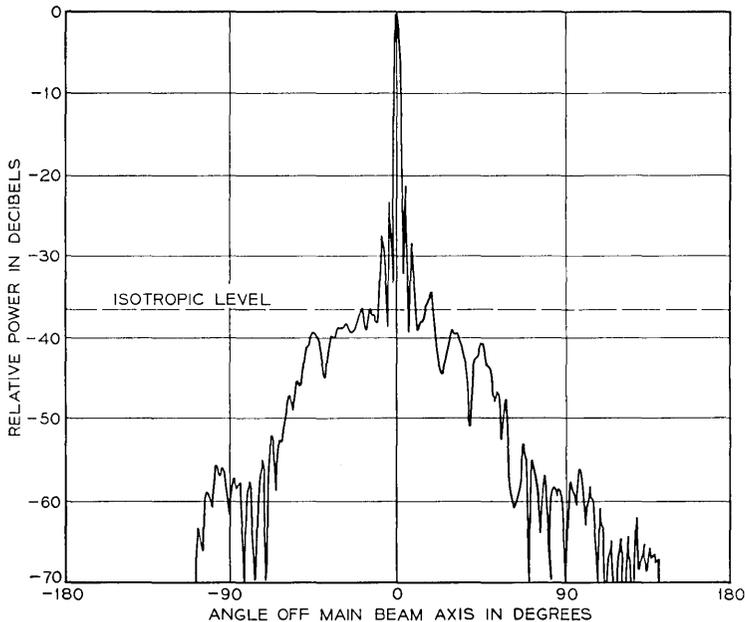


Fig. 9 — Far field radiation pattern in the horizontal plane with vertical polarization and with the sandwich weather cover installed.

TABLE I—SUMMARY OF ANTENNA MEASUREMENTS WITH SANDWICH WEATHER COVER INCLUDED

| Measurement plane | Antenna 1 | | | | Antenna 2 | | | |
|---|---|---|---|---|---|---|---|---|
| | Azimuth | | Elevation | | Azimuth | | Elevation | |
| Polarization | Vertical | Horizontal | Vertical | Horizontal | Vertical | Horizontal | Vertical | Horizontal |
| −3 dB beam width in degrees | 2.37 | 2.46 | 2.42 | 2.39 | 2.43 | 2.35 | 2.39 | 2.43 |
| 1st side-lobe level (dB) | −18.5 | −20.5 | −22.0 | −23.5 | −21.5 | −22.0 | −22.0 | −20.0 |
| Gain (average, vertical, and horizontal polarization) | 36.7 | | | | 36.6 | | | |

Area gain, $4\pi A/\lambda^2$, for a 30-inch circular aperture at 11.2 GHz is 39.02 dB.
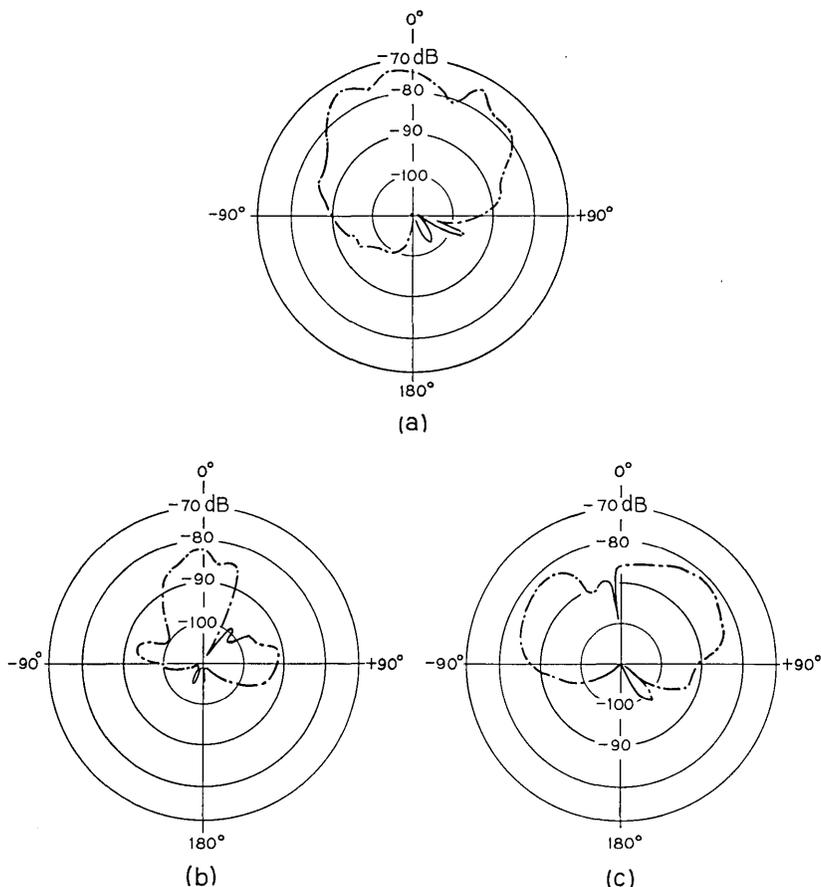
Fig. 10 — Measured mutual coupling between stacked antenna packages for various polarization combinations: (a) both vertically polarized, (b) both horizontally polarized, and (c) one horizontally and the other vertically polarized.

foot of projected area. (Two antenna packages have a projected area of about 20 square feet.) Thus the deflection of the radio beam, which is the same as the end slope of the mast, should be no more than a half beamwidth, 1.2 degrees, for steady winds up to about 65 miles per hour. The mast is expected to survive winds in excess of 200 miles per hour. A stiffer mast of the same weight could be obtained by increasing the diameter and decreasing the wall thickness.

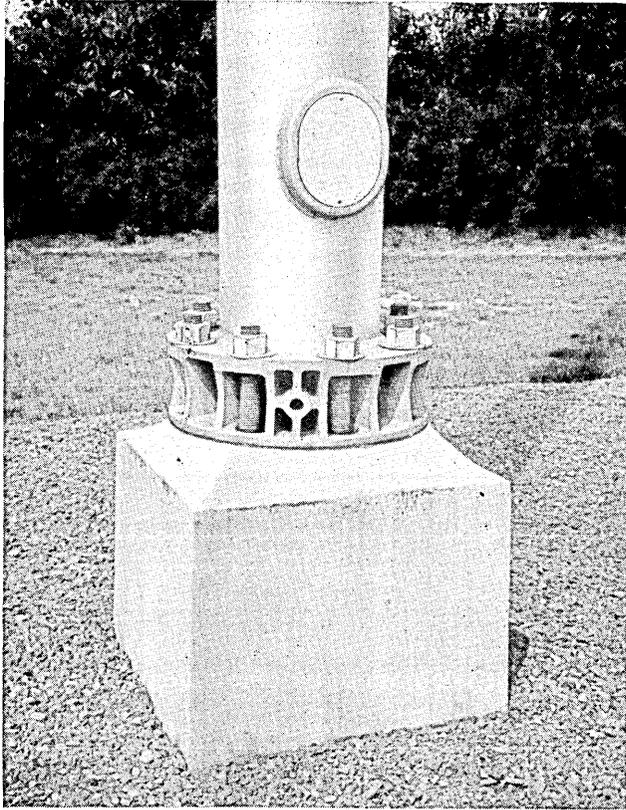The top of the mast is fitted with a cast aluminum spider flange to

Fig. 11 — View of base flange and above ground portion of the foundation. The mounting bolts are 2 inches in diameter.

which the antenna package is fastened by means of the girth band and clamp described in Section II. The base of the mast is welded to a heavy cast aluminum flange for bolting to the foundation as shown in Fig. 11. The overturning moment is about 100,000 pound-feet for steady winds of 100 miles per hour. Figure 12 shows a complete repeater installation with support mast and two stacked antenna packages. Beneath the shallow conical top which overhangs the upper cannister is a screened opening which, with the large holes in the bases of the cannisters, provides a free flow of air for ventilation of the repeater packages. The rectangular box, about 8 feet above the base of the tower, houses a propane gas-fueled thermoelectric generator which provides electric power for the repeater.[2]
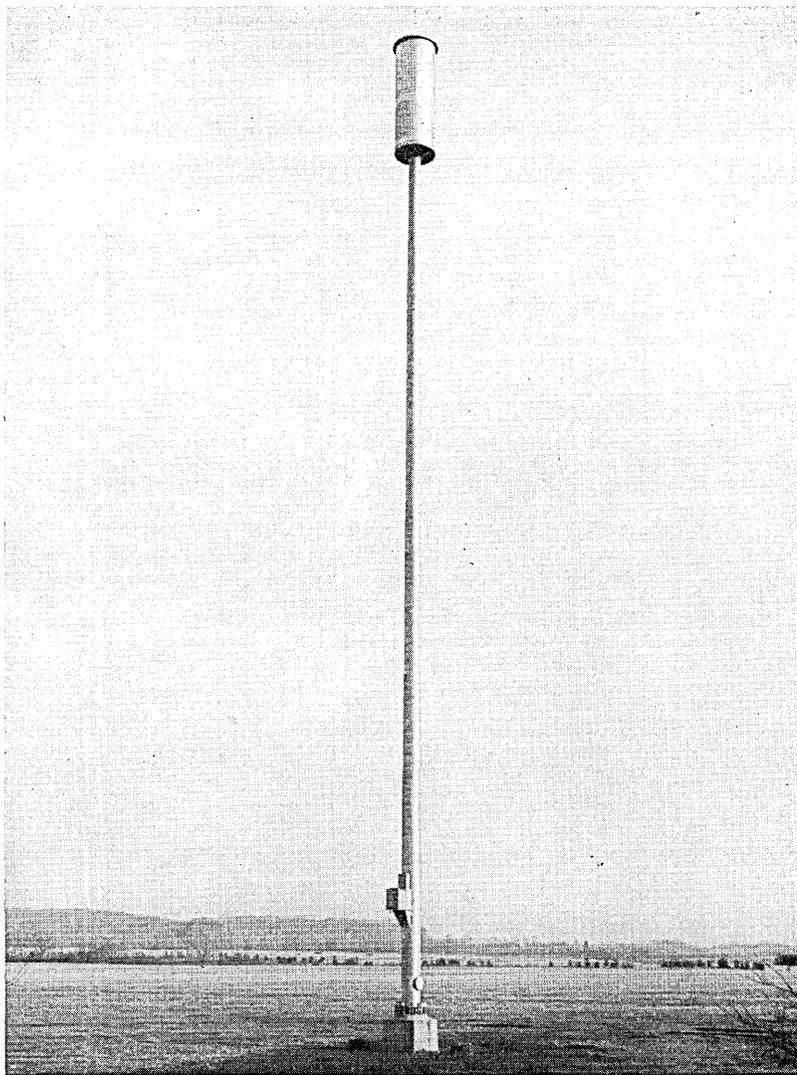
Fig. 12 — A simulated repeater installation showing two antenna packages stacked on the 60-foot aluminum mast. The mast is tapered from 16 inches in diameter at the base to 10 inches at the top. The antenna packages are 34 inches in diameter and each is 42 inches high.

8 – 2 INCH  DIAMETER
ANCHOR BOLTS

GROUND
LEVEL

GAS TANK

10 INCH  DIAMETER
STEEL PIPE
15 FEET LONG

REINFORCING  RODS

(a)

8 – 2 INCH DIAMETER
ANCHOR BOLTS

GROUND
LEVEL

GAS TANK

REINFORCING RODS

ELECTRICAL GROUND  ROD
BONDED TO  STUDS

(b)

Fig. 13 — Schematic cross sections of the two foundations constructed for the experimental short-hop installation. The above ground pedestals are about 3 feet square and include an access duct for the gas line. The cylindrical gas tanks are also shown. (a) A reinforced concrete post-type foundation 18 inches in diameter and extending about 15 feet into the ground. (b) A reinforced concrete pad foundation 9 feet square and 2 feet thick.

For the experimental installation, two different foundations were constructed. For the repeater installation at Holmdel, where the bearing quality of the soil is only fair because of poor drainage, a reinforced concrete pad foundation, 9 feet square and 2 feet thick was constructed. At Crawford Hill, where the bearing strength of the soil is very good, a reinforced concrete post foundation was used. Schematic cross-sections in Figs. 13a and b show both of these installations. For the post-type foundation, a hole 15 feet deep was drilled with an 18-inch auger. A steel pipe, 10 inches in diameter, with reinforcing rods at right angles to the axis, was inserted and the pipe as well as the hole was filled with concrete. This type of foundation is preferred where possible because it is simple to install and requires a minimum of right-of-way area.

The reinforced concrete pedestals are 3 feet square and are identical for both foundations. Eight threaded rods, 2 inches in diameter, secure the base flange of the mast to the foundation. Adjustable nuts below and above the flange are used to plumb the mast while elongated holes in the flange permit a small adjustment of the mast (and antenna packages) in azimuth. Conduits, not shown in the drawings, are provided for gas lines from the gas storage tanks through the centers of the pedestals to the mast.

For the experimental installation, it was possible to determine the angles of elevation and the bearings of the sites by surveying techniques so that the masts, with antenna packages attached, were erected and aligned mechanically with good accuracy using theodolites and alignment marks on the antennas and mast flanges. This was shown later when it was found that the antenna beams were within 0.05 degree in elevation and 0.3 degree in azimuth of true bearing.

V. ACKNOWLEDGMENTS

REFERENCES

1. Tillotson, L. C., "Use of Frequencies Above 10 GHz for Common Carrier Applications," B.S.T.J., this issue, pp. 1563–1576.
2. Ruthroff, C. L., Osborne, T. L., and Bodtmann, W. F., "Short Hop Radio System Experiment," B.S.T.J., this issue, pp. 1577–1604.

3. Ohm, E. A., "A Broad-Band Microwave Circulator," IRE Trans. Microwave theory and techniques *MTT-4,* No. 4, (October 1956), pp. 210–217.
4. Crawford, A. B., Hogg, D. C., and Hunt, L. E., "A Horn-Reflector Antenna for Space Communication," B.S.T.J., No. 4, (July 1961), pp. 1095–1116.
5. Jasik, H., *Antenna Engineering Handbook,* Section 32.3, "Sandwich Radomes," New York: McGraw-Hill, 1961.
6. Minnett, H. C., and Thomas, B. MacA., "A Method of Synthesizing Radiation Patterns with Axial Symmetry," IEEE Trans. Antennas and Propagation, *AP-14,* No. 5 (September 1966), pp. 654–656.
7. Turrin, R. H., "Dual Mode Small Aperture Antennas," IEEE Trans. Antennas and Propagation, *AP-15,* No. 2 (March 1967), p. 307.

# Design of Efficient Broadband Varactor Upconverters

By T. L. OSBORNE

*This paper describes a design procedure for optimizing the performance of a varactor upconverter microwave power amplifier with respect to maximum pump efficiency; the procedure gives explicit diode and circuit parameters and operating levels. An optimum diode is selected by inclusion of an empirical relation between diode breakdown voltage and cutoff frequency. A fully driven abrupt junction diode is assumed.*

*The bandwidth of an upconverter is limited by the intermediate frequency input circuit which typically has a 3 dB bandwidth of about 10 percent. We describe a method of obtaining broadband operation where the interface between the intermediate frequency source and the varactor diode is mismatched in an optimum way. Analysis shows that the frequency variation in mismatch loss just compensates for the frequency variation in upconverter gain predicted by the Rowe–Manley relations.*

*The design procedure is illustrated by a 300 MHz to 10.960 GHz varactor upconverter built for use in the transmitter of the short hop radio system experiment. A bandwidth of 120 MHz between 1.4 dB points represents a bandwidth of more than 40 percent at the intermediate frequency. An output power of +16 dBm was obtained using a pump power of +20 dBm giving a pump efficiency of 40 percent. The normal input to the driver amplifier is +3 dBm giving an overall gain of +13 dB. The upconverter operates over a temperature range of −40°F to +140°F with only a small change in bandshape and output power.*

## I. INTRODUCTION

In many microwave radio relay systems the transmitted microwave signal is obtained by upconverting the IF signal in a mixer or modulator and amplifying it in a microwave power amplifier, such as a traveling wave tube, to a power level set by system requirements. By use of a varactor upconverter as a microwave power amplifier the func-

tions of upconversion and amplification can be performed by one device. In fact, for all-solid-state radio systems the varactor upconverter is the only broadband microwave power amplifier available. Therefore, in radio systems such as those described in companion papers, broadband, efficient, rugged varactor upconverters are required as power amplifiers.[1,2]

Although several analyses of varactor upconverters have been published, design has required a trial and error procedure with its inherently uncertain results. In this paper we describe a design procedure for optimizing the performance of an upper sideband varactor upconverter with respect to maximum pump efficiency. This procedure gives diode and circuit parameters explicitly in terms of operating power levels. We also describe a method of obtaining broadband operation. Although the procedure is derived using specific frequencies and the maximum efficiency optimization, the basic method can be used for any frequency and optimization.

The design procedure and broadbanding method are illustrated by the design of an upconverter for the experimental radio system described by Ruthroff and others.[3] The system requires a bandwidth of 120 MHz between 1 dB points and has an intermediate frequency of 300 MHz. A maximum of 100 mW (+20 dBm) of pump power was available; the available IF power was +3 dBm. The objective was maximum output power consistent with the foregoing constraints and the limitations imposed by the environment. The performance of the two upconverters required in the experiment was essentially the same, so the discussion is primarily of the 300 MHz to 10.960 GHz upconverter.

II. SUMMARY

The limited pump power available at 11 GHz requires maximum pump efficiency upconverter operation. From the large signal analysis of Penfield and Rafuse and an empirical relation between diode cutoff frequency and breakdown voltage, the output power versus pump power for fully driven maximum efficiency operation is found.[4-6] Then, given minimum output power and maximum pump power available, an operating range and all diode circuit parameters and diode characteristics are specified. The details of this procedure are given in Section III.

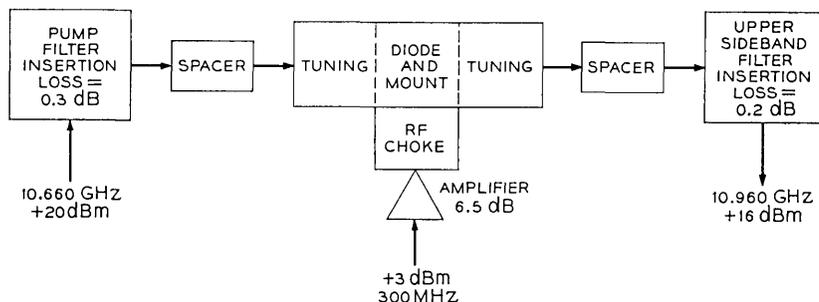The upconverter is a straight-through type with a single diode (see

Fig. 1 — Block diagram of upconverter.

the block diagram in Fig. 1). The RF circuit uses a standard WR90 X-band waveguide and consists of properly located signal and pump filters for providing the required terminations, tuning sections for matching input and output signals, and the diode mount and impedance transforming capacitor. The diode-to-waveguide impedance transformation circuit consists of the diode inductance and an adjustable shunt capacitance formed by an adjustable sleeve concentric with the diode. This method of matching eliminates the need for reduced height waveguide and associated transitions and allows the filters to be close to the diode with a corresponding reduction of bandshape ripples. The details of the RF circuit are discussed in Section IV.

The +3 dBm IF signal from the main IF amplifier is amplified by a driver amplifier attached to the upconverter and applied to the diode through a low capacitance coaxial RF choke in the broad wall of the waveguide. The bandwidth of an upconverter is limited by the IF input circuit which typically has a 3 dB bandwidth of about 10 per cent. Analysis of the IF circuit including the effect of the frequency variable input resistance shows that if the varactor is driven from a constant resistance source the frequency variation in mismatch loss will compensate for the frequency variation in upconverter gain. Consequently the driver amplifier-upconverter interface is purposely mismatched to get the required 40 per cent bandwidth at the cost of reduced IF-to-RF gain. The IF circuit analysis and the amplifier performance is discussed in Section V.

Complete data on the upconverter performance is given in Section VI and summarized in Table I. The operating temperature range is $-40°F$ to $+140°F$ with small changes in the frequency response as shown in Fig. 2. Figure 3 is a photograph of the upconverter.

TABLE I—10.960 GHz BROADBAND UPCONVERTER PERFORMANCE

| | |
|---|---|
| Pump power into filter | +20 dBm |
| IF power into amplifier | +3 dBm |
| Output power at band center | +16.2 dBm at +140°F and 76°F |
| | +16.9 dBm at −40°F |
| Frequency response at 76°F | 1.4 dB down at 10.900 GHz |
| | 1.0 dB down at 11.020 GHz |
| | ripples <0.2 dB peak-to-peak |

Diode: Bell Telephone Laboratories L2280 with
$V_B = 32$ V, $C_o = 1.091$ pF, $f_{co} = 137$ GHz.

### III. MAXIMUM EFFICIENCY OPERATION AND DIODE SELECTION

If the output power is to be maximum with a given pump power, the upconverter must operate at maximum pump efficiency. In this section the operating point for maximum efficiency and the diode parameters are derived using the large signal analysis of Penfield and Rafuse, assuming a fully driven abrupt-junction silicon diode with $S_{\min} = 0$ and output load tuned.[4]

From Penfield and Rafuse, the pump efficiency, $\epsilon$, and upper-side-band power, $P_u$ in milliwatts, are

$$\epsilon = \frac{f_u}{f_p} \frac{m_s - \dfrac{f_u}{f_c}\dfrac{m_u}{m_p}}{m_s + \dfrac{f_p}{f_c}\dfrac{m_p}{m_u}}, \tag{1}$$

and

$$P_u = 16\pi(V_B + \Phi)^2 f_c C_{\min}\left(\frac{f_u}{f_c} m_s m_p m_u - \frac{f_u^2}{f_c^2} m_u^2\right), \tag{2}$$

where

$f_s, f_p, f_u$ are the signal (IF), pump, and upper-sideband frequencies in GHz, and are known constants in this analysis,

$m_s, m_p, m_u$ are the magnitudes of the normalized elastance coefficients,

$f_c$ is the diode cutoff frequency in GHz defined as

$$f_c = \frac{10^3}{2\pi R_s C_{\min}} \quad \text{GHz}, \tag{3}$$

$V_B$ is the diode reverse breakdown voltage,

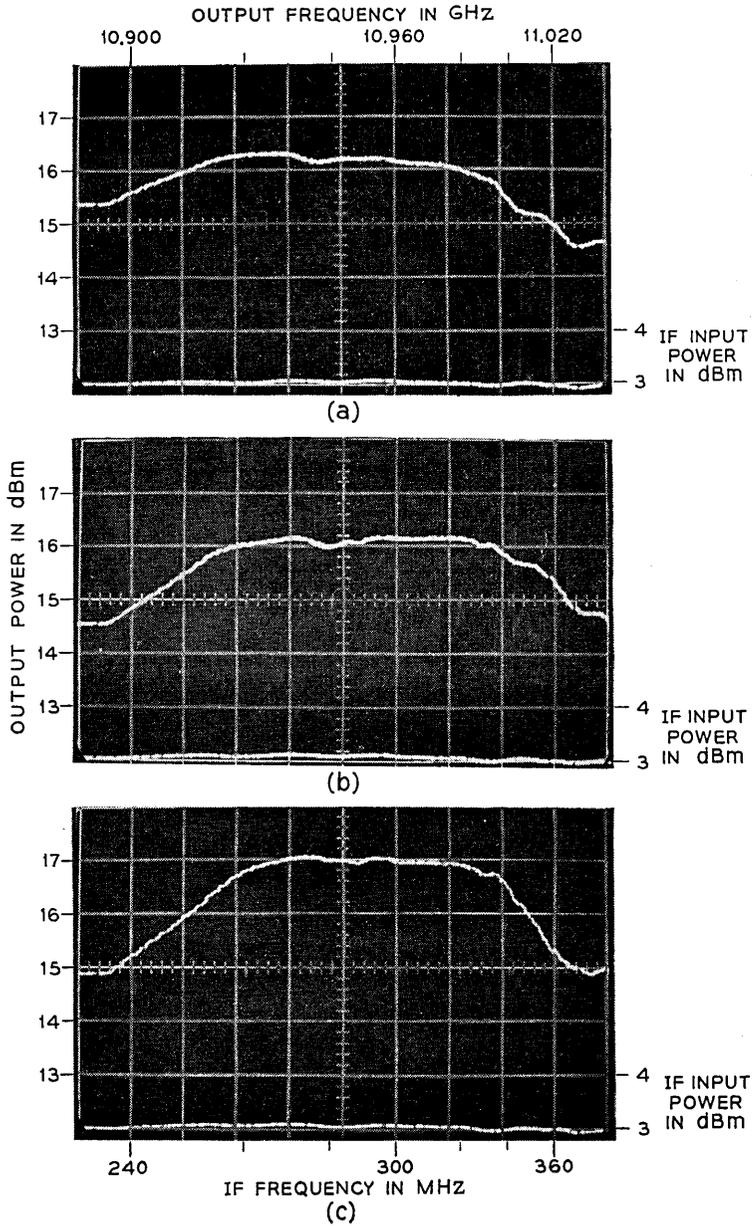$\Phi$ is the contact potential (0.7 volts for silicon),

and

Fig. 2 — Upconverter frequency response at (a) +140°F, (b) +76°F, and (c) −40°F, with +3 dBm IF input power and +20 dBm pump power.
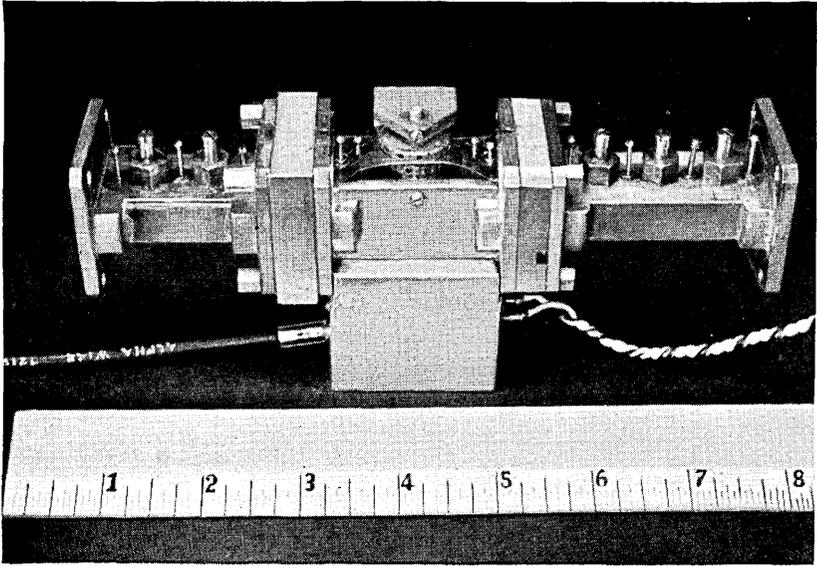
Fig. 3 — 10.960 GHz upconverter.

$C_{\min}$ is the minimum or breakdown capacitance in picofarads. For finite $f_c$ and fully driven operation, that is,

$$m_s + m_p + m_u = 0.25,\tag{4}$$

$\epsilon$ has an absolute maximum when $m_s = 0.25$ and $m_u = m_p = 0$, a useless condition because the output power is zero. For $m_s < 0.25$, the ratio $m_u/m_p$ which gives maximum efficiency can be found by differentiating (1). The result is

$$R = \frac{m_u}{m_p}\bigg|_{\substack{\max \\ \text{eff}}} = \frac{1}{m_s}\frac{f_p}{f_c}\left[\left(\frac{f_c^2}{f_p f_u}\,m_s^2 + 1\right)^{\frac{1}{2}} - 1\right].\tag{5}$$

when $f_c$ and $m_s$ are given, $m_p$ and $m_u$ can be calculated from (4) and (5):

$$m_p = \frac{1}{R+1}\,(0.25 - m_s)\tag{6}$$

$$m_u = \frac{R}{R+1}\,(0.25 - m_s).\tag{7}$$

Substitution of (5), (6), and (7) into (2) eliminates $m_u$ and $m_p$ leaving only $m_s$, $f_c$, $V_B$, and $C_{\min}$. For any set of these variables the pump efficiency is maximum with respect to $m_p$ and $m_u$.

In practice the diode parameters $f_c$, $V_B$, and $C_{\min}$ are not independent because of the experimentally determined dependence of breakdown voltage on material resistivity. Irvin has calculated $f_c$ and $f_{co}$, the zero bias cutoff frequency, as a function of $V_B$ for abrupt-junction silicon diodes using easily realizable breakdown versus resistivity values.[5] His curve for $f_{co}$ agrees well with a curve by Lee giving the state of the art for epitaxial diffused silicon diodes (see Fig. 3).[6]

The equation,

$$f_c = \frac{14 \times 10^3}{C_o^{\frac{1}{2}}(V_B + \Phi)} \qquad \text{GHz,} \tag{8}$$

where $C_o$ is in picofarads, closely approximates Irvin's calculated values of $f_c$ versus $V_B$ as shown in Fig. 4. The $C_o$ dependence was obtained by plotting $1/f_c$ versus $C_o$ for constant $V_B$ for several commercially available diodes. Also shown in Fig. 4 are the $f_{co}$ and $V_B$ for several 1.0 pF Bell Telephone Laboratories L2280 silicon diodes obtained for this and related experiments.

Substituting (8) into (2), using the relation

$$C_{\min} = C_o \frac{\Phi^{\frac{1}{2}}}{(V_B + \Phi)^{\frac{1}{2}}}, \tag{9}$$

and setting $\Phi = 0.7$ gives an equation for $P_u$ normalized to $C_o^{\frac{1}{2}}$:

$$\frac{P_u}{C_o^{\frac{1}{2}}} = \frac{6.97 \times 10^7}{f_c^{\frac{1}{2}}} \left( \frac{f_u}{f_c} m_s m_p m_u - \frac{f_u^2}{f_c^2} m_u^2 \right). \tag{10}$$

The following equations can be derived from the equations in Penfield and Rafuse in a similar way. In maximum pump efficiency operation, as defined by equation (5), the load and pump input resistances, $R_u$ and $R_{\mathrm{in},\,p}$, are equal.

$$\frac{P_{\mathrm{in},p}}{C_o^{\frac{1}{2}}} = \frac{6.97 \times 10^7}{f_c^{\frac{1}{2}}} \left( \frac{f_p}{f_c} m_s m_p m_u + \frac{f_p^2}{f_c^2} m_p^2 \right). \tag{11}$$

$$\frac{P_{\mathrm{in},s}}{C_o^{\frac{1}{2}}} = \frac{6.97 \times 10^7}{f_c^{\frac{1}{2}}} \left( \frac{f_s}{f_c} m_s m_p m_u + \frac{f_s^2}{f_c^2} m_s^2 \right). \tag{12}$$

$$C_o^{7/6} R_u = C_o^{7/6} R_{\mathrm{in},p} = \frac{22.5 \times 10^3}{f_c^{\frac{3}{2}}} \left( \frac{f_c}{f_p} \frac{m_s m_u}{m_p} + 1 \right). \tag{13}$$

$$C_o^{7/6} R_{\mathrm{in},s} = \frac{22.5 \times 10^3}{f_c^{\frac{3}{2}}} \left( \frac{f_c}{f_s} \frac{m_p m_u}{m_s} + 1 \right). \tag{14}$$

All these equations are valid for maximum pump efficiency operation of fully driven abrupt junction silicon diodes and any frequencies.
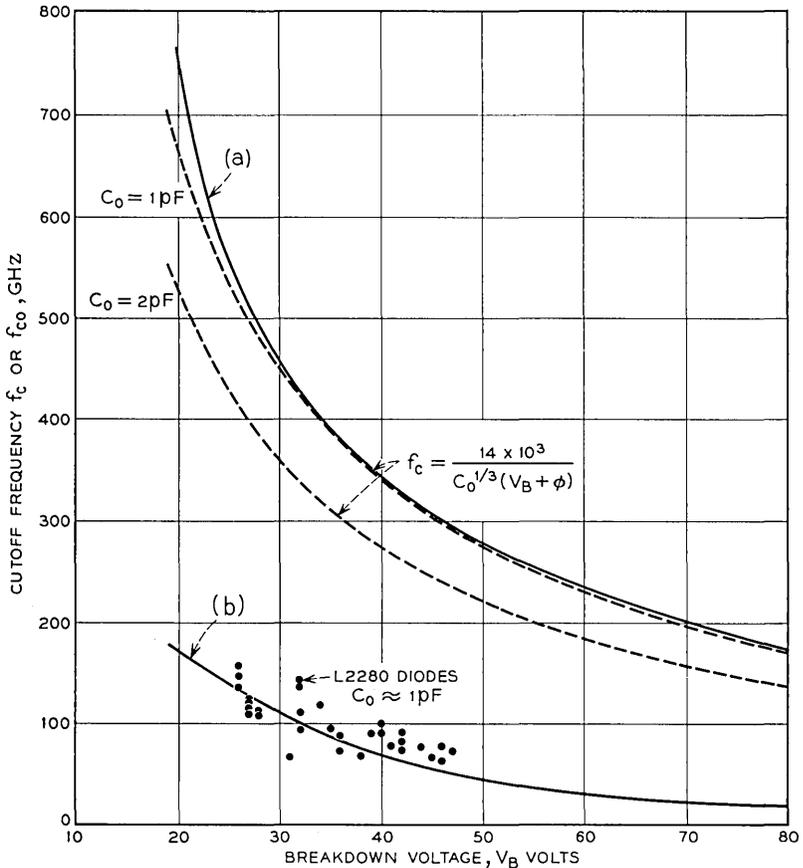
Fig. 4 — Cutoff Frequencies, $f_c$ and $f_{co}$, versus breakdown voltage for silicon abrupt junction diodes: (a) calculated $f_c$ for silicon abrupt junction (Ref. 5), $C_o = 1$ pF, (b) calculated $f_{co}$ for silicon abrupt junction (Ref. 5), $C_o = 1$ pF and state of art for epitaxial diffused silicon (Ref. 6), $C_o = 0.8$ pF.

In Figs. 5, 6, and 7 the equations are plotted as a function of $m_s$ with $f_c$ as a parameter with $f_s = 300$ MHz, $f_p = 10.660$ GHz, and $f_u = 10.960$ GHz. The ratio $R$ is plotted in Fig. 8 for use when calculating the values of $m_u$ and $m_p$.

In these equations $m_s$ is used as an independent variable even though it is not really independent but rather an intermediate variable which relates the independent variables, power and resistance. In Fig. 5 the powers are plotted as functions of $m_s$ so that the relations between them are shown. Any point in Fig. 5 defines an operat-

ing point in the sense that if the indicated diode (defined by $f_c$), $P_{in, s}$, and $P_{in, p}$ are used and the diode terminated in the resistances shown in Fig. 7, the output power theoretically will be that value shown on the ordinate of Fig. 5 and $m_s$ will have the value shown on the abscissa of Fig. 5. The values of $m_s$ and $f_c$ are used to locate the operating point on the graphs in Figs. 6, 7, and 8.

Even though the ratio of $m_u$ to $m_p$ was chosen to give maximum pump efficiency it is still possible to choose the remaining variables to get further maximization of pump efficiency or other optimizations. The required operating point can be found from Fig. 5 as follows.

The solid lines are lines of constant $f_c$ and show the variation of normalized output power with $m_s$ for a given diode. The peaks show that for a given diode maximum output power would be obtained by operating near $m_s = 0.10$.
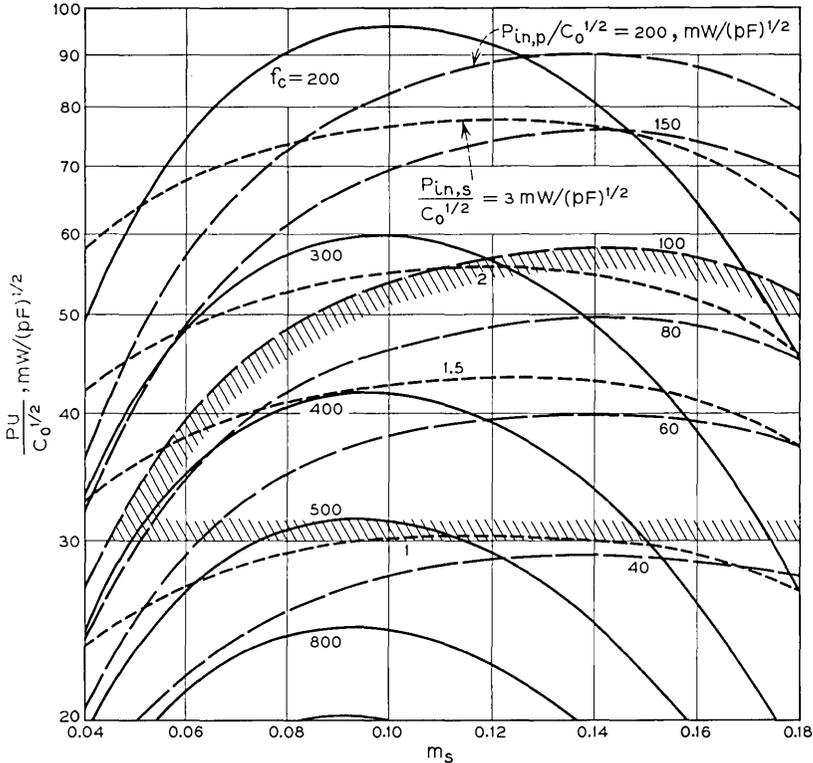


Fig. 5 — Upper-sideband output power versus $m_s$ with $f_c$ as a parameter and with contours of constant pump and signal power superimposed.
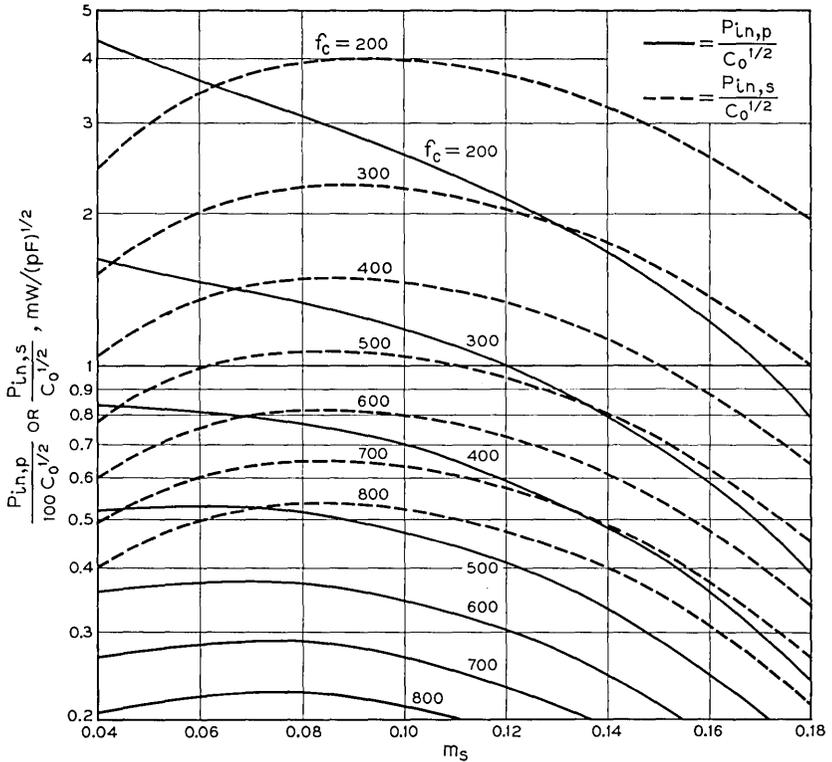
Fig. 6 — Normalized pump and signal power versus $m_s$ with $f_c$ as a parameter.

The broken lines are contours of constant normalized pump power. The peak in these curves shows that a given output power would be obtained with the least amount of pump power by operating near $m_s$ = 0.14. Thus, for a given output power, kept constant by letting $f_c$ vary, this is a point of maximum pump efficiency. However if $f_c$ is constant the pump efficiency increases monotonically if the operating point is chosen at a larger $m_s$ and the output power goes to zero.

The dashed lines are contours of constant normalized signal power. The peak in these curves shows that a given output power would be obtained with the least amount of signal power by operating near $m_s$ = 0.12. Thus for a given output power this is a point of maximum gain. If $f_c$ is kept constant, maximum gain occurs when operating near $m_s$ = 0.15.

Figure 5 can also be used to define the set of operating points which meet the power requirements of a given problem. Using the radio system upconverter as an example, the shaded area defines a set of operating points each satisfying the requirements, $P_u \geqq 30$ mW and $P_{\text{in},p} \leqq 100$ mW, when the diode capacitance is 1.0 pF. From Figs. 4 through 7 a theoretical operating point and diode were chosen. Since the output power theoretically obtainable is safely above the minimum needed, a lower value of $m_s$ was chosen to avoid large values of $R_u$ and $R_{\text{in},p}$ . The theoretical diode characteristics and upconverter diode performance are summarized in Table II. The average junction capacitance,



Fig. 7 — Normalized resistances versus $m_s$ with $f_c$ as a parameter.

Fig. 8 — Plot of $R$, the ratio of $m_u$ to $m_p$, versus $m_s$ with $f_c$ as a parameter.

$C_j$ , is calculated from

$$C_j = 2C_{min} = 2C_o\left(\frac{\Phi}{V_B + \Phi}\right)^{\frac{1}{3}} \text{pF} \tag{15}$$

assuming current pumping and cubic capacitance variation.*

IV. WAVEGUIDE CIRCUIT

Figures 9a and b show a cross section of the L2280 U-package diode and its equivalent circuit.[6] Above about 7 GHz the diode operates

---

\* The actual capacitance variation of the L2280 diode is approximately cubic. Even though a square root variation was assumed for ease of analysis, use of a cubic variation in calculating $C_j$ should give better agreement between theory and practice.

TABLE II—THEORETICAL UPCONVERTER DIODE
CHARACTERISTICS AND PERFORMANCE

| $f_c = 350$ GHz | $R_u = R_{in,p} = 10.6$ ohms | $P_{in,s} = 1.8$ mw |
|---|---|---|
| $V_B = 39$ volts | $R_{in,s} = 280$ ohms | $\epsilon = 52.4\%$ |
| $f_{co} = 75$ GHz | $P_u = 49.2$ mw | $G = 27.2 = 14.3$ dB |
| $C_o = 1.0$ pF | $P_{in,p} = 93.9$ mw | $C_j = 0.521$ pF |
| $m_s = 0.09$ | | |

above its self-resonant frequency and appears as a frequency depend-
ent inductance in series with $R_{in,\,p}$ as shown in Fig. 9c, where

$$L(\omega) = L_s\left(1 - \frac{1}{\omega^2 L_s C_j}\right). \tag{16}$$

At $\omega_o$ the equivalent parallel resistance and inductance, Fig. 9d, are

$$R_p(\omega_o) = R_{in,p}\left[1 + \frac{\omega_o^2 L_s^2}{R_{in,p}^2}\left(1 - \frac{1}{\omega_o^2 L_s C_j}\right)^2\right], \tag{17}$$
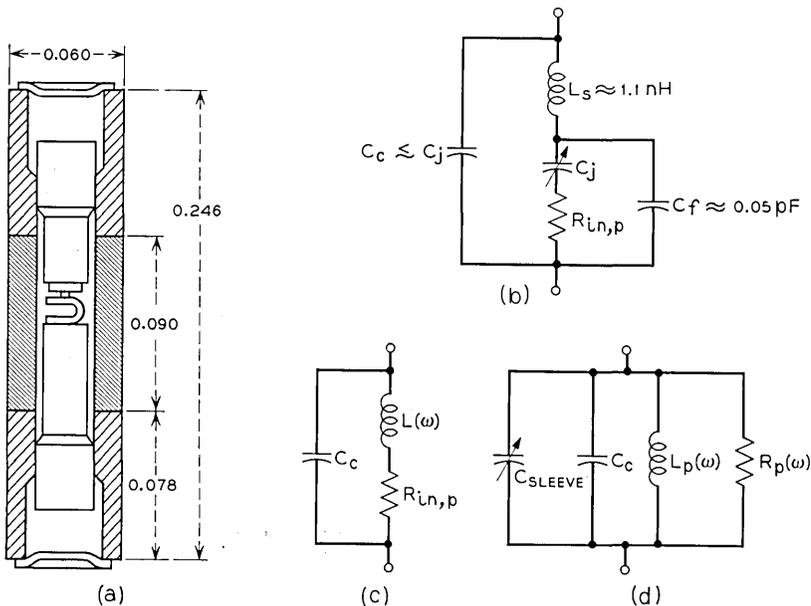


Fig. 9 — (a) L2280 U-package epitaxial diode (Ref. 6). (b) Equivalent circuit
of diode. (c) Equivalent circuit of diode at 10.7 GHz. (d) Parallel equivalent of
(c) with adjustable capacitance of the sleeve in parallel.

and

$$L_p(\omega_o) = L_s\left(1 - \frac{1}{\omega_o^2 L_s C_j}\right)\left[1 + \frac{1}{\frac{\omega_o^2 L_s^2}{R_{in,p}^2}\left(1 - \frac{1}{\omega_o^2 L_s C_j}\right)^2}\right]. \tag{18}$$

The total capacitance in parallel with the diode, including $C_c$, required to tune $L_p(\omega_o)$ at $\omega_o$ is

$$C_T = \frac{L_s\left(1 - \frac{1}{\omega_o^2 L_s C_j}\right)}{R_{in,p}^2 + \omega_o^2 L_s^2\left(1 - \frac{1}{\omega_o^2 L_s C_j}\right)^2}. \tag{19}$$

Substituting the values of $R_{in,\,p}$ and $C_j$ from Table II and $L_s = 1.1$ nH and $f_o = 10.900$ GHz into equations (17) and (19), gives $R_p = 222$ ohms and $C_T = 0.294$ pF. This value of $C_T$ can be easily obtained by adding a small capacitance, $C_{sleeve}$, in parallel with the diode, thus using the diode self-inductance and parallel capacitance as a tuned impedance transformer. The $R_p$ of 222 ohms across a waveguide with $Z_o = 420$ ohms can be matched easily with three $\lambda/8$ tuning screws. This method of impedance matching eliminates the need for stepped impedance transformers and results in a significant simplification of the RF circuit.

Physically, the diode is mounted across a gap formed by two cylindrical posts protruding into standard height X-band waveguide from opposite broad walls; see Fig. 10. The lower post is a diode chuck and also the center conductor of the coaxial RF choke through which the IF signal is applied. The top post is a diode chuck surrounded by an adjustable cylindrical sleeve. The diode tuning capacitor is formed by the end of the sleeve and the lower diode chuck. The capacitance can be varied from approximately 0.1 pF to some large value when the sleeve touches the lower chuck.

The $Q$ and the bandwidth of the RF impedance matching circuit can be calculated from (17) and (19):

$$Q(f_o) = 2\pi f_o C_T(f_o) R_p(f_o) = 4.47.$$

Ignoring the frequency dependence of $R_{in,\,p}$ and $R_p$, the 3 dB bandwidth is approximately

$$BW = \frac{f_o}{Q(f_o)} = 2.4 \text{ GHz},$$

which is much larger than the required 120 MHz. Therefore the diode will be matched to the waveguide over the 120 MHz band and the RF bandwidth will be determined by the output filter.

The pump and upper-sideband output filters are direct coupled filters with multiple posts as the shunt inductive reactances. A two-section maximally flat filter is used for the pump filter and a three-section 0.1 dB Tchebyscheff filter is used for the upper-sideband output filter. The insertion loss and return loss of the output filter are shown in Fig. 11; the measured performance of both filters is summarized in Table III.

## V. IF CIRCUIT

The main consideration in the IF circuit is the bandwidth requirement. Since the required 120 MHz bandwidth is a very small percentage bandwidth at RF, the RF circuit presents no problem as far as bandwidth is concerned. At the 300 MHz IF, however, the 120 MHz represents a 40 percent bandwidth. An estimate of the severity



Fig. 10 — Cross section of upconverter.

Fig. 11 — Return loss and insertion loss of the 3 section 0.1 dB Tchebyscheff output filter.

TABLE III—MEASURED PERFORMANCE OF PUMP
AND OUTPUT FILTERS

|  | Pump filter | Upper sideband output filter |
|---|---|---|
| Type | Maximally flat | 0·1 dB Tchebyscheff |
| Number of sections | 2 | 3 |
| Center frequency | 10.658 GHz | 10.960 GHz |
| Bandwidth | 134 MHz at 3 dB | 152 MHz at 0.1 dB |
| Skirt attenuation | 20 dB at 10.900 GHz | 26.5 dB at 10.660 |
| Midband loss | 0.3 dB | 0.2 dB |
| Return loss | 27 dB at 10.658 GHz | >18 dB in pass band |

of the problem can be found by calculating the $Q$ of the IF input circuit shown in Fig. 12. The input impedance at the RF choke terminals is

$$Z = \frac{R_{in,s}}{\left(1 + \dfrac{C_s}{C_i}\right)^2 \left[1 + \omega^2 R_{in,s}\left(\dfrac{C_s C_i}{C_s + C_i}\right)^2\right]}$$
$$+ \frac{\left[1 + \omega^2 R_{in,s} C_i\left(\dfrac{C_s C_i}{C_s + C_i}\right)\right]}{j\omega(C_s + C_i)\left[1 + \omega^2 R_{in,s}\left(\dfrac{C_s C_i}{C_s + C_i}\right)^2\right]} \quad (20)$$

where the frequency variation of $R_{in,s}$ has been neglected. For the values given in Table II, the impedance is approximately

$$Z = \frac{R_{in,s}}{\left(1 + \dfrac{C_s}{C_i}\right)^2} + \frac{1}{j\omega(C_s + C_i)}. \quad (21)$$



Fig. 12 — IF input circuit.

The $Q$ and 3 dB bandwidth at $f_o = 300$ MHz are

$$Q = \frac{\left(1 + \frac{C_s}{C_j}\right)^2}{\omega_o R_{\text{in},s}(C_s + C_j)} = 15.2$$

$$BW = \frac{2\pi f_o^2 R_{\text{in},s}(C_s + C_j)}{\left(1 + \frac{C_s}{C_j}\right)^2} = 19.8 \text{ MHz.}$$
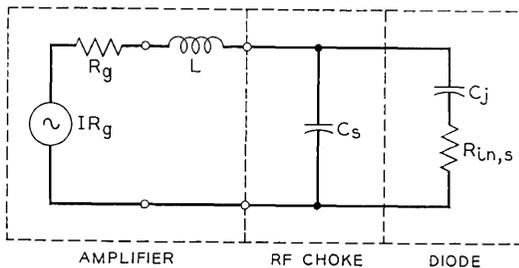
Simple resistance loading to lower the $Q$ would require a total resistance about 6 times the unloaded input resistance, would result in 15.7 dB of loss, and the frequency response would still be down 3 dB at the band edge. Some improvement on this loss could be achieved by using a more elaborate filter, but the filter itself could contribute considerable loss.

The inherent narrow bandwidth is caused by the shunt capacitance, $C_s$, which limits the high frequency response and the junction capacitance, $C_j$, which in turn limits the low frequency response giving a bandpass type of characteristic. To widen the bandwidth, $C_s$ should be smaller and $C_j$ larger. The value of $C_s$ is primarily determined by the capacitance of the RF choke which is required to isolate the IF circuit from the RF. Considerable effort was required to design a choke with a capacitance of 1.23 pF; it is not likely that this capacitance can be reduced significantly. To increase $C_j$ would require a higher capacitance diode resulting in a complete redesign of the RF circuit. Therefore the situation cannot be improved easily by changing the circuit parameters.

The solution is in the inverse frequency variation of $R_{\text{in},s}$, equation (14), and upconverter IF gain, the ratio of equation (10) to (12). Neglecting $L$, the power transfer of the input circuit of Fig. 12 with $R_{\text{in},s}$ constant is a bandpass characteristic with band-edge slopes of 6 dB per octave. The frequency variation of $R_{\text{in},s}$ causes the slopes to be only 3 dB per octave. If the IF band lies on the low frequency side of the bandpass characteristic, the 3 dB per octave increase in power transfer will just compensate for the 3 dB per octave decrease in IF gain giving an overall constant gain relative to the available power from the generator. Since the IF band must lie below the peak of the power transfer characteristic there must be some mismatch loss. However, the bandwidth can be very wide and is limited practically by the available power from the generator and filter requirements.

The equation for the power into $R_{\text{in},s}$ relative to the generator

available power is derived in the appendix and the response in dB is plotted in Fig. 13 as a function of normalized frequency, $\omega/r\omega_o$ , for several representative combinations of $C_j$ and $L$. The diode and operating point is that given in Table II. The location of a 120 MHz band centered at 300 MHz with $r = 5$ is shown by the shaded area. The inductance $L$ is used to decrease the loss at the peak of the response and to present an open circuit to frequencies in the diode which are below the rejection band of the RF choke. The value of $L$ cannot be too large or it will resonate with the equivalent input capacitance and cause a very narrow band response as shown in Fig. 13 for $L = 80$ nH. Figure 13 also shows that the amount of mismatch loss is determined by the value of $C_j$ . The location of the power transfer peak is determined by the value of the inductance, $L$, and the shunt capacitance, $C_s$ , the latter being determined primarily by the choke capacitance.

The IF signal incident at the choke terminal of Fig. 12 must be large enough to supply the required power to the diode through the mismatched input circuit and must be supplied by a source of about 50 ohms impedance. The source must be physically close to the choke to minimize stray reactances because of the relatively high input frequency. The return loss requirement on the transmitter input
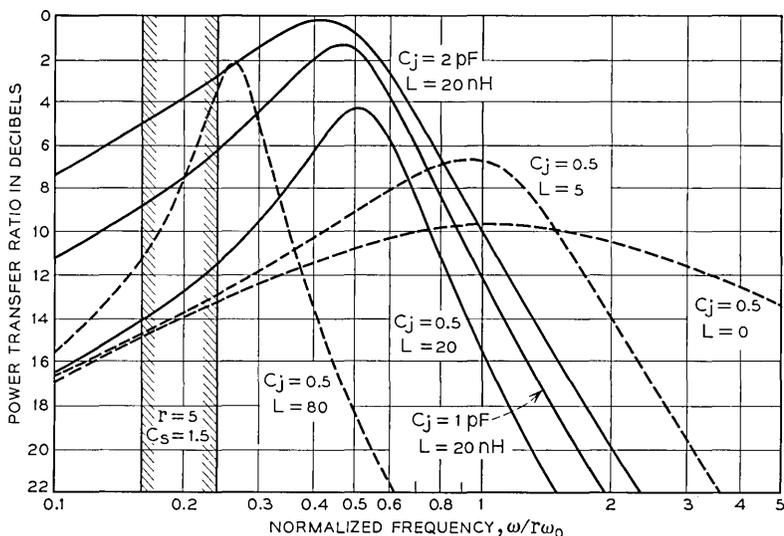


Fig. 13 — Ratio of power into $R_{in, s}$ to the generator available power versus $\omega/r\omega_o$ for several values of $C_s$, $C_j$, and $L$. The shaded band is the location of the 240 to 360 MHz band when $r = 5$ and $f_o = 300$ MHz.

in the short hop radio system experiment is approximately 18 dB over the IF band to prevent echoes in the long cable between the receiver and transmitter. All these requirements are met by using an IF driver amplifier as an integral part of the upconverter.

Figure 14 is a schematic diagram of the driver amplifier. Two Western Electric 45B transistors, used as common-base transformer coupled stages, are required because of the relatively high output power and frequency. With 0 dB input power the amplifier has 6.5 dB gain between 50 ohm terminations; the gain-frequency response is down 0.4 dB at 230 MHz and 360 MHz as shown in Fig. 15a. The input return loss, shown in Fig. 15b, is larger than 20 dB from 200 to 400 MHz. The output return loss, Fig. 15c, is 17 dB at 340 MHz and drops to 3 dB at 240 MHz as expected from the single-tuned output circuit.

## VI. UPCONVERTER MEASUREMENTS

Initial measurements were made on the 10.760 GHz upconverter using a constant 300 MHz IF signal of +5 dBm from a 50 ohm source. The IF signal was applied to a coupling capacitor, inductor, and bias circuit which were similar to $C_8$, $L_3$, and the bias circuit of Fig. 14. Then the RF circuit was tuned for maximum output power. The actual diode performance, given in Table IV compares well with the the-
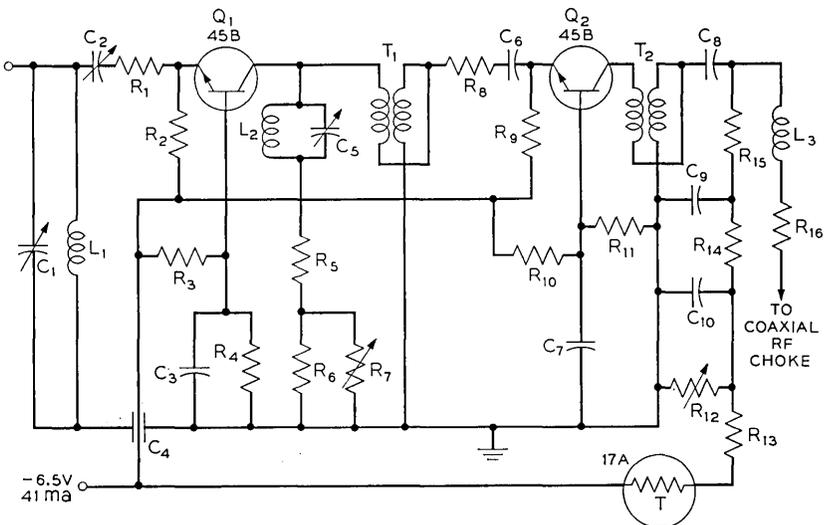


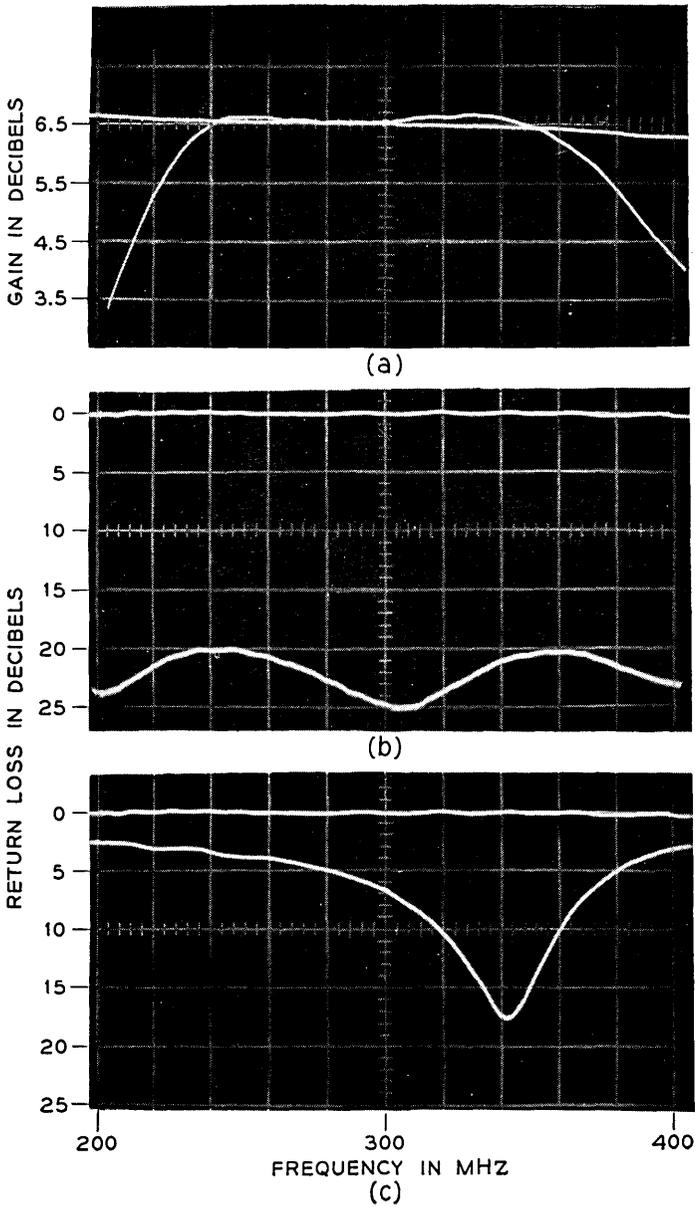Fig. 14 — Schematic diagram of upconverter driver amplifier.

Fig. 15 — IF driver amplifier (a) insertion gain, (b) input return loss, and (c) output return loss versus frequency.

TABLE IV—10.760 GHz UPCONVERTER PERFORMANCE AND DIODE
CHARACTERISTICS WITH A CONSTANT 300 MHz IF SIGNAL

| | | |
|---|---|---|
| $V_B = 32$ V | $f_u = 10.760$ GHz | $P_{in,s} = 2.5$ mw (+4 dBm) |
| $f_{co} = 112$ GHz | $f_s = 300$ MHz | $\epsilon = 44\%$ |
| $C_o = 1.075$ pF | $P_u = 42$ mw | $G = 12.3$ dB |
| $f_p = 10.460$ GHz | $P_{in,p} = 96$ mw | |

oretical performance in Table II. The output power was maximum
when the diode was self-biased at about 1.0 volt. The low bias voltage
indicates that the diode is being overdriven causing its average
junction capacitance to be larger than the theoretical value which,
from Fig. 13, accounts for the low IF circuit mismatch loss. This
initial performance showed that the 9.5 dBm available power from the
driver amplifier should be sufficient even for broadband operation.

Since the upconverters were intended to work directly into an
antenna, the effect of a poor antenna return loss on the upconverter
was measured on the 10.760 GHz upconverter by terminating it in a
variable attenuator followed by a sliding short. To determine the
effect of the phase of the reflection, several sweeps were recorded on
the same photograph as the position of the sliding short was changed.
The result is an envelope of the possible gain-frequency responses
for each return loss magnitude. The envelopes for a return loss of
> 23, 20, 10, and 3 dB are shown in Fig. 16. The upconverter was
stable and did not produce extraneous tones for any magnitude or
phase of reflection.

Figures 17a and b show the variation of output power and fre-
quency response as the pump and IF levels were varied with the
upconverter at 76°F. In Fig. 17a the IF power incident on the driver
amplifier was kept at +3 dBm and the pump power was varied in
1 dBm steps from +16 dBm to +22 dBm without retuning. There is
very little compression and little change in band shape. In Fig. 17b
the pump power was kept at +20 dBm and the incident IF power was
varied in 1 dB steps from −3 dBm to +4 dBm without retuning.
There is considerable compression and change in band shape. This
shows that the diode is being driven very hard by the IF signal and in
a sense is being pumped by the IF signal.

Figure 17c shows the return loss at both pump port and driver
amplifier input port as the input frequency was swept. The pump
return loss does not have the conventional meaning because the
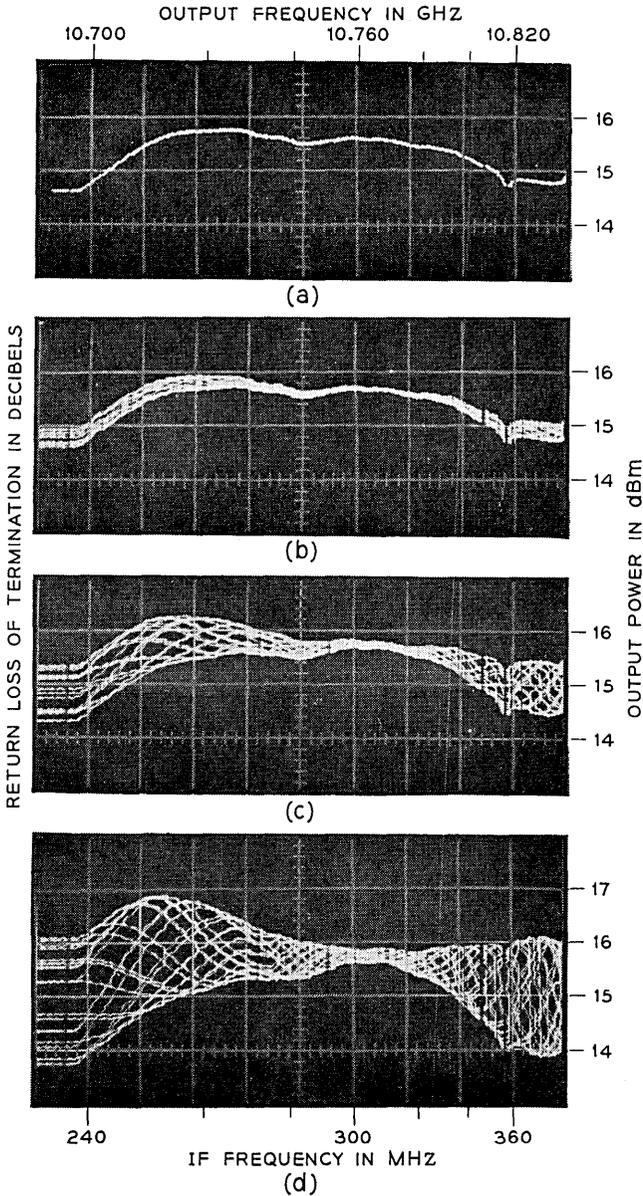pump signal was not swept; it was constant in frequency, the pump

Fig. 16 — Frequency response of the upconverter as the phase of the terminating return loss is varied for return loss magnitudes of (a) > 23 dB, (b) 20 dB, (c) 10 dB, (d) 3 dB.
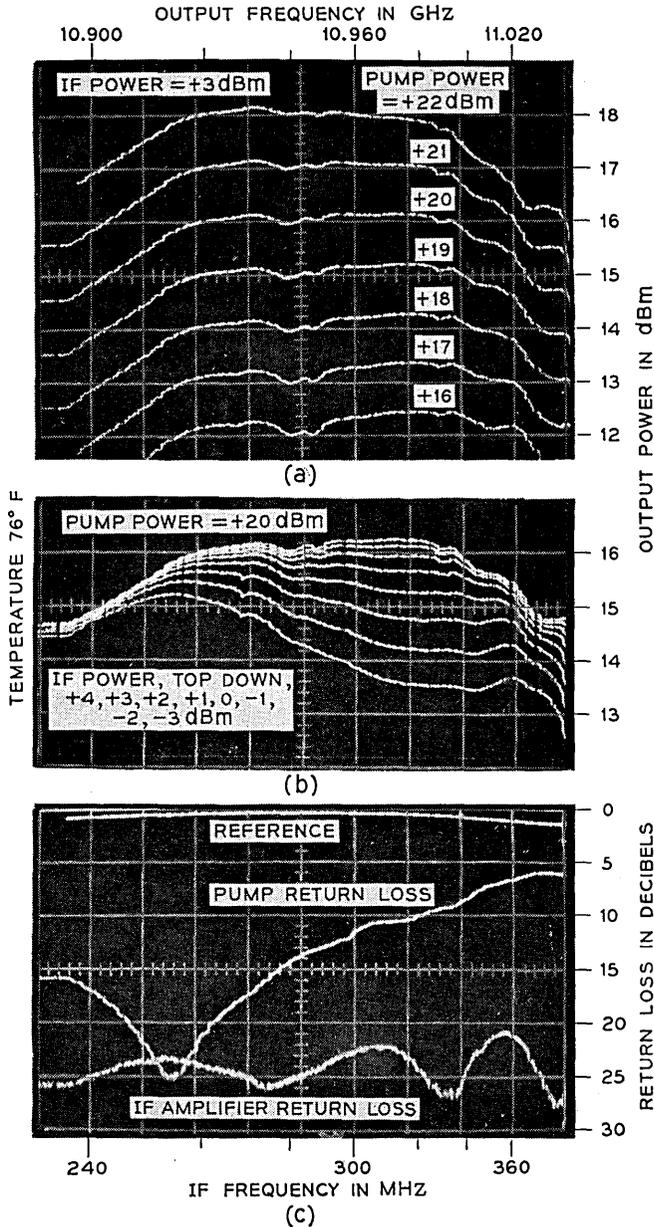
Fig. 17 — (a) Frequency response for different values of pump power. (b) Frequency response for different values of IF power. (c) Return loss at pump and IF amplifier input ports.

reflection varying because of the swept IF signal. The amplifier input return loss is not the same as in Fig. 15b because of the imperfect isolation between amplifier output and input, and because the output was terminated in the variable upconverter impedance rather than 50 ohms.

The diode bias was measured at the top of R14 in Fig. 14 with +20 dBm pump power and a +3 dBm IF signal at 300 MHz. The bias voltage was −1.5 volts and did not vary as the pump power varied but decreased 0.1 volt per dB as the incident IF power was decreased again indicating that the diode is being pumped primarily by the IF.

VII. DISCUSSION

Some bandshape and stability problems were caused by parametric oscillations in the IF circuit at frequencies below the RF choke rejection band and above the self-resonance of the series inductance in the IF amplifier output circuit. A pumped varactor diode can produce signals at any pair of frequencies whose sum is equal to the pumping frequency, if the terminating impedances are favorable at both frequencies. Since the RF choke rejects frequencies above about 10 GHz, reactive termination of all frequencies below 5 GHz reactively terminates one of the frequencies in any pair and prevents any parametric oscillation. However, this requires an inductance with a self-resonance above 5 GHz which is difficult to obtain. In retrospect a better solution would be to use a miniature 50 ohm low-pass filter, which has no spurious responses to above about 25 GHz, to replace both the inductance and RF choke. If the amplifier output were 50 ohms the filter could be considered part of the 50 ohm source and the IF circuit could be operated as discussed in Section V except that the major part of the shunt capacitance would be eliminated. The output element of the filter would have to be a shunt capacitor to provide an RF short circuit at the waveguide wall.

APPENDIX

*Derivation of Power Transfer Characteristic of IF Input Circuit*
    Referring to Fig. 12, the available power from the generator is

$$P_A = \frac{I^2 R_g}{4}.$$ (22)

The power into $R_{\text{in, }s}$ is

$$P_L = |I_2|^2 R_{\text{in,}s},$$ (23)

where $I_2$ is the loop current flowing through $R_{in, s}$ and is given by

$$I_2 = \frac{IR_g}{(1 - \omega^2 LC_s + j\omega C_s R_g)\left(R_{in,s} + \dfrac{1}{j\omega C_s} + \dfrac{1}{j\omega C_i}\right) - \left(\dfrac{1}{j\omega C_s}\right)} . \quad (24)$$

The ratio of load power to available power is

$$\frac{P_L}{P_A} = \frac{4R_g R_{in,s}}{\left|(1 - \omega^2 LC_s + j\omega C_s R_g)\left(R_{in,s} + \dfrac{1}{j\omega C_s} + \dfrac{1}{j\omega C_i}\right) - \left(\dfrac{1}{j\omega C_s}\right)\right|^2} . \quad (25)$$

The frequency variation of $R_{in,s}$ is given by equation (14). Since $f_c \gg f_s$,

$$R_{in,s} \approx \frac{22.5 \times 10^3}{C_o^{7/6} f_c^{1/2}} \frac{m_p m_s}{m_u} \left(\frac{1}{f}\right)\left(\frac{f_o}{f_o}\right) = \frac{R_{in,s}(\omega_o)}{\omega/\omega_o} , \quad (26)$$

where $\omega$ is the variable IF and $\omega_o$ is the center frequency of the IF band. From Table II the value of $R_{in,s}(\omega_o)$ is 280 ohms. Let the ratio of

$$R_{in,s}(\omega) = \frac{\omega_o}{\omega} r R_g . \quad (27)$$

Substituting (27) into (25) and taking the absolute value of the denominator gives

$$\frac{P_L}{P_A} = 4\frac{\omega}{r\omega_o} \bigg/ \left\{ \left[1 + \left(\frac{\omega}{r\omega_o}\right)\left(1 + \frac{C_s}{C_i}\right) - \left(\frac{\omega}{r\omega_o}\right)(r^2\omega_o^2 LC_s)\right]^2 \right.$$
$$\left. + \left[\frac{1}{\omega_o r R_g C_i} - \frac{\omega}{r\omega_o}(r\omega_o R_g C_s) - \left(\frac{\omega}{r\omega_o}\right)^2 \frac{r\omega_o L}{R_g}\left(1 + \frac{C_s}{C_i}\right)^2\right] \right\}. \quad (28)$$

The following frequencies are defined and substituted into (28) to get (29):

$$\omega_L = \frac{R_g}{L} , \qquad \omega_c = \frac{1}{R_g C_s} , \qquad \omega_j = \frac{1}{R_g C_i}.$$

$$\frac{P_L}{P_A} = 4\frac{\omega}{r\omega_o} \bigg/ \left\{ \left[1 + \frac{\omega}{r\omega_o}\left(1 + \frac{C_s}{C_i}\right) - \left(\frac{\omega}{r\omega_o}\right)^2 \frac{r^2\omega_o^2}{\omega_L \omega_c}\right]^2 \right.$$
$$\left. + \left[\frac{\omega_j}{r\omega_o} - \left(\frac{\omega}{r\omega_o}\right)\left(\frac{r\omega_o}{\omega_c}\right) - \left(\frac{\omega}{r\omega_o}\right)^2 \left(\frac{r\omega_o}{\omega_L}\right)\left(1 + \frac{C_s}{C_i}\right)\right]^2 \right\}. \quad (29)$$

Now let $x$ be the frequency normalized to $r\omega_o$, for example,

$$x_L = \frac{\omega_L}{r\omega_o} = \frac{R_g}{r\omega_o L}$$

giving finally

$$\frac{P_L}{P_A} = \frac{4x}{\left[1 + x\left(1 + \frac{x_i}{x_c}\right) - x^2\frac{1}{x_c x_L}\right]^2 + \left[x_i - \frac{x}{x_c} - \frac{x^2}{x_L}\left(1 + \frac{x_i}{x_c}\right)\right]^2}.$$

$$(30)$$

Equation (30) is plotted in dB in Fig. 13 as a function of the normalized frequency $x$ for several representative combinations of $C_i$ and $L$; $r$ is assumed to be 5 which approximates the ratio of $R_{in,s}(\omega_o) = 280$ ohms to a generator impedance of about 50 ohms.

REFERENCES

1. Tillotson, L. C., "Use of Frequencies Above 10 GHz for Common Carrier Applications," B.S.T.J., this issue, pp. 1563–1576.
2. Ruthroff, C. L., and Tillotson, L. C., "Interference in a Dense Radio Network," B.S.T.J., this issue, pp. 1727–1743.
3. Ruthroff, C. L., Osborne, T. L., and Bodtmann, W. F., "Short Hop Radio System Experiment," B.S.T.J., this issue, pp. 1577–1604.
4. Penfield, P., and Rafuse, R. P., *Varactor Applications,* Cambridge, Massachusetts: M.I.T. Press, 1962, p. 484.
5. Irwin, J. C., unpublished work.
6. Lee, T. P., "Evaluation of Voltage Dependent Series Resistance of Epitaxial Varactor Diodes at Microwave Frequencies," IEEE Trans. Electron Devices, *ED-12,* No. 8 (August 1965), pp. 457–470.

# Low Noise Receiving Downconverter

By T. L. OSBORNE, L. U. KIBLER, and
W. W. SNELL

*Significant improvements in noise figure and band-width have been obtained in down conversion from frequencies in the 11 GHz range to frequencies of several hundred megacycles. An average noise figure of 5.6 dB and greater than 50 percent bandwidth have been obtained using a Schottky-barrier diode balanced mixer and a low noise transistor pre-amplifier. We discuss basic design criteria, and present a complete circuit description and measurements for one of the downconverters used in the short hop radio system experiment.*

## I. INTRODUCTION

In microwave radio relay systems the receiver thermal noise is one of the main contributors to the total noise of the radio line which is one of the fundamental limitations on many important system parameters, such as the number of voice circuits per radio frequency channel, repeater spacing, and system length. Since in heterodyne systems the receiver noise is essentially the noise of the receiving downconverter, the noise figure of the downconverter is one of the most important factors in setting the entire system performance. In radio systems such as those described by Tillotson, low noise figures and wide bandwidths are required in conversion from the higher microwave frequencies above 10 GHz to higher than normal intermediate frequencies in the range of several hundred MHz.[1] The system concept also places a premium on circuit and mechanical simplicity, low power consumption, and environmental stability. The downconverter described in this paper was designed to meet the requirements of such a system.

Before the Schottky-barrier diode, noise figures below about 7 dB could only be obtained by using tunnel diodes or parametric microwave amplifiers with atendant sacrifice in receiver simplicity, stability,

and dynamic range. With the perfection of high quality GaAs Schottky-barrier diodes and low noise transistors, noise figures below 7 dB are possible without RF amplification. The 5.6 dB average noise figure and greater than 50 per cent 3 dB bandwidth achieved with the Schottky-barrier diode mixer and low noise preamplifier described here represent a significant improvement over those previously obtained in down conversion from frequencies in the 11 GHz range. Other important characteristics of the downconverter are the temperature stability, 19 dB RF to IF gain, low pump power, single dc supply, low dc power consumption, and circuit and mechanical simplicity. Table I summarizes the results achieved.

In the following sections, basic considerations in the design of a low noise Schottky-barrier diode mixer are discussed, followed by a description of the mixer and preamplifier circuitry, measurements, and performance.

## II. BASIC DESIGN CONSIDERATIONS

### 2.1 Noise Figure

The basic consideration in the design of the mixer and preamplifier is the receiver noise figure. Friis has shown that the noise figure of a receiver can be expressed as a function of three parameters,[2]

$$n_R = l_x(t_x + n_{IF} - 1), \tag{1}$$

where $n_R$ is the receiver noise figure, $l_x$ is the overall conversion loss of the mixer, $t_x$ is the equivalent output temperature of the mixer and source normalized to 290°K, and $n_{IF}$ is the noise figure of the IF amplifier following the mixer. Equation (1) can be written in terms of the normalized equivalent diode temperature, $t_{av}$,[3]

$$n_R = 1 + l_x \left[ t_{av} \left( 1 - \frac{1}{l_x} \right) + n_{IF} - 1 \right]. \tag{2}$$

The equivalent diode temperature has the following meaning. The shot noise current produced by the diode, a function of the diode conduction current, is equated to the thermal noise from a conductance equal to the diode conductance. The temperature required to make the thermal noise equal the shot noise is the equivalent diode temperature. In a pumped diode the conduction current is a periodic function of time and thus the equivalent diode temperature is a periodic function of time; $t_{av}$ is its average value. The conditions

TABLE I — SUMMARY OF PERFORMANCE

| | |
|---|---|
| Radio frequency | 10.760 GHz |
| Intermediate frequency | 300 MHz |
| Pump power (total for 2 diodes) | 10 mW |
| RF to IF gain at band center | 18.6 dB |
| Frequency response | flat ±0.1 dB, 240 to 350 MHz |
| | 0.3 dB down at 230 and 360 MHz |
| Noise figure at 75°F | 5.3 dB at band center |
| | 5.6 dB average 240 to 360 MHz |
| Temperature Range | −40 to +140°F |

under which this representation is valid have been established rigorously by Dragone.[4] In a practical diode, because of its series resistance, $R_s$, $t_{av}$ will be between 0.5 and 1.0. Equation (2) and the small range of values of $t_{av}$ show that the receiver noise figure is primarily determined by the mixer conversion loss, especially when the loss is small. Therefore, to obtain the minimum receiver noise figure, the mixer should be designed for minimum conversion loss.

## 2.2 Conversion Loss

For a resistive mixer to have minimum conversion loss, the terminations at the image frequency and all pump harmonics must have zero admittance, that is, the diode is pumped with a sinusoidal current.[5] The conversion loss is then a function of the pumping voltage and the terminating resistances at the input and output frequencies. In general, the conversion loss decreases as the pump voltage increases, but the minimum loss values of the terminating resistances increase rapidly to impractical values. If the actual terminating resistances are smaller than the resistances required for minimum loss, the conversion loss, as derived by Dragone, is

$$ l_x = 1 + \frac{2}{R} \left( R_s + \frac{KT}{qI_{co}} \right), \tag{3} $$

where $R$ is the terminating resistance, $R_s$ is the diode series resistance, and $I_{co}$ is the dc component of the diode current. Equation (3) shows that minimum values of $R_s$ and large values of $I_{co}$ are required.

## 2.3 Mixer Preamplifier Interface

In addition to the conversion loss of the mixer diode, there are circuit losses in the filters, coupler, and impedance matching circuits which must be minimized. The IF output termination is a special problem in the case of a mixer followed by a low noise preamplifier

because the noise figure of the input transistor depends on the source impedance and there is an optimum value for minimum noise. The interface between the mixer and preamplifier should be a lossless transformer to provide optimum terminations for both mixer and transistor. However, the loss of a nonideal broadband transformer can easily exceed the loss caused by nonoptimum terminating impedances.

In preliminary experiments, we found that the best performance was obtained with no impedance matching circuit at the mixer-preamplifier interface. A spot noise figure of 4.8 dB and an average noise figure of 5.4 dB over a 40 per cent bandwidth were obtained with 24 mW of local oscillator power and the diodes connected directly to the preamplifier. These results show that excellent performance can be obtained without the use of an impedance matching circuit at this interface; the downconverter reported here uses this technique.

III. DOWNCONVERTER CIRCUITRY

The downconverter consists of signal and local oscillator filters, a magic tee type of balanced mixer, and a two-transistor IF preamplifier connected as an integral part of the mixer. Figure 1 is a block diagram and Fig. 2 is a photograph of the completed unit.

3.1 *Filters*

The signal filter rejects input signals at the image frequency, reactively terminates the diode at the image frequency, and prevents leakage of local oscillator power to the antenna. The characteristics of the signal filter are important in an image rejection mixer because the filter loss at the signal frequency adds directly to the mixer conversion loss. Thus low insertion loss is required and at the same time good skirt selectivity is required to provide the image frequency termination. The signal filter is a three-section, 0.1 dB Tchebyscheff filter with a 120 MHz bandwidth and 0.2 dB midband insertion loss. Its insertion loss and return loss are shown in Fig. 3.

The local oscillator filter is a narrowband filter which attenuates undesired harmonics generated in the local oscillator multiplier chain which otherwise could enter the mixer at a higher level than the received signal. It is a three-section, 0.1 dB Tchebyscheff filter with a 30 MHz bandwidth and 1.0 dB midband insertion loss.

The electrical distance from the signal filter to the diode determines the admittance which terminates the diodes at the image frequency.

Fig. 1 — Block diagram of the receiving downconverter.



Fig. 2 — The downconverter.

Fig. 3 — Insertion loss and return loss of the signal filter.

The distance required for minimum noise figure was determined by using a continuously adjustable waveguide line stretcher as the spacer between the signal filter and mixer. Because of imperfect balance the local oscillator filter spacer has a small effect on the noise figure and also requires adjustment.

## 3.2 Mixer

A Microwave Associates *Orthotee®* balanced mixer was used because it is a compact folded magic-tee with stepped impedance transformers for matching the diode impedance to the waveguide impedance. The original diode mounts were removed and the U package diodes were mounted across the waveguide in a tunable shorted coaxial line.* The RF choke in the IF output circuit was modified to reduce the shunt capacitance to 2.8 pF with the RF rejection remaining better than 25 dB.

The diodes are matched pairs of Bell Laboratories L-2486 GaAs Schottky-barrier diodes with these typical characteristics: total capacitance at zero volts bias, $C_{To}$, 0.42 pF; series resistance, $R_s$, 0.7

---

* The U package is described in Ref. 6.

ohms; cutoff frequency, 560 GHz; and breakdown voltage, 15 volts. Operating the diodes with self bias eliminates the need for an extra bias supply ordinarily required because of the reverse polarity of the two diodes. The diodes are connected in parallel directly to the IF preamplifier input transistor to avoid loss in a matching circuit or connecting line which would increase the noise figure.

3.3 *Preamplifier*

The preamplifier is a two stage stagger-tuned transformer coupled amplifier using low noise KMC Corporation 5002 and 5003 transistors in the common-emitter configuration. Noise figures at 450 MHz for these transistors are typically 1.7 and 2.2 dB, respectively, and the small signal power gain at 450 MHz is 22 dB. Figure 4 is a schematic diagram of the preamplifier. The bandwidth is obtained by stagger tuning the collector circuits of $Q_1$ and $Q_2$. Transformer $T_1$ is tuned by $C_6$ and the collector capacitance of $Q_1$; loading is provided by $R_5$. Transformer $T_2$ is similarly tuned by $C_9$ and $Q_2$, and loading is provided by the 50 ohm input to the main IF amplifier. The transformers are 4:1 bifilar transmission line transformers wound on a Rexolite core which provides good temperature stability.

In the circuit construction, a planar layout on a copper-clad epoxy-glass circuit board, and miniature components such as $\frac{1}{10}$-watt resistors and chip capacitors were used. In Fig. 5 the preamplifier is shown with a test connector (right side) replacing one of the normal inputs from the mixer.



Fig. 4 — IF preamplifier schematic diagram.

Fig. 5 — The IF preamplifier.

IV. MEASUREMENTS AND RESULTS

4.1 *Frequency Response*

4.1.1 *Preamplifier*

For initial tuning and measurements on the preamplifier, a BNC connector was connected directly to the input transistor in the same way the connection to the mixer would be made (see Figs. 2 and 5). Using a 50 ohm sweep frequency generator, an automatic noise measuring set, and a 50 ohm load, the preamplifier was adjusted for minimum noise figure and best band shape.

At 75°F and −22 dBm output power, the gain of the preamplifier alone when operating between a 50 ohm source and a 50 ohm load is 24 dB. The gain-frequency response, shown in Fig. 6, is flat to within 0.1 dB and the gain is down 0.3 dB at 240 and 370 MHz. At −40°F and +140°F the gain decreases by 0.4 dB and the gain-frequency response has a 0.1 dB positive slope across the 130 MHz band. The output power for which there is 0.2 dB gain compression is −14 dBm.

4.1.2 *Downconverter*

With the preamplifier attached to the mixer, the filter spacers, tuning screws, and diode mounts were adjusted for minimum overall

TEMPERATURE 75°F



Fig. 6 — Gain-frequency response of the IF preamplifier.

noise figure and flat bandshape. Frequency response measurements were made with a leveled X-band backward wave oscillator sweep frequency generator and the IF output was displayed in dB on an Alfred sweep network analyzer. Except as noted, all measurements were made with a signal power of −40 dBm into the signal filter and a pump power of +10 dBm into the pump filter in accordance with system requirements.

At 75°F, the overall gain of the downconverter is 18.6 dB. The gain-frequency response, Fig. 7, is flat to within 0.1 dB and the gain is down 0.3 dB at 230 and 360 MHz. As shown in Fig. 8, the gain remains at 18.6 dB at −40°F but decreases to 17.7 dB at +140°F.

In system use, the pump power will change somewhat with temperature and from one pump source to another. The effect on the band shape is shown in Fig. 9. For a decrease in pump power from +10 dBm to +8 dBm there is a decrease in gain of 0.7 dB and a positive slope of 0.8 dB across the 120 MHz band. This slope can be corrected by minor retuning of the preamplifier.

The overload characteristics of the downconverter are shown in Fig. 10. The normal output power level with −40 dBm input is −21.4 dBm. At −16 dBm some slope is evident and at −14 dBm output the gain has decreased by 0.2 dB and there is considerable slope in the band shape. This compression primarily is due to the preamplifier since this is the same amount of compression measured for the preamplifier alone.

TEMPERATURE 75°F



Fig. 7 — Gain-frequency response of the downconverter.

The input return loss at the input to the signal filter is more than 15 dB across the 120 MHz band. The return loss at the narrow band input to the pump filter at the pump frequency is more than 30 dB. The isolation from the pump port to the signal port of the hybrid, a measure of the mixer balance, is 13 dB.

### 4.2 Noise Figure

Noise figure measurements on the preamplifier were made using a Hewlett-Packard 343A temperature limited diode as the noise source.



Fig. 8 — Downconverter frequency response at −40, +75, and +140°F.

Fig. 9 — Downconverter frequency response as a function of pump power.



Fig. 10 — Downconverter frequency response as a function of IF output power.

Fig. 11 — Preamplifier and downconverter noise figures as a function of frequency and temperature.

The preamplifier was followed by an Avantek amplifier, an AN/APR-4 radar receiver with 1 MHz bandwidth for converting the UHF signal to 30 MHz, and a Hewlett-Packard 342A noise figure meter. A Hewlett-Packard X347A argon gas discharge tube was used as a noise source for the downconverter. For both preamplifier and downconverter, the spot noise figure in a 1 MHz band was measured as a function of frequency from 200 to 400 MHz and at temperatures of −40, +75, and +140°F. The results are plotted in Fig. 11.

At 75°F, the preamplifier spot noise figure at 300 MHz is 1.98 dB and the preamplifier average noise figure for the 240 to 360 MHz band is 2.05 dB. The corresponding numbers for the downconverter are 5.3 and 5.6 dB. At −40°F the downconverter average noise figure decreases to 4.9 dB and at +140°F increases to 6.5 dB.

V. DISCUSSION

The noise figures obtained in the downconverter reported here are equal to or better than those obtainable with tunnel diode amplifiers and are approaching those obtainable with uncooled parametric amplifiers. However, Dragone's analysis shows that the noise figures reported here are still far above those theoretically possible with Schottky-barrier diode mixers. Theoretically the mixer conversion loss can be as small as a few tenths of a dB and, since noise figures of 1.5 dB have been obtained for transistor amplifiers, an overall noise figure under 2 dB is theoretically possible. The practical problems involved in achieving extremely low conversion loss center around the terminating impedances, particularly at the pump harmonics. In a waveguide circuit such as the magic tee, which has the advantage of low loss, control of harmonic terminations is difficult because of multimoding. On the other hand, in single mode TEM lines, line loss has prevented low conversion losses. Therefore, much lower conversion losses will require careful attention to the RF circuit construction.

When calculating the mixer conversion loss from equation (2) assuming $t_{av} = 1$, its largest value, the conversion loss is the difference between the overall noise figure and the preamplifier noise figure. Using 300 MHz spot noise figures gives $5.3 - 2.0 = 3.3$ dB.

REFERENCES

1. Tillotson, L. C., "Use of Frequencies Above 10 GHz for Common Carrier Applications," B.S.T.J., this issue, pp. 1563–1576.
2. Friis, H. T., "Noise Figures of Radio Receivers," Proc. IRE, *32*, No. 7 (July 1944), pp. 419–423.
3. Messinger, G. C., and McCoy, C. T., "Theory and Operation of Crystal Diodes as Mixers," Proc. IRE, *45*, No. 9 (September 1957), pp. 1269–1283.
4. Dragone, C., "Analysis of Thermal and Shot Noise in Pumped Resistive Diodes," B.S.T.J., *47*, No. 9 (November 1968), pp. 1883–1902.
5. Dragone, C., "Amplitude and Phase Modulations in Resistive Diode Mixers," B.S.T.J., this issue, pp. 1967–1998.
6. Lee, T. P., "Evaluation of Voltage Dependent Series Resistance of Epitaxial Varactor Diodes at Microwave Frequencies," IEEE Trans. Electron Devices, *ED-12*, No. 8 (August 1965), pp. 457–470.

# Broadband 300 MHz IF Amplifier Design

### By W. F. BODTMANN and F. E. GUILFOYLE

(Manuscript received October 29, 1968)

*Short hop radio systems of the type described in this issue require broadband intermediate frequency amplifiers in the range of a few hundred MHz which have a large gain control range and stable characteristics with respect to temperature; they must also make efficient use of bias power and have well behaved noise properties.*

*This paper describes the design of such amplifiers. We give extensive measurements for an amplifier with a 1 dB bandwidth of 120 MHz centered at 300 MHz and an automatic gain control range of 43 dB.*

## I. INTRODUCTION

The incentive for this work arises from the need for IF amplifiers in short hop radio systems. Such systems may be powered by thermoelectric generators driven by propane gas and have the electronic repeater mounted at the top of a pole. Amplifiers used for this application must:

(*i*) Be sufficiently broadband to amplify, with low distortion, large index analog FM or PSK-PCM signals.

(*ii*) Have good temperature stability since they operate at outdoor ambient temperatures.

(*iii*) Be small and lightweight to fit into the limited space of the pole mounted enclosure and to keep the weight loading at the top of the mounting to a minimum.

(*iv*) Have low power consumption to operate from a thermoelectric generator which is the primary power source.

(*v*) Be designed, if possible, with a thought toward the use of integrated circuitry for future amplifiers.

In the amplifier designed to satisfy these requirements, most of the gain is obtained in broadband common emitter stages. Transistors with a low current gain are used so that the interstage is simple,

stable, and has good power handling capabilities. Broadband vario-
lossers, separating the early gain stages, provide the necessary auto-
matic gain control range. Band shaping is obtained in the transformer
coupled output stages.

## II. LOW LEVEL STAGES

Most of the gain and all of the gain control range are accomplished
in the low level section of the amplifier. It is important to do as much
as possible at low signal levels so that the total power consumption—
an important system consideration—is minimized. The low level sec-
tion consists of eight transistors connected alternately in common
base and common emitter configurations. The common emitter stage
provides current gain, while the interstage network, coupling the
common emitter common base stage, is part of a broadband variolosser
used in the automatic gain control circuits. The common emitter com-
mon base gain-variolosser stages are coupled together by a simple
peaking circuit to distribute both the gain and variable loss through-
out the amplifier. This gain-variolosser-coupler arrangement illus-
trated schematically in Fig. 1 was chosen for the following reasons:

(*i*) Gain is achieved through the current gain of the common emit-
ter transistor, obviating the need for transformers and lending the
interstage to future circuit integration.



Fig. 1 — Common emitter-common base gain-variolosser stage and peaking
circuit.

(*ii*) The output impedance of the common emitter stage is readily adapted to a simple broadband variolosser circuit.

(*iii*) The zero frequency current gain of the transistor $\beta_o$, can be tailored to the requirements on gain and bandwidth. For the low $\beta_o$ used in this amplifier, the result is a stable gain-frequency characteristic achieved without the use of feedback circuits.

(*iv*) The common base stage provides isolation between the variolosser and the peaking circuit. This configuration is a satisfactory load for the broadband variolosser and provides a stable broadband high impedance generator for the common emitter stage.

(*v*) A low intermediate frequency favors low power consumption and increases the gain of the upper-sideband varactor amplifier which is driven by the IF amplifier.[1] The intermediate frequency of 300 MHz chosen for this application is sufficiently high to obtain a 120 MHz IF bandwidth.

## 2.1 *Common Emitter Gain Stage*

The common emitter transistor is a Bell Telephone Laboratories L-2526 microwave transistor chosen to have a low frequency current gain between 8 and 10 and a minimum gain bandwidth product $F_T = 2{,}500$ MHz. With this combination of high $F_T$ and low $\beta_o$ substantial current gain and a reasonable high frequency characteristic are obtained without using external feedback. In addition, the low $\beta_o$ transistor has higher power capability than a high $\beta_o$ transistor with external feedback unless the feedback can be maintained down to dc. The difference in power handling capability can be seen by referring to Fig. 2. This plot shows a transistor with a high $\beta_o$ which has sufficient feedback to reduce the high frequency current gain to $\beta_1$. The feedback is not maintained down to dc. Curve 2 shows the frequency characteristics of a low $\beta$ transistor with a low frequency current gain equal to $\beta_1$. Consider first the high $\beta$ transistor. The ratio of collector current to base current at dc is given by

$$\frac{I_{c_{dc}}}{I_{b_{dc}}} = \beta_o \ . \tag{1}$$

At some high frequency $f_1$ the ratio of collector current to base current is

$$\frac{I_c}{I_b} \cong \beta_1 \ . \tag{2}$$

Fig. 2 — Common emitter transistor characteristics.

If the base is not driven beyond cutoff, then the peak current swing in the base at frequency $f_1$ is limited to

$$I_{bpeak} = I_{bd_o} = \frac{I_{cd_o}}{\beta_o}. \tag{3}$$

Combining (2) and (3) yields the peak collector current swing at frequency $f_1$ for a high $\beta$ transistor

$$I_{cpeak} = I_{cd_o} \left(\frac{\beta_1}{\beta_o}\right). \tag{4}$$

Thus, the peak collector current is reduced by the ratio of the high frequency to low frequency current gain. For the low $\beta$ transistor, $\beta_1/\beta_o \cong 1$ and

$$I_{cpeak} \cong I_{cd_o} .$$

The peak collector current of a high $\beta$ unit is $(\beta_1/\beta_o)$ less than that obtainable from the low $\beta$ unit, and the power output is reduced by the square of this ratio. This is important since the variolosser must be operated at a reasonably high power level to maintain a good amplifier noise figure. The power dissipation is minimized by the proper choice of $\beta_o$.

2.1.1 *Common Emitter Frequency Characteristic*

The frequency characteristic of the common emitter stage as a result of the variation of $\beta$ is given by[2]

$$\beta = \frac{\beta_o}{1 + \dfrac{jf}{(1 - \alpha_o)f_{ab}}} \tag{5}$$

where

$\alpha_o$ = low frequency grounded base current gain
$\beta_o$ = low frequency grounded emitter current gain

and

$f_{ab}$ = alpha cutoff frequency.

Assuming a $\beta_o$ of 9 and an $f_{ab}$ of 3.8 GHz for the 2526 transistor, the frequency characteristic given by equation (5) is shown in Fig. 3, curve 1. The 240 to 360 MHz band indicated by the arrows is the IF band of interest. There is a negative slope across the IF band of 1.3 dB. Because this slope is small, it can be readily compensated for by a peaking circuit, and the IF band will be flat and stable.

### 2.1.2 Common Base Common Emitter Coupling and Peaking Circuit

As shown in Fig. 1, the output of the common base stage Q3 is a capacitor shunted by a resistor of approximately 1K. The input impedance of the common emitter stage Q4 is also capacitive, and has approximately 35 ohms in shunt. The addition of one inductor $L_3$, forms a simple series peaking circuit which has been used extensively in tube amplifiers.[3] This interstage circuit can be adjusted to have a positive slope across the band which compensates for the negative

Fig. 3 — Low power gain stage characteristics.

slope of the $\beta$ characteristic. The slope can be varied by changing $R_4$ or $L_3$. The peaking circuit shown in Fig. 1 has been redrawn in Fig. 4 where

$$Y_1 = \frac{1}{R_3} + j\omega C_3 , \qquad (6)$$

$$Y_2 = \frac{1}{j\omega L_3} , \qquad (7)$$

$$Y_3 = j\omega C_4 , \qquad (8)$$

$$Y_4 = \frac{1}{R_4}. \qquad (9)$$

The current transfer ratio for this network is

$$\frac{I_4}{I_o} = \frac{y_4}{\left[ y_1 + y_3 + y_4 + \dfrac{y_1(y_3 + y_4)}{y_2} \right]}. \qquad (10)$$

Substituting (6), (7), (8), and (9) into (10), the transfer function becomes

$$\left| \frac{I_4}{I_o} \right|^2 = \frac{1}{\left[ 1 + K - \left( \dfrac{C_3 + K C_4}{C} \right)\gamma^2 \right]^2 + \left[ \left( \dfrac{1 + Q_3 Q_4}{Q_4} \right)\gamma - \dfrac{C_3 C_4}{C^2 Q_4} \gamma^3 \right]^2} \qquad (11)$$

where

$$\gamma = \frac{\omega}{\omega_o} , \qquad \omega_o = \frac{1}{(LC)^{\frac{1}{2}}} , \qquad C = C_3 + C_4 ,$$

$$K = \frac{R_4}{R_3} , \qquad Q_3 = \frac{1}{\omega_o C R_3} , \qquad Q_4 = \frac{1}{\omega_o C R_4}.$$



Fig. 4 — Common base-common emitter coupling and peaking circuit. $Y_1 = 1/R_3 + j\omega C_3$; $Y_2 = 1/j\omega L_3$; $Y_3 = j\omega C_4$; $Y_4 = 1/R_4$.

Equation (11) includes the loading from the common base as well as the common emitter stage and expresses the amplitude characteristic in terms of the currents which is a useful form for transistor circuits. Equation (11) has been plotted in Fig. 3, curve 2, for the case where $C_3 = C_4 = 2.5\rho F$, $L_3 = 25$ nH, $R_3 = 1K$, and $R_4 = 35\Omega$. This curve shows the positive slope needed to compensate for the fall off in the $\beta$ characteristic. The peaking circuit is adjusted for a flat amplitude by adjusting $L_3$.

Curve 3 of Fig. 3 shows the complete transfer function of the gain stage; it is the product of the functions shown in curves 1 and 2. The amplitude characteristic is flat from 200 to 400 MHz which is considerably wider than the desired IF band, and it also shows that a peak will occur at approximately twice the center frequency of the IF band. Filtering proceeding the low level stages, and filtering resulting from the IF transformers following the low level stages, attenuate this peak so that it has a small effect on the overall amplifier characteristic.

## 2.2 Variolosser

The electronically variable attenuator used in this amplifier uses one Western Electric 2480 PIN diode to vary the loading on the primary of a double tuned circuit and in this manner changes the interstage loss. This attenuator is described in detail in Ref. 4 where it has been shown that at the midband, the interstage loss is given by

$$\text{loss in dB} = 20 \log \frac{1}{1 + G_1 R_2}. \tag{12}$$

$R_2$ is the total resistance in the emitter of transistor Q3 of Fig. 1 and $G_1$ is the parallel combination of the output resistance of Q2 and the PIN diode $R_D$. Thus, the interstage loss can be changed over a wide range by varying $G_1$ or $R_2$.

The parallel—series tuned circuit incorporates the output capacity of the common emitter stage and the input inductance of the common base stage as part of the interstage circuit. The interstage is designed to have a bandwidth considerably wider than the IF band and to have flat transmission characteristics at the minimum loss and at the nominal operating loss points. As the variolosser is varied from the maximum attenuation of 17 dB to the minimum attenuation of 2.7 dB there is little change in the band shape, that is, less than 0.1 dB across the band. Three variolossers, used as part of the low level stages, distrib-

ute the loss throughout the amplifier so that the attenuation can be varied over a 43 dB range and the main amplifier will still have a reasonable noise figure at the high attenuation point.

## III. DRIVER AND OUTPUT STAGES

The driver and the output stages are common base transformer coupled stages as shown in Fig. 16. Two Western Electric 2254 transistors are used in parallel in the output stage to provide +7.5 dBm of power into a 50 ohm load. The transformers are parallel-series tuned transformers designed to have 0.25 dB loss at 240 and 360 MHz. They provide the primary band shaping for the amplifier and are also used to provide:

(*i*) 6 dB current gain in the interstage and to match the common emitter transistor Q8 and the common base transistor Q9;

(*ii*) 4.5 dB current gain between transistor Q9 and the output stage. This transformer is designed as a mismatched transformer with the loading on the secondary;

(*iii*) An impedance match between the output transistors and the 50 ohm load.

The output stage can supply +7.5 dBm of power into a 50 ohm load with 0.5 dB of compression as shown by Fig. 5. This curve is a plot of the ouput power into a 50 ohm load as a function of the input power to the amplifier.

## IV. NOISE FIGURE

The noise figure of a main IF amplifier used in a radio system should not appreciably increase the overall noise figure of the system. Thus, a satisfactory noise figure for the main amplifier is affected by preamplifier gain and down converter preamplifier noise figure. The main amplifier spot noise figure (at 300 MHz) varies from 8.6 dB at maximum gain to 13.6 dB at the normal operating level, and has a value of 14.8 dB at minimum gain. At the normal operating level, the main amplifier will increase the receiver noise figure by 0.4 dB.

The noise figure of the amplifier as a function of frequency with gain as a parameter is shown in Fig. 6. This noise figure is measured in a 1 MHz band.

### 4.1 *Variolosser Placement*

The amplifier noise figure is affected by the number of variolossers and their distribution throughout the amplifier. Variolossers placed

Fig. 5 — Compression characteristic of amplifier.



Fig. 6 — Noise figure verses frequency and amplifier gain.

in the early stages of the amplifier must have a limited range if a low noise figure is required. Variolossers placed near the output of the main amplifier can have a large attenuation range; however, the early stages of the amplifier and the variolosser will then operate at high power levels which results in greater dc power consumption.

The noise figure of several networks in cascade is given by the expression[5,6]

$$F_T = F_1 + \frac{F_2 - 1}{G_1} + \frac{F_3 - 1}{G_1 G_2} \cdots , \tag{13}$$

where

$F_T$ is the overall noise figure,

$F_1, F_2, F_3$ are the noise figures of each cascade section,

and

$G_1, G_2, \ldots$ are available gains of each cascade section.

If $G_1$ can be made large, then the noise figure is essentially $F_1$. If $G_1$ is small, $F_2, F_3$, and so on, contribute to the overall noise figure.

As used in the radio system, the main IF amplifier is preceded by a preamplifier with 24 dB gain. The available gain from this amplifier greatly reduces the noise contribution to the system noise figure from the main amplifier and places the first stages of the main amplifier at an intermediate power level.

Three variolossers are placed in the first three stages of the amplifier, and the attenuation range of each one is limited to 14 dB, primarily by noise figure considerations. Several variolossers, each with a moderate attenuation range, placed at an intermediate power level, provide a good compromise between radio system noise figure, the number of variolossers required, and the dc power consumption. The available gain of each amplifier-variolosser section varies with attenuation setting and approaches a value slightly greater than one at the nominal amplifier gain of 33 dB. Thus, from equation (13), the noise figure should vary with amplifier gain. The variation in noise figure as a function of amplifier gain and frequency is shown in Fig. 6.

### 4.2 Distribution of Gain and Loss

The distribution of the gain and loss throughout the amplifier is shown in Fig. 7. A gain stage is followed by a variable attenuator stage so that each low level section has a maximum gain of 14.5 dB. Three such sections with a maximum gain of 43.5 dB are followed by a common emitter transformer coupled stage which provides 22 dB

Fig. 7 — Gain-loss distribution of IF amplifier at maximum gain.

of fixed gain. The driver stage has a net gain of 4 dB while the output stage matches the amplifier to 50 ohms with 1.5 dB loss.

## V. AUTOMATIC GAIN CONTROL

The automatic gain control samples the output power and provides an error signal to the variolossers to correct for changes in amplifier input power or for amplifier gain as a function of temperature. As shown by Figs. 8 and 16, the automatic gain control consists of a detector diode D4, two dc amplifiers Q12 and Q13, and variolosser diodes D1, D2, and D3. The amplifier output is sampled by diode D4 and any error signals are amplified by the dc amplifiers to change the



Fig. 8 — Automatic gain control schematic.

Fig. 9 — PIN diode current verses attenuation.

current through the three PIN diodes. The resistance of the PIN diodes varies with the diode current to form a variable current divider in the IF path. The use of a PIN diode to vary the interstage loss is discussed in Section 2.2 and in Ref. 4. As used in this amplifier, the PIN diode will provide a large change in attenuation for a small variation in diode current. This is shown in Fig. 9 which also shows the attenuation of three variolossers as a function of the total diode current. A 43 dB change requires slightly over 2 milliamperes. This small variation in current makes the power requirements of the amplifier virtually constant at all gain settings.

The loop gain of the AGC circuit varies with variolosser setting; however, it is greater than 25 dB over most of the attenuation range. The output power decreases 1.2 dB for a 30 dB decrease in signal level.

## 5.1 Manual Gain Control

A manual gain control is necessary to adjust the gain of the IF amplifier when some radio system measurements are made. In order for the amplifier to reproduce a swept input signal accurately, the response time of the AGC must be long compared with the sweep rate of the generator. For a generator using a 60 cycle sweep rate, the 0.4 second time constant of the AGC loop is satisfactory; however, in some system measurements, the sweep rate can be as low as one cycle per second. The manual gain control clamps the gain of the amplifier at a selected level; the IF amplifier characteristic can then be measured without being distorted by the AGC. This control is shown in

Fig. 10. The main requirement for this control is that it can supply a stable negative voltage that can be varied from $-3.8$ to $-4.1$ volts.

### 5.2 *AGC Monitor*

The AGC monitor is used to indicate the input power to the main IF amplifier. As shown in Fig. 8, this monitor measures the voltage across the collector-emitter junction of transistor Q13. The meter has been calibrated as shown in Fig. 11 to indicate the input power to the amplifier in dBm with temperature as a parameter. In the system experiment the AGC monitor is used as a remote indicator approximately 60 feet from the main IF amplifier.

### VI. AMPLIFIER PERFORMANCE

The performance of the amplifier as a function of gain control setting, temperature, and power supply variations is shown by the series of pictures in Figs. 12, 13, and 14. The input and output impedance match is shown in Fig. 15.

### 6.1 *Amplifier Performance with Gain Changes*

The performance of the amplifier as a function of the gain control setting is shown in Fig. 12. This series of curves show three characteristics of the amplifier:



Fig. 10 — Manual gain control schematic (J1-battery charge jack; S1-battery switch; S2-manual or AGC selector; S3-battery test).

(*i*) The basic amplitude characteristic over the 240 to 360 MHz range (−26 dBm input signal level),

(*ii*) The distortion contributed to the amplitude characteristic by the three variolossers as the gain is varied from 26 dB (−18 dBm input) to the maximum gain of 67 dB (−62 dBm input),

(*iii*) The performance of the AGC loop by showing the output power changes as a function of the input power.

Each photograph shows the amplitude characteristic as a function of frequency with input power as a parameter. The vertical scale is 1 dB per division. The input power level is shown for each photo as well as an output power level of +7.5 dBm. Figure 12a shows the characteristic when the amplifier is driven with a −18 dBm signal level. This is 8 dB greater than the normal signal level. Figure 12b shows the characteristic for a normal input signal level of −26 dBm.

Fig. 11 — Automatic gain control moniter calibration.

Fig. 12 — Amplitude characteristic and output power as a function of input power at (a) —18 dBm, (b) —26 dBm, (c) —36 dBm, (d) —46 dBm, (e) —56 dBm, and (f) —62 dBm.

Fig. 13 — Amplitude characteristic verses temperature at (a) 140°F, (b) 100°F, (c) 75°F, (d) 0°F, (e) −25°F, and (f) −40°F.

Fig. 14 — Amplitude characteristic and output power as a function of dc supply voltage at (a) 6.5 V, (b) 6.175 V, and (c) 6.825 V.

Fig. 15 — Return loss verses frequency for amplifier input and output circuits: (a) IF input, and (b) IF output.

This is the point where the amplifier was adjusted for a flat transmission characteristic with an output power of +7.5 dBm. Figures 12c through f show the frequency characteristic and the output power delivered by the amplifier for decreasing signal levels.

### 6.2 Amplifier Performance with Temperature Changes

The amplifier must operate in a temperature range from −40°F to +135°F. The variation in the amplitude characteristic as a function of temperature is shown by the series of curves in Fig. 13. The input signal level was held constant at −26 dBm and the output power was set at +7.5 dBm at 75°F. The vertical scale is 1 dB per division. In addition to the amplitude-frequency characteristic, the curves show how the output power changes with temperature from the value set at 75°F.

### 6.3 *Amplifier Performance with Power Supply Variations*

The −6.5 volts supplied by the radio system power supply will vary less than ±2 percent for expected input voltage variations, load variations, and temperature changes.[7] The three curves in Fig. 14 show the frequency characteristic and the change in output power for supply voltage that are normal, 5 percent high, and 5 percent low. The vertical scale is 1 dB per division. These curves show that for ± 2 percent change in power supply voltage, the output power will vary ±0.4 dB and the frequency amplitude characteristic will remain constant.

### 6.4 *Input Impedance*

The input transformer, used to match a 50 ohm source to the impedance of the common base stage, consists of a low pass circuit as shown in Fig. 16. The input impedance of transistor Q1 is inductive and resistive. A resistor in series with the input of transistor Q1 adjusts the $Q$ of the network to provide a flat transmission characteristic over a wide frequency range.

The measured return loss of the input circuit from 240 to 360 MHz is shown in Fig. 15a. The return loss is 17 dB at the center of the band and approximately 15 dB at the band edges.

### 6.5 *Output Impedance*

The amplifier has an output impedance of 50 ohms, and it is matched by means of a double tuned transformer. This circuit is shown in the collector stage of Q10 and Q11 of Fig. 16. The load resistor for the parallel-tuned series-tuned transformer is also used as a dc return for the automatic gain control diode detector D4.

The measured return loss of the output circuit from 240 to 360 MHz is shown by the photo in Fig. 15b. The return loss in the center of the band is 22 dB; it decreases to 10 dB at the band edges.

### VII. MECHANICAL DESIGN

The IF amplifier is constructed on a $\frac{1}{16}$ inch thick epoxy glass board, laminated on both sides with two ounce copper and enclosed in a brass box with removable top and bottom covers. The enclosure with covers in place measures $1\frac{3}{4}$ inches × $1\frac{3}{4}$ inches × $10\frac{1}{2}$ inches. The transistor bias and RF circuits are mounted on the top of the board and the AGC circuits on the bottom as shown in Figs. 17 and 18. This type of construction is used to keep lead lengths in the RF circuits to a minimum. To insure a common ground plane, the board is fastened

Q₁,₃,₅,₇ TIX 2N2999    Q₂,₄,₆,₈  WE-2526    Q₁₂ 2N844(FET)    D₁,₂,₃ WE-2480

Q₉,₁₀,₁₁  WE-2254    Q₁₃ 2N329A    D₄ WE-404B

Fig. 16 — IF amplifier schematic.

Fig. 17 — Top view of IF amplifier.

to a $\frac{3}{16}$ inch wide support on the inner walls of the box with screws spaced approximately $1\frac{1}{4}$ inches apart. The covers are constructed to overlap and to be screwed to the outside walls of the box. Silver-plated phosphor bronze fingers, attached to the inside of the covers, make contact with the inner walls of the box. This shielding is necessary because the amplifier is normally operated close to a 320 MHz source with an output power 76 dB greater than the minimum input power to the IF amplifier. With the covers in place, no interference from the power source was observed in the amplifier.



Fig. 18 — Bottom view and cover of IF amplifier.

VIII. DISCUSSION

It has been demonstrated that a high performance wide band amplifier can be made using a combination of common emitter-common base stages to obtain current gain and a stable wideband variolosser. The use of low $\beta_o$ transistors leads to a simple interstage with good power handling capabilities and a stable temperature characteristic. The low power consumption of 0.65 watts for an amplifier with a maximum gain of 67 dB meets the short hop radio system requirements.

REFERENCES

1. Osborne, T. L., "Design of Efficient Broadband Varactor Upconverters," B.S.T.J., this issue, pp. 1623–1649.
2. Phillips, A. B., *Transistor Engineering,* New York: McGraw-Hill, pp. 300–304.
3. Rideout, V. C., *Active Networks,* Englewood Cliffs, New Jersey.: Prentice-Hall, 1954, pp. 99–102.
4. Bodtmann, W. F., "Design of Efficient Broadband Variolossers," B.S.T.J., this issue, pp. 1687–1702.
5. Friis, H. T., "Noise Figure of Radio Receivers," Proc. IRE, *32,* No. 7 (July 1944), pp. 419–422.
6. Roberts, S., "Some Considerations Governing Noise Measurements on Crystal Mixers," Proc. IRE, *35,* No. 3 (March 1947), pp. 257–265.
7. Bodtmann, W. F., unpublished work.

# Design of Efficient Broadband Variolossers

By W. F. BODTMANN

*This paper describes new variolosser circuits which are suitable for use in high quality broadband IF amplifiers operating at several hundred MHz, and offering substantial improvements in power consumption, temperature stability, bandshape stability, minimum attenuation and attenuation range. Three variolossers of this type are used in the automatic gain control for the IF amplifier used in the short hop radio system experiment.*

*We present an analysis, a design procedure, and extensive measurements. As an example of the performance to be expected, one experimental vario- loser has a minimum attenuation of 1.35 dB, a maximum attenuation of 19.5 dB, and is flat to ±0.15 dB over a 40 percent band centered at 300 MHz.*

## I. INTRODUCTION

The gain control of broadband high quality transistor IF amplifiers has been a difficult problem for many years. Control of gain by varia- tion of bias on the gain stages causes severe bandshape distortion, so it is customary to use one or more semiconductor diode variolossers for this purpose. The use of a variolosser, which introduces a large attenuation to the signal, can impair the noise figure of the amplifier; for this reason several variolossers of limited range are imbedded be- tween suitable gain stages to minimize the effect on the noise figure. Several types of variolossers have been used but all of them have one or more of the following disadvantages.

($i$) Complicated circuitry, making adjustment difficult and thermal stability hard to achieve.

($ii$) Large minimum loss, which requires additional gain stages and adds to the noise figure problem.

($iii$) Large bias power.

($iv$) Parasitic reactances degrading the performance at high fre- quencies.

1687

For some applications, power efficiency, thermal stability, and size
are important parameters. This is true, for example, of the short hop
radio system which has a pole-mounted repeater operating from a
primary power source.[1] The circuits described in this paper were de-
signed for such applications.

The basic variolosser circuit is a double-tuned IF stage in which
a PIN diode is inserted to provide variable attenuation. At frequencies
of several hundred MHz the parasitic reactances of the diode are
important; it is a feature of the circuit that the parasitics can be
absorbed in a simple way to improve the performance of the vario-
losser. For example, for a given bandwidth and flatness requirement,
the attenuation range is doubled in dB, (see Fig. 4).

Bandwidths of 50 percent, flat to a few tenths of a dB over an at-
tenuation range of 20 dB have been achieved with these circuits.

## II. ANALYSIS OF BASIC CIRCUIT

The basic variolosser circuit is shown in Fig. 1 without bias circuits.
The loss of the double-tuned interstage circuit is controlled by a PIN
diode as indicated in Fig. 1a. An equivalent circuit including the out-
put impedance of the transistor Q1 and the input impedance of Q2
is shown in Fig. 1b; $L_2$ is the sum of $L_s$ and the input inductance of
Q2, $R_2$ is the input resistance of Q2 plus the resistance $R_s$, and $R_1$ is
the parallel resistance combination of the PIN diode and the output



Fig. 1 — Basic interstage circuit.

resistance of Q1. The interstage has been further simplified in Fig. **1c** where,

$$Z_2 = R_2 + j\left(\omega L_2 - \frac{1}{\omega C_2}\right), \tag{1}$$

$$Y_1 = G_1 + j\left(\omega C_1 - \frac{1}{\omega L_1}\right), \tag{2}$$

$$G_1 = \frac{1}{R_1}.$$

The design approach is to make the transfer function as flat as possible for a prescribed attenuation range and bandwidth. From Fig. 1c, the current transfer function is

$$\frac{I_2}{I_0} = \frac{1}{1 + Z_2 Y_1}. \tag{3}$$

After substituting (1) and (2) into (3), the transfer function can be written:

$$\left|\frac{I_2}{I_0}\right|^2 = \frac{1}{\left[A - B\gamma^2 - \dfrac{C}{\gamma^2}\right]^2 + \left[D\gamma - \dfrac{E}{\gamma}\right]^2} \tag{4}$$

where,

$$A = 1 + G_1 R_2 + \frac{L_2}{L_1} + \frac{C_1}{C_2},$$

$$B = \omega_0^2 C_1 L_2,$$

$$C = \frac{1}{\omega_0^2 L_1 C_2},$$

$$D = \omega_0 R_2 C_1 + \omega_0 G_1 L_2,$$

$$E = \frac{R_2}{\omega_0 L_1} + \frac{G_1}{\omega_0 C_2},$$

$$\gamma = \frac{\omega}{\omega_0}.$$

The transfer function is maximally flat when

$$B = C, \tag{5}$$

$$D = E, \tag{6}$$

and

$$D^2 = 2B(A - 2B). \tag{7}$$

From (5) and (6)

$$\omega_0^2 C_1 L_1 = \omega_0^2 C_2 L_2 = 1, \tag{8}$$

so the primary and secondary circuits both resonate at frequency $f_o = \omega_0/2\pi$.

Substituting (8) into (4), the transfer function is conveniently expressed as

$$\left| \frac{I_2}{I_0} \right|^2 = \frac{1}{(A - 2B - BF^2)^2 + (DF)^2}, \tag{9}$$

where $F = (\gamma - 1/\gamma) = (\omega/\omega_0 - \omega_0/\omega)$. Equation (9) can also be written in the form:

$$\left| \frac{I_2}{I_0} \right|^2 = \frac{1}{\left( 1 + G_1 R_2 - \dfrac{C_1}{C_2} F^2 \right)^2 + (\omega_0 R_2 C_1 + \omega_0 G_1 L_2)^2 F^2}. \tag{10}$$

The midband loss is given when $F = 0$

$$\left| \frac{I_2}{I_0} \right|^2_{F=0} = \frac{1}{(1 + G_1 R_2)^2}. \tag{11}$$

The attenuation of the circuit can be controlled by the choice of $G_1$ and $R_2$. As shown in Fig. 1, the PIN diode provides variation in $G_1$. Assuming that the parameters have been chosen for a maximally flat transmission for some particular $G_1$ and $R_2$, and that only $G_1$ or $R_2$ will be changed, then (10) gives the amplitude response as a function of $G_1$ and $R_2$.

If (7) is also satisfied, then the current transfer ratio will be maximally flat and (9) becomes

$$\left| \frac{I_2}{I_0} \right|^2 = \frac{1}{(A - 2B)^2 + B^2 F^4}. \tag{12}$$

Equation (12) also can be expressed in the form:

$$\left| \frac{I_2}{I_0} \right|^2 = \frac{1}{[1 + G_1 R_2]^2 + \left( \dfrac{C_1}{C_2} \right)^2 F^4}. \tag{13}$$

The midband loss of the interstage is determined when $F = 0$:

$$\left| \frac{I_2}{I_0} \right|^2_{F=0} = \frac{1}{[1 + G_1 R_2]^2}.$$ (14)

The midband loss expressed in dB is

$$L = 20 \log \left| \frac{I_2}{I_0} \right| = 20 \log \frac{1}{1 + G_1 R_2} \text{ in dB.}$$ (15)

### 2.1 Design Equations and Procedure

The value of $G_1 R_2$ can be determined by writing (15) in the following form and using the maximum value of attenuation, which the designer must choose.

$$R_2 G_{1_{\max}} = \frac{1}{\log^{-1} \left( \frac{-L_{\max}}{20} \right)} - 1.$$ (16)

The required value of $R_2$ can be determined by choosing a maximum value of $G_1$ from the PIN diode characteristic. Such a characteristic is shown in Fig. 2 for a Western Electric L-2480 PIN diode.

The minimum attenuation is determined by using $R_2$ and the minimum value of $G_1$ . $G_{1_{\min}}$ is the parallel combination of the PIN diode at zero bias, the output impedance of transistor $Q1$, and diode bias resistors. From (15) the minimum midband attenuation if given by

$$L = 20 \log \frac{1}{1 + G_{1_{\min}} R_2}.$$ (17)

From (7), which must be satisfied to obtain a maximally flat transmission characteristic at the minimum loss point, we obtain

$$C_1 C_2 = \frac{1 \pm [1 - (R_2 G_{1_{\min}})^2]^{\frac{1}{2}}}{\omega_0^2 R_2^2},$$ (18)

where $\omega_0 = (\omega_1 \omega_2)^{\frac{1}{2}}$, and $\omega_1$ and $\omega_2$ are the band edges. From (13) the response is down 3 dB when

$$F_{3\,\text{dB}} = \left[ (1 + R_2 G_{1_{\min}}) \frac{C_2}{C_1} \right]^{\frac{1}{2}}.$$ (19)

Letting this value of $F_{3\text{dB}} = \beta$ where

$$\beta = \frac{\omega_2 - \omega_1}{(\omega_2 \omega_1)^{\frac{1}{2}}} = \frac{\omega_2 - \omega_1}{\omega_0},$$

Fig. 2 — Impedance of W. E. 2480 PIN diode as a function of frequency and diode current.

yields

$$\beta = \left[ (1 + R_2 G_{1\min}) \frac{C_2}{C_1} \right]^{\frac{1}{2}}. \tag{20}$$

From (8)

$$L_1 = \frac{1}{\omega_0^2 C_1} \tag{21}$$

and

$$L_2 = \frac{1}{\omega_0^2 C_2}. \tag{22}$$

## 2.2 PIN *Diode Characteristic*

The PIN diode used in this variolosser is a Western Electric L2480 diode.[2] The diode is essentially resistive over a wide frequency range and over a large range of diode currents as shown by the Smith chart plot of Fig. 2. This data shows the impedance of the diode at 200, 300, and 400 MHz as a function of the diode current. The data were obtained using the circuit shown schematically 'in Fig. 3a. The mechanical arrangement of the diode mounting with its bias circuit, as shown in Fig. 3b, consists of a cap, a clip, and a low inductance capacitor. This capacitor is chosen to resonate with the inductance of the diode and its mount at 300 MHz; the circuit has a $Q$ of 0.3 when the diode has a resistance of 10 ohms.

## 2.3 *Design Example of a Common Base, Common Base Variolosser*

A variolosser operating in the 200 to 400 MHz band was designed as outlined in Section 2.1. The bandwidth was chosen narrow compared with the maximum bandwidth obtainable so that the variation in the band shape as a function of attenuation could be observed. The conditions were that:

($i$) Attenuation of the interstage be 20 dB for a PIN diode current of 0.68 mA.
($ii$) $\beta = 0.725$.
($iii$) $f_o = 283$ MHz.

From the PIN diode characteristic of Fig. 2, the conductance at a diode current of 0.68 mA was found to be 0.0715 mhos. From (16), $R_2$ was computed to be 126 ohms. From the diode characteristic and interstage circuit, $G_{min}$ was found to be 1.33 millimhos. Using (17), the minimum loss was computed to be 1.35 dB. From (18) and (20),



Fig. 3 — L-2480 PIN diode mounting.

the values of $C_1$ and $C_2$ required to satisfy the bandwidth and maximally flat transmission characteristic at the minimum loss point are

$$C_1 = 9.38 \text{ pF}, \qquad C_2 = 4.22 \text{ pF}.$$

From (21) and (22), $L_1 = 34$ nH, $L_2 = 75$ nH.

The computed frequency response of a variolosser using the above design values and equation (10) is shown by the dotted curves in Fig. 4. These data show the loss relative to the loss at $\omega_o$ for various interstage losses. Figure 4a shows the frequency response for an interstage loss of 1.35 dB which is the minimum loss. Figures 4b and 4c show the response for a loss of 9.5 and 19.5 dB, respectively. These figures show the undercoupled response that is obtained as $G_1$ is varied. It also shows that by designing for a 3 dB bandwidth large compared with the desired bandwidth, that the amplitude distortion over a given attenuation range can be small. This design provides 19.5



Fig. 4 — Common base, common base variolosser characteristics.

dB attenuation over a 30 percent band with a variation in amplitude of ±0.15 dB. This bandwidth can be increased by providing a compensating circuit as discussed in Section III.

### III. ANALYSIS OF COMPENSATED CIRCUIT

Variolossers used to control the gain of an amplifier are usually operated near the maximum attenuation point and varied toward the minimum loss point. In this application it is desirable to have the frequency response flat near the maximum loss point. The addition of one inductor in series with the PIN diode allows the $Q$ of this circuit to be chosen so that the impedance of this arm varies with frequency in the correct way to compensate the under-coupled characteristic of the double tuned circuit. As the interstage attenuation is decreased, the $Q$ of the series resonant circuit is also decreased and thus the compensating circuit has little effect on the minimum loss point.

Figure 5a shows the interstage circuit with $L_3$ and $C_3$ in series with $R_3$ to form a series resonant circuit. This interstage has been simplified as shown in Fig. 5b. The current transfer function for the circuit in Fig. 5b is

$$\frac{I_2}{I_0} = \frac{Z_3}{(1 + y_1 Z_2)Z_3 + Z_2} \tag{23}$$



(a)



(b)

Fig. 5 — Variolosser with compensating circuit $Z_3$.

where

$$y_1 = G_1 + j\left(\omega C_1 - \frac{1}{\omega L_1}\right), \qquad G_1 = \frac{1}{R_{T1}}, \tag{24}$$

$$Z_2 = R_2 + j\left(\omega L_2 - \frac{1}{\omega C_2}\right), \tag{25}$$

$$Z_3 = R_3 + j\left(\omega L_3 - \frac{1}{\omega C_3}\right). \tag{26}$$

The magnitude of the transfer function is

$$\left|\frac{I_2}{I_0}\right|^2 = \frac{1+(Q_3 F)^2}{\left[1+G_1 R_2+\frac{R_2}{R_3}-\left(\frac{C_1}{C_2}+DQ_3\right)F^2\right]^2+\left[\left(1+G_1 R_2-\frac{C_1}{C_2}F^2\right)Q_3+\frac{R_2}{R_3}Q_2+D\right]^2 F^2} \tag{27}$$

where

$$G_1 = \frac{1}{R_{T1}}, \qquad Q_2 = \frac{1}{\omega_0 C_2 R_2}, \qquad Q_3 = \frac{1}{\omega_0 C_3 R_3},$$

$$D = \omega_0 G_1 L_2 + \omega_0 R_2 C_1 .$$

The frequency characteristic of (27) will have a flat amplitude response over a wide band at the maximum loss point if $Q_3$ is properly chosen. As the loss is adjusted from maximum to minimum, the value of $R_3$ must increase and thus $Q_3$ decreases. At the minimum loss point $Q_3$ is very small and the amplitude response is the maximally flat interstage characteristic that was used as the design point.

The common base, common base variolosser described in Section 2.3, whose characteristic is shown by the dotted curves in Fig. 4, was compensated at the 19.5 dB loss point by setting $Q_3$ equal to 0.832. The results are shown by the solid curves in Figs. 4b and 4c, and the dotted curve in 4a which applies to both the compensated and uncompensated case. Over the 40 percent band indicated by the arrows, the compensated characteristic deviates $\pm 0.15$ dB from a flat transmission characteristic for a variolosser that has a minimum loss of 1.35 dB and a change in attenuation of 18 dB. The interstage characteristic is flat at the maximum loss point.

IV. COMPUTED AND EXPERIMENTAL RESULTS

4.1 *Measured Performance of a Common Base, Common Base Variolosser*

A common base, common base variolosser was constructed using the design values computed in Section 2.3. The measured characteristic, without compensation, is shown by the series of pictures, Figs. 6a, b,

Fig. 6 — Frequency characteristic versus attenuation for 300 MHz common base-common base variolosser without compensation at attenuations of (a) 1.35 dB, (b) 9.5 dB, (c) 19.5, and with compensation at attenuations of (d) 1.35 dB, (e) 9.5 dB, and (f) 19.5 dB.

and c, for the minimum loss point, at an interstage loss of 9.5 dB, and at an interstage loss of 19.5 dB. The measured characteristic is also shown by the solid line in Fig. 4a at the minimum loss point. As described in Section III, the interstage was compensated to provide a flat transmission characteristic at the 19.5 dB loss point over the 229 to 352 MHz range. The effect of this compensation is shown by the series of pictures, Figs. 6d, e, and f. At 19.5 dB, the frequency characteristic is virtually flat over the 229 to 352 MHz range. At 9.5 and 1.35 dB, the compensating circuit has little effect on the frequency characteristic over the same frequency range. It was necessary to decrease the loss approximately 0.7 dB at 229 MHz to provide the desired frequency characteristics. The variolosser is shown schematically in Fig. 7.

## 4.2 *A Wideband Common Base, Common Base Variolosser*

A variolosser was designed for a large bandwidth and low amplitude distortion to demonstrate the wideband possibilities of the circuit. The essential difference between this design and the design presented in Section 2.3 is that the minimum loss is greater, the variation in attenuation is smaller, and the interstage has been designed to have



Fig. 7 — Common base, common base variolosser schematic.

a large 3 dB bandwidth. A compensating circuit with a $Q$ equal to 0.458 was used to provide a flat amplitude characteristic over a 50 percent band at the 17 dB loss point. The amplitude versus frequency characteristic was computed using equation (27); the results are shown by the series of curves in Fig. 8. This figure shows the amplitude characteristic at the 17 dB attenuation point, where the interstage was compensated, and the characteristic at 5 dB intervals down to the minimum loss of 2 dB. The arrows indicate a bandwidth slightly greater than 50 percent. Over this band the amplitude distortion does not exceed 0.1 dB for any attenuation level.

### 4.3 Common Emitter, Common Base Variolosser

A transistor operating in the common emitter configuration has an output admittance which is capacitive and resistive. Thus, transistor Q1 of Fig. 1 could be operated in the common base configuration, or in the common emitter configuration, without requiring a change in the basic circuit. The value of $C_1$ and the maximum value of $R_1$ depends upon the configuration. The equations in Section 2.1 can be



Fig. 8 — Wideband variolosser characteristics.

used equally well to design a common base, common base, or a common emitter, common base variolosser. The main difference between the two configurations is that the output conductance is considerably larger for common emitter operation, resulting in a higher minimum loss. Measurements at 300 MHz for the WE 2526 transistor show that the common emitter configuration has an output capacity of 1.7 pF shunted by 500 ohms. Using the combined shunt resistance for the PIN diode, a bias resistor of 1200 ohms, and equation (17), the minimum loss for a common emitter, common base variolosser is computed to be 2.7 dB. The maximum attenuation for $R_2 = 126$ ohms would remain at 20 dB.

### 4.4 *Performance of Common Emitter, Common Base Variolosser*

Figure 9 shows a variolosser operating between 240 and 360 MHz using the common emitter, common base configuration, designed to provide automatic gain control in a main IF amplifier.[3] This figure shows one section of the amplifier that contains the variolosser. By connecting point A of an identical section to point B of the one shown, gain-variolosser sections are cascaded.

The performance of three such sections cascaded to provide an attenuation range of 43 dB and operating from 240 to 360 MHz is shown in Fig. 10. The attenuation shown in the photos does not in-



Fig. 9 — Common emitter, common base variolosser schematic.

Fig. 10 — Frequency characteristic versus attenuation for three variolossers operating in 300 MHz IF amplifier at attenuations of (a) 39 dB, (b) 36 dB, (c) 26 dB, (d) 16 dB, (e) 6 dB, and (f) 0 dB.

clude the minimum loss of each section of approximately 2.7 dB. The falloff at the band edges is the result of several fixed gain stages following the variolossers. The frequency characteristic versus attenuation is shown when the three sections are providing attenuations of 39, 36, 26, 16, 6, and zero dB. There is an additional 4 dB range in attenuation which is used to correct for the variation in gain of the amplifier over the −40°F to +135°F temperature range. This series of photographs show the large range in attenuation and the small change in band shape that can be realized using the design method described in this paper.

V. CONCLUSIONS

A high quality, wideband variolosser can be made by coupling the common base, common base or the common emitter, common base configuration with a double tuned circuit, and varying the shunt impedance with a PIN diode. A minimum attenuation of 1.35 dB for the common base, common base configuration and 2.7 dB for the common emitter, common base configuration is realizable at 300 MHz. The minimum losses would be lower at lower frequencies. A maximum attenuation greater than 20 dB is possible for both configurations.

By designing the interstage circuit for a 3 dB bandwidth that is large compared with the desired bandwidth and compensating the circuit to provide a flat transmission characteristic at the maximum loss point, a variolosser with a loss variation of 15 dB can be made to operate over a 50 per cent bandwidth with little distortion in the frequency characteristic as a function of attenuation.

VI. ACKNOWLEDGMENTS

REFERENCES

1. Ruthroff, C. L., Osborne, T. L., and Bodtmann, W. F., "Short Hop Radio System Experiment," B.S.T.J., this issue, pp. 1577–1604.
2. Granna, N. G. and Forster, J. H., and Leenov, D., "PIN Diodes for Protective Limiter Applications," Solid-State Circuits Conf. Digest, February 1961, pp. 84–85.
3. Bodtmann, W. F., and Guilfoyle, F. E., "Broadband 300 MHz IF Amplifier Design," B.S.T.J., this issue, pp. 1665–1686.

# Microwave and Millimeter Wave Hybrid Integrated Circuits for Radio Systems

By M. V. SCHNEIDER, BERNARD GLANCE, and
W. F. BODTMANN

*Hybrid integration of microwave and millimeter wave circuits is essential for achieving future communication objectives in radio systems. Hybrid integrated circuits are circuits which are manufactured on a single planar substrate. Passive elements are fabricated by partial metallization of the substrate; active devices are inserted by bonding semiconductor diodes or bulk devices to the metal conductors.*

*We discuss the electrical properties of passive line elements on insulating substrates. We also compare the design formulas given with measurements made at 30 GHz, and present the results obtained at 30 GHz with wideband transitions from waveguide to microstrip and the measurements obtained with microstrip IMPATT oscillators and high order varactor multipliers in the same frequency range.*

*There are advantages of scaling for building hybrid integrated circuits which we discuss. Oversize models can be built and tested in a relatively short time and substantial savings in turnaround time, required manpower, and cost can be achieved.*

## I. INTRODUCTION

Advances in solid-state technology and related processing techniques in recent years make it possible to produce small size, minimum weight, low production cost circuits for communication systems. These advances became possible because of simultaneous improvements in the fields of planar diffusion, photolithographic pattern delineation, and vacuum and sputtering techniques. A typical end product of this process is the beam-leaded, sealed junction monolithic integrated circuit. Although the process is extremely complex, its main advantage is that many identical circuits, untouched by human hands, can be produced in a short time at low cost.

In spite of these advances only limited progress has been achieved in applying similar concepts and techniques to circuits at microwave and millimeter wave frequencies. There are two reasons for this. One is that the planar silicon technology is not fully applicable in the microwave and millimeter wave frequency range because many types of solid-state devices are manufactured from group III-V intermetallic compounds. The other reason is that stray reactances at high frequencies are extremely important, that is, metallic overlays on semiconductors have the electrical properties of short sections of low impedance transmission lines which can create special circuit problems.

We show that such problems can be solved and that hybrid integrated circuits for solid-state radio systems can be built with solid-state devices mounted on or bonded to suitable metallized dielectric substrates. Circuits of this type are not complex; because they can be produced economically, they are essential in radio systems where duplicate functons are required at many repeater locations such as for the radio pole line, for high-capacity domestic satellite systems, and for *Picturephone®* see-while-you-talk-service distribution.[1,2]

The substrate, the metallization process, and the ensuing electrical properties are important factors in the design of hybrid integrated microwave circuits. We demonstrate the usefulness of scaling, which serves as an analog computer in circuit design, with data obtained from microstrip transmission line measurements, band-pass filters built at 800 MHz and 30 GHz, and measurements obtained with oversize frequency multipliers. Scaling can also be used for designing impedance transformers, microstrip circulators, and frequency converters. Excitation of hybrid modes, radiation problems, and box resonances can be easily detected and corrected in an oversize model of the final circuit. Scaling also insures substantial savings in turnaround time, required manpower, and design cost of the RF section of a repeater.

## II. MICROSTRIP LINES AND SUBSTRATES

### 2.1 *Types of Microstrips*

Strip lines or microstrips provide the linear circuit functions in microwave integrated circuits. The lines are made by partial metallization of an insulating substrate which is supported by a metal ground plane or spaced between two ground planes with a suitable support. Two basic configurations which can be used are the standard microstrip shown in Fig. 1a and the triplate line in Fig. 1b. Both types

Fig. 1 — Microstrip transmission lines: (a) standard microstrip, (b) triplate microstrip embedded in dielectric, (c) inverted microstrip, (d) triplate microstrip with suspended dielectric, (e) shielded microstrip with suspended dielectric.

are used at frequencies up to a few GHz. It is often necessary to use modifications of both configurations at higher frequencies in order to obtain a low effective dielectric constant. This insures small total attenuation, improved relative mechanical tolerances, and little influence by nonuniformities of the substrate on the line parameters. The modified lines are the inverted microstrip shown in Fig. 1c and the triplate with a symmetrically suspended substrate in Fig. 1d which has been used successfully by Engelbrecht and Kurokawa for hybrid

integrated microwave transistor amplifiers.[3] The inverted microstrip and the triplate are special cases of the partially filled and completely shielded microstrip line in Fig. 1e where the major part of the field energy of a TEM mode is concentrated between one ground plane and the strip conductor, while the other ground plane and the side-walls are used for RF shields.

## 2.2 Solid-State Devices in Microstrips

Solid-state devices in microstrips can be shunt-mounted or series-mounted. Device terminals are connected between the ground plane and the strip conductor for shunt-mounting. For series-mounting, one interrupts the strip conductor and bonds the beam-leaded device terminals to the free ends of the strip conductor. Shunt-mounts are preferred in circuits which require a good heat sink. Series-mounts are suitable if the device terminal has to be connected to a high impedance circuit. Series and shunt mounts are equivalent in most other respects with the important exception that dc biasing requirements normally determine the solution which is more economical.

Diodes and bulk devices create different circuit problems in inverted and in triplate microstrips. Radiation losses occur in standard and inverted microstrips. They can be reduced by partial or full shielding as shown in Fig. 1e. The impedance of the shielded line does not change substantially because the major part of the RF energy is still concentrated between the bottom ground plane and the strip conductor.

Another problem occurs if a solid-state device is shunt-mounted in the triplate microstrip of Fig. 1d. If the device is inserted between the lower ground plane and the strip conductor it will be shunted by the capacitance between the strip conductor and the upper ground plane. A series reactance is also added because the RF energy is about equally divided between the lower and upper half of the triplate section which does not contain a device, and is not equally divided in the section which contains the diode or the bulk device. Symmetrical distribution of the RF energy can be obtained by mounting two identical solid-state devices, one in the lower and one in the upper half of the microstrip.

## 2.3 Computation of Line Parameters

The electrical parameters of microstrips are characterized by the impedance $Z$, the attenuation $\alpha$, the guide wavelength $\lambda$, and the unloaded $Q$. They can be calculated by numerical approximation as

TABLE I—MICROSTRIP IMPEDENCE

| Line parameter | Standard microstrip | Triplate microstrip |
|---|---|---|
| Ratio $\dfrac{w}{h}$ or $\dfrac{w}{b}$ | $\dfrac{w}{h} = \dfrac{2}{\pi} \dfrac{\partial \ln \vartheta_4(\zeta, \kappa)}{\partial \zeta}$ $\mathrm{dn}^2(2K\zeta) = E/K$ | $\dfrac{w}{b} = -\dfrac{\ln k}{\pi}$ $m = k^2$ |
| Impedance ohms | $Z_o = 60\pi\kappa$ $\kappa = K'/K$ | $Z_o = \dfrac{60\pi}{\kappa}$ $\kappa = K'/K$ |
| Narrow strip approximation | $w/h \ll 1$ $Z_o = 60 \ln \dfrac{8h}{w}$ ohm | $w/b \ll 1$ $Z_o = 60 \ln \dfrac{8b}{\pi w}$ ohm |
| Wide strip approximation | $w/h \gg 1$ $Z_o = \dfrac{240\pi}{2\dfrac{w}{h} + 5 - \dfrac{h}{w}}$ ohm | $w/b \gg 1$ $Z_o = \dfrac{30\pi}{\dfrac{w}{b} + \dfrac{\ln 4}{\pi}}$ ohm |

| | |
|---|---|
| $K(m)$, $K'(m)$, $E(m)$ | Complete elliptic integrals |
| $m = k^2$ | Modulus $m$ |
| $\mathrm{dn}^2(2K\zeta) = 1 - \mathrm{sn}^2(2K\zeta)$ | Jacobian elliptic function |
| $\vartheta_4(\zeta, \kappa)$ | Theta function |

shown by Brenner or by computation of Schwarz-Christoffel integrals if one assumes that the TEM mode is dominant.[4] The exact solution for the standard microstrip and the triplate microstrip of Fig. 1a and b with thickness $t = 0$ and $\epsilon_r = 1$ is given in Table I. The lines are defined by the ratios $w/h$ or $w/b$ which are functions of the complete elliptic integrals $K(m)$, $K'(m)$, $E(m)$ with modulus $m$ and the logarithmic derivative of the theta function $\vartheta_4$ defined by Tölke.[5]

The line parameters for $\epsilon_r \neq 1$ can be calculated from an effective dielectric constant $\epsilon_{eff}$ and from the line parameters for the free space microstrip $Z_o$, $\alpha_o$, and $\lambda_o$ as:

$$Z = \frac{Z_o}{(\epsilon_{eff})^{\frac{1}{2}}} \qquad \text{impedance} \qquad (1)$$

$$\lambda = \frac{\lambda_o}{(\epsilon_{eff})^{\frac{1}{2}}} \qquad \text{wavelength} \qquad (2)$$

$$\alpha = (\epsilon_{eff})^{\frac{1}{2}}\alpha_o \qquad \text{attenuation} \qquad (3)$$

$$Q = Q_o = \frac{20\pi}{\ln 10}\frac{1}{\alpha_o\lambda_o} \qquad \alpha_o\lambda_o \text{ in dB.} \qquad (4)$$

The free space impedance $Z_o$ has to be computed from the equations given in Table I. The following approximations can be derived for the standard microstrip of Fig. 1a with $(\mu_o/\epsilon_o)^{\frac{1}{2}} = 120\pi$ ohm

$$Z_o = 60 \ln\left(\frac{8h}{w} + \frac{w}{4h}\right) \text{ ohm} \qquad \frac{w}{h} \leqq 1 \qquad (5)$$

$$Z_o = \frac{120\pi \text{ ohm}}{\frac{w}{h} + 2.42 - 0.44\frac{h}{w} + \left(1 - \frac{h}{w}\right)^6} \qquad \frac{w}{h} \geqq 1. \qquad (6)$$

The accuracy obtained for strips with $w/h \leqq 10$ is $\pm0.25$ percent and for $w/h > 10$, the accuracy is $\pm1$ percent.

The effective dielectric constant $\epsilon_{eff}$ depends on the ratio $w/h$, the relative dielectric constant $\epsilon_r$, and the geometrical shape of the boundary between the air and the dielectric support material. An approximation for $\epsilon_{eff}$ can be derived with an accuracy of $\pm2$ percent from Wheeler's theory for the standard microstrip of Fig. 1a,[6]

$$\epsilon_{eff} = \frac{\epsilon_r + 1}{2} + \frac{\epsilon_r - 1}{2}\left(1 + \frac{10h}{w}\right)^{-\frac{1}{2}}. \qquad (7)$$

The conductor attenuation $\alpha_o$ for the standard microstrip, with the same resistivity for the ground plane and the strip conductor, can be computed from the skin resistance $R_s$ and the partial derivatives $\partial Z_o/\partial(w/h)$ and $\partial w/\partial t$. The skin resistance is

$$R_s = (\pi\mu_o f\rho)^{\frac{1}{2}} \text{ ohm} \qquad (8)$$

where $f$ is the frequency in Hz, $\rho$ the conductor resistivity in ohm·cm, and $\mu_o = 4\pi\cdot10^{-9}$ H/cm. The partial derivative $\partial w/\partial t$ is obtained from an equation $Z_o = Z_o(w, h, t)$ with $Z_o = $ constant. The attenuation $\alpha_o$ in dB per unit length for the standard microstrip is

$$\alpha_o = -\frac{R_s}{6\pi \ln 10}\frac{\partial Z_o}{\partial\left(\frac{w}{h}\right)}\frac{1 + \frac{w}{h} + \frac{\partial w}{\partial t}}{hZ_o} \qquad (9)$$

with the partial derivative $\partial w/\partial t$ given by

$$\frac{\partial w}{\partial t} = \frac{1}{\pi} \ln \frac{4\pi w}{t} \qquad \frac{w}{h} \leqq \frac{1}{2\pi} \qquad (10)$$

$$\frac{\partial w}{\partial t} = \frac{1}{\pi} \ln \frac{2h}{t} \qquad \frac{w}{h} \geqq \frac{1}{2\pi}. \qquad (11)$$

Equations (10) and (11) are valid for $t \ll h$, $t < w/2$, and $\partial w/\partial t > 1$. These conditions are fulfilled in most practical microstrip circuits. The attenuation obtained for narrow strips $w/h \leqq 1$ and for wide strips $w/h \geqq 1$ from equations (5), (6), and (9) in dB per unit length is

$$\alpha_o = \frac{10R_s}{\pi \ln 10} \frac{\left( \frac{8h}{w} - \frac{w}{4h} \right)\left( 1 + \frac{h}{w} + \frac{h}{w} \frac{\partial w}{\partial t} \right)}{hZ_o \exp\left( \frac{Z_o}{60} \right)} \qquad \frac{w}{h} \leqq 1 \qquad (12)$$

$$\alpha_o = \frac{R_s Z_o}{720\pi^2 h \ln 10} \left[ 1 + \frac{0.44h^2}{w^2} + \frac{6h^2}{w^2}\left( 1 - \frac{h}{w} \right)^5 \right]$$

$$\cdot \left( 1 + \frac{w}{h} + \frac{\partial w}{\partial t} \right) \qquad \frac{w}{h} \geqq 1. \qquad (13)$$

The unloaded $Q$ can be computed from equations (4), (12), and (13).

For wide strips $w/h \gg 1$ one can compute an approximate value for $\alpha_o$ based on the assumption of a uniform current density in the ground plane and on the opposite bottom face of the strip conductor. By assuming that the current is uniform over the width $w$ one obtains

$$\alpha_o = \frac{20}{\ln 10} \frac{R_s}{wZ_o} \qquad \frac{w}{h} \gg 1. \qquad (14)$$

One can show that the attenuation computed from the uniform current model is always higher than the attenuation obtained from the exact theory given by equation (9). Notice that the conductor resistivity $\rho$ at dc is always lower than the RF conductor resistivity. This explains the fact that good agreement of theory and measurement is often obtained by assuming uniform current density distribution and a skin resistance based on $\rho_{dc}$ because the two errors tend to compensate for each other. Section 2.5 discusses a typical example.

2.4 *Dielectric Loss*

The attenuation of a microstrip is increased if the dielectric support material is lossy. If the complex dielectric constant of the substrate is

$\epsilon_r = \epsilon'_r - j\epsilon''_r$ with $\epsilon''_r \ll \epsilon'_r$, one obtains, for the unloaded dielectric quality factor $Q_D$ of a microstrip which is fully embedded in the dielectric,

$$\frac{1}{Q_D} = \tan \delta = \frac{\epsilon''_r}{\epsilon'_r}. \tag{15}$$

The line properties for partial dielectric filling are characterized by an effective complex dielectric constant $\epsilon_{eff} = \epsilon'_{eff} - j\epsilon''_{eff}$. For $\epsilon''_{eff} \ll \epsilon'_{eff}$ we obtain by partial derivation

$$\frac{\epsilon''_{eff}}{\epsilon'_{eff}} = \frac{\partial \epsilon_{eff}}{\partial \epsilon_r} \frac{\epsilon''_r}{\epsilon'_{eff}}. \tag{16}$$

The unloaded quality factor for partial dielectric filling is $Q_D = \epsilon'_{eff}/\epsilon''_{eff}$ or from equations (15) and (16)

$$\frac{1}{Q_D} = (\tan \delta)_{eff} = \frac{\epsilon_r}{\epsilon_{eff}} \frac{\partial \epsilon_{eff}}{\partial \epsilon_r} \tan \delta. \tag{17}$$

Equation (17) can be simplified if the air dielectric interface is parallel to an electric field line. It is possible to define a filling factor $q$ for such a line as

$$q = \frac{\partial \epsilon_{eff}}{\partial \epsilon_r} = \frac{\epsilon_{eff} - 1}{\epsilon_r - 1}. \tag{18}$$

One can show that the effective dielectric constant given by equation (7) for the standard microstrip satisfies equation (18). This can be explained by the fact that the boundary between dielectric and air is parallel to the electric field at both corners of the strip conductor where field intensities reach their maximum. For the standard microstrip one obtains from equations (17) and (18)

$$\frac{1}{Q_D} = \frac{\dfrac{1}{\epsilon_{eff}} - 1}{\dfrac{1}{\epsilon_r} - 1} \tan \delta. \tag{19}$$

The total unloaded Q of the line is

$$\frac{1}{Q_T} = \frac{1}{Q} + \frac{1}{Q_D} \tag{20}$$

with $Q$ being the contribution from conductor loss and $Q_D$ the contribution from dielectric loss. The attenuation $\alpha_D$ in dB per unit length

resulting from the dielectric is

$$\alpha_D = \frac{20\pi}{\ln 10} \frac{\dfrac{1}{\epsilon_{eff}} - 1}{\dfrac{1}{\epsilon_r} - 1} \frac{\tan \delta}{\lambda} \qquad (21)$$

with $\lambda$ being the guide wavelength $\lambda_0/(\epsilon_{eff})^{\frac{1}{2}}$.

### 2.5 Measurements of Microstrips

The impedance, the guide wavelength, and the effective dielectric constant of microstrips can be measured with a time domain reflectometer, a capacitance bridge, and modified slotted line techniques. The attenuation is measured with a microstrip resonator or a low loss transition from microstrip to waveguide. The measurement of line parameters is necessary because conformal mapping is often too complex and also because the skin resistance $R_s$ is not exactly known because sufficient data for the conductor RF conductivity is lacking.

The effective dielectric constant $\epsilon_{eff}$ is an important line parameter as seen from equations (1) through (4). Computation of $\epsilon_{eff}$ is difficult for the inverted or suspended microstrip shown in Figs. 1c and d. It can be measured with a time domain reflectometer or a capacitance bridge. In the latter case, $\epsilon_{eff}$ is given by $\epsilon_{eff} = C/C_o$ with $C$ being the capacitance with the dielectric and $C_o$ the capacitance without the dielectric. Measurements of $(\epsilon_{eff})^{\frac{1}{2}}$ with a time domain reflectometer are shown in Fig. 2. Quartz substrates are used in the oversize microstrip transmission line. One concludes from the results of Fig. 2 that the effective dielectric constant of inverted microstrips is very low. This is also true for the shielded microstrip with suspended dielectric shown in Fig. 1e and the triplate microstrip with suspended dielectric shown in Fig. 1d.[4] The dielectric losses are small because the substrate material is located in the low field region of the microstrip.

The attenuation of a standard microstrip on a quartz substrate has been measured from 26.5 to 32 GHz with two low-loss transitions from waveguide to microstrip. The transitions are described in Section III. Line dimensions, measured and computed attenuation are given in Table II.

The calculated $Q$ for nonuniform current distribution is based on the measured dc resistivity of the composite conductor material and a partial derivative $\partial w/\partial t = 2.1$ obtained from equation (11). A

Fig. 2 — Square root of effective dielectric constant for inverted microstrip.

measurement of $\partial w/\partial t$ has been performed with an oversize micro-strip line with conductor dimensions shown in Fig. 2. The partial derivative is obtained by changing the thickness $t$ by $\Delta t$ and the width $w$ by $\Delta w$ in such a way that the characteristic impedance $Z_o$ remains constant. Satisfactory agreement has been obtained between calculated and measured derivatives if the conductor thickness exceeds several skin depths.

## 2.6 Substrate Material

The electrical properties of the insulating substrate determine the effective dielectric constant and the dielectric losses of the micro-strip. A small loss and a low $\epsilon_{eff}$ are obtained by choosing a configuration with a low RF field in the dielectric and by using a small $\epsilon_r$. Other important design parameters are the thermal expansion coefficient and the thermal conductivity. Table III lists the properties of some substrates which can be used for microwave and millimeter wave circuits.

Clear fused and polished quartz substrates have a small surface roughness and a low $\epsilon_r$ which is independent of frequency up to a few

TABLE II—MICROSTRIP ATTENUATION AT 30 GHz

| Microstrip type | Standard of Fig. 1a |
|---|---|
| Line dimensions | $w = h = 0.75$ mm, $t = 2$ $\mu$m |
| Clear fused polished quartz substrate* | $\epsilon_r = 3.78$ |
| Evaporated and photoetched conductor material | 150 Å nichrome (80% Ni, 20% Cr), 2 $\mu$m gold |
| Composite dc conductor resistivity | $\rho = 3.0 \cdot 10^{-6}$ ohm·cm |
| Measured attenuation 30 GHz, line length $l = 7.60$ cm | $\alpha l = 0.78$ dB |
| Measured unloaded $Q$ | $Q = 450$ |
| Calculated $Q$ for uniform current distribution, equation (14) | $Q = 514$ |
| Calculated $Q$ nonuniform curren t distribution, equation (9) | $Q = 840$ |

* 99.8 per cent $SiO_2$, Amersil Inc., Hillside, New Jersey 07205

hundred GHz. It also has a small thermal expansion coefficient, and it is more uniform than ceramics which have to be fabricated by firing metal oxides with a suitable binder.

The least expensive material is epoxy glass which consists of a glass fiber texture with an epoxy binder. It is available as a laminate with copper cladding. Special precautions are necessary because of some moisture absorption and also because of the relatively high loss tangent.

## 2.7 *Metal Deposition*

Conductor materials used in microwave integrated circuits should have the following properties

TABLE III—SUBSTRATE PROPERTIES
FREQUENCY $f = 30$ GHz, ROOM TEMPERATURE $T = 25°C$

| Substrate | Dielectric constant $\epsilon_r$ | Loss tangent $10^{4} \cdot \tan \delta$ | Thermal expansion coefficient $10^{6} \cdot \alpha$ °C$^{-1}$ | Thermal conductivity watts/cm °K |
|---|---|---|---|---|
| Clear fused quartz, $SiO_2$ | 3.78 | 1 | 0.4 | $1.4 \cdot 10^{-2}$ |
| Epoxy glass, Thiokol Panelite G10 | 4.4 | 80 | 10 | $1.6 \cdot 10^{-1}$ |
| Borosilicate glass, Corning 7040 | 4.5 | 73 | 4.8 | $1.1 \cdot 10^{-1}$ |
| Berillia BeO, Alsimag 754 | 6.0 | 40 | 6.0 | 2.3 |
| Alumina $Al_2O_3$, Alsimag 772 | 9.5 | 1 | 6.0 | $3.7 \cdot 10^{-1}$ |

(i) Low RF sheet resistance,

(ii) Good adherence and high stability,

(iii) Bondable top surface, and

(iv) Must be compatible with devices mounted into circuits.

The metal can be applied by vacuum deposition followed by electroplating. Adherence of the metal film to the substrate depends on the cleanliness of the substrate and choice of the base metal. Chromium, titanium, or nickel-chromium alloys are good base metals because of their high affinity for oxygen. A film several hundred Angstroms thick topped by an evaporated gold or copper film gives satisfactory results. The composite film must exceed several skin depths in thickness for low RF loss. The composite film should also be uniform over the total area of the circuit in order to obtain circuits which can be reproduced. Figure 3 shows a measurement of the film thickness of a nichrome-gold film over a total length of 7.5 cm. The thickness variation measured with a Tolansky-interferometer is less than ±5 percent. This result is obtained by using a three-section, three-strand tungsten coil with dimensions given in Fig. 3. The uniform film thickness leads to consistent results in the photolithographic process and in the final conductor delineation step.

III. WAVEGUIDE TO MICROSTRIP TRANSITION AT 30 GHZ

3.1 *Design of Microstrip Transitions*

The use of microstrip structures at frequencies near 30 GHz has generated a need for broadband waveguide to microstrip transitions.



Fig. 3 — Film thickness of evaporated gold film. Tungsten coil with three sections. Gold charge: 0.060 inch diameter wire, center piece 9/16 inch long, outer pieces ¾ inch long. Distance: 2.5 inches from substrate.

Most work with microstrips has been done at frequencies below 10 GHz where coaxial-microstrip transitions have proven satisfactory;[7,8] however, at higher frequencies test equipment is usually supplied with waveguide connectors, and waveguide-to-microstrip transistions are needed to determine the performance of microstrip components. The desirable features of such a transition are that it:

(*i*) Has a high return loss so that it does not appreciably effect loss or reflection measurements made on microstrip components.

(*ii*) Has a low transmission loss.

(*iii*) Can be connected to the microstrip with reproducable results and is easily connected or removed.

(*iv*) Be an in-line design for ease of connecting to test equipment.

(*v*) Is mechanically easy to reproduce.

The design approach has been to transform from the waveguide impedance to the microstrip impedance by use of a broadband stepped ridgeline transformer which is mechanically connected to the microstrip by a tab and a single pressure screw. The reactance introduced at the ridgeline-microstrip junction can be made sufficiently small so that the final transition satisfies the five desirable features just listed.

### 3.2 *Transformer*

A broadband four-step transformer is used to match the impedance of the RG-96/U waveguide to that of the microstrip. For a given bandwidth and return loss, the desired impedance at each step is computed by the design method outlined in Ref. 9. A stepped single ridgeline was chosen for this transformer since it is readily adapted to the unsymmetrical microstrip line; this type of structure is easily machined and attached to the waveguide. The mechanical dimensions of a ridgeline required to achieve the desired impedance at each step can be computed from the extensive data on ridgeline given in Ref. 10. By setting the height of the last step so that the substrate will hit and stop against the ridgeline, the microstrip ridgeline junction is made in a reproducible way. For a given waveguide size and impedance, choosing the height of the ridgeline also determines the width.[10] The ridgeline used in this transition is 0.096 inch wide.

### 3.3 *Performance of Transition*

The frequency band of interest is the band between 27.5 and 31.3 GHz. The return loss of a four-step ridgeline transformer designed for

this band, as shown in Fig. 4, curve 1, is substantially greater than 30 dB. The return loss of the complete transition is shown by curve 2 of Fig. 4. The reactance at the microstrip-ridgeline junction has been decreased by tapering the edge of the transformer at the junction. As curve 2 shows, the junction reactance decreased the return loss over the center of the band, improved the performance of the transformer at higher frequencies, and extended the useful bandwidth of the transition. This transition has more than 30 dB return loss over a 17 percent bandwidth and an insertion loss less than 0.1 dB.

### 3.4 *Mechanical Design*

The mechanical arrangement is shown by Fig. 5a. The waveguide and substrate are clamped to a common base for rigidity; a single insulated pressure screw (not shown) connects the 0.020-inch tab to the 0.030-inch-wide microstrip and applies pressure to the microstrip-waveguide ground planes. The edge of the ridgeline is used as a stop to position the substrate. This transition can be readily connected or removed without damaging the microstrip ground plane or line. Figure 5b shows the dimensions of the ridgeline transformer.

### 3.5 *Microstrip Line*

The microstrip conductor is deposited by evaporating a 2 $\mu$m thick gold film on a 0.030-inch clear fused and polished quartz substrate. A nichrome film 150 Å thick is deposited from a tungsten coil in vacuum prior to the gold evaporation in order to obtain good adherence. The strip is made by photoresist application with exposure through a contact photographic film, and by subsequent etching steps in a



Fig. 4 — Return loss of transformer and transition from waveguide to microstrip from 26.5 GHz to 32 GHz.

(a)



(b)

Fig. 5 — Dimensions of waveguide to microstrip transition: (a) mechanical design, (b) ridgeline transformer.

potassium iodine solution for the gold and ferric chloride for the nichrome. The impedance calculated from equations (1), (5), and (7) with $\epsilon_r = 3.78$ is $Z = 75.5$ ohm. Table II gives the results of an attenuation measurement performed with two waveguide-to-microstrip transitions for a 3-inch long microstrip line.

## IV. MICROSTRIP IMPATT OSCILLATORS

### 4.1 *Hybrid Integrated* IMPATT *Circuits*

IMPATT oscillators and bulk sources are required for frequency conversion in repeaters or terminals of solid-state RF systems. Major efforts in the past have been directed towards building sources with high efficiency in waveguide and coaxial circuits. These circuits give satisfactory performance, but present some problems in phase locked applications at higher frequencies or in circuits requiring solid-state devices for power amplifiers of angle-modulated signals. The widest possible locking bandwidth in such a circuit is obtained if the device

is connected to a short at the shortest possible distance from the diode required for achieving resonance with the effective device reactance. This can be achieved with a microstrip circuit in which one diode terminal is bonded to the ground plane and the other terminal is connected to the strip conductor which is shorted a small distance from the diode. Such a circuit resonates with a minimum of stored electromagnetic energy. It therefore has a low $Q$ and the maximum locking bandwidth for a given injected power.

## 4.2 Microstrip Circuit

Figure 6a shows the microstrip circuit built to obtain the maximum locking range with a given diode. The dimensions of a circuit built at X-band are given in Fig. 6b. The IMPATT oscillator at X-band consists of a silicon avalanche diode in a V package which is mounted in

QUARTER WAVE TRANSFORMER
COUPLING CAPACITOR
MICROSTRIP CAVITY
DIELECTRIC SHEET
(BOTH SIDES
METALLIZED)
IMPATT DIODE

HEAT SINK

(a)

COUPLING CAPACITOR        DIELECTRIC
IMPATT DIODE        |←--0.150"--→|    MICROSTRIP CAVITY
        COPPER CLAD
    0.100"   TEFLON SHEET
        HEAT SINK

(b)

Fig. 6 — Microstrip mount for X-band oscillator: (a) IMPATT oscillator (b) microstrip mount for X-band oscillator.

a short microstrip cavity. A similar circuit built at 30 GHz uses the same passive elements with all physical dimensions smaller. The circuit built at 30 GHz uses an unpackaged diode which is thermocompression-bonded to the heat sink. A loop or small cavity is obtained by mounting the diode in the center of a short microstrip line. The loop acts to a first approximation as a lumped inductance in parallel with the diode. It is observed that the frequency of oscillation multiplied by the square root of the loop length is constant, which means that the oscillation frequency can be changed by varying a single external circuit parameter. This important fact facilitates the circuit design.

Circuits with devices that have a negative impedance characteristic over part of the frequency spectrum present special problems in applying the dc bias. A simple solution has been found for applying the dc bias to devices mounted in microstrip circuits. A thin dielectric sheet with a hole for the diode is mounted on a solid metal block which is also the heat sink for the diode. The dielectric sheet is metallized on both sides and provides an RF bypass with two terminals for applying the dc bias.

The RF coupling between the circuit and the output is obtained with a quarter wavelength microstrip line in series with a coupling capacitor C as shown in Fig. 6. This solution gives dc isolation between circuit and load. The capacitance and the strip line impedance are adjusted for maximum output power. A transition from microstrip to coaxial transmission line is used for measuring the output power at X-band and a transition to waveguide is used for power measurements at 30 GHz.

### 4.3 RF Measurements and Locking Experiments

The circuit parameters and the RF performance obtained with an X-band and a V-band microstrip oscillator are given in Table IV. Locking measurements performed with a V-band oscillator operating at 31.6 GHz and 33.03 GHz are shown in Fig. 7. The diode used in this oscillator has a breakdown voltage of 17 V and a breakdown capacitance of 0.24 pF. The $Q$ of the circuit is defined by

$$\frac{\Delta f}{f_o} = \frac{2}{Q}\left(\frac{P_i}{P_o}\right)^{\frac{1}{2}} \tag{22}$$

where $\Delta f$ is the locking bandwidth, $f_o$ the frequency of the free-running oscillator, $P_o$ the output power without locking and $P_i$ the in-

TABLE IV — MICROSTRIP IMPATT OSCILLATOR PERFORMANCE

| Circuit and diode parameters, RF output | X-band oscillator | V-band oscillator |
|---|---|---|
| Microstrip cavity (inches) length, width, height | 0.150 0.120 0.100 | 0.070 0.070 0.020 |
| Silicon IMPATT diode mount | V-type pill package | Epoxy coated for mechanical support |
| Junction capacitance at breakdown (pF) | 0.44 | 0.22 |
| Breakdown voltage (V) | 80 | 24 |
| Inductance across diode (nH) | 0.32 | 0.16 |
| Coupling capacitance (pF) | 0.1 | 0.01 |
| Oscillator frequency (GHz) | 9.6 | 28.5 |
| RF output (mW) | 225 | 135 |
| Efficiency (%) | 2.3 | 2.2 |

jected locking power. The value of $Q$ is optimized by adjusting the impedance of the quarter wave transformer and the coupling capacitance. A circuit $Q$ between 3.5 and 5.0 is obtained for a gain from 17 to 30 dB.

### 4.4 Hybrid Integration With Other RF Circuits

The IMPATT oscillator described in this section is built with an external air line microstrip circuit. Hybrid integration of microwave oscillators with other circuits in a solid-state radio system requires a common substrate such as quartz or alumina on which all conductors are simultaneously deposited and then photoetched by photoresist and mask delineation processes. A circuit built with a common substrate consisting of a standard microstrip transmission line requires a hole in the dielectric for the diode. This type of mount is difficult and expensive to make. A better solution is to use the inverted microstrip of Fig. 1c or the triplate microstrip of Fig. 1d with a shunt-mounted diode which is bonded to the ground plane to obtain a good heat sink. The RF bypass shown in Fig. 8 is deposited on the ground plane and two metal studs are bonded to the metallized part of the RF bypass. Additional hybrid circuits may be deposited on the same microstrip line in order to provide other functions, such as downconversion, upconversion, or phase locked operation of oscillators.

### V. MICROSTRIP HIGH ORDER VARACTOR MULTIPLIERS

### 5.1 Hybrid Integrated Multiplier Circuits

Frequency multipliers in repeaters and terminals of solid-state radio systems are necessary for driving upconverters and downcon-

Fig. 7 — Locking range of V-band microstrip oscillator. The oscillation frequency is 31.6 GHz for $I_o = 90$ mA and 33.03 GHz for $I_o = 118$ mA.

verters or for phase-locking solid-state sources such as IMPATT and LSA oscillators. Major efforts in the past have been directed towards building doublers, triplers, and quadruplers with optimized efficiency at specified power levels in waveguide or coaxial circuits. Recently, work has been reported on hybrid integrated multipliers. A quadrupler from 2.25 to 9.0 GHz with two series-mounted silicon diodes on an alumina substrate has given a 4 dB conversion loss with an output power of 800 mW.[11] A hybrid integrated quadrupler from 15 to 60 GHz has been built with gallium arsenide Schottky barrier diodes;[12] a conversion loss of 18 dB has been measured. Quadruplers in the same frequency range have also been built with waveguide circuits. A gallium arsenide Schottky barrier quadrupler from 12.6 to 50.4 GHz, built in a waveguide circuit, has given a conversion loss of 17 to 18 dB. Notice, however, that the same waveguide quadrupler built with diffused junction gallium arsenide diodes has given a highly improved conversion loss of approximately 10 dB because of the beneficial effects of minority carrier charge-storage in the diffused junction.[13]

Fig. 8 — IMPATT oscillator cavity with output on inverted microstrip line.

Harmonic generators of order 4 or higher are used in RF systems where high order frequency multiplication from a crystal controlled power source from a few hundred MHz to 10 GHz and above is required. An example is the multiplier used in the short hop radio system which uses a chain consisting of only two stages, an octupler followed by a quadrupler. Lumped elements and waveguide components are used in the chain for all passive circuit elements. It is possible, with present integrated circuit technology, to build the passive elements for the idler circuit and input and output filters on one single insulating substrate. Results obtained with a hybrid integrated octupler from 3.8 to 30.4 GHz are reported in sections 5.2 and 5.3. Results are also given for multipler circuits built on an oversize substrate with an input frequency of 100 MHz and an output frequency of 800 MHz. By linear scaling, which means reducing all physical dimensions by the same factor, and by proper scaling of the diode package and the diode characteristics, one can obtain designs for all hybrid integrated octuplers by specifying the input or output frequency of the multiplier. This approach is also useful for designing other components and circuits of an RF repeater such as impedance transformers, directional couplers, filters, upconverters, and downconverters.

## 5.2 Microstrip Conductor Configuration for Idlers and Filters

Idler resonators, and low-pass, and high-pass filters for high order multipliers can be built on insulating substrates with transmission line elements shown in Fig. 1. Open-ended or shorted strip conductor stubs are used for the idlers, alternating sections of high impedance and low impedance microstrip lines for the low-pass filters, and parallel-coupled strip-line resonators for the band-pass filter. The design of these filters has been treated by Matthaei and others.[14] Figure 9 shows two micro-

strip conductor configurations which can be used for building a hybrid integrated octupler. It shows a top view of the microstrip conductor pattern which is deposited on an insulating substrate such as quartz or alumina. The varactor diode is shunt-mounted at point $D$ between the strip conductor and the ground plane for obtaining a good heat sink. Quarter wave stubs are used in Fig. 9a for providing the idlers at the second and fourth harmonic of the pump frequency $\omega_p$ .

The shortest stub of the $F$ structure is $\lambda/4$ long at the eighth harmonic and is transformed into a short at point A and an open circuit across the diode at $D$. This prevents power at the frequency $8\omega_p$ from flowing back into the pump circuit. The second stub, $\lambda/4$ long at the frequency $4\omega_p$ , is transformed into a short at the fourth harmonic at point B and into an open circuit across $A$. This prevents fourth harmonic power



(a)



(b)

Fig. 9 — Microstrip conductor configuration for high order varactor multipliers: (a) with waveguide output (b) with band-pass filter and microstrip output.

from flowing back into the pump circuit. The shortest stub of the $F$, the section of the microstrip between $A$ and $D$, and the shunt reactance of the output circuit, provide the idler circuit at the fourth harmonic. The remaining stub at $C$, $\lambda/4$ long at $2\omega_p$, has a similar function at the second harmonic.

The conductor pattern shown in Fig. 9b is made with a single stub resonator and a shunt-mounted diode at D. The stub length is designed so that idler currents are flowing at the second harmonic and at least one other higher order harmonic. Microminiature coils are thermocompression bonded between two low impedance microstrip sections for obtaining a two-section low-pass input filter. The two-section parallel-coupled bandpass filter has its center frequency at the eighth harmonic and is designed to reject at least 25 dB at the seventh and ninth harmonic. A transition from microstrip to waveguide for the circuits of Figs. 9a and b is used for measuring the RF performance of the octupler.

### 5.3 RF Performance of High Order Multipliers

The RF output obtained with a high order multiplier built with the microstrip conductor pattern shown in Fig. 9a is given in Table V. A coaxial-to-microstrip transition is used at the input of 3.8 GHz and a microstrip-to-RG-96/U waveguide transition is used for measuring the output power. The varactor is a planar diffused gallium arsenide diode which is shunt-mounted at point D between the conductor pattern and the ground plane. This is achieved by using the standard microstrip of Fig. 1a with an additional upper ground plane which is spaced at a distance $h$ from the insulating substrate. Tuning screws are located above points A, B, and C in the upper ground plane in order to maximize the idler currents. The diode terminals are connected between the microstrip pattern and the upper ground plane.

TABLE V — HYBRID INTEGRATED HIGH ORDER MULTIPLIER
FROM 3.8 TO 30.4 GHz

| | |
|---|---|
| Varactor diode | Planar diffused epitaxial GaAs diode |
| Junction diameter | 22 $\mu$m |
| Zero bias capacitance | 0.21 pF |
| Breakdown voltage | 18.5 V |
| Input power at 3.8 GHz coaxial transmission line | 200 mW |
| Output power measured in RG-96/U waveguide | 9.6 mW |
| Conversion loss | 13.2 dB |

This solution facilitates the assembly and avoids fabrication of a hole in the substrate which is normally done for shunt-mounting solid-state devices. The conductor pattern is deposited on a clear fused quartz subtrate with a thickness of 0.5 mm.

The waveguide provides sufficient rejection for all low order harmonics. The rejection at the seventh and ninth harmonic is lower and a band-pass filter in microstrip would be required in this circuit to suppress the unwanted high order harmonics.

A band-pass filter with a rejection of 28 dB at the seventh harmonic and a 31 dB rejection at the ninth harmonic has been built in an oversize microstrip transmission line and in a reduced microstrip line at 30 GHz. The measured 3 dB bandwidth for the oversize and for the reduced filter is ±2.4 percent. The filter consists of two parallel-coupled half-wavelength resonant microstrip lines on a standard quartz microstrip line. The insertion loss of the oversize model is 0.3 dB at 800 MHz; the measured insertion loss of the filter after linear reduction by a factor of 37.5 is 1.8 dB at 30 GHz. From this measurement one concludes that conductor losses are predominant because the insertion loss increases as the square root of the frequency because of the skin effect in the conductor material.

The RF performance of a completely shielded oversize high order multiplier from 100 to 800 MHz with this band-pass filter and with the conductor configuration of Fig. 9b is given in Table VI.

Linear scaling of the measured power levels to higher frequencies is not applicable because no attempt has been made to scale the RC-product of the varactor diode. The main purpose of the experiment is to obtain optimum conductor dimensions in order to build efficient integrated multipliers at much higher frequencies.

VI. CONCLUSION

Hybrid integrated microwave and millimeter wave circuits and components can be built on planar insulating quartz substrates. We

TABLE VI—HIGH ORDER OVERSIZE MULTIPLIER FROM
100 TO 800 MHz

|  |  |
| --- | --- |
| Varactor diode | Planar diffused epitaxial Si diode |
| Zero bias capacitance | 11.5 pF |
| Breakdown voltage | 80 volts |
| Input power at 100 MHz | 975 mW |
| Output power at 800 MHz | 600 mW |
| Conversion loss | 2.1 dB |

report the results on the attenuation of microstrip line elements, on transitions from waveguide to microstrip, on IMPATT oscillators, and on high order frequency multipliers. Linear scaling is used for obtaining optimum conductor configurations and for making essential measurements on passive microstrip components. A similar approach can be used for development of other hybrid integrated circuits required for RF systems at 10 GHz and above.

## VII. ACKNOWLEDGMENT

## REFERENCES

1. Ruthroff, C. L., Osborne, T. L., and Bodtmann, W. F., "Short Hop Radio System Experiment," B.S.T.J., this issue, pp. 1577–1604.
2. Tillotson, L. C., "A Model of a Domestic Satellite Communication System," B.S.T.J., *47*, No. 10 (December 1968), pp. 2111–2137.
3. Engelbrecht, R. S., and Kurokawa, K., "A Wideband Low Noise L-Band Balanced Transistor Amplifier," Proc. IEEE, *53*, No. 3 (March 1965), pp. 237–247.
4. Brenner, H. E., "Computer Design of Suspended-Substrate Integrated Circuits," Microwaves, *7*, No. 9 (September 1968), pp. 38–45.
5. Tölke, F., *Praktische Funktionenlehre, Zweiter Band, Theta-Funktionen und spezielle Weierstrass'sche Funktionen,* Berlin: Springer-Verlag, 1966, pp. 1–83.
6. Wheeler, H. A., "Transmission-Line Properties of Parallel Strips Separated by a Dielectric Sheet," IEEE Trans. Microwave Theory and Techniques, *MTT-13*, No. 2 (March 1965), pp. 172–185.
7. Grieg, D. D. and Engelmann, H. F., "Microstrip—A New Transmission Technique for the Kilomegacycle Range," Proc. IRE, *40*, No. 12 (December 1952), pp. 1644–1650.
8. Arditi, M., "Characteristics and Applications of Microstrip for Microwave Wiring," IEEE Trans. on Microwave Theory and Techniques, *MTT-3*, No. 2 (March 1955), pp. 31–56.
9. Cohn, S. B., "Optimum Design of Stepped Transmission-Line Transformers," IRE Trans. on Microwave Theory and Techniques, *MTT-3*, No. 3 (April 1955), pp. 16–21.
10. Hopfer, S., "The Design of Ridged Waveguides," IRE Trans. on Microwave Theory and Techniques, *MTT-3*, No. 5 (October 1955), pp. 20–29.
11. Johnson, K. M., "A High-Performance Integrated Microwave Circuit Frequency Quadrupler," IEEE Trans. on Microwave Theory and Techniques, *MTT-16*, No. 7 (July 1968), pp. 420–424.
12. Mao, S., Jones, S., and Vendelin, G. D., "Millimeter-Wave Integrated Circuits," IEEE Trans. on Microwave Theory and Techniques, *MTT-16*, No. 7 (July 1968), pp. 455–461.
13. Lee, T. P., and Burrus, C. A., "A Millimeter-Wave Quadrupler and an Up-Converter Using Planar-Diffused Gallium Arsenide Varactor Diodes," IEEE Trans. on Microwave Theory and Techniques, *MTT-16*, No. 5 (May 1968), pp. 287–296.
14. Matthaei, G. L., Young, L., and Jones, E. M. T., *Microwave Filters, Impedance Matching Networks, and Coupling Structures,* New York: Mc-Graw-Hill, 1964, pp. 355–477.

# Interference in a Dense Radio Network

By CLYDE L. RUTHROFF and LeROY C. TILLOTSON

*Radio systems operating at frequencies above 10 GHz are limited to short hops by rain attenuation. A severe interference problem arises from the fact that in a given area many repeaters may mutually interfere through through the back and side lobe responses of their antennas. Introducing the concept of frequency spectrum conservation by maximizing the communication flow through an area, a model of a dense network of radio systems has been studied to determine the effect of antenna discrimination and bandwidth expanding modulation methods upon the total communication capacity in a given frequency band. We conclude that for efficient performance, bandwidth expansion is required; we also conclude that communication capacity can be considerably increased by improving the near side lobe performance of antennas.*

## I. INTRODUCTION

The major factor influencing the design of radio systems above 10 GHz is the attenuation caused by rainfall; the principal result, from which many other prameters are determined, is that the repeaters must be closely spaced compared with lower frequency systems.[1,2] Consequently, the number of exposures to co-channel interference may be large—requiring interference resistant modulation. Such solutions generally use considerable bandwidth, thus raising questions about the efficiency of frequency occupancy.

The uses for short hop systems are such that many systems may exist in the same area on the same frequency assignment. Among the possible applications are broadband exchange area service, *Picturephone®* visual telephone distribution, and communication between mobile telephone concentrators. The important property with respect to efficient spectrum usage is the total "communication capacity through the area" in a given bandwidth. It is conceivable that, in terms of the concept of capacity through an area, interference resistant wideband modulation is more efficient than the narrowband

methods now used. The purpose of this work is to explore and answer this question. The approach is to set up a two-dimensional system model, compute the total interference at each repeater site, and explore the consequences of the RF signal-to-interference ratio (S/I) on the RF bandwidth required to meet the baseband interference criteria.

## II. INTERFERENCE PARAMETERS

The three most important system parameters affecting interference are antenna discrimination, the modulation method used to transmit information, and the frequency assignment plan.

### 2.1 *Antenna Discrimination*

Increased antenna discrimination reduces interference, but a finite limit on antenna size and performance is set by practical considerations such as antenna cost, stability of the supporting structure, and the repeater antenna environment. The short wavelengths resulting from the very high radio frequency greatly reduce the overall size and bulk of an antenna having a given gain and beamwidth, and thus make feasible new, more precise construction techniques. This is one of the few areas where the problems become easier as frequency increases. Hence in the calculation of interference it is admissible to assume well-shielded antennas with very small back and far sidelobe radiations and an aperture illumination which is tapered to control off-axis radiation near the main beam; this requires increased sophistication in antenna design. Some loss of antenna gain and a slight increase in repeater cost results; the increase in efficiency with which the frequency spectrum and the area of the countryside are used are well worthwhile.

Since there are practical but no natural limits on antenna beamwidth, the performance criteria selected for use in an interference study is, in part, a matter of judgment. One possibility is to demand only what has been obtained with antennas currently in production. This appears much too restrictive since it makes no allowance for new designs, changed requirements, and new construction techniques.

In the past antenna gain has been a dominant design parameter. For short hop systems operating in an environment of severe interference this is no longer justified; the angular response is much more important. For maximum communication capacity through an area, systems should be designed to be interference limited under conditions

of normal propagation; thermal noise is controlling only during heavy rain. Hence, it is reasonable to expect new antennas having better off-axis discrimination; the question is how much better. In this paper three antenna patterns are used, ranging from the theoretical limit, which can be approached but not achieved, to the patterns of existing antennas.

In addition to the antenna pattern, the environment in which the antenna is placed is also of importance. Reflections from buildings, trees, rocks, and so forth, can distort an otherwise ideal antenna pattern and considerably reduce the discrimination obtained in an ideal environment. This problem is well recognized in that most antenna testing ranges are chosen to be free of such obstacles. However, it is unlikely that many repeater locations for short hop systems of the type considered here will be so ideal. For these reasons there is a practical upper limit to the amount of discrimination which can be obtained even with an "ideal" antenna. In our calculations we assume the antenna response falls off in accordance with its diffraction pattern until it reaches the level limited by the environment and levels off at this value.

## 2.2 *Modulation Methods*

A second important aspect of system design which affects our ability to achieve maximum use of a frequency assignment in the area in which the systems are located is the modulation method used. Signals passing through a radio system are subjected to degradations caused by a variety of sources of noise and interference from the same and other systems. If maximum traffic is to be carried, the system must be resistant to interference so that systems can be densely packed. Part of this resistance can be obtained by using well-designed narrow beam antennas.

Additional resistance to interference can be obtained by using modulation methods such as PCM and large index phase or frequency modulation which exchange bandwidth for improved signal-to-interference ratio. Increased resistance to interference will allow more systems to be built in a given area but will require more bandwidth per channel so that an optimum is eventually reached beyond which total capacity begins to decrease.

Communication capacity is an involved function of many parameters including some related to the cost of providing a given level of sophistication in apparatus performance. These parameters change

with time as technology evolves and our understanding improves, so we do not attempt to calculate a grand maximum. Instead we study a somewhat arbitrary but sophisticated network which we believe demonstrates that, by properly exchanging bandwidth for interference resistance, the maximum communication flow through an area can be achieved.

In this paper only large index phase modulation is studied, although both digital and analog methods of expanding bandwidth are of interest.[1] The relations between bandwidth expansion and interference must be known for this work. For large index phase (or frequency) moduation, these relations are well known.[3,4] However, multilevel phase-shift keyed modulation, the digital method of interest, has not been so completely analyzed. The deterioration in performance produced by one or more co-channel interferers can be evaluated, and while some excellent work has been done on the spectral properties of phase-shift keying, we do not yet have sufficient information to determine the bandwidth required for multilevel phase-shift keying systems with realizable baseband pulses.[5-7]

In some circumstances it is possible to use each frequency assignment twice by use of orthogonal polarizations. The efficacy of this approach to doubling the capacity of the network depends upon the cross-polarization discrimination of the antenna, the tower stability, and the bandwidth expansion factor used in the method of modulation. In this paper the frequency assignment is used only once.

### 2.3 Frequency Plan

A two-frequency plan is assumed, at each site one frequency is transmitted and the other received; two-way transmission is accomplished. This is believed to be the most efficient plan and is the one used for most microwave systems. A single frequency plan is not feasible because the maximum repeater gain exceeds the coupling between the transmitting and receiving antennas. Use of the same frequency in the receiver and transmitter could result in an unstable repeater.

### III. ANALYSIS

### 3.1 The Network Model

There is no completely satisfactory way of choosing a network of systems on which to base an interference study. Regular geometric net-

works are artificial and unrealistic. A closer approach to field conditions could be obtained by choosing a "typical" area of the country and by using topographical maps, laying out several interesting systems, and then studying the resulting interference. The trouble comes in choosing a typical area; for every situation, favorable or not, that one chooses someone else can find a counter example. There is no "typical" area. A compromise approach suggested by R. Kompfner is to assume that possible repeater sites occur at random locations and that on the average, investigation of two sites is required to locate one repeater. One result of this approach is shown in Fig. 1.

The points shown on Fig. 1 were located by using two successive random 5-digit entries from the "Table of a Million Random Digits," as X, Y coordinates of repeater locations.[8] This approach has the advantage that sites are located randomly but their locations are precisely known so that interference can be computed. Figure 1 has a total of 11 north-south systems and 11 east-west systems, using 160 sites from 250 possible locations. Each of these systems is two-way; a two-frequency plan is assumed.



Fig. 1 — Network of radio systems with randomly located sites.

### 3.2. *Signal-to-Interference Ratios*

Let the antenna have gain $G(\theta)$ where $\theta$ is the angle in azimuth measured from the beam axis. The transmission between any two antennas can be written[9]

$$P_R = P_T\left(\frac{\lambda}{4\pi r_{TR}}\right)^2 G(\theta_{TR})G(\theta_{RT}) \tag{1}$$

where

$P_T$ is the power transmitted from one antenna,
$P_R$ is the power received in the other,
$\lambda$ is the free space wavelength,
$r_{TR}$ is the distance between antennas $T$ and $R$, and
$\theta_{TR}$ is the angle between the axis of the beam of antenna $T$ and the
   direction of antenna $R$.

In an interference computation, the quantity of interest is the ratio of the received signal power to the total interference power in each receiver. From (1) the received signal power in the $i$th receiver is

$$S_i = P_T\left(\frac{\lambda}{4\pi r_{TR}}\right)^2 G(0)^2, \tag{2}$$

where $r_{\mathrm{TR}}$ is the distance between transmitter and receiver. The total interference power in the $i$th receiver is given by

$$I_i = P_T\left(\frac{\lambda}{4\pi}\right)^2 \sum_j \frac{G(\theta_{ij})G(\theta_{ji})}{r_{ij}^2}, \tag{3}$$

and the summation is over all sources of interference. The ratio of interference-to-signal power in the $i$th receiver is therefore

$$\frac{I_i}{S_i} = \sum_j \frac{G(\theta_{ij})G(\theta_{ji})}{G(0)^2} \left(\frac{r_{TR}}{r_{ij}}\right)^2. \tag{4}$$

As expected, the signal-to-interference ratio is a function of relative antenna response and the ratios of antenna separations. The signal-to-interference ratios discussed in this paper were computed for the network of Fig. 1 from expression (4), and in accordance with the following comments.

(*i*) Each receiver receives interference from the adjacent sites in the same system because two signals are transmitted from each adjacent site on the same frequency; one is the desired signal, the other is interference.

(*ii*) There is no interference in the *i*th receiver from the $i \pm 2$ sites because these sites do not transmit on the interfering frequency.

(*iii*) When systems cross at the same site the distance is $r_{ij} = 0$ and expression (4) cannot be used. In this instance we assume separate antennas with coupling limited by the environment.

(*iv*) Only half of the repeaters transmit on the same frequency. If it is known which repeaters transmit on a given frequency, then only those sites need be considered in the computation. By the proper choice of frequency assignment within each system, the interference can be minimized. Such a procedure requires extensive and, perhaps, unmanageable frequency coordination. We assume that no frequency coordination is required and, for computational purposes, all repeaters transmit on both frequencies except as just noted. The resulting interference is therefore somewhat greater than would otherwise be the case.

The three antenna patterns used are shown in Fig. 2. Curves A, B,



Fig. 2 — Antenna patterns used in the computations.

and C represent the envelopes of the responses. The patterns are for a beamwidth of two degrees and an environment limitation of −60 dB as discussed in Section 2.1. These responses were used for all computations; the beamwidths and environmental limitations are parameters.

Antenna A is the limiting case described by Dolph and approximated by Taylor,[10,11] Antenna C is an approximation to the type of antenna described by Crawford and Turrin;[12] the envelope of the measured response is approximated by the function

$$\frac{G(\theta)}{G(0)} = \frac{1}{1 + \left(\dfrac{\theta}{\theta_0}\right)^4} + \frac{10^{-3}}{1 + \left(\dfrac{\theta}{25\theta_0}\right)^4}.$$ (5)

As shown in Fig. 2 this function levels off at the assumed environment limit. Antenna B was chosen as a compromise between the theoretical limiting antenna A and the existing antenna C. Its response is given by

$$\frac{G(\theta)}{G(0)} = \frac{1}{1 + \left(\dfrac{\theta}{\theta_0}\right)^4}.$$ (6)

## 3.3 *Bandwidth Expansion and Interference Resistance*

As mentioned in Section 2.1, resistance to interference can be improved by suitably increasing RF bandwidth. At the present time there is enough knowledge about the bandwidth requirements of analog angle-modulated systems to determine the necessary relations between interference and bandwidth expansion. For this reason, large index analog phase modulation is discussed in this paper.

Prabhu and Enloe have derived the relationships between the signal-to-interference ratio at a receiver input and the corresponding signal-to-interference ratio in the baseband output.[3] In their work the modulating baseband signal was a flat band of gaussian noise extending from 0 to $W$ Hz. For large modulation index and for one co-channel interferer, they have shown that the worst baseband interference occurs at $f = 0$, and that this minimum signal-to-interference ratio $S_o/I_o$ is given by

$$\frac{S_o}{I_o} \approx 2\varphi^3(S/I), \qquad \varphi \geqq 2.5,$$ (7)

where

$\varphi$ is the rms phase index in radians, and
$S/I$ is the signal-to-interference ratio at the receiver input.

Ref. 13 shows that when there is more than one co-channel interferer, the problem of estimating the baseband interchannel interference can be reduced to that resulting from a single equivalent co-channel interferer. The power in this one interferer can be expressed in terms of the powers of the individual interferers and the total number of interferers.

The RF bandwidth required for a large-index, phase-modulated baseband signal with uniform spectral density in the range 0 to $W$ Hz, and with a gaussian amplitude distribution, is given by Carson's rule. (The efficacy of Carson's rule is demonstrated in Ref. 4.)

$$B = 2W\left[1 + \frac{4\varphi}{\sqrt{3}}\right]. \tag{8}$$

The relationship between the bandwidth expansion factor and the RF signal-to-interference ratio is found by eliminating $\varphi$ from (7) and (8).

$$B/W \approx 2 + \frac{8}{\sqrt{3}}\left[\frac{1}{2}\frac{(S_o/I_o)}{(S/I)}\right]^{\frac{1}{3}}, \tag{9}$$

where

$B/W$ is the bandwidth expansion factor,

$S_o/I_o$ is the baseband signal-to-interference ratio at the receiver output, and

$S/I$ is the RF signal-to-interference ratio at the receiver input.

Figure 3, which has been drawn from (9), illustrates the relationship between the bandwidth expansion factor and $S/I$ with the baseband signal-to-interference ratio as a criterion.

### 3.4 Communication Capacity

We assume that if any repeater in a system does not meet the requirement on $S_o/I_o$, the system is inoperative. For any $S/I$, and therefore for any $B/W$, the percentage of systems in operation is designated by $\epsilon$. The relative communication capacity of the network, $C_R$, is defined as

$$C_R = \frac{\epsilon W}{B}. \tag{10}$$

This definition satisfies our intuitive ideas concerning communication capacity; it is proportional to the baseband bandwidth $W$ and the percentage of operating systems and is inversely proportional to the

Fig. 3 — Bandwidth expansion as a function of signal-to-interference ratio for large index phase modulation.

RF bandwidth $B$ required for successful transmission. The computations performed on the network model were to determine $C_R$ as a function of antenna response and the bandwidth expansion factor.

## IV. COMPUTATIONS AND RESULTS

### 4.1 *Computations*

In accordance with the foregoing analysis the total $S/I$ at each of the 222 receivers was computed for the network of Fig. 1, for the antenna responses of Fig. 2, for six 3 dB beamwidths, and for three values of environmental limit. The results, presented in several ways in the following sections, illustrate the interference behavior of the network and demonstrate the effects of antenna response, beamwidth, and bandwidth expansion on the total capacity of the network.

4.2 *Maximum Interference*

Figures 4, 5, and 6 show the lowest total $S/I$ computed for the network for various bandwidths and environment limits and for the three antennas of Fig. 2. If all systems in the network are to operate successfully, the modulation method must meet the baseband interference criterion at the $S/I$ given in Figs. 4, 5, and 6.

For large beamwidths the results are insensitive to the environment limit and antenna response. Conversely, for small beamwidths, the results depend strongly on antenna response and the environment limit. For antenna A there is no reason to have a beamwidth of less than 0.02 radians, whereas, a smaller beamwidth is beneficial for antenna C. These relationships can be seen more clearly in Fig. 7 where the data for the three antenna patterns and an environment limit of −60 dB are plotted from Figs. 4, 5, and 6.

Figure 7 shows that narrowing the antenna beamwidth improves the network performance as expected. As the beamwidth becomes smaller, the antenna response pattern dominates the behavior; as the beamwidth decreases further the network performance becomes independent of antenna response and is determined by the environment limit. The best result occurs for the narrowest beamwidth for all



Fig. 4 — Minimum signal-to-interference ratio as a function of antenna beamwidth and environment.

Fig. 5 — Minimum signal-to-interference ratio as a function of antenna beamwidth and environment.

antenna patterns. However, there are other constraints such as antenna size and tower stability which limit the minimum beamwidth of antennas for practical systems. Figure 7 also illustrates the important result that for a given $S/I$ criterion, the antenna with the best response pattern would allow a wider beamwidth to be used, thus



Fig. 6 — Minimum signal-to-interference ratio as a function of antenna beamwidth and environment.

Fig. 7 — Worst site signal-to-interference ratio for three antennas and an environment limit of —60 dB.

reducing the requirements on tower stability and permitting the use of a smaller antenna.

### 4.3 *Distributions of the Number of Sites With Respect to Interference*

The previous remarks concerned only the $S/I$ at the worst site in the network. This section illustrates the way in which the interference is distributed throughout the network. A beamwidth of 0.03 radians and an environment limit of −60 dB were chosen.

Figure 8 shows the number of receiver sites having a given $S/I$ in intervals of one dB. For example, 6 percent of the sites have $S/I$ in the interval 35 dB $\leq$ $S/I$ < 36 dB for antenna $C$. The same data are plotted in the form of distribution curves in Fig. 9.

It is probably safe to speculate on the basis of Figs. 8 and 9 that simple modifications in the network would result in substantial improvement in communication capacity for all antenna patterns.

### 4.4 *Communication Capacity of the Network*

For computation of the relative communication capacity, the distribution of the number of systems as a function of $S/I$ is required. As discussed in Section III and illustrated in Fig 1, there are 11 east-west and 11 north-south systems. Transmission is in both directions in each system. The distributions for the three antenna patterns are shown in Fig. 10 for an antenna beamwidth of 0.03 radians. In con-

Fig. 8 — Density of the number of sites as a function of signal-to-interference ratio for three antennas with 0.03 radian beamwidth and an environment limit of −60 dB.

trast to the distributions of sites, the distribution of systems has a clear and specific interpretation. If, for example, a modulation index is satisfactory for a $S/I = 39$ dB or greater, then Fig. 10 indicates all systems meet requirements with antenna A, 50 percent with antenna B, and none with antenna C. The quality, $\epsilon$, in expression (10) is found from Fig. 10.

The relative communication capacity is shown in Fig 11 for the three antenna patterns of Fig. 2 and a beamwidth of 0.03 radians. The points on these curves are computed from (10) with the aid of Figs. 3 and 10. The baseband requirement is $S_o/I_o = 60$ dB. (A baseband signal-to-interference ratio of 60 dB per hop would allow systems with about 100 hops to meet Bell System transmission objectives.)

It was mentioned in Section 3.3 that all computations were based on a single interferer. It has been shown that for a given interference power the single interferer results in the least baseband interference;[3] the worst case occurs when the interference power is divided equally among the maximum number of interferers. If the exact number and power distribution of the interferers were used in the computations the effect would be to shift the results with respect to the $S/I$ scales. We believe that none of the general conclusions are affected by the computing for a single interferer.

## V. CONCLUSIONS

For future applications of radio systems above 10 GHz it may be desirable to serve a large number of customers in a limited area. To make the best use of the available frequency spectrum, many systems in the same frequency band must co-exist in this area. Optimum de-



Fig. 9 — Distribution of the number of sites as a function of signal-to-interference ratio for three antennas with 0.03 radian beamwidth and an environment limit of −60 dB.

sign of these systems must include the concept of communication capacity through an area. This paper analyzes a dense radio network of this type and describes the effects of antenna response patterns and bandwidth expansion on communication capacity.

The limitations of the network model are clear: the earth isn't plane, sites are not chosen at random, and such a square network is highly artificial. We believe, however, that important conclusions are demonstrated by the results and that the directions for further research are plainly indicated.

(*i*) The importance of antennas with good off-axis discrimination is amply demonstrated. The potential improvements in communication capacity is large and offers incentive for better antenna design.

(*ii*) Modulation methods which expand bandwidth are not, in themselves, wasteful of frequency spectrum in the sense of total communication through an area. To the contrary, the results of computations in a simulated field environment indicate that bandwidth expansion is required to achieve efficient results.



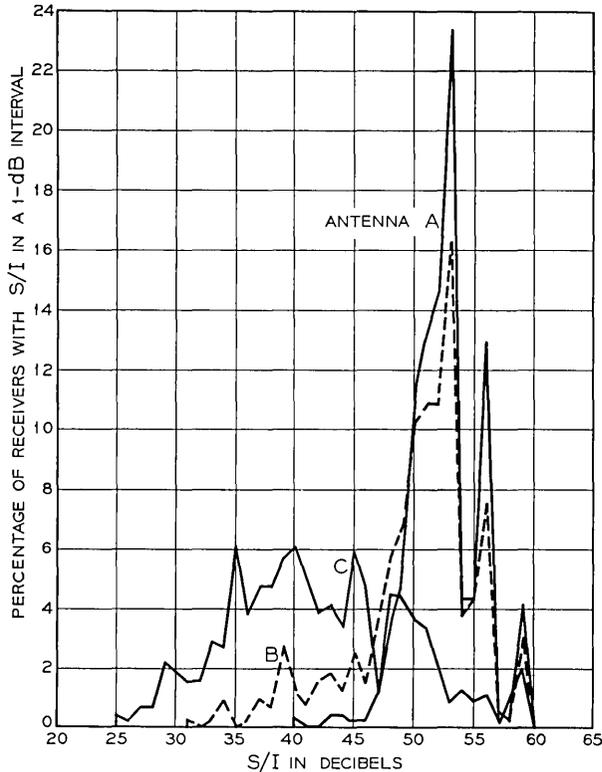Fig. 10 — Distribution of the number of systems as a function of signal-to-interference ratio for three antennas with 0.03 radian beamwidth and an environment limit of —60 dB.

Fig. 11 — Network efficiency as a function of bandwidth expansion ratio for three antennas with 0.03 radian beamwidths, an environment limit of —60 dB, and a baseband signal-to-interference criterion of 60 dB.

REFERENCES

1. Tillotson, L. C., Use of Frequencies Above 10 GHz for Common Carrier Applications," B.S.T.J., this issue, pp. 1563–1576.
2. Ruthroff, C. L., Osborne, T. L., and Bodtmann, W. F., "Short Hop Radio System Experiment," B.S.T.J., this issue, pp. 1577–1604.
3. Prabhu, V. K., and Enloe, L. H., "Interchannel Interference Considerations in Angle-Modulated Systems," scheduled for B.S.T.J., September 1969.
4. Ruthroff, C. L., "Computation of FM Distortion in Linear Networks for Bandlimited Periodic Signals," B.S.T.J., *47*, No. 6 (July–August 1968), pp. 1043–1063.
5. Prabhu, V. K., "Error Rate Considerations for Coherent Phase-Shift Keyed Systems With Co-Channel Interference," B.S.T.J., *48*, No. 3 (March 1969), pp. 743–767.
6. Rosenbaum, A. S., PSK Error Performance with Gaussian Noise and Interference B.S.T.J., *48*, No. 2 (February 1969), pp. 413–442.
7. Anderson, R. R. and Salz, J., "Spectra of Digital FM," B.S.T.J., *44*, No. 6 (July–August 1965), pp. 1165–1189.
8. The Rand Corporation, *A Million Random Digits With 100,000 Normal Deviates*, New York: The Free Press, 1966, pp. 231.
9. Schelkunoff, S. A., and Friis, H. T., *Antennas: Theory and Practice*, New York: John Wiley & Sons, Inc., 1952.
10. Dolph, C. L., "A Current Distribution for Broadside Arrays Which Optimizes the Relationship Between Beamwidth and Side-Lobe Level," Proc. IRE, *34*, No. 6 (June 1946), pp. 335–347.
11. Taylor, T. T., "Design of Line Source Antennas for Narrow Beamwidth and Low Side Lobes," IRE. Trans. on Antennas and Propagation, *AP-3*, No. 1 (January 1955), pp. 16–28.
12. Crawford, A. B., and Turrin, R. H., "A Packaged Antenna for Short-Hop Microwave Radio Systems," B.S.T.J., this issue, pp. 1605–1622.

# A Dense Network for Rapid Measurement
# of Rainfall Rate*

### By R. A. SEMPLAK and H. E. KELLER

(Manuscript received December 20, 1968)

*We discuss the design and operation of a dense rain gauge system for obtaining statistical data on both the temporal and spatial distribution of heavy rainfall. The data are collected every ten seconds from a network of approximately one hundred rain gauges spaced 1.34 kilometer in a square array which covers an area of 130 square kilometers. The rain gauge is a continuous, flow type with a response time of the order of one second. We describe the system used for recording the data on magnetic tape and give typical computer-generated rain maps for large area storms and for localized showers.*

## I. INTRODUCTION

Studies of the effects of rain on propagation at centimeter and millimeter wavelength have shown that attenuation by rain becomes more severe as the wavelength decreases.[1-5] The attenuation becomes objectionable at a wavelength of about 5 cm and increases rapidly. Information on the temporal and spatial characteristics of a rainfall at the surface of the earth is scant and further knowledge is needed to permit design of systems for the radio frequency bands above 10 GHz.

An experiment has been designed to obtain statistical data on both the temporal and spatial distribution of heavy rainfall in both time and space. The rainfall data are collected from a network consisting of approximately one hundred gauges about 1.34 km from each other and covering an area of 130 square km, as shown in Fig. 1.[†] The dot within each grid indicates the physical location of a rain gauge mounted at the top of a telephone pole, approximately 7.6 m above

---

Fig. 1 — Location of rain gauges are shown by small dots in each grid. The heavy lines denote central office boundaries. The dashed line from grid 9 to grid 33 indicates the 18.5 GHz propagation path.

the ground. The heavy lines denote the boundaries of various New Jersey Bell Telephone Company central office areas. The locations of each central office is indicated by a heavy circle; the Red Bank central office serving the lower right hand area is located some distance off the map.* Also shown by the dashed line extending between grids 9 and 33, is the transmission path used in the 18.5 GHz propagation studies which are discussed in a companion paper.[4]

---

* To reduce the number of telephone lines which carry information from the rain gauges, the lines are concentrated in each central office and sampling is accomplished there. This is discussed in detail in Section II.

The continuously measuring rainfall rate gauge is discussed fully in Ref. 6; therefore we give only a brief description. The gauge uses the high dielectric constant of water (approximately 80 at low frequencies) to change the frequency of an oscillator. Rain is collected by a funnel in the same manner as that used for other types of gauges; however, the collected water is directed down an inclined plane and flows between the insulated electrodes of a capacitor. Since the capacitance is a function of the amount of water flowing between the electrodes, the frequency of the oscillator is a function of rainfall rate. This capacitor is the variable element in a simple audio frequency resistance-capacitance relaxation oscillator. After the oscillator there is a multivibrator which divides the oscillator frequency by two and provides a square wave output with a fifty percent duty cycle. An emitter follower is used to couple the output to the telephone line; thus a single telephone pair can be used both for supplying dc power and for carrying the signal to the central office. A rain gauge and typical calibration curve are shown in Fig. 2.

In a manner discussed in Section II, the entire network of one hundred gauges is scanned every ten seconds; that is, each gauge is sampled for one-tenth of a second every ten seconds. This choice of scan rate (one scan per ten seconds) was limited in part by economic considerations, but, based on experience with conventional rain gauges, it was also deemed rapid enough to provide sufficient temporal resolution of rain-rate. However, these experiments have shown that one scan per second is needed to resolve some of the intense showers.* Some discussion of the errors introduced by sampling only every ten seconds is given in Ref. 4.

## II. DATA ACQUISITION SYSTEM

The apparatus for data acquisition is shown in Fig. 3. The signal conditioning units, one for each line, supply dc power to the gauge from the central office 130 volt battery through a current adjustment and monitor section. The audio frequency signal from the gauge, is separated from the dc and equalized to within about $\pm$ 2-$\frac{1}{2}$ dB between one and twelve kHz. All circuits from gauges to counter input transformers are carried and switched on balanced pairs.

The multiplexing is done in two stages. In each of the four central offices, a multiplex samples the lines coming into that office and concentrates them on 15 kHz program quality line, complete with equal-

---

* The time constant of the capacitor-type gauge is about one second.

Fig. 2 — (a) Rain gauge removed from container; notice that the electronic unit is fastened to the lid for ease of inspection; (b) Typical calibration data for a rain gauge: the points represent the measured values, and the curve is empirical.

izers and amplifiers, to the master multiplex located at Crawford Hill. The master control synchronizes both the central office and master multiplexes; the latter transfers the incoming samples on the four program lines to the appropriate counter for conversion to digital data. Figure 4 shows the signal flow from a typical gauge to the recorder.

The basic part of the multiplexer is the matrix switch unit; it consists of ten vertical lines by five horizontal lines forming 50 crosspoints each of which consists of a diode in series with a two contact reed relay. Each central office multiple uses only part of its matrix switch since there are less than 50 gauges in any central office area. The master multiplex uses two matrices and each forms the base for a 50 channel multiplex. Thus two gauges can be sampled at the same time. The program for scanning the gauges (an interlaced pattern) is wired into both the central office and master multiplex.

A submultiplex consisting of one horizontal line of ten crosspoints is also synchronized with the master matrix switch and recycles every ten steps; it is used to feed auxiliary data to counter 3. The auxiliary data consists of received level on the microwave propagation path, wind velocity at Crawford Hill, and detailed data from rain gauge 33.

The master control (Fig. 3) receives a pulse every 0.2 seconds from a one hundred kHz time base; this pulse is fed to counters with counts of ten, five, three, and two times. The count-of-ten counter furnishes the vertical information for all matrix switches through a binary coded decimal to ten-line converter. The vertical information is also converted into a "two of five" code and transmitted to each central office multiplex over three pairs where it is reconstructed to ten-line for the matrix vertical lines. The five times counter, through a binary coded decimal to ten-line converter, drives the horizontal matrix lines.



Fig. 3 — Data acquisition system.

Fig. 4 — Flow diagram showing the direct current feed to a gauge and the signal path from a gauge to the tape recorder.

The count-of-three counter is used to group three scans into a record.

There are two modes of operation, slow and fast scan. In the former, a timing pulse fed to the master control produces a scan every 100 seconds; this is the "rain-watch" mode. When significant rain is indicated on any two gauges, the system switches to fast-scan. Low pass filters with cutoff frequencies of 9.40 kHz* determine the point at which the fast scan begins. When rain is no longer indicated the system returns to slow-scan.

Three frequency counters, used as in Fig. 3, provide the timing pulse for scanner and recorder. In the master control a digital scanner combines, in sequence, the output of counter 3, information from a time clock, and day-of-the-year. The output of each counter is four digits of binary code decimal which is serialized and fed to the tape recorder.

## III. DATA PROCESSING

The adjunctive data discussed in Section II are recorded on the A and B tracks of the magnetic tape; both these and the gauge data are expressed as four digit numbers (16 bits). This means that all six of the tracks carry data and that an unpacking routine is necessary.

The raw data on the tape are in the form of frequencies (for example, data from the rain gauges, microwave receiver, and so on). The 200 bits per inch raw data tape is copied to 800 bits per inch. The latter is read into the computer with an unpacking program which sorts the intermixed bits and completes all partial numerical digits. The tape resulting from the unpacking procedure is used for all data reduction. The manner in which the data are analyzed is discussed in companion papers.[4,7]

## IV. OPERATIONAL DATA

The system has been operational since late 1966 and data for all rainfalls since then have been recorded,† but because of possible errors in both gauges and recording apparatus, only data taken since May

---

* For the average gauge this frequency corresponds to a rain rate of 15 mm per hour (see Fig. 2b).

† Since the rain gauges in the network are not equipped to handle snow (that is, heaters to melt snow) and in order to avoid any ambiguous data because of freezing rain, it was decided not to operate the system for outside temperatures less than 40°F. However, data obtained from both tipping bucket and weighing gauges indicate that with two exceptions, the rain rate during periods of low temperature was less than 25 mm per hour. The first exception occurred December 3, 1967 when a peak rate of 150 mm per hour was recorded for a period of less than two minutes. The second exception occurred March 12, 1968 with a peak rate of 75 mm per hour.

1967 are considered to be a high reliability; numerous heavy rains occurred during the summers of 1967 and 1968. Several examples of the output of the system are given below.

Figure 5 shows eight computer-generated rain intensity photographs. Each photograph shows the instantaneous rain rate as measured on each rain gauge in one scan. The sequence of photos, Fig. 5a thru 5d, occurred in 10 second intevals (the regular scan period). The dark grey areas indicate absence of rain whereas, the brightest areas correspond to 200 mm per hour rainfall. These photos portray a general, large area storm; it is evident that although it was raining over the entire network, cells of considerable intensity occur at various places;



(a)                                    (b)

(c)                                    (d)

(e)

(f)

(g)

(h)

Fig. 5 — Computer-generated rainfall-rate patterns. Small "patches" correspond to geographic area associated with each gauge. Each patch can have one of 48 intensities: dark for no rain, and gradually brightening for increasing rainfall. Frames (a) through (d) map a general storm over the area at ten-second intervals. The sequence (e) through (h) portrays a localized shower system moving diagonally across the system at 5 minute intervals.

thus "shower" activity is often embedded in a general rain environment. For reasons such as this, spaced-path diversity (discussed in Refs. 5 and 7) provides a considerable advantage in millimeter-wave radio-relay systems.

The sequence of photos, of Figs. 5e thru 5h, have been selected at 5 minute intervals; they show a localized storm, or shower, moving

from southwest to northeast through the network; this shower moved about ten miles in 20 minutes (30 mph). In the northwest quadrant of Figs. 5a thru 5d, the 6.4 km propagation path operating at 18.5 GHz (Ref. 4) is shown by a heavy dark line.* The isolated square in the upper right hand corner of the photos of this sequence represents the received 18.5 GHz signal. The magnitude of the signal is indicated by the brightness of this square. In Figs. 5e, f, and h, the square is fairly bright indicating little attenuation, but in Fig. 5g the square is dark, indicating an attenuation of about 30 dB; this coincides with heavy rain on three of the gauges on the path. Entire storms have been processed by computer in this manner to produce moving pictures of the storm activity.

In Fig. 6, data from the rain gauge network has been used to generate contour maps. For this presentation, a linear interpolation is used to produce the contours between adjacent gauges. Elapsed time between the upper and lower maps was 20 minutes. In Fig. 6a, the rain rates are very low, indicating a fairly light and steady rain over the network. However, in Fig. 6b, taken twenty minutes later than 6a, several intense cells have formed within a few miles of each other. The contour interval is 20 mm per hour in Fig. 6b; even with this large an interval, the contours become crowded in some places simply because the rain rate on a given gauge is so much higher than that of its nearest neighbor.

## V. MAINTENANCE, TESTING, AND CALIBRATION

A daily record is made of the network performance. Difficulties are caused by malfunction and noise on telephone lines, gauges damaged by inadvertent application of high voltage, clogged capacitors (mainly caused by small spiders), and component failures caused by lightning and undetermined reasons. To date, more than 90 percent of the network has always been operational.

After analysis of preliminary data, the rain gauges were calibrated in the field; this involved taking a source of water to each gauge position and measuring the output frequency of the gauge for three water-flow rates, namely 45, 75, and 254 mm per hour. Due to pollution in the gauge capacitor, some departures from the original calibrations were observed. However, when the gauges were cleaned, they faithfully returned to their original calibration. Thus the present pro-

---

* A 1.9 km path at 30 GHz was in operation during 1968. The author plans to discuss these data in a future paper.

Fig. 6 — (a) Plot of rainfall-rate contours in millimeters per hour showing several rain cells on the network; (b) contours 20 minutes later. Notice that there is only one gauge per square and linear interpolation was used in obtaining the detailed contours.

cedure is to clean the gauges bimonthly and to use the original calibrations in all data reduction.

## VI. ACKNOWLEDGMENTS

The efforts of the many people at Bell Telephone Laboratories and New Jersey Bell Telephone Company, especially H. A. Gorenflo and J. Batton, are greatly appreciated.

REFERENCES

1. Gunn, K. L. S., and East, T. W. R., "The Microwave Properties of Precipitation Particles," J. Royal Meteorological Soc., *80*, No. 1954, pp. 522–545.
2. Hathaway, S. D., and Evans, H. W., "Radio Attenuation at 11 Kmc and Some Implications Affecting Relay System Engineering," B.S.T.J., *38*, No. 1 (January 1959), pp. 73–98.
3. Medhurst, R. G., "Rainfall Attenuation of Centimeter Waves: Comparison of Theory and Measurement," IEEE Trans. on Antennas and Propagation, *AP-13*, No. 4 (July 1965), pp. 550–564.
4. Semplak, R. A., and Turrin, R. H., "Some Measurements of Attenuation by Rainfall at 18.5 GHz," B.S.T.J., this issue, pp. 1767–1787.
5. Hogg, D. C., "Millimeter-Wave Communication Through the Atmosphere," Science, *159*, No. 3810 (January 5, 1968), pp. 39–46.
6. Semplak, R. A., "Gauge for Continuously Measuring Rate of Rainfall," Rev. of Sci. Instruments, *37*, No. 11 (November 1966), pp. 1554–1558.
7. Freeny, A. E., "Statistical Treatment of Rain Gauge Calibration Data," B.S.T.J., this issue, pp. 1757–1766.

# Statistical Treatment of Rain Gauge Calibration Data

By MRS. A. E. FREENY

(Manuscript received December 2, 1968)

*This paper describes the statistical treatment of the calibration data of the capacitance gauges used for measurement of rain rates in the rain gauge network set up in a 160 square kilometer area surrounding Crawford Hill, Holmdel, New Jersey. The expression* $R = \alpha e^{-\beta f}$, *where* R *is the rain rate in millimeters per hour, f is the frequency of the rain gauge oscillator, and* $\alpha$ *and* $\beta$ *are fitted coefficients, allows calibration curves to be developed for all 99 of the devices in the system with an overall equivalent standard error of 7.3 millimeters per hour. This paper discusses the distribution of parameters and residuals and gives a refinement which corrects for a fitting bias.*

## I. INTRODUCTION

This paper describes the statistical treatment of the calibration data of capacitance gauges used for measurement of rain rate in the rain gauge network set up in a 160 square km area surrounding Crawford Hill, Holmdel, New Jersey. Section II describes the original calibration data plus added observations, Section III discusses the least squares fits to the individual oscillators and their results, and Section IV describes the correction made for bias.

## II. CALIBRATION DATA

The calibration data consisted of approximately 20 pairs of readings for each of the 100 gauges.[1,2] These readings, which consist of frequency in KHz and rain rate in mm per hour, were obtained by pouring water via a flowmeter through the gauge at a given rate and recording the corresponding oscillator frequency reading. The data used includes the original calibration readings plus supplementary readings made in the spring of 1967 when a field check on the original

Fig. 1 — Rain rate versus frequency, oscillator 47.

calibrations was made. For oscillator 47, which shows a typical configuration, rain rate $R$ is plotted against frequency $f$ on both a linear and a semilog scale in Fig. 1. Points plotted with a "•" are the original calibration measurements; those with a "o" are added measurements.

## III. EXPONENTIAL FITS TO INDIVIDUAL OSCILLATORS

The linear character of the semilog plots show that $R = \alpha e^{-\beta f}$, that is, a simple exponential curve, is a reasonable description of the data.

After some consideration the exponential form, rather than the linear alternative $\ln R = \ln \alpha - \beta f$, was chosen so that the residuals would be equally weighted over the frequency range. This minimizes the relative errors at high rain rates. The data for each of the oscillators were put through a nonlinear least squares fitting program; estimates $a$ and $b$ of the parameters $\alpha$ and $\beta$ were obtained. Figure 2a shows the fit of the data to oscillator 47.

For some of the oscillators, it was clear from the additional readings that the calibration at this time was not the same as it had been at the time the original measurements were made. For these oscillators, the original data were not used and fitting was done on the supplementary data only. Two of these cases were left with only the three additional points, and one had only five points, so the quality of the estimates for these oscillators is poor. In all, six oscillators had less than ten points.

Inspection of the residuals from the fits for the individual oscillators showed very good agreement between the calculated values and the observed values in most cases; however, a significant number of the oscillators showed two types of nonrandom scatter about the fitted curve.

The first type is a pronounced lack of fit of the exponential curve to low rain rate data ($R \leq 30$ mm per hour) for oscillators 13, 21, and 35.



Fig. 2 — Exponential fits: (a) oscillator 47 $[R = 883.84 \exp(-0.516f)]$; (b) oscillator 76 $[R = 849.86 \exp(-0.534f)]$.

Fig. 3 — Normal probability plots of estimates from exponential fits to each oscillator: (a) alpha, (b) beta.

The fitted curve overestimates the rain rate by as much as 20 mm per hour at very low rates; but the overestimate disappears by about 30 mm per hour.

The second type of scatter is a tendency of the observed rain rate to lie below the curve for $4 \lesssim f \lesssim 5$ KHz, and above the curve for $5 \lesssim f \lesssim 9$ KHz for some oscillators. Figure 2b shows the fit of the curve to the data from oscillator 76 which follows this pattern. Several other straightforward functions were tested on a few of these oscillators but did not result in improvement of the residual patterns. However, the simple exponential equation expresses the relationship between rain rate and frequency quite well.

Of the 100 oscillators, 89 had an error mean square of less than 100, and 11 had an error mean square greater than 100. Inspection of the individual fits to the oscillators indicated that there was an outlier in the data of the oscillators having two of the 11 worst fits. These observations were removed and new fits were calculated. These have error mean squares less than 100. Of the remaining nine worst fits, one was to oscillator 56 which was not being used in the rain gauge network, so the data from this oscillator were excluded from further calculations. After removal of this data, the average value of the error mean square is 57.75. The mean square of the pooled residuals is 53.55, which is equivalent to a standard error of about 7.3.

An informal assessment of the distribution of $\alpha$ and $\beta$ may be made from normal probability plots of $a$ and $b$ which show the values of $a$ and $b$ plotted against standard normal quantiles for a sample of size 99.[3] These plots are shown in Fig. 3. The plot of $a$ shows an upper tail which is too long to be normally distributed. The plot of $b$ has both tails too short for a normal distribution. It is also possible to interpret these plots as showing a slight noticeable curvature upward for $a$ and downward for $b$. This suggests that the estimates from the equation in the form $\ln R = A - \beta f$ would have been more nearly normally distributed.

Figure 4 is a normal plot of the pooled residuals for all oscillators. The outstanding feature of this plot is a change in the slope toward the upper tail of the distribution which could be produced by the residuals being from different distributions (hopefully of the same shape but with different parameters). This could occur if the oscillators were not all from the same population.

The plot of residuals versus rain rate (Fig. 5a) shows the tendency of the residuals to be positive for $10 \lesssim R \lesssim 30$ mm per hour, nega-

Fig. 4 — Normal probability plot of residuals from exponential fits.

tive for $50 \lesssim R \lesssim 100$ mm per hour, and positive again for $120 \lesssim R \lesssim 200$ mm per hour. This agrees roughly with the tendency toward nonrandom scatter of the residuals about the individual exponential fits and is not unexpected.

The plot of residuals versus frequency (Fig. 5b) shows the non-homogeneity of fit over the frequency range of the oscillators, that is, that the spread of the residuals about the fitted curve is not the same for all values of the frequency. The plot also demonstrates the tendency of the residuals toward oscillation about 0.

IV. CORRECTION FOR CALIBRATION BIAS

It was decided to try to remove the calibration bias shown as a cyclic trend in Fig. 5a. Accordingly a function of the form

$$r_1 = A \sin (\alpha R^\beta + \varphi),$$

where $r_1$ is the residual from the exponential fit and $R$ is the associated rain rate, was fit to the pooled residuals.

The resulting equation is

$$r_1 = 3.813 \sin (1.450 R^{0.3894} - 2.560)$$

with an error mean square of 46.33.

Fig. 5 — Residuals from exponential fits versus (a) rain rate and (b) frequency for all oscillators.

Fig. 6 — Sine fit of residuals from exponential fits to rain rate.

Figure 6 plots this fit and shows half of the residuals selected at random. The residuals from this fit ($r_2$) are plotted against rain rate in Fig. 7a and against frequency in Fig. 7b.

Figure 7 shows that the cyclic trend has indeed been removed by the sine function fit leaving a more random scatter of the residuals $r_2$. The spread of these residuals is different at different values of rain rate and frequency as was the spread of the residuals $r_1$. A normal plot of the residuals $r_2$ looks substantially like that of Fig. 4, with the same change in slope toward the upper tail, and somewhat more curvature in the extreme tails.

V. CONCLUSIONS

A sine corrected exponential fit to this calibration data produces a pooled standard error of less than 7 mm per hour in the rain rate. This is not very much better than the pooled standard error of about 7.3 mm per hour associated with the exponential fits alone. However, the sine correction has removed the average bias present in the exponential fits and will produce a small average error in the rain rates between

Fig. 7 — Residuals from sine fit versus (a) rain rate and (b) frequency.

10 and 100 mm per hour which is where the majority of the data of interest lies. The normal plots of the pooled residuals from the exponential fits and the residuals from the sine corrected fit are sufficiently similar so as not to give preference to either the plain or corrected exponential fit.

In processing the rainfall data collected by this network in 1967, the individual exponential coefficients for each gauge were used to calculate rain rates from the recorded frequencies. The sine correction was then applied to the rain rates from all the gauges. If the correction produced a negative value, 0 was substituted. No further corrections were made to the fits to the three gauges which showed overestimation of rain rates $\leq$ 30 mm per hour. These rain-rate values were recorded for future analysis.

## VI. ACKNOWLEDGMENT

REFERENCES

1. Semplak, R. A., "A Gauge for Continuously Measuring Rate of Rainfall," Rev. Sci. Instruments, *37*, No. 11 (November 1966), pp. 1554–1558.
2. Semplak, R. A., and Keller, H. E., "A Dense Network for Rapid Measurement of Rainfall Rate," B.S.T.J., this issue pp. 1745–1756.
3. Wilk, M. B., and Gnanadesikan, R., "Probability Plotting Methods for the Analysis of Data," Biometrika, *55*, No. 1 (June 1968), pp. 1–18.

# Some Measurements of Attenuation by Rainfall at 18.5 GHz*

## By R. A. SEMPLAK and R. H. TURRIN

(Manuscript received November 7, 1968)

*We discuss 18.5 GHz attenuation measurements on a 6.4 km path within the Holmdel rain gauge network. The period of measurement includes the summer of 1967 during which many very heavy showers occurred. We examine the data for individual storms separately. There is a marked variability. For example, one shower shows strong evidence of an updraft. The composite results show that $\gamma = 0.041 \sum_i d_i R_i^{1.04}$ where $\gamma$ is the attenuation per unit length, $R_i$ is the rainfall rate in millimeters per hour measured at the ith rain gauge on the path, and $d_i$ is the distance in kilometers over which the rain rate $R_i$ applies. When examined in detail, this relationship is satisfactory for attenuations $< 3$ dB per kilometer; however, the higher attenuations exceed this prediction and agree with the relationship $\gamma = 0.055 \langle R \rangle_{av}^{1.09}$ where $\langle R \rangle_{av}$ is the path-average rain rate. Percent-of-time distributions are given for $\langle R \rangle_{av}$, the attenuation, and the duration of attenuation. For this sample of data, the path attenuation exceeds thirty dB 0.03 percent of the time; thus 6.4 kilometers is probably too long for a conventional 18 GHz radio-relay path in New Jersey. All fades observed to date have been associated with rainfall; thus no "selective fading" has as yet been observed on this 6.4 kilometer path.*

## I. INTRODUCTION

One of the problems encountered in the utilization of the millimeter-wave bands for radio-relay systems is attenuation resulting from precipitation. There is a continued effort to overcome the fundamental inadequacy concerning the knowledge of the propagation environment. Recently a paper by Medhurst, which collects most of the published measurements on microwave attenuation by rain, has shown that there is wide disagreement between observed and calculated values.[1]

---

This paper discusses measurements at 18.5 GHz on a 6.4 kilometer path that passes above four capacitor-type rain gauges.[2] The gauges are part of a network centered at Bell Laboratories, Crawford Hill, Holmdel, New Jersey.

## II. EXPERIMENTAL ARRANGEMENT

Figure 1, which is a portion of the rain-gauge-network map, indicates the location of the 6.4 kilometer propagation path with a dashed line.[3] The dot within each grid indicates the physical location of a rain gauge mounted at the top of a telephone pole, approximately 7.6m above the ground. In particular, the four rain gauges located in grids 9, 17, 25, and 33 are closely associated with the propagation path. The rainfall rates obtained from these gauges are used in the following discussions.

## III. EQUIPMENT

The transmitting and receiving antennas are identical 76-cm diameter parabolic reflectors with a 36.8-cm focal length. To minimize waveguide losses, a spash plate fed by circular waveguide is used. The waveguide extends through the vertex of the reflector where a sidewall coupler provides the transition to rectangular $k$-band guide with vertical polarization in the dominant mode. The antenna gain measured from the rectangular waveguide port, is 39.9 dB. A klystron operating cw with a power output of about 50 mW is used as the source. The transmitter and antenna are in a small equipment house (a slanted fiberglas weather cover window provides the exit for the transmitted beam) at Cliffwood, New Jersey on the Garden State Parkway Authority right-of-way, (grid 9 of Fig. 1). The line-of-sight path has good foreground clearance at both ends. The path loss $(4\pi R/\lambda)^{-2}$ is 134 dB.

The receiver and its antenna are situated in a housing at the northern end of Crawford Hill (grid 33 of Fig. 1). Entry of the transmitted beam is provided by a sloping Mylar window. The receiver consists of a standard balanced converter followed by a low noise FET preamplifier. The single side band noise figure of the converter-preamplifier is 13 dB and the intermediate frequency is 70 MHz. To overcome effects of frequency drifts associated with the klystrons, the receiver is operated in a swept mode. However, this method of operation penalizes the dynamic measuring range by about 5 dB.

Fig. 1 — An area map showing the physical location of the rain gauges as closed dots. The dashed line is the 18.5 GHz propagation path.

The main IF amplifier is somewhat unique in that it provides an output linear in decibels, over a 50 dB dynamic range. The circuit is British in origin and uses successive transistorized twin gain stages to achieve the log-linear response.[4] The bandwidth of the IF amplifier is about 1.5 MHz and the beating oscillator is swept approximately ± 3 MHz. Visual indication of receiver tuning and sweep is obtained with an oscilloscope monitor of the video detector output. Peak detection of the video pulses provides a dc output (time constant ≈ 0.1 sec) which drives an Esterline Angus pen recorder through an FET isolating stage.

Since both rainfall and the resulting attenuation often change rapidly it is almost essential, for purposes of data reduction, to include the propagation data on the same magnetic tape that contains the information from the rain gauges;[3] thus the output of the receiver must be in a form suitable for transmission over a standard unloaded telephone pair, from the receiving site on Crawford Hill to the data collection room in the Crawford Hill laboratory. For this purpose, a voltage-to-frequency converter is used to convert the receiver dc signal into a frequency with a range of 1 KHz to 10 KHz which is suitable for telephone line transmission. The receiver signal is then sampled and recorded on the magnetic tape at the rate of one sample per second.[3]

A permanently installed 18.5 GHz precision variable attenuator is used to calibrate the receiving system. System and path-loss considerations indicate a maximum received signal at the receiving antenna port of the order −38 dBm. The receiver sensitivity is −94 dBm including a 5 dB penalty for sweeping the beating oscillator. Thus the available dynamic range is 56 dB. However, the log-linear amplifier and precision attenuator have a useful dynamic range of 50 dB.

IV. RESULTS

The data discussed here were obtained for the recording period June 23 through October 31, 1967, during which there were twenty-three storms associated with the propagation path.* Complete propagation and rainfall data are available for twenty-one of these storms.

Three storms have been selected to show the marked variability between measured and predicted attenuation on an individual-storm

---

* The summer of 1967 in New Jersey was prolific in heavy showers. During this period, of 3144 hours, a total of 21 inches of rain was measured at Crawford Hill. The annual average rainfall for New Jersey is 40 inches.

basis. Figures 2a, 3a, and 4a show strip chart records obtained during rains on July 21, July 28, and October 25, respectively. An examination of these figures reveals a lack of similarity in both the duration and absolute value of the attenuation.

The measured attenuation is averaged along the path. A comparison is made of the observed attenuations with values predicted from the rainfall rates. For all data to be discussed, the attenuations were read



Fig. 2 — (a) Propagation data July 21, 1967. (b) Observed versus predicted path attenuation for July 21, 1967, where the weighing factor $w_i = d_i/6.4$.

Fig. 3 — (a) Propagation data July 28, 1967. These data are considered as two separate storms and the remarks of the text refer to the second storm beginning about 4:50 P.M. (b) Observed versus predicted path attenuation for July 28, 1967, where the weighing factor $w_i = d_i/6.4$.

once from the magnetic tape during every 30-second interval; the corresponding rainfall rates recorded every ten seconds on the tape were averaged over this same 30-second interval. Plots on individual storms showed that averaging in this manner reduced the scatter in the calculated data. For the predictions, the empirical expression proposed by Gunn and East[5]

$$\gamma = k \sum_i d_i R_i^\alpha \tag{1}$$

are used. $R_i$ is the rainfall rate measured at the $i$th rain gauge, $d_i$ the distance in km over which the rain rate $R_i$ applies, $\Sigma_i d_i = d = 6.4$ km, and $i = 1, 2, 3, 4$ is the gauge position.* $k = 0.055$ dB per km and $\alpha = 1.09$ are the values suggested by Gunn and East, based on the Laws and Parsons drop-size distribution.

Heavy rainfalls occur over small areas; it is believed that this area decreases (in general) with increasing rainfall rate. Therefore, unless the gauge spacing is small it is less likely that a heavy rainfall will occur at a gauge but more likely that it will occur between gauges. In such a case, equation (1) tends to underestimate the attenuation. When a heavy rainfall does occur at one gauge (assuming that gauges are not sufficiently closely spaced) equation (1) tends to over-estimate; this is because $d_i$ is fixed and because we assume that the heavy rain rate measured at the $i$th gauge extends uniformly along $d_i$. From experience with the network, we have learned that the gauge spacing is somewhat too large to permit resolving the distance associated with some of the very heavy rain rates. Since the rain gauges are sampled every ten seconds, there is also some sampling error in the time domain; this is discussed in the appendix. In particular, note that the four gauges associated with the propagation path are not read consecutively. There is an approximate two second interval between successive readings from the four gauges on the path.

Comparison of the measured attenuation with that calculated using equation (1) for the three selected storms is shown in Figs. 2b, 3b, and 4b. In Fig. 2b, the predicted values agree with the measurements. The agreement is not nearly as good for the rainstorm corresponding to Fig. 3b. The spread of the calculated points is fairly well contained but the indications are that both $k$ and $\alpha$ are too large, that is, we predict an attenuation higher than is actually measured. Both of these constants depend on the drop-size distribution.† The comparison for the last

* The appropriate values are $d_1 = 1.933$ km, $d_2 = 1.792$ km, $d_3 = 1.741$ km, and $d_4 = 0.934$ km.
† $k$ is also a function of the fall-velocity of the drops.

DATA OF OCTOBER 25, 1967
METEOROLOGICAL DATA–FOR 18.5 GHz
6.4 km PATH–CLIFFWOOD–CRAWFORD HILL

(a)

(b)

Fig. 4 — (a) Propagation data for October 25, 1967 along with average rain rate; the lower charts show the pressure and wind at Crawford Hill. (b) Observed versus predicted path attenuation for October 25, 1967, where the weighting factor $w_i = d_i/6.4$.

storm to be considered on an individual basis is shown in Fig. 4b. It is readily apparent that here is a case where the predictions are under-estimates; these data are discussed in detail later in this section. For the moment, it is sufficient to say that equation (1) with the above values for $k$ and $\alpha$ tends to overestimate the attenuation in most cases; this is readily seen in Fig. 5 where data for the period of June 23, 1967 to October 25, 1967 have been plotted. Using a nonlinear least squares regression, the values $0.041 \pm 0.007$ dB per km and $1.042 \pm 0.031$ are found for $k$ and $\alpha$ respectively.

There are at least three factors which enter into predicting values higher than the measured ones, as in Fig. 5. First, the drop-size distribution often may differ from that commonly used in computation. Secondly, the fall velocity of the drops has a strong effect. The third factor is polarization; since falling rain drops tend to be oblate, attenuation for a vertically polarized electromagnetic wave (as used here) is several percent less than values calculated assuming spherical drops.[5,6]

Fig. 5 — Plot of predicted values $\gamma = k \ \Sigma_i \ d_i R_i^a$ where $k = 0.055$ dB per kilometer and $\alpha = 1.09$ for all storms during the period June 23 to October 31, 1967.

## 4.1 A Special Case

Let us return to the storm of October 25 (Figs. 4a and b) and examine Fig. 6 which is a plot of measured attenuation versus average rain rate. The dashed lines on this plot are theoretical maximum and minimum attenuations for hypothetical rains containing drops which are all of the same diameter;[1] the diameters specifically chosen are those which produce the maximum and minimum attenuation.* For this case the measured values of attenuation lie well above the theoretical maximum. The spread of the data is small and the solid line is a calculated least squares regression. Since these data differ considerably from the rest of our measurements, an attempt to account for this behavior was made by examining synoptic data for this period.

Large-scale synoptic data indicated the passage of a cold front

---

* Commonly accepted terminal velocities are used in these estimates.

system; this was confirmed by our local wind and pressure measurements shown in Fig. 4a. There is a pressure drop and a wind reversal near the time the attenuation occurred. The classical picture for vertical wind movement at the interface of the cold front is updrafts associated with the retreating warm system and downdrafts in the advancing cold front.[7] Since an updraft decreases the terminal velocity of the rain drops, one has a condition where the density of rain drops in the radio path is greater than is indicated by the rain gauges.

Assume for the moment a modest updraft of two miles per hour, that



Fig. 6 — Plot of measured attenuation obtained for October 25, 1967 versus average rain rate.

is, 0.95 m per second. Using this value to correct the rain drops terminal velocities (about 7 m per second) to account for the effect of the updraft, the three points indicated by crosses on Fig. 6 are calculated. If an updraft of 1.9 m per second is assumed, the points indicated by the open circles on Fig. 6 are obtained. A line drawn through the crosses lies close to the calculated regression line and indicates that the updraft is of the order 1 m per second; thus with this assumption, the data are readily explained. Presumably the converse can take place; that is, with a downdraft of 1 m per second the measured attenuation would be about 15 per cent lower than one would estimate.

## 4.2 *Comparison with Path-Averaged Rain Rate*

Figure 7 shows the totality of the attenuation data plotted as a function of average rain rate $\langle R \rangle_{av}$ (and not the summation of attenuations discussed earlier in Section IV). Since the data obtained from the storm of October 25, (shown by crosses) are reasonably accounted for by an updraft, they are not included in the following discussion. The dashed lines in Fig. 7 are, as in Fig. 6, theoretical maximum and minimum boundaries and the line labeled $\gamma = 0.048 \langle R \rangle_{av}$ dB per km is a linear least square regression for these data;[1] this gives an attenuation directly proportional to average rain rate. To remove any possible influence of improper readings from the rain gauges at low rates, that is, small drops lodging in the trough-type capacitor, these measurements at average rain rates of 10 mm per hour and less have been removed; this accounts for the absence of data points in this range in Fig. 7.

The bulk of the data is well contained within the theoretical boundaries; however, there is a cluster of points near the origin that lies below the boundary. At higher rain rates neither the linear regression line nor the nonlinear regression line labeled $\gamma = 0.04 \langle R \rangle_{av}^{1.04}$ dB per kilometers provide suitable prediction. The line labeled $\gamma = 0.055 \langle R \rangle_{av}^{1.09}$ dB per kilometers, from Gunn and East appears to represent the higher attenuations fairly well; these higher values of attenuation and the associated rainrates are discussed further in Section 4.3.

## 4.3 *Distributions*

Percent of time distributions for attenuation and average rain rate are plotted in Figs. 8a and b respectively. In addition to the distribution (Fig. 8a) for the summer period, distributions for the period covering spring, fall, and winter and for the full year are shown.

Fig. 7 — Plot of observed attenuation versus average path rain rate, $\langle R \rangle_{av}$, as measured by four gauges with an intergauge spacing of 1.6 km, for all storms during the period June 23 through October 31, 1967. The dashed lines are theoretical maxima and minima. The solid line is a least squares regression. The data points represented by crosses pertain to October 25, 1967.

Fig. 8 — (a) Distributions of the percent of time the attenuation exceeds the abscissa for periods of measurements of 3144 hours (summer), 5616 hours winter), and 8760 hours (full year). (b) The percent of time the average rain rate exceeds the abscissa for a total period of 3144 hours (summer).

That the summer distribution is very conservative is readily apparent.

Bussey proposed that there was a relationship between point and space-averaged rainfall rates.[8] If certain, not fully defined restrictions are placed on the length of the path over which the average is taken, system designs based on the point rainfall rates could provide reliable operation. A comparison of point and path averaged distributions is shown in Fig. 8b. The three dashed curves are for point rainfall rate distributions from gauge 25 on the propagation path; averages over time intervals of 10 seconds, 3 minutes, and 12 minutes were used. One sees that, as the time of averaging interval is increased, the tail of the distribution obtained from the point gauge approaches that obtained for the path-average distribution (the solid curve). For rain rates less than about 100 mm per hour, the equivalence between the point and path-average distributions is fairly good; however, for higher rates, the point distribution, even for a 12 minute averaging interval, indicates higher percentage times than the true path-average distribution.

Rather than plotting instantaneous attenuation versus instantaneous path-average rain rate, as has been done in Fig. 7, one can examine the relationship between these two variables by means of the cumulative distributions in Fig. 8. Since the same samples was used in obtaining the distributions of attenuation (Fig. 8a) and average rain rate (Fig. 8b), one can draw a correspondence between the attenuation and the average rain rate for a given probability. For example, at 0.03 per cent probability the average rain rate is 80 mm per hour and the attenuation 30 dB; at the 0.003 per cent value, the average rate is 113 mm per hour and the attenuation 45 dB. For these low probabilities, both distributions are quite linear on the semilogarithmic plot (see Fig. 8). Thus for the 6.4 km path, one obtains

$$\gamma \text{ path} = 0.38 \langle R \rangle_{av} \text{ dB}, \langle R \rangle_{av} > 80 \text{ mm per hour},$$

which results in

$$\gamma = 0.059 \langle R \rangle_{av} \text{ dB/km}, \langle R \rangle_{av} > 80 \text{ mm per hour}.$$

Let us compare this result with the plots of instantaneous data for the 6.4 km path given in Fig. 7. At a path-average rate of 80 mm per hour the Gunn-East curve shows an attenuation of 43 dB; the method of comparison of distributions just discussed, results in an attenuation of 30 dB; and the least squares fit, 26 dB. However the rain gauges are sensors only of points on the path; and since it is more probable

that intense rain cells occur between gauges than on a gauge, the points that fall along the upper end of the Gunn and East curve (the high attenuations) are believed to be associated with actual average rain rates higher than those plotted in the figure. The effect of inadequate sampling of rain rate in the time domain, discussed in the appendix, also leads to conservative estimates of attenuation. Suppose an intense cell produces a rapid peak in the rain rate. If, as a result of sampling error, this peak were not recorded, a given attenuation would be plotted against a rain rate lower than that which actually existed. Again, this means that the high attenuation points are actually associated with average rain rates higher than those plotted in Fig. 7.

The above discussion leads to the question: Based on the present experiment, what is the best relationship between the attenuation and the path-average rain rate? One can only say that at a rate of 100 mm per hour, the attenuation is between five and eight dB per km. However,

$$\gamma = 0.06\langle R\rangle_{\mathrm{av}} \qquad \mathrm{dB/km}$$

is considered a working relationship for path-average rates exceeding 50 mm per hour.

Distributions of the duration of fades are plotted in Fig. 9. The ordinate is the percentage of the total number of fades for the 17.8 hours of heavy rain that occured during the period June 23 through October 31. (An example: at the −20 dB level, 5.8 percent of the total number of fades have a duration exceeding five minutes. The total number of fades is 182, therefore 10 fades have durations greater than five minutes at the −20 dB level.

Attenuation by rain at a rate of 100 mm per hour is shown in Fig. 10 for the 8 to 80 GHz band. Most of the plotted points are measurements of the Bell Telephone Laboratories; the measurements at 8 and 15 GHz were made in Canada.[9] All of the points correspond to the heavy average rain rate of 100 mm per hour. In most cases the points were actually measured at that rate; for the others, as indicated by arrows, some extrapolation was made. The solid line is a faired-in fit to the data points.

The broken line is a weighted mean of measured data from many parts of the world assembled by Medhurst for 100 mm per hour rate.[1] The two dashed lines, also shown, represent the maxima and minima of these observations. Clearly, the solid line lies well below the "world" average over all of the band.

Fig. 9 — Percent of total number of fades for which the fade length exceeds the abscissa. Based on 17.8 hours of rain data and a total of 182 fades.

V. CONCLUSIONS

Calculations using the commonly accepted expression $\gamma = k \sum_i d_i R_i^\alpha$ where $k = 0.055$ dB per kilometer and $\alpha = 1.09$ overestimate the present measurements of 18.5 GHz attenuation by rain. Our data, obtained from four rain gauges with an intergauge spacing of 1.6 km, result in values $k = 0.041$ and $\alpha = 1.04$ over much of the range of observations. The discrepancy may be caused in part by formation of ellipsoidal drops at high rainrates; these would produce lower attenuation for vertical polarization. An incorrect assumption for the fall velocity has a strong influence; as shown by the data for October 25 (Fig. 4) an updraft can produce attenuations above the bounds of conventional predictions. On the other hand, a drop velocity greater than normally assumed in the theory would modify the estimates in a direction toward better agreement with most of the storms we have measured. On the 6.4 km path, the 18.5 GHz attenuation exceeded 30 dB 0.03 percent of the time for this summer season; this path is therefore

Fig. 10 — The "world" weighted mean from Medhurst is shown as a broken line.[1] The solid line represents data appropriate to an average rainfall rate of 100 mm per hour and is considered the best available basis for prediction in system design. Data plotted:

    8 and 15 GHz—Blevis and others[9]
    11 GHz—Hathaway and Evans[11]
    18.5 GHz—Semplak and Turrin (present data).
    48 and 70 GHz—Hogg.[12]

too long for the repeater spacing in a conventional system at this frequency in the New Jersey climate.* For the present sample of data, the average rain rate on the 6.4 km path exceeded 100 mm per hour 0.008 percent of the time; however, point rates in excess of 280 mm per hour have been observed for very short periods. All fades observed to date have been associated with rainfall; thus no "selective fading" has as yet been observed on this 6.4 km path.

VI. ACKNOWLEDGMENTS

The authors are grateful to D. C. Hogg for helpful discussions and to Mrs. E. Kerschbaumer for the computations. The assistance of H. A.

---

* For a discussion of path-diversity systems see Ref. 10.

Fig. 11 — Comparison of one and ten second sampling of rain gauge data.

APPENDIX A

*Effect of Temporal Sampling Error in the Raingauge Data*

As discussed in a companion paper, the output of the individual rain gauges in the network are sampled and recorded once every ten seconds.[3] In addition, one of the gauges (grid 33 of Fig. 1) is sampled and recorded every second. A portion of the recorded one second sampling data from this gauge for the storm of July 21, 1967 (Fig. 2a) are shown by the solid line of Fig. 11. The dashed line connects the normal 10 second sampling data points from this gauge; it is apparent that 10 second sampling misses many of the peaks shown by the one second sampling. However cumulative distributions of both the one (solid line) and ten second samplings (dashed line) for this storm, whose duration was thirteen minutes, are shown as Fig. 12; one can readily see that the distributions are identical out to rain rates of the order 100 mm per hour. For rates greater than 100 mm per hour, the ten second sampling tends to be conservative. Another cumulative distribution based on a ten second sampling period of this data and



Fig. 12 — Cumulative distributions of rain fall rates based on one and ten second sampling.

selected in a manner such that a minimum of one second peaks would be encountered is shown by the broken line on this figure. Here again this result is not too different from the one second sampling in Fig. 12. Both ten second distributions straddle the one second distribution.

REFERENCES

1. Medhurst, R. G., "Rainfall Attenuation of Centimeter Waves: Comparison of Theory and Experiment," IEEE Trans. Antennas and Propagation, *AP-13*, No. 4 (July 1965), pp. 550–564.
2. Semplak, R. A., "Gauge for Continuously Measuring Rate of Rain Fall," Rev. of Sci. Instruments, *37*, No. 11 (November 1966), pp. 1554–1558.
3. Semplak, R. A., and Keller, H. E., "A Dense Network for Rapid Measurement of Rainfall Rate," B.S.T.J., this issue, pp. 1745–1756.
4. Croney, J., and Worowcow, A., "A True IF Logarithmic Amplifier Using Twin-Gain Stages," REE, *32*, No. 3 (September 1966), pp. 149–155.
5. Gunn, K. L. S., and East, T. W. R., "The Microwave Properties of Precipitation Particles," J. Royal Meteorological Soc., *80*, 1954, pp. 522–545.
6. Oguchi, T., "Attenuation of Electromagnetic Wave Due to Rain with Distorted Rain Drops," J. of the Radio Res. Labs (Tokyo), *7*, No. 33 (September 1960), pp. 467–485.
7. Berry, F. A., Jr., Bollay, E., and Beers, N. R., *Handbook of Meterology*, New York: McGraw-Hill, 1945, pp. 638ff.
8. Bussey, H. E., "Microwave Attenuation Statistics Estimated from Rainfall and Water Vapor Statistics," Proc. IRE, *38*, No. 7 (July 1950), pp. 781–785.
9. Blevis, B. C., Dohoo, R. M., and McCormick, K. S., "Measurements of Rainfall Attenuation at 8 and 15 GHz," IEEE Trans. Antennas and Propagation, *AP-15*, No. 3 (May 1967), pp. 394–403.
10. Hogg, D. C., "Path Diversity in Propagation of Millimeter Waves through Rain," IEEE Trans. *T-AP 15*, No. 3 (May 1967), pp. 410–415.
11. Hathaway, S. D., and Evans, H. W., "Radio Attenuation at 11 Kmc and Implications Affecting Relay System Engineering," B.S.T.J., *38*, No. 1 (January 1959), pp. 73–98.
12. Hogg, D. C., "Millimeter-Wave Communication through the Atmosphere," Science, *159*, No. 3810 (January 1968), pp. 39–46.

# A Statistical Description of Intense Rainfall

By MRS. A. E. FREENY and J. D. GABBE

*This paper contains a statistical summary of the 14,000,000 measurements taken during 27 rainfalls in a six-month period in 1967 from a 96-station, rapid-response rain guage network spread over a rectangular area 13 by 14 kilometers centered near Crawford Hill, New Jersey. The analysis emphasizes rain rates greater than 50 millimeters per hour, which interfere with radio transmission in the 10 to 30 GHz frequency range.*

*Heavy rain rates are relatively rare events, come in irregular bursts, and do not appear amenable to description by simple analytic distributions. This paper presents statistics concerning the behavior of rain rates at a point in space, the relationship of rain rates separated in space or time, and the relationship of average rain rates on pairs of paths in various configurations.*

## I. INTRODUCTION

This paper presents some statistics from the rainfall data collected on a rain gauge network during the period from June 1 to November 30, 1967. The network consists of 96 gauges spaced approximately 1.3 km apart on a rectangular grid centered near Crawford Hill, Holmdel, New Jersey. The design of the rain gauges and the equipment for recording data from the network are described elsewhere.[1,2]

In communications, interest in rain-rate data arises from the relationship of attenuation of radio signals in the 10 to 30 GHz frequency range to the number and size of raindrops present in the transmission path. The quantity of water that the signal penetrates is directly related to the average rain rate on the path. Thus a major direction of our analysis was toward a statistical description of the behavior of rain rates at a point in space, the relationship of two rain rates separated in space or in time, and the relationship of average rain rates on pairs of paths in various configurations. Knowledge about these relationships, particularly for rain rates greater than 50 mm per hour,

where substantial attenuation occurs,[3] is important for the design of microwave radio transmission systems.

A major characteristic of the rain-rate data taken in this experiment is its extreme variability, which increases as the rain rate increases. The separations between successive readings in time (10 s) and neighboring readings in space ($\sim$ 1.3 km) are too large to provide continuous representation of rain-rate behavior through time and space. Because the time series at the measuring stations cannot be considered even piecewise stationary during intense rainfalls, the usual time series techniques are not applicable. This leads to large numbers of descriptive statistics rather than a concise representation of the characteristics of rainfall.

A brief summary of the results and general conclusions given in Section VIII is: On the basis of 14,000,000 measurements obtained during 27 rainfalls that occurred in the Crawford Hill locale during the 1967 recording season, the empirical probabilities of observing point rain rates above 50, 100, 150, and 200 mm per hour are found to be $4.3 \times 10^{-4}$, $1.3 \times 10^{-4}$, $4.2 \times 10^{-5}$, and $1.0 \times 10^{-5}$ respectively; the joint probability that the rain rate exceeds a given value at both of two stations simultaneously decreases rapidly at short distances as the separation between the stations increases, and goes through a minimum at a separation of about 12 km; the joint probability that the average rain rate on both of a pair of parallel paths exceeds a given value decreases as the path length increases, and shows a minimum for paths separated by 9 km; and the probability that the average rain rate will exceed 150 mm per hour on a single path 6.5 km long is 200 times greater than the probability that the average rain rate will simultaneously exceed 150 mm per hour on both of two parallel paths 6.5 km long and 6.5 km apart.

Detailed descriptions of the components of the analysis are:

Section II—treatment of the data, selection of a subset for analysis and procedures used for the detection and removal of spurious data;

Section III—general characteristics of observed rain, rainfall behavior at individual stations and in individual rainfalls, and the partitioning of the data made necessary by the variability of rain;

Section IV—statistics on point rain rates, both for the complete subset of the data and under the condition that the rain rate is greater than 50 mm per hour;

Section V—conditional and joint probabilities of various events for

pairs of stations at selected separations in space or time, and the relationship of the probabilities with distance and time;

Section VI—results of an analysis of average rain rates on pairs of paths in various configurations;

Section VII—engineering calculations and "benchmarks" for translating the relative probabilities calculated from the selected sample to "annual" probabilities based on the duration of the experiment, comparison of the results of various analyses using these data, and presentation of a result from another body of data.

## II. TREATMENT OF THE DATA

Most bodies of data containing several million data points can be separated into segments of primary and secondary interest in the context of a particular analysis. Furthermore, raw data in such quantities inevitably contains some fraction of specious readings. In this section a brief description of the data collection system is followed by several presentations of rain-rate data and an outline of the data selection and screening procedures. The closing subsection contains remarks concerning the data retained for analysis.

### 2.1 *The Basic Data*

The data were acquired with a network of rapid-response rain-rate gauges most of which were mounted on telephone poles about 15 meters above the ground and well clear of all obstructions.[1,2] The collecting surface of the gauge has an area of 478 $cm^2$; the response characteristics of the gauge are such that the rain rate represents an average over less than one second. Although occasional anomalies in the data appear to be traceable to the gauges, they seem to give satisfactory readings for rain rates greater than about 10 mm per hour.

The network consisted of 96 stations, about 85 percent of which were operational at any given time. The area around Crawford Hill was divided into squares 1.29 km on a side by grid lines oriented north-south and east-west; one of the gauges was emplaced as close to the center of each square as practicable. The actual positions of the devices are shown on the map in Fig. 1. The data from the rain gauges, in the form of oscillator frequencies, is telemetered via telephone lines into the telephone central office serving the gauge location. At the central offices, which act as collecting points, the readings are commutated and forwarded to the Crawford Hill Laboratory, where the frequencies are detected. These frequencies and some auxiliary data are then automatically recorded on magnetic tape.

Fig. 1 — The rain gauge network, showing the station numbers, scanning order (small circled numbers), and row numbers (in squares). Station 33 (CH) is at Crawford Hill, New Jersey.

During periods of rain, the network is scanned once every ten seconds to produce a "scan" of readings. The gauges are not read simultaneously, but are scanned at the rate of ten gauges per second in a sequence dictated by the telemetry arrangements and indicated by the numbers in the lower left corners of the grid squares in Fig. 1. The gauge at Crawford Hill (station 33) was read once per second separately from, and in addition to, the regular scan. This high frequency sample is recorded as part of the auxiliary data. The solid line between the transmitter (T) and receiver (R) shows an experimental microwave transmission path. A pair of the parallel paths ($P_{11}$—$P_{12}$, $P_{21}$—$P_{22}$) and a pair of the adjoining paths ($I_1$—$I_4$—$I_2$) used in the path analysis are indicated by dashed lines. During the time covered by this report, the devices were not interchanged among grid locations, so that a station number always refers to the same device. The telemetry system is exposed to interfering signals which may produce spurious readings, creating some of the data-screening problems discussed in Section 2.2.

The data recorded at Crawford Hill are transferred from the original tapes onto tapes compatible with the computers used in the data analysis. For a variety of reasons this process could not be carried out successfully with about ten percent of the tapes; these data were lost. The next major processing programs transformed the data from oscillator frequencies to rain rates by applying calibration functions.[4] As noted in Ref. 4 these calibrations have been optimized for the higher rain rates at the expense of some loss of relative accuracy below 25 mm per hour. The programs also (*i*) discarded time periods during which rain rates greater than 15 mm per hour were recorded at fewer than four stations during each scan for at least 20 consecutive minutes, (*ii*) indicated some kinds of questionable data, (*iii*) appended additional information, such as a list of stations known to be inoperative, (*iv*) organized the data into a format suitable for subsequent analysis, and (*v*) produced the isometric plots described in Section 2.2.

## 2.2 *Data Selection and Screening*

It is neither practical nor desirable to analyze all the data acquired from the rain gauge network. This is partly because more than 70 percent of the data were taken when the rain was too light to seriously affect microwave transmission in the frequency range of interest, partly because a tiny but important fraction of the data are spurious, and partly because the total number of data points is so large (over 14,000,000).

The first selection procedure operating on the data was "natural" selection, that is, owing to various equipment failures we have no data for a number of rainfalls. The second selection took place during the data processing where lengthy periods, during which virtually no rain was falling, were discarded. At this point about 430 hours of processable data remained. The procedures used to identify data of interest from the available 430 hours and remove spurious values are based on several graphical presentations of the data.

The isometric plots, which were produced for each rainfall in the final stage of the data processing, are the most basic graphical presentation. These plots give one a spatial and temporal appreciation of the rainfall as a whole, in addition to displaying particular features of the data in an illuminating manner.

Three hours of data from the rainfall of July 28, and one hour from July 11, 1967 are displayed in Figs. 2a and b. Each solid trace represents the rain-rate measurements from a single station, and is constructed by connecting six-reading averages in Fig. 2a and two-reading averages in Fig. 2b. The time in minutes is measured from the first data set. The traces are arranged in a horizontal and vertical grid which is isomorphic to the geographic grid in Fig. 1. Each trace is labeled with the station number. No trace is plotted if the station is known to be inoperative.

On the basis of the isometric plots, a sample of about 110 hours of rain was chosen for analysis in the following somewhat subjective but operationally convenient manner. An hour's worth of rain was included or rejected as a unit, with the exception of one isolated 20-minute shower which was included. A unit during which any six-reading average rain rate exceeded 50 mm per hour was included. A unit during which no six-reading average exceeded 30 mm per hour was excluded unless it fell between two units which were included on the basis of the 50-mm-per-hour criterion. Units with occasional six-reading averages greater than 30 mm per hour but less than 50 mm per hour were included if there were other units on the same day chosen on the basis of the 50-mm-per-hour criterion, and excluded otherwise. Subject to these constraints, the beginning and end of each time period was determined somewhat arbitrarily in terms of such operational conveniences as accepting an "entire rainfall" or eliminating a few bad magnetic tape records near the beginning or end of a rainfall. For example the entire three-hour period shown in Fig. 2a was accepted, and no effort was made to delete the first 20 minutes, during which there is little rain. On the whole, we were generous in the

Fig. 2 — Isometric plots of rainfalls for (a) July 28, 1967 and (b) July 11, 1967.

# A Statistical Description of Intense Rainfall

By MRS. A. E. FREENY and J. D. GABBE

*This paper contains a statistical summary of the 14,000,000 measurements taken during 27 rainfalls in a six-month period in 1967 from a 96-station, rapid-response rain guage network spread over a rectangular area 13 by 14 kilometers centered near Crawford Hill, New Jersey. The analysis emphasizes rain rates greater than 50 millimeters per hour, which interfere with radio transmission in the 10 to 30 GHz frequency range.*

*Heavy rain rates are relatively rare events, come in irregular bursts, and do not appear amenable to description by simple analytic distributions. This paper presents statistics concerning the behavior of rain rates at a point in space, the relationship of rain rates separated in space or time, and the relationship of average rain rates on pairs of paths in various configurations.*

## I. INTRODUCTION

This paper presents some statistics from the rainfall data collected on a rain gauge network during the period from June 1 to November 30, 1967. The network consists of 96 gauges spaced approximately 1.3 km apart on a rectangular grid centered near Crawford Hill, Holmdel, New Jersey. The design of the rain gauges and the equipment for recording data from the network are described elsewhere.[1,2]

In communications, interest in rain-rate data arises from the relationship of attenuation of radio signals in the 10 to 30 GHz frequency range to the number and size of raindrops present in the transmission path. The quantity of water that the signal penetrates is directly related to the average rain rate on the path. Thus a major direction of our analysis was toward a statistical description of the behavior of rain rates at a point in space, the relationship of two rain rates separated in space or in time, and the relationship of average rain rates on pairs of paths in various configurations. Knowledge about these relationships, particularly for rain rates greater than 50 mm per hour,

where substantial attenuation occurs,[3] is important for the design of microwave radio transmission systems.

A major characteristic of the rain-rate data taken in this experiment is its extreme variability, which increases as the rain rate increases. The separations between successive readings in time (10 s) and neighboring readings in space ($\sim$ 1.3 km) are too large to provide continuous representation of rain-rate behavior through time and space. Because the time series at the measuring stations cannot be considered even piecewise stationary during intense rainfalls, the usual time series techniques are not applicable. This leads to large numbers of descriptive statistics rather than a concise representation of the characteristics of rainfall.

A brief summary of the results and general conclusions given in Section VIII is: On the basis of 14,000,000 measurements obtained during 27 rainfalls that occurred in the Crawford Hill locale during the 1967 recording season, the empirical probabilities of observing point rain rates above 50, 100, 150, and 200 mm per hour are found to be $4.3 \times 10^{-4}$, $1.3 \times 10^{-4}$, $4.2 \times 10^{-5}$, and $1.0 \times 10^{-5}$ respectively; the joint probability that the rain rate exceeds a given value at both of two stations simultaneously decreases rapidly at short distances as the separation between the stations increases, and goes through a minimum at a separation of about 12 km; the joint probability that the average rain rate on both of a pair of parallel paths exceeds a given value decreases as the path length increases, and shows a minimum for paths separated by 9 km; and the probability that the average rain rate will exceed 150 mm per hour on a single path 6.5 km long is 200 times greater than the probability that the average rain rate will simultaneously exceed 150 mm per hour on both of two parallel paths 6.5 km long and 6.5 km apart.

Detailed descriptions of the components of the analysis are:

Section II—treatment of the data, selection of a subset for analysis and procedures used for the detection and removal of spurious data;

Section III—general characteristics of observed rain, rainfall behavior at individual stations and in individual rainfalls, and the partitioning of the data made necessary by the variability of rain;

Section IV—statistics on point rain rates, both for the complete subset of the data and under the condition that the rain rate is greater than 50 mm per hour;

Section V—conditional and joint probabilities of various events for

pairs of stations at selected separations in space or time, and the relationship of the probabilities with distance and time;

Section VI—results of an analysis of average rain rates on pairs of paths in various configurations;

Section VII—engineering calculations and "benchmarks" for translating the relative probabilities calculated from the selected sample to "annual" probabilities based on the duration of the experiment, comparison of the results of various analyses using these data, and presentation of a result from another body of data.

## II. TREATMENT OF THE DATA

Most bodies of data containing several million data points can be separated into segments of primary and secondary interest in the context of a particular analysis. Furthermore, raw data in such quantities inevitably contains some fraction of specious readings. In this section a brief description of the data collection system is followed by several presentations of rain-rate data and an outline of the data selection and screening procedures. The closing subsection contains remarks concerning the data retained for analysis.

### 2.1 *The Basic Data*

The data were acquired with a network of rapid-response rain-rate gauges most of which were mounted on telephone poles about 15 meters above the ground and well clear of all obstructions.[1,2] The collecting surface of the gauge has an area of 478 cm$^2$; the response characteristics of the gauge are such that the rain rate represents an average over less than one second. Although occasional anomalies in the data appear to be traceable to the gauges, they seem to give satisfactory readings for rain rates greater than about 10 mm per hour.

The network consisted of 96 stations, about 85 percent of which were operational at any given time. The area around Crawford Hill was divided into squares 1.29 km on a side by grid lines oriented north-south and east-west; one of the gauges was emplaced as close to the center of each square as practicable. The actual positions of the devices are shown on the map in Fig. 1. The data from the rain gauges, in the form of oscillator frequencies, is telemetered via telephone lines into the telephone central office serving the gauge location. At the central offices, which act as collecting points, the readings are commutated and forwarded to the Crawford Hill Laboratory, where the frequencies are detected. These frequencies and some auxiliary data are then automatically recorded on magnetic tape.

Fig. 1 — The rain gauge network, showing the station numbers, scanning order (small circled numbers), and row numbers (in squares). Station 33 (CH) is at Crawford Hill, New Jersey.

During periods of rain, the network is scanned once every ten seconds to produce a "scan" of readings. The gauges are not read simultaneously, but are scanned at the rate of ten gauges per second in a sequence dictated by the telemetry arrangements and indicated by the numbers in the lower left corners of the grid squares in Fig. 1. The gauge at Crawford Hill (station 33) was read once per second separately from, and in addition to, the regular scan. This high frequency sample is recorded as part of the auxiliary data. The solid line between the transmitter (T) and receiver (R) shows an experimental microwave transmission path. A pair of the parallel paths ($P_{11}$—$P_{12}$, $P_{21}$—$P_{22}$) and a pair of the adjoining paths ($I_1$—$I_4$—$I_2$) used in the path analysis are indicated by dashed lines. During the time covered by this report, the devices were not interchanged among grid locations, so that a station number always refers to the same device. The telemetry system is exposed to interfering signals which may produce spurious readings, creating some of the data-screening problems discussed in Section 2.2.

The data recorded at Crawford Hill are transferred from the original tapes onto tapes compatible with the computers used in the data analysis. For a variety of reasons this process could not be carried out successfully with about ten percent of the tapes; these data were lost. The next major processing programs transformed the data from oscillator frequencies to rain rates by applying calibration functions.[4] As noted in Ref. 4 these calibrations have been optimized for the higher rain rates at the expense of some loss of relative accuracy below 25 mm per hour. The programs also (*i*) discarded time periods during which rain rates greater than 15 mm per hour were recorded at fewer than four stations during each scan for at least 20 consecutive minutes, (*ii*) indicated some kinds of questionable data, (*iii*) appended additional information, such as a list of stations known to be inoperative, (*iv*) organized the data into a format suitable for subsequent analysis, and (*v*) produced the isometric plots described in Section 2.2.

## 2.2 *Data Selection and Screening*

It is neither practical nor desirable to analyze all the data acquired from the rain gauge network. This is partly because more than 70 percent of the data were taken when the rain was too light to seriously affect microwave transmission in the frequency range of interest, partly because a tiny but important fraction of the data are spurious, and partly because the total number of data points is so large (over 14,000,000).

The first selection procedure operating on the data was "natural" selection, that is, owing to various equipment failures we have no data for a number of rainfalls. The second selection took place during the data processing where lengthy periods, during which virtually no rain was falling, were discarded. At this point about 430 hours of processable data remained. The procedures used to identify data of interest from the available 430 hours and remove spurious values are based on several graphical presentations of the data.

The isometric plots, which were produced for each rainfall in the final stage of the data processing, are the most basic graphical presentation. These plots give one a spatial and temporal appreciation of the rainfall as a whole, in addition to displaying particular features of the data in an illuminating manner.

Three hours of data from the rainfall of July 28, and one hour from July 11, 1967 are displayed in Figs. 2a and b. Each solid trace represents the rain-rate measurements from a single station, and is constructed by connecting six-reading averages in Fig. 2a and two-reading averages in Fig. 2b. The time in minutes is measured from the first data set. The traces are arranged in a horizontal and vertical grid which is isomorphic to the geographic grid in Fig. 1. Each trace is labeled with the station number. No trace is plotted if the station is known to be inoperative.

On the basis of the isometric plots, a sample of about 110 hours of rain was chosen for analysis in the following somewhat subjective but operationally convenient manner. An hour's worth of rain was included or rejected as a unit, with the exception of one isolated 20-minute shower which was included. A unit during which any six-reading average rain rate exceeded 50 mm per hour was included. A unit during which no six-reading average exceeded 30 mm per hour was excluded unless it fell between two units which were included on the basis of the 50-mm-per-hour criterion. Units with occasional six-reading averages greater than 30 mm per hour but less than 50 mm per hour were included if there were other units on the same day chosen on the basis of the 50-mm-per-hour criterion, and excluded otherwise. Subject to these constraints, the beginning and end of each time period was determined somewhat arbitrarily in terms of such operational conveniences as accepting an "entire rainfall" or eliminating a few bad magnetic tape records near the beginning or end of a rainfall. For example the entire three-hour period shown in Fig. 2a was accepted, and no effort was made to delete the first 20 minutes, during which there is little rain. On the whole, we were generous in the

inclusion of data. This selection procedure reduced the data base to about 110 hours of "interesting" rain.

Attention was next directed to removing spurious values from the data. It is clearly impractical to attempt a point-by-point examination of so large a body of data ($3.4 \times 10^6$ data points remained after the selection process), so various anomalies in the data were identified and studied with the aid of three additional presentations of information on the distribution of rain rate.

The most primitive displays of distributional information consisted of histograms giving the frequency of occurrence of rain rates for each station for each rainfall. Figure 3 contains four examples of the histograms from the rainfall of July 28, 1967, comprising 3840 readings per station, or about 11 hours of rain. The ordinate is the number of occurrences in 2 mm per hour intervals plotted on a logarithmic scale so that the frequencies of occurrence of the higher rain rates are visible in the diagram. Note the very substantial differences in the frequency of rain rates greater than 50 mm per hour at the four stations.

The frequency distributions were also displayed as probability plots of the rain-rate observations (ordinate) against quantiles of the standard normal distribution (abscissa).[5] (In this presentation, a normal distribution would plot as a straight line with slope equal to the standard deviation.) Figure 4 contains four examples of the distributions of data pooled for an entire rainfall and two examples of the distribution of rain rates from individual stations. Adjacent curves are translated vertically by 100 mm per hour to increase the clarity of the plot. Curves $a$ and $b$ are from the rainfall of July 28, 1967 pooled over all operational stations. The data base for curve $a$ contains 314,021 observations; that of curve $b$ is smaller by 22 observations identified as outliers. Curves $c$ and $d$ are from the rainfall of July 11, 1967. Curve $c$ is based on 285,098 observations pooled over all stations, curve $d$ on the 281,869 observations that remained after deletion of the data from a "malfunctioning" station. Curves $e$ and $f$ are based on the 3840-observation sets collected during the rainfall of July 28, 1967 from stations 10 and 84, respectively. The distributions are far from normal, having either an overwhelming surfeit of low (less than 40 mm per hour) values or too long a high-rain-rate tail, depending on the point of view. The presentation has been retained because: ($i$) of the large number of analytic distributions with fixed parameters that we have tried, none provides a good fit to even a minority of the empirical distributions; ($ii$) the rain rates higher than 40 mm per hour, in

Fig. 3 — Histograms of the frequency of occurrence of rain rates in 2 mm per hour intervals for (a) station 1, (b) station 10, (c) station 27, and (d) station 84.

Fig. 4 — Plots of empirical cumulative probabilities.

which we are particularly interested, are approximately linear on the normal plots (implying that these observations may be considered to be part of separate subpopulations); and (*iii*) the normal plots are readily produced and provide a convenient standard for comparisons among the empirical distributions.

Percentile information was also extracted from the frequency distributions and used to make pseudogeographic percentile plots. Figure 5 displays four of these. Each set of five dots connected by a

Fig. 5 — Pseudogeographic percentile plots for rainfall of (a) July 10, (b) July 11, (c) July 28, and (d) November 17, all in 1967. Individual stations may be identified by the correspondence between sets of ten stations and the rows of Fig. 1. A single dot on the abscissa indicates an inoperative (or nonexistent) station.

vertical line represents the five points of the empirical cumulative rain-rate distribution corresponding to 99.9, 99.5, 99.0, 98.0, and 97.0 percent in descending vertical order. Each block along the abscissa corresponds to a row of stations on the grid of Fig. 1. Thus these plots display the geographic distribution of the high-rain-rate tail of the rain-rate distribution.

The first step in the data screening process was to remove data from malfunctioning stations. The term station covers all the equipment for measuring, telemetering, and recording data associated with a single rain gauge. Malfunctioning stations were identified by looking at the isometric plots for stations behaving differently from the surrounding stations, and at the pseudogeographic percentile plots for stations having distributional tails that appeared peculiar in the context of distributions of the surrounding stations. When a station was judged to have been malfunctioning, it was treated as though it

Fig. 6 — Partial isometric plots of rainfalls for (a) July 11, 1967 and (b) November 17, 1967. The positions of the base lines of the traces correspond to the geographic positions of the stations, Fig. 1. Traces are six-reading averages.

had been inoperative throughout that rainfall; no more refined attempt was made to separate the spurious from the valid data.

Figure 5b shows that the tail of the distribution of data from station 67 during the rainfall of July 11, 1967, is much longer than the tails of the distributions recorded by the other stations. Part of the corresponding isometric plot, Fig. 6a, shows a rain rate in excess of 400 mm per hour at station 67 starting at minute 410 and lasting for nearly an hour, while the surrounding stations show no more than a light drizzle. Such a rainfall event is very unlikely in New Jersey and we discarded the data from station 67 for this storm. However the decision is not always so easy.

An earlier portion of the same storm is shown enlarged (two-reading-averages, rather than six-reading averages) in Fig. 2b. The behavior of station 67 after minute 25 is peculiar, but not outstandingly so, and without the subsequent evidence this data would have been retained. Another example of anomalous data is the behavior of station 29 during the storm of November 17, 1967. Figures 6b and 5d show this. The high-rain-rate tail of the distribution from station 29 is certainly very different from that of any other station. However, the isometric plot shows that the burst of intense rainfall (∼ 400 mm per hour) lasted less than five minutes. As station 29 is on the edge of the network and there are only two stations which can be classed as close neighbors, there is little context in which to form a judgment of the validity of the data. In this case the data were retained.

These data-screening decisions may have substantial effects on the high-rain-rate tail of the distributions. Curves $c$ and $d$ in Fig. 4 are probability plots showing the distribution of the data from the rainfall of July 11, 1967 with and without the data from station 67, respectively. The effect on the shape of the rain-rate distribution for high rain rates is very evident. The probability of observing a rain rate greater than 100 mm per hour is changed by a factor of about 3; this factor increases rapidly at still higher rain rates (a factor of more than ten at 200 mm per hour). These results are typical for malfunctioning stations. Substantially all the retained data for rain rates greater than 350 mm per hour comes from station 29 during the November 17 rainfall, so the shape of the distribution in this region depends entirely upon this event.

After the malfunctioning stations had been disregarded, a number of individual bad data points were removed from the retained stations. Some of these data points are definitely not connected with the rainfall measurements. For example, the source of the numerous spikes

at time 530 on Fig. 2a is a faulty magnetic tape record. In other cases, such as the spikes at time 425 on station 36, 405 on station 45, and 370 and 435 on station 46, one cannot identify the extraneous sources. Outliers have been eliminated by examining the histograms and deleting any isolated points more than 50 mm per hour above what would otherwise be the last point in the tail of the frequency distribution. Such a point appears at 256 mm per hour in Fig. 3b. This point is also shown as the highest point on curve *e* in Fig. 4 (just under curve *d*) and is obviously inconsistent with the rest of the distribution.

Although the consequences of discarding outliers are less drastic than those of disregarding stations, the effect on the shape of the tail of the distribution is significant. Curves *a* and *b* in Fig. 4 show two probability plots of the data from the rainfall of July 28, 1968 pooled over all the operating stations. Curve *a* includes the outliers; curve *b* is without the 22 points identified as outliers by examining the station histograms individually. The five highest readings were among those judged outliers, and their deletion changes the shape of the distributions for rain rates greater than 250 mm per hour substantially. Extrapolations to rain rates above 300 mm per hour would give very different results in the two cases. The fate of spurious data in the more common rain rates is of less concern. They are overwhelmed by the valid readings and have little effect on the observed distribution.

In retrospect it is clear that the magnitude of the interference problem was not fully appreciated by the designers of the telemetry system. When one is interested in infrequent events (a probability of $10^{-6}$ corresponds to 30 seconds per year) cleanliness of the raw data is very important. We emphasize this point for future experimenters.

## 2.3 *The Retained Data*

The selection and screening processes leave a body of data (the "retained data") containing 3,418,623 observations at 96 grid locations on about 110 hours of rainfall. The observations are distributed roughly as in Table I.

These data were collected over 27 rainfalls of which nine may be classed as heavy (that is, at least 12 gauges showed rain-rate readings greater than 125 mm per hour). Parts a, c, and d of Fig. 5 show heavy rainfalls; part b shows a light rainfall.

The question of what population of rain rates the data represent is more fully discussed in Section 3.6. For the present, notice that the selection of the retained data is operationally defined and does not lead to a sample whose relationship to the sampled population

TABLE I—DISTRIBUTION OF RAIN RATES IN THE RETAINED DATA

| Rain rate (millimeters per hour) | min. max. | 0 30 | 30 50 | 50 100 | 100 150 | 150 200 | 200 250 | 250 300 | 300 350 | 350 400 | 400 ∞ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Number of observations | | 3,280,000 | 90,600 | 35,300 | 9440 | 3600 | 948 | 124 | 32 | 12 | 12 |
| Percent of observations | | 95.9 | 2.7 | 1.0 | 0.3 | 0.1 | 0.03 | 0.004 | 0.001 | 0.0003 | 0.0003 |

may be precisely defined in statistical terms. Moreover, the selection procedure severely decimates the low-rain-rate part of the rain-rate distribution. This is intentional because the experiment is specifically oriented toward measurements of rain rates high enough to interfere seriously with the transmission of microwave radio signals, and the rain gauges are not designed to provide reliable readings for rain rates below about 10 mm per hour (in fact, once wet the rain gauges record some low rain rate for long periods of time).[6]

The retained data may be used to calculate empirical probabilities regarding the occurrence of various events given the condition that the rain rate at some station is greater than 50 mm per hour. Viewed this way, the sample contains only 50,000 points, a number which must be considered small in the context of the enormous variations observed in rainfall. The top 1100 of these (rain rate greater than 200 mm per hour) are in addition particularly subject to the screening problems discussed in Section 2.2. Some evidence presented in Section VII suggests that the distribution of rain rates above 50 mm per hour derived from the retained data is not atypical of summer and fall rains in littoral northern New Jersey. Broader interpretations of the data should be viewed with appropriate caution.

III. GENERAL CHARACTERISTICS OF OBSERVED RAIN

The outstanding characteristics of the rainfalls observed are the variety of behavior, both within and among individual rainfalls, and the extremely rapid fluctuations in the measured rain rates. First we discuss a possible demographic model for rain rate and the poolings of the data necessitated by the variability. Then we discuss sampling rates at one station; typical spatial, temporal, and distributional behavior; and a systematic difference among the observations from different stations. Finally we consider the matter of the relationship between this sample and the general population of rain rates.

3.1 *A Demographic Model*

It is convenient to have a demographic model in terms of which the measurements can be discussed. The very elementary phenomenology offered in the following paragraph is only a convenient hypothesis which should not be regarded as being confirmed in any sense by the data, some of which is discussed later in this section. (Of course, the model would not be presented if it were contradicted.)

A large variety of populations of rain rates is produced by differ-

Fig. 7 — Typical section of the time series from station 33 for the rainfall of July 25, 1967. The origin of the time axis has been arbitrarily set to zero.

ent, very local meteorological conditions. Each population has a characteristic distribution. A rain gauge records a sample of rain rates from the populations associated with those rain clouds which happen to pass over it. Thus a rain gauge observes a composite distribution which is the result of a two-stage sampling process, the first stage of which is the "natural" selection of the cloud, the second, the discrete scanning interval. Many rain-rate populations are inextricably intermingled as the data are recorded; indeed, as it may rain simultaneously from several strata of clouds, even a single rain-rate measurement may represent a mixture of populations.

Our extensive attempts to classify the heavy rain rates by distributional characteristics have been defeated by the large variety of distributions observed, and so subsets of the data have been pooled in various ways in attempts to synthesize some typical mixtures of populations. The unit of pooling is the data from one station during one rainfall. Some pools are: the data from each station pooled over all rainfalls for which the station was operational; the data from each rainfall pooled over all operational stations; the heavy-rainstorm pool, which contains all the data from the nine heavy rainfalls; and the grand pool which contains all $3.4 \times 10^6$ retained observations.

## 3.2 Behavior at Stations

Figure 7 shows two time series obtained from station 33 during the rainfall of July 25, 1967. The thin trace connects successive readings from the every-one-second (fast) sample. The thick trace represents the every-ten-second (slow) sample and is drafted as if the rain rate remained at the sampled value for ten seconds. (The equipment design prevents the every-ten-second and every-one-second readings from ever being coincident in time.)

The plot shows two notable features: the rain rate changes very rapidly from second to second, on occasion by a factor of three; and the magnitude of the fluctuations varies with the rain rate. These large, rapid fluctuations are characteristic of rain rates measured on a short (less than one-second) time scale with a small-area gauge.[7] It is clear that the every-ten-second frequency of sampling the stations is not nearly fast enough to provide an accurate time-series representation of the rain rate at any gauge, but instead acts as a low-pass filter.

The rain-rate distributions of the slow and fast samples for an hour of the rainfall of July 25 are plotted against quantiles of the normal distribution in Fig. 8. Curve $a$, the slow sample, is based on 360 observations; curve $b$, the fast sample, is based on 3600 observations. The curves are nearly identical below the $+2$ quantile, indicating that the slow sample may be regarded as an unbiased sample of the fast sample and the distribution can be expected to be the same for large enough samples. The last three points in the upper tail of the slow sample affect the shape of the distribution above the second quantile significantly. One of these points is shown at time 370 seconds on Fig. 7. It is almost surely spurious, as are the other two outliers; the distributional effect of these three points should be discounted.



Fig. 8 — Normal plots of the fast and slow samples from Station 33 for one hour of the rainfall of July 25, 1967.

Another characteristic of rainfall, which is apparent in Fig. 2, is that heavy rain comes in bursts. Although some of the showers last for almost 30 minutes, longer showers usually contain a number of short bursts. In general, these periods of rain cannot (at our sampling rate) be considered piecewise-stationary time series. Our many attempts to analyze periods during which the rain rate was greater than 35 mm per hour yielded results badly confounded by the nonstationarity produced by occurrences of these bursts. Thus one cannot naively apply time-series analysis techniques to the heavy rain rates; such common statistics as correlation and autocorrelation coefficients, and power spectra are not directly meaningful.

### 3.3 An Overview of the Grid

At about time 400 on July 28, 1967 it began to rain heavily on two separate portions of the grid, the southeastern portion and the north central portion (Fig. 2a). The southeastern rain built up slowly (in general) for about half an hour, decreased abruptly, and then resumed at lower intensities and intermittent intervals for another two hours. The north central rain was heavy for about ten minutes and was followed by 40 minutes of light rain. At about time 460, very heavy rain started in the northwest portion of the grid and continued for between 20 minutes and an hour before dying out. The general behavior of this rainfall is typical. It rains heavily first on one part of the grid then on another, the regions of heavy rainfall are fairly local (often only a few stations register heavy rain), and it seldom rains heavily for as long as an hour. The downpour is not continuous but has substantial variations in intensity throughout its lifetime.

In view of the rapid fluctuations of rain rate at the stations, it is not surprising that the correlation among the fine structures observed at different stations is poor. It would require an extraordinary mechanism to synchronize one-second rain-rate fluctuations over an area of several square kilometers. The temporal relationship among rain rates at the different stations may be seen on Fig. 2b. Even allowing for the fact that the stations are not sampled simultaneously (see Section 2.1) and that there may be a propagation lag, the structure of the rain-rate traces is not especially similar from station to station.

Of the 27 rainfalls, that of July 11 is one of the best examples of systematic motion of a rainstorm. The traces in Fig. 2b behave as though a rain cloud roughly 3 km wide and many kilometers long, with the long axis oriented in the WNW-ESE direction, moved NNE across the grid with a velocity of about 15 km per hour.

3.4 *Distributional Characteristics*

In discussing the distributional characteristics of rain rates, our attention is focused on rain rates greater than 50 mm per hour. Over our sample this is the upper $1\frac{1}{2}$ percent of the empirical distribution, although for particular stations and rainfalls the amount of data above 50 mm per hour varies from zero to five percent.

The purpose of this subsection is to demonstrate that there seem to be many different kinds of rainfall, so the data contain observations from many different populations with different rain-rate distributions and mixture ratios. Many common analytic distributions (the normal, log normal, gamma family, and so on) have been tried in attempts to find a simple description of the observed composite distributions which behaves reasonably in the tail region (rain rate $\geq$ 50 mm per hour); but none have been found that can serve even approximately.

We conclude that the large sampling variation and the complexity of the mixtures precludes obtaining reliable estimates of distributional parameters of rain rates above 50 mm per hour from any small sample. Some samples we consider small in this context are single stations for a season and the entire network for a single intense rainfall. For still higher rain rates the present overall sample may be inadequate (there are only 1100 observations for rain rates above 200 mm per hour). In the absence of a distribution on which to base precise estimates of formal statistics, such as confidence limits, discussion of these matters can best await the analysis of the 1968 data.

The remainder of this subsection indicates some of the distributional variety observed. The arrival of heavy rainfall in bursts, as indicated in Fig. 2, makes it obvious that the distribution of rain rates is very different for different segments of the same rainfall measured at the same station. The same holds true for the distributions observed at different stations for the same rainfall. This may be seen from the two probability plots for stations 10 and 84 (curves *e* and *f* in Fig. 4), from the differences in the tails of the distributions shown in Fig. 5, and the four histograms in Fig. 3. The rain-rate distributions differ greatly among stations even for the rainfall of July 11, whose active portion, shown in Fig. 2a, appears deceptively similar for many of the traces. Apparently many populations of rain rates coexist within rainstorms.

In an attempt to find a common mixture of populations, the stations have been pooled within rainfalls. Curves *b* and *d* in Fig. 4 are examples of such pooled distributions, and the difference between them is typical of what is observed. A glance at the percentile plots of Fig.

5 indicates that it is unlikely that pooling across these collections of stations can produce distributions with similar shapes above 50 mm per hour. The next pooling is within stations, across rainfalls, that is, at each station for the whole season. Naively this seems to be the procedure most likely to produce a "typical" sample and thus similar distributions. Cloudbursts, however, are limited in extent. Some of the stations were never hit by cloudbursts, whereas others appear to be deluged quite often. Figure 9a shows normal probability plots of the rain rates from each of four stations pooled for the season, and Fig. 9b shows the matching upper five percent of the empirical cumulative distributions. Adjacent curves have been translated vertically by 100 mm per hour. Curves $aa$ and $ba$ are for station 29, $ab$ and $bb$ for station 1, $ac$ and $bc$ for station 98, and $ad$ and $bd$ for station 76. The data base for each station contains about 39,000 observations. The distributions are seen to be quite different.

The overall pseudogeographic percentile plot, Fig. 10, shows the



Fig. 9 — Distributions of the rain rates at each of four stations pooled for the 1967 season: (a) plots against quantiles of the standard normal distribution; (b) the upper five percent of the empirical cumulative distribution.

Fig. 10 — Pseudogeographic percentile plot. Individual stations may be identified by the correspondence between sets of ten stations and the rows of Fig. 1. The 99.9, 99.5, 99.0, 98.0 and 97.0 rain-rate percentiles of the empirical cumulative distributions are indicated for each station by the dots in decreasing vertical order. A single dot on the abscissa indicates a grid square that has no station.

upper tails of the distribution of rain rates from each station pooled across all the rainfalls (for which the station was operational) in the final sample. This plot covers the upper tail down to a rain rate of 50 mm per hour for most of the stations. The figure indicates the variety of distributions, and shows that the four samples selected for Fig. 9 are not atypical. It also confirms the notion that data from a single station over a season are not a sufficient sample on which to base rain-rate statistics for rain rates above 50 mm per hour.

3.5 *Systematic Behavior*

Figure 10 gives an idea of the areal distribution of the intense rain rates during the 1967 season. It rained most intensely on the north-western part of the network and least intensely in the northeastern and southerly central portions. Intense rain was also observed by the stations in the southeastern section of the grid. A more detailed study

of the data reveal that station 1 always records more intense rain than most of the network and station 76 always records rain of lower intensity than most stations. Three possible explanations are confounded in the present data: first, there may be a systematic geographic effect; second, the effect may be inherent in the equipment of the stations (the calibrations of the gauges may have shifted); and third, one may be witnessing a sampling fluctuation. Action has been initiated to determine whether the effect is connected with the equipment. Analysis of data from 1968 may help to distinguish between the first and third possibilities.

### 3.6 *What Does the Sample Represent?*

The sample contains virtually all of the available data for rain rates greater than 50 mm per hour and about a quarter of the data below 30 mm per hour. As such it can be expected to be representative (with appropriate proportionality adjustments) of rainfall with rain rates greater than 10 mm per hour within a small area of northern littoral New Jersey over a relatively short period. Neither the year nor the locality are known to be exceptional in any respect.

The distributional stability of the observations above 50 mm per hour is poor, as has been indicated, which shows that the sample is small in the context of the demographic composition. In Section VII, the data are compared with the 1958 data taken at Island Beach, New Jersey by Mueller and Sims and found to be reasonably similar.[8] Some informal indication of the dispersion of various statistics is given in the following sections by noting the quartiles and interquartile ranges of the distributions. (The quartiles are the 25 and 75 percent points of the distribution; the interquartile range is the difference between them.) While there is nothing to indicate that the results of this study are in any way atypical of heavy rainfall in this part of New Jersey, it is unlikely that they permit an accurate assessment of the extremes of rainfall that may not infrequently be encountered.

Extension of these results outside the immediate locality will, of course, be subject to even greater uncertainties.

### IV. POINT-RAIN-RATE STATISTICS

### 4.1 *The Grand Pool*

Any interpretation of this section must take into account the fact that the retained data include all of the available data for rain rates

greater than 50 mm per hour but only about 25 percent of the available data for rain rates below 30 mm per hour.

Figure 11 is the grand histogram based on the 3,418,623 observations (about 110 hours of rain) retained from 1967. Log frequency of occurrence is plotted linearly on the ordinate in order to accommodate the range of frequencies. The interval of the histogram is 2 mm per hour. Although only about 25 percent of the data below 30 mm per hour were selected, the relative frequencies in this range are unlikely to have been substantially affected by this selection. However, the general tendency of the rain gauges to give biased readings which greatly overestimate the rain rates below about 15 mm per hour make this histogram unsuitable for hydrological use.

The empirical cumulative distribution is shown in Fig. 12a, and the upper four percent of the distribution is shown on an expanded scale in Fig. 12b. The shape of these curves may be adjusted to approximately represent the distribution of the entire 430 hours of rain recorded by scaling the probabilities below 0.96 by the factor 0.99/0.96. Above 0.96 the new *probability*, $P_N$ , would be given approximately by the expression $P_N = 0.75 + 0.25 P_B$ , where $P_B$ is the *probability* on the plot. The cumulative distribution points up the small fraction of the data which lies above about 50 mm per hour. The empirical probabilities, Fig. 12a, are transformed to quantiles of the standard normal distribution and the rain rates as ordinates are plotted against the quantiles in Fig. 12c. (Only the data above the mean, quantile 0, rain rate about 8 mm per hour, are shown.) If the data were a random sample from a normal distribution, the points would lie along a straight line. The data are not even approximately normally distributed, which must of course be the case as the frequency distribution (Fig. 11) decreases monotonically from the lowest interval.

As the distribution has such a long tail and is skew (no negative rain rates have been recorded although all that water must have gotten up there somehow), it seems reasonable to look at various long-tailed skew distributions such as the gamma, Wiebull, and extreme value. These distributions yield probability plots which do have an appreciably better appearance than the normal probability plot of Fig. 12c; however, they all show substantial deviations from linear behavior near rain rates of 100 mm per hour. Efforts to improve the linearity by adjusting the proportion of the data below rain rates of 30 mm per hour were not successful. Since a distributional description in such circumscribed conditions is of only very limited usefulness the details of these efforts are not reported here. Typical of the results is

Fig. 11 — The grand histogram. There are 12 additional observations scattered above 400 mm per hour.

Fig. 12 — Distribution of the grand pool: (a) cumulative distribution, (b) enlargement of the upper tail of cumulative distribution, (c) normal probability plot, and (d) log-normal probability plot. Only the part of the distribution above the mean is shown in parts (c) and (d).

Fig. 12d, which is a plot of log rain rate vs. quantiles of the standard normal distribution (that is, a probability plot of the "log normal" distribution). As in Fig. 12c, only data above the mean are shown. Those familiar with probability plots will recognize that the inflection at the 2.5 quantile cannot be removed by adjusting the fraction of data assumed to lie above 30 mm per hour.[5]

In the present data-analytic situation, where there is no simple analytic distributional description of the data, and indeed the data may represent a drawing in unknown proportions from many populations, it is not possible to compute formal confidence limits for the various estimates of probabilities. An informal indication of the

Fig. 13 — Histograms showing the distribution of the (a) 95.0th, (b) 97.0th. (c) 98.0th, (d) 99.0th, (e) 99.5th, and (f) 99.9th percentiles from the rain-rate distributions of the 96 stations where the data from each station has been pooled over the 1967 season. Intervals are 5 mm per hour.

dispersion in the data may be obtained by examining the behavior
of the distribution from individual stations, some of which are dis-
played in Fig. 9. The behavior of all the stations is summarized in
histograms showing the frequency distribution of rain rates corre-
sponding to 95, 97, 98, 99, 99.5 and 99.9 percent for the data from each
station pooled over all the rainfalls. These histograms are presented
in Fig. 13. All 96 stations were used. While 11 stations were operational
only about half of the time and four were operational less than one-
third of the time, inspection of the distributions from some of these
stations indicates that an assumption of random loss of observations
is not unreasonable. Summary statistics comprising the value of rain
rate from the grand pool, together with mean, upper, and lower 25
percent points (quartiles), and the interquartile range (spread of the
middle 50 percent) of the histograms, are given in Table II.

The overall rain-rate percentiles are larger than the means of the
corresponding histograms. The amount by which they are larger in-
creases as the overall percent (and thus the rain rate) increases until
the largest overall percentile corresponds to the upper quartile of the
histograms. This indicates that most of the high-rain-rate data came
from only a few stations. The interquartile ranges are between 50
and 75 percent of the means of the histograms and tend to increase
with rain rate. The histograms of the 99th and higher percentiles sug-
gest a bimodal distribution, as if the data contained a component of
very high rain rates which (because it is infrequent) is never observed
at a large number of the stations. Thus Fig. 13 indicates a substantial
sampling uncertainty above the 99th percentile; this uncertainty would
be much greater for data collected at only a few stations.

While it is possible to make histograms similar to those of Fig. 13

TABLE II—STATISTICS ON THE DISTRIBUTION OF VARIOUS
PERCENTILES FROM THE RAIN-RATE DISTRIBUTION

| Percent of rain rate distribution | Value from the grand pool | Rain rates in mm per hour | | | |
|---|---|---|---|---|---|
| | | The 96 stations, each station pooled for the season | | | |
| | | Lower quartile | Upper quartile | Interquartile range | Mean |
| 95.0 | 28 | 17 | 29 | 12 | 25 |
| 97.0 | 34 | 22 | 38 | 16 | 32 |
| 98.0 | 41 | 27 | 47 | 20 | 41 |
| 99.0 | 62 | 37 | 75 | 38 | 59 |
| 99.5 | 91 | 50 | 109 | 59 | 81 |
| 99.9 | 162 | 79 | 160 | 81 | 127 |

Fig. 14 — Cumulative distribution of the grand pool subject to the condition that the rain rate is greater than 50 mm per hour. The data base contains 48,481 observations.

for other poolings of the data (say across stations for each rainfall) this has not been done because the rain-rate percentiles are affected by the dilution of the high-rain-rate data when periods of light rain are included. As a result of the selection procedure (Section II) this dilution varies greatly (by a factor of five) from rainfall to rainfall. Thus percentiles for different rainfalls are not really comparable, but percentiles from poolings across rainfalls are. This remark does not apply to the data when the condition in Section 4.2 is applied.

4.2. *Conditioned Probabilities*

All the probabilities in this section are based on the condition that the rain rate is greater than 50 mm per hour. This simple condition may be applied straight-forwardly to the point-rain-rate data because in arriving at the point-rain-rate statistics each data point is

treated as independent, and no relationships in time or space are taken into account. The 50-mm-per-hour condition is chosen because our selection procedure accepts virtually all data meeting this criterion. Thus at the cost of severely limiting the range of rain rates examined, one may work with probabilities that are well-defined and independent of the subjective criteria used to determine the beginning and end of periods of "interesting" rain.

The frequency distribution of rain rates greater than 50 mm per hour (in 2 mm per hour intervals) is given by the part of Fig. 11 which lies above this value. The corresponding empirical cumulative distribution appears in Fig. 14, which shows the full range from 0 to 1.0, as well as the upper tail on an expanded scale from 0.9 to 1.0. Once more there is no obvious overall correspondence between the data and the more common distributions. We have however fitted the function

$$N = \tfrac{1}{2}\alpha \exp\,(-\beta R), \qquad R \geqq 65 \text{ mm/hr}, \tag{1}$$

where $N$ is the number of observations in a 1 mm per hour interval, $R$ is the rain rate, and $\alpha$ and $\beta$ are fitted coefficients, to the tail of the frequency distribution of Fig. 11. The fitting procedure was weighted least squares, with an *ad hoc* weighting that emphasizes the goodness



Fig. 15 — Residuals from the fit to the tail of the frequency distribution. The interval is 1 mm per hour.

Fig. 16 — Distributions of rain rates at each of four stations pooled for the 1967 season subject to the condition that the rain rate is greater than 50 mm per hour. Adjacent curves have been translated vertically by 100 mm per hour to increase the clarity of the plot. Curves *a, b, c,* and *d* are for stations 29, 1, 98, and 76, respectively and the respective data bases contain 1384, 1817, 415, and 45 observations.

Fig. 17 — Histograms showing the distribution for (a) 50.0th, (b) 75.0th, (c) 90.0th, (d) 95.0th, and (e) 99.0th percentiles from the rate-rate distributions of the 96 stations, for rain rates greater than 50 mm per hour. Intervals are 5 mm per hour.
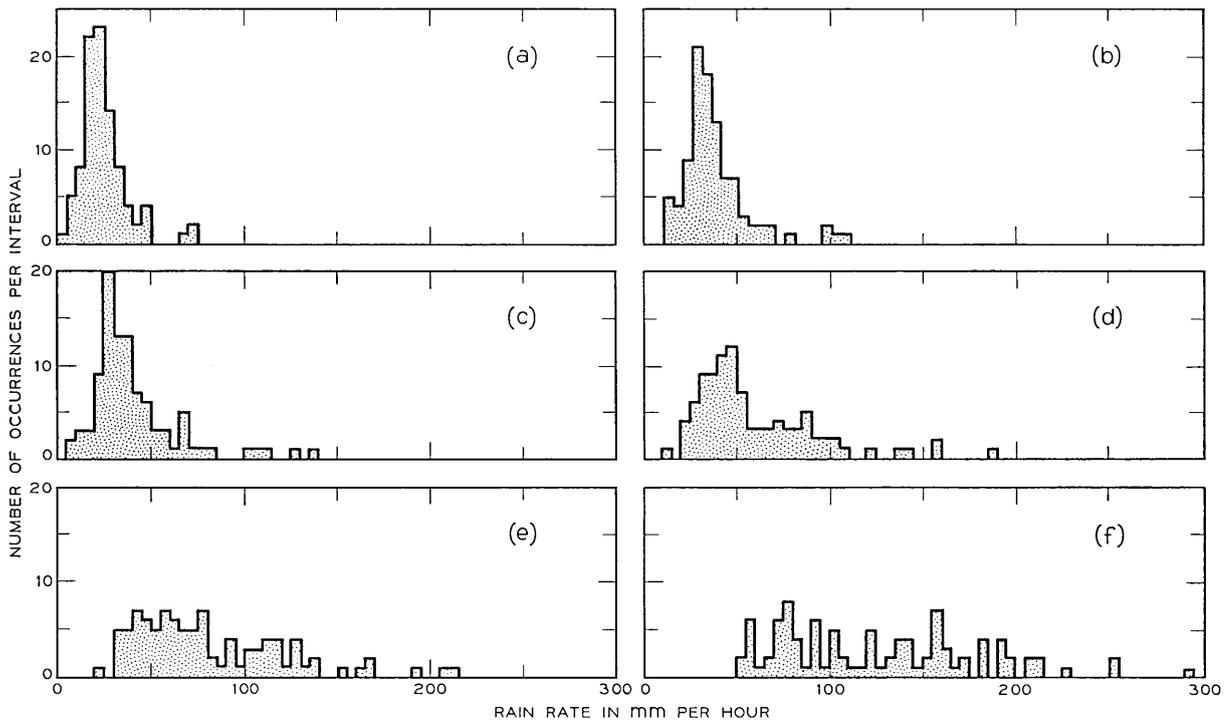
Fig. 18 — Histograms showing the distribution for (a) 50.0th, (b) 75.0th, (c) 90.0th, (d) 95.0th, and (e) 99.0th percentiles from the rain rate distributions of the 27 rainfalls, for rain rates greater than 50 mm per hour. Intervals are 5 mm per hour.

TABLE III—STATISTICS ON THE DISTRIBUTION OF VARIOUS PERCENTILES
FROM RAIN-RATE DISTRIBUTIONS*

| | | Rain rates in mm per hour | | | | | | | |
| Percent of rain rate above 50 mm per hour | Value from the grand pool | The 96 stations, each pooled for the season | | | | The 27 rainfalls, each pooled over all stations | | | |
| | | Lower quartile | Upper quartile | Inter-quartile range | Mean | Lower quartile | Upper quartile | Inter-quartile range | Mean |
|---|---|---|---|---|---|---|---|---|---|
| 50 | 74 | 65 | 78 | 13 | 73 | 64 | 77 | 13 | 71 |
| 75 | 106 | 80 | 112 | 32 | 95 | 81 | 105 | 24 | 93 |
| 90 | 148 | 97 | 152 | 55 | 125 | 100 | 139 | 39 | 121 |
| 95 | 174 | 110 | 167 | 57 | 142 | 108 | 167 | 59 | 139 |
| 99 | 222 | 132 | 204 | 72 | 174 | 141 | 208 | 67 | 179 |

* For rain rates greater than 50 mm per hour.

of fit for the higher rain rates. The estimates $a$ and $b$ of $\alpha$ and $\beta$ are 7754 and 0.02408, respectively. The properties of the model are best displayed by the residuals, which are plotted against rain rate in Fig. 15. The fit is seen to be moderately good for rain rates greater than 200 mm per hour, and fair for rain rates between 65 and 200 mm per hour. Below 65 mm per hour the function and the observations diverge rapidly.

The model provides a useful smoothing function between rain rates of 65 and 400 mm per hour, and also represents the best currently available extrapolation of these data above 400 mm per hour. However, the systematic behavior of the residuals shows that an exponential function is not really a suitable distributional description of the data.

Figure 16 shows empirical cumulative distributions for rain rates above 50 mm per hour for the same four stations whose unconditioned cumulative distributions appear in Fig. 9. This once more illustrates the large differences among stations. Histograms of the frequency with which certain percentiles fall in various rain-rate intervals are once more used to quantify the dispersion in the data. Figure 17 contains the results obtained when data are pooled across the season's rainfalls for each station; Figure 18 contains analogous information for data pooled across stations for individual rainfalls. The outliers on the two 99 percentile histograms are both caused by the same event, a very high rain rate at one station during one rainfall (see Section 2.2). Smaller poolings generally contain too few data points above 50 mm per hour to be informative. The results are summarized in Table III.

Most notable in Table III is the great similarity between the distributions of upper-tail percentiles for the two poolings. There is no obvious reason why this should be so, and the result may be peculiar to this set of data. As in Table II, the interquartile range increases with the mean (ranging in value from about 20 to about 40 percent of the mean) indicating the increasing uncertainty associated with the small sample of high rain rates. Above 90 percent, the value from the grand pool lies above all other values in the same line of the table, showing once again the influence of a few local cloudbursts.

## V. SPATIAL AND TEMPORAL RELATIONSHIPS

Figure 2 shows the general space-time relationship of rain rates: When it is raining heavily at a station it is likely to be raining heavily

at nearby stations, and it is also likely to be raining heavily a short time later. This section quantifies these general observations.

If the time series representing the rain-rate observations were stationary the relationships would be given by the cross correlation as a function of distance and the autocorrelation as a function of lag (time difference). However as indicated in Section 3.1 the time series are not even piecewise stationary in the regions containing the high rain rates. As the situation is more complex than that represented by stationary time series, the statistical summary is somewhat less succinct. The spatial relationships are described in terms of the distribution of the rain rate at a point *there* given certain conditions of rain at point *here*. The distance between *here* and *there* and the rain rate at point *here* are varied as parameters. The temporal relationships are described analogously in terms of distributions of the rain at time *hence* given certain conditions of the rain at time *now*. The lag, that is, the interval between *now* and *hence*, and the rain rate at time *now* are varied as parameters.

## 5.1 *Spatial Relationships*

The frequency distributions underlying this presentation were collected as follows. The readings from each scan of the network were assumed to be simultaneous. The very rapid variations in rain rate (Fig. 7) and lack of detailed correlation over distances of the magnitude of the station separations (Section III) provide some statistical justification for this treatment. Each station (*here*) was then paired with every other station (*there*), giving $(96)^2$ pairs. Each pair was classed as belonging to one of 238 possible histograms according to the value of the rain rate *here*, which is sorted into 17 intervals, and the nominal distance (stations are assumed to be at the center of the grid square) between the stations, which is sorted into 14 separation intervals. Results are labeled with the mean of the rain-rate interval, and the mean of the nominal distances in the separation interval. For the rain rate *here* the interval is increased as the rain rate increases; for the separation, the interval is about 20 percent of the mean. Once the histogram was selected, the frequency of the rain rate *there* was accumulated in intervals of 5 mm per hour.

The 110 hours of rain data contain about $4 \times 10^4$ scans of the network and each scan yields about $10^4$ pairs giving a total data base of $4 \times 10^8$ pairs. The data are accumulated separately for each of the 27 rainfalls and then pooled for many presentations.

In spite of the large number of pairs, there are relatively few pairs at the high rain rates. Thus the high-rain-rate results show noticeable sampling fluctuations and remain sensitive to the data-screening procedures. Probabilities involving rain rates of less than 50 mm per hour are distorted by the data-selection processes as already discussed.

Three typical histograms (from the set of 238) are shown in Fig. 19. Figures 19a and b show the smooth behavior produced by the large numbers of pairs at low and intermediate rain rates; Figure 19c shows the fluctuations that result from the relatively small numbers of pairs at high rain rates. These fluctuations are reflected in occasional irregular behavior seen in other displays in this section.

Three sets of empirical probabilities for the rainfall of July 25, 1967 are shown in Fig. 20. The curves give the probability of the rain rate *there* being less than $R$ mm per hour when the rain rate *here* falls in the interval indicated on the curve. Figure 20a shows that when it is raining lightly at station *here*, it is unlikely to be raining heavily at stations 1.3 km away. As the rain rate at station *here* increases it becomes increasingly more likely to be raining heavily at nearby stations.

The strength of this relationship decreases with increasing separation until at about 8 km separation the probability of observing a rain rate less than $R$ mm per hour is largely independent of the rain rate observed at the distant station; this is demonstrated by the contraction of curves in Fig. 20b into a narrow band. As the separation increases still further the curves reverse (that is, the probabilities associated with curve $5 \pm 5$ mm per hour are generally lower than those associated with curve $200 \pm 20$ mm per hour for the same $R$. This may be seen on Fig. 20c which shows the results for stations separated by 12 km. The reversal means that if it is raining heavily at station *here* it is less likely to be raining heavily at a station 12 km away than if it were raining lightly at station *here*. This effect is noticeable at separations up to 16 km, but few pairs of stations are so widely separated (the diagonal measurement of the network is 18 km). This behavior suggests that regions of heavy rainfall have a "range of influence" of about 16 km (twice the 8 km distance corresponding to Fig. 22b) and that the "centers" of these regions are separated by more than the dimension of the network, that is, more than about 16 km.

The probabilities may also be plotted against distance, with the rain-rate-*here* intervals as parameters and different values of $R$ appearing on different plots. Figures 21a shows the behavior of the prob-

Fig. 19 — Histograms showing the distribution of the rain rate *there* for three separation, rain-rate-*here* combinations pooled over all rainfalls. (a) 6.3 ± 0.6 km, 45 ± 5 mm per hour; (b) 14.1 ± 1.4 km, 155 ± 15 mm per hour; (c) 6.3 ± 0.6 km, 300 ± 20 mm per hour. The interval of the histogram is 5 mm per hour.

Fig. 20 — Relative empirical probability that the rain rate *there* is less than $R$ mm per hour for various intervals of the rain rate *here* during the rainfall of July 25, 1967, for stations separated by (a) 1.3 ± 0.2, (b) 7.5 ± 0.8, and (c) 12.3 ± 1.2 km, respectively. Points shown on this and following plots indicate actual calculated values.

Rain Rate *Here* (mm per hour)

|        | A  | B  | C  | D  | E   | F   | G   |
|--------|----|----|----|----|-----|-----|-----|
| Min.   | 0  | 20 | 40 | 70 | 110 | 170 | 200 |
| Max.   | 10 | 30 | 50 | 90 | 140 | 200 | 240 |

ability that the rain rate *there* is less than 35 mm per hour (the threshold) as a function of distance for various intervals of the rain rate *here*. The probability that the rain rate *there* is less than 35 mm per hour increases monotonically with distance when the rain rate *here* is greater than 30 mm per hour. For rain rates *here* that are less

than 30 mm per hour, the probability for rain rate *there* being less than 35 mm per hour goes through a minimum at a separation that depends on the rain rate *here*. Figure 21a also displays the crossover of the probability curves at a separation of about 8 km. For smaller separation it is less probable for the rain rate *there* to be below 35



Fig. 21 — Relative empirical probability that the rain rate *there* is less than R as a function of separation with various intervals of the rain rate *here* as a parameter, for the rainfall of July 25, 1967. R is (a) 35, (b) 60, and (c) 125 mm per hour.

Rain Rate *Here* (mm per hour)

|      | A  | B  | C  | D  | E   | F   | G   | H   |
|------|----|----|----|----|-----|-----|-----|-----|
| Min. | 0  | 20 | 40 | 70 | 110 | 170 | 200 | 240 |
| Max. | 10 | 30 | 50 | 90 | 140 | 200 | 240 | 280 |

mm per hour if it is raining heavily *here* than if it is raining lightly *here*. At separations exceeding 8 km, the inverse is true. Figures 21b and 21c show the analogous plots for rain rates less than thresholds of 60 mm per hour and 125 mm per hour, respectively. The pattern is the same, and if 70 and 125 mm per hour respectively are substituted for 30 mm per hour in the previous discussion it remains applicable. In particular the crossover still occurs at 8 km separation. As might be expected, the probability of finding the rain rate *there* lower than the threshold increases with increasing threshold.

For completeness we include four plots in Fig. 22, which show the behavior with separation of the probability that the rain rate *there*



Fig. 22 — Relative empirical probability that the rain rate *there* is as indicated on the curve as a function of separation for various intervals of the rain rate *here* for the rainfall of July, 25, 1967. The intervals of rain rate *here* are (a) 20 to 30, (b) 50 to 70, (c) 90 to 110, and (d) 200 to 240 mm per hour. Curves are labeled in units of millimeters per hour.

is less than certain values for the rainfall of July 25, 1967. Each plot is for a given interval of the rain rate *here*. The behavior of the curves is generally smooth without any abrupt changes of slope. The direction of curvature and position of the maxima and minima depend on the particular interval of rain rate *here* displayed. We are unable at this time to evaluate the significance of the fine structure that appears occasionally.

The patterns of the other heavy rainfalls have been examined and are similar to that of the July 25, 1967 rainfall. In particular the "range of influence" of the rainfalls is approximately 16 to 20 km. When the data from all 27 rainfalls (or even for the nine heavy rainfalls) are aggregated, the differences among the rainfalls conceal the individual patterns, and the reversal of the order of the curves that take place between Figs. 20a and 20c disappears. However, the concept of a "range of influence," that is, a separation at which the rain rate *there* is relatively independent of the rain rate *here,* is still valid for the aggregate when the rain rate *there* is below about 125 mm per hour. The "range of influence" is about 20 km.

Of particular engineering interest is the probability that the rain rate is *greater* than a particular value at two stations simultaneously. Results in this form are presented in Fig. 23, which is a plot of the joint probability that the rain rate is greater than $R$ mm per hour at both of two stations against the separation in km. These probabilities change smoothly with distance, and have a minimum at a separation of about 12 km for rain rates with sufficient data to establish the shape of the curves. The joint probabilities were also computed for each of the 27 rainfalls separately; the upper quartiles of the 27-rainfall distributions (a distribution for each separation) are shown for minimum rain rates of 50 and 110 mm per hour by the crosses in Figs. 23a and b, respectively. The lower quartiles are indistinguishable from the abscissa on the scale of the plots. The greater-than-50-mm-per-hour curve lies close to its upper quartiles; the greater-than-110-mm-per-hour curve lies far above the corresponding upper quartiles. These results again demonstrate the sensitivity of the upper tail of the aggregated distribution to a small number of "cloudbursts." Although the plots for individual rainfalls differ substantially in behavior and statistical stability, about one third of the rainfalls display a minimum in the joint probabilities at a separation of about 12 km.

Because of the bias introduced by the method of selecting the data, the numerical values of joint probabilities for rain rates below 50 mm per hour are distorted relative to those for rain rates above that value.

Fig. 23 — Relative empirical joint probability that the rain rate at both stations exceeds the value with which the curve is labeled versus separation in kilometers. Crosses in (a) and (b) indicate the upper quartiles of the distributions of joint probabilities that the rain rate is greater than 50 and 110 mm per hour, respectively, when the 27 rainfalls are considered separately.

Minimum Rain Rate (mm per hour)

| A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|
| 30 | 40 | 50 | 70 | 90 | 110 | 140 | 170 | 200 | 240 |

## 5.2 Temporal Behavior

The frequency distributions on which descriptions of the temporal behavior are based are accumulated in a manner completely analogous to that used to collect the spatial distributions. Instead of distance intervals, time intervals (lags) are used, and each station is "paired"

only with itself. Thirteen lags were used and there are about 100 stations and 40,000 scans. This leads to $5 \times 10^7$ pairs which (as 17 intervals of rain rate *now* are used) fall within $13 \times 17 = 221$ histograms. The value of the rain rate *hence* is sorted into intervals of 5 mm per hour to form the histograms. All the data presented in this section are pooled across the 27 rainfalls.

The probabilty that the rain rate *hence* is less than $R$ mm per hour for various intervals of the rain rate *now* is given by the curves in Fig. 24. While the curves shift toward lower rain rates at the longer lags, the pattern remains essentially the same with time. If it is raining heavily *now* it is more likely to be raining heavily a short time *hence* than if the rain *now* were light. Curves G and H, which show rain rate *now* greater than 240 mm per hour are much more variable than the others because of the small number of observations in this range.

Figure 25 shows curves giving the empirical probability that the rain rate *hence* will be less than the value indicated on the curve versus lag for rain rates *now* in the ranges 20 to 30, 50 to 70, 90 to 110, and 200 to 240 mm per hour, respectively. These plots show that the probability that the rain rate will fall at or below its present value is about 0.7 (increasing to 0.9 at the highest rain rate) and relatively independent of the lag (the increase is about 0.1 over the 360 seconds).

The tendency of the curves to move in the direction of the line corresponding to a probability of 0.8 is noticeable at the small lags. This means that the probability, that the rain rates a short time hence are very different from the current rain rate, is small. The relationship has both short and long term components. The short term component dies down within about five minutes, as evidenced by the decrease in the slopes of the curves with increasing lag, and is probably associated with the fine structure of the precipitation. The long term component is present after five minutes as evidenced by the different values for the same curves on the four plots. This component is probably associated with some average property of the local precipitation. Except for the highest rain rates (Fig. 25d), the spacing between the curves after a 360-second lag is larger for the lower rain rate *hence* limits (in spite of the fact that the rain rate interval is also smaller for the lower rain rate *hence* limits), indicating that the probability density for the rain rate *hence* peaks at the low rain rates after a lag of 360 seconds for values of the rain rate *now* below about 200 mm per hour.

An interesting feature appears in Fig. 25d. All the curves with rain

rate *hence* limits of less than 240 mm per hour show a dip at a lag of 20 seconds, and the curve labeled 240 mm per hour shows a double dip. This indicates a tendency for the heavy rain rates to have a fine structure oscillation with a period of about 20 seconds. The conjecture is confirmed by examination of plots (not presented here) for rain rate *now* intervals with lower bounds greater than 200 mm per hour.



Fig. 24 — Relative empirical probability that the rain rate *hence* is less than *R* mm per hour for various intervals of the rain rate *now* for lags of (a) 10, (b) 120, and (c) 360 seconds.

Rain rate (mm per hour)

|       | A  | B  | C  | D  | E   | F   | G   | H   |
|-------|----|----|----|----|-----|-----|-----|-----|
| Min.  | 0  | 20 | 40 | 70 | 110 | 170 | 240 | 290 |
| Max.  | 10 | 30 | 50 | 90 | 140 | 200 | 280 | 320 |

Fig. 25 — Relative empirical probability that the rain rate *hence* is as indicated on the curve as a function of lag for intervals of the rain rate *now* of (a) 20 to 30, (b) 50 to 70, (c) 90 to 110, and (d) 200 to 240 mm per hour. Curves are labeled in units of mm per hour.

Another presentation, Fig. 26, plots the empirical probability that the rain rate *hence* is less than $R$ mm per hour (the threshold) for various intervals of the rain rate *now* as a parameter. In agreement with the previous plots, one notes that the probability that the rain rate exceeds its present value decreases as both lag and rain rate *now* increase and is smaller than about 0.2. There is a tendency for the rain-rate intervals below the threshold (for example, the 40 to 50 mm per hour curve in Fig. 26b) to be slightly concave upward, with a minimum at about 180 seconds. This might be associated with a characteristic growth time for showers; but the averaged effect is very weak.

Lastly, the empirical probability that the rain rate exceeds a given

value both now and after various lags is given in Fig. 27 for a number
of rain rates. The probabilities are seen to decrease slowly and
monotonically with time up to lags of 360 seconds. The crosses in
Fig. 27 are analogous to those in Fig. 23 and indicate the upper quar-
tiles of the distributions obtained (a distribution for each lag) when
the joint probabilities of minimum rain rates of 50 and 110 mm per
hour are calculated separately for each of the 27 rainfalls. (The lower
quartiles are too close to zero to show on the plots.) That the joint
probabilities in time are also very sensitive to small numbers of heavy-
rain-rate events is demonstrated by the fact that the mean value is
higher than the crosses for a large fraction of the data. However the



Fig. 26 — Relative empirical probability that the rain rate *hence* is less than
*R* mm per hour as a function of lag with intervals of the rain rate *now* as param-
eter. *R* is (a) 35, (b) 60, (c) 125, and (d) 220 mm per hour. Curves are labeled
in units of mm per hour.
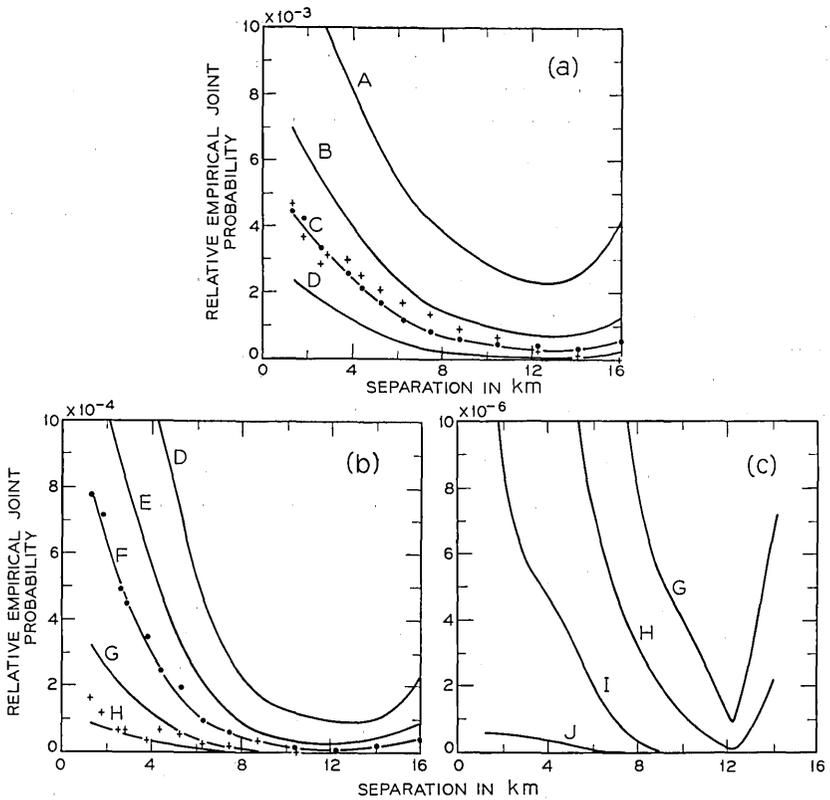
Fig. 27 — Relative empirical joint probability that the rain rate exceeds the value with which the curve is labeled both *now* and *hence* (but not necessarily between the two times).

Minimum Rain Rate (mm per hour)

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| 20 | 30 | 40 | 50 | 70 | 90 | 110 | 140 |

| I | J | K | L | M | N | O | |
|---|---|---|---|---|---|---|---|
| 170 | 200 | 240 | 280 | 320 | 360 | 400 | |

situation is less sensitive than in the case of joint probabilities in space (Fig. 23).

This presentation does not provide direct information on the probability that the rain rate will *remain* above a certain value for the time given by the abscissa. We have not attempted to provide sta-

tistics on the length of time the rain rate exceeds given values because
(as shown in Fig. 7) the rain rate changes rapidly within the ten
second sampling interval and the form of distribution of which the
rain rates within the interval are a sample is not known.

The remark on sampling bias at the end of Section 5.1 applies to the
temporal as well as the spatial probabilities.

VI. PATH STATISTICS

The average rain rate along the path between a microwave trans-
mitter and receiver affects the amplitude of the received signal. This
section presents the probability that the average rain rate exceeds a
given value on a single path or on two paths simultaneously for a
variety of conditions. We calculate the probabilities directly from the
data without recourse to any of the statistics and intermediate re-
sults previously presented. Because of the very large amount of
computation involved, it is essential to reduce the data base and regu-
larize the geometry as much as possible. We chose to examine only
the nine heavy rainfalls (see Section 2.3), and from these we selected
only those 6725 scans during which at least one station showed a rain
rate in excess of 50 mm per hour.

The geometry has been regularized by rectangularizing the network
to 11 east-west rows of ten stations each. The averages of the values
from the surrounding stations have been substituted for values missing
because stations were inoperative (or fictional, such as the northeast
corner) or where the values were judged to be outliers. Substituting for
missing values has only a small effect (statistically) on the distribution
of average rain rates but a large effect on the population of paths in
various categories. This matter has been examined in some detail
and, because of the configuration of heavy-rain-rate occurrences, it
has been concluded that omission of paths containing missing observa-
tions introduces artifactual sampling variations that distort the re-
sults of the path analysis very seriously in unpredictable ways, and
that this approach is unacceptable. The distortion introduced by
using averages of surrounding values to estimate missing observations
is that of generally lowering the probability of occurrence of high rain
rates. As noted in Section 7.3, the effect is most severe on "single
station" paths.

Omission of data from rows 1, 2, and 11 of the grid, in an attempt
to make the grid rectangular without filling in corners, seriously biases
the results by throwing out too much of the heavy-rain-rate data and

many of the long paths. As justified in Section V, stations were treated as if they were located at the centers of their grid squares and scanned simultaneously. The paths were taken as starting and ending at stations; in computing the average rain rate the values at the stations at the ends of the path were given the weight $\frac{1}{2}$ and the values at the intermediate stations, the weight 1.

The diameter of a rain gauge is only 30 cm and the distance between grid square centers is 1.3 km, so that much less than 1/4000 of the path is sampled. As the rain rate has a great deal of fine structure and the number of samples per path is small (less than 11), these averages are themselves subject to very substantial fluctuations about the "true" path average.

We considered many possible paths and pairs of paths through the network. Because a single station is generally included in several paths of the same length and orientation, the separate path averages are not statistically independent. This factor has not been taken into account explicitly in our analysis.

### 6.1 *Parallel Paths*

Parallel paths are characterized by three parameters: length, separation and orientation. (To allow consistency of notation, length zero indicates paths which include only a single station and separation zero indicates a single path.) We consider four orientations: two square, north-south, and east-west; and two diagonal, northeast-southwest and northwest-southeast. A pair of parallel north-south paths of length 5.2 km and separation 2.6 km is indicated by the two dashed lines, $P_{11}$—$P_{12}$ and $P_{21}$—$P_{22}$ toward the western side of the grid in Fig. 1.

The relative probability that the rain rate on both paths exceeds the value on the abscissa is given by the ordinate for nine separations as parameters and six lengths of path as parts of Fig. 28 for nine heavy rainfalls in 1967. (The abscissa is then the minimum average rain rate for the pair.) The results for the north-south and east-west paths have been pooled. The largest number of pairs of paths entering into the data base for an individual curve is 1,338,275 (Fig. 28a curve $B$) the smallest is 80,700 (12 per scan, Fig. 28f curve $I$).

Except for a small amount of straggling in the lowest decade, the curves are smooth and well behaved. In order to display the small probabilities and high rain rates effectively, the logarithm of the probability has been used as the ordinate, and rain rate has been

Fig. 28 — Relative empirical probability that the average rain rate on both of two parallel paths will exceed the value on the abscissa. Path lengths are (a) 0 km, (b) 1.3 km, (c) 2.6 km, (d) 5.2 km, (e) 7.8 km, and (f) 10.4 km. Points indicate actual calculated values.

Separation (km)

| A | B | C | D | E | F | G | H | I |
|---|-----|-----|-----|-----|-----|-----|-----|------|
| 0 | 1.3 | 2.6 | 3.9 | 5.2 | 6.5 | 7.8 | 9.1 | 10.4 |

plotted linearly as the abscissa. The curves are generally quite straight on this semilog plot, but tend to droop somewhat at the high rain rates. This indicates that for purposes of extrapolation at the high rain rates the customary log-probability versus log-rain-rate plot is unsuitable and that improvement should be sought in the opposite direction of scaling (for example, by raising the rain rate to a power slightly greater than 1). For paths up to 10 km long, the joint probabilities decrease with increasing separation until a minimum is reached at a separation of about 9 km.

For larger separations the joint probability increases again. This result was observed for single stations in Section V. Comparison among the parts of Fig. 28 shows that for fixed probabilities the minimum average rain rate decreases as the path length increases at all separations. Thus the product of the minimum average rain rate and the path length increases more slowly than the path length itself. (This product is approximately proportional to the quantity of water along the path.)

Results from the diagonal paths generally confirm the results from the square paths. Pooling has not been carried out, however, because of the differences in path lengths and separations, and the fact that the absence of long diagonal paths at large separations tends to weight the stations in the center of the network very heavily.

The curve for zero length and zero separation corresponds to single stations, that is, point-rain-rate statistics. Similarly the curves for zero length and various separations are subsets of the space pairs of Section V. However, the data bases are not the same; missing observations have been estimated for the subset of the data used for the path analysis, making direct comparison of these results very difficult. Section VII presents a rough comparison and a discussion of the differences.

In order to minimize edge effects (the stations on the edges are included in fewer paths, that is, less heavily weighted than the central section even for the square paths) the maximum length of path considered should be short compared with the dimensions of the network. This is certainly not the case for the paths more than 5 km long. Furthermore, the results appear quite sensitive to the particular patterns of rainfall observed. The relative values, while subject to some distortion introduced by the estimation of missing points, are probably fairly good; however, we do not believe that the numerical values of the results in Fig. 28 should be considered better than an order of magnitude for the high rain rates at this stage of the analysis.

This conservative view is vindicated by the large dispersion among the results for the four orientations considered individually. For example, Fig. 29 shows the results for all parallel paths about 7.5 km long and 3.7 km apart pooled over the nine heavy rainfalls categorized by orientation. The "orientation" effect exceeds a factor of 2 at a minimum average rain rate of 60 mm per hour and a factor of 10 above about 100 mm per hour. Part of the dispersion may be real, that is, a result of systematic factors (such as the direction of the prevailing winds or the location of local updrafts) and part may be a result of differences in the population of stations which go to make up the different sets of paths. The dispersion also indicates the additional caution with which results from single lines of stations should be viewed.



Fig. 29 — "Orientation" effect. The four curves represent the two square and two diagonal orientations in the network.

## 6.2 *Adjoining Paths*

The three parameters characterizing adjoining paths (two paths having an end point in common) are the length of the arms, the angle between the arms, and the orientation of the apex. We consider three angles, $\pi/4$ (acute), $\pi/2$ (right), and $3\pi/4$ (obtuse) and the eight orientations corresponding to rotations through an angle of $\pi/4$. Right angles containing diagonal paths (**D**) are differentiated from right angles containing square paths (**R**) to avoid confounding the different lengths of arm, thus there is a total of four categories: acute (**A**), obtuse (**O**), **D**, and **R**. Two paths with arms about 3 km long joined in an obtuse angle are indicated by the two dashed lines $I_1$—$I_A$—$I_2$ toward the eastern side of the grid in Fig. 1.

The relative probability that the average rain rate on both paths exceeds the value on the abscissa is given by the ordinate for the four categories as parameters and six lengths of arm in Figs. 30a through f. The data, which are from the nine heavy rainfalls in 1967, have been pooled across orientations. The smallest and largest data bases for individual curves contain 53,800 and 4,842,000 points, respectively. For the same length of path, the results for adjoining paths lie between those for single paths (zero separation) and those for two paths separated by 1.3 km; and show the same decrease in probability with increasing length and rain rate.

A given rain rate is exceeded most frequently for the acute angles and least frequently for obtuse angles. The factor between the two increases with both path length and rain rate, and ranges from less than two (Fig. 30a) to about ten (Fig. 30d). Except for the acute angle, long paths cannot transit the center of the network so the results for various path lengths are based on different populations of stations, an effect which may introduce some bias. Very long path pairs with obuse and diagonal angles cannot be formed in the network, and are missing from Figs. 30e and f. Once again the concatenation of uncertainties suggests that these results should not be considered better than an order of magnitude in absolute values.

Figure 31 contains results for the four pairs of adjoining paths of the square family (**R**) whose arms are 5.2 km long. The paths whose apex angles point to the northwest yield probabilities substantially higher than the other three orientations. This result is not peculiar to this category of angle or length of path. Symmetry considerations would lead one to expect the results for diagonally opposed pairs (that is, NW, SE and NE, SW) to be the same. The symmetry expectation,

RELATIVE EMPIRICAL JOINT PROBABILITY

MINIMUM AVERAGE RAIN RATE IN mm PER HOUR

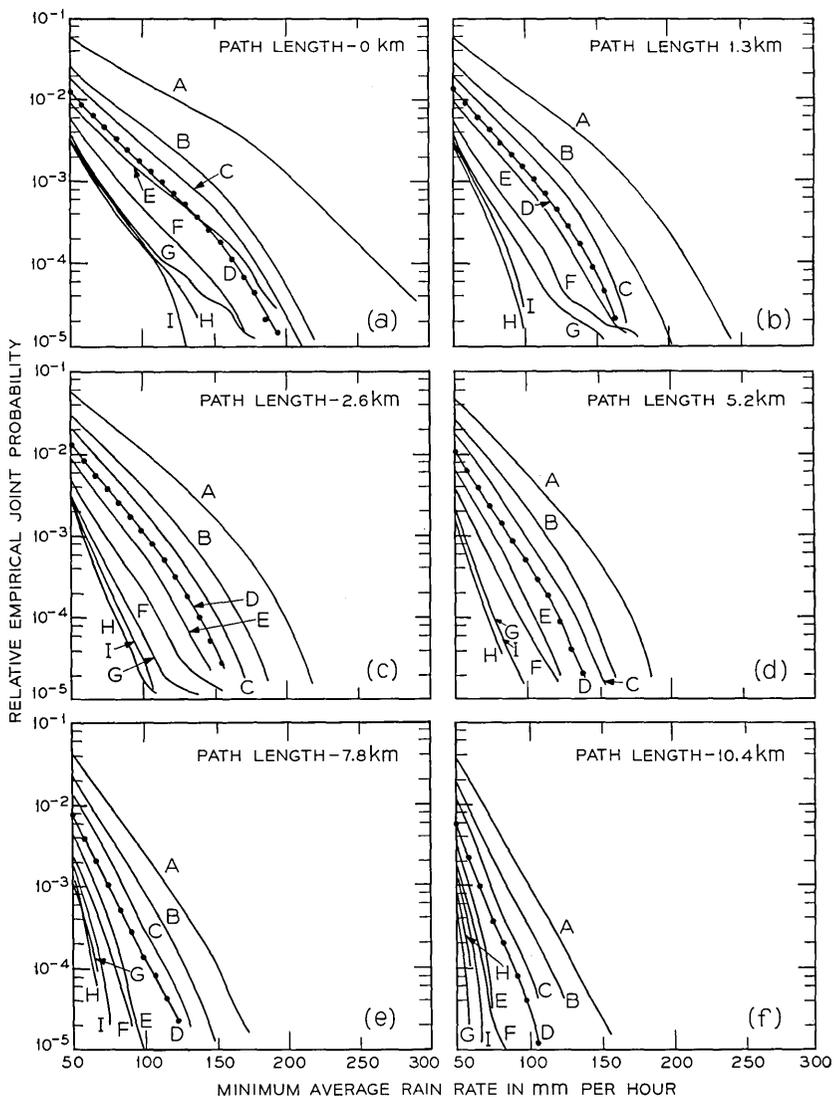Fig. 30 — Relative empirical probability that the average rain rate on both of two adjoining paths will exceed the value on the abscissa. Curves are labeled with the apex angle; see Section 6.2.

however, is quite generally violated. The explanation is that for a small grid the population of path pairs is very different for the two sets of pairs. In this case the SE, NE, and SW path pairs can never have their apexes in the northwest corner of the grid where heavy rain rates are most frequent. This sensitivity to sampling bias further confirms our reluctance to regard the numerical probabilities derived from the results of this section as better than an order of magnitude.

Some further work on statistical dependence, edge and other geometric effects, and dispersion in the data is being considered.



Fig. 31 — Asymmetry effect. Relative empirical probabilities that the average rain rate on both of two adjoining paths which constitute one of the four sets of pairs that make up the square family (R) with arms 5.2 km long will exceed the value on the abscissa. Inset: the four sets of pairs represented by the four curves.

VII. BENCHMARKS FOR ENGINEERING APPLICATIONS

All the probabilities presented so far have been calculated with respect to their own particular data bases. Thus the results do not contain assumptions about the representative properties of the data or subsets thereof, correction factors for missing data, or other adjustments which may require future revision. In this section we supply, for those who must make engineering calculations, factors which allow conversion of the relative probabilities to annual probabilities. The annual probabilities should be used only in the light of the various caveats and discussions regarding representativity, precision, and accuracy contained in this paper.

All benchmarks are computed at rain rates of 50 mm per hour. They may be applied to higher rain rates directly, but can be applied to lower rain rates only if they are scaled (see Section 4.1) to compensate for the bias introduced by excluding low-rain-rate data in the selection process.

### 7.1 A Point-Rain-Rate Benchmark

During the time covered by these data the average probability of observing a rain rate greater than 50 mm per hour was $4.3 \times 10^{-4}$. This number may be used to transform the scale of Fig. 14 as follows:

$$P(R) = 4.3 \times 10^{-4}[1 - p(R)],$$

where $P(R)$ is the annual probability that the rain rate exceeds $R$ and $p(R)$ (the empirical probability that the rain rate is less than $R$) is given by Fig. 14. This benchmark is based on a nominal observation period of $158 \times 10^4$ scans (June 1, 1967 through November 30, 1967), an allowance of 15 percent for data known to have been lost, and a correction of all the stations to 40,069 scans (110 hours of rain) on the assumption that the loss of data is random over the data base.

For convenience, some of the point-rain-rate data have been replotted in Fig. 32a on the log-log scales often used in engineering. The figure shows the Crawford Hill 1967 data and the 1958 Island Beach, New Jersey, data of Mueller and Sims.[8] (The curve in Fig. 32a was derived from Muellers and Sims Figs. 25 and 26 by D. C. Hogg.) Rain rates greater than 100 mm per hour are six times more probable in the Crawford Hill than the Island Beach measurements. However, rain rates above 250 mm per hour are equally frequent indicating

Fig. 32 — (a) Annual probabilities for point rainfall rates. (b) Annual joint probabilities for point rainfall rates at two stations separated in space.

| Label | | Separation (km) |
|:---:|:---:|:---:|
| A | (●) | 1.3 |
| B | (○) | 3.9 |
| C | (▲) | 6.5 |
| D | (△) | 9.0 |

somewhat different distributions of the rain rates measured. Fig. 32a also shows the results from stations 1 and 76. The spread between them is more than two orders of magnitude and encompasses most of the stations in the network. In view of this spread there is no reason to regard either the Crawford Hill aggregate or the Island Beach results as atypical. The bars in the figure are the interquartile ranges from the distributions of the results from the 96 stations (taken from Table III). The aggregate curve lies above the upper end of the bar when the rain rate is higher than 150 mm per hour, emphasizing the importance of rare events at these rain rates. If essential, extrapolation of the Crawford Hill curve to higher rain rates may best be accomplished by using the formula in Section 4.3 to construct the necessary ratios.

### 7.2 *A Spatial-Temporal Benchmark*

For rain rates above 50 mm per hour the relative probability scales of Figs. 23 and 27 may be converted to annual probabilities by multiplying by the factor $3.0 \times 10^{-2}$. This factor contains a correction of about 15 percent for rainfalls known to have been missed.

### 7.3 *Path Probabilities Benchmark*

The path probabilities benchmark is calculated on the premise that the 6725 scans analyzed contain all the pertinent data. The consequences of this premise are further discussed in this section. The probability scales of Figs. 28 through 31 should be multiplied by the factor $5.0 \times 10^{-3}$ to convert to annual probabilities. This factor also contains the correction (about 15 percent) for rainfalls known to have been missed.

Curve *A* of Fig. 28a, which represents single stations, is replotted using this benchmark on Fig. 32 as the thin solid line labeled "path analysis." For single stations the corrected probabilities from the path analysis are about 25 percent below those obtained from the point-rain-rate analysis for rain rates over 100 mm per hour and 35 percent lower for rain rates near 50 mm per hour. This suggests that about 30 percent of the data above 50 mm per hour have been excluded from the data base for the path analysis. However, this has been confounded by the inclusion of estimates for missing observations in the path analysis. We feel that most of the data above 50 mm per hour were included, and the lower result for the path analysis is mainly artifactual.

Relative joint probabilities of pairs of stations for four different separations are taken from Fig. 28a, corrected to annual probabilities, and replotted as the solid curves in Fig. 32b, together with the corrected relative probabilities for the same nominal separations (the points) taken from Fig. 23. The geometry of the pooling in the two sets of data differs for all except for curve *A* of Fig. 32b. The path analysis included only separations along the square grid, while the spatial analysis included all stations (both square and diagonal) within appropriately separated concentric circles. As previously noted, the data base for the path analysis is a subset of that used for the joint probabilities, and the probabilities in the path analysis have been lowered by the inclusion of estimates for missing observations. Figure 32b demonstrates the results; the probabilities from both

analyses are in good agreement for low rain rates and diverge as the rain rate increases. The values from the spatial analysis (the points in Fig. 32b) are clearly more reliable.

Results for high rain rates on paths consisting of single gauges are most affected by the estimation of missing observations, so that Fig. 32b presents the worst cases for the path analysis. As agreement is good near 50 mm per hour and the missing-data effect is more serious at the higher rain rates and less severe for the longer paths, no single compensatory adjustment of the benchmark can be made.

### 7.4 *Seasonal Correction*

Data concerning fading on the microwave transmission path (T-R in Fig. 1) indicate that during 1967 most of the heaviest rain rates, which are usually generated by summer thunderstorms, occurred between June 1 and December 1.[6] This seasonal effect decreases the probability of observing the highest rain rates by a factor of two, without materially affecting the probability of observing light rain. Insufficient data are available to correct our distributions for the seasonal variation. However, an approximate adjustment may be made by multiplying the annual probabilities associated with rain rates greater than about 100 mm per hour by the factor $\frac{2}{3}$, if the probabilities are to be applied on an annual basis. Correction for the seasonal effect would improve some aspects of the comparison between the Island Beach and Crawford Hill data (Fig. 32) and worsen others.

### VIII. SUMMARY AND CONCLUSIONS

This article presents some statistical summaries of rainfall data acquired on the Crawford Hill, New Jersey, rain gauge network between June 1 and November 30, 1967. The emphasis of the analysis is on rain rates greater than 50 mm per hour, which cause substantial attenuation of electromagnetic waves in the 10 to 30 GHz frequency range. The objects of analysis are: the behavior of rain rates at a point in space, the relationship of two rain rates separated in space or in time, and the relationships of average rain rates on pairs of paths in various configurations.

The network consisted of 96 stations approximately uniformly distributed over a rectangular area 13 by 14 km.[2] The gauges are small in area (478 cm[2]), measure the rain rate averaged over less than one seconds, and are read out in sequence every ten seconds.[1] About 430

hours of data were recorded; the 110 hours of greatest interest have been examined in detail. The data are generally satisfactory for rain rates greater than about 10 mm per hour, although some substantial data-screening problems were encountered. There is no indication that these data are atypical of rain rates in northern littoral New Jersey, but they constitute a very small sample at the high rain rates; there is only informal indication of the large range of variation to be expected.

Heavy rainfall is found to come in irregular bursts. The time series containing these bursts are nonstationary (not even piecewise stationary at the sampling rate used), and no simple description of the wide variety of observed distributions of the rain rates has been found. Several poolings (for example, one station for the season, all stations for a rainfall) of the data were examined in an unsuccessful effort to find subsets of the data that behaved consistently enough to be described by a single distribution.

Figures 11 and 12 summarize the distribution of the point-rain-rate data. The high-rain-rate tail ($i$) is longer than would be expected on the basis of the normal distribution, ($ii$) can be approximated by an exponential function [equation (1)], ($iii$) depends heavily on a small number of cloudburst-like events (and thus shows a lot of scatter, see Fig. 13 and Table II), and ($iv$) is sensitive to the data-screening procedures used in the data processing. Restricting the data to rain rates greater than 50 mm per hour gives a better-defined subset of the data than does the original data-selection process and allows a further examination of the dispersion in the data (Table III), but fails to reduce the distributional variety observed.

Empirical probabilities for point rain rates above 50 mm per hour are given in Table IV. These measurements are somewhat higher (up to a factor of 6 at 50 mm per hour) for rain rates below 250 mm per

TABLE IV—EMPIRICAL PROBABILITIES FOR POINT RAIN RATES*

| $R$ mm per hour | Probability (rain rate $> R$) | Total time (rain rate $> R$) minutes per year |
|---|---|---|
| 50 | $4.3 \times 10^{-4}$ | 230 |
| 100 | $1.3 \times 10^{-4}$ | 70 |
| 150 | $4.2 \times 10^{-5}$ | 22 |
| 200 | $1.0 \times 10^{-5}$ | 5 |

\* Data are for Crawford Hill network during 1967.

hour than the Island Beach, New Jersey measurements of Mueller and Sims.[8] However, better agreement would be obtained if a seasonal correction were applied.

When it is raining heavily at a station, nearby stations are also likely to show heavy rain. More detailed examination of the spatial characteristic of the data in terms of the joint probability of the rain rate being greater than $R$ mm per hour at each of two stations simultaneously shows a rapid decrease of the joint probability with increasing separation for short separations and a minimum in the joint probability when the stations are separated by about 12 km (Fig. 23). Other spatial relationships also change qualitatively when the separations exceeds about 12 km. The temporal behavior may be characterized by the behavior of the joint probability of the rain rate at a station exceeding $R$ mm per hour at both the beginning and end of a given time period (Fig. 27). This probability declines smoothly and relatively slowly as the interval increases from 10 to 360 seconds.

Of particular engineering interest are the joint probabilities that the average rain rate will simultaneously exceed a particular value on all of a number of paths. Figure 28 summarizes results of analysis of various sets of paths within the network for parallel paths; Figure 30 summarizes the same results for adjoining paths. Figures 29 and 31 give some indication of the dispersion among subsets of the paths.

For parallel paths up to 10 km long, the joint probability is minimum for paths separated by about 9 km. The joint probabilities for parallel paths 6.5 km long (and ordinary probabilities for single paths) are tabulated in Table V. As seen in Table V, for paths of this length the joint probability may be more than 200 times lower than the ordinary probability. This indicates the improved reliability of microwave transmission that may be obtained with appropriate choice of redundant paths. While the relative values of the probabilities obtained in this analysis are likely to be reasonably good, for reasons given in Section VI the numerical values should be regarded as correct only within an order of magnitude. The joint probabilities associated with average rain rates on adjoining paths are somewhat less than the probabilities associated with single paths through the network.

Section VII supplies benchmarks for converting the relative probabilities presented in the body of the paper to the annual probabilities needed for engineering calculations and discusses an additional correction for seasonal effects. Annual probabilities for point rain rates are plotted in Fig. 32a, and results of various analyses are compared and found to be consistent.

TABLE V—JOINT PROBABILITY THAT THE AVERAGE RAIN RATE ON BOTH
OF TWO PARALLEL PATHS EXCEED $R^*$

| $R$ mm per hr \ Separation km | 0 (Single path) | 1.3 | 3.9 | 6.5 | 9.1 |
|---|---|---|---|---|---|
| **50** | $2.5 \times 10^{-4}$ (2 hr per yr) | $1.5 \times 10^{-4}$ (75 min per yr) | $5 \times 10^{-5}$ (25 min per yr) | $1.5 \times 10^{-5}$ (7.5 min per yr) | $7.5 \times 10^{-6}$ (4 min per yr) |
| 100 | $2.0 \times 10^{-5}$ (10 min per yr) | $7.5 \times 10^{-6}$ (4 min per yr) | $1.3 \times 10^{-6}$ (30 s per yr) | $1.0 \times 10^{-7}$ (3 s per yr) | $< \sim 10^{-8}$ |
| 150 | $1.5 \times 10^{-6}$ (30 s per yr) | $1.0 \times 10^{-7}$ (3 s per yr) | $< \sim 10^{-8}$ | $< \sim 10^{-8}$ | $< \sim 10^{-8}$ |

* The paths are 6.5 km long. Probability is tabulated by the separation between paths. (Ordinary probability for the single paths.)

## IX. ACKNOWLEDGMENT

REFERENCES

1. Semplak, R. A., "Gauge for Continuously Measuring Rate of Rainfall," Rev. Sci. Instruments, *37*, No. 11 (November 1966), pp. 1554–1558.
2. Semplak, R. A., and Keller, H. E., "A Dense Network for Rapid Measurement of Rainfall rate," B.S.T.J., this issue, pp. 1745–1756.
3. Semplak, R. A., and Turrin, R. H., "Some Measurements of Attenuation by Rainfall at 18.5 GHz," B.S.T.J., this issue, pp. 1767–1787.
4. Freeny, A. E., "Statistical Treatment of Rain Gauge Calibration Data," B.S.T.J., this issue, pp. 1757–1766.
5. Wilk, M. B., and Gnanadesikan, R., "Probability Plotting Methods for the Analysis of Data," Biometrika, *55*, No. 1 (March 1968), pp. 1–17.
6. Semplak, R. A., unpublished work.
7. Hogg, D. C., unpublished work.
8. Mueller, E. A., and Sims, A. L., "Investigation of the Quantitative Determination of Point and Areal Precipitation by Radar Echo Measurements," Technical Report ECOM-00032-F, U. S. Army Electronic Command, Fort Monmouth, New Jersey, December 1966. (Contract DA-28-043 AMC-00032(E) with the Illinois State Water Survey, University of Illinois, Urbana, Illinois.)

# Overload Stability Problem in Submarine Cable Systems

By CLEO D. ANDERSON

*Modern submarine cable systems usually provide bidirectional trans-
mission over a single cable with repeaters which use a common amplifier
for both directions of transmission. The signals for the two directions of
transmission occupy separate frequency bands. Under normal loads, the
amplifier is highly linear and there is negligible interaction of the signals
in the two bands. However, when approaching overload, intermodulation
in the repeaters may transfer appreciable power between bands. It has been
discovered that this feedback, especially in the presence of large misalign-
ments, can result in a system that maintains itself in overload. Such a
system, once excited by a momentary signal or noise peak, generates suffi-
cient intermodulation noise to keep itself overloaded, even in the absence of
any further external signal. This paper describes the occurrence which drew
attention to this phenomenon, presents an analytical approach used to
predict stability margins for any given repeater, and describes the action
taken to ensure the stability of the SF Submarine Cable System.*

## I. INTRODUCTION

Experience gained while laying the Oahu-Guam SD Submarine
Cable System shows that a bidirectional single cable system whose
repeaters have a common amplifier for both directions of transmission
is potentially unstable. So far, this instability has occurred only when
an abnormally high-gain repeater section is present. When the in-
stability occurs, noise levels corresponding to repeater overload make
communication in either direction completely impossible. The in-
stability results from the feedback of intermodulated signals or noise
in overloaded repeaters. Low-band power traveling in one direction
through the system is partially converted to high-band power which
propagates in the opposite direction, and vice versa. In this manner,
a configuration of feedback loops develops. Loops in which repeaters
are separated by less than nominal loss are especially critical.

## II. THE FEEDBACK MECHANISM

The SD and SF repeaters consist of a common amplifier and directional filters. This configuration, shown in Fig. 1, is well suited to undersea applications because of the reduced number of components, greater reliability, and lower power drain compared with the more conventional twin amplifier arrangement. With typical pre-emphasis, the low-band signal levels at the amplifier output are about 10 dB below those of the high band. Thus, the low-band signals comprise a small portion of the total multichannel repeater load.

The common amplifier configuration has a severe limitation when one considers the effect of a very large overload. Assume that the high-band signal, which enters port A of Fig. 1, overloads the amplifier. This will produce, by intermodulation, not only a distorted output at B, but will also return power in the low band at A. Thus the repeater, through its nonlinearity, partially redistributes and reflects the spectrum of the incident power.

The simplest type of system overload instability occurs when a single high-gain repeater section is present as shown in Fig. 2. Assume that no signal is being transmitted by either terminal so that the load results entirely from noise generated by the repeaters. Because of the low loss equalizer section, repeaters $n + 1$ to $N$ will be loaded mainly with low-band power, while repeaters $n$ to 1 will be loaded mainly with high-band power. If sufficient loss is removed from the equalizer, the repeaters on one side of the equalizer will become overloaded with noise generated by repeaters on the opposite side and a self-sustaining overload results.

## III. SD SYSTEM EXPERIENCE

### 3.1 *Wake Island*

The conditions shown in Fig. 2 were present for several days while the Oahu-Guam SD Submarine Cable System was being laid. This system has intermediate stations at Midway and Wake Islands which allow channels to be bridged onto the system at these two points. A modified ocean block equalizer is included in the through circuit at the intermediate stations. In the normal setting, this equalizer's loss is equal to 10 nm of cable.

While the second shipload of cable was being laid, a dc fault occurred which was apparently in the immediate vicinity of Wake Island. At that time 144 of the 200 repeaters had been laid. Wake Island is

Fig. 1 — SF and SD repeater configuration.

between repeaters 121 and 122 counting from Oahu. To eliminate the Wake Island equalizer as a possible source of the dc fault, it was completely bypassed resulting in an excess gain equal to the loss of 10 nm of cable. After the equalizer was bypassed, system power was turned back up and the dc fault was no longer present. In the meantime, four additional repeaters had been laid without power.

Following power turn up, laying continued, but after a short time the system became noisy to the extent that order wire communication between the ship and Oahu was impossible. Power was turned down and when it was raised again the noise was no longer present. This was repeated several times. However, soon after repeater 156 was laid the noise reappeared and subsequent efforts to squelch it by lowering and raising power were unsuccessful. Then it was decided to install the spare equalizer on Wake Island. After this was accomplished, no further noise developed and the laying was completed. Subsequent examination of the original equalizer revealed no imperfections. The location or cause of the temporary dc fault has never been established.

3.2 *Description and Explanation of Noise Conditions*

The noise which existed after bypassing the equalizer appeared to result from repeater overload. The noise covered both bands and no



Fig. 2 — Interband power flow during self overload.

discrete frequencies were observed. At the time, the noise condition was believed to be related to the previously observed dc fault. Several months elapsed before the phenomenon was attributed to the excess gain resulting from by-passing the equalizer.

The excess gain between repeaters 121 and 122 was about 25 dB in the high band and 17 dB in the low band. Thus, any large burst of noise generated in the high band west of Wake Island tended to overload those repeaters east of Wake. They, in turn, generated modulation noise, some of which fell into the low band. This low-band noise was transmitted west, encountered excess gain, and tended to overload the repeaters west of Wake, regenerating the high-band noise. Thus, a potential feedback loop was established through the intermodulation of noise. With the excess gain present and a sufficient number of repeaters west of Wake converting power from low to high band, normal background noise alone was sufficient to initiate a self-sustained overload.

Evidently the stability of the system was marginal with the excess gain of 10 nm of cable at Wake and 22 repeaters laid west of Wake. When 34 repeaters had been laid west of Wake the system was definitely unstable. This means that 34 repeaters shifted enough power from the low to the high band so that, when this shifted power was subjected to the 25 dB of excess gain at Wake Island, it reached a level high enough to cause a significant shift of power into the low band by the repeaters between Wake Island and Oahu. This low-band power, enhanced by the 17 dB of excess gain at Wake, was sufficient to "close the loop" and the overload was self-sustaining. The repeaters that had not yet been laid did not contribute significantly in the process because of the excess loss of the relatively warm cable aboard ship. This excess loss attenuated the low-band power that caused the power shift as well as any high-band power that might have been generated.

### 3.3 Two-Repeater Experiment

To test the preceeding hypothesis, an experiment was conducted using two SD pilot model repeaters connected to each other with an adjustable cable shape equalizer. It was found that by reducing the loss between the repeaters to an amount equivalent to 2.5 nm of cable a "sing condition" could be established with the repeaters sustaining each other in noise overload. The reduction to 2.5 nm in the experiment rather than the 10 nm at Wake Island was necessary because in

the experiment only two repeaters shifted power between bands. Thus, the Wake Island experience was very valuable in revealing a source of instability which is inherent in an equivalent four-wire system using a common amplifier for both directions of transmission. An intensive study of this problem was begun to evaluate the stability of the SF system which was then in the final design stage.

## IV. SF SYSTEM STABILITY ANALYSIS

### 4.1 *Repeater Characterization*

Although the SD overload condition was analyzed and the result agreed with the experimental evidence, it was obvious that a direct application of those methods would not be valid for predicting the stability of a normal system where no large abrupt changes in repeater levels occur. The fundamental problem is one of characterizing the repeater input-output relations under severe overload conditions. The SF repeater overload performance was determined by noise loading both bands simultaneously and measuring the apparent compression or expansion of the power in each band as a function of the average driving power in each band. The instrumentation of the repeater measurements is shown in Fig. 3.

Such a characterization is obviously not rigorous in that it neglects the spectral distribution of either the input or output noise power. The input was band-limited white noise. Figures 4 and 5 show the three dimensional characterization. The nonlinear behavior is described simply as an apparent compression or expansion of the signal relative to the transmission of a truly linear repeater. Compression indicates a reduction in gain while expansion, which is the source of instability, is an apparent increase in gain.

Figures 4 and 5 show that the power transferred from low to high band is much greater than the power flow from high to low band. Since the mechanism of interband power flow is intermodulation noise produced by one band falling into the other, it seems reasonable that the band with greater feedback would undergo least expansion. The difference in feedback between the two bands is about 12 dB, which is roughly the same as the difference between the maximum high- and low-band expansion.

### 4.2 *Analytical Model*

In line with the empirical nature of the repeater characterization, the analysis of system stability is based upon an iterative scheme.

Fig. 3 — Instrumentation of repeater measurements.

Fig. 4 — SF repeater high band expansion versus drive. (0 dB is input power for 0 dBm output power.)



Fig. 5 — SF repeater low band expansion versus drive. (0 dB is input power for 0 dBm output power.)

Simply stated, the method seeks to determine whether a system would "recover" from a gross overload. By "recover" is meant that the signal levels at the repeaters would return to their normal value once the overloading signal is removed from the transmitting terminals. If any repeater remains overloaded in the steady state after the overloading signal is removed from both terminals, then clearly the system is unstable.

The quantities which must be determined numerically are the steady state values of output power in each band of every repeater, as indicated in Fig. 6.

Let $P_H(n)$   = high-band power delivered by repeater $n$,

$P_L(n)$   = low-band power delivered by repeater $n$,

$G_H(n)$   = excess high-band gain of $n$th repeater section, (that is, misalignment),

$G_L(n)$   = excess low-band gain of $n$th repeater section,

$F_H(P_L , P_H)$ = repeater high-band expansion with driving powers $P_L$ and $P_H$ in the low and high bands, respectively, and

$F_L(P_L , P_H)$ = repeater low-band expansion with driving powers $P_L$ and $P_H$ as above.

$F_H$ and $F_L$, which were determined experimentally by noise loading a repeater, are plotted in Figs. 4 and 5. These functions describe the nonlinear input-output relationship of the repeaters. The two equations which must be satisfied are:

$$P_H^i(n) = P_H^{i-1}(n + 1) + G_H(n)$$

$$+ F_H[P_L^{i-1}(n - 1) + G_L(n), P_H^{i-1}(n + 1) + G_H(n)]$$

$$P_L^i(n) = P_L^{i-1}(n - 1) + G_L(n)$$

$$+ F_L[P_L^{i-1}(n - 1) + G_L(n), P_H^{i-1}(n + 1) + G_H(n)]$$



Fig. 6 — Quantities which describe the state of the system.

where the index $i$ is the number of the iteration while $n$ ranges over all repeaters.

These two equations simply say that the state of a repeater on the $i$th iteration is determined by the states of the adjacent repeaters on the $(i - 1)$th iteration. The initial conditions are set by assuming the system is linear and applying a large overload at the two transmitting terminals. The overload signal is then removed and the system "relaxed" by successive applications of the two preceeding equations to every repeater. When the maximum difference between the states of each repeater on two successive iterations is sufficiently small, say 0.1 dB, the process is halted and the equilibrium or steady state is assumed to be reached. Then if any repeater carries a load greater than that caused by the normally applied signals, an unstable condition exists.

This sequence of operations was programmed in FORTRAN on the IBM 7094 computer. The input data characterizing the repeater non-linear behavior was in the form of two 20 × 20 matrices which covered the entire repeater power range of interest. Linear interpolation was used between input data points. Critical misalignment conditions were found by successive trials with different values of misalignment. Stability margins were determined by finding that amount of mis-alignment in the system which would just cause instability. The stability margins are, by definition, equal to this misalignment.

Various forms of misalignment (net gain or loss) can be assumed. The simplest is lumped gain at a particular point in the system. (This corresponds to the Wake Island experience.) By successive trials, that point in the system is found where the smallest amount of excess gain results in instability. That value of gain at this point is defined as the lumped gain margin. This situation does not correspond to an event that is likely to occur in practice. It is simply a figure of merit characterizing the stability of the system.

Another type of margin studied corresponds to within-block mis-alignment. Submarine cable systems are divided into ocean blocks. Each block has associated with it a block equalizer which is expected to compensate for most of the misalignment (net gain or loss) accumu-lated in that block. In the SF system, for example, a block consists of 20 repeaters, 192 miles of cable, and an equalizer. The within-block misalignment stability margin is the amount of misalignment that could occur in every block, equalized perfectly at every equalizer, and result in a system which is marginally stable. This situation cor-responds more closely than the Wake Island experience to the type

of misalignment encountered in practice. It has been found that excess loss within a block compensated by gain in the equalizer section is more critical than excess gain in a block compensated by loss in the equalizer section.

A third type of misalignment margin assumes a uniform gain across the whole system compensated only in the terminals. This situation corresponds to misalignment resulting from aging or temperature where equalization by the ocean block equalizers is not possible. For this condition, only excess gain can cause instability; excess loss tends to prevent instability.

### 4.3 Computed SF Stability Margins

Based on the measured repeater nonlinear performance and the preceeding analysis, these marginal stability conditions were computed:

    (i) 10 dB excess gain per block completely equalized at the ocean block equalizer,

    (ii) 6 dB loss per block completely equalized at the ocean block equalizer,

    (iii) 6 dB of lumped gain in both bands at one ocean block equalizer and

    (iv) 13 dB of uniform positive misalignment in both bands.

Although the expected misalignment was considerably less than these critical values, it was felt that the margins were dangerously low, especially for within-block loss compensated by equalizer gain. Block losses of 3 dB or more may be expected under normal conditions. There was also a large uncertainty in the amount of variation in the overload response of repeaters and the effect of aging. What a "safe" stability margin would be was not known. However, from SD system experience, we knew that its stability margin, about 15 dB lumped gain, was adequate. Consideration of these points led to the study of a practical means of improving the SF system stability margin.

### 4.4 Improving the SF Stability Margin

Two different approaches were taken to improve the stability margin. The first was to modify the repeater feedback. No improvement resulted, probably because in the process the repeater feedback phase was changed from its optimum 90° value. The second attempt involved

using a diode limiter to prevent the low band signal from severely overloading a repeater.

Figure 4 shows that if the low band repeater output power remains below about 12 dBm, then no power is transferred from the low band to the high band and the feedback loop is opened. Since a limiter by itself generates intermodulation noise, it must be placed somewhere in the transmission path where the two bands are physically separated and are isolated by filters. Otherwise the limiter itself would transfer power between bands.

The limiter was placed in the low band branch of the ocean block equalizer as shown in Fig. 7. Since the average low band signal power at the repeater output is about −10 dBm, limiting the maximum rms repeater low band output power to about 12 dBm involves very little signal-to-noise penalty. System stability studies assuming diode limiters in each ocean block equalizer indicated that it would be possible to increase the stability margin by at least 6 dB with a resulting noise contribution of less than 20 dBrnC0 on a 3500 nm system. Figure 8 compares the measured low band compression of the limiter with that of a repeater. The driving signal was white noise, band limited to cover only the low band.

V. CONCLUSIONS

Equivalent four-wire transmission systems which use repeaters with a common amplifier for both directions of transmission are po-

Fig. 7 — Location of low-band limiter.

Fig. 8 — Effect of limitation on low-band compression. (0 dB is input power for 0 dBm output power.)

tentially unstable. Although this instability has never yet been experienced under normal operating conditions, it has occurred when an unusual amount of excess gain was present. A method of calculating a first order approximation of the stability margin has been described. This method has been used in evaluating the stability of the SF Submarine Cable System. The nonlinear stability problem appears to become more acute as system bandwidths increase. This effect may limit, or even possibly preclude, the use of a common amplifier in future broadband submarine cable repeaters. Although only the surface of this problem has been probed we hope that publishing these results will stimulate a deeper investigation.

VI. ACKNOWLEDGEMENTS

# A Continuously Adaptive Equalizer for General-Purpose Communication Channels

By H. R. RUDIN, JR.

(Manuscript received July 18, 1968)

*This paper describes and analyzes a technique for the automatic and continuous minimization of the linear distortion in a communication channel. The channel can be used for the transmission of information in any format while the minimization process continues simultaneously. The equalizer which implements this strategy thus responds adaptively to changes in the channel's frequency-response characteristics.*

*Recent work in the field of automatic equalization has been along two main avenues. The first has led to a series of special-purpose equalizers, each designed to function with a particular information signalling format. This first class of equalizers is readily made adaptive by use of the known signal structure. The second approach is the optimization of the channel's frequency response in the general sense so that its transmission capability is improved for any signal. Equalizers of the latter type have, to date, been of the preset type, requiring adjustment in a dedicated period prior to information transmission.*

*This paper describes a technique which retains the advantage of the general approach yet adds that of adaptive operation. The result is an equalizer which can be used with any information transmission scheme and which continuously strives to compensate for any change which occurs in the channel's characteristics.*

## I. INTRODUCTION

Automatic equalization techniques have done much to alleviate the deleterious effects of linear distortion in communication channels.[1] In particular, equalization schemes such as that suggested in Ref. 2 have permitted significantly increased transmission speeds and better error performance for linear, synchronous data transmission systems. The same improvements can be obtained for arbitrary information transmission systems (not necessarily synchronous or linear) through the

use of the generalized equalization techniques described in Ref. 3.

The equalization problem becomes more complicated when the linear distortion in the communication channel becomes time varying. In the time-varying case an adaptive equalizer is needed; such an equalizer for the synchronous, linear information transmission case has been discussed in the literature.[4] This paper describes a continuously adaptive equalizer which can be used with any information transmission system: synchronous or asynchronous, using linear or nonlinear modulation.

The continuously adaptive equalizer, described here, has the advantage over that described in Ref. 3 that the communication channel is under continuous surveillance by the equalizer controller. If the frequency characteristic of the transmission path should change, the equalizer would immediately begin to compensate for that change. Such a correction would be made regardless of the format being used for information transmission in the channel and independent of whether the channel were being used for information transmission.

The operation of the system relies on the transmission of a low-level test signal sent simultaneously with the information-bearing signal. Despite the low energy level of the test signal, powerful correlation detectors extract the necessary information for continuous, adaptive equalization of the channel.

## II. THE TECHNIQUE

The preset version of a general-purpose channel equalizer is shown in Fig. 1. A test signal is transmitted and the channel equalized prior to use of the channel for information transmission purposes.[3] The test signal used for equalization is removed before information transmission begins.



Fig. 1 — Preset mean-square channel equalizer.

The adaptive channel equalizer is shown in Fig. 2. Here the test signal is added to the information signal. The equalizer controller operates continuously so that the channel is always being monitored and the equalizer immediately responds to a change in the channel characteristic.

Because the test signal appears in the high "noise" environment created by the information-bearing signal, steps must be taken to ensure that the test signal and not the "noise" (or information signal) dictates the behavior of the equalizer. The preset equalizer responded to both test signal and noise.[3] This is an advantageous mode of operation for the preset case because if the signaling statistics are known *a priori* the equalizer could be designed to maximize the total received signal-to-noise ratio where the noise consists of both random noise and a component resulting from residual linear distortion. The same approach cannot be used in the adaptive equalizer because the much larger effective noise would dominate in the control of the equalizer. A stunt is used to make the equalizer controller blind to the information signal.

## 2.1 *Review of Preset Equalizer Operation*[3]

In both preset and adaptive mean-square channel equalizers, the desire is to equalize the channel transmission characteristic so that it best resembles an ideal transmission characteristic. The fit is optimized under a mean-squared error criterion. The distortion to be minimized is

$$E_1 = \int_{-\infty}^{\infty} |H(\omega) - G(\omega)|^2 \, d\omega \qquad (1)$$



Fig. 2 — Adaptive mean-square channel equalizer.

where $H(\omega)$ is the equalized channel characteristic and $G(\omega)$ is the ideal channel characteristic. The error criterion given in equation (1) can be made more general by adding information concerning the relative importance of errors at various frequencies by the inclusion of a real, nonnegative weighting function $\mid W(\omega) \mid^2$ which assigns a relative weight $\mid W(\omega) \mid^2$ to the equalization error at each frequency $\omega$, as shown in Fig. 3. The resultant expression for the distortion to be minimized is

$$E = \int_{-\infty}^{\infty} \mid H(\omega) - G(\omega) \mid^2 \mid W(\omega) \mid^2 d\omega. \tag{2}$$

The ideal characteristic $G(\omega)$ would normally have flat amplitude and linear phase responses in the band of interest while the error weighting function $\mid W(\omega) \mid^2$ would be chosen to resemble the spectral density of the signal most likely to be transmitted. This choice of $\mid W(\omega) \mid^2$ would ensure that the best equalization would occur at frequencies where most of the information signal energy occurs.

Parseval's theorem allows equation (2) to be rewritten

$$E = \int_{-\infty}^{\infty} [h(t)*w(t) - g(t)*w(t)]^2 dt \tag{3}$$

where $h(t)$, $g(t)$, and $w(t)$ are the time functions corresponding to $H(\omega)$, $G(\omega)$, and $W(\omega)$, respectively, and * denotes convolution. Convergence of the algorithm to be stated shortly can be guaranteed for the case where the equalizer impulse response can be written as a linear,



Fig. 3 — Mean-square equalizer.

weighted sum of the form[3,5-7]

$$h(t) = \sum_{n=-N}^{N} c_n y_n(t) * x(t). \qquad (4)$$

Here $x(t)$ is the distorting channel's response, and $c_n$ is the weight (which can be either positive or negative) associated with the impulse response $y_n(t)$ of the elemental network $Y_n(\omega)$ shown in Fig. 4. When equation (4) is substituted into equation (3) the resulting distortion measure is

$$E = \int_{-\infty}^{\infty} \left\{ \sum_{n=-N}^{N} [c_n y_n(t) * x(t) * w(t)] - g(t) * w(t) \right\}^2 dt. \qquad (5)$$

Partial differentiation of this last relation with respect to the $j$th attenuator, $c_j$, yields

$$\frac{\partial E}{\partial c_j} = 2 \int_{-\infty}^{\infty} \{h(t) * w(t) - g(t) * w(t)\} \{y_j(t) * x(t) * w(t)\} \, dt. \qquad (6)$$

Since it can be easily demonstrated that $E$ is a convex function of the weights $c_n$ $(n = -N, N)$, the set of $2N + 1$ equations of the form of equation (6) provides the information necessary to obtain the desired optimum adjustment of the weights under the specified mean-squared error criterion. The polarity of equation (6) determines the directions in which the attenuator weights must be incremented to



Fig. 4 — Generalized mean-square equalizer.

minimize the distortion $E$; the minimum is achieved when all the partial derivatives are zero.

Equation (6) is simply the cross-correlation between the error signal and the input signal to the $j$th attenuator. Thus the system can be implemented as shown in Fig. 4 where only one correlator is shown for clarity.

If the general parallel structure of Fig. 4 is replaced by a series tapped delay line structure so that

$$h(t) = \sum_{n=-N}^{N} c_n x(t - n\tau) \tag{7}$$

where $\tau$ is the tap spacing, a relation for the attenuator settings identical to equation (8) of Ref. 3 is obtained.

If noise is introduced, the optimum settings of the attenuator weights are changed. If noise $\eta(t)$ with spectrum $N(\omega)$ is added in the channel shown in Fig. 3, the received signal is

$$y(t) = w(t)^* x(t) + \eta(t). \tag{8}$$

If the measure of distortion in the presence of noise, $E_n$, is again taken as the average mean-square error between the equalized received signal $z(t)$ and the transmitted signal passed through the ideal, noiseless channel $G(\omega)$,

$$E_n = \langle [z(t) - w(t)^* g(t)]^2 \rangle. \tag{9}$$

The brackets $\langle \ \rangle$ denote a time average. The equalized signal $z(t)$ is now given by

$$z(t) = \sum_{n=-N}^{N} c_n [\eta(t - n\tau) + w(t)^* x(t - n\tau)], \tag{10}$$

where the equalizer is of the tapped delay line type. Substitution of equation (10) into (9) and partial differentiation of the result yields:

$$\frac{\partial E_n}{\partial c_i} = 2\langle [z(t) - w(t)^* g(t)][\eta(t - j\tau) + w(t)^* x(t - j\tau)] \rangle. \tag{11}$$

The noise $\eta$ now appears in both the first (or error) term, as a component of $z(t)$, and in the second term so that the noise does in general affect the optimum attenuator settings.

In the preset equalizer this phenomenon can be used to advantage. If the spectral density of the information signal to be transmitted over the equalized channel is known *a priori*, then the $W(\omega)$ error spectral weighting function can be selected such that the equalizer will mini-

mize the total expected mean-squared error including both noise and linear distortion components over the frequency band of interest.

## 2.2 *Operation of the Adaptive Channel Equalizer*

While the dependence of the optimum equalizer attenuator settings on the noise could be used to advantage in the preset equalizer, such is not the case in the adaptive equalizer. In the adaptive equalizer the information-bearing signal is added to the test signal. Insofar as the equalization process is concerned, the information- bearing signal acts as noise and would interfere with the equalization process. Since this "noise" is apt to be larger than the test signal, the equalizer controller must be made insensitive to it.

Because it is necessary to generate a replica of the transmitted test signal at the location of the equalizer controller, it is convenient to use a particular pseudorandom sequence for the test signal. Such a test signal has a periodic autocorrelation function.[3] Over a short period of time the clocked ones and zeroes of the sequence look quite random but seen over a longer period of time it is clear that the sequence is in fact a periodic sequence.

The periodic property of the pseudorandom sequence permits the use of a simple stunt to achieve the desired independence of the attenuator settings from the high-noise environment. In the adaptive version of the channel equalizer the error signal is passed through a delay of $T_D$ seconds before reaching the correlator. The delay $T_D$ is chosen equal to a multiple of the period of the pseudo-random sequence. The fully expanded version of equation (11) is

$$\frac{\partial E_n}{\partial c_j} = 2\left\langle\left\{\sum_{n=-N}^{N} c_n[\eta(t - n\tau) + w(t)^*x(t - n\tau)] - w(t)^*g(t)\right\}\right.$$
$$\left. \cdot[\eta(t - j\tau) + w(t)^*x(t - j\tau)]\right\rangle. \qquad (12)$$

If the error signal is now passed through the delay $T_D$, the relation becomes

$$\frac{\partial E_n}{\partial c_j} = 2\left\langle\left\{\sum_{n=-N}^{N} c_n[\eta(t - n\tau - T_D) + w(t)^*x(t - n\tau - T_D)]\right.\right.$$
$$\left.\left. - w(t)^*g(t - T_D)\right\}[\eta(t - j\tau) + w(t)^*x(t - j\tau)]\right\rangle. \qquad (13)$$

Since the test signal and information signal are generated in completely different fashions it is certainly reasonable to assume that

$x(t)$ and $\eta(t)$ are uncorrelated as are $g(t)$ and $\eta(t)$. Also, because $T_D$ is selected large compared with the channel's dispersion (or length of the channel's impulse response), time translations such as the pair $[\eta(t - j\tau), \eta(t - n\tau - T_D)]$ are also independent. Thus, assuming zero mean values for the signals involved, equation (13) may be rewritten

$$\frac{\partial E_n}{\partial c_i} = 2\left\langle\left\{ \sum_{n=-N}^{N} c_n[w(t)^*x(t - n\tau - T_D)] - w(t)^*g(t - T_D)\right\}\right.$$
$$\left. \cdot [w(t)^*x(t - j\tau)]\right\rangle. \qquad (14)$$

Because of the periodicity of the test waveform, $g(t - T_D) = g(t)$ and $x(t - T_D) = x(t)$. Equation (14) then becomes

$$\frac{\partial E_n}{\partial c_i} = 2\langle[h(t)^*w(t) - g(t)^*w(t)][x(t - j\tau)^*w(t)]\rangle \qquad (15)$$

which is entirely equivalent to equation (6) for the noise-free case and where the equalizer is of the tapped delay line or transversal filter type. The same arguments can be made for the parallel form of equalizer shown in Fig. 4. Use of delay to achieve the kind of independence needed here has been described elsewhere in the literature.[8]

So far the discussion has dealt with the effect of the information-bearing signal on the test signal. Insofar as the information receiver is concerned, the test signal looks just like noise. If the input to the information receiver is taken as the output of the error differencing amplifier as shown in Fig. 5, the noise component resulting from the test signal can be substantially reduced. The better the equalization, the better the match between desired and equalized test signals and the lower the effective noise for the information-bearing signal. In the laboratory the effective noise of the test signal is reduced by about 20 dB.

## 2.3 An Estimate of Settling Time

In the preset mean-square channel equalizer the settling time (the time required for the equalizer to reach equilibrium) is dependent largely on the nature of the dispersion in the channel. One reason for this is that the preset equalizer functions typically in a low-noise environment. Such is not the case for the adaptive equalizer, where the test signal is effectively buried in the noise of the information signal. In this instance the effective noise determines the response time of the equalizer.

Fig. 5 — Adaptive channel equalizer block diagram.

The implication in many of the equations written thus far is that the integration or averaging necessary for determining the various cross-correlations is carried out over an infinite period of time. In the implementation, however, integration is carried out only for as long as needed to determine cross-correlation to the desired statistical accuracy. A discussion of settling time, then, hinges upon an estimate of just how long an integration time is necessary.

In the implementation of the automatic transversal filter used here (described in greater detail in Ref. 3) the digitally controlled attenuators each can take on $2^M$ values separated by constant increments. $M$ is the number of bits in the attenuator's memory (see Fig. 6). Specifically then, one wants to know how long an integration time is necessary to determine whether or not an undesirable component in the received signal waveform warrants a correction by a $\pm 2/2^M$ portion of the signal from an elemental network. Such a decision must be made with a satisfactorily large probability of being correct.

Thus the problem can be viewed as a problem in signal detection theory. Specifically this question must be answered for the $j$th tap: "Is there or is there not a component of the error signal present which is of the form $2^{1-M} [w(t)^*x(t - j\tau)]$?" This question is

Fig. 6 — Attenuator setting errors.

frequently asked in statistical decision theory.[9] As discussed above in the nonadaptive case the component signal $w(t)^*x(t - j\tau)$ is directly available as the $j$th tap voltage. Thus a maximum likelihood estimator of the component to be measured can be constructed by correlating the tap signal with the error signal. There is a new wrinkle in the adaptive mode, however: although the component signal is still available, it is effectively buried in a very noisy environment (created by the information signal). Thus the problem at hand can be called the detection of a known signal by correlation with a noisy reference signal. The details of the mathematical argument are given in the appendix and essentially follow Helstrom's development but with modifications for the noisy reference case.[9]

In the appendix an approximate relation is obtained which relates the equalizer's maximum residual signal-to-noise ratio to the correlators' integration time. This residual signal-to-noise ratio considers only the error or noise resulting from improper setting of the attenuators. In the analysis there is one noise component resulting solely from the granularity of the attenuators. For an equalizer of specified length and composed of attenuators with $M$-bit memories this residual noise (called the noise floor) is shown as a number of horizontal lines in Fig. 7. In the appendix the increase in this noise resulting from finite correlator integration time (which causes further inaccuracies in the setting of the attenuators) is determined. The resultant total noise as a function of integration time is displayed in Fig. 7 as curved lines. The results are shown for the case of a 19-tap equalizer with information signal and test signal transmitted at the

same power level and a useful channel bandwidth of 2,400 Hz. The increase in noise should be maintained at a small value, for if this is not the case the same residual noise could have been obtained from attenuators with fewer bits accuracy and smaller cost. If the increase in noise is held to, say, ¼ of the value of the noise floor (that is, an increase of 1 dB), the equalizer design is restricted to points on the dashed line in Fig. 7. The probability $P$ that an attenuator is incorrectly set can then be determined from equations (25) and (26) in the appendix by setting the permitted fraction of the noise floor equal to the increase in noise, that is,

$$\frac{1}{4} \frac{(2N + 1)}{(2)^{2M}} = \frac{(2N + 1)}{(2^{M-1})^2} P. \tag{16}$$

For the case stated the probability of incorrect setting should be 1/16, independent of equalizer length or attenuator granularity. Because of the assumed gaussian statistics, the value of $P$ equal to 1/16 determines the ratio of the mean [equation (21)] to the standard deviation [equation (23)]. An approximate relation for the integra-



Fig. 7 — Integration time.

tion time $T$ can then be found for the region of interest, that is,

$$2^{-2[M-1]}N_o < V_o \, .$$

In fact, in this region it follows that

$$T \cong \frac{2^{2M}}{\Omega} \left[ \frac{V_0}{S} + \frac{N_o V_o}{S^2} \right] \tag{17}$$

where $N_o/S$ is the received information to test signal power ratio, $V_o/S$ is the equalized information to test signal power ratio, $\Omega$ is the channel bandwidth, and $M$ is the number of bits in each attenuator's memory.

The noise levels shown in Fig. 7 were calculated for received and equalized reference to information signal ratios of unity. An increase in the information signal to reference noise ratio can be obtained by reducing the level of the transmitted reference signal. Thus, more of the power capacity of the channel can be used for the information-bearing signal. The penalty which must be paid is an increase in the necessary integration time, as shown by equation (17).

It may be wise to emphasize that the argument here has dealt with the time, $T$, required for the equalizer to make its finest adjustment reliably. The most abrupt possible change in the transmission channel's characteristics would force at least one of the equalizer's attenuators to move through its entire range of $2^M$ steps. The time required for this is long, in fact $2^M T$ seconds, ignoring the interaction among attenuator settings. This period could easily be reduced by the incorporation of a scheme for accelerated operation of the equalizer when the reference error signal exceeds some threshold. The equalizer could first run quickly to an approximate setting and then slowly to the exact optimum setting.

III. EXPERIMENTAL IMPLEMENTATION AND PERFORMANCE

In order to see what difficulties might arise in the implementation of the adaptive equalizer described here, such a device was constructed in the laboratory. The heart of the implementation is the automatic transversal filter as shown in Fig. 8. This has been described in considerable detail in Ref. 3 as are the necessary timing and carrier recovery circuits. The function of the timing circuitry is to establish the synchronization of the two pseudorandom sequences while the carrier recovery circuitry compensates for any net carrier offset which may have occurred during transmission. A Bell System *Data-Phone*® data communications set 201B was used as the information transmitter

Fig. 8 — Adaptive equalizer block diagram.

and receiver. This modem happens to be of the four-phase, differentially detected type and is capable of 2400 bits per second transmission on well-equalized (C2 conditioned) private lines. It should be emphasized that the design and utility of the adaptive equalizer is independent of the information transmission format.

In the laboratory experiment the channel was simulated and had a frequency response given by $X(\omega) = 1 - ae^{-i\omega 2\tau}$, where $\tau$ is the tap spacing of the equalizer and is equal to 150 microseconds. The performance of the equalizer on this simulated channel (with the constant $a$ assuming the value of about 0.5) is shown in the series of eye patterns in Fig. 9. Figure 9a shows the closed eye-pattern for the received, distorted signal in the absence of the equalizer. Figure 9b shows the improved eye pattern resulting from the adaptively equalized signal.

Fig. 9 — Eye patterns: (a) received, distorted signal, (b) adaptively equalized signal, test signals removed, (c) adaptively equalized signal, (d) preset-equalized signal.

The correlators used in the adaptive equalizer were designed for use in the preset equalizer. These correlators imposed stringent restrictions on the relative levels of the test and data signals for optimum performance. Specifically the test signal had to be transmitted at a higher level than would be necessary, given correlators with a greater dynamic range. Comparison of the second and third eye patterns shown in Figs. 9b and c shows the resulting noise. The two patterns are identical except that the second includes the effect of the imperfectly matched equalized and locally generated test signals. The final eye pattern (Fig. 9d) shows the results which would be expected from an adaptive equalizer using correlators with greater dynamic range than those now available.

IV. OBSERVATIONS AND CONCLUSION

The experimental performance of the adaptive channel equalizer encourages further development. At the same time, however, the experimental work has indicated several problems which require special attention.

First, the correlators used to establish the attenuator settings in the adaptive channel equalizer perform in a very hostile environment. They must search out small components of a signal which itself is buried in noise; this requires tremendous dynamic range. The observation times involved are long; this imposes constraints on the integrators necessary for performing cross-correlation.

Second, the delay line of length $T_D$ used to prevent domination of the equalizer by the information signal also poses a problem. The operation of the equalizer can be severely impaired by any irregularities in this long delay line. The correlators have no means of compensating for distortion which occurs in this delay, and even an amount of distortion which would be considered small in many applications can make the equalizer virtually useless.

These two observations indicate that considerable advantage might well be obtained in an all-digital implementation. In such a device a small penalty would initially be suffered in analog-to-digital conversion, but this is the only source of distortion. Delay of arbitrary length can easily be obtained by inexpensive shift registers, and correlators can be constructed to arbitrary accuracy.

Third, a substantial component of the noise seen by the information signal receiver is the very small difference between the nearly identical desired and equalized transmitted test signals. Even a small difference

or jitter in timing or carrier phase is thus capable of causing a very large increase in the difference signal or noise seen by the information receiver. Therefore the carrier and timing recovery circuits must be designed with care and the reference signal should be transmitted at as low a level as is practicable. It should also be stated that this technique does nothing to eliminate the effects of nonlinear distortion in the channel.

In conclusion, there are problems which need further attention but the indications are that this technique holds great promise for the continuous adaptive equalization of communication channels where the designer has no knowledge of the information format used on the channel.

APPENDIX

*Calculation of Correlator Integration Time*

Let the signal, which is the input to the $j$th attenuator, be broken down into two parts: the signal $s_j(t)$ resulting from the test signal, and the noise $\eta_j(t)$ resulting from the information signal. As the equalization process nears completion, the correctable components of the error signal become more and more difficult to detect. Finally the $j$th correlator must determine whether or not a component of the error signal of the form $[2/2^M]$ $[s_j(t - T_D)]$ is present. This signal is also buried in noise $\nu(t - T_D)$ which is a result almost exclusively of the information signal. Thus the other waveform delivered to the correlator is $a[2/2^M]$ $[s_j(t - T_D)] + \nu(t - T_D)$ where the variable $a$ is either zero (signal not present) or one (signal present). The $j$th correlator thus computes

$$\rho_j = \left\langle \left\{ \frac{2a}{2^M} [s_j(t - T_D)] + \nu(t - T_D) \right\} \{ s_j(t) + \eta_j(t) \} \right\rangle \qquad (18)$$

to ascertain whether or not the signal is present, that is, whether $a$ is zero or one.

Following Helstrom's outline for the noiseless reference case, it is necessary to calculate the mean and variance of the statistic $\rho_j$ since in the physical system integration is carried out only over a range of 0 to $T$ seconds.[9] The mean and variance of $\rho_j$ suffice because we assume at this point in the argument that the noise created by the information signaling system is gaussian. For convenience we also assume that the noise is white. This is clearly not the case but the search here is only for an estimate of the settling time and an analysis based on the exact statistics of the noise would yield little for the greatly increased effort necessary.

The mean of the various statistics $\rho_j$ is

$$\langle \rho \rangle_{\mathrm{av}} = \int_0^T \left\langle \left\{ \frac{2a}{2^M} \left[ s_i(t - T_D) \right] + \nu(t - T_D) \right\} \{ s_i(t) + \eta_i(t) \} \right\rangle_{\mathrm{av}} dt \quad (19)$$

where $\langle \ \rangle_{\mathrm{av}}$ indicates an ensemble average. Because the noise and signal sources are independent and have zero mean values and further because $T_D$ is chosen large enough to make $\nu(t - T_D)$ and $\eta_j(t)$ independent,

$$\langle \rho \rangle_{\mathrm{av}} = \frac{2a}{2^M} \int_0^T \langle s_i(t - T_D) s_j(t) \rangle_{\mathrm{av}} dt. \quad (20)$$

However, $s_j(t - T_D) = s_j(t)$ so that

$$\langle \rho \rangle_{\mathrm{av}} = \frac{2a}{2^M} \int_0^T \langle s_i^2(t) \rangle_{\mathrm{av}} dt = \frac{2a}{2^M} TS \quad (21)$$

where $T$ is the integration time and $S$ is the power level of the test signal.

The variance of the statistics $\rho_j$ can be shown to be essentially independent of $a$ under the assumptions made and is given by

$$\sigma_\rho^2 = \int_0^T \int_0^T \left\langle \left\{ \frac{2a}{2^M} \left[ s_i(t_1 - T_D) \right] + \nu(t_1 - T_D) \right\} \right.$$
$$\left. \cdot \left\{ \frac{2a}{2^M} \left[ s_i(t_2 - T_D) \right] + \nu(t_2 - T_D) \right\} \right\rangle_{\mathrm{av}}$$
$$\cdot \langle [s_i(t_1) + \eta_i(t_1)][s_i(t_2) + \eta_i(t_2)] \rangle_{\mathrm{av}} \, dt_1 \, dt_2 \ . \quad (22)$$

After some manipulation the nonzero terms are

$$\sigma_\rho^2 = \left[ \frac{2a}{2^M} \right]^2 \int_0^T \int_0^T \langle s_i(t_1) s_i(t_2) \eta_j(t_1) \eta_j(t_2) \rangle_{\mathrm{av}} \, dt_1 \, dt_2$$

$$+ \int_0^T \int_0^T \langle \nu(t_1)\nu(t_2)s_j(t_1)s_j(t_2)\rangle_{\mathrm{av}} \, dt_1 \, dt_2$$

$$+ \int_0^T \int_0^T \langle \nu(t_1)\nu(t_2)\eta_j(t_1)\eta_j(t_2)\rangle_{\mathrm{av}} \, dt_1 \, dt_2 \; .$$

$$\sigma_\rho^2 = \left[\frac{2a}{2^M}\right]^2 \int_0^T \int_0^T \langle s_j(t_1)s_j(t_2)\rangle_{\mathrm{av}} R_\eta(t_1 - t_2) \, dt_1 \, dt_2$$

$$+ \int_0^T \int_0^T \langle s_j(t_1)s_j(t_2)\rangle_{\mathrm{av}} R_\nu(t_1 - t_2) \, dt_1 \, dt_2$$

$$+ \int_0^T \int_0^T R_\nu(t_1 - t_2) R_\eta(t_1 - t_2) \, dt_1 \, dt_2 \; .$$

$R_\nu$ and $R_\eta$ are autocorrelation functions. After evaluation the relation for $\sigma_\rho^2$ is of the following form:

$$\sigma_\rho^2 \cong \left\{\left[\frac{2a}{2^M}\right]^2 \frac{2N_{j0}}{\Omega} + \frac{2V_0}{\Omega}\right\}TS + \frac{2N_{j0}V_0 T}{\Omega}. \tag{23}$$

$N_{j0}$ is the power level of the signal $\eta_j(t)$ and $V_0$ is the power level of the signal $\nu(t)$; the effective channel bandwidth is $\Omega$. Had the reference not been noisy, only the variance of the statistic $\rho_j$ would be changed and would in fact be

$$\sigma_\rho^2 = \frac{2V_0}{\Omega} \, (TS) \tag{24}$$

which agrees with the standard result.[9]

With the assumptions mentioned above, equations (21) and (23) completely specify the statistics necessary to calculate the probability of making a correct decision.

Figure 6 helps explain the operation of the system. The range of the attenuator $(-1, +1)$ is divided into $2^M$ equal increments of width $2^{-(M-1)}$. A worst case assumption that the desired attenuator setting is midway between two possible settings is shown in Fig. 6 and leads to a minimum attainable signal-to-noise ratio (resulting from the granularity of the attenuator) of

$$\frac{S_\eta}{N_S} = 10 \log_{10} \frac{(2)^{2M}}{(2N + 1)} \, \mathrm{dB} \tag{25}$$

where $(2N + 1)$ is the number of attenuators. Relation (25) was obtained in Ref. 3 and is the ratio of the information signal power to the test-signal-generated noise, assuming that the two signals are

transmitted at the same level. This is approximately the signal-to-noise ratio seen by the information receiver.

Incorrect setting of the attenuator decreases the effective signal-to-noise ratio from this value. Referring to Fig. 6, assume that the attenuator is set at level $c$. At the end of the correlator integration period, the correlator output, $\rho_j$, is sampled. If $\rho_j$ is positive, the attenuator setting increases to level $d$; if $\rho_j$ is negative, the setting decreases to level $b$. Assume that the optimum attenuator value corresponds to level $P$, which is not achievable because of the attenuator granularity. Then the expected value of $\rho_j$, $\langle \rho_j \rangle_{av}$, is positive by an amount proportional to the difference between levels $P$ and $c$. If the actual attenuator settings fluctuate between levels $c$ and $d$, the distortion is that given by equation (25); if, however, the attenuator setting assumes the value of level $b$, the result is an increase in the distortion power of the amount $V_o/(2^{M-1})^2$. The probability that the attenuator setting goes to level $b$ is then the probability that $\rho_j$ is less than zero. Since $\langle \rho_j \rangle_{av}$ is positive in this example, and assuming gaussian statistics as above, this probability of incorrect attenuator setting is given by the relation

$$P = \text{erfc} \left( \langle \rho_i \rangle_{av}/\sigma_P \right).$$

Further, assuming independence of the attenuator weights, the increase in noise to the information signal is

$$(2N + 1)(P) \frac{V_0}{(2^{M-1})^2}. \tag{26}$$

The results of the evaluation of equation (26) are plotted in Fig. 7 in various regions of interest. For Fig. 7 a channel bandwidth of 2,400 Hz is assumed.

REFERENCES

1. Rudin, H. R. Jr., "Automatic Equalization Using Transversal Filters," IEEE Spectrum, 4, No. 1 (January 1967), pp. 53–59.
2. Lucky, R. W., "Automatic Equalization for Digital Communication," B.S.T.J., 44, No. 4 (April 1965), pp. 547–588.
3. Lucky, R. W., and Rudin, H. R., Jr., "An Automatic Equalizer for General-Purpose Communication Channels," B.S.T.J., 46, No. 9 (November 1967), pp. 2179–2208.
4. Lucky, R. W., "Techniques for Adaptive Equalization of Digital Communication," B.S.T.J., 45, No. 2 (February 1966), pp. 255–286.
5. Narendra, K. S., and McBride, L. E., "Multiparameter Self-Optimizing Systems Using Correlation Techniques," IEEE Trans. on Automatic Control, AC-12, No. 1 (January 1967), pp. 53–59.
6. Widrow, Bernard, "Adaptive Filters I: Fundamentals" Stanford Electronics Laboratory Report, SEL-66-126, December 1966.

7. Rudin, H. R., Jr., "Automatic Filter Synthesis Under A Frequency-Weighted Mean-Squared Error Criterion," Proc. First Hawaii Int. Conf. on Syst. Sci., Honolulu, Hawaii, January 1968, pp. 213–215.
8. Narendra, K. S., and Streeter, D. N., "An Adaptive Procedure for Controlling Undefined Linear Processes," IEEE Trans. Automatic Control, AC-9, No. 10 (October 1964), pp. 545–548.
9. Helstrom, C. W., Statistical Theory of Signal Detection, New York: Pergamon Press, 1960 (see especially pages 84–88).

# An Imaging System Exhibiting Wavelength-Dependent Resolution

## By S. Y. CHAI

(Manuscript received January 20, 1969)

*The Christiansen filter consists of particles of a transparent solid immersed in a suitable liquid such that the two materials have different dispersions but equal indexes of refraction at some specified wavelength. This paper analyzes such a filter as an optical device for providing wavelength-dependent resolution.*

*We derive the modulation transfer function of this filter by using the Huygens–Fresnel diffraction principle. General agreement is obtained with experimental measurements made using monochromatic light.*

*One proposed single-tube color camera system for picture telephone service requires the three primary color component images to have different resolutions. An optical filter illustrating the characteristics required by this application was designed using the formulas developed in this paper. Subjective evaluation of the experimental results obtained with this filter indicates no unforeseen degradation of the composite image of the primary color components.*

## I. INTRODUCTION

When homogeneous isotropic particles, such as crushed optical glass, are immersed in a liquid which has, at a specified wavelength, the same index of refraction as the particles but a different dispersion, light at the specified wavelength is unaffected in passing through the mixture. On the other hand, light of other wavelengths is dispersed. A cell containing such a mixture of solid particles in an appropriate liquid constitutes a Christiansen filter (see Figure 1).[1-4]

The wavelength dependent action of the filter can be used to construct a low-pass spatial filter with a wavelength dependent cut-off. In particular, with the proper glass particles and liquid, the filter acts as weakly diffusing ground glass for the blue and red images while not affecting the green images.

An optical device using a color-selective spatial low-pass filter of this type is desirable in order to realize a proposed single pickup-tube color camera.[5,6] In this camera, color separation is achieved by modulating the amplitude of the color signals onto different carrier frequencies. These carrier frequencies are in turn generated by striped filters that have different spatial frequencies for the different primary colors. It is necessary to reduce the spatial bandwidth of the composite color image in order to meet the resolution limitation of the camera tube.

This bandwidth reduction can be used because the human eye does not readily discriminate between composite images composed of three high resolution primary color images, and a composite image composed of a high resolution green image and low resolution red and blue images.

By considering the filter as a random phase-changing screen, one can determine its spatial frequency response as a function of glass particle size, thickness of the filter, and position of the filter in relation to the image plane.[4,7]

The quantitative understanding thus gained is used to design an optical filter with characteristics which are roughly those which would be required by an experimental color picture telephone system. The experimentally observed effect of this filter indicated no unforeseen degradation of the subjective quality of the final composite image transmitted through it.

## II. OPTICAL TRANSFER FUNCTION OF THE FILTER SYSTEM

Figure 1 shows an optical system transforming an object plane into an image plane. If the object plane is incoherently illuminated, the system is linear in intensity and therefore:

$$i(x, y) = \iint o(x', y')h(x - x', y - y') \, dx' \, dy' \qquad (1)$$

where $i(x, y)$ and $o(x', y')$ are the intensity distribution in the image plane and object plane, respectively; and the point spread function $h$ of the optical system represents the intensity distribution in the image plane for a point object.

In terms of spatial frequency $(f_x, f_y)$ response, we have

$$I(f_x, f_y) = O(f_x, f_y)H(f_x, f_y) \qquad (2)$$

where the optical transfer function $H(f_x, f_y)$ is the Fourier transform

Fig. 1 — Optical system with Christiansen filter.

of the point spread function $h$. If the transfer function is real and only the modulus is used, it is then called the modulation transfer function.[8]

The point spread function $h$, by definition, is:

$$h(x, y) = u(x, y) \cdot u(x, y)^* \qquad (3)$$

where $u(x, y)$ is the complex amplitude distribution of a point object in the image plane. This is called the "point image function," and $u(x, y)^*$ is its complex conjugate.

Using the Huygens-Fresnel diffraction principle, the point image function $u(x, y)$ can be expressed in terms of the complex field distribution $A(\xi, \eta)$ over the exit plane of the filter (see Fig. 1).[9-11] We should pay particular attention to the fact that we are assuming a point object and an isoplanatic region. That is, the optical system is assumed linear such that equation (1) holds throughout the entire object plane. Thus, for a point object, the point image function is given by

$$u(x, y) = \frac{1}{j\lambda} \frac{e^{jkd}}{d} \iint_{-\infty}^{+\infty} A(\xi, \eta) \exp\left\{\frac{jk}{2d}[(\xi - x)^2 + (\eta - y)^2]\right\} d\xi \, d\eta.$$

$$(4)$$

The field distribution in the exit plane of a thin Christiansen filter can be written as

$$A(\xi, \eta) = E_o(\xi, \eta)e^{j\Phi(\xi, \eta)} \cdot \exp\left[-\frac{jk}{2d}(\xi^2 + \eta^2)\right] \qquad (5)$$

where $\Phi(\xi, \eta)$ is a random phase function produced by scattering (see

the appendix), and the term $E_o(\xi, \eta) \exp[-jk/2d \ (\xi^2 + \eta^2)]$ represents the spherical wavefront produced by the lens. Since the field outside the lens aperture is suppressed,

$$E_o(\xi, \eta) = \begin{cases} \text{constant} & (\xi^2 + \eta^2)^{\frac{1}{2}} \leqq b \\ 0, & (\xi^2 + \eta^2)^{\frac{1}{2}} > b \end{cases}$$

and $b$ is proportional to the lens aperture as shown in Fig. 1.

Since the phase function $\Phi(\xi, \eta)$ was assumed to be a two-dimensional random process with isotropic fluctuation, and the experimental test object contains only one-dimensional variations in the transmission, the problem can be reduced to one dimension for simplicity. Thus, by neglecting the constant factor outside the integrals, we may rewrite equation (4) as

$$u(x) = \int_{-\infty}^{+\infty} E_o(\xi) e^{j\Phi(\xi)} \exp\left(-\frac{jk}{2d}\xi^2\right) \cdot \exp\left\{\frac{jk}{2d}(\xi - x)^2\right\} d\xi$$

$$= \exp\left(\frac{jk}{2d}x^2\right) \int_{-\infty}^{+\infty} E_o(\xi) e^{j\Phi(\xi)} \exp\left(-\frac{jk}{d}\xi x\right) d\xi. \tag{6}$$

The optical transfer function $H(f_x)$ of the system is now obtained by taking the Fourier transformation of $u(x)u(x)^*$ to give

$$H(f_x) = \int_{-\infty}^{+\infty} u(x)u(x)^* \exp(j2\pi f_x x) \, dx$$

$$= \int_{-\infty}^{+\infty} E_o(\xi) e^{j\Phi(\xi)} \cdot \{E_o(\xi - d\lambda f_x) \exp[j\Phi(\xi - d\lambda f_x)]\}^* d\xi. \tag{7}$$

Normalizing (dropping the subscript $x$) this equation results in

$$H(f) = \frac{\int_{-\infty}^{+\infty} E_o(\xi) e^{j\Phi(\xi)} \{E_o(\xi - d\lambda f) \exp[j\Phi(\xi - d\lambda f)]\}^* d\xi}{\int_{-\infty}^{\infty} |E_o(\xi)|^2 d\xi}. \tag{8}$$

The statistics of the phase fluctuation $\Phi(\xi)$ are assumed to be such that the average value with respect to $\xi$ can be replaced by an ensemble average. By denoting an ensemble average by $\langle \ \rangle_{av}$, the normalized transfer function for this case becomes

$$H(f) = D(f)\langle \exp\{j[\Phi(\xi) - \Phi(\xi - d\lambda f)]\} \rangle_{av} \tag{9}$$

where

$$D(f) = 1 - \frac{d\lambda f}{2b} \cong 1.$$

Because the range of the spatial frequency $f$ under discussion is much lower than the cutoff frequency produced by diffraction at the lens aperture, $d\lambda f$ is much less than $2b$. Therefore, $D(f)$ can be approximated by unity.

Under the condition that $\Phi(\xi)$ is a random gaussian variable of zero mean and variance $\Phi_o^2$, equation (9) is shown to be[12]

$$H(f) = \exp\{-\Phi_o^2[1 - \rho(d\lambda f)]\} \qquad (10)$$

where $\rho(d\lambda f)$ is the autocorrelation function for the phase fluctuations at the exit plane of the filter and is expressed as

$$\rho(d\lambda f) = \frac{\displaystyle\int \Phi(\xi)\Phi(\xi - d\lambda f)\, d\xi}{\displaystyle\int \Phi(\xi)^2\, d\xi}. \qquad (11)$$

Now, the autocorrelation function $\rho(d\lambda f)$ must be determined. The behavior of the function $\rho$ is complicated; however, it is known that

$$\rho(d\lambda f) = 1 \qquad \text{when} \qquad d\lambda f = 0$$

and

$$\rho(d\lambda f) = 0 \qquad \text{when} \qquad \begin{aligned}&d\lambda f \text{ is greater than the linear}\\&\text{dimension of the scattering}\\&\text{particles.}\end{aligned}$$

In this case there are many possible types of behavior that the function $\rho(r)$ may exhibit. Two common types are.[13–16]

$$\rho(r) = \exp(-|r/q|) \qquad (12a)$$

and

$$\rho(r) = \exp[-(r/q)^2] \qquad (12b)$$

where the $q$ represents a correlation length.

These are both particularly important forms to us because we are interested primarily in small $r$. Thus, equations (12a) and (12b) reduce to

$$\rho(r) \cong 1 - \left|\frac{r}{q}\right| \qquad (13a)$$

and

$$\rho(r) \cong 1 - \frac{r^2}{q^2}. \qquad (13b)$$

The autocorrelation function for any particles, regardless of shape, reduces to one of these two forms for small $r$. In our experiments we assumed spherical particles and equation (13a) is found to be the more realistic approximation. Thus, let us assume[15,16]

$$\rho(d\lambda f) = \exp \left\{ - \left| \frac{d\lambda f}{K2a} \right| \right\} \tag{14}$$

where the diameter $2a$ of the particle indicates the order of the size of the irregularities. The constant $K$ is a factor expressing the actual correlation length to the diameter of the particle and can be obtained from the function $\rho(d\lambda f)$ which is empirically determined. Only the small values of $d\lambda f$ are of interest so that by approximating equation (14)

$$\rho(d\lambda f) \cong 1 - \left| \frac{d\lambda f}{K2a} \right|. \tag{15}$$

Finally, substituting equations (15) and (27) into equation (10), the modulation transfer function obtained is

$$H(f) = \exp \left\{ - \frac{3.29CL \, \triangle n^2 \, df}{K\lambda} \right\} \tag{16}$$

where $C$ is the ratio of particle volume to total volume, $L$ is the thickness of the filter, and $\triangle n$ is the difference in refractive index between particles and liquid.

Notice that when $f$ is very large so that $\rho(d\lambda f) \cong 0$, the modulation transfer function reaches a constant value of

$$H(f) = \exp \{ -\Phi_o^2 \} = \exp \left\{ - \frac{6.58CLa \, \triangle n^2}{\lambda^2} \right\}. \tag{17}$$

This means that there is always finite transmission at high spatial frequencies, but that the loss can be made as large as desired by controlling the parameters of equation (17).

## III. EXPERIMENTAL

### 3.1 Preparation of Filters

The filters were made of crown glass type K5 ($N_d = 1.52320$) in a solution of ethyl salicylate at room temperature. This combination provides a nearly transparent filter in the green region as shown in Fig. 2. The dispersion curves were determined by a Carl Zeiss, Model A Abbe-refractometer.

Fig. 2 — Dispersion curves of optical glass (*K5*) and ethyl salicylate.

The glass particles were prepared by pulverizing them with a ball mill; they were then cleaned with aqua regia and distilled water. The size of the particles was estimated by averaging the dimensions of the wire screen used in sorting. The radii of the prepared samples ranged from $37\mu$ to $44\mu$, $44\mu$ to $62\mu$, $75\mu$ to $125\mu$, and $125\mu$ to $150\mu$. The cells were made by epoxy-cementing glass plates to 1-inch diameter glass tubing of lengths 2mm, 2.5mm, 3.5mm, and 4mm.

The cells were about half filled with the liquid, then the glass particles were poured in slowly so that air bubbles were not carried down. The volume of the glass particles was determined from their weight and known density of 2.59 grams per cubic centimeter. The volume ratio $C$ was found to be approximately 0.42, independent of particle size.

### 3.2 *Measurement of Modulation Transfer Function*

The apparatus used to measure the transfer function of the optical system has been described previously.[17] Figure 3 is a rough sketch of this apparatus. A target film with sinusoidal intensity distribution was used as a test object. The film was imaged through the Christiansen filter onto a slit in front of the photomultiplier. Rotation of the drum containing the test film caused the image to move across the sampling slit. The sinusoidal gratings of the test film varied logarithmically from 4 to 300 cycles per inch and were arranged on the target film in sequence to allow rapid measurement or direct recording of the modulation transfer function.

Fig. 3 — Apparatus for measuring the modulation transfer function of the filter.

The measurements were usually carried out in the region of the blue mercury arc light to maximize the response of the S-11 photomultipler. The slit width of the photomultiplier was 100$\mu$, which corresponds to a sampling interval equal to one-third the period of highest frequency. An $f/2.3$ lens of 5-cm focal length was used to image the target film onto the slit. The target-lens spacing was 6 cm and the lens-slit spacing was 30 cm. This resulted in a magnification of 5 and gave a range of spatial frequencies from about 0.03 cycles per millimeter to 3 cycles per millimeter. This range of spatial frequencies is of interest for potential color picture telephone camera applications.

The target film was first focused with a cell filled with only the liquid placed between the lens and the slit. This compensates for the increased optical path length resulting from the nonunity index of refraction of the Christiansen filter.

IV. RESULTS AND DISCUSSION OF THE MODULATION TRANSFER FUNCTION

The modulation transfer functions of the filters were measured for various positions of the filters as well as for various thicknesses and particle sizes of the filter. The position of the filter was defined by the distance $d$ between the slit plane and the inside of the exit glass plate of the cell (see Fig. 1).

The experimentally determined modulation transfer functions, for

several different filter positions with a filter thickness of 2.5 mm and particle radius of 100 $\mu$, are plotted in Fig. 4. Figure 5 shows the variation of the modulation transfer functions with filter thickness for a constant particle radius of 100 $\mu$ and a $d$ of 1/8-inch. From Fig. 4, the linear dependence of the exponent $\Phi_o^2(1 - \rho)$ on $d$ can be demonstrated, thus validating equation (10). In addition, Fig. 5 shows the approximately linear dependence of $\Phi_o^2$ on $L$, indicating the validity of equation (27). From these results, the assumption of equation (15) for the auto-correlation function of the phase irregularities also appears reasonable.

The dependence of the modulation transfer function on particle size is shown in Fig. 6. From this figure, the factor $K$ representing the ratio of the particle diameter to correlation length [see equation (14)] was obtained and plotted in Fig. 7 as a function of the particle size $a$. A quantitative analysis of the effect of the particle size on the transfer function is unavailable because the exact expression for the phase autocorrelation function in terms of the particle size is not known. It is believed that the empirically determined factor $K$ provides a qualitative dependence of the particle size on the phase autocorrelation function because only small values of $d\lambda f/2aK$ [see equation (14)] are of interest. A theoretical analysis of the phase autocorrelation function may not be of great value unless the exact shape of the particle is considered.



Fig. 4 — Experimental modulation transfer functions for various filter positions ($a = 100\mu$, $d = 2.5$ mm, $\lambda = 4358$ Å).

Fig. 5 — Experimental modulation transfer functions for various filter thickness ($a = 100\mu$, $d = 3.18$ mm, $\lambda = 4358$ Å).



Fig. 6 — Experimental modulation transfer functions for various particle sizes of the filter ($L = 2.5$ mm, $d = 3.18$ mm, $\lambda = 4358$ Å).

Fig. 7 — Factor $K$ of phase correlation length versus particle sizes.

The linear dependence of the factor $K$ on the particle size $a$, shown in Fig. 7, cannot be extended much further than the radius of about $150\mu$ because the gaussian phase fluctuation imposes an upper limit on $a$ for which the central limit theorem is applicable (see the appendix). In practice, it is not possible to determine the exact correlation length from Fig. 7, since it is impossible to obtain a large value of the correlation length $2aK$ with sufficient accuracy from extrapolation of the small values of $d\lambda f$.

In order to simulate a situation of independent resolution control of primary color images, measurements of the transfer functions at the blue, green, and red regions were made with a Christiansen filter 2.5mm thick and with a $100\mu$ particle radius. In this case, the filter was set $\frac{1}{8}$-inch away from the image plane. As light sources of primary colors, a mercury arc lamp was filtered with Optics Technology, Inc., monopass filters, No. 433 for blue (4350 Å), No. 566 for green (5600 Å), and No. 633 for red (6250 Å).* Figure 8 shows the results. The modulation transfer function for the imaging lens itself also is plotted in Fig. 8 to show the effect of the lens aberrations.

---

* These narrow band filters were used in order to give quantitative agreement with theory. In its proposed application to resolution control the true "wide band" primaries should be used. The evaluation of this filter for resolution control in a TV color camera is more involved and is the subject of the following sections.

Fig. 8 — Experimental modulation transfer functions for different wavelengths and for imaging lens with $a = 100\mu$, $L = 2.5$ mm, and $d = 3.18$ mm.

Figure 9 shows the comparison between the experimental results at different wavelengths and the theoretical modulation transfer function given by equation (16). The values of $\Delta n$ used in this equation, $1.5 \times 10^{-2}$ at $\lambda = 4350$ Å, $10^{-3}$ at $\lambda = 5600$ Å and $6 \times 10^{-3}$ at $\lambda = 6250$ Å, were obtained from Fig. 2. The factor $K$ of 1.82 obtained from Fig. 7 and the measured value $C$ of 0.42 are also used in equation (16). The experimental results are plotted from Fig. 8 after correcting for the effects of the aberrations caused by the imaging lens. This correction becomes necessary at the higher spatial frequencies encountered with red and green light.

The deviation between the experimental and analytical results probably result from thermally induced changes in the refractive index of ethyl salicylate, which has a thermal coefficient of index of refraction of about $5 \times 10^{-4}$ per C°. A slight change in $\Delta n$ causes a significant change in the analytical results for the green region, since $\Delta n$ for green is a small value. On the other hand, the effect of a small change in $\Delta n$ has very little effect on the calculated blue or red responses. Choosing the correct value of $\Delta n$ to be used in equation (16) is further complicated by the 100 Å spectral bandwidth of the primary colors used.

The theoretical and experimental results shown in Fig. 9 agree sub-

stantially considering the many approximations made. These approximations include the representation of the actual particles by spheres with no variation in particle size and assuming a filter thin enough to allow the difference in path length through the filter between normally and obliquely incident rays to be neglected.

## V. APPLICATION OF THE FILTER

As Section IV shows, the cutoff frequency of the modulation transfer function (spatial frequency response) of the Christiansen filter can be controlled by the filter thickness, the particle size, and the position of the filter with respect to the image plane. Using this quantitative understanding, a filter for independently controlling the resolutions of the three superposed primary color images was designed. The operation of the filter and its suitability for the intended application were then experimentally verified.

In order to evaluate the filter in its proposed application to the color picture telephone system, the true NTSC primary color "taking" filters were used.

### 5.1 *Evaluation of the Filter*

For the color camera, the desirability of having the shortest possible distance between the filter and the image plane (that is camera



Fig. 9 — Comparison of experimental and theoretical modulation transfer functions for selected wavelengths with $a = 100\mu$, $L = 2.5$ mm, and $d = 3.18$ mm.

tube target) dictates choosing a glass particle with about a $40\mu$ radius. The filter had a circular cross section 20 mm in diameter, and it was 2.5 mm thick.

The experimental arrangement used is shown in Fig. 10. Incandescent (white) light illuminated a ground glass diffuser which was positioned in front of the 35 mm color slide used to simulate a real object. In order to be able to separate the three primary color images, the incandescent (white) light was actually derived by adding the primary colors with the half-silvered mirrors as shown in Fig. 10. The primary colors were generated by projecting tungsten filament light through the appropriate Wratten filters: No. 25 for red, No. 47B for blue, and No. 58 for green.

A Bell and Howell $f/1.3$ Angeniux television camera lens with a 15 mm focal length was used as an object lens. The image size was about 10 mm by 10 mm to simulate the raster size of a typical picture telephone camera. For viewing, the image formed at the image plane (that is, vidicon tube target) was magnified about ten times by a second lens. This magnified view of the projected image was used to evaluate the filter.

The Christiansen filter was positioned in front of the image plane; the distance between the exit plane of the filter and the image plane was adjusted to obtain the desired spatial frequency response.

In color picture telephone operation, a figure for picture bandwidth has not yet been determined. However, at the present time, approximations of 1 MHz for the green, 500kHz for the red, and 300kHz for the blue images may be assumed for a single pickup camera system. By assuming a horizontal line scan time of 100 $\mu$s and a raster



Fig. 10 — Experimental arrangement for evaluation of the Christiansen filter as a resolution controlling device.

size of 12 mm (horizontal), the relationship between the image resolution and required electrical bandwidth is determined.

The spatial frequency response was measured by using a pie-chart test target disk with 30-cycle gratings along the circumference. Since the spatial frequency of the grating varies continuously along the radial direction of the disk, the cutoff frequency of the image can be measured by knowing the point where resolution fails. The image size of the test target was 7.5 mm in diameter.

In order to determine the effectiveness of the filter as a resolution controlling device, an image without the filter was viewed for comparison. In this case the Christiansen filter was replaced by a cell the same size as the filter, filled with only the ethyl salicylate solution. This maintains the same optical path lengths within the imaging system.

### 5.2 *Image Quality through the Filter*

In order to determine the resolution, each pictoral subject imaged through the system was associated with the corresponding image of the resolution test pattern. This was done by replacing the object slide with the test target object when photographing the system output. The resolutions were determined simply by observing the images of the test pattern resolution chart. The spatial frequency of the test image was 1.27 cycles per mm at the outer edge of the circular pattern. This spatial frequency corresponded to 153 kHz by assuming 100 microseconds for the line scan time and 12 mm for the raster size of the camera tube.

Figure 11a is a photograph of the image formed by incandescent (white) light without the Christiansen filter; Figure 11b is the corresponding photograph of the image of the resolution test pattern. Notice that the same quality images were obtained with separate illumination of each primary color since achromatic lenses were used in the experiment.

Figure 12a is a photograph of the image formed by incandescent (white) light with the Christiansen filter placed 1 mm in front of the image plane. Figure 12b is the corresponding image of the resolution test chart. The distance of 1 mm between the filter and the image plane was chosen to demonstrate the particular case for which the filter will be used in the proposed color camera system. Some of the dark spots in the photographs were caused by the presence of dirt in the filter. Blurring of the image by low-pass filtering is noticeable.

Fig. 11 — (a) Portrait and (b) resolution test chart showing images formed by incandescent (white) light without the Christiansen filter.

This was introduced deliberately by the Christiansen filter and should not be considered as a defect of the optical system. The quality of the image seems acceptable in all other respects.

In order to illustrate to what extent each primary color image in the composite image (Fig. 12) is spatially filtered, three separate



Fig. 12 — (a) Portrait and (b) resolution test chart showing images formed by incandescent (white) light with the Christiansen filter inserted in the optical system.

photographs were taken with each primary color by blocking the other two color sources. Figures 13a and b show photographs of the images formed with green primary light only. The image would be as sharp as that in Fig. 11 if the filter was completely transparent at green. However, the spectral bandwidth of the NTSC green primary color is wider than the matched bandwidth of the filter in green, so that the quality of Fig. 13 is not quite that of Fig. 11. This illustrates that it is advantageous to limit that bandwidth of the green image so that it does not extend beyond the required minimum value.

Figures 14a and b show the images formed with red primary light only. Blurring of the image caused by low-pass spatial filtering by the Christiansen filter is noticeable. The cutoff frequency of the image, determined from Fig. 14b, is about 4 cycles per mm (the resolution falls about one-third of the radius of the pattern) corresponding to a 500 kHz electrical signal.

Similarly, Figs. 15a and b show photographs of the images formed with blue primary light. Blurring of the image is much more notice-able than that of Fig. 14, as expected from characteristics of the Christiansen filter. The cutoff frequency of the image is about 2.5 cycles per mm (the resolution falls about half of the radius of the pattern) corresponding to a 300 kHz electrical signal.

It should be emphasized that although the blurring of the red and blue images is very noticeable when they are viewed alone (as in Figs. 14 and 15), the composite image of all three primaries is of much better quality. This indicates that a significant amount of the detail in the composite image is provided by the green component.[18]



Fig. 13 — Images formed by green primary component color with the filter inserted in the optical system.

Fig. 14 — Images formed by red primary component color with the filter inserted in the optical system.

In order to illustrate the subjective quality of a composite color image transmitted through this system, Figs. 11a and 12a are printed in color as Figs. 16 and 17, respectively. These photographs were taken on Ektacolor L film.

Much work remains to be done to determine precisely what resolution values should be used, but such a study is beyond the scope of this paper. The point we are trying to make is that we can control the resolution of the three primary color component images and that there is no appreciable degradation beyond that deliberately introduced in the form of spatial bandwidth reduction. Although this system does not have a sharp cutoff characteristic, that is, it does not have a rectan-



Fig. 15 — Images formed by blue primary component color with the filter inserted in the optical system.

gular passband [see equation (16)], it can nevertheless reduce the image detail of a selected color component by any desired amount, while still permitting coarser image structure to remain.

## VI. CONCLUSIONS

Formulas for the modulation transfer function of the Christiansen filter have been derived. The filter thickness, particle size, position of the filter from the image plane, and difference between the particle and liquid indices of refraction have all been shown to control the cutoff frequency but not the shape of the modulation transfer function. The modulation transfer functions as measured with essentially monochromatic sources were in quantitative agreement with those predicted from theory.

The proposed use of the Christiansen filter to reduce the red and blue color detail in a composite image is illustrated using the true "wide band" primary colors. Subjective evaluation of the composite image shows that there is no appreciable degradation beyond that deliberately introduced in the form of spatial bandwidth reduction. The result confirms the applicability of the filter in color picture telephone-like schemes as a means of independently controlling the spatial resolution of the component color images.

## VII. ACKNOWLEDGMENTS

## APPENDIX

### Random Phase Variation in the Christiansen Filter

When a plane wave is normally incident on the entrance plane of the Christiansen filter, the wave emerges from the filter with the phase varying across its wave front but with no appreciable change in amplitude. The particle size is on the order of hundreds of a wavelength; the difference in refractive index between the particles and liquid is a few parts per thousand. Therefore, the distribution of phase over the emerging wave-front can be described approximately in terms of the distribution of the path length of each ray through the particles.

Since the particles are close-packed, it is assumed that the average

number of particles traversed by a ray as it crosses the filter is uniform. The particles are considered to be spherical with constant radii. The statistical fluctuation in the phase across the exit plane of the filter is caused by the ray passing through different parts of the particles.

Assume the ratio of particle volume to total volume is $C$ and the cross-sectional area of the filter is $S$. The average number of particles intersected by a single ray which passes through a filter of the thickness $L$ can then be calculated as

$$N = \left(\frac{CSL}{\frac{4}{3}\pi a^3}\right)\left(\frac{\pi a^2}{S}\right) = \frac{3}{4}\frac{CL}{a} \tag{18}$$

where $a$ is radius of the particle. The first parenthesis represents the total number of particles in the filter.

If the ray is traveling in the $z$ direction in the Cartesian coordinates system $(x, y, z)$, the ray path length through an individual particle



Fig. 16 — Ektacolor print of Fig. 11a.

Fig. 17 — Ektacolor print of Fig. 12a.

$l$ will be a function of the lateral displacement of the ray from the particle center in direction $x$ and $y$. The average path length through a particle is

$$\langle l \rangle_{av} = \int_{-a}^{a} \int_{-(a^2-y^2)^{\frac{1}{2}}}^{(a^2-y^2)^{\frac{1}{2}}} 2(a^2 - x^2 - y^2)^{\frac{1}{2}} P(x, y) \, dx \, dy \qquad (19)$$

where $P(x, y)$ is the probability density function. The lateral displacements in directions $x$ and $y$ are uniformly distributed random variables in the range of the area $\pi a^2$. Therefore, $P(x, y)$ is $1/\pi a^2$ and equation (19) becomes

$$\langle l \rangle_{av} = \int_{-a}^{a} \int_{-(a^2-y^2)^{\frac{1}{2}}}^{(a^2-y^2)^{\frac{1}{2}}} 2(a^2 - x^2 - y^2)^{\frac{1}{2}} \frac{1}{\pi a^2} \, dx \, dy = \tfrac{4}{3}a. \qquad (20)$$

This result differs from that obtained by R. H. Clarke, whose probability density function (equation A4 of Ref. 4) depends upon one of the spatial coordinates.

The second moment of $l$ is

$$\langle l^2 \rangle_{\text{av}} = \int_{-a}^{a} \int_{-(a^2-y^2)^{\frac{1}{2}}}^{(a^2-y^2)^{\frac{1}{2}}} 4(a^2 - x^2 - y^2) \frac{1}{\pi a^2} \, dx \, dy = 2a^2. \qquad (21)$$

Therefore, the variance of the random variable $l$ is

$$\langle l^2 \rangle_{\text{av}} - \langle l \rangle^2_{\text{av}} = 2a^2 - (\tfrac{4}{3}a)^2 = \tfrac{2}{9}a^2. \qquad (22)$$

If the plane wave $E_o(x, y)$ (traveling in direction $z$) is normally incident to the entrance plane of the filter, the field distribution at the exit plane of the filter can be expressed by the equation

$$E(x, y) = E_o(x, y) \exp \{j\Phi(x, y)\} \qquad (23)$$

where $\Phi(x, y)$ is a random phase function.

The phase of the rays at the exit plane of the filter with reference to the average phase across the plane can be written as

$$\Phi = \sum_{m=1}^{N} k \, \triangle n(l_m - \langle l_m \rangle_{\text{av}}) \qquad (24)$$

where $k = 2\pi/\lambda$ ($\lambda$ is a wavelength in free space) and $\triangle n$ is the difference in refractive index between particles and liquid. The process corresponding to each $m$ is an independent random variable so that the summation according to the central limit theorem should be gaussian. This gives

$$P(\Phi) = \frac{1}{(2\pi)^{\frac{1}{2}}\Phi_o} \exp\left(-\Phi^2/2\Phi_o^2\right). \qquad (25)$$

Thus, the variance of the phase at the exit plane is

$$\Phi_o^2 = \left\langle \left[ \sum_{m=1}^{N} k \, \triangle n(l_m - \langle l_m \rangle_{\text{av}}) \right]^2 \right\rangle_{\text{av}} = N(k \, \triangle n)^2 \langle (l - \langle l \rangle_{\text{av}})^2 \rangle_{\text{av}}$$

$$+ N(N-1)(k \, \triangle n)^2 \langle (l - \langle l \rangle_{\text{av}}) \rangle^2_{\text{av}}. \qquad (26)$$

Noting that the second term is zero and using equations (18) and (22), equation (26) becomes

$$\Phi_o^2 = N(k \, \triangle n)^2 \tfrac{2}{9}a^2 = \frac{6.58 CLa \, \triangle n^2}{\lambda^2}. \qquad (27)$$

REFERENCES

1. Christiansen, C., "Untersuchungen uber die optischen Eigenschaften von feinvertheilten Korpern," Annalen der Physik und Chemie, *23* No. 10 (1884), pp. 298–306.

2. McAlister, E. D., "The Christiansen Light Filter: Its Advantages and Limitations," Smithsonian Miscellaneous Collections, *93,* No. 7 (April 1935), pp. 1–11.
3. Raman, C. V., "The Theory of the Christiansen Experiment," Proc. India Acad. Sci., *A29,* (1943), pp. 381–390.
4. Clarke, R. H., "A Theory for the Christiansen Filter," Appl. Opt., *7,* No. 5 (May 1968), pp. 861–868.
5. Kell, R. D., "Color Television Camera," U. S. Patent No. 2,733,291, applied for July 29, 1952 issued January 31, 1956.
6. Takagi, T., and Nagahara, S., "Single-Pickup-Tube Color Camera System," Japan Elec. Eng., No. 11, (1967), pp. 41–44.
7. Ratcliffe, J. A., "Some Aspects of Diffraction Theory and Their Application to the Ionosphere," Rep. Progress Phys., *19,* (1956), pp. 188–267.
8. Brown, E. B., *Modern Optics,* New York: Reinhold Publishing, 1965, p. 484.
9. Born, M., and Wolf, E., *Principles of Optics,* New York: Pergamon Press, 1965, 3rd ed., p. 370.
10. Leith, E. N., and Palermo, C. J., "Introduction to Optical Data Processing," University of Michigan Eng. Summer Conf., May 24–June 4, 1965, pp. 2–12 to 2–28.
11. Van der Lugt, A., "Operational Notation for the Analysis and Synthesis of Optical Data-Processing Systems," Proc. IEEE, *54,* No. 8 (August 1966), pp. 1055–1063.
12. O'Neill, E. L., *Introduction to Statistical Optics,* Reading, Massachusetts: Addison-Wesley, 1963, p. 100.
13. Rowe, H. E., *Signals and Noise in Communication Systems,* New York: Van Nostrand, 1965, pp. 125–135.
14. Tatarski, V. I., *Wave Propagation in a Turbulent Medium,* New York: McGraw-Hill, 1961), p. 6.
15. Booker, H. G., and Gordon, W. E., "A Theory of Radio Scattering in the Troposphere," Proc. Inst. Radio Engineers, *38,* No. 4 (April 1950), p. 401.
16. Pekeris, C. L., "Note on the Scattering of Radiation in an Inhomogeneous Medium," Phys. Rev., *71,* No. 4 (February 1947), p. 268.
17. Herriott, D. R., "Recording Electronic Lens Bench," J. Opt. Soc. Amer., *48,* No. 12 (December 1958), pp. 968–971.
18. Baldwin, M. W., Jr., "Subjective Sharpness of Additive Color Picture," Proc. IRE, *39,* No. 10 (October 1951), pp. 1173–1176.

# Step Correction of Misaligned Beam Waveguides

By J. C. DALY

*This paper studies the steady-state performance of a system that stabilizes the beam position in optical waveguides. We have constrained the system to make only a single fixed amount of correction in the beam position at any given lens. We consider a symmetric corrector capable of correcting both positive and negative errors at any given lens and an unsymmetric corrector capable of correcting only positive errors at any given lens. We give the results from our studies of the performance of these systems when the lens misalignment forms a wave at the guide resonant spatial frequency, $\omega_0$ , and from our simulation of 5,000 confocal guides which were subjected to uncorrelated lens misalignment. We also derive an approximate statistical theory relating the root mean square beam displacement to the root mean square lens misalignment. We relate systems where correction occurs at every lens to systems where correction can occur only at every nth lens.*

## I. INTRODUCTION

Both theoretical and experimental studies of the guided transmission of optical beams indicate that some sort of active control of the beam position is required in order to keep the beam within the guide when it is transmitted over long distances.[1-7]

The optical waveguide considered here consists of a sequence of identical lenses of focal length $f$ and spacing $d$, as shown in Fig. 1. The system operates by sensing the position of the beam and making discrete adjustments in the transverse positions of the optical centers of the lenses. The system maintains the beam within a given region of the guide axis. We consider steady-state performance. Various schemes that eliminate the possibility of bothersome overshoot in the transient response have been proposed.[8-10]

We consider the problem in two dimensions. It has been shown that the three-dimensional problem can be split into two separate two-

Fig. 1 — Optical waveguide.

dimensional problems.[1] Although linear systems have been shown to be capable of providing highly accurate beam position control, the desired positional stability can be obtained using simpler more efficient nonlinear mechanisms.[11] Christian, Goubau, and Mink have demonstrated experimentally and with computer simulations that a marked improvement in optical transmission can be achieved through the use of nonlinear self-aligning beam waveguides.[8]

In an effort to achieve the simplest physical design, we have constrained the system to make only a fixed amount of correction in the beam position at any given lens. When the beam displacement is large the system does not attempt to totally correct at any lens but rather accumulates correction in the same manner that lens misalignment causes beam displacement to accumulate with distance of propagation. If the beam displacement is small the system does nothing until the beam has propagated through enough misaligned lenses so that its displacement exceeds a threshold; then the beam displacement is reduced by an amount equal to the threshold. The lenses are moved in a manner that suppresses the component of the lens displacement at the guide resonant spatial frequency, $\omega_o = (1/d) \cos^{-1} (1 - d/2f)$.

The correction of the beam position at the $n$th lens is introduced by changing the slope of the beam at the $(n - 1)$th lens. One way to accomplish this is to induce a corrective displacement of the $(n - 1)$th lens.*

II. SYSTEMS

2.1 *The Three-State Corrector*

The three-state corrector system is capable of making either positive or negative corrections in the beam position at any lens. The

---

* The slope of the beam can also be changed by inserting prisms into the beam, by changing the focal length of the lenses, or by changing the distance between lenses.

amount of the correction is equal to the threshold. The name "three state corrector" arises from the number of positions available to the lenses. For this system each lens occupies either its center position or one of two alternate positions. If at the $n$th lens the beam displacement exceeds a threshold $T$, the $(n-1)$th lens is displaced an amount $-Tf/d$ causing the beam displacement at the $n$th lens to be reduced by $T$. If the beam displacement at the $n$th lens passes the threshold in the negative direction the $(n-1)$th lens would be moved an amount $Tf/d$ to its other alternate position. If either of these corrections does not result in the beam displacement being less than the threshold no further correction can be made at this lens. However, additional correction is added at other lenses as the beam propagates.

## 2.2 *The Two-State Corrector*

The two-state corrector system differs from the three state system in that it is capable of making corrections in only one direction at any given lens. A lens is displaced $-Tf/d$ to its alternate position when the beam displacement at the next lens exceeds the positive threshold. Although the system is capable of correcting only positive errors at any given lens, negative errors do not become large because the beam oscillates about the guide axis as it propagates. Negative errors become positive errors after the beam has propagated a distance $d = \pi/\omega_o$ .

The two-state system is capable of reducing any errors in the beam position provided the beam remains within the aperture of the guide. The correction is distributed over more lenses than with the three-state corrector.

## III. WORST CASE PERFORMANCE IN A CONFOCAL GUIDE

Let us consider the worst case situation to be a sequence of equal lens displacements, $D$, forming a square wave at the spatial resonant frequency $\omega_o$ . The effect of each displacement adds directly to the effect of previous lens displacements. The beam displacement is proportional to the number of lenses through which it has passed. The beam oscillates about the guide axis as it propagates.

The response of the control mechanism is to generate a correction in the beam displacement which subtracts from the beam displacement caused by the lens misalignment. The increase in the beam displacement at any lens is $Dd/f$ where $D$ is the amount that the lens is misaligned and $d/f$ equals two for a confocal guide.

For the three-state system the amount of correction in the beam displacement that can be introduced at any lens is equal to the

threshold level, $T$. In order for the beam to remain within a finite region of the guide axis the increase in the beam displacement at each lens must be less than the amount of correction that is possible, that is $D < T/2$. As long as the lens misalignment is less than $T/2$ the beam displacement can not exceed the threshold.

For the two-state system the amount of correction that can be introduced at any given lens is either $T$ or zero, depending on whether the beam displacement is positive or negative. If, on the average, the amount of the increase in the beam displacement is to be less than the amount of the correction, $D$ must be less than $T/4$. As long as the lens misalignment is less than $T/4$ the beam displacement can not exceed $1.5T$.

IV. RESULTS OF SIMULATIONS

Computer experiments were performed in order to evaluate the performance of the systems. The two-state and three-state systems were each simulated in 5,000 confocal guides. Each guide was subjected to a different set of uncorrelated gaussian amplitude distributed, transverse lens displacements. Quantities that were observed were: $\sigma_r$, the rms beam displacement as a function of $\sigma_L$ (the rms lens misalignment), and the distribution of beam displacements at the 25th lens for various values of rms lens misalignment. All quantities are measured in units of $T$, the threshold displacement.

From the rms beam displacement the mean square beam displacement $\langle r^2 \rangle$ was determined and plotted in Figs. 2 and 3 as a function of the distance of propagation through the guides for the three-state and



Fig. 2 — Mean square beam deviation averaged over 5,000 confocal guides with the three-state correctors at each lens.

Fig. 3 — Mean square beam deviation averaged over 5,000 confocal guides with the two-state correctors.

two-state systems, respectively. The staircase appearance of the plotted points is because of the effects of the displacements of alternate lenses are independent in a confocal guide, for example, displacing even numbered lenses causes the beam displacement only at odd numbered lenses. The plots of $\langle r^2 \rangle$ in Figs. 2 and 3 approach straight line asymptotes. The increase in $\langle r^2 \rangle$ is linear until it has increased to the point where the control mechanism begins to act to maintain $\langle r^2 \rangle$ at a constant equilibrium value. In the appendix, an approximate expression relating the equilibrium value of $\sigma_r$ to $\sigma_L$ is derived for both the two-state and three-state position control systems. Figure 4 con-



Fig. 4 — The rms beam displacement versus the rms lens misalignment.

tains plots of this approximation. The experimental points marked on the plot were obtained from computer simulations of the two-state mechanism.

Figures 5 and 6 are the distribution of beam displacements from the confocal guide axis at the 25th lens, when the guide is controlled with a three-state system. The lens displacements are uncorrelated, gaussian, and have a standard deviation, $\sigma_L$, of 0.1 in Fig. 5 and 0.5 in Fig. 6. In Figs. 5 through 9, $(\Delta n/N)/\Delta r$ is the fraction of beam displacements per unit displacement. When $\sigma_L$ equaled 0.1 (Fig. 5) in the three-state system, the mechanism kept the beam displacement below the threshold in all of the 5,000 guides simulated. Figures 7, 8, and 9 are the distributions of beam displacements from the confocal guide axis at the 25th lens when the guide is stabilized using a two-state position corrector. The lack of symmetry of the corrector is evident in the distributions. A comparison of Figs. 5 and 7 shows that when the rms lens deviation is small (0.1 of the threshold) the two-state system performs nearly as well as the three-state system.

The system can be further simplified by allowing only every $n$th lens to be movable. Then for both the two- and three-state systems, it follows from equation (1) that the guide misalignment produces an increase in the mean square beam displacement of $2n\sigma_L^2$ between each corrector. Therefore the performance of a system where only every $n$th lens is adjustable is equivalent to the performance of a system where



Fig. 5 — Distribution of beam displacements for a three-state corrector with the rms lens misalignment $\sigma_L$ equal to 0.1.

Fig. 6 — Distribution of beam displacements for a three-state corrector with the rms lens misalignment $\sigma_L$ equal to 0.5.

each lens is adjustable, provided the lens displacements are uncorrelated and the mean square lens displacement is less by a factor of $1/n$.

Since, in a confocal guide, corrections in the beam position are introduced at the even lenses by displacing the odd lenses and at the odd lenses by displacing the even lenses, the distribution of correctors must be evenly spread between the even and the odd numbered lenses.

V. CONCLUSION

The use of two- and three-state beam position controllers in optical waveguides stabilizes the beam position in the guide.

When the misalignment of the lenses in a confocal guide forms waves at $\omega_o$, the resonant spatial frequency of the guide, the rms displacement of the beam reaches an equilibrium value beyond which it does not



Fig. 7 — Distribution of beam displacements for a two-state corrector with the rms lens misalignment $\sigma_L$ equal to 0.1.

Fig. 8 — Distribution of beam displacements for a two-state corrector with the rms lens misalignment $\sigma_L$ equal to 0.3.

grow, as long as the amplitude of the wave of lens displacements is less than $T/2$ for the three-state system or $T/4$ for the two-state corrector.

When the lens displacements are gaussian, uncorrelated, and have an rms value, $\sigma_L$, that is less than 0.1 of the threshold the two-state system works nearly as well as the three-state system. If $\sigma_L < 0.1$, then $\sigma_r < 0.5$ and the distribution of the beam displacements is approximately gaussian. The probability of the beam exceeding the threshold at any given lens is less than $10^{-4}$ if $\sigma_L < 0.1$.

APPENDIX

*Approximate Statistical Theory*

It has been shown that in a beam waveguide with misaligned lenses the probability distribution of the beam displacement from the guide axis is gaussian with a standard deviation that increases with increasing distance of propagation through the guide.[12] In a confocal guide the average value of this increase is given by

$$\triangle_g \sigma_r^2 = 2\sigma_L^2 \tag{1}$$

where $\triangle_g \sigma_r^2$ is the change in the mean square beam deviation resulting from guide misalignment and $\sigma_L$ is the rms lens displacement.

When a discrete-state beam position control system is used, the probability distribution is stationary and the mean square beam deviation, $\sigma_r^2$, does not increase with increasing distance of propagation. The increase in $\sigma_r^2$ resulting from lens misalignment is counteracted by a decrease in $\sigma_r^2$ resulting from the action of the controller.

Assume that $P(r)$ is the probability that beam displacement equals $r$ before the controller acts. When the controller is a three-state cor-

rector, the mean square beam deviation after the controller acts is

$$\langle r^2 \rangle_a = 2 \int_0^1 r^2 P(r) \, dr + 2 \int_1^\infty (r-1)^2 P(r) \, dr; \tag{2}$$

expanding the term $(r-1)^2$ in the second integral yields the equation

$$\langle r^2 \rangle_a = \langle r^2 \rangle_b - 4 \int_1^\infty rP(r) \, dr + P(\mid r \mid > 1) \tag{3}$$

where $\langle r^2 \rangle_b$ and $P(\mid r \mid > 1)$ are the mean square beam deviation before the controller acts and the probability that the magnitude of the deviation before the controller acts is greater than one, respectively.

When $P(r)$ drops off rapidly for $r > 1$ an approximation to the integral in equation (3) can be obtained by setting $r$ equal to one in the integral. This results in the following expression for the change in the mean square beam deviation resulting from the controller

$$\triangle_c \sigma_r^2 = \langle r^2 \rangle_a - \langle r^2 \rangle_b \approx -P(\mid r_b \mid > 1) \tag{4}$$

where $P(\mid r_b \mid > 1)$ is the probability that the magnitude of the beam deviation before the controller acts is greater than one. It is also the probability that the controller acts. Since $\triangle_g \sigma_r^2 = -\triangle_c \sigma_r^2$ it follows from equations (1) and (4) that

$$P(\mid r_b \mid > 1) \approx 2\sigma_L^2 . \tag{5}$$

Equation (5) is the probability that the three-state system has responded to a threshold crossing. Assuming that $P(r)$ is gaussian it follows from equation (5) that

$$\sigma_L^2 = \text{erf} \, (-1/\sigma_r) \tag{6}$$

where $\sigma_r$ is the rms beam displacement and $\text{erf}(-1/\sigma_r)$ is the error function of $-1/\sigma_r$.



Fig. 9 — Distribution of beam displacements for a two-state corrector with the rms lens misalignment $\sigma_L$ equal to 0.4.

Equation (6) is approximate. In order to apply it to the guide one must assume (i) that the probability distribution of beam displacements at any given lens is gaussian, and (ii) that the change in the distribution at any lens resulting from the guide misalignment is small. These assumptions are more accurately satisfied as the rms misalignment decreases. A comparison of the distribution obtained from the simulations when $\sigma_L = 0.1$ and the gaussian used to approximate it is shown in Fig. 5.

An analysis of the two-state system that proceeds in the same manner as above yields the following relationship between the rms beam and rms lens displacements

$$\sigma_L^2 = \tfrac{1}{2} \operatorname{erf} (-1/\sigma_r); \tag{7}$$

the probability that at any given lens the control system will respond to a threshold crossing is

$$P(r > 1) = 2\sigma_L^2. \tag{8}$$

The relationships between $\sigma_r$ and $\sigma_L$ given by equations (6) and (7) are plotted in Fig. 4 along with experimental points obtained from simulations of the two-state system.

REFERENCES

1. Hirano, J., and Fukatsu, Y., "Stability of a Light Beam in a Beam Waveguide," Proc. IEEE, *52*, No. 11 (November 1964), pp. 1284–1292.
2. Steier, W. H., "The Statistical Effects of Random Variations in the Components of a Beam Waveguide," B.S.T.J., *45*, No. 3 (March 1966), pp. 451–471.
3. Marcuse, D., "Physical Limitations on Ray Oscillation Suppressors," B.S.T.J., *45*, No. 5 (May–June 1966), pp. 743–751.
4. DeLange, O. E., "Losses Suffered by Coherent Light Redirected and Refocused Many Times in an Enclosed Medium," B.S.T.J., *44*, No. 2 (February 1965), pp. 283–302.
5. Gloge, D., "Experiments With an Underground Lens Waveguide," B.S.T.J., *46*, No. 4 (April 1967), pp. 721–735.
6. Christian, J. R., Goubau, G., and Mink, J. W., "Further Investigations With an Optical Beam Waveguide for Long Distance Transmission," IEEE Trans. on Microwave Theory and Techniques, *MTT-15*, No. 4 (April 1967), pp. 216–219.
7. Beck, A. C., "An Experimental Gas Lens Optical Transmission Line," IEEE Trans. on Microwave Theory and Techniques, *MTT-15*, No. 7 (July 1967), pp. 433–434.
8. Christian, J. R., Goubau, G., and Mink, J. W., "Self-Aligning Optical Beam Waveguides," paper 5–6, 1967 IEEE Conf. on Laser Eng. and Applications, Washington, D. C., June 6–8, 1967.
9. Ring, D. H., unpublished work.
10. Richter, P. S., unpublished work.
11. Daly, J. C., "Linear Beam Position Control in Optical Waveguides," B.S.T.J., *47*, No. 5 (May–June 1968), pp. 783–799.
12. Marcuse, D., "Probability of Ray Position in Beam Waveguides," IEEE Trans. on Microwave Theory and Techniques, *MTT–15*, No. 3 (March 1967), pp. 167–171.

# Eigenmodes of an Asymmetric Cylindrical Confocal Laser Resonator with a Single Output-Coupling Aperture

By D. E. McCUMBER

*Previous calculations of the low-loss modes of a symmetric cylindrical confocal laser resonator have been extended to the asymmetric case. Diffraction losses are governed by the geometric-mean Fresnel number $N_m$ of the two end mirrors and, in the system we consider, by the Fresnel number $N_o$ of an output coupling aperture in one of the mirrors. Loss factors and mirror field distributions have been calculated numerically for different $N_o$ for $N_m$ in the range $0.6 \leq N_m \leq 2$.*

## I. INTRODUCTION

In a previous paper[1] we described the diffraction losses and the field distributions at the reflectors of the low-loss modes of a symmetric cylindrical confocal resonator for Fresnel numbers $0.6 \leq N_m \leq 2$. We considered the effect of output-coupling apertures in the reflectors but we assumed that both reflectors were identical, each with the same output aperture and the same maximum radius. In this paper we again consider the cylindrical confocal geometry, but we do not require identical reflectors and, in particular, we assume that only one reflector is pierced by an output-coupling aperture, as in the coupling scheme proposed by Patel and others.[2]

Figure 1 shows an axial section of the confocal resonator in question. The cavity is bounded at its two ends by confocal spherical mirrors (more exactly, confocal paraboloids[3]). The first mirror is perfectly reflecting over the annular region $0 \leq a_{1o} \leq \rho \leq a_{1m}$, the second over the circular section $0 \leq \rho \leq a_{2m}$. The maximum radii ($a_{1m}$, $a_{2m}$) are both much less than the mirror separation $b$.

Expressions for the eigenvalues and eigenfunctions of asymmetric rectangular confocal resonators with output coupling slits have been

Fig. 1 — Axial section of cylindrical confocal laser cavity. The cavity is bounded at its two ends by confocal spherical mirrors with radius of curvature $b$. One mirror is perfectly reflecting over the annular region $0 \leqq a_{1o} \leqq \rho \leqq a_{1m}$, the other over the circular section $0 \leqq \rho \leqq a_{2m}$. Both $a_{1m}$ and $a_{2m}$ are much less than $b$.

derived by Boyd and Kogelnik.[4] Properties of symmetric resonators without coupling apertures are summarized with an extensive list of references in the review article by Kogelnik.[5] Equivalence relations relating asymmetric and symmetric resonators with circular mirrors have been derived by Gordon and Kogelnik.[6]

Our analysis of the resonator of Fig. 1 closely parallels that of the symmetric resonator.[1] Assuming that all dimensions are large compared with the optical wavelength $\lambda$, we again use a scalar formulation of Huygens' principle.[3-5] For the cylindrical confocal geometry the field amplitude at reflector $j$, $j = 1$ or 2, for a typical mode can be written in the form

$$F_{lp}^{(j)}(\rho, \varphi) = f_{lp}^{(j)}(\rho) \exp(-il\varphi), \qquad (1)$$

where $(\rho, \varphi)$ are radial and angular coordinates in a plane perpendicular to the resonator axis and where $(l, p)$ are angular and radial quantum integers (transverse quantum numbers). For this asymmetric system with nonidentical mirrors, we cannot require for an eigenmode that the field amplitude distribution $F_{lp}^{(j)}(\rho, \varphi)$ at one mirror be a constant multiple of that at the other. Rather, we must require for eigenmodes that after a round-trip transit of the resonator the field amplitude at one mirror be a constant multiple of the initial field amplitude at the same mirror. This more elaborate self-reproducing requirement together with Huygens' principle gives the following pair of simultaneous integral equations which must be satisfied by the radial eigenfunctions $f_{lp}^{(j)}(\rho)$ and eigenvalues $\kappa_{lp}^{(i)}$ [compare equation (2) of Ref. 1].

$$\kappa_{lp}^{(2)} f_{lp}^{(2)}(\rho_2) = \frac{2\pi}{b\lambda} \int_{a_{1o}}^{a_{1m}} d\rho_1 \, \rho_1 J_l(2\pi\rho_2\rho_1/b\lambda) f_{lp}^{(1)}(\rho_1), \qquad (2a)$$

$$\kappa_{lp}^{(1)} f_{lp}^{(1)}(\rho_1) = \frac{2\pi}{b\lambda} \int_{o}^{a_{2m}} d\rho_2 \; \rho_2 J_l(2\pi \rho_1 \rho_2 / b\lambda) f_{lp}^{(2)}(\rho_2). \qquad (2b)$$

$J_l(z)$ is the Bessel function of order $|l|$. The loss factor is

$$\alpha_{lp} = 1 - |\kappa_{lp}^{(1)} \kappa_{lp}^{(2)}|, \qquad (3)$$

which is the fractional energy of a mode lost per reflection (or during the one-way transit time $b/c$, where $c$ is the velocity of light in the resonator).[5] The phase of the eigenvalue product $\kappa_{lp}^{(1)} \kappa_{lp}^{(2)}$ determines the resonant wavelength:

$$\text{resonant } \lambda = 4\pi b/[(l+1)\pi - \text{Arg } \kappa_{lp}^{(1)} \kappa_{lp}^{(2)} - 2\pi n], \qquad (4)$$

where $n$ is an arbitrary integer (longitudinal quantum number).

It is useful to introduce the mirror Fresnel numbers

$$N_m^{(1)} = a_{1m}^2 / \lambda b, \qquad N_m^{(2)} = a_{2m}^2 / \lambda b, \qquad (5a)$$

and their geometric mean

$$N_m \equiv [N_m^{(1)} N_m^{(2)}]^{\frac{1}{2}} = a_{1m} a_{2m} / \lambda b. \qquad (5b)$$

In place of the variables $\rho_i$ and the functions $f_{lp}^{(i)}(\rho_i)$ in equations (2), we introduce new variables

$$r_i = r_m \rho_i / a_{im} \qquad (6a)$$

and functions

$$g_{lp}^{(i)}(r_i) = f_{lp}^{(i)}(r_i a_{im} / r_m), \qquad (6b)$$

where $r_m^2 = N_m$. We characterize the size of the radius-$a_{1o}$ hole in the first mirror by the Fresnel number

$$N_o \equiv r_o^2 = (a_{1o}/a_{1m})^2 N_m = a_{1o}^2 a_{2m} / \lambda b a_{1m}. \qquad (7)$$

With no significant loss of generality we can define the functions $f_{lp}^{(i)}(\rho_i)$ such that

$$a_{2m} \kappa_{lp}^{(2)} / a_{1m} = a_{1m} \kappa_{lp}^{(1)} / a_{2m} = \kappa_{lp} \qquad (8)$$

and

$$\delta_{pq} = 2\pi \int_{r_o}^{r_m} dr_1 \; r_1 g_{lp}^{(1)}(r_1) g_{lq}^{(1)}(r_1), \qquad (9a)$$

$$\delta_{pq} = 2\pi \int_{o}^{r_m} dr_2 \; r_2 g_{lp}^{(2)}(r_2) g_{lq}^{(2)}(r_2). \qquad (9b)$$

With this notation the eigenvalue equations (2) become

$$\kappa_{lp} g_{lp}^{(2)}(r_2) = 2\pi \int_{r_o}^{r_m} dr_1\, r_1 J_l(2\pi r_2 r_1) g_{lp}^{(1)}(r_1), \tag{10a}$$

$$\kappa_{lp} g_{lp}^{(1)}(r_1) = 2\pi \int_{o}^{r_m} dr_2\, r_2 J_l(2\pi r_1 r_2) g_{lp}^{(2)}(r_2). \tag{10b}$$

Equation (3) simplifies to

$$\alpha_{lp} = 1 - \mid \kappa_{lp} \mid^2, \tag{11}$$

which depends upon the parameters $(a_{1o}, a_{1m}, a_{2m}, b, \lambda)$ only through the Fresnel numbers $N_m = r_m^2$ and $N_o = r_o^2$. In what follows we describe how the loss factor $\alpha_{lp}$ and the amplitudes $g_{lp}^{(i)}(r)$ change with $N_o$ for $N_m$ in the interval $0.6 \leq N_m \leq 2$. Solutions for $N_o = 0$ are described elsewhere.[1,3,7-9]

Our numerical method is similar to that previously described for the symmetric resonator.[1] We expand the Bessel-function kernels in equation (10) as power series, truncate the series after a finite number $M = \max(10\, N_m + 1, 10)$ of terms, and reduce the integral equations (10) to $M$-dimensional matrix equations which are solved numerically with standard matrix routines. [The reduction of equation (10) to matrix form is described in the Appendix.] The merits of this technique remain as described before.[1]

## II. ANALYTIC METHODS FOR SMALL $N_o$

The eigenvalues $\kappa_{lp}$ and field amplitudes $g_{lp}(r)$ for $N_o = 0$ are described elsewhere.[1,3,7-9] The loss factor $\alpha_{lp} = 1 - \mid \kappa_{lp} \mid^2$ for the four lowest-loss modes (compare with Fig. 2 of Ref. 1) are tabulated for $0.6 \leq N_m \leq 2$ in Table I. For $l = 0$, the field amplitude $g_{lp}(r)$ at $r = 0$ is finite; for $l \neq 0$, it vanishes as $r^{|l|}$. As before, we expect the modes with angular quantum number $l = 0$ to be more sensitive to the coupling aperture $(N_o > 0)$ than the $l \neq 0$ modes.[1]

If we use a superscript "0" to lable the eigenvalues and field amplitudes for $N_o = 0$, then for small $r$, to within corrections of relative order $r^2$,

$$g_{lp}^0(r) = g_{0p}^0(0) \quad \text{for} \quad l = 0 \tag{12a}$$

$$= c_{lp}^0 r^{|l|} \quad \text{for} \quad l \neq 0, \tag{12b}$$

where coefficients $g_{0p}^0(0)$ are listed in Table II (an expanded version of Table III, Ref. 1) and coefficients $c_{lp}^0$ in Table III. The forms (12) are

TABLE I—LOSS FACTOR $\alpha_{lp}$ FOR $N_o = 0$

| $N_m$ | $\alpha_{00}{}^0$ | $\alpha_{01}{}^0$ | $\alpha_{10}{}^0$ | $\alpha_{20}{}^0$ |
|---|---|---|---|---|
| 0.6 | $3.614 \times 10^{-2}$ | 0.7931 | 0.2724 | 0.6679 |
| 0.7 | $1.301 \times 10^{-2}$ | 0.6131 | 0.1371 | 0.4712 |
| 0.8 | $4.448 \times 10^{-3}$ | 0.4131 | $6.103 \times 10^{-2}$ | 0.2900 |
| 0.9 | $1.471 \times 10^{-3}$ | 0.2411 | $2.477 \times 10^{-2}$ | 0.1563 |
| 1.0 | $4.759 \times 10^{-4}$ | 0.1233 | $9.417 \times 10^{-3}$ | $7.505 \times 10^{-2}$ |
| 1.1 | $1.515 \times 10^{-4}$ | $5.651 \times 10^{-2}$ | $3.424 \times 10^{-3}$ | $3.285 \times 10^{-2}$ |
| 1.2 | $4.767 \times 10^{-5}$ | $2.382 \times 10^{-2}$ | $1.206 \times 10^{-3}$ | $1.343 \times 10^{-2}$ |
| 1.3 | $1.485 \times 10^{-5}$ | $9.462 \times 10^{-3}$ | $4.151 \times 10^{-4}$ | $5.225 \times 10^{-3}$ |
| 1.4 | $4.59 \times 10^{-6}$ | $3.601 \times 10^{-3}$ | $1.403 \times 10^{-4}$ | $1.961 \times 10^{-3}$ |
| 1.5 | $1.41 \times 10^{-6}$ | $1.328 \times 10^{-3}$ | $4.674 \times 10^{-5}$ | $7.164 \times 10^{-4}$ |
| 1.6 | $4.3 \times 10^{-7}$ | $4.78 \times 10^{-4}$ | $1.539 \times 10^{-5}$ | $2.561 \times 10^{-4}$ |
| 1.7 | $1.3 \times 10^{-7}$ | $1.69 \times 10^{-4}$ | $5.01 \times 10^{-6}$ | $9.000 \times 10^{-5}$ |
| 1.8 | $4 \times 10^{-8}$ | $5.89 \times 10^{-5}$ | $1.62 \times 10^{-6}$ | $3.116 \times 10^{-5}$ |
| 1.9 | $1 \times 10^{-8}$ | $2.02 \times 10^{-5}$ | $5.2 \times 10^{-7}$ | $1.066 \times 10^{-5}$ |
| 2.0 | $4 \times 10^{-9}$ | $6.86 \times 10^{-6}$ | $1.7 \times 10^{-7}$ | $3.60 \times 10^{-6}$ |

useful for estimating the perturbations induced by a small finite $N_o$. To first order, the perturbed field amplitudes are

$$g_{lp}^{(1)}(r) = g_{lp}^0(r)\left\{1 + \pi \int_o^{r_o} dr_1\, r_1 [g_{lp}^0(r_1)]^2\right\}$$

$$- \sum_{q \neq p}{}' g_{lq}^0(r) \frac{(\kappa_{lq}^0)^2}{(\kappa_{lp}^0)^2 - (\kappa_{lq}^0)^2} 2\pi \int_o^{r_o} dr_1\, r_1 g_{lq}^0(r_1) g_{lp}^0(r_1), \qquad (13a)$$

TABLE II—FIELD AMPLITUDE AT $r = 0$ FOR $l = 0$
MODES WITH $N_o = 0$

| $N_m$ | $g_{00}{}^0(0)$ | $g_{01}{}^0(0)$ | $g_{02}{}^0(0)$ |
|---|---|---|---|
| 0.6 | 1.2770 | 1.2251 | 1.6159 |
| 0.7 | 1.3021 | 1.1541 | 1.4851 |
| 0.8 | 1.3213 | 1.1254 | 1.3735 |
| 0.9 | 1.3354 | 1.1286 | 1.2760 |
| 1.0 | 1.3457 | 1.1511 | 1.1930 |
| 1.1 | 1.3536 | 1.1807 | 1.1291 |
| 1.2 | 1.3597 | 1.2093 | 1.0890 |
| 1.3 | 1.3647 | 1.2336 | 1.0742 |
| 1.4 | 1.3688 | 1.2532 | 1.0812 |
| 1.5 | 1.3723 | 1.2688 | 1.1025 |
| 1.6 | 1.3752 | 1.2814 | 1.1302 |
| 1.7 | 1.3778 | 1.2918 | 1.1580 |
| 1.8 | 1.3800 | 1.3006 | 1.1826 |
| 1.9 | 1.3820 | 1.3081 | 1.2033 |
| 2.0 | 1.3838 | 1.3146 | 1.2203 |
| $\infty$ | $\sqrt{2} = 1.4142$ | 1.4142 | 1.4142 |

TABLE III—FIELD-AMPLITUDE COEFFICIENT FOR
$l \neq 0$ MODES AT $r = 0$ FOR $N_o = 0$

| $N_m$ | $c_{10}{}^0$ | $c_{11}{}^0$ | $c_{20}{}^0$ |
|---|---|---|---|
| 0.6 | 2.6428 | 4.4649 | 3.9948 |
| 0.7 | 2.7222 | 4.0086 | 3.9131 |
| 0.8 | 2.8269 | 3.7024 | 3.9979 |
| 0.9 | 2.9266 | 3.5247 | 4.1770 |
| 1.0 | 3.0094 | 3.4618 | 4.3917 |
| 1.1 | 3.0746 | 3.4932 | 4.5999 |
| 1.2 | 3.1255 | 3.5889 | 4.7809 |
| 1.3 | 3.1659 | 3.7148 | 4.9305 |
| 1.4 | 3.1988 | 3.8431 | 5.0526 |
| 1.5 | 3.2260 | 3.9580 | 5.1531 |
| 1.6 | 3.2491 | 4.0549 | 5.2371 |
| 1.7 | 3.2690 | 4.1351 | 5.3087 |
| 1.8 | 3.2862 | 4.2017 | 5.3704 |
| 1.9 | 3.3014 | 4.2578 | 5.4244 |
| 2.0 | 3.3148 | 4.3059 | 5.4721 |
| $\infty$ | $2\pi^{1/2} = 3.5449$ | $2^{3/2}\pi^{1/2} = 5.0133$ | $2\pi = 6.2832$ |

$$g_{lp}^{(2)}(r) = g_{lp}^0(r) - \sum_{q \neq p}{}' \, g_{lq}^0(r) \frac{\kappa_{lp}^0 \kappa_{lq}^0}{(\kappa_{lp}^0)^2 - (\kappa_{lq}^0)^2}$$

$$\cdot 2\pi \int_o^{r_o} dr_1 \, r_1 g_{lq}^0(r_1) g_{lp}^0(r_1). \qquad (13b)$$

To second order, the loss factor is

$$\alpha_{lp} = \alpha_{lp}^0 + (1 - \alpha_{lp}^0) 2\pi \int_o^{r_o} dr_1 \, r_1 [g_{lp}^0(r_1)]^2$$

$$- \sum_{q \neq p}{}' \frac{(1 - \alpha_{lq}^0)(1 - \alpha_{lp}^0)}{\alpha_{lq}^0 - \alpha_{lp}^0} \left[ 2\pi \int_o^{r_o} dr_1 \, r_1 g_{lq}^0(r_1) g_{lp}^0(r_1) \right]^2. \qquad (14)$$

In deriving (14), we used the fact that for the cylindrical confocal geometry the eigenvalues $\kappa_{lp}$ are real and $\kappa_{lp}^2 = |\kappa_{lp}|^2$.

The last terms in (13) and (14) describe mode mixing by the aperture. The amount of mixing depends upon the separation of the eigenvalues as well as upon the strength of the perturbation. Degenerate or nearly degenerate modes are much more sensitively coupled than are modes with greatly different losses. In the symmetric resonator the even-$p$ and odd-$p$ modes of a particular angular quantum number $l$ do not mix; such modes do mix in the asymmetric geometry.[1]

For $N_o$ sufficiently small, we can neglect the second-order or mode-

mixing term in (14). If we use the small-$r$ approximations (12) in the remaining integral, we find for $l = 0$ that

$$\alpha_{0p} = \alpha_{0p}^0 + (1 - \alpha_{0p}^0)\pi N_0[g_{0p}^0(0)]^2 \tag{15a}$$

and for $l \neq 0$ that

$$\alpha_{lp} = \alpha_{lp}^0 + (1 - \alpha_{lp}^0)\pi (c_{lp}^0)^2 N_o^{|l|+1}/(|l| + 1). \tag{15b}$$

These expressions confirm our previous conjecture that modes with $l = 0$ are more sensitive to aperture loss than are modes with $l \neq 0$.

A quantity of interest in the design of lasers with aperture output coupling is that value $N_{oc}$ of $N_o$ for which the losses of the longitudinal (00) mode equal the losses of the (least lossy) transverse (10) mode.[2] All other things being equal, the laser will operate in the (00) mode for $N_o < N_{oc}$, in the (10) mode for $N_o > N_{oc}$, and in still another mode for larger values of $N_o$. Using Eqs. (15), we estimate that

$$N_{oc} = (\alpha_{10}^0 - \alpha_{00}^0)/\pi(1 - \alpha_{00}^0)[g_{00}^0(0)]^2. \tag{16}$$

III. NUMERICAL RESULTS FOR $N_o$ FINITE

Let us compare estimates based upon the approximate expressions (15) and (16) with accurate numerical results. Using the numerical technique outlined at the end of the introduction and in greater detail in the Appendix, we have computed the loss factor $\alpha_{lp}$ for $N_o$ finite and $N_m$ in the range $0.6 \leq N_m \leq 2$. Results for $N_m = 0.8$ and $N_o$ variable are shown in Fig. 2 and similar results for $N_m = 1.6$ in Fig. 3. Results for $N_o = 0.001$ and $N_m$ variable are shown in Fig. 4. These examples were chosen to facilitate comparison with the two-aperture symmetric geometry of Ref. 1. The results are qualitatively similar to those obtained before, except that here odd-$p$ and even-$p$ modes interact whereas before they did not.

Predictions based upon the first-order expressions (15) and Tables I to III are shown as dashed lines in Figs. 2 and 3. The fit to the exact results is good for sufficiently small $N_o$, but deviations are large for those $N_o$'s for which the interaction between the $(l, p)$ and $(l, p + 1)$ modes is evident as a repulsion in the calculated loss curves.

The critical single-aperture Fresnel number $N_{oc}$ for which the loss factor $\alpha_{lp}$ of the longitudinal (00) mode equals that of the lowest transverse (10) mode is shown as a function of mirror Fresnel number $N_m$ in Fig. 5. This is an important parameter in the design of aperture-

Fig. 2 — Loss factor $\alpha_{lp}$ versus single-aperture Fresnel number $N_o$ for low-loss modes when $N_m = 0.8$. (Compare with Fig. 9 of Ref. 1 for the symmetric two-aperture geometry.) The dashed lines are estimates based on equation (15) and the data from Tables I through III.

out-put-coupled cavities having good mode selection.[2] Also shown in Fig. 5 is the estimate of $N_{oc}$ derived from Eq. (16) and the data from Tables I through III. The agreement is reasonably good, and consistent with what one would expect from the accuracy of the estimates derived from Eqs. (15) in Figs. 2 and 3.

The effect of mode coupling is apparent in the amplitude $g_{lp}^{(i)}(r)$ of the field at the two mirrors. Consider the $l = 0$ modes, which are those most sensitive to finite $N_o$. The field amplitudes and intensities at the mirrors of the two lowest-loss modes are shown for $N_m = 0.8$ in Figs. 6 and 7 and for $N_m = 1.6$ in Figs. 8 through 11. Figures 6 and 8 show the distributions for $N_o = 0$; they are the same on both mirrors. [Distributions for other $(lp)$ modes are shown for $N_o = 0$ in Ref. 1.] The other figures show how the distributions change for $N_o > 0$. In each figure the dashed lines indicate the field distributions on the mirror pierced by the aperture, the solid lines those on the intact mirror. The radius of the mirrors and the radius of the aperature are indicated on the plots.

One should distinguish two effects apparent in the field plots as $N_o$ increases. First, there is a change in the magnitude of the field amplitude $g_{lp}^{(1)}(r)$ on the pierced mirror. This is a simple renormalization correction implicit in the requirement (9) that the fields be normalized

over the reflecting areas of the mirrors. Second, there are changes in the shape of the field distributions on both mirrors. This is a consequence of mode mixing, which becomes appreciable for those $N_o$'s for which significant mode repulsion is apparent in the curves of Figs. 2 and 3. In each case the effect is to reduce the intensity of the less-lossy mode in regions where the mirrors are not reflecting, at the expense of the more lossy of the two interacting modes. This is apparent, for example, in Fig. 9 with $N_m = 1.6$ and $N_o = 0.0001$, for which there is appreciable interaction between the (00) and (01) modes (compare with Fig. 3). The amplitude of the (00) mode at the aperture is decreased below that in Fig. 9; that of the (01) mode is increased. In Fig. 11 with $N_m = 1.6$



Fig. 3 — Loss factor $\alpha_{lp}$ versus single-aperture Fresnel number $N_o$ for low-loss modes when $N_m = 1.6$. (Compare with Fig. 11 of Ref. 1 for the symmetric two-aperture geometry.) The dashed lines are estimates based on equation (15) and the data from Tables I through III.

Fig. 4 — Loss factor $\alpha_{lp}$ versus Fresnel number $N_m$ for low-loss modes with $N_o$ fixed at 0.001. (Compare with Fig. 16 of Ref. 1 for the symmetric two-aperture geometry.)

and $N_o = 0.01$, the amplitude of the (01) mode at the aperture is reduced as a consequence of interaction with the (02) mode (Fig. 3).

IV. DISCUSSION

Our previous treatment[1] of a symmetric cylindrical confocal laser cavity is extended here to the asymmetric cylindrical confocal geometry of Fig. 1 and specific numerical values for the loss factor and for the field distributions at the mirrors have been calculated for Fresnel numbers $0.6 \leqq N_m \leqq 2$. The transformations outlined in Section I show that the loss factor of a cavity with different-sized mirrors $(N_m^{(1)} \neq N_m^{(2)})$ equals that of a cavity having two mirrors with the same outer dimensions $[N_m = (N_m^{(1)} N_m^{(2)})^{\frac{1}{2}}]$. The field distributions

scale accordingly [see equations (6)]. Similar results also obtain in the rectangular confocal geometry.[4]

For sufficiently small aperture Fresnel numbers $N_o$ the cavity losses associated with a single output coupling aperture can be estimated as in equations (15) from first-order perturbation theory. Just as in the symmetric two-aperture case considered previously, the value of $N_o$ for which such first-order calculations fail decreases rapidly as the Fresnel number $N_m$ increases, because the field distributions distort through mode mixing to minimize the aperture losses in the lowest-loss modes.[1] As before, this distortion occurs at approximately those values of $N_o$ and $N_m$ for which an observer at one reflector, using light of the relevant wavelength and optics limited by the radius $r_m = N_m^{\frac{1}{2}}$, can resolve the aperture of radius $r_o = N_o^{\frac{1}{2}}$ from the other reflector.[1,10]

APPENDIX

*Reduction of Integral Equations to Matrix Equations*

We express the Bessel-function kernels in power-series form. ($l = |l|$ throughout this appendix.)



Fig. 5 — Critical single-aperture Fresnel number $N_{oc}$ for which diffraction losses of longitudinal (00) mode equal those of the lowest transverse (10) mode versus Fresnel number $N_m$. (Compare with Fig. 18 of Ref. 1 for the symmetric two-aperture geometry.) The dashed line is an estimate based on equation (16) and the data from Tables I through III.

Fig. 6 — (a) Field amplitude $g_{lp}(r)$ and (b) field intensity $|g_{lp}(r)|^2$ for modes $(lp) = (00)$ and $(01)$ with $N_m = 0.8$ and $N_o = 0$. The field distributions are identical on both mirrors.

$$J_l(2\pi r_1 r_2) = (\pi r_1 r_2)^l \sum_{m=1}^{\infty} \frac{(-1)^{m-1}(\pi r_1 r_2)^{2(m-1)}}{(m+l-1)!\,(m-1)!}. \tag{17}$$

Truncating this series after $M$ terms and substituting the result into (10a), we obtain

$$\kappa_{lp} g_{lp}^{(2)}(r_2) = 2\pi \int_{r_o}^{r_m} dr_1\, r_1 \sum_{m-1}^{M} \frac{(-1)^{m-1}(\pi r_1 r_2)^{l+2(m-1)}}{(m+l-1)!\,(m-1)!}\, g_{lp}^{(1)}(r_1) \tag{18a}$$

$$= \left[\frac{(\pi r_2^2)^l}{l!}\right]^{\frac{1}{2}} \sum_{m=1}^{M} \frac{(-1)^{m-1}(\pi r_2^2)^{m-1}}{[(m-1)!\,(m+l-1)!/l!]^{\frac{1}{2}}}\, G_m^{(1)}(l,\,p), \tag{18b}$$

where

$$G_m^{(1)}(l,\,p) = \frac{2\pi}{[(m+l-1)!\,(m-1)!]^{\frac{1}{2}}} \int_{r_o}^{r_m} dr_1\, r_1(\pi r_1^2)^{m-1+l/2}\, g_{lp}^{(1)}(r_1). \tag{19}$$

Likewise we obtain from (10b)

$$\kappa_{lp}g_{lp}^{(1)}(r_1) = \left[\frac{(\pi r_1^2)^l}{l!}\right]^{\frac{1}{2}} \sum_{m=1}^{M} \frac{(-1)^{m-1}(\pi r_1^2)^{m-1}}{[(m-1)!\,(m+l-1)!/l!]^{\frac{1}{2}}} G_m^{(2)}(l,\,p), \quad (20)$$

where

$$G_m^{(2)}(l,\,p) = \frac{2\pi}{[(m+l-1)!\,(m-1)!]^{\frac{1}{2}}} \int_o^{r_m} dr_2\,r_2(\pi r_2^2)^{m-1+l/2}g_{lp}^{(2)}(r_2). \quad (21)$$

Substituting the expression (20) for $g_{lp}^{(1)}(r_1)$ into the right hand side of (19) and the expression (18b) for $g_{lp}^{(2)}(r_2)$ into (21), we obtain after simple manipulations

$$\kappa_{lp}G_m^{(1)}(l,\,p)$$

$$= \sum_{k=1}^{M} \frac{(-1)^{k-1}[(\pi N_m)^{l+m+k-1} - (\pi N_o)^{l+m+k-1}]}{[(m-1)!\,(m+l-1)!\,(k-1)!\,(k+l-1)!]^{\frac{1}{2}}(l+m+k-1)} G_k^{(2)}(l,\,p),$$

$$(22a)$$



Fig. 7 — (a) Field amplitude $g_{lp}^{(i)}(r)$ and (b) field intensity $|g_{lp}^{(i)}(r)|^2$ for modes $(lp) = (00)$ and $(01)$ with $N_m = 0.8$ and $N_o = 0.03$. The dashed lines refer to mirror 1 (Fig. 1) and the solid lines to mirror 2. The radius of the aperture in mirror 1 is $r_o = 0.173$.

Fig. 8 — (a) Field amplitude $g_{lp}(r)$ and (b) field intensity $\mid g_{lp}(r) \mid^2$ for modes $(lp) = (00)$ and $(01)$ with $N_m = 1.6$ and $N_o = 0$. The field distributions are identical on both mirrors.

$$\kappa_{lp} G_m^{(2)}(l,\,p)$$

$$= \sum_{k=1}^{M} \frac{(-1)^{k-1}(\pi N_m)^{l+m+k-1}}{[(m-1)!\,(m+l-1)!\,(k-1)!\,(k+l-1)!]^{\frac{1}{2}}(l+m+k-1)}\, G_k^{(1)}(l,\,p).$$

(22b)

We have used the definitions (5b) and (7) to replace $(r_m^2,\,r_o^2)$ by the Fresnel numbers $(N_m,\,N_o)$.

It is convenient to view $G_m^{(i)}(l,\,p)$ as the $m$th component of an $M$-dimensional vector $\mathbf{G}^{(i)}(l,\,p)$. We define $\mathbf{B}$ to be the $M \times M$ diagonal matrix with elements $B_{mm} = (-1)^{m-1}$. We also define real symmetric matrices $\mathbf{S}^{(1)}(l)$ and $\mathbf{S}^{(2)}(l)$ with elements

$$S_{mk}^{(1)}(l) = \frac{(\pi N_m)^{l+m+k-1}}{[(m-1)!\,(m+l-1)!\,(k-1)!\,(k+l-1)!]^{\frac{1}{2}}(l+m+k-1)},\quad (23a)$$

$$S_{mk}^{(2)}(l) = \frac{[(\pi N_m)^{l+m+k-1} - (\pi N_o)^{l+m+k-1}]}{[(m-1)!\,(m+l-1)!\,(k-1)!\,(k+l-1)!]^{\frac{1}{2}}(l+m+k-1)}.\quad (23b)$$

With these definitions equations (22) can be written more compactly

$$\kappa_{lp}\mathbf{G}^{(1)}(l, p) = \mathbf{S}^{(2)}(l)\cdot\mathbf{B}\cdot\mathbf{G}^{(2)}(l, p), \qquad (24a)$$

$$\kappa_{lp}\mathbf{G}^{(2)}(l, p) = \mathbf{S}^{(1)}(l)\cdot\mathbf{B}\cdot\mathbf{G}^{(1)}(l, p). \qquad (24b)$$

Eliminating $\mathbf{G}^{(2)}$, we obtain

$$\kappa_{lp}^2\mathbf{G}^{(1)}(l, p) = \mathbf{S}^{(2)}(l)\cdot\mathbf{B}\cdot\mathbf{S}^{(1)}(l)\cdot\mathbf{B}\cdot\mathbf{G}^{(1)}(l, p), \qquad (25)$$

which is a single $M$-dimensional matrix eigenvalue equation.

It is generally useful to transform equation (25) such that the matrix to be diagonalized is real symmetric (or Hermitian). Because $\mathbf{S}^{(2)}(l)$ is real symmetric with nonnegative eigenvalues, we can find a real lower-triangular matrix $\mathbf{P}(l)$ such that

$$\mathbf{S}^{(2)}(l) = \mathbf{P}(l)\cdot\mathbf{P}(l)^T. \qquad (26)$$



Fig. 9 — (a) Field amplitude $g_{lp}{}^{(j)}(r)$ and (b) field intensity $|g_{lp}{}^{(j)}(r)|^2$ for modes $(lp) = (00)$ and $(01)$ with $N_m = 1.6$ and $N_o = 0.0001$. The dashed lines refer to mirror 1 (Fig. 1) and the solid lines to mirror 2. The radius of the aperture in mirror 1 is $r_o = 0.01$.

If we define a new vector

$$\mathbf{F}(l,\,p) \;=\; \mathbf{P}(l)^{-1} \cdot \mathbf{G}^{(1)}(l,\,p), \tag{27}$$

then (25) can be written

$$\kappa_{lp}^{2}\mathbf{F}(l,\,p) \;=\; \mathbf{P}(l)^{T} \cdot \mathbf{B} \cdot \mathbf{S}^{(1)}(l) \cdot \mathbf{B} \cdot \mathbf{P}(l) \cdot \mathbf{F}, \tag{28}$$

for which the matrix on the right hand side is obviously real symmetric. If $\mathbf{U}(l)$ is the real orthogonal matrix which diagonalizes this matrix, then

$$\mathbf{U}(l) \cdot \mathbf{K}(l) \;=\; \mathbf{P}(l)^{T} \cdot \mathbf{B} \cdot \mathbf{S}^{(1)}(l) \cdot \mathbf{B} \cdot \mathbf{P}(l) \cdot \mathbf{U}(l) \tag{29}$$

where $\mathbf{K}(l)$ is diagonal with elements $K_{pp}(l) = \kappa_{lp}^{2}$, $p = 1$ to $M$. The eigenvector $\mathbf{F}(l,\,p)$ of (28) corresponds to the $p$th column of $\mathbf{U}(l)$ and, from (27),

$$\mathbf{G}^{(1)}(l,\,p) \;=\; \mathbf{P}(l) \cdot \mathbf{F}(l,\,p). \tag{30}$$



Fig. 10 — (a) Field amplitude $g_{lp}^{(j)}(r)$ and (b) field intensity $|\,g_{lp}^{(j)}(r)\,|^{2}$ for modes $(lp) = (00)$ and $(01)$ with $N_m = 1.6$ and $N_o = 0.001$. The dashed lines refer to mirror 1 (Fig. 1) and the solid lines to mirror 2. The radius of the aperture in mirror 1 is $r_o = 0.0316$.
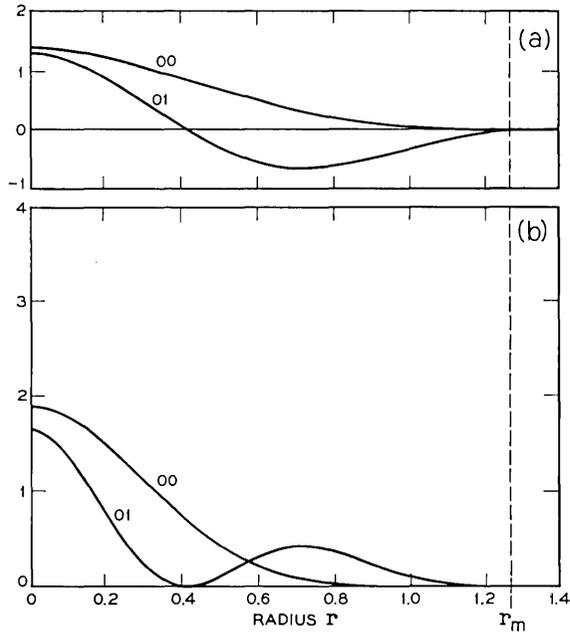
Fig. 11 — (a) Field amplitude $g_{lp}^{(i)}(r)$ and (b) field intensity $|g_{lp}^{(i)}(r)|^2$ for modes $(lp) = (00)$ and $(01)$ with $N_m = 1.6$ and $N_o = 0.01$. The dashed lines refer to mirror 1 (Fig. 1) and the solid lines to mirror 2. The radius of the aperture in mirror 1 is $r_o = 0.1$.

The elements of $\mathbf{U}(l)$ and $\mathbf{K}(l)$ are easily computed numerically; the vectors $\mathbf{G}^{(1)}(l, p)$ follow from (30); the vectors $\mathbf{G}^{(2)}(l, p)$ follow from (24b); and the amplitudes $g_{lp}^{(i)}(r)$ follow from (18b) and (20).

The program used to compute the results reported in this paper required a nominal 0.0003 hr. of GE 635 processor time to compute the $M$ different eigenvalues $|\kappa_{lp}|$ and eigenvectors $[\mathbf{G}^{(1)}(l, p), \mathbf{G}^{(2)}(l, p)]$ for $M = 20$. Timing for other values of $M$ varies roughly as $M^3$.

REFERENCES

1. McCumber, D. E., "Eigenmodes of a Symmetric Cylindrical Confocal Laser Resonator and their Perturbation by Output-Coupling Apertures," B.S.T.J., *44*, No. 2 (February 1965), pp. 333–363.
2. Patel, C. K. N., Faust, W. L., McFarlane, R. A., and Garrett, C. G. B., "Laser Action up to 57.355μ in Gaseous Discharges (Ne, He-Ne)," Appl. Phys. Letters, *4*, No. 1 (January 1, 1964), pp. 18–19.
3. Fox, A. G., and Li, T., "Resonant Modes in a Maser Interferometer," B.S.T.J., *40*, No. 2 (March 1961), pp. 453–488.

4. Boyd, G. D., and Kogelnik, H., "Generalized Confocal Resonator Theory," B.S.T.J., *41*, No. 4 (July 1962), pp. 1347–1369.
5. Kogelnik, H., "Modes in Optical Resonators," in *Lasers*, vol. 1, ed. A. K. Levine, New York: Marcel Dekker, Inc., 1966, pp. 295–347.
6. Gordon, J. P., and Kogelnik, H., "Equivalence Relations Among Spherical Mirror Optical Resonators," B.S.T.J., *43*, No. 6 (November 1964), pp. 2873–2886.
7. Goubau, G., and Schwering, F., "On the Guided Propagation of Electromagnetic Wave Beams," IRE Trans. Antennas Propagation, *AP-9*, No. 3 (May 1961), pp. 248–256.
8. Beyer, J. B., and Scheibe, E. H., "Higher Modes in Guided Electromagnetic-Wave Beams," IRE Trans. Antennas Propagation, *AP-10*, No. 3 (May 1962), pp. 349–350.
9. Slepian, D., "Prolate Spheroidal Wave Functions, Fourier Analysis and Uncertainty—IV: Extensions to Many Dimensions; Generalized Prolate Spheroidal Functions," B.S.T.J., *43*, No. 6 (November 1964), pp. 3009–3057.
10. Faust, W. L., unpublished work.

# Some Basic Characteristics of Broadband Negative Resistance Oscillator Circuits

By K. KUROKAWA

(Manuscript received January 14, 1969)

*This paper discusses the behavior of oscillators with multiple-resonant circuits. It discusses the condition for free-running stable oscillations, the injection locking phenomena, the stable locking range, the noise of free-running and injection-locked oscillators, and a condition for parasitic oscillations in detail, and presents a graphical interpretation of this study for clarity. Finally, this paper shows how broadbanding of oscillators can be achieved with a double-resonant circuit. This provides a systematic guide for the design of broadband frequency deviators and broadband injection-locked oscillators for numerous applications.*

## I. INTRODUCTION

With the advent of high-frequency negative resistance elements such as tunnel diodes, Gunn diodes, and IMPATT's, it now appears practical, at microwave frequencies, to build negative resistance oscillators which perform various functions other than that of a fixed frequency oscillator.[1] Among the proposals recently made are variable frequency oscillators as FM deviators, and locked oscillators as FM amplifiers, limiters, and FM demodulators.[2-4] In the past, the analysis of these oscillators was primarily based on a simple model with a single-resonant circuit.[5-8] In practice, however, the designer of such a microwave circuit provides a number of tuning elements and adjusts them by trial and error until a desired bandwidth of locking or frequency deviation is obtained. During this adjustment, the designer observes numerous phenomena, including sudden changes in noise or oscillating frequency, as well as various hysteresis effects which could not be explained by the simple oscillator model.

This paper presents a more realistic model of oscillators in which the load is separated from the active device by a multiple-resonant circuit. This gives us a better understanding of the oscillator behavior

and, hopefully, a more systematic approach toward an oscillator design for broadband applications.

## II. EQUATIONS FOR AMPLITUDE AND PHASE OF OSCILLATING CURRENT

Consider the circuit shown in Fig. 1. The active device is represented by $-\bar{R} + j\bar{X}$ and the load by $Z_L$. The box between the active device and the load represents a multiple-resonant circuit; $Z(\omega)$ is the impedance looking into the box from the active device. As far as the active device is concerned, the entire circuit can be expressed as a series connection of $Z(\omega)$ and $-\bar{R} + j\bar{X}$, as shown in Fig. 2 where $e(t)$ represents noise or locking signal voltages that may be present. Let the current flowing through the active device be

$$i(t) = A \cos (\omega t + \varphi), \tag{1}$$

where $A$ and $\varphi$ are assumed to be slowly varying functions of time. From the beginning, the effects of higher harmonics are neglected since they disappear in the process of averaging over one period of the oscillation which we perform later.[7,8] The voltage drop across the active devices is given by

$$v = -\bar{R}A \cos (\omega t + \varphi) - \bar{X}A \sin (\omega t + \varphi) \tag{2}$$

where $\bar{R}$ and $\bar{X}$ are functions of the current amplitude $A$. In the following discussion the frequency dependence of $\bar{R}$ and $\bar{X}$ are neglected, which is generally justifiable over the frequency range of interest. Furthermore, the magnitude of $\bar{X}$ is assumed to be small, which is usually justifiable after lumping any large constant portion of the device reactance in with $Z(\omega)$.

In order to calculate the voltage drop across $Z(\omega)$, consider the time derivative of $i(t)$. The first derivative is given by

$$\frac{di}{dt} = -A\left(\omega + \frac{d\varphi}{dt}\right) \sin (\omega t + \varphi) + \frac{dA}{dt} \cos (\omega t + \varphi)$$



Fig. 1 — Schematic diagram for an oscillator.

Fig. 2 — An equivalent circuit of the oscillator in Fig. 1.

$$= \mathrm{Re}\left\{\left[j\left(\omega + \frac{d\varphi}{dt}\right) + \frac{1}{A}\frac{dA}{dt}\right]Ae^{j(\omega t + \varphi)}\right\}.$$

Similarly, to first order approximation, the $n$th derivative is given by

$$\frac{d^n i}{dt^n} \cong \mathrm{Re}\left\{\left[j\left(\omega + \frac{d\varphi}{dt}\right) + \frac{1}{A}\frac{dA}{dt}\right]^n Ae^{j(\omega t + \varphi)}\right\}.$$

In ac circuit theory, the time derivative corresponds to multiplication by $j\omega$. Therefore, replacing $\omega$ by $\omega + d\varphi/dt - j(1/A)(dA/dt)$ everywhere in $Z(\omega)$, Re $\{ZI\}$ should give the desired voltage drop, where $I$ is the expression for $i(t)$ in the form of $Ae^{j(\omega t + \varphi)}$. Assuming $d\varphi/dt \ll \omega$ and $(1/A)(dA/dt) \ll \omega$, we have

$$Z\left(\omega + \frac{d\varphi}{dt} - j\frac{1}{A}\frac{dA}{dt}\right)$$

$$\cong Z(\omega) + \frac{dZ(\omega)}{d\omega}\left(\frac{d\varphi}{dt} - j\frac{1}{A}\frac{dA}{dt}\right)$$

$$= R(\omega) + jX(\omega) + [R'(\omega) + jX'(\omega)]\left(\frac{d\varphi}{dt} - j\frac{1}{A}\frac{dA}{dt}\right) \quad (3)$$

where $R(\omega) = \mathrm{Re}\,[Z(\omega)]$, $X(\omega) = \mathrm{Im}\,[Z(\omega)]$, and the primes indicate the derivatives with respect to $\omega$. Using the above approximate expression, Re $\{ZI\}$ is calculated to be

$$\mathrm{Re}\,\{ZI\} = \left[R(\omega) + R'(\omega)\frac{d\varphi}{dt} + X'(\omega)\frac{1}{A}\frac{dA}{dt}\right]A\,\cos\,(\omega t + \varphi)$$

$$- \left[X(\omega) + X'(\omega)\frac{d\varphi}{dt} - R'(\omega)\frac{1}{A}\frac{dA}{dt}\right]A\,\sin\,(\omega t + \varphi).$$

Referring to Fig. 2, we have $v + \mathrm{Re}\,\{ZI\} = e(t)$ which is equivalent to

$$\left[ R(\omega) - \bar{R} + R'(\omega) \frac{d\varphi}{dt} + X'(\omega) \frac{1}{A} \frac{dA}{dt} \right] A \cos (\omega t + \varphi)$$

$$- \left[ X(\omega) + \bar{X} + X'(\omega) \frac{d\varphi}{dt} - R'(\omega) \frac{1}{A} \frac{dA}{dt} \right] A \sin (\omega t + \varphi) = e(t).$$

After multiplying this equation by cos $(\omega t + \varphi)$ and sin $(\omega t + \varphi)$, respectively and integrating over one period of oscillation, we obtain

$$R(\omega) - \bar{R} + R'(\omega) \frac{d\varphi}{dt} + X'(\omega) \frac{1}{A} \frac{dA}{dt} = \frac{1}{A} e_c(t) \qquad (4)$$

$$-X(\omega) - \bar{X} - X'(\omega) \frac{d\varphi}{dt} + R'(\omega) \frac{1}{A} \frac{dA}{dt} = \frac{1}{A} e_s(t) \qquad (5)$$

where

$$e_c(t) = \frac{2}{T_o} \int_{t-T_o}^{t} e(t) \cos (\omega t + \varphi) \, dt \qquad (6)$$

$$e_s(t) = \frac{2}{T_o} \int_{t-T_o}^{t} e(t) \sin (\omega t + \varphi) \, dt \qquad (7)$$

and $T_o$ is the oscillation period. Equations (4) and (5) both contain $d\varphi/dt$ and $dA/dt$. However, by multiplying equation (4) by $X'(\omega)$ and equation (5) by $R'(\omega)$ and adding, an equation with $dA/dt$ alone is obtained. Similarly, multiplying equation (4) by $R'(\omega)$ and equation (5) by $-X'(\omega)$ and adding, gives an equation with $d\varphi/dt$ alone. These are

$$[R(\omega) - \bar{R}]X'(\omega) - [X(\omega) + \bar{X}]R'(\omega) + |Z'(\omega)|^2 \frac{1}{A} \frac{dA}{dt}$$

$$= \frac{1}{A} [X'(\omega)e_c(t) + R'(\omega)e_s(t)] \qquad (8)$$

$$[R(\omega) - \bar{R}]R'(\omega) + [X(\omega) + \bar{X}]X'(\omega) + |Z'(\omega)|^2 \frac{d\varphi}{dt}$$

$$= \frac{1}{A} [R'(\omega)e_c(t) - X'(\omega)e_s(t)]. \qquad (9)$$

Equations (8) and (9) are the basic equations for the amplitude and phase of an oscillating current.

For a steady-state free-running oscillation, we assume $e_c(t) = e_s(t) = 0$, $dA/dt = 0$ and $d\varphi/dt = 0$. Thus, from equations (8) and (9)

we have

$$R(\omega) - \bar{R} = 0, \quad X(\omega) + \bar{X} = 0 \tag{10}$$

which determine the amplitude $A_o$ and the frequency $\omega_o$ of the oscillation. Suppose $A$ somehow deviates from its steady-state value $A_o$ by a small amount $\delta A$. We then have

$$R(\omega_o) - \bar{R} = \frac{\delta A}{A_o} sR_o , \qquad X(\omega_o) + \bar{X} = \frac{\delta A}{A_o} rR_o \tag{11}$$

where $R_o$ indicates the value of $\bar{R}$ at $A = A_o$ [which is also equal to $R(\omega_o)$], and $sR_o$ and $rR_o$ represent $A_o(-\partial\bar{R}/\partial A)$ and $A_o(\partial\bar{X}/dA)$, respectively, as indicated in Fig. 3.

From equations (8) and (11), we obtain a differential equation for $\delta A$. To first-order approximation, it is

$$sR_o X'(\omega_o) \frac{\delta A}{A_o} - rR_o R'(\omega_o) \frac{\delta A}{A_o} + | Z'(\omega_o) |^2 \frac{1}{A_o} \frac{d\ \delta A}{dt} = 0.$$

If $\delta A$ decays with time, we have

$$sR_o X'(\omega_o) - rR_o R'(\omega_o) > 0. \tag{12}$$

An operating point determined by equation (10) is stable if and only if it satisfies condition (12).

## III. CONDITIONS FOR INJECTION LOCKING

Let us next investigate injection locking of the oscillator. The injection signal is represented by



Fig. 3 — $\bar{R}$ versus $A$; $\bar{X}$ versus $A$.

$$e(t) = a_o \cos \omega_i t.$$

Then, from (6) and (7), we have

$$e_c(t) = a_o \cos \varphi, \quad e_s(t) = a_o \sin \varphi. \tag{13}$$

Since $\omega_i$ may differ from $\omega_o$, and $A$ from $A_o$, we write

$$R(\omega_i) - \bar{R} = \frac{\triangle A}{A_o} sR_o + \triangle R, \quad X(\omega_i) + \bar{X} = \frac{\triangle A}{A_o} rR_o + \triangle X \tag{14}$$

where $\triangle R = R(\omega_i) - R(\omega_o)$ and $\triangle X = X(\omega_i) - X(\omega_o)$. Substituting (13) and (14) into (8) and (9) gives

$$\left( \frac{\triangle A}{A_o} sR_o + \triangle R \right) X'(\omega_i)$$

$$- \left( \frac{\triangle A}{A_o} rR_o + \triangle X \right) R'(\omega_i) + \mid Z'(\omega_i) \mid^2 \frac{1}{A_o} \frac{d \triangle A}{dt}$$

$$= \frac{a_o}{A_o} \{ X'(\omega_i) \cos \varphi + R'(\omega_i) \sin \varphi \} \tag{15}$$

$$\left( \frac{\triangle A}{A_o} sR_o + \triangle R \right) R'(\omega_i) + \left( \frac{\triangle A}{A_o} rR_o + \triangle X \right) X'(\omega_i) + \mid Z'(\omega_i) \mid^2 \frac{d\varphi}{dt}$$

$$= \frac{a_o}{A_o} \{ R'(\omega_i) \cos \varphi - X'(\omega_i) \sin \varphi \}. \tag{16}$$

When the oscillator is locked $d \triangle A/dt = 0$, $d\varphi/dt = 0$. Substituting these conditions into (15) and (16) and eliminating $\triangle A$ we have

$$\frac{a_o}{A_o} (s^2 + r^2)^{\frac{1}{2}} \sin (\varphi_o + \theta) = r \triangle R - s \triangle X \tag{17}$$

where $\varphi_o$ is the phase of the steady-state current and $\theta$ is defined by

$$\tan \theta = r/s. \tag{18}$$

Since the magnitude of $\sin(\varphi_o + \theta)$ must be less than one, (17) leads to the condition

$$\frac{A_o}{a_o} \frac{\mid r \triangle R - s \triangle X \mid}{(s^2 + r^2)^{\frac{1}{2}}} \leqq 1 \tag{19}$$

which determines the possible locking range.

Next, assume that $\triangle A$ and $\varphi$ deviate from the steady-state values $\triangle A_o$ and $\varphi_o$ and write $\triangle A = \triangle A_o + \delta A$ and $\varphi = \varphi_o + \delta\varphi$. Then, (15)

and (16) become

$$\frac{\delta A}{A_o} [sR_oX'(\omega_i) - rR_oR'(\omega_i)] + \mid Z'(\omega_i) \mid^2 \frac{1}{A_o} \frac{d \, \delta A}{dt}$$

$$= \frac{a_o}{A_o} [-X'(\omega_i) \sin \varphi_o + R'(\omega_i) \cos \varphi_o] \, \delta \varphi \qquad (20)$$

$$\frac{\delta A}{A_o} [sR_oR'(\omega_i) + rR_oX'(\omega_i)] + \mid Z'(\omega_i) \mid^2 \frac{d \, \delta \varphi}{dt}$$

$$= \frac{a_o}{A_o} [-R'(\omega_i) \sin \varphi_o - X'(\omega_i) \cos \varphi_o] \, \delta \varphi. \qquad (21)$$

Eliminating $\delta A$ from equations (20) and (21) gives

$$\mid Z'(\omega_i) \mid^2 \frac{d^2 \, \delta \varphi}{dt^2} + \left\{ sR_oX'(\omega_i) - rR_oR'(\omega_i) \right.$$

$$\left. + \frac{a_o}{A_o} [R'(\omega_i) \sin \varphi_o + X'(\omega_i) \cos \varphi_o] \right\} \frac{d \, \delta \varphi}{dt}$$

$$+ \frac{a_o}{A_o} R_o(s^2 + r^2)^{\frac{1}{2}} \cos (\varphi_o + \theta) \, \delta \varphi = 0.$$

For stable operation, $\mid \delta \varphi \mid$ must decrease with time if it is not already zero. Thus, we have two conditions

$$sR_oX'(\omega_i) - rR_oR'(\omega_i) + \frac{a_o}{A_o} [R'(\omega_i) \sin \varphi_o + X'(\omega_i) \cos \varphi_o] > 0$$
$$\qquad (22)$$

$$\frac{a_o}{A_o} R_o(s^2 + r^2)^{\frac{1}{2}} \cos (\varphi_o + \theta) > 0.$$

The first is usually satisfied if $sR_oX'(\omega_i) - rR_oR'(\omega_i)$ is positive, as required for a stable free-running oscillation at $\omega_i$ (not $\omega_o$). The second condition is equivalent to

$$\cos (\varphi_o + \theta) > 0. \qquad (23)$$

The term $\varphi_o$ is uniquely determined from equations (17) and (23). Using this $\varphi_o$, the left-hand side of equation (22) must then be calculated to check whether or not the locking is stable.

   The above discussion was in terms of voltage and current. Since power is a more meaningful quantity at microwave frequencies, let us rewrite equations (17), (19), and (22). For simplicity, we assume that the box in Fig. 1 is lossless. Since the available power is invariant

under a lossless transformation the injected signal power $P_i$ is given by

$$P_i = \frac{a_o^2}{8(R_o + \triangle R)}.$$

The free-running output power is equal to

$$P_o = R_o \frac{A_o^2}{2}.$$

Using the above relations, equation (17) can be rewritten as

$$\triangle R \sin \theta - \triangle X \cos \theta \cong 2R_o(P_i/P_o)^{\frac{1}{2}} \sin (\varphi_o + \theta) \qquad (24)$$

where higher order terms of $\triangle R$ are neglected. The locking range is determined from

$$| \triangle R \sin \theta - \triangle X \cos \theta | \leq 2R_o(P_i/P_o)^{\frac{1}{2}} \qquad (25)$$

$$(s^2 + r^2)^{\frac{1}{2}}R_o \sin (\Theta - \theta) + 2R_o(P_i/P_o)^{\frac{1}{2}} \sin (\Theta + \varphi_o) > 0 \qquad (26)$$

where

$$\tan \Theta = \frac{X'(\omega_i)}{R'(\omega_i)}. \qquad (27)$$

Equations (25) and (26) correspond to equations (19) and (22), respectively. To the same order of approximation, the output power of the locked oscillator is given by

$$P = P_o \left\{ 1 + \frac{\triangle R}{R_o} - \frac{2}{R_o} \frac{\triangle R X'(\omega_i) - \triangle X R'(\omega_i)}{s X'(\omega_i) - r R'(\omega_i)} \right.$$
$$\left. + 4\left(\frac{P_i}{P_o}\right)^{\frac{1}{2}} \left[ \frac{X'(\omega_i) \cos \varphi_o + R'(\omega_i) \sin \varphi_o}{s X'(\omega_i) - r R'(\omega_i)} - \frac{1}{2} \cos \varphi_o \right] \right\} \qquad (28)$$

where the power wave concept has been used in the calculation.[9]

## IV. NOISE AND PARASITIC OSCILLATIONS

Oscillator noise near the carrier can be discussed in terms of fluctuations in $A$ and $\varphi$, just as in a single-tuned oscillator. Let $e(t)$ indicate a noise voltage in equations (6) and (7); then we have (see Ref. 5)

$$e_c(t) = n_1(t), \quad e_s(t) = n_2(t)$$

where $n_1(t)$ and $n_2(t)$ represent the cosine and sine components of the noise voltage $e(t)$ and

$$\overline{n_1(t)^2} = \overline{2e(t)^2}, \qquad \overline{n_2(t)^2} = \overline{2e(t)^2}, \qquad \overline{n_1(t)n_2(t)} = 0.$$

Substituting the above expressions into equations (8) and (9), we obtain

$$sR_oX'(\omega_o)\ \delta A\ -\ rR_oR'(\omega_o)\ \delta A\ +\ |\ Z'(\omega_o)\ |^2\ \frac{d\ \delta A}{dt}$$

$$=\ X'(\omega_o)n_1(t)\ +\ R'(\omega_o)n_2(t)$$

$$sR_oR'(\omega_o)\ \delta A\ +\ rR_oX'(\omega_o)\ \delta A\ +\ |\ Z'(\omega_o)\ |^2\ A_o\frac{d\varphi}{dt}$$

$$=\ R'(\omega_o)n_1(t)\ -\ X'(\omega_o)n_2(t)$$

from which the frequency spectra of $\delta A$ and $\varphi$ are calculated as

$$|\ \delta A(f)\ |^2\ =\ \frac{2\ |\ Z'(\omega_o)\ |^2\ |\ e\ |^2}{\omega^2\ |\ Z'(\omega_o)\ |^4\ +\ [sR_oX'(\omega_o)\ -\ rR_oR'(\omega_o)]^2} \qquad (29)$$

$$|\ \varphi(f)\ |^2\ =\ \frac{2\ |\ e\ |^2}{\omega^2A_o^2}\ \frac{\omega^2\ |\ Z'(\omega_o)\ |^2\ +\ (s^2\ +\ r^2)R_o^2}{\omega^2\ |\ Z'(\omega_o)\ |^4\ +\ [sR_oX'(\omega_o)\ -\ rR_oR'(\omega_o)]^2}. \qquad (30)$$

These equations indicate that the oscillation becomes very noisy as we approach the boundary of the stable region defined by equation (12). This is what we expect. However, equation (29) and (30) are not valid near the boundary since $\delta A$ is no longer small, as was initially assumed.

Sometimes we need to consider additional current components at frequencies $\omega \pm \triangle\omega$, so far away from the oscillating frequency $\omega$ that, although $\triangle\omega \ll \omega$, equation (3) is no longer a valid approximation. In this case, we express the current in the form

$$i(t)\ =\ A\ \cos\ (\omega t\ +\ \varphi)\ +\ A_+\ \cos\ \{(\omega\ +\ \triangle\omega)t\ +\ \varphi_+\}$$

$$+\ A_-\ \cos\ \{\omega\ -\ \triangle\omega)t\ +\ \varphi_-\} \qquad (31)$$

where $A$, $\varphi$, $A_+$, $\varphi_+$ $A_-$, and $\varphi_-$ are all slowly varying functions of time. Assuming that $A_+ \ll A$ and $A_- \ll A$, and neglecting higher order terms of small quantities, the current can be expressed as

$$i(t)\ =\ \tilde{A}\ \cos\ (\omega t\ +\ \tilde{\varphi})$$

where

$$\tilde{A}\ =\ A\ +\ A_+\ \cos\ (\triangle\omega t\ +\ \varphi_+\ -\ \varphi)\ +\ A_-\ \cos\ (\triangle\omega t\ -\ \varphi_-\ +\ \varphi)$$

$$\tilde{\varphi}\ =\ \varphi\ +\ \frac{1}{A}\ \{A_+\ \sin\ (\triangle\omega t\ +\ \varphi_+\ -\ \varphi)\ -\ A_-\ \sin\ (\triangle\omega t\ -\ \varphi_-\ +\ \varphi)\}.$$

The voltage drop across the active device is then given by

$$
\begin{aligned}
v = {}& -\bar{R}(\tilde{A})\tilde{A}\,\cos\,(\omega t + \tilde{\varphi}) - \bar{X}(\tilde{A})\tilde{A}\,\sin\,(\omega t + \tilde{\varphi}) \\
\cong {}& -\bar{R}A\,\cos\,(\omega t + \varphi) - \bar{X}A\,\sin\,(\omega t + \varphi) \\
& - \left(2\bar{R} + \frac{\partial\bar{R}}{\partial A}\,A\right)\frac{A_+}{2}\,\cos\,[(\omega + \triangle\omega)t + \varphi_+] \\
& - \frac{\partial\bar{R}}{\partial A}\,A\,\frac{A_-}{2}\,\cos\,[(\omega + \triangle\omega)t + 2\varphi - \varphi_-] \\
& - \left(2\bar{X} + \frac{\partial\bar{X}}{\partial A}\,A\right)\frac{A_+}{2}\,\sin\,[(\omega + \triangle\omega)t + \varphi_+] \\
& - \frac{\partial\bar{X}}{\partial A}\,A\,\frac{A_-}{2}\,\sin\,[(\omega + \triangle\omega)t + 2\varphi - \varphi_-] \\
& - \left(2\bar{R} + \frac{\partial\bar{R}}{\partial A}\,A\right)\frac{A_-}{2}\,\cos\,[(\omega - \triangle\omega)t + \varphi_-] \\
& - \frac{\partial\bar{R}}{\partial A}\,A\,\frac{A_+}{2}\,\cos\,[(\omega - \triangle\omega)t + 2\varphi - \varphi_+] \\
& - \left(2\bar{X} + \frac{\partial\bar{X}}{\partial A}\,A\right)\frac{A_-}{2}\,\sin\,[(\omega - \triangle\omega)t + \varphi_-] \\
& - \frac{\partial\bar{X}}{\partial A}\,A\,\frac{A_-}{2}\,\sin\,[(\omega - \triangle\omega)t + 2\varphi - \varphi_+] \qquad (32)
\end{aligned}
$$

where $\bar{R}$ and $\bar{X}$ express the values at $A$, that is, $\bar{R}(A)$ and $\bar{X}(A)$, respectively. To derive equation (32), an assumption was made that the values of $\bar{R}$ and $\bar{X}$ change with time following the envelope of the current, which is generally a valid assumption when $\triangle\omega \ll \omega$, as assumed here. The voltage drop across the impedance $Z(\omega)$ can be calculated as follows.

$$
\begin{aligned}
\mathrm{Re}\ \{ZI\} = {}& \left[AR(\omega) + AR'(\omega)\frac{d\varphi}{dt} + X'(\omega)\frac{dA}{dt}\right]\cos\,(\omega t + \varphi) \\
& - \left[AX(\omega) + AX'(\omega)\frac{d\varphi}{dt} - R'(\omega)\frac{dA}{dt}\right]\sin\,(\omega t + \varphi) \\
& + \left(A_+R_+ + A_+R'_+\frac{d\varphi_+}{dt} + X'_+\frac{dA_+}{dt}\right)\cos\,[(\omega + \triangle\omega)t + \varphi_+] \\
& - \left(A_+X_+ + A_+X'_+\frac{d\varphi_+}{dt} - R'_+\frac{dA_+}{dt}\right)\sin\,[(\omega + \triangle\omega)t + \varphi_+]
\end{aligned}
$$

$$+ \left( A_- R_- + A_- R'_- \frac{d\varphi_-}{dt} + X'_- \frac{dA_-}{dt} \right) \cos \left[ (\omega - \triangle\omega)t + \varphi_- \right]$$

$$- \left( A_- X_- + A_- X'_- \frac{d\varphi_-}{dt} - R'_- \frac{dA_-}{dt} \right) \sin \left[ (\omega - \triangle\omega)t + \varphi_- \right]$$

where $R_+ = R(\omega + \triangle\omega)$, $X_+ = X(\omega + \triangle\omega)$, $R_- = R(\omega - \triangle\omega)$ and $X_- = X(\omega - \triangle\omega)$.

Let $n$ be defined by $2\pi/T_o \triangle\omega$. Multiplying Re $\{ZI\} + v = e(t)$ by $\cos (\omega t + \varphi)$ and $\sin (\omega t + \varphi)$, and integrating the results over a period of $nT_o$, we obtain approximate equations for $A$ and $\varphi$. These are identical to equations (4) and (5), except that $e_c(t)$ and $e_s(t)$ are now defined by

$$e_c(t) = \frac{2}{nT_o} \int_{t-nT_o}^{t} e(t) \cos (\omega t + \varphi) \, dt,$$

$$e_s(t) = \frac{2}{nT_o} \int_{t-nT_o}^{t} e(t) \sin (\omega t + \varphi) \, dt.$$

This means that $A$ and $\varphi$ of free running and of injection-locked oscillators behave exactly in the same manner. For example, the amplitude of the free-running oscillation is determined by equation (10) and hence $(\partial \bar{R}/\partial A)A$ and $(\partial \bar{X}/\partial A)A$ in equation (32) can be replaced by $-sR_o$ and $rR_o$, respectively. The integration of Re $\{ZI\} + v = e(t)$ over a period of $nT_o$ after multiplying by $\cos [(\omega + \triangle\omega)t + \varphi_+]$, and. so on, gives

$$(R_+ - \bar{R} + \tfrac{1}{2}sR_o)A_+ + R'_+ A_+ \frac{d\varphi_+}{dt} + X'_+ \frac{dA_+}{dt}$$

$$+ (sR_o \cos \delta - rR_o \sin \delta) \frac{A_-}{2}$$

$$= \frac{2}{nT_o} \int_{t-nT_o}^{t} e(t) \cos \{ (\omega + \triangle\omega)t + \varphi_+ \} \, dt, \tag{33}$$

$$- (X_+ + \bar{X} + \tfrac{1}{2}rR_o)A_+ - X'_+ A_+ \frac{d\varphi_+}{dt} + R'_+ \frac{dA_+}{dt}$$

$$- (rR_o \cos \delta + sR_o \sin \delta) \frac{A_-}{2}$$

$$= \frac{2}{nT_o} \int_{t-nT_o}^{t} e(t) \sin \{ (\omega + \triangle\omega)t + \varphi_+ \} \, dt \tag{34}$$

and two similar equations, in which the subscripts $+$ and $-$ are interchanged and $+\triangle\omega$ is replaced by $-\triangle\omega$. These are the basic equations

determining the behavior of $A_+$, $A_-$, $\varphi_+$, and $\varphi_-$. In equations (33) and (34), $\delta$ stands for $2\varphi - \varphi_- - \varphi_+$.

When $e(t)$ equals zero, the above equations usually give $A_+ = A_- = 0$. However, if the condition

$$\frac{R_+ - \bar{R} + \tfrac{1}{2}sR_o}{X_+ + \bar{X} + \tfrac{1}{2}rR_o} = \frac{R_- - \bar{R} + \tfrac{1}{2}sR_o}{X_- + \bar{X} + \tfrac{1}{2}rR_o} \tag{35}$$

is satisfied, the steady-state values of $A_+$ and $A_-$ can become finite even when the integrals in equations (33), and (34) (and the other two equations) are all zero. Condition (35) is not satisfied by a single-tuned oscillator. With a multiple-resonant circuit, however, the locus of $Z(\omega)$ on the complex plane may form a loop (or loops) which intersects with itself; the above condition is often satisfied at a certain oscillation frequency (or frequencies). This means that as the oscillation frequency approaches this particular value the noise output power at $\omega \pm \triangle\omega$ increases enormously, resulting in a noisy parasitic oscillation; this first-order discussion then becomes invalid.

The noise of locked oscillators with a multiple-resonant circuit can be treated in a way similar to that of a single-tuned locked oscillator. For simplicity, let us assume that the locking signal is noise-free. Then, we have

$$|\, \delta A \,|^2 = [(\omega^2 A^2 \,|\, Z' \,|^2 + a_o^2)2 \,|\, e \,|^2]$$

$$\cdot \{[\omega^2 A \,|\, Z' \,|^2 - a_o(sR_o \cos \varphi_o - rR_o \sin \varphi_o)]^2$$

$$+ \omega^2 [A(sR_o X' - rR_o R') + a_o(R' \sin \varphi_o + X' \cos \varphi_o)]^2\}^{-1} \tag{36}$$

$$|\, \delta\varphi \,|^2 = \{[\omega^2 \,|\, Z' \,|^2 + (s^2 + r^2)R_o^2]2 \,|\, e \,|^2\}$$

$$\cdot \{[\omega^2 A \,|\, Z' \,|^2 - a_o(sR_o \cos \varphi_o - rR_o \sin \varphi_o)]^2$$

$$+ \omega^2 [A(sR_o X' - rR_o R') + a_o(R' \sin \varphi_o + X' \cos \varphi_o)]^2\}^{-1} \tag{37}$$

where $Z'$, $R'$ and $X'$ indicate their values at $\omega_i$.

The parasitic oscillations also take place in the injection-locked oscillators at the same frequency as before. This is because the injection-locked signal does not contribute to the values of the integrals in equation (33), and the condition for parasitic oscillations remains the same.

## V. GRAPHICAL INTERPRETATION

An excellent discussion of oscillator behaviors using a graphical method was already presented by Slater in his book *Microwave Elec-*

*tronics.* This section will extend his method first by including the case in which the locus of $Z(\omega)$ forms a loop (or loops),* and then by discussing injection locking.

Suppose that the locus of $Z(\omega)$ and the line representing $\bar{R} - j\bar{X}$ are drawn on the complex plane as shown in Fig. 4. The parameters are $\omega$ and $A$; the arrows indicate the directions of increasing $\omega$ and $A$. According to equation (10), the intersections of these two curves give possible operating points. For intersection (a), $sR'(\omega) > 0$ and

$$R_o\left(\frac{X'(\omega)}{R'(\omega)} - \frac{r}{s}\right) = R_o(\tan \Theta - \tan \theta) > 0.$$



Fig. 4 — $Z(\omega)$ and $\bar{R} - j\bar{X}$ on the complex plane.

From equation (12), (a) is a stable operating point. For intersection (b), $sR'(\omega) < 0$ and hence equation (12) becomes

$$R_0\left(\frac{X'(\omega)}{R'(\omega)} - \frac{r}{s}\right) = R_o(\tan \Theta - \tan \theta) < 0.$$

This inequality is satisfied and intersection (b) represents another stable operating point. However, intersection (c) is unstable since $sR'(\omega) < 0$ and $\tan \Theta - \tan \theta > 0$.

Which of the two stable operating points the oscillation selects depends on the oscillation's history. Suppose an oscillation corresponding to (a) is taking place and $\bar{X}$ is gradually decreased, then the curve $\bar{R} - j\bar{X}$ moves up and the intersection (a) moves toward the right until it reaches $(a')$. If we further decrease $\bar{X}$ the intersection $(a')$ disappears and the operating point jumps to $(b')$. Just before it jumps the oscillation becomes very noisy [since

$$sR_oX'(\omega) - rR_oR'(\omega) = (s^2 + r^2)^{\frac{1}{2}}R_o \mid Z'(\omega) \mid \sin (\Theta - \theta)$$

in the denominator of the noise expression (29) becomes zero and

---

* In this case $\omega = \omega_1 + j\omega_2$ becomes a multivalued function of impedance.

$| \delta A(f) |$ becomes infinite as $\omega \to 0$]. After the jump, if $\bar{X}$ is increased to its original value, operating point (b) is reached. Similarly, with the operating point initially at (b) we can realize operating point (a) by increasing $\bar{X}$ sufficiently and then bringing it back to the original value. If $\bar{X}$ cannot be changed, we might still be able to change the circuit so as to move the locus of $Z(\omega)$ up or down, thereby producing the same effect. During these adjustments, if three points $\bar{R} - \frac{1}{2}sR_o - j(\bar{X} + \frac{1}{2}rR_o)$, $R_+ + jX_+$, and $R_- + jX_-$ happen to fall on a straight line as illustrated in Fig. 5, then condition (35) is satisfied and we observe simultaneous oscillations at three different frequencies (sometimes more than three at $\omega \pm m\triangle\omega$, with $m$ representing integers).

Let us next consider the injection locking case. Suppose that the free-running operating point is at (a) in Fig. 6. In order to lock the oscillator at $\omega$ corresponding to point (e) in Fig. 6, the conditions for stable locking, equations (25) and (26), must be satisfied. Since $| \triangle R \sin \theta - \triangle X \cos \theta |$ corresponds to the distance from (e) to the line $\bar{R} - j\bar{X}$ as illustrated in Fig. 6a, we see that if this distance is less than $2R_o \cdot (P_o/P_o)^{\frac{1}{2}}$ the first condition is satisfied. Then, from equations (23) and (24) $\varphi$ can be graphically obtained as shown in Fig. 6b. In this particular case $\varphi$ is negative, but when (e) is located sufficiently below (a), $\varphi$ becomes positive. Using this $\varphi$, $2R_o(P_i/P_o)^{\frac{1}{2}} \sin (\Theta + \varphi)$ can be evaluated; for example, as shown in Fig. 6c, and if condition (26) is satisfied, (e) is in the stable locking range. This condition is generally satisfied when $P_i \ll P_o$, unless (e) is located near the boundary point of the stable free-running oscillation such as $(a')$ in Fig. 4. (Notice that $s$ is equal to 2 if the circuit is adjusted for maximum power.)



Fig. 5 — Condition for parasitic oscillations.

Fig. 6 — Graphical explanation of locking conditions.

The apparent impedance presented to the active device is given by

$$Z(\omega) - \frac{E}{I} = Z(\omega) - \frac{a_o e^{i\omega t}}{A e^{i(\omega t + \varphi)}} = Z(\omega) - 2R_o(P_i/P_o)^{\frac{1}{2}} e^{-i\varphi}$$

where $E$ and $I$ are the complex representations of $e(t)$ and $i(t)$, respectively. The last term on the right, including the minus sign, corresponds to a vector drawn from (e) to (f) in Fig. 6b. Consequently, the effect of the injection signal can be considered as that of adding to $Z(\omega)$ an impedance corresponding to this vector; the active device is operating at (f). This can be confirmed by calculating the amplitude of the oscillating current from equations (20) and (21) under steady-state conditions.

VI. BROADBAND OSCILLATOR CIRCUITS

Let us consider the frequency deviator. The oscillating frequency can be deviated by changing $\bar{X}$ of the active device, which may be a function

of supply voltage or current. For an oscillator with a single-resonant circuit, the relative frequency deviation is given by

$$\frac{\triangle f}{f} = \frac{-1}{2Q_{\text{ext}}} \frac{\triangle \bar{X}}{\bar{R}} \tag{38}$$

where $\triangle \bar{X}$ indicates the variation in $\bar{X}$. In order to increase the magnitude of the frequency deviation for a given $\triangle \bar{X}/\bar{R}$, the external $Q$ of the circuit, $Q_{\text{ext}}$ , must be reduced. However, it cannot be made smaller than the $Q$ of the active device, and in practice it tends to be larger by as much as an order of magnitude because of the difficulty in obtaining strictly lumped-constant elements at microwave frequencies.

Suppose that everything has been done to reduce $Q_{\text{ext}}$ to its practical limit; the next thing to consider is the possibility of adding another resonant circuit as shown in Fig. 7. The series arm of $L_1$ and $C_1$ represents the original resonant circuit; the parallel resonant circuit $L_2$ , $C_2$ has been added. For simplicity, let us assume

$$\omega_0 = (L_1 C_1)^{-\frac{1}{2}} = (L_2 C_2)^{-\frac{1}{2}}.$$

The locus of the impedance $Z(\omega)$, looking from the active device, should then look like Fig. 8, depending on the relative magnitude of the $Q$'s of the resonant circuits (which are defined by $Q_1 = \omega_o L_1/R_L = Q_{\text{ext}}$ , $Q_2 = \omega_o C_2 R_L$). The markers indicate equally spaced frequencies. Figure 8a shows the limiting case of $L_2 \rightarrow \infty$ and $C_2 \rightarrow 0$, which corresponds to the original single-tuned oscillator. From Section IV, it is now obvious how the oscillator behaves in each case. Thus, the conditions of Figs. 8c and d are to be avoided if smooth frequency deviation is desired. Consequently, we should concentrate on the case where $Q_2 < Q_1$ . In the vicinity of the operating point shown in Fig. 8b,

$$X(\omega) \cong R_L(Q_1 - Q_2)\frac{2 \triangle \omega}{\omega_o} \tag{39}$$

where $\triangle \omega = \omega - \omega_o$ , so that the frequency deviation corresponding to a



Fig. 7 — An oscillator with a double-resonant circuit.

Fig. 8 — $Z(\omega)$ for the oscillator shown in Fig. 7.

small $\triangle\bar{X}$ is given by

$$\frac{\triangle f}{f} \cong \frac{1}{2(Q_1 - Q_2)} \frac{\triangle\bar{X}}{\bar{R}}. \tag{40}$$

As $Q_2$ approaches $Q_1$, $\triangle f/f$ increases for a given small $\triangle\bar{X}/\bar{R}$.

However, there are several disadvantages in bringing $Q_2$ too close to $Q_1$. The first is noise: since $R'(\omega_o) = 0$ and $X'(\omega_o)$ becomes small, from equation (30) the oscillator output becomes noisy. The second is non-linearity: unless $Q_2(2\triangle\omega/\omega_o) \ll 1$, equation (39) and hence (40) become poor approximations. In fact, the frequency markers will be crowded near $\omega_o$; they rapidly get sparser away from it. This means that the relation $\triangle f$ versus $\triangle\bar{X}$ is highly nonlinear, which is objectionable in certain applications. Finally, the output power $\frac{1}{2}\bar{R}A^2$ may change rapidly as we deviate from $\omega_o$. This is because the intersection between $\bar{R} - j\bar{X}$ and $Z(\omega)$ moves to the left; hence, both $A$ and $\bar{R}$ change. If the operating point is selected to maximize the output power, then the increase in $A$ and the decrease in $\bar{R}$ may compensate for each other to produce approximately the same output power over a wide range of deviation, which is, however, highly dependent on the saturation charac-teristic of the device itself.

If the double-resonant circuit does not offer a satisfactory result be-cause of its nonlinearity or power variation, one should consider triple or even higher order resonant circuits. If all the resonant circuits in the ladder structure shown in Fig. 9 are tuned to the same frequency $\omega_o$,

and if the $Q$'s of the circuits satisfy

$$Q_1 + Q_2 + Q_4 + \cdots > Q_1 + Q_3 + \cdots > Q_2 + Q_4 + \cdots \,,$$

then a broadbanding effect can be obtained near $\omega_o$. However, the rate of return diminishes as we increase the complexity of the circuit.

Let us next consider injection locking of oscillators. It is clear from the discussion in Section IV that the circuits corresponding to Figs. 8c and d are not suitable. We, therefore, concentrate on Figs. 8a and b. The total locking range $\Delta f$ can be calculated from

$$\frac{\Delta f}{f} = \frac{2}{Q_{\text{ext}}} \left(\frac{P_i}{P_o}\right)^{\frac{1}{2}} \frac{1}{\cos \theta}$$

$$\frac{\Delta f}{f} = \frac{2}{Q_1 - Q_2} \left(\frac{P_i}{P_o}\right)^{\frac{1}{2}} \frac{1}{\cos \theta}$$

for a single-resonant and a double-resonant oscillator, respectively. As $Q_2$ approaches $Q_1$, the locking range increases for given $(P_i/P_o)^{\frac{1}{2}}$. However, the output becomes noisy if we bring $Q_2$ too close to $Q_1$. Thus, some compromise must be made.

If a precise limiter action is desired, equation (28) should be closely investigated. If good linearity of $\varphi$ versus $\Delta\omega$ is desired, multiple-resonant circuits are worth investigating. In any case, once the desired characteristic is given, it determines the necessary locus of $Z(\omega)$. Then, well-known techniques developed in connection with the design of filters are available to optimize the circuit parameters until the desired locus is sufficiently approximated.

In the above discussions, impedances were used exclusively. The same discussions can be carried out using admittances, producing no new results. However, when the active device has a large parallel capacitance, using admittances may have certain practical advantages since the capacitance can be more easily included in the multiple resonant circuit.



Fig. 9 — An oscillator with a multiple-resonant circuit.

VII. CONCLUSION

We have discussed, in detail, the expected behavior of negative resistance oscillators with multiple-resonant circuits. The systematic realization of microwave hardware which simulates these equivalent circuits still remains a problem for the future. Also, an accurate method to evaluate the device impedance $-\bar{R} + j\bar{X}$ must be developed before this discussion becomes more than just a guidance. Efforts in these directions are presently undertaken.

VIII. ACKNOWLEDGEMENT

Although they are not coauthors, N. D. Kenyon and J. P. Beccone contributed experimental investigations which were indispensable for developing the ideas presented here. Acknowledgment is due to R. S. Engelbrecht for his support, encouragement, and comments.

REFERENCES

1. Huntoon, R. D., and Weiss, A., "Synchronization of Oscillators," Proc. IRE, *35*, No. 12 (December 1947), pp. 1415–1423.
2. Hubbard, W. M., and others, "A Solid-State Regenerative Repeater for Guided Millimeter-Wave Communication Systems," B.S.T.J., *46*, No. 9 (November 1967), pp. 1977–2018.
3. Amoss, J. W., and Gsteiger, K. E., "Frequency Modulation of Avalanche Transit Time Oscillators," IEEE Trans., *MTT-15*, No. 12 (December 1967), pp. 742–747.
4. Ruthroff, C. L., "Injection-Locked-Oscillator FM Receiver Analysis," B.S.T.J., *47*, No. 8 (October 1968), pp. 1655–1661.
5. Minorsky, N., *Nonlinear Oscillators*, Princeton, N. J.: Van Nostrand, 1962.
6. Stoker, J. J., *Nonlinear Vibration*, New York: John Wiley, 1950.
7. Krylov, N., and Bogoliubov, N., *Introduction to Nonlinear Mechanics*, Princeton, N. J.: Princeton University Press, 1943.
8. Kurokawa, K., "Noise in Synchronized Oscillators," IEEE Trans., *MTT-16*, No. 4 (April 1968), pp. 234–240.
9. Kurokawa, K., "Power Waves and the Scattering Matrix," IEEE Trans., *MTT-13*, No. 2 (March 1965), pp. 194–202.

Note added in proof:

Since the submission of this article, the following book has been published. It discusses the power wave concept in detail. Also, chapter 9 is devoted to discussing oscillators with a single tuned circuit.

Kurokawa, K., *An Introduction to the Theory of Microwave Circuits*, New York: Academic Press, 1969.

# A Method for Digitally Simulating Shorted Input Diode Failures

## By HERBERT Y. CHANG

(Manuscript received September 26, 1968)

*Most existing digital fault simulators can simulate only a somewhat restrictive class of failures; namely, stuck at "1" and stuck at "0" faults. The problem seems to be the lack of techniques for properly treating the effects due to "backward propagation" of errors. This paper describes a method for illustrating how shorted input diode failures, which previously could not be simulated, can be handled by digital methods. The method is applicable to all other modes of failures describable by truth tables or Boolean expressions. Furthermore, we examine the problem of circuit oscillations caused by backward propagating errors. We conclude that failure induced oscillations can only occur under very restrictive conditions.*

## I. INTRODUCTION

Digital fault simulation is the method of predicting the behavior under failure of a logical circuit by a computer program: that is, the use of a computer to aid in computing the output(s) of a logical circuit for a given input or a given set of inputs.[1] One of the drawbacks of existing digital fault simulators is that the class of failures capable of being simulated is somewhat restrictive. Most simulators, such as IBM's Saturn Fault Simulator and Seshu's Sequential Analyzer (see Refs. 2 and 3) consider only those failures which cause some connection in the logic circuit to appear stuck at "1" (stuck-at-1 )or stuck at "0" (stuck-at-0).* Shorted input diode failures of low level logic (LLL) gates, for example, therefore cannot be simulated because they are not describable by stuck-at-1 or stuck-at-0 types of faults.

In this paper we describe a method illustrating how shorted input diode failures can be simulated digitally, using a technique somewhat

---

* Other common assumptions made are (*i*) the fault-free circuit is well behaved; (*ii*) the class of faults considered is finite and nonintermittent; and (*iii*) the single-fault assumption is used.

different from the conventional approaches. The method can be easily extended to simulate other types of failures describable by means of truth tables or Boolean expressions. The oscillatory circuit behavior resulting from "backward propagation" of errors is also analyzed.

## II. NATURE OF THE PROBLEM

The effectiveness of any simulator depends largely on how accurately the structure and the behavior of a physical model can be described by the model used in the simulator. In the case of a fault simulator, one must ascertain that the presence of any fault in the physical model can be accurately reflected in the simulator by appropriately changing the structural description of the simulated model. Experience has shown that stuck-at-1 and stuck-at-0 faults are quite common in most logic circuit realizations, and that these faults can be conveniently represented in the simulator by a technique called the failure-injection-word approach (see Ref. 3).† However, for shorted input diode failures, and possibly for many other unforeseen integrated circuit failure modes, the failure-injection-word approach will no longer be applicable. Some means for digitally simulating faults other than those of stuck-at-1 or stuck-at-0 types must therefore be devised.

Consider the three-NAND low level logic gate circuit shown in Fig. 1 and investigate, respectively, the possible malfunctions of this circuit caused by (i) shorting diode $D_1$ between input terminal $a$ and transistor Q1 and (ii) shorting diode $D_2$ between transistors Q1 and Q2.

Case (i): When $D_1$ is shorted, point $x$ will always be at the same potential as point $a$. Whenever the potential at terminal $b$ goes "low," the level at point $x$, and therefore terminal $a$, will also go "low." However, since there is no fanout at terminal $a$, other gates will not be affected. Furthermore, it can be easily verified that for all combinations of inputs at terminals $a$ and $b$, the output of Q1 remains normal; no trouble symptoms will be propagated to other parts of the circuit. Hence, we may conclude that the shorting of diode $D_1$ does not affect the operation of this circuit.

Case (ii): When diode $D_2$ is shorted, we can verify that no erroneous signals will be observed at the output terminal of Q2. However, because of the presence of fanout at point $y$, whenever the potential level at terminal $c$ goes "low," the level at point $y$ will also go low, causing a possible erroneous output at terminal $t$. Thus although the fault is

---

† Other methods of analysis can be found in Refs. 4 and 5.

Fig. 1 — A logic circuit with the NANDs.

associated with gate Q2 (from a logic viewpoint), trouble symptoms caused by the fault are not observable at the output of Q2, but are observable at the fanout gate Q3. This is the "backward propagation" of errors. In the case where there is fanout to many gates, all these gates would be affected.

At this point, we may summarize that the behavior of a logic circuit under a shorted input diode failure is inconsistent. If the input terminal of the faulty diode has no fanout, the circuit outwardly behaves in the same manner as a fault-free circuit; if the input terminal has fanout, malfunctions then propagate to other fanout gates but do not affect the gate to which this faulty diode is connected. At first glance it may appear that faults, such as shorted input diode failures, causing errors propagating backwards are difficult to analyze and therefore unsuitable for digital simulation. However, further study shows that if the nonstuck at "1" and nonstuck at "0" types of failures

are describable by truth tables, a systematic procedure can easily be devised to cope with these situations.

III. A SOLUTION

Consider a three-input NAND gate with input terminals $x$, $y$, $z$, and output terminal Q. As previously mentioned, any time an input diode is shorted, the output of a NAND behaves in the same way a fault-free NAND behaves. However, if any one of the input terminals goes "low," the terminal associated with the shorted diode will also go "low." The condition can be described by constructing a truth table for the faulty three-input NAND as shown in Table I. Part ($i$) enumerates all the input combinations; part ($ii$) shows the effect of input diode $x$ shorted on all input and output terminals. Equivalently, by realizing Table I, we may represent the condition using the following expressions:

$$\left. \begin{aligned} Q_{\text{NAND}}(x \text{ shorted}) &= (xyz)' \\ x(x \text{ shorted}) &= Q' \end{aligned} \right\}. \tag{1}$$

Equation (1) indicates that to simulate a shorted input diode fault, say at terminal $x$, of a NAND gate, one first computes the output Q as usual; one then changes the logic state of terminal $x$ to be the complement of the gate output Q. This must be done before other gates which are fanouts from terminal $x$ are simulated. In other words, simulation is accomplished in two passes (over the gate), rather than in one pass, as was the case with the failure-injection-word approach.

In order to avoid the complication resulting from backward error

TABLE I—TRUTH TABLE REPRESENTATION OF SHORTED INPUT
DIODE CONDITIONS FOR NAND GATE

| ($i$) Inputs | | | ($ii$) Diode $x$ shorted | | | |
|---|---|---|---|---|---|---|
| $x$ | $y$ | $z$ | $x$ | $y$ | $z$ | $Q$ |
| 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | ◻ | 0 | 0 | 1 |
| 1 | 0 | 1 | ◻ | 0 | 1 | 1 |
| 1 | 1 | 0 | ◻ | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 | 1 | 1 | 0 |

propagations, all gates in a logic circuit must be organized by "level," as was done in the conventional approach.[3] By "level" we mean that all primary inputs and feedback inputs (in the sense of Huffman, see Ref. 6) are of level zero; a gate is of level $n$ if all the inputs at the gate are of level $n - 1$ or less and at least one input is of level $n - 1$. To begin simulation, we apply a particular combination of primary inputs and feedback inputs to all gates of level 1; we compute the logic states of gate outputs and then all input terminals that have fanout using equation (1) to reflect the shorted input diode failure mode (of NAND gates). We proceed in a similar fashion to compute the outputs of all gates of level 2, all gates of level 3, and so on.

Notice that faults can still be simulated several at a time if, as in the failure-injection-word approach, each fault is represented by a "bit" of a computer word and these bits are properly packed. Finally, we resimulate the whole circuit once more, in a way similar to the failure-injection-word approach. The reason for resimulating the circuit is that there are gates of level $j$ whose input(s) may be fed from some gate of level $k$, and $j > k$ (see Fig. 2). When gates of level $j$ are being simulated, the states of all gates of lower levels, such as those of level $k$, have already been computed. Thus, to insure that the effect of shorted input diode fault(s) at any level properly propagates backward to gates of lower level(s), such as shown in Fig. 2, another pass must be computed over the whole circuit to obtain the correct results. With this approach, backward error propagations resulting from fanouts can be properly handled.

The technique can easily be extended to simulate the same failure for other types of gates, as long as they are describable by Boolean expressions. The truth tables describing shorted input diode failures for AND and OR gates are given in Table II.*

The equivalent Boolean expressions for Table II are

$$Q_{AND}(x \text{ shorted}) = x \cdot y$$
$$x(x \text{ shorted}) = Q_{AND} \tag{2a}$$

and

$$Q_{OR}(x \text{ shorted}) = x$$
$$y(x \text{ shorted}) = x \cdot y. \tag{2b}$$

---

* We assume the driving element is much more potent in the zero-going (positive logic) direction than in the positive-going direction, as in the case of the collector of a saturating n-p-n transistor with resistive pull-up.

Fig. 2 — Backward propagation of error resulting from shorted input diode failure.

It is apparent that the truth table or Boolean expression approach, for describing failure modes of gates for digital simulation, can also be used to represent stuck-at-1 and stuck-at-0 faults. It should also be mentioned that methods for handling certain nonstuck at "1" and nonstuck at "0" types of failures have been treated by Roth.[5] However, the technique for constructing the singular cover of the faulty gate(s) for shorted input diode types of failures is not discussed.

TABLE II — TRUTH TABLE REPRESENTATIONS OF SHORTED
INPUT DIODE CONDITIONS

| 2—Input AND | | | | |
|---|---|---|---|---|
| (*i*) Inputs | | (*ii*) Diode *x* shorted | | |
| *x* | *y* | *x* | *y* | $Q_{AND}$ |
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | ⊡ | 0 | 0 |
| 1 | 1 | 1 | 1 | 1 |



| 2—Input OR | | | | |
|---|---|---|---|---|
| (*i*) Inputs | | (*ii*) Diode *x* shorted | | |
| *x* | *y* | *x* | *y* | $Q_{OR}$ |
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | ⊡ | ⊡ |
| 1 | 0 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 | 1 |

IV. INDUCED OSCILLATIONS

In some cases the backward propagation of errors can produce oscillations in a normally well-behaved circuit. For example, consider the exclusive-OR circuit (realized in terms of NAND gates) as shown in Fig. 3. With inputs $a = 1$, $b = 1$ and the input diode from $b$ to NAND gate Q2 shorted, the error will propagate around the "loop" Q2 → $b$ → Q1 → Q2 and thus makes Q1 oscillate, which in turn causes Q4 to oscillate.† At first glance, it seems that the problem presented by induced oscillation can be extremely complicated. However, the situation is still manageable.

One of the necessary conditions under which oscillations caused by backward propagation of errors may result is that the site of failure be located at a place where reconvergent fanout paths converge. Reconvergent fanout paths are defined to be those fanout paths of some gate of level $k$ that reconverge at some gate of level $k + i$ ($i \geq 1$). This is because if there are no reconvergent fanout paths in a circuit, no "loop" can be formed when an input diode of any gate is shorted. Consequently, no oscillation is induced. Thus in circuits having no reconvergent fanouts, the two-pass technique accurately simulates the circuit behavior.

Now, suppose a circuit has some reconvergent fanout paths. Consider gate $Q_a$ of level $k$ whose fanout paths reconverge at some gate $Q_b$ of level $k + i$ (see Fig. 4). Call the number of gates traversed by a reconvergent fanout path from $Q_a$ to $Q_b$ the degree of that path. It is not difficult to see that another necessary condition for the existence of induced oscillation (caused by a shorted input diode at $Q_b$) is that there is at least one path of degree zero from $Q_a$ to $Q_b$. In other words, if every reconvergent fanout path from $Q_a$ to $Q_b$ is of degree one or greater, errors at $Q_b$ cannot propagate backwards to $Q_a$, as they will be "blocked" by gate(s) lying in paths between $Q_a$ and $Q_b$. Consequently, induced oscillation is not possible under these conditions.

Furthermore, in the case where one of the reconvergent fanout paths from $Q_a$ to $Q_b$ is of degree zero, oscillation will not result if there is no "sensitized" path(s), under the given input condition, from $Q_a$ to $Q_b$.[4,5] This is because if there are no sensitized path, errors originated at $Q_b$ and propagated backwards to $Q_a$ will not be able to propagate forward to $Q_b$. Even in cases where there is some sensitized path(s), if the number of inversions of logic signal along the sensitized path(s), from

---

† Physically, if the temporal length of the loop is short with respect to circuit operation time, oscillation will not occur.

Fig. 3 — An exclusive-or circuit.

$Q_a$ to $Q_b$ is even, no oscillation will result either, since all those inputs of $Q_b$ belonging to the reconverging fanout paths will have the same parity (as that of the particular input where the diode is shorted).

As a result, we conclude that oscillations caused by backward propagation of errors can occur, but only under a very restrictive set of conditions. In the case of shorted input diode failures, oscillation can occur if (*i*) the site of failure is located at a place where some reconvergent fanout paths converge, and is on at least one reconvergent fanout path of zero degree (see Fig. 4), and (*ii*) for the given input condition, the number of inversions of logic signal along some sensitized, reconvergent, and gain producing path(s) terminating at the site of failure is odd.

There are a number of ways to handle the problem presented by induced oscillations. One method is to first locate, from the circuit description, all groups of reconvergent fanout paths in which at least one



Fig. 4 — A circuit with a reconvergent path of degree zero.

path is of degree zero. Then, by flagging those gates that are members of these groups, we can repeatedly interrogate those flagged gates during simulation runs to determine whether there is any oscillation. As multipass simulation can be very time consuming for large circuits, it may not be economically justifiable in view of the stringent set of constraints that must be met for oscillations to occur. The other method is to modify the logic design by eliminating all zero-degree reconvergent fanout paths. This can be done by either inserting a gate in series in the zero-degree path, or by adding a gate in parallel at the site of fanout to replace the zero-degree path, or by using some other design tricks. The choice of method depends largely on the cost of its implementation and the reliability requirement of the particular application.

V. CONCLUSION

A method has been described illustrating how shorted input diode failures, which previously could not be simulated by digital fault simulators, can be handled by digital methods. It has also been shown that the method can be easily extended to simulate all other modes of failures, including stuck-at-1 and the stuck-at-0 types of faults, describable by truth tables or Boolean expressions.

Furthermore, it has been pointed out that the backward propagation of errors can sometimes produce oscillatory behavior in a normally well-behaved circuit. Failure induced oscillations, although occurring only when a rare combination of conditions exists, can nevertheless complicate the simulation process, and make the circuit behavior unpredictable. It is therefore desirable to eliminate these problems by modifying the logic circuit during the design stage.

VI. ACKNOWLEDGMENT

REFERENCES

1. Manning, E. G., and Chang, H. Y., "A Comparison of Methods for Simulating Faults of Digital Systems," Digest of the First Annual Computer Conf., Chicago, Illinois, September 1967, pp. 10–13.
2. Hardie, F., and Suhocki, R., "Design and Use of Fault Simulation for Saturn Computer Design," IEEE Trans. Elec. Computers, EC-16, No. 4 (August 1967), pp. 412–429.

3. Seshu, S., and Freeman, D. N., "The Diagnosis of Asynchronous Sequential Switching Systems," IRE Trans. Elec. Computers, *EC-11,* No. 4 (August 1962), pp. 459–465.
4. Armstrong, D. B., "On Finding a Nearly Minimal Set of Fault Detection Tests for Combinational Logic Nets," IEEE Trans. Elec. Computers, *EC-15,* No. 1 (February 1966), pp. 66–73.
5. Roth, J. P., "Diagnosis of Automata Failures: A Calculus and a Method," IBM J. Res. Development, *10* (July 1966), pp. 278–291.
6. Huffman, D. A., "Synthesis of Sequential Switching Circuits," J. Franklin Inst., *257,* Nos. 3 and 4 (March and April 1954), pp. 161–190, 275–303.

# Amplitude and Phase Modulations in Resistive Diode Mixers

## By C. DRAGONE

(Manuscript received April 30, 1968)

*This paper presents some new aspects of mixer operation and shows that the behavior of a mixer, using a resistive diode, can in general be represented by means of an equivalent circuit consisting of two transducers connected in cascade. The first of these two transducers transforms the input signals into amplitude modulations; the second is an AM detector which transforms these amplitude modulations into output signals. An important feature of this equivalent circuit is that it gives the dependence of the conversion loss upon certain important mixer parameters.*

*We also show that extremely low (< 0.3 dB) conversion losses can be achieved from a Schottky barrier diode, if the pump frequency $\omega_0$ is much smaller than the cutoff frequency of the diode. To achieve such low conversion losses the diode must be open-circuited at the harmonics $2\omega_0$, $3\omega_0$, $4\omega_0$, and so on, of the pump frequency.*

## I. INTRODUCTION AND SUMMARY OF THE PRINCIPAL RESULTS

The process of frequency conversion and its applications are well known and are extensively treated in the literature.[1-12] This paper considers the special case of a resistive diode frequency converter. We assume that the diode is pumped periodically by a strong source, the pump, which generates power at a single frequency $\omega_0$. In most of the analysis we assume that the frequency converter is required to transform small input signals occurring in the vicinity of $\omega_0$ into low-frequency output signals.

In such a frequency converter the large signals occurring at the frequencies 0, $\omega_0$, $2\omega_0$, $3\omega_0$, and so on, are perturbed by small amplitude and phase modulations caused by the input signals. The occurrence of these modulations can be shown by using the amplitude-phase representation (sometimes called AM-PM representation) which describes the small signals occurring in the vicinity of the $k$th harmonic

of $\omega_0$ in terms of the corresponding amplitude and phase modulations. When this representation is used, a convenient way of determining the small-signal terminal behavior of the frequency converter is provided by the incremental method. This method consists of perturbing the large-signals by introducing small and stationary variations in the pump level and in the dc bias. Then, from the relations between these variations and the perturbations they cause in the large signals, the small-signal terminal behavior of the frequency converter can be readily determined. In fact, we show that the connection between the foregoing relations and the desired terminal behavior can be represented by means of a simple equivalent circuit.

This equivalent circuit represents the frequency converter as a cascade connection of two transducers, as shown in Fig. 1. The first of these two transducers transforms the input signals into amplitude modulations with carrier frequency $\omega_0$ . The second is an AM detector which transforms these amplitude modulations into the desired low-frequency output signals.

In other words, when a small signal generator at some frequency $\omega_0 + p$, with $p \ll \omega_0$ , is connected to the input terminals of the frequency converter, small signals are produced in the diode at $\omega_0 + p$ and $\omega_0 - p$. These signals contain two distinct types of components, the components which are produced by amplitude modulations, and the components which are produced by phase modulations. The second transducer of Fig. 1 shows that there is a one-to-one correspondence between the output signals of the frequency converter and the amplitude modulation components. The first transducer represents the relations between these components and the input signals.

An important feature of the foregoing equivalent circuit is that it shows the dependence of the conversion loss upon certain important parameters of the frequency converter. This follows from the fact that the first transducer can be represented by means of a network consisting of two ideal transformers and two admittances. One of these two admittances is the admittance $y_{\beta 1}$ terminating the diode at the



Fig. 1 — Equivalent circuit representing a frequency converter as a cascade connection of two transducers.

Fig. 2 — First transducer of Fig. 1 for $y_{\beta 1} = 0$.

image frequency; the other is the nonlinear admittance $g_p$ presented by the diode at the pump frequency $\omega_0$. In particular, if $y_{\beta 1} = 0$ then the first transducer can be represented as shown in Fig. 2. An analogous equivalent circuit is obtained when $y_{\beta 1} = \infty$.

Finally, we examine the conversion loss of a Schottky barrier diode frequency converter. Its minimum value can be expressed as

$$L \cong 1 + 4 \exp\left(-\frac{q}{KT}\frac{\Delta V}{4}\right), \tag{1}$$

if the parasitic capacitance of the diode can be neglected. In equation (1) $\Delta V$ is the difference between the maximum and minimum values of the large-signal voltage of the diode. In order to achieve this conversion loss the diode must be open-circuited at all frequencies except $\omega = 0$, $\omega = \omega_0$, and the input and output frequencies.

Practical considerations almost always require that the input and output impedances be much smaller than the impedances required for achieving the conversion loss given by equation (1). If these practical considerations are taken into account, one finds that the frequency converter can be represented by means of the equivalent circuit shown in Fig. 3. The conversion loss of this equivalent circuit depends upon: the impedance $R$ of the input signal generator, the diode series resistance $R_s$, and the dc component $I_{c0}$ of the current flowing through the diode. In fact, the conversion loss can be expressed as

$$L = 1 + \frac{2}{R}\left(R_s + \frac{KT}{qI_{c0}}\right). \tag{2}$$

Notice that the operating frequencies at which this conversion loss



Fig. 3 — Low frequency equivalent circuit of a Schottky barrier diode frequency converter under optimum circuit conditions.

can be achieved are limited by the junction capacitance $C_j$ of the diode. Finally, notice that equation (2) implies that $R_s$ does not set, by itself, any limit to the conversion loss.

## II. SMALL-SIGNAL TERMINAL BEHAVIOR OF PUMPED NONLINEAR RESISTORS

Consider a nonlinear element in which the terminal voltage is related to the current by

$$i = f(v) \tag{3}$$

and let

$$g_d(v) = \frac{df(v)}{dv}. \tag{4}$$

Then, if the voltage $v$ is the sum of a large component $v_c(t)$ and a small component $\delta v(t)$, that is,

$$v = v(t) = v_c(t) + \delta v(t), \tag{5}$$

the current can be written to a first approximation as

$$i \cong (t) = i_c(t) + \delta i(t), \tag{6}$$

where

$$i_c(t) = f[v_c(t)] \tag{7}$$

and

$$\delta i(t) = g(t)\ \delta v(t) \tag{8}$$

with

$$g(t) = g_d[v_c(t)]. \tag{9}$$

Equation (8) completely describes the small-signal terminal behavior of a nonlinear resistor pumped by a large-signal voltage $v_c(t)$, in the absence of internal noise sources.

In equation (8) $g(t)$ is the time-varying differential conductance of the nonlinear resistor, as shown by equations (4) and (9). In the following analysis we assume that $v_c(t)$ is periodic with some frequency $\omega_0$. Therefore $i_c(t)$ and $g(t)$ also are periodic. Furthermore, we assume that the circuit connected to the nonlinear resistor is resistive at the harmonics $2\omega_0$, $3\omega_0$, $4\omega_0$, and so on, and at these frequencies, does not contain generators. Under these conditions it is always possible to choose the origin of time in such a way as to make $v_c(t)$, $i_c(t)$, and $g(t)$ even

functions of time.[4-6] Thus let

$$v_c(t) = v_c(-t), \qquad i_c(t) = i_c(-t), \qquad g(t) = g(-t). \qquad (10)$$

This allows one to write $v_c(t)$, $i_c(t)$, and $g(t)$ in the form

$$v_c(t) = V_{c0} + 2 \sum_{k=1}^{\infty} V_{ck} \cos k\omega_0 t \qquad (11)$$

$$i_c(t) = I_{c0} + 2 \sum_{k=1}^{\infty} I_{ck} \cos k\omega_0 t \qquad (12)$$

$$g(t) = g_0 + 2 \sum_{k=1}^{\infty} g_k \cos k\omega_0 t. \qquad (13)$$

III. AMPLITUDE-PHASE REPRESENTATION OF $\delta v(t)$ AND $\delta i(t)$

Normally, it is convenient to represent the perturbations $\delta v(t)$ and $\delta i(t)$ in terms of the Fourier coefficients of their frequency components. The discussion of certain properties of the behavior of a pumped diode is often simplified by the use of an alternative representation, the so-called amplitude-phase representation. The main feature of this representation is that it emphasizes the amplitude and phase modulations occurring, because of $\delta v(t)$ and $\delta i(t)$, in the various quasisinusoidal harmonic components of $v(t)$ and $i(t)$.

This representation is well known and has already been shown to be particularly useful in simplifying the analysis of certain pumped nonlinear systems.[13-15] In this section a complete derivation of it is given, in a form suitable for the discussion of the behavior of a frequency converter. We regard $\delta v(t)$ and $\delta i(t)$ as representing small perturbations produced on $v(t)$ and $i(t)$, about the condition $v(t) = v_c(t)$ and $i(t) = i_c(t)$. In accordance with this point of view, which is clarified by the following discussion, $\delta v(t)$ and $\delta i(t)$ are often referred to as perturbations.

An arbitrary voltage perturbation $\delta v(t)$ can always be expressed in the form

$$\delta v(t) = \delta v_{a0}(t) + \sqrt{2} \sum_{k=0}^{\infty} [\delta v_{ak}(t) \cos k\omega_0 t + \delta v_{pk}(t) \sin k\omega_0 t] \qquad (14)$$

where $\delta v_{a0}(t)$, $\delta v_{ak}(t)$, and $\delta v_{pk}(t)$ are low-pass functions limited to the band $|\omega| < \omega_0/2$. Therefore in equation (14) the time function

$$\delta v_k(t) = \sqrt{2}[\delta v_{ak}(t) \cos k\omega_0 t + \delta v_{pk}(t) \sin k\omega_0 t] \qquad (15)$$

represents the components of $\delta v(t)$ occurring in the frequency range between $k\omega_0 - \omega_0/2$ and $k\omega_0 + \omega_0/2$ and between $-k\omega_0 - \omega_0/2$ and $-k\omega_0 + \omega_0/2$. In particular, if $\delta v(t)$ contains frequency components at only the side-frequencies $r\omega_0 \pm p(|\ r\ | = 0, 1, 2$, and so on; $0 \leqq p < \omega_0/2)$, then $\delta v_k(t)$ and $\delta v_{a0}(t)$ can be written as

$$\delta v_k(t) = 2(\text{Re})\{V_{\alpha k} \exp\ [j(p + k\omega_0)t]$$

$$+ V_{\beta k} \exp\ [j(p - k\omega_0)t]\}\qquad (k = 1, 2, \text{and so on}) \qquad (16)$$

$$\delta v_{a0}(t) = 2(\text{Re})(V_{\alpha 0}e^{jpt}) \qquad (17)$$

where $V_{\alpha 0}$, $V_{\alpha k}$ and $V_{\beta k}$ are the complex amplitudes of the components of $\delta v(t)$ occurring at $p$, $p + k\omega_0$, and $p - k\omega_0$, respectively ($k = 1$, $2, 3, \cdots$). Equation (16) can be rewritten in the form

$$\delta v_k(t) = 2(\text{Re})[(V_{\alpha k} + V_{\beta k})e^{jpt}\ \cos\ k\omega_0 t$$

$$+ j(V_{\alpha k} - V_{\beta k})e^{jpt}\ \sin\ k\omega_0 t]. \qquad (18)$$

A comparison of this expression and equation (15) shows that the low-pass functions $\delta v_{ak}(t)$ and $\delta v_{pk}(t)$ are sinusoidal of frequency $p$. That is, if equation (16) is satisfied one can write

$$\delta v_{ak}(t) = 2(\text{Re})(V_{ak}e^{jpt}) \qquad (k = 0, 1, 2, \cdots) \qquad (19)$$

$$\delta v_{pk}(t) = 2(\text{Re})(V_{pk}e^{jpt}) \qquad (k = 1, 2, 3, \cdots). \qquad (20)$$

Furthermore, by comparing equation (18) and the expressions that one obtains by substituting equation (19) and (20) into equation (15), one obtains the following relations between the Fourier coefficients of $\delta v_{ak}(t)$ and $\delta v_{pk}(t)$ and those of $\delta v(t)$

$$\begin{bmatrix} V_{ak} \\ V_{pk} \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ j & -j \end{bmatrix} \begin{bmatrix} V_{\alpha k} \\ V_{\beta k} \end{bmatrix} \qquad (k = 1, 2, 3, \cdots). \qquad (21)$$

Furthermore,

$$V_{a0} = V_{\alpha 0}. \qquad (22)$$

The two sets of low-pass functions $\delta v_{a0}(t)$, $\delta v_{a1}(t)$, and so on, and $\delta v_{p1}(t)$, $\delta v_{p2}(t)$, and so on, defined by equation (14) provide a complete representation of the voltage perturbation $\delta v(t)$. Their significance can be readily obtained by noticing that if $v(t)$ is written as

$$v(t) = V_0(t) + \sqrt{2} \sum_{k=1}^{\infty} V_k(t)\ \cos\ [k\omega_0 t + \delta\varphi_k(t)], \qquad (23)$$

where $V_0(t)$, $V_k(t)$ and $\delta\varphi_k(t)$ are low-pass functions limited to the band $|\omega| < \omega_0/2$, from equations (5), (11), and (14) one obtains

$$V_0(t) = V_{c0} + \delta v_{a0}(t) \tag{24}$$

$$V_k(t) = \sqrt{2}\, V_{ck} + \delta v_{ak}(t) \qquad (k > 0) \tag{25}$$

$$\delta\varphi_k(t) = \frac{-1}{\sqrt{2}\, V_{ck}}\, \delta v_{pk}(t) \qquad (k > 0). \tag{26}$$

That is, if one decomposes the voltage $v(t)$ into a sum of harmonic terms, one finds that, because of the presence of the perturbation $\delta v(t)$, each harmonic component is modulated. The $k$th harmonic component is amplitude modulated by $\delta v_{ak}(t)$ and phase modulated by $\delta v_{pk}(t)$. Because of this property the representation of $\delta v(t)$ in terms of the low-pass functions $\delta v_{ak}(t)$ and $\delta v_{pk}(t)$ (or of their Fourier coefficients $V_{ak}$ and $V_{pk}$) is called the amplitude-phase representation of $\delta v(t)$. The direct representation of $\delta v(t)$ by means of the Fourier coefficients $V_{ak}$ and $V_{\beta k}$ is called the $\alpha$-$\beta$ representation.[13]

The foregoing discussion can be extended to $\delta i(t)$ by replacing $v$ with $i$ throughout. That is, $\delta i(t)$ can be represented in terms of low-pass functions $\delta i_{ak}(t)$ and $\delta i_{pk}(t)$

$$\delta i(t) = \delta i_{a0}(t) + \sqrt{2} \sum_{k=1}^{\infty} [\delta i_{ak}(t) \cos k\omega_0 t + \delta i_{pk}(t) \sin k\omega_0 t]. \tag{27}$$

Furthermore, in the special case where $\delta i(t)$ of the type

$$\delta i(t) = 2(\mathrm{Re})\Bigg\{ \sum_{k=0}^{\infty} I_{\alpha k} \exp\,[j(p + k\omega_0)t]$$
$$+ \sum_{k=1}^{\infty} I_{\beta k} \exp\,[j(p - k\omega_0)t] \Bigg\}, \tag{28}$$

the low-pass functions $\delta i_{ak}(t)$ and $\delta i_{pk}(t)$ are sinusoidal of frequency $p$ and, if $I_{ak}$ and $I_{pk}$ are their Fourier coefficients, one has

$$\begin{bmatrix} I_{ak} \\ I_{pk} \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ j & -j \end{bmatrix} \begin{bmatrix} I_{\alpha k} \\ I_{\beta k} \end{bmatrix} \qquad (k = 1, 2, 3, \cdots) \tag{29}$$

and

$$I_{a0} = I_{\alpha0}. \tag{30}$$

Equations (21) and (29) can be interpreted, because of their particular form, as the relations describing the terminal behavior of a four-terminal-pairs lossless network consisting of ideal transformers, such as the network $T$ shown in Fig. 4.

Fig. 4 — Network $T$ describing the relations between the $\alpha - \beta$ signals at $p \pm k\omega_0$ and the corresponding amplitude and phase modulations.

IV. SMALL-SIGNAL TERMINAL BEHAVIOR OF THE DIODE IN THE AMPLITUDE-PHASE REPRESENTATION

According to equations (14) and (27)

$$\delta v(t) = \delta v_a(t) + \delta v_p(t) \tag{31}$$

$$\delta i(t) = \delta i_a(t) + \delta i_p(t) \tag{32}$$

where $\delta v_a(t)$ and $\delta i_a(t)$ are the amplitude modulation components of $\delta v(t)$ and $\delta i(t)$, and $\delta v_p(t)$ and $\delta i_p(t)$ are the phase modulation components. Thus

$$\delta v_a(t) = \delta v_{a0}(t) + \sqrt{2} \sum_{k=1}^{\infty} \delta v_{ak}(t) \cos k\omega_0 t \tag{33}$$

$$\delta v_p(t) = \sqrt{2} \sum_{k=1}^{\infty} \delta v_{pk}(t) \sin k\omega_0 t \tag{34}$$

$$\delta i_a(t) = \delta i_{a0}(t) + \sqrt{2} \sum_{k=1}^{\infty} \delta i_{ak}(t) \cos k\omega_0 t \tag{35}$$

$$\delta i_p(t) = \sqrt{2} \sum_{k=1}^{\infty} \delta i_{pk}(t) \sin k\omega_0 t. \tag{36}$$

Since from equations (13), (35), and (33) one has that $\delta v_a(t) g(t)$ only contains amplitude modulation components and $\delta v_p(t) g(t)$ only contains phase modulation components, by substituting equations (31) and (32) into equation (8) one obtains

$$\delta i_a(t) = g(t) \, \delta v_a(t) \tag{37}$$

$$\delta i_p(t) = g(t) \, \delta v_p(t). \tag{38}$$

The relations between the low-pass functions $\delta i_{a0}(t)$, $\delta i_{a1}(t)$, and so

on and $\delta v_{a0}(t)$, $\delta v_{a1}(t)$, $\delta v_{a2}$, ... may now be obtained by substituting equations (13), (33), and (35) into equation (37). One obtains

$$\delta i_{a0}(t) + \sqrt{2} \sum_{k=1}^{\infty} \delta i_{ak}(t) \cos k\omega_0 t = g_0 \, \delta v_{a0}(t)$$

$$+ \sqrt{2} \sum_{r=1}^{\infty} g_r \, \delta v_{ar}(t) + \sum_{k=1}^{\infty} \left[ 2g_k \, \delta v_{a0}(t) \right.$$

$$\left. + \sqrt{2} \sum_{r=1}^{\infty} (g_{r+k} + g_{|r-k|}) \, \delta v_{ar}(t) \right] \cos k\omega_0 t, \qquad (39)$$

which gives

$$
\begin{bmatrix} \delta i_{a0}(t) \\ \delta i_{a1}(t) \\ \delta i_{a2}(t) \\ \delta i_{a3}(t) \\ \vdots \end{bmatrix} = \begin{bmatrix} g_0, & \sqrt{2}\, g_1, & \sqrt{2}\, g_2, & \sqrt{2}\, g_3, & \cdots \\ \sqrt{2}\, g_1, & g_0 + g_2, & g_1 + g_3, & g_2 + g_4, & \cdots \\ \sqrt{2}\, g_2, & g_1 + g_3, & g_0 + g_4, & g_1 + g_5, & \cdots \\ \sqrt{2}\, g_3, & g_2 + g_4, & g_1 + g_5, & g_0 + g_6, & \cdots \\ \vdots & \vdots & \vdots & \vdots & \end{bmatrix} \begin{bmatrix} \delta v_{a0}(t) \\ \delta v_{a1}(t) \\ \delta v_{a2}(t) \\ \delta v_{a3}(t) \\ \vdots \end{bmatrix} \cdot
$$

$$(40)$$

In a completely similar way, from equations (13), (34), (36) and (38) one obtains

$$
\begin{bmatrix} \delta i_{p1}(t) \\ \delta i_{p2}(t) \\ \delta i_{p3}(t) \\ \vdots \end{bmatrix} = \begin{bmatrix} g_0 - g_2, & g_1 - g_3, & g_2 - g_4, & \cdots \\ g_1 - g_3, & g_0 - g_4, & g_1 - g_5, & \cdots \\ g_2 - g_4, & g_1 - g_5, & g_0 - g_6, & \cdots \\ \vdots & \vdots & \vdots & \end{bmatrix} \begin{bmatrix} \delta v_{p1}(t) \\ \delta v_{p2}(t) \\ \delta v_{p3}(t) \\ \vdots \end{bmatrix} \cdot \qquad (41)
$$

Equations (40) and (41) provide the amplitude-phase admittance-matrix representation of the small-signal terminal behavior of the diode, in the absence of internal noise sources. If the various low-pass functions $\delta v_{ak}(t)$ and $\delta v_{pk}(t)$ are sinusoidal of frequency $p$, then equations (40) and (41) give

$$I_a] = [G_a] V_a]$$

$$I_p] = [G_p] V_p] \qquad (42)$$

where the matrix notation is defined as

$$[G_a] = \begin{bmatrix} g_0, & \sqrt{2}\,g_1, & \sqrt{2}\,g_2, & \sqrt{2}\,g_3, & \cdots \\ \sqrt{2}\,g_1, & g_0 + g_2, & g_1 + g_3, & g_2 + g_4, & \cdots \\ \sqrt{2}\,g_2, & g_1 + g_3, & g_0 + g_4, & g_1 + g_5, & \cdots \\ \sqrt{2}\,g_3, & g_2 + g_4, & g_1 + g_5, & g_0 + g_6, & \cdots \\ \vdots & \vdots & \vdots & \vdots & \end{bmatrix} \qquad (43)$$

$$[G_p] = \begin{bmatrix} g_0 - g_2, & g_1 - g_3, & g_2 - g_4, & \cdots \\ g_1 - g_3, & g_0 - g_4, & g_1 - g_5, & \cdots \\ g_2 - g_4, & g_1 - g_5, & g_0 - g_6, & \cdots \\ \vdots & \vdots & \vdots & \end{bmatrix} \qquad (44)$$

and

$$I_a] = \begin{bmatrix} I_{a0} \\ I_{a1} \\ \vdots \end{bmatrix}, \qquad V_a] = \begin{bmatrix} V_{a0} \\ V_{a1} \\ \vdots \end{bmatrix}, \qquad I_p] = \begin{bmatrix} I_{p1} \\ I_{p2} \\ \vdots \end{bmatrix}, \qquad V_p] = \begin{bmatrix} V_{p1} \\ V_{p2} \\ \vdots \end{bmatrix}. \qquad (45)$$

Notice that $I_{ak}$, $V_{ak}$, $I_{pr}$, and $V_{pr}$ are the Fourier coefficients of $\delta i_{ak}(t)$, $\delta v_{ak}(t)$, $\delta i_{pr}(t)$, and $\delta v_{pr}(t)$, as shown by equations (19) and (20).

Equation (42) shows that the terminal behavior of the diode has the following special property: the diode does not produce any coupling between the amplitude modulation components and the phase modulation components of the perturbations $\delta v(t)$ and $\delta i(t)$; or equivalently, the diode does not produce amplitude $\rightleftarrows$ phase conversion. This property can be useful in simplifying the discussion of the terminal behavior of the diode because it allows the two types of signals (amplitude and phase) to be treated separately.

According to equation (42) the small-signal terminal behavior of the diode can be represented by two separate linear and time-invariant networks $D_a$ and $D_p$, as shown in Fig. 5. In these two equivalent networks the terminal voltages and currents occur at the same frequency $(p)$ and their Fourier coefficients are equal to those of the various low-pass functions of $\delta v(t)$ and $\delta i(t)$. Notice that from equations (43) and (44) one has that $[G_a]$ and $[G_p]$ are real symmetric matrices and

Fig. 5 — Time-invariant equivalent networks $D_a$ and $D_p$ providing the amplitude-phase representation of the terminal behavior of a resistive pumped diode.

therefore the two networks $D_a$ and $D_p$ can be realized by means of ordinary resistive networks.

## V. TERMINAL BEHAVIOR OF A FREQUENCY CONVERTER IN THE AMPLITUDE-PHASE REPRESENTATION

In the following part of this paper the results obtained in the preceding sections are applied to the study of a very common type of frequency converter, a frequency converter which is required to transform input signals occurring at some frequency $\omega_0 + p$, with $p \ll \omega_0$, into output signals occurring at $\omega = p$. We assume that $\delta v(t)$ and $\delta i(t)$ contain components at only the pairs of side-frequencies $k\omega_0 + p$ and $k\omega_0 - p(|k| = 0, 1, 2,$ and so on; $2p < \omega_0)$.

In this section attention is focused on the signals occurring at $p$ and $p \pm \omega_0$. At set of relations among these signals is derived, from the equations describing the small-signal terminal behavior of the diode, by taking into account the constraints imposed by the external circuit on the signals occurring at the side-frequencies of $2\omega_0$, $3\omega_0$, $4\omega_0$, and so on. We show that these relations establish the existence of a one-to-one correspondence between the output signals ($V_{a0}$ and $I_{a0}$) and the amplitude modulation coefficients ($V_{a1}$ and $I_{a1}$) of the signals occurring

at $p \pm \omega_0$. That is, they do not impose any relation between the output signals and the phase modulation coefficients ($V_{p1}$ and $I_{p1}$) of the signals occurring at $p \pm \omega_0$.

Let $y_k$, $y_{\alpha k}$ and $y_{\beta k}$ ($k > 0$) denote the admittances terminating the diode at $k\omega_0$, $p + k\omega_0$, and $p - k\omega_0$, respectively. We assume that

$$y_k = y_{\alpha k} = y_{\beta k}, \qquad k > 1, \qquad (46)$$

a condition satisfied in many practical applications. Notice that this condition is satisfied for $p \to 0$ because $y_k$ is real.

Consider the equivalent network $T$ shown in Fig. 4. One can verify that if one terminates the two terminal pairs relative to the ($\alpha$, $\beta$) signals with two admittances equal to $y_k$ then its behavior at the remaining two terminal pairs is described by the relations

$$I_{\alpha k} = -y_k V_{\alpha k}, \qquad I_{pk} = -y_k V_{pk} \qquad (k > 1), \qquad (47)$$

which show that, if condition (46) is satisfied, then the external circuit does not produce amplitude $\rightleftarrows$ phase conversion. That is, if $y_{\alpha k} = y_{\beta k}$, then in the amplitude-phase representation the behavior of the external circuit at $p \pm k\omega_0$ can be represented by means of two separate one-terminal-pair networks, Furthermore, the admittances of these two networks are equal to the admittance presented by the external circuit at $p \pm k\omega_0$.

Now, consider the relations among the signals occurring at $p \pm \omega_0$ and $p$ when the remaining signals are constrained by the conditions imposed by the external circuit at $p \pm 2\omega_0$, $p \pm 3\omega_0$, and so on. If the two networks $D_a$ and $D_p$ shown in Fig. 5 are terminated as specified by equations (47) one obtains the two networks $A$ and $P$ shown in Figs. 6 and 7. The relations between the terminal voltages and currents of these two networks can be written in the form

$$I_{p1} = g_p V_{p1} \qquad (48)$$

$$\begin{bmatrix} I_{a0} \\ I_{a1} \end{bmatrix} = [G'_a] \begin{bmatrix} V_{a0} \\ V_{a1} \end{bmatrix} \qquad (49)$$

where $g_p$ and $[G'_a]$ can be derived from equations (42) and (47). For example, in the special case where $y_k = \infty$ for $k > 1$, from equations (42) and (47) one obtains

$$[G'_a] = \begin{bmatrix} g_0 & \sqrt{2}\, g_1 \\ \sqrt{2}\, g_1 & g_0 + g_2 \end{bmatrix}, \qquad g_p = g_0 - g_2. \qquad (50)$$

Fig. 6 — Network $A$ representing the relations in a frequency converter between the output signals and the amplitude modulation signals at $p \pm \omega_0$.



Fig. 7 — Network $P$ representing the behavior of the diode at $p \pm \omega_0$ with respect to phase modulation.

Notice that equation (49) establishes the existence, in a frequency converter satisfying equation (47), of a one-to-one correspondence between the output signals $V_{a0}$ and $I_{a0}$ and the amplitude modulation signals $V_{a1}$ and $I_{a1}$.

The significance of the two networks $A$ and $P$ is best understood by considering the following alternative method of deriving $g_p$ and $[G'_a]$. Suppose, for the moment, that $\delta v(t) = \delta i(t) = 0$ and that the level of the pump and the dc bias applied to the diode are variable. That is, assume that the amplitudes of the various harmonic components of $v(t)$ and $i(t)$ can be varied, by changing the level of the pump and the dc bias, while the terminations $y_2$, $y_3$, $y_4$, $\cdots$ presented by the external circuit to the diode are kept constant. The functions $v_c(t)$ and $i_c(t)$ must satisfy the constraint $i_c(t) = f[v_c(t)]$ imposed by the diode. Furthermore, the harmonic components of $v_c(t)$ and $i_c(t)$ must satisfy the constraints $I_{ck} = -y_k V_{ck}$ $(k > 1)$ imposed by the external circuit. Clearly, from these constraints two independent nonlinear relations can be derived, among the four variables $V_{c0}$, $V_{c1}$, $I_{c0}$ and $I_{c1}$. That is, if one considers the variables $\eta_0 = I_{c0}$, $\eta_1 = \sqrt{2}I_{c1}$, $x_0 = V_{c0}$ and $x_1 = \sqrt{2}V_{c1}$, one can write

$$\eta_0 = F_0(x_0, x_1)$$
$$\eta_1 = F_1(x_0, x_1)$$
(51)

where the two nonlinear functions $F_0(x_0, x_1)$ and $F_1(x_0, x_1)$ are determined by the particular function $f(v)$ describing the diode nonlinearity and by the parameters $y_2$, $y_3$, $y_4$, $\cdots$. Notice that the variables $\eta_0$, $\eta_1$, $x_0$ and $x_1$ represent the root mean square values of the harmonic components of order zero and one of $i(t)$ and $v(t)$.

Now, suppose $\delta v(t)$ and $\delta i(t)$ differ from zero and are produced by small generators occurring in the external circuit at the frequencies $p \pm \omega_0$ and $p$. We want to derive from equation (51) the relations among the coefficients $V_{a1}$, $I_{a1}$, $V_{a0}$, $I_{a0}$, $V_{p1}$ and $I_{p1}$ of $\delta v(t)$ and $\delta i(t)$. Since the terminal behavior of the two networks $A$ and $P$ (shown in Figs. 6 and 7) is frequency independent, the relations in question are independent of $p$ and therefore one can set $p = 0$ without loss of generality. Thus, let $\delta v(t)$ and $\delta i(t)$ be produced by small-signal generators occurring at $\omega = \omega_0$ and $\omega = 0$. These small-signal generators can be interpreted as follows. A small-signal amplitude modulation generator occurring at $\omega_0$ can be regarded as representing a small change of the level of the pump; a small-signal phase modulation generator occurring at $\omega_0$ can be interpreted as a small change $\delta\varphi$ of the phase of the pump.

Finally, a small-signal generator occurring at dc can be interpreted as a small change produced in the dc bias circuit. When small changes of the above three types occur in the external circuit, the components of $\delta v(t)$ and $\delta i(t)$ at $\omega = 0$ and $\omega = \omega_0$ can be evaluated as follows.

First, consider the amplitude modulation components $\delta v_a(t)$ and $\delta i_a(t)$ of $\delta v(t)$ and $\delta i(t)$. Clearly, the harmonic components of order zero and one of $v_c(t) + v_a(t)$ and $i_c(t) + i_a(t)$ must satisfy equation (51). From equations (11), (12), (33), and (35) and from the fact that $\delta v_{ak}(t) = 2V_{ak}$ and $\delta i_{ak}(t) = 2I_{ak}$ because $p = 0$, the root mean square values of the harmonic components in question are $x_0 = V_{c0} + 2V_{a0}$, $x_1 = \sqrt{2}V_{c1} + 2V_{a1}$, $\eta_0 = I_{c0} + 2I_{a0}$ and $\eta_1 = \sqrt{2}I_{c1} + 2I_{a1}$. Therefore, from equation (51),

$$I_{c0} + 2I_{a0} = F_0(\sqrt{2}\ V_{c1} + 2V_{a1}\ ,\ V_{c0} + 2V_{a0}) \qquad (52)$$

$$\sqrt{2}\ I_{c1} + 2I_{a1} = F_1(\sqrt{2}\ V_{c1} + 2V_{a1}\ ,\ V_{c0} + 2V_{a0}),$$

and since $V_{ak} \ll V_{ck}$, from equations (49) and (52) one obtains

$$[G'_a] = \begin{bmatrix} \left(\dfrac{\partial F_0}{\partial x_0}\right)_c & \left(\dfrac{\partial F_0}{\partial x_1}\right)_c \\[2mm] \left(\dfrac{\partial F_1}{\partial x_0}\right)_c & \left(\dfrac{\partial F_1}{\partial x_1}\right)_c \end{bmatrix}, \qquad (53)$$

where $(\ )_c$ indicates that the partial derivatives are calculated for

$$x_0 = V_{c0}\ , \qquad x_1 = \sqrt{2}V_{c1}\ . \qquad (54)$$

Next, consider the phase modulation components $\delta v_p(t)$ and $\delta i_p(t)$ of $\delta v(t)$ and $\delta i(t)$. They are produced by a small change $\delta\varphi$ of the phase of the pump. Such a change is equivalent to a small change $\delta\varphi/\omega_0$ of the origin of time. Therefore

$$v_c(t) + \delta v_p(t) = v_c\left(t + \frac{\delta\varphi}{\omega_0}\right) \qquad (55)$$

$$i_c(t) + \delta i_p(t) = i_c\left(t + \frac{\delta\varphi}{\omega_0}\right). \qquad (56)$$

From these relations and from equations (11), (12), (34), and (36) one can calculate the phase modulation coefficients $V_{p1}$ and $I_{p1}$ of $\delta v_p(t)$ and $\delta i_p(t)$. One obtains

$$V_{p1} = -\delta\varphi V_{c1}/\sqrt{2}, \qquad I_{p1} = -\delta\varphi I_{c1}/\sqrt{2}.$$

Therefore, from these relations and equation (48),

$$g_p = \frac{I_{c1}}{V_{c1}} = \left[ \frac{F_1(x_0, x_1)}{x_1} \right]_c . \tag{57}$$

In words, $g_p$ is the nonlinear conductance presented at $\omega_0$ by the diode to the external circuit under normal operating conditions.

The foregoing discussion presents two alternative methods of evaluating the terminal behavior of the two networks $A$ and $P$. One method consists of analyzing these two networks by using equations (42) and (47), as suggested by the equivalent circuits of Figs. 6 and 7. Therefore, it requires that the small-signal terminal behavior of the diode be determined. The other method consists of deriving $[G_a]$ and $g_p$ directly from the large-signal terminal behavior, at $\omega = 0$ and $\omega = \omega_0$, of the nonlinear network consisting of the diode terminated with $y_2, y_3, y_4, \cdots$, at $2\omega_0, 3\omega_0, 4\omega_0, \cdots$. If these admittances are non-zero and finite, considerable analytical difficulties often arise when one tries to analyze the equivalent circuits $A$ and $P$ by using the representation shown in Figs. 6 and 7. In these cases the second method may be used. Furthermore, it is important to point out that a well known technique of measuring frequency converters is based on the second method, which is often called the incremental method.[4-6]

In the special case where $y_k = \infty$ for $k > 1$, one can readily verify that

$$F_0(x_0, x_1) = \frac{1}{2\pi} \int_0^{2\pi} f(x_0 + \sqrt{2}\, x_1 \cos \omega_0 t)\, dt,$$

and

$$F_1(x_0, x_1) = \frac{1}{\sqrt{2}\,\pi} \int_0^{2\pi} f(x_0 + \sqrt{2}\, x_1 \cos \omega_0 t)\, \cos \omega_0 t\, dt, \tag{58}$$

and by using equations (53) and (57) one obtains equation (50).

VI. COMPLETE EQUIVALENT CIRCUIT OF A FREQUENCY CONVERTER

In Section V a one-to-one correspondence between the output signals of a frequency converter and the amplitude modulation coefficients $V_{a1}$ and $I_{a1}$ was derived. In this section the terminal behavior of the frequency converter is completely determined by deriving the relations between these coefficients and the input signals. The relations in question can be obtained with the help of the equivalent circuits T and P shown in Figs. 4 and 7. In fact, they are given by the terminal behavior of the network shown in Fig. 8.

Fig. 8 — Network $T'$ describing the transformation of the input signals of a frequency converter into amplitude modulation signals at $p \pm \omega_0$.

Now, by connecting the two networks of Figs. 6 and 8 in cascade, as shown in Fig. 9, one obtains a two-terminal-pairs network, which provides a complete representation of the terminal behavior of the frequency converter. This network corresponds to the equivalent circuit shown in Fig. 1. Notice that in Fig. 9 the input terminals of the frequency converter are connected to a small-signal generator with short-circuit terminal current $I_s$ and with internal admittance $y_{a1}$, and that the output terminals are connected to a load $y_{a0}$.

From equation (15) the signals occurring at the terminals of the diode at $p \pm \omega_0$ can be expressed as

$$\delta v_1(t) = \sqrt{2} \ \delta v_{a1}(t) \ \cos \omega_0 t + \sqrt{2} \ \delta v_{p1}(t) \sin \omega_0 t$$

$$\delta i_1(t) = \sqrt{2} \ \delta i_{a1}(t) \ \cos \omega_0 t + \sqrt{2} \ \delta i_{p1}(t) \sin \omega_0 t. \tag{59}$$



Fig. 9 — Equivalent circuit representing a frequency converter as a cascade connection of two transducers.

From these relations one can verify that

$$P_1 = P_{a1} + P_{p1},\tag{60}$$

where

$$P_1 = \langle \delta v_1(t)\ \delta i_1(t)\rangle_{av}\tag{61}$$

$$P_{a1} = 2\langle \delta v_{a1}(t)\ \delta i_{a1}(t)\ \cos^2 \omega_0 t\rangle_{av} = \langle \delta v_{a1}(t)\ \delta i_{a1}(t)\rangle_{av}$$

$$P_{p1} = 2\langle \delta v_{p1}(t)\ \delta i_{p1}(t)\ \sin^2 \omega_0 t\rangle_{av} = \langle \delta v_{p1}(t)\ \delta i_{p1}(t)\rangle_{av}\tag{62}$$

and $\langle\ \rangle_{av}$ indicates the time average. Equation (59) shows that $\delta v_1(t)$ and $\delta i_1(t)$ contain two types of signals: the amplitude modulation components $\sqrt{2}\ \delta v_{a1}(t)\ \cos\ \omega_0 t$ and $\sqrt{2}\ \delta i_{a1}(t)\ \cos\ \omega_0 t$, and the phase modulation components $\sqrt{2}\ \delta v_{p1}(t)\ \sin\ \omega_0 t$ and $\sqrt{2}\ \delta i_{p1}(t)\ \sin\ \omega_0 t$. Equations (59) through (62) show that the total power $P_1$ flowing into the diode at $p \pm \omega_0$ is the sum of the individual powers $P_{a1}$ and $P_{p1}$ carried by these two types of components. Furthermore, $P_{a1}$ and $P_{p1}$ can be calculated directly from the modulating functions $\delta v_{a1}(t)$ and $\delta i_{a1}(t)$ as shown by equations (62).

Now, consider the equivalent circuit shown in Fig. 9. The network $T'$ can be regarded as a transducer which transforms the input signals of the frequency converter into amplitude modulations of the fundamental harmonic components of $v_c(t)$ and $i_c(t)$. These amplitude modulations are then transformed by the network $A$ into the output signals of the frequency converter. The network $T'$ is dissipative because it contains the two admittances $y_{\beta 1}$ and $g_p$. The power dissipated in $g_p$, $P_{p1}$, represents the power dissipated in the frequency converter because of the generation of phase modulation signals. The power flowing into the network $A$, at its left terminals, is $P_{a1}$. Figure 9 clearly shows that $P_{a1} + P_{p1}$, the total power flowing into the diode at $p \pm \omega_0$, is in general less than the power $P_{a1}$ flowing into the input terminals of the frequency converter. In fact the difference $P_{a1} - (P_{a1} + P_{p1})$ between these two powers is lost in the admittance $y_{\beta 1}$ of Fig. 9. This is the power flowing from the diode into the external circuit at $p - \omega_0$. In general, only if either one of the two conditions

$$y_{\beta 1} = \infty\tag{63}$$

$$y_{\beta 1} = 0\tag{64}$$

is satisfied, is this power equal to zero. These two conditions are examined in Section VIII.

In the following section the special and very important case of a

frequency converter in which

$$y_{\beta 1} = y_{\alpha 1} \tag{65}$$

is considered. The discussion of the terminal behavior of this frequency converter will also provide the background needed in Section VIII for the discussion of the two cases given by equations (63) and (64).

## VII. FREQUENCY CONVERTER WITH EQUAL TERMINATIONS AT $\omega_0 + p$ AND $p - \omega_0$

One can verify that if $y_{\beta 1} = y_{\alpha 1}$, then the Thevenin representation of the one-terminal-pair network connected to the left terminals of the network $A$ of Fig. 9 has a short-circuit current $I_s/\sqrt{2}$ and an internal admittance $y_{\alpha 1}$, as shown in Fig. 10a. One can also verify that the one-terminal-pair network connected to the terminals of $g_p$, in Fig. 9, has the Thevenin representation shown in Fig. 10b. Notice that the only significant difference between the two equivalent generators of Fig. 10 and the generator connected to the input terminals of the frequency converter (Fig. 9) is that the available power of this generator is twice the power available from each of the two equivalent generators of Fig. 10. Notice, furthermore, that Fig. 10 shows that if $y_{\alpha 1} = y_{\beta 1}$, in the amplitude-phase representation the terminal behavior of the external circuit at $p \pm \omega_0$ can be represented by means of two separate one-terminal-pair networks, the two equivalent generators shown in Fig. 10. This property is in accordance with the remarks in Section VI about the significance of Equation (47).

Thus, if $y_{\alpha 1} = y_{\beta 1}$, then only half of the power available from the input generator shown in Fig. 9 is used by the frequency converter to



(a)                                              (b)

Fig. 10 — Equivalent generators in the amplitude-phase representation for the external circuit at $p \pm \omega_0$, when $y_{\beta 1} = y_{\alpha 1}$.

produce output signals. The remaining half is lost because of the generation of phase modulation signals, as shown by Fig. 10b. This explains why the minimum conversion loss of a frequency converter of the type considered here is always greater than two and can be expressed as

$$L = 2L_a , \tag{66}$$

where $L_a$ is the minimum conversion loss of the network $A$ shown in Fig. 10a. The conversion loss given by equation (66) is achieved when

$$y_{\alpha 1} = Y_{a1} , \qquad y_{\alpha 0} = Y_{a0} , \tag{67}$$

where $Y_{ak}$ $(k = 0, 1)$ are the image admittances of the network $A$.

A frequency converter of the type considered in this section is characterized by the following special behavior. If one minimizes its conversion loss, the frequency converter reflects, at its input terminals, a part of the power incident from the input generator. The reflected power $P_\rho$ can be calculated with the help of the equivalent circuit shown in Fig. 8, by terminating with $Y_{a1}$ the terminals relative to the signals $V_{a1}$ and $I_{a1}$, connecting the remaining terminals to the input generator shown in Fig. 9, and setting $y_{\alpha 1} = y_{\beta 1} = Y_{a1}$. One finds that $P_\rho$ is equal to the power dissipated in $y_{\beta 1}$ and

$$P_\rho = \frac{1}{4}\left(\frac{g_p - y_{\alpha 1}}{g_p + y_{\alpha 1}}\right)^2 P_s = \frac{1}{4}\left(\frac{g_p - Y_{a1}}{g_p + Y_{a1}}\right)^2 P_s \tag{68}$$

where $P_s$ is the available power of the input generator. Therefore, the following conclusion can be drawn: in a frequency converter of the type considered in this section, minimum conversion loss is achieved only when the power reflected at the input terminals of the frequency converter is equal to the power dissipated in the admittance terminating the diode at $p - \omega_0$ and is given by equation (68).

This result can be given the following interpretation. The available power of the equivalent generator shown in Fig. 10a represents the amplitude modulation power available from the external circuit at $p \pm \omega_0$ . This power is entirely absorbed by the diode, when equations (67) are satisfied. Only a part of the power available from the generator of Fig. 10b, on the other hand, is generally absorbed by the diode, because normally $g_p \neq Y_{a1} = y_{\alpha 1}$ . From Fig. 10b, the power reflected at the terminals of the generator in question is

$$\frac{P_s}{2}\left(\frac{g_p - y_{\alpha 1}}{g_p + y_{\alpha 1}}\right)^2 = \frac{P_s}{2}\left(\frac{g_p - Y_{a1}}{g_p + Y_{a1}}\right)^2 .$$

Since half of this power is reflected at $p + \omega_0$ and the remaining half is reflected at $p - \omega_0$, one obtains equation (68) and the conclusion following it.

Notice that, in the special case where $y_k = \infty$ for $k > 1$, from equations (50) one obtains the well known expressions:[4]

$$L_a = \frac{1 - [1 - 2g_1^2/g_0(g_0 + g_2)]}{1 + [1 - 2g_1^2/g_0(g_0 + g_2)]^{\frac{1}{2}}} \tag{69}$$

$$Y_{a1} = \frac{g_0 + g_2}{g_0}, \qquad Y_{a0} = g_0 + g_2\left[1 - \frac{2g_1^2}{g_0(g_0 + g_2)}\right]^{\frac{1}{2}}. \tag{70}$$

## VIII. FREQUENCY CONVERTER WITH EITHER $y_{\beta 1} = \infty$ OR $y_{\beta 1} = 0$

One can readily verify that in the two cases corresponding to conditions (64) and (63), the equivalent circuit of Fig. 9 reduces to the two equivalent circuits shown in Figs. 11 and 12, respectively. In both cases the minimum conversion loss can be expressed as

$$L = 4L_a\{[(L_a + 1)^2 + \gamma(L_a^2 - 1)]^{\frac{1}{2}}$$
$$- [(L_a - 1)^2 + \gamma(L_a^2 - 1)]^{\frac{1}{2}}\}^{-2} \tag{71}$$

where

$$\gamma = \frac{g_p}{Y_{a1}} \qquad \text{if} \quad y_{\beta 1} = \infty \tag{72}$$

$$\gamma = \frac{Y_{a1}}{g_p} \qquad \text{if} \quad Y_{\beta 1} = 0. \tag{73}$$

The values of $y_{a1}$ and $y_{a0}$ required for achieving this conversion loss are given in the appendix.

Of great practical importance is the problem of determining which



Fig. 11 — Equivalent circuit of a frequency converter when $y_{\beta 1} = 0$.

Fig. 12 — Equivalent circuit of a frequency converter when $y_{\beta 1} = \infty$.

of the two conditions $y_{\beta 1} = \infty$ and $y_{\beta 1} = 0$ yields the lowest conversion loss. Since equation (71) shows that $L$ decreases with decreasing $\gamma$, for $\gamma \geqq 0$, one has that of the two foregoing conditions the one which gives $\gamma < 1$ yields the lowest $L$. Therefore one obtains the fundamental result: Of the two conditions $y_{\beta 1} = 0$ and $y_{\beta 1} = \infty$ the one which yields the lowest conversion loss is $y_{\beta 1} = 0$ if $g_p/Y_{a1} > 1$ and is $y_{\beta 1} = \infty$ if $g_p/Y_{a1} < 1$; the two conditions are equivalent if $g_p = Y_{a1}$.

Therefore, if $y_{\beta 1}$ is so chosen as to minimize the conversion loss, then $\gamma \leqq 1$ and from equation (71) one obtains

$$L_a \leqq L \leqq \frac{1}{L_a - (L_a^2 - 1)^{\frac{1}{2}}} \tag{74}$$

where the two equal signs occur when $\gamma = 0$ and $\gamma = 1$, respectively. The first inequality is a direct consequence of the fact that in both equivalent circuits of Figs. 11 and 12 the conversion loss is always greater than $L_a$, the minimum conversion loss of the network $A$. The second inequality shows that in a frequency converter of the type considered here the optimum conversion loss is always less than $2L_a$, the lowet conversion loss obtainable when $y_{a1} = y_{\beta 1}$. In fact

$$1 \leqq \frac{1}{L_a - (L_a^2 - 1)^{\frac{1}{2}}} \leqq 2L_a \tag{75}$$

where the equal signs occur when $L_a = 1$ and $L_a = \infty$.

Often practical considerations require that the admittance $y_{a1}$ of the input generator be equal to a prescribed value $G$, not necessarily equal to the value required for optimum performance. In such a case it is important to point out that one can readily determine, with the help of the two equivalent circuits of Figs. 11 and 12, which one of the two conditions $y_{\beta 1} = 0$ and $y_{\beta 1} = \infty$ is to be chosen, in order to minimize the conversion loss. Clearly the choice depends upon the values of $G$, $Y_{a1}$ and $L_a$.

IX. SCHOTTKY BARRIER DIODE

Among the various resistive diodes presently available, the Schottky barrier diode offers the most suitable electrical characteristics for microwave frequency conversion.[9,16-18] The main features of this diode are that its low-frequency terminal behavior is described by the equation

$$i = i_s \left\{ \exp \left[ \frac{q}{KT} (v - iR_s) \right] - 1 \right\} \tag{76}$$

over a very wide range of voltages, and that its frequency response is not limited by minority carrier lifetime. In equation (76) $i_s$ is the saturation current, $q$ the electronic charge, $K$ the Boltzmann constant, $T$ the absolute temperature, and $R_s$ the series resistance. Therefore,

$$\frac{q}{KT} \cong 40 \quad \text{for} \quad T = 290°\text{K}. \tag{77}$$

The useful frequency range of the diode is limited primarily by the junction capacitance $C_j$. In some instances the lead inductance and the case capacitance also have to be considered. However, in Section 9.1 consideration is restricted to the range of frequencies over which the foregoing parasitics can be neglected.

9.1 *Analysis*

Assume that

$$i = i_s \exp \left[ \frac{q}{KT} (v - R_s i) \right] \tag{78}$$

and that

$$y_k = 0 \quad \text{for} \quad k > 1. \tag{79}$$

Then from equation (78) one obtains

$$\frac{dv}{di} = \frac{KT}{q} \frac{1}{i} + R_s \tag{80}$$

and, from equations (12) and (79),

$$i_c(t) = I_{c0} + 2I_{c1} \cos \omega_0 t. \tag{81}$$

Therefore, by setting $i = i_c(t)$ in equation (80), the differential resistance of the diode can be expressed as

$$r(t) = R_s + \frac{KT}{2qI_{c1}} \frac{1}{\xi + \cos \omega_0 t} \tag{82}$$

where

$$\xi = \frac{I_{c0}}{2I_{c1}}. \tag{83}$$

From equation (82) one can readily calculate the Fourier coefficients of $r(t)$. One obtains

$$r_n = \delta(n)R_s + \frac{KT}{2qI_{c1}} \frac{1}{\pi} \int_0^\pi \frac{\cos n\theta}{\xi + \cos \theta}\, d\theta = \delta(n)R_s + \frac{KT}{2qI_{c1}} (-1)^n$$

$$\cdot \left(\frac{\xi + 1}{\xi - 1}\right)^{\frac{1}{2}} (1 + \xi)^{n-1} \left[\frac{\xi}{\xi + 1} - \left(\frac{\xi - 1}{\xi + 1}\right)^{\frac{1}{2}}\right]^n \tag{84}$$

where

$$\delta(n) = \begin{cases} 1, & n = 0 \\ 0, & n \neq 0 \end{cases}. \tag{85}$$

Because of the formal analogy between the impedance-matrix representation and the admittance-matrix representation, from equation (50) one has that the impedance matrix $[R'_a]$ of the network $A$ is

$$[R'_a] = \begin{bmatrix} r_0 & \sqrt{2}\, r_1 \\ \sqrt{2}\, r_1 & r_0 + r_2 \end{bmatrix} \tag{86}$$

and the impedance $r_p$ presented at $\omega_0$ by the diode to the pump is

$$r_p = r_0 - r_2. \tag{87}$$

Therefore, from equations (84), (86), and (87) one obtains

$$[R'_a] = \frac{KT}{2qI_{c1}} \begin{bmatrix} (\xi^2 - 1)^{-\frac{1}{2}} & -\sqrt{2}\, [\xi(\xi^2 - 1)^{-\frac{1}{2}} - 1] \\ -\sqrt{2}\, [\xi(\xi^2 - 1)^{-\frac{1}{2}} - 1] & 2\xi[\xi(\xi^2 - 1)^{-\frac{1}{2}} - 1] \end{bmatrix}$$

$$+ \begin{bmatrix} R_s & 0 \\ 0 & R_s \end{bmatrix} \tag{88}$$

and

$$r_p = \frac{KT}{qI_{c1}} [\xi - (\xi^2 - 1)^{\frac{1}{2}}]. \tag{89}$$

Equation (88) shows that the network $A$ can be represented by means of the equivalent circuit shown in Fig. 13, with

$$R_1 = \frac{KT}{qI_{c0}} \tag{90}$$

Fig. 13 — Network $A$ for a Schottky barrier diode frequency converter with $y_k = 0$ for $k > 1$.

$$R_2 = \frac{KT}{qI_{c0}} \frac{\xi - (\xi^2 - 1)^{\frac{1}{2}}}{(\xi^2 - 1)^{\frac{1}{2}}} \qquad (91)$$

$$n = \sqrt{2}\, \xi. \qquad (92)$$

Now, let $v_M$ and $v_m$ denote the maximum and minimum values, respectively, of $v_c(t) - R_s i_c(t)$ and let

$$\triangle V = v_M - v_m. \qquad (93)$$

Then from equations (78), (81), and (93) one obtains

$$I_{c0} + 2I_{c1} = (I_{c0} - 2I_{c1}) \exp\left(\frac{q}{KT}\, \triangle V\right)$$

which, by making use of equation (83), gives

$$\xi = 1 + 2\left[\exp\left(\frac{q}{KT}\, \triangle V\right) - 1\right]^{-1}. \qquad (94)$$

At this point the assumption is made that

$$\triangle V > 0.5 \text{ volts,} \qquad (95)$$

a condition satisfied in most cases of practical interest. Then, if one substitutes equation (94) into equations (89), (91), (92) and examines the behavior of $n$, $R_2$ and $r_p$ for large $q\,\triangle V/KT$, one finds that

$$n \cong \sqrt{2} \qquad (96)$$

$$r_p \cong 2\frac{KT}{qI_{c0}} + R_s \qquad (97)$$

$$R_2 \cong 2\frac{KT}{qI_{c0}} \exp\left(\frac{q}{KT}\, \frac{\triangle V}{2}\right). \qquad (98)$$

Furthermore, one can verify that the fractional errors in these approximate expressions decrease exponentially with $\triangle V$ and are less

than 0.01 per cent, because of condition (95). Therefore the equivalent circuit of Fig. 13 can be replaced by that shown in Fig. 14.

Notice that

$$R_s , R_1 \ll R_2 , \tag{99}$$

because of condition (95) and of the fact that typically $R_s < 10\Omega$ and

$$I_{c0} \ll 100 \text{ mA.} \tag{100}$$

Then, with the help of the equivalent circuit of Fig. 14 and by making use of condition (95), from standard network theory one obtains

$$Z_{a1} \cong Z_{a0} \cong \left(2 + \frac{2R_s q I_{c0}}{KT}\right)^{\frac{1}{2}} \frac{KT}{q I_{c0}} \exp\left(\frac{q}{KT} \frac{\triangle V}{4}\right) \tag{101}$$

$$L_a \cong 1 + 2\left(2 + \frac{3R_s q I_{c0}}{KT}\right)^{\frac{1}{2}} \exp\left(- \frac{q}{KT} \frac{\triangle V}{4}\right) \tag{102}$$

where $Z_{a0}$ and $Z_{a1}$ are the image impedances of the network $A$ (that is, $Z_{a0} = 1/Y_{a0}$ and $Z_{a1} = 1/Y_{a1}$).

Equations (97) and (101) show that $r_p \ll Z_{a1}$. Therefore, in accordance with the discussion in Section VIII, condition (64) yields better performance than condition (63). Thus let $y_{\beta 1} = 0$. Then, from the equivalent circuits shown in Figs. 11 and 14 and from equation (97) one obtains the equivalent circuit shown in Fig. 15.

9.2 *Discussion of the Equivalent Circuit*

From the equivalent circuit shown in Fig. 15 one obtains the following expression for the minimum conversion loss of a Schottky barrier diode

$$L \cong 1 + 4\left(1 + \frac{R_s q I_{c0}}{KT}\right)^{\frac{1}{2}} \exp\left(- \frac{q}{KT} \frac{\triangle V}{4}\right) \tag{103}$$

which is achieved when

$$z_{\alpha 1} = z_{\alpha 0} \cong \left(1 + \frac{R_s q I_{c0}}{KT}\right)^{\frac{1}{2}} \frac{KT}{q I_{c0}} \exp\left(\frac{q}{KT} \frac{\triangle V}{4}\right) \tag{104}$$

where $z_{\alpha k} = 1/y_{\alpha k}$. From equation (103), by selecting $I_{c0}$ in such a way that

$$\frac{R_s q I_{c0}}{KT} \ll 1, \tag{105}$$

one obtains equation (1).

Fig. 14 — Equivalent circuit of Fig. 13 when $\triangle V > 0.5$ volts.

Notice that it is always desirable that $\triangle V$ be as large as possible, be-cause this increases the value of the resistance $R_2$ shown by Fig. 15. On the other hand, equations (104) and (105) show that it $\triangle V$ is too large, then very high input and output impedances are required in order to achieve minimum conversion loss. For instance, if $\triangle V > 0.85$ (which is a condition that can be readily satisfied in most practical cases since the breakdown voltage of the diode typically is greater than 1 volt), $I_{c0} = 10$ mA and $R_s = 2\Omega$, then from equations (103) and (104) one obtains

$$L < 1.0015, \quad \text{but} \quad z_{\alpha 1} = z_{\alpha 0} > 17000\,\Omega. \tag{106}$$

Since normally it is required that the input and output impedances be much smaller than this value, one concludes that practical considera-tions almost always require that

$$z_{\alpha 1} = z_{\alpha 0} \ll \frac{KT}{qI_{c0}} \exp\left(\frac{q}{KT}\frac{\triangle V}{4}\right). \tag{107}$$

Because of this restriction, the equivalent circuit of Fig. 15 reduces to that of Fig. 16 which shows that, under practical conditions, the termi-nal behavior of the frequency converter only depends upon the two parameters $I_{c0}$ and $R_s$ . Furthermore, Fig. 16 clearly shows that low conversion losses require high dc currents $(I_{c0})$; that is, they require



Fig. 15 — Equivalent circuit of a Schottky barrier diode frequency converter with $y_k = 0$ for $k > 1$, $y_{\beta 1} = 0$, and $\triangle V > 0.5$ volts.

Fig. 16 — Equivalent circuit of a Schottky barrier diode frequency converter under certain optimum circuit conditions.

that the impedance $r_p$ presented by the diode at $\omega_0$ be small. Notice that from the equivalent circuit of Fig. 16 one readily obtains equation (2).

The pump power $P_p$ absorbed by the diode at $\omega_0$ is

$$P_p = 2I_{c1}^2 r_p \ . \tag{108}$$

Therefore, by using equations (77) and (97) and the fact that $I_{c0} \cong 2I_{c1}$ one obtains

$$P_p \cong \frac{I_{c0}}{40} + \frac{R_s I_{c0}^2}{2000} \qquad \text{(mW)}, \tag{109}$$

where $R_s$ and $I_{c0}$ are measured in ohms and milliamperes, respectively.

9.3 *Optimum Terminations at $2\omega_0$ , $3\omega_0$ , $4\omega_0$ , $\cdots$*

The preceding analysis assumes that the diode is open-circuited at the frequencies $p \pm k\omega_0$ and $k\omega_0$ , $k > 1$. In this section we show that the conversion losses obtainable when the diode is short-circuited at these frequencies are appreciably higher than those given by equation (1).

Assume that the diode is short-circuited at $p \pm k\omega_0$ and $k\omega_0 (k > 1)$. Reference 10 shows that then the optimum termination at $p - \omega_0$ is[*]

$$y_{\beta 1} = 0. \tag{110}$$

With this optimum termination the minimum conversion loss can be expressed as follows[4-10]

$$L' = \left\{ 1 + \left[ \frac{1 + \dfrac{g_2}{g_0} - 2\left(\dfrac{g_1}{g_0}\right)^2}{\left[1 - \left(\dfrac{g_1}{g_0}\right)^2\right]\left(1 + \dfrac{g_2}{g_0}\right)} \right]^{\frac{1}{2}} \right\}^2 \frac{\left[1 - \left(\dfrac{g_1}{g_0}\right)^2\right]\left(1 + \dfrac{g_2}{g_0}\right)}{\left(\dfrac{g_1}{g_0}\right)^2\left(1 - \dfrac{g_2}{g_0}\right)} , \tag{111}$$

[*] The analysis of Ref. 10 is valid only if the diode is short-circuited at $p \pm k\omega_0$, $k > 1$. Therefore, even though the analysis covers all cases where such an assumption can be made, it cannot be applied to the case $(y_k = y_{\alpha k} = y_{\beta k} = 0)$.

where

$$\frac{g_2}{g_0} = \frac{J_2\left(\frac{q}{KT}\frac{\triangle V}{2}\right)}{J_0\left(\frac{q}{KT}\frac{\triangle V}{2}\right)}, \tag{112}$$

$$\frac{g_1}{g_0} = \frac{J_1\left(\frac{q}{KT}\frac{\triangle V}{2}\right)}{J_0\left(\frac{q}{KT}\frac{\triangle V}{2}\right)}, \tag{113}$$

provided $R_s$ can be neglected.[10] From these equations, by making use of well known asymptotic expansions of the modified Bessel functions $J_0$, $J_1$, and $J_2$, one obtains

$$L' = 1 + 2\left(\frac{KT}{q\,\triangle V}\right)^{\frac{1}{2}} + 2\left(\frac{KT}{q\,\triangle V}\right) + (5 + \sqrt{2})\left(\frac{KT}{q\,\triangle V}\right)^{\frac{3}{2}} + \cdots .$$

$$\tag{114}$$

Thus, if one compares $L'$ with the conversion loss $L$ of equation (1) one finds

$$\frac{L' - 1}{L - 1} \to \infty \quad \text{as} \quad \frac{KT}{q\,\triangle V} \to 0.$$

If for example $\triangle V = 1$, then equations (1) and (114) give

$$L = 1.00018 \tag{115}$$

and

$$L' = 1.39, \tag{116}$$

respectively. Clearly, condition (79) gives better performance than the condition $Y_k = \infty$ for $k > 1$.

## X. CONCLUSIONS

It has been shown that, under certain circuit conditions, a resistive diode frequency converter can be represented by the equivalent circuit shown in Fig. 9. The general case where the diode cannot be regarded as a nonlinear resistor and the terminations $y_2$, $y_3$, $y_4$, $\cdots$ are not resistive has not been considered, in order to prevent the mathematical difficulties inherent in it from obscuring the significance of the results.

However, the equivalent circuit of Fig. 9 is valid under conditions much more general than those considered in the preceding sections.

In fact, in this section we show that it is valid also in the general case of a frequency converter using an arbitrary nonlinear element which is pumped at $\omega_0$ and is arbitrarily terminated at $2\omega_0$, $3\omega_0$, $4\omega_0$, and so on, provided that three conditions are satisfied. ($i$) $p \ll \omega_0$, ($ii$) small and slow changes of the level of the pump and of the dc bias do not cause phase variations in the fundamental harmonic components of $v_c(t)$ and $i_c(t)$, and ($iii$) the nonlinear admittance ($g_p$) seen by the pump be real.

Assume first that $iii$ is satisfied. Then, the fundamental harmonic components of $v_c(t)$ and $i_c(t)$ have the same phase $\varphi_1$. Let the origin of time be chosen in such a way that $\varphi_1 = 0$. If one now defines the amplitude modulation and phase modulation components of $\delta v_1(t)$ and $\delta i_1(t)$ by means equation (59), the relations between $V_{a1}$, $I_{a1}$, $V_{p1}$, $I_{p1}$ and $V_{\alpha 1}$, $I_{\alpha 1}$, $V_{\beta 1}$, $I_{\beta 1}$ are provided by the terminal behavior of the network $T$ shown in Fig. 4.

Next, assume that $p \ll \omega_0$ and let the incremental method be used to derive the relations between the output signals $V_{a0}$ and $I_{a0}$ and the coefficients $V_{\alpha 1}$, $I_{\alpha 1}$, $V_{p1}$ and $I_{p1}$. By using equations (51), which are valid even if the aforementioned three conditions are not satisfied, one finds that there is a one-to-one correspondence between the output signals of the frequency converter and the amplitude modulation coefficients $V_{\alpha 1}$ and $I_{\alpha 1}$. More precisely,

$$
\begin{bmatrix} I_{a0} \\ I_{a1} \end{bmatrix} = \begin{bmatrix} \left(\dfrac{\partial F_0}{\partial x_0}\right)_c & \left(\dfrac{\partial F_0}{\partial x_1}\right)_c \\ \left(\dfrac{\partial F_1}{\partial x_0}\right)_c & \left(\dfrac{\partial F_1}{\partial x_1}\right)_c \end{bmatrix} \cdot \begin{bmatrix} V_{a0} \\ V_{a1} \end{bmatrix} \tag{117}
$$

where $F_0(x_0, x_1)$ and $F_1(x_0, x_1)$ have the significance given in Section V.

Finally, assume that all of the three conditions are satisfied. Then, by using the incremental method one finds that $I_{p1}/V_{p1}$ is independent of $I_{a1}$, $V_{a1}$, $I_{a0}$ and $V_{a0}$, because of $ii$. The ratio $I_{p1}/V_{p1}$ can therefore be determined by applying a small variation $\delta\varphi$ to the phase of the pump; and one finds

$$
I_{p1} = g_p V_{p1}. \tag{118}
$$

Now, from equations (117) and (118) and from the network $T$ of Fig. 4 one obtains the equivalent circuit of Fig. 9, where the terminal behavior of the network $A$ is now specified by equation (117).

Notice that the validity of the equivalent circuit shown in Fig. 9

can be extended to the general case where only the condition $p \ll \omega_0$ is satisfied, by properly modifying the first transducer. This follows from the fact that $p \ll \omega_0$ is sufficient to guarantee the validity of equation (117).

The special case of a Schottky varrier diode is examined in Section IX, which shows that the optimum terminations at the image frequency and at the harmonics $2\omega_0$, $3\omega_0$, $4\omega_0$, $\cdots$, are open-circuits. Furthermore, at low frequency and under optimum circuit conditions, the terminal behavior of a Schottky barrier diode frequency converter is quite simple. It can be adequately represented by means of the equivalent circuit shown in Fig. 16.

Finally, the fact that consideration has been restricted to the frequency range where the Schottky barrier diode can be regarded as a nonlinear resistor should not be interpreted as an indication that $C_j$ can generally be neglected. However, the simple results obtained for the low frequency case serves as a guide for treating the high frequency case, which is reserved to a future article.

APPENDIX

The conversion loss given by equation (71) is obtained when

$$\hat{y}_{a1} = Y_{a1}\left(1 + \gamma \frac{L_a + 1}{L_a - 1}\right)^s \left(1 + \gamma \frac{L_a - 1}{L_a + 1}\right)^s \tag{119}$$

$$\hat{y}_{a0} = 2Y_{a0}\left(1 + \gamma \frac{L_a + 1}{L_a - 1}\right)^s \left(1 + \gamma \frac{L_a - 1}{L_a + 1}\right)^s \tag{120}$$

where

$$\gamma = g_p/Y_{a1}, \quad s = 1 \quad \text{if} \quad y_{a1} = \infty$$

$$\gamma = Y_{a1}/g_p, \quad s = -1 \quad \text{if} \quad y_{\beta 1} = 0.$$

REFERENCES

1. Strutt, M. J. O., "Mixing Values," Wireless Eng., *12*, 1935, pp. 59–64.
2. Peterson, E., and Hussey, L. W., "Equivalent Modulator Circuits," B.S.T.J., *18*, No. 1 (January 1939), pp. 32–48.
3. Caruthers, R. S., "Copper Oxide Modulators in Carrier Telephone Systems," B.S.T.J., *18*, No. 2 (April 1939), pp. 315–337.
4. Peterson, L. C., and Llewellyn, F. B., "The Performance of Mixers in Terms of Linear-Network Theory," Proc. IRE, *33*, No. 7 (July 1945), pp. 458–476.
5. Torrey, H. C., and Whitmer, C. A., Crystal Rectifiers, *15*, Rad. Lab. Series, New York: McGraw-Hill, 1948.
6. Pound, R. V., Microwave Mixers, *16*, Rad. Lab. Series, New York: McGraw-Hill, 1948.

7. Southworth, G. C., *Principles and Applications of Waveguide Transmission,* Princeton, N. J.: D. van Nostrand, 1950, pp. 626–636.
8. Messenger, G. C., and McCoy, C. T., "Theory and Operation of Crystal Diodes as Mixers," Proc. IRE, *45,* No. 9 (September 1957), pp. 1269–1283.
9. Barber, M. R., and Ryder, R. M., "Ultimate Noise Figure and Conversion Loss of the Schottky Barrier Mixer Diode," Int. Microwave Symp. Digest, (May 1966), pp. 13–17.
10. Barber, M. R., "Noise Figure and Conversion Loss of the Schottky Barrier Mixer Diode," IRE Trans., *MTT-15,* No. 11 (November 1967), pp. 629–635.
11. Rafuse, R. P., "Low Noise and Dynamic Range in Symmetric Mixer Circuits," Conf. on High Frequency Generation and Amplification, Cornell University, Ithaca, New York, August 1967.
12. Rafuse, R. P., "Symmetric MOSFET Mixers of High Dynamic Range," 1968 Int. Solid-State Circuits Conf., Philadelphia, Pennsylvania, February 1968.
13. Penfield, P., Jr., "Circuit Theory of Periodically Driven Nonlinear Systems," Proc. IEEE, *54,* No. 2 (February 1966), pp. 266–280.
14. Dragone, C., "Phase and Amplitude Modulation in High Efficiency Varactor Frequency Multipliers—General Scattering Properties," B.S.T.J., *46,* No. 4 (April 1967), pp. 775–796.
15. Dragone, C., and Prabhu, V. K., "Scattering Relations in Lossless Varactor Frequency Multipliers," B.S.T.J., *46,* No. 8 (October 1967), pp. 1699–1731.
16. Kahng, D., and D'Asaro, L. A., "Gold-Epitaxial Silicon High-Frequency Diodes," B.S.T.J., *43,* No. 1 (January 1964), pp. 225–232.
17. Kahng, D., "Conduction Properties of the Au-n-Type-Si Schottky Barrier," Solid State Elec., *6,* No. 3 (May–June 1963), pp. 281–295.
18. Sze, S. M., and Ryder, R. M., "The Nonlinearity of the Reverse Current-Voltage Characteristics of a p-n Junction Near Avalanche Breakdown," B.S.T.J., *16,* No. 6 (July–August), pp. 1135–1139.

# On Conditions Under Which It Is Possible To Synchronize Digital Transmission Systems

## By I. W. SANDBERG

*J. R. Pierce has recently proposed a system for synchronizing an arbitrary number of geographically separated oscillators, and, under the assumption of zero transmission delays between stations, has shown that a certain linear model of the system is stable in the sense that all of the station frequencies approach a common final value as $t \to \infty$.*

*This paper reports on some results concerning the dynamic behavior of Pierce's linear model. The results take into account transmission delays. More explicitly, we prove that if a certain set of simple inequalities involving the delays is satisfied, then the system is stable and the oscillator frequencies approach their common final value at an exponential rate. These inequalities have the property that they are always satisfied for sufficiently small delays. This paper presents a simple example showing that the system can be unstable when the inequalities are not met. In addition, we present some information concerning the rate of decay of the natural modes of stable systems, discuss an alternative stability criterion not involving the transmission delays and derive an explicit expression for the final frequency. Finally, we discuss the mathematical relationship between Pierce's model and earlier models of synchronization systems.*

## I. INTRODUCTION, DISCUSSION, AND SUMMARY OF RESULTS

### 1.1 *The Model*

It is well known that the problem of synchronizing the frequencies of geographically separated oscillators is of importance in connection with the detection and switching of pulse-code-modulated signals.

J. R. Pierce has recently proposed a system for synchronizing digital

1999

transmission networks* The system uses oscillators of adjustable rate of change of frequency, buffers which accept pulses at an incoming rate and which produce corresponding output pulses at the local clock rate, and adequate delay at each station to enable PCM frames to be properly aligned in time.

In Pierce's model the content $b_{ij}$ of the buffer at station $i$ which accepts pulses from station $j$ is assumed to satisfy the equation,[†]

$$\dot{b}_{ij} = f_j - f_i \tag{1}$$

in which $f_j$ and $f_i$ are the frequencies at stations $j$ and $i$, respectively, and the overall system of coupled oscillators is assumed to be governed by the equations

$$\dot{f}_i = c_i(f_{0i} - f_i) + \sum_{j \neq i} a_{ij} b_{ij} \qquad i = 1, 2, \cdots, n \tag{2}$$

for $t \geq 0$, where $n$ is the number of stations, each $c_i$ is a positive constant, $f_{0i}$ is the center frequency of the oscillator at station $i$, and each $a_{ij}$ is nonnegative and satisfies $a_{ij} = a_{ji}$.

In his discussion of equations (2), Pierce assumes that the transmission delay between stations $i$ and $j$ can be neglected for all $i \neq j$, in which case equation (1) can be interpreted as

$$\dot{b}_{ij}(t) = f_j(t) - f_i(t), \qquad t \geq 0 \tag{3}$$

instead of as

$$\dot{b}_{ij}(t) = f_j(t - \tau_{ij}) - f_i(t), \qquad t \geq 0 \tag{4}$$

with each $\tau_{ij}$ a nonnegative constant. Under the natural assumption that there is some path from each station to every other station, Pierce has shown, by directing attention to a passive RLC network analog of the equations,[‡] that the system is stable in the

---

sense that each frequency $f_i$ approaches the same final frequency as $t \to \infty$.[*]

## 1.2 *Example of An Unstable System*

We wish to show now that if transmission delays are taken into account, then the system described above need not be stable. Consider the equations of a two-station system with $c_1 = c_2 = c$, $a_{12} = a_{21} = a$, and $\tau_{12} = \tau_{21} = \bar{\tau}(c, a, \bar{\tau} > 0)$:

$$\dot{f}_1 + cf_1 + a \int_0^t [f_1(\tau) - f_2(\tau - \bar{\tau})] \, d\tau = ab_{12}(0) + cf_{01}$$

$$\dot{f}_2 + cf_2 + a \int_0^t [f_2(\tau) - f_1(\tau - \bar{\tau})] \, d\tau = ab_{21}(0) + cf_{02}$$

for $t \geq 0$. It is not difficult to verify that these equations posses a solution

$$f_1(t) = \mathrm{Re} \, [v_1 e^{i\omega t}] \qquad\qquad f_2(t) = \mathrm{Re} \, [v_2 e^{i\omega t}]$$

for $t \geq -\bar{\tau}$ with $v_1$, $v_2$, and $\omega \neq 0$ constants if and only if

$$-a \, \mathrm{Re} \, [(v_1 - v_2 e^{-i\omega\bar{\tau}})(i\omega)^{-1}] = ab_{12}(0) + cf_{01} \qquad (5)$$

$$-a \, \mathrm{Re} \, [(v_2 - v_1 e^{-i\omega\bar{\tau}})(i\omega)^{-1}] = ab_{21}(0) + cf_{02} \qquad (6)$$

and $Mv = \theta$, in which $v$ is the transpose of $(v_1, v_2)$, $\theta$ is the zero 2-vector, and

$$M = \begin{bmatrix} -\omega^2 + i\omega c + a & -ae^{-i\omega\bar{\tau}} \\ -ae^{-i\omega\bar{\tau}} & -\omega^2 + i\omega c + a \end{bmatrix}.$$

Thus if $\det M = 0$ for some real value of $\omega \neq 0$, then the system is unstable in the sense that for some values of the right sides of equations (5) and (6) there exists a pair of real-valued functions $f_1(\cdot)$ and $f_2(\cdot)$ such that the equations are satisfied for all $t \geq 0$, and at least one of these functions does not approach a limit as $t \to \infty$. Moreover, if for example $\omega^2 = a$, $c^2 = a$, and $\omega\bar{\tau} = \frac{1}{2}\pi$, then $\det M = 0$, which shows that the two-station system can be unstable if additional restrictions are not imposed on $c$, $a$, and $\bar{\tau}$.[†]

---

[*] He has also exploited the network analogy further in order to obtain an expression for the final frequency and to make assertions concerning the behavior of the system when certain elements are nonlinear.

[†] Note that $f_1(t)$ or $f_2(t)$ can be negative in this example. While it is certainly true that instantaneous frequencies are not negative, our analysis is intended to show only that the solution of the equations can possess a sinusoidal mode. Thus, the conclusion and the essential details of the analysis are unchanged if we add a constant $u$ to $f_1(t)$ and $f_2(t)$ provided that we subtract the constant $cu$ from the right side of equation (5) and the right side of equation (6).

### 1.3 Summary of Results

The purpose of this paper is to report on some results concerning the dynamic behavior of Pierce's model. The results take into account transmission delays. In particular, we prove that the (assumed connected) system is stable whenever

$$c_i > \sum_{j \ne i} a_{ij} \tfrac{1}{2}(\tau_{ij} + \tau_{ji}) \qquad i = 1, 2, \cdots, n. \tag{7}$$

In fact we prove that if condition (7) is satisfied, then the frequencies $f_i$ approach their final value $\rho$ at an exponential rate.* Concerning the final frequency $\rho$, we derive the explicit expression

$$\rho = \frac{\sum_i c_i f_{0i} + \sum_i \sum_{j \ne i} a_{ij} b_{ij}(0) + \sum_i \sum_{j \ne i} a_{ij} \int_{-\tau_{ij}}^{0} f_i(u)\, du}{\sum_i c_i + \sum_i \sum_{j \ne i} a_{ij} \tau_{ij}}. \tag{8}$$

This paper also presents some material concerning bounds on the rate of decay of transients in stable systems. More explicitly, we derive a lower bound on the rate of decay of all complex natural modes.

We prove also that the system is stable whenever

$$c_i \geqq (2 \sum_{j \ne i} a_{ij})^{\frac{1}{2}} \qquad i = 1, 2, \cdots, n \tag{9}$$

and that if condition (9) is satisfied, then [as in the case of condition (7)] the frequencies $f_i$ approach $\rho$ at an exponential rate.

Unlike condition (7), condition (9) is obviously not always satisfied for sufficiently small delays. On the other hand, if condition (9) is satisfied, then the system is stable independent of the values of the $\tau_{ij}$. In this sense the results corresponding to conditions (7) and (9) are complementary.

Our results described above are stated in a more precise manner in Section II, and proofs are given in Section III.

### 1.4 The Relation to Earlier Work

The system governed by equations (2) and (4) is closely related to synchronization systems which have been studied earlier. This relation is made clear as follows. From equations (2) and (4)

---

* Of course in the vast majority of cases it is reasonable to assume that $\tau_{ij} = \tau_{ji}$ for all $i \ne j$. We have proceeded without this assumption in order to show that our stability result is not critically dependent on it.

$$\dot{f}_i = c_i(f_{0i} - f_i) + \sum_{j \neq i} a_{ij} \int_0^t [f_i(\tau - \tau_{ij}) - f_i(\tau)] \, d\tau$$

$$+ \sum_{j \neq i} a_{ij} b_{ij}(0) \qquad i = 1, 2, \cdots, n \qquad (10)$$

for all $t \geqq 0$. We write equations (10) as

$$\dot{f}_i + c_i f_i = \sum_{j \neq i} a_{ij} \int_0^t [f_i(\tau - \tau_{ij}) - f_i(\tau)] \, d\tau + c_i f_{0i}$$

$$+ \sum_{j \neq i} a_{ij} b_{ij}(0) \qquad i = 1, 2, \cdots, n \qquad (11)$$

for all $t \geqq 0$. Then, treating the right sides of equations (11) as "driving functions," we can solve the $n$ first-order differential equations (11) to obtain

$$f_i = \int_0^t \exp\left[-c_i(t - u)\right] \sum_{j \neq i} a_{ij} \int_0^u [f_i(\tau - \tau_{ij}) - f_i(\tau)] \, d\tau \, du$$

$$+ f_i(0) \exp\left(-c_i t\right) + \int_0^t \exp\left[-c_i(t - \tau)\right][c_i f_{0i} + \sum_{j \neq i} a_{ij} b_{ij}(0)] \, d\tau$$

for all $i$ and all $t \geqq 0$. But

$$\int_0^u f_i(\tau - \tau_{ij}) \, d\tau = \int_0^{(u - \tau_{ij})} f_i(\tau) \, d\tau + \int_{-\tau_{ij}}^0 f_i(v) \, dv$$

for all $j$ and all $u \geqq 0$. Thus, with

$$p_i(t) = \int_0^t f_i(\tau) \, d\tau$$

for all $i$, we have (after some simplification)

$$\dot{p}_i = \int_0^t k_i(t - \tau) \sum_{j \neq i} a_{ij}[p_i(\tau - \tau_{ij}) - p_i(\tau)] \, d\tau$$

$$+ v_i + u_i(t), \qquad t \geqq 0 \quad i = 1, 2, \cdots, n \qquad (12)$$

in which each $v_i$ is a constant and each $u_i(t)$ approaches zero as $t \to \infty$. More explicitly, $k_i(t) = \exp(-c_i t)$,

$$v_i = f_{0i} + (c_i)^{-1} \sum_{j \neq i} a_{ij}\left[\int_{-\tau_{ij}}^0 f_i(u) \, du + b_{ij}(0)\right],$$

and $u_i(t) = [f_i(0) - v_i] \exp(-c_i t)$.

Equations of the same type as (12) have been studied extensively in connection with linear models of synchronized networks, and many

interesting and informative results have been obtained (see Refs. 2–6). In particular, Beneš has derived a sufficient condition for stability (see Refs. 2 and 3). His condition does not depend on the delays and is therefore of a very different type than condition (7). When applied to our system of equations (10) via the relation between equations (10) and (12), Beneš' condition reduces to our condition (9).* It therefore does not yield Pierce's result obtained via the network analog.

As is often true of pioneering work, Beneš' proof is long and involved. Our derivation of what amounts to his condition applied to our system is relatively simple and is much less involved compared with the rather long proof of our main result which asserts that if condition (7) is satisfied, then the system is stable. To a considerable extent, the methods of proof used here are very different from those of earlier writers concerned with the synchronization problem, and they provide some grip on the problem of estimating the rate of decay of the natural modes of the system. On the other hand, in this paper we do not consider several other important practical problems such as that of obtaining a useful upper bound on the contents of the buffers or of predicting the effects of variable transmission delays (resulting from temperature changes). There is a clear need for much more work in this area, especially in connection with models which take into account purposely-inserted nonlinearities.

## II. STATEMENT OF RESULTS

### 2.1 *Definitions and Assumptions*

Let $M$ denote an arbitrary complex matrix. We denote by $M^{tr}$ and $M^*$ the transpose of $M$ and the complex-conjugate transpose of $M$, respectively. If $M$ is not a row vector or a column vector, then $(M)_{ij}$ denotes the $ij$th element of $M$. If $M$ is a column vector then $(M)_j$ denotes the $j$th element of $M$. If $M$ is square, then $M^{(j, j)}$ denotes the matrix obtained from $M$ by deleting the $j$th row and column. The zero element of complex Euclidean $n$-space is indicated by $\theta$, and $1_n$ denotes the identity matrix of order $n$.

---

* Beneš proves that a suitable equilibrium state is approached. He makes no assertions concerning the rate at which it is approached. However, it is possible to modify the alternative proof[3] of Beneš' result due to Gersho and Karafin to show that the equilibrium state is approached exponentially under reasonable additional assumptions. For example, it suffices to assume that each transfer function $H_i(s)$ of Ref. 3 is a rational function of $s$. See the appendix.

The statement "for all $i$" means for all $i = 1, 2, \cdots, n$ in which $n$ denotes an arbitrary fixed positive integer (the number of geographically separated stations) such that $n \geqq 2$. If $f_i$ or $F = (f_1, f_2, \cdots, f_n)^{tr}$ denotes a differentiable function of $t$ or a differentiable $n$-vector-valued function of $t$, then $\dot{f}_i$ or $\dot{F}$ indicates the derivative with respect to $t$ of $f_i$ or $F$, respectively.

We assume throughout the paper that

(*i*) $A$ is the real $n \times n$ matrix defined by

$$(A)_{ii} = \sum_{j \neq i} a_{ij} \qquad \text{for all } i \tag{13}$$

$$(A)_{ij} = -a_{ij} \qquad \text{for all } i \neq j \tag{14}$$

in which $a_{ij} = a_{ji} \geqq 0$, for all $i \neq j$.

(*ii*) $$\det A^{(n,n)} > 0 \tag{15}$$

(*iii*) $$\tau_{ij} \geqq 0 \quad \text{for all} \quad i \neq j \tag{16}$$

(*iv*) the operator $\underline{A}$ is a mapping of the set of $n$-vector-valued functions of $t$ into itself defined by the condition that if $F(t) = [f_1(t), f_2(t), \cdots, f_n(t)]^{tr}$, then

$$[(\underline{A}F)(t)]_i = \sum_{j \neq i} a_{ij}[f_i(t) - f_i(t - \tau_{ij})] \tag{17}$$

for all $i$.

(*v*) $C = \operatorname{diag}(c_1, c_2, \cdots, c_n)$ with $c_i > 0$ for all $i$
(*vi*) $\tilde{A}$ is the complex $n \times n$ matrix defined by

$$(\tilde{A})_{ii} = \sum_{j \neq i} a_{ij} \qquad \text{for all } i \tag{18}$$

$$(\tilde{A})_{ij} = -a_{ij} \exp(-s\tau_{ij}) \qquad \text{for all } i \neq j \text{ and all complex } s \tag{19}$$

(*vii*) $\Delta(s)$ denotes the determinant $\det[s^2 1_n + sC + \tilde{A}]$.

Assumption (*ii*) is satisfied if there exists a path (not necessarily a direct path) from each station to every other station.*

---

* This proposition is a special case of a known result in probability theory. A simple proof is given in Ref. 3. A "network theoretic" proof is as follows. The matrix $A$ may be interpreted as the indefinite conductance matrix of a non-negative element $n$-node resistance network, and $A^{(n, n)}$ is the node-pair conductance matrix of the network obtained by grounding node $n$. If the original $n$-node network is connected, then the common-ground network possesses an open-circuit resistance matrix, which means that $\det A^{(n, n)} \neq 0$. But, since the network contains only nonnegative elements, $\det A^{(n, n)} > 0$.

The equations of our system are

$$\dot{f}_i = c_i(f_{0i} - f_i) + \sum_{j \neq i} a_{ij} \int_0^t [f_j(\tau - \tau_{ij}) - f_i(\tau)]\, d\tau$$

$$+ \sum_{j \neq i} a_{ij} b_{ij}(0), \qquad t \geq 0$$

for all $i$. Therefore, with $F = (f_1, f_2, \cdots, f_n)^{tr}$, we have

$$\ddot{F} + C\dot{F} + \underline{A}F = \theta, \qquad t \geq 0.$$

### 2.2 Results

Theorem 1 is the principal result of this paper:

*Theorem 1:   Let $F(\cdot)$ be a twice differentiable n-vector-valued function defined on $[-\bar{\tau}, \infty)$, in which $\bar{\tau} = \max_{i \neq j}\{\tau_{ij}\}$, such that*

$$\ddot{F} + C\dot{F} + \underline{A}F = \theta, \qquad t \geq 0^*.$$

*If*

$$c_i > \sum_{j \neq i} a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji})$$

*for all $i$, then there exist a constant $\rho$, positive constants $\beta$ and $\gamma$, and an n-vector-valued function $G(\cdot)$ defined for $t \in [0, \infty)$ such that*

$$|\, [G(t)]_i \,| \leq \beta e^{-\gamma t}, \qquad t \geq 0$$

*for all $i$, and*

$$F(t) = G(t) + \rho(1, 1, \cdots, 1)^{tr}$$

*for all $t \geq 0$.*

Some information concerning the rate of decay of the complex natural modes of the system, assuming that the stability condition of Theorem 1 is satisfied, is provided by the following theorem.

*Theorem 2:   Let $\tau_{ij} = \tau_{ji}$ for all $i \neq j$. If there exists a positive constant $\delta$ such that*

$$c_i - \sum_{j \neq i} a_{ij}\tau_{ij} \geq \delta$$

*for all $i$, and if $\triangle(s) = 0$ with $\mathrm{Im}\,[s] \neq 0$, then $\mathrm{Re}\,[s] \leq -\alpha_0$ in which $\alpha_0$ is the solution of*

$$-2\alpha_0 + \max_i c_i = [\max_i c_i - \delta]\exp(\alpha_0\bar{\tau})$$

---

* Questions concerning the existence and uniqueness of solutions of equations of this type are discussed in Ref. 7.

*where* $\bar{\tau} = \max_{i \neq j} \{\tau_{ij}\}$.

An expression for the final frequency $\rho$ of the system is given in Theorem 3:

*Theorem 3: If $F(\cdot)$ is differentiable on $[-\bar{\tau}, \infty)$ with $\bar{\tau} = \max_{i \neq j} \{\tau_{ij}\}$ such that for all $i$*

$$\dot{f}_i = c_i(f_{0i} - f_i) + \sum_{j \neq i} a_{ij} \int_0^t [f_j(\tau - \tau_{ij}) - f_i(\tau)] \, d\tau$$

$$+ \sum_{j \neq i} a_{ij} b_{ij}(0), \qquad t \geqq 0$$

*in which the $f_{0i}$ and the $b_{ij}(0)$ are constants, and if there exists a constant $\rho$ such that for all $i$: $(f_i - \rho) \to 0$ as $t \to \infty$, and $(f_i - \rho)$ is (absolutely) integrable on $[0, \infty)$, then*

$$\rho = \frac{\sum_i c_i f_{0i} + \sum_i \sum_{j \neq i} a_{ij} b_{ij}(0) + \sum_i \sum_{j \neq i} a_{ij} \int_{-\tau_{ij}}^0 f_j(u) \, du}{\sum_i c_i + \sum_i \sum_{j \neq i} a_{ij} \tau_{ij}}.$$

Our final result is as follows.

*Theorem 4: The statement obtained from the statement of Theorem 1 by replacing the condition that*

$$c_i > \sum_{j \neq i} a_{ij} \tfrac{1}{2}(\tau_{ij} + \tau_{ji}) \qquad \text{for all } i$$

*by the condition that*

$$c_i \geqq (2 \sum_{j \neq i} a_{ij})^{\frac{1}{2}} \qquad \text{for all } i$$

*is a theorem.*

III. PROOFS

3.1 *Proof of Theorem 1*

Our proof consists of proving the following four lemmas.

*Lemma 1: $\Delta(s)$ has a simple zero at $s = 0$.*

*Lemma 2: If*

$$c_i > \sum_{j \neq i} a_{ij} \tfrac{1}{2}(\tau_{ij} + \tau_{ji})$$

*for $i = 1, 2, \cdots, n$ then*

(i)  $\triangle(s) \neq 0$ for Re $[s] > 0$

(ii) there exists a positive constant $\alpha_0$ such that if $\triangle(s) = 0$ with $s = -\alpha + i\omega$ ($\alpha, \omega$ real; $\alpha \geqq 0$; $\omega \neq 0$), then $\alpha \geqq \alpha_0$.

*Lemma 3*:  Let $F(\cdot)$ be a twice differentiable n-vector-valued function defined on $[-\bar{\tau}, \infty)$, in which $\bar{\tau} = \max_{i \neq j} \{\tau_{ij}\}$, such that

$$\ddot{F} + C\dot{F} + \underline{A}F = \theta, \qquad t \geqq 0.$$

Let $\triangle(s)$ have a simple zero at $s = 0$, and with the exception of this zero assume that $\triangle(s) \neq 0$ for all Re $[s] \geqq -\alpha_0$ for some positive constant $\alpha_0$. Then there exist a constant n-vector $K$, positive constants $\beta$ and $\gamma$, and an n-vector-valued function $G(\cdot)$ defined for $t \; \varepsilon \; [0, \infty)$, such that

$$| \; [G(t)]_i \; | \leqq \beta e^{-\gamma t}, \qquad t \geqq 0$$

for all i, and $F(t) = G(t) + K$ for all $t \geqq 0$.

*Lemma 4*:  If $F(t) = G(t) + K$ satisfies

$$\ddot{F} + C\dot{F} + \underline{A}F = \theta, \qquad t \geqq 0$$

with $G(t) \to \theta$ as $t \to \infty$, and K a constant vector, then, for some constant $\rho$,

$$K = \rho(1, 1, \cdots, 1)^{tr}.$$

*Proof of Lemma 1*:  The determinant $\triangle(s)$ is analytic throughout the s-plane. It can be written as some power series

$$\sum_{n=0}^{\infty} \xi_n s^n$$

which converges for all s. Since $\triangle(0) = \det A = 0$, $\xi_0 = 0$. To prove Lemma 1, it suffices to show that

$$\lim_{s \to 0} \frac{\triangle(s)}{s} \neq 0,$$

for then $\xi_1 \neq 0$. We show this as follows.

We write $s^{-1} \det [s^2 1_n + sC + \tilde{A}]$ as $\det M$ in which

$$M = \text{diag} \; (s^2 + c_1 s, \; s^2 + c_2 s, \; \cdots, \; s^2 + c_{n-1} s, \; s + c_n) + \tilde{A}'$$

where $\tilde{A}'$ is obtained from $\tilde{A}$ by dividing each element of the nth row of $\tilde{A}$ by s. But, with $M^{(n,n)}$ the submatrix obtained from $M$ by deleting the nth row and column,

$$\det M = (s + c_n) \det M^{(n,n)} + s^{-1} \det P$$

in which

$$P = \tilde{A} + \text{diag} \ (s^2 + c_1 s, \ s^2 + c_2 s, \ \cdots, \ s^2 + c_{n-1} s, \ 0).$$

Since, by assumption, $\det A^{(n,n)} > 0$, we have

$$0 < \lim_{s \to 0} (s + c_n) \det M^{(n,n)}.$$

The determinant of $P$ is an analytic function of $s$ for all $s$. For $s = 0$, it vanishes. Therefore its power-series expansion at $s = 0$ is of the form

$$\sum_{n=1}^{\infty} p_n s^n.$$

We note that for $s$ real and positive, $\det P > 0$ because for $s$ real and positive $P$ is strongly dominant in its first $(n - 1)$ rows and at least weakly dominant in its last row.* Thus $p_1$ is not negative, and

$$\lim_{s \to 0} s^{-1} \det P \geqq 0.$$

Therefore

$$\lim_{s \to 0} s^{-1} \det [s^2 1_n + sC + \tilde{A}] > 0. \quad \square$$

*Proof of Lemma 2:* If $s$ is such that $\triangle(s) = 0$, then there exists a nonzero complex $n$-vector $x$ such that

$$(s^2 1_n + sC + \tilde{A})x = \theta$$

and consequently

$$x^*(s^2 1_n + sC + \tilde{A})x = 0.$$

With $\tilde{A}_H = \frac{1}{2}(\tilde{A} + \tilde{A}^*)$ and $\tilde{A}_S = \frac{1}{2}(\tilde{A} - \tilde{A}^*)$, we have

$$s^2 + s \frac{x^*Cx}{\| x \|^2} + \frac{x^* \tilde{A}_H x}{\| x \|^2} + \frac{x^* \tilde{A}_S x}{\| x \|^2} = 0 \tag{20}$$

in which $\| x \| = (x^*x)^{\frac{1}{2}}$,

$$(\tilde{A}_H)_{ij} = -\tfrac{1}{2}a_{ij}[\exp (-s\tau_{ij}) + \exp (-s^*\tau_{ji})], \qquad \text{for all } i \neq j$$

$$(\tilde{A}_H)_{ii} = \sum_{j \neq i} a_{ij}, \qquad \text{for all } i$$

and

$$(\tilde{A}_S)_{ij} = -\tfrac{1}{2}a_{ij}[\exp (-s\tau_{ij}) - \exp (-s^*\tau_{ji})], \qquad \text{for all } i \neq j$$

$$(\tilde{A}_S)_{ii} = 0, \qquad \text{for all } i.$$

---

* In other words, for such values of $s$, $\det P > 0$ because (for each such value of $s$) there exists a diagonal matrix $D = \text{diag}(1, 1, \ldots, k)$ with $k > 1$ such that $PD$ is strongly row-sum dominant.

Notice that both $\tilde{A}_H$ and $(i)^{-1}\tilde{A}_S$ $[i \triangleq (-1)^{\frac{1}{2}}$ when not used as an index] are hermitian matrices, and that

$$\frac{x^*\tilde{A}_H x}{\|\,x\,\|^2} \quad \text{and} \quad \frac{x^*\tilde{A}_S x}{\|\,x\,\|^2}$$

are real and pure imaginary, respectively.

*Sublemma 1:* If $\triangle(s) = 0$ with $s = \alpha + i\omega$ and $\alpha \geqq 0$, then

$$\omega^2 \geqq \alpha^2 + \alpha \min_i c_i \ . \tag{21}$$

*Proof:* From

$$s^2 + s\frac{x^*Cx}{\|\,x\,\|^2} + \frac{x^*\tilde{A}_H x}{\|\,x\,\|^2} + \frac{x^*\tilde{A}_S x}{\|\,x\,\|^2} = 0 \tag{22}$$

for some nonzero $x$ corresponding to the assumed value of $s = \alpha + i\omega$, we have

$$\alpha^2 - \omega^2 + \alpha\frac{x^*Cx}{\|\,x\,\|^2} + \frac{x^*\tilde{A}_H x}{\|\,x\,\|^2} = 0.$$

But for $\alpha \geqq 0$, $\tilde{A}$ is both at least weakly row-sum dominant and weakly column-sum dominant. Thus $\tilde{A}_H$ is at least weakly dominant and hence nonnegative definite (that is, $x^*\tilde{A}_H x \geqq 0$). Therefore

$$\omega^2 \geqq \alpha^2 + \alpha\frac{x^*Cx}{\|\,x\,\|^2}$$

$$\geqq \alpha^2 + \alpha \min_i c_i \ . \quad \square$$

*Sublemma 2:* If

$$c_i > \sum_{j \neq i} a_{ij} \max(\tau_{ij}, \tau_{ji})$$

*for all $i$, then $\triangle(s) \neq 0$ for* Re $[s] > 0$.

*Proof:* Sublemma 1 implies that $\triangle(s)$ has no zeros on the positive-real axis. Assume now that $s = \alpha + i\omega$ with $\alpha > 0$ and $\omega \neq 0$. Then, using equation (20),

$$2\alpha i\omega + i\omega\frac{x^*Cx}{\|\,x\,\|^2} + \frac{x^*\tilde{A}_S x}{\|\,x\,\|^2} = 0$$

or

$$2\alpha + \frac{x^*Cx}{\|\,x\,\|^2} + \frac{x^*[(i\omega)^{-1}\tilde{A}_S]x}{\|\,x\,\|^2} = 0. \tag{23}$$

Let $m_{ij}$ denote the $ij$th element of $(i\omega)^{-1}\tilde{A}_S$ . Then $m_{ii} = 0$ for all $i$, and for $i \neq j$:

$$m_{ij} = -\tfrac{1}{2}a_{ij}(i\omega)^{-1}[\exp(-\alpha\tau_{ij})\exp(-i\omega\tau_{ij}) - \exp(-\alpha\tau_{ji})\exp(i\omega\tau_{ji})]$$

$$= -\tfrac{1}{2}a_{ij}\exp(-\alpha\tau_{ij})(i\omega)^{-1}[\exp(-i\omega\tau_{ij}) - \exp(i\omega\tau_{ji})]$$

$$\quad + \tfrac{1}{2}a_{ij}(i\omega)^{-1}[\exp(-\alpha\tau_{ji})\exp(i\omega\tau_{ji}) - \exp(-\alpha\tau_{ij})\exp(i\omega\tau_{ji})]$$

$$= a_{ij}\exp(-\alpha\tau_{ij})\exp[i\omega\tfrac{1}{2}(\tau_{ji} - \tau_{ij})](\omega)^{-1}\sin[\tfrac{1}{2}(\tau_{ij} + \tau_{ji})\omega]$$

$$\quad + \tfrac{1}{2}a_{ij}\exp(i\omega\tau_{ji})(i\omega)^{-1}[\exp(-\alpha\tau_{ji}) - \exp(-\alpha\tau_{ij})].$$

It follows that

$$|m_{ij}| \leqq a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji}) + \tfrac{1}{2}a_{ij}|\omega^{-1}[\exp(-\alpha\tau_{ji}) - \exp(-\alpha\tau_{ij})]|.$$

But, from inequality (21), $\omega^2 \geqq \alpha^2$, so that

$$|\omega^{-1}[\exp(-\alpha\tau_{ji}) - \exp(-\alpha\tau_{ij})]|$$

$$\leqq |\alpha^{-1}[\exp(-\alpha\tau_{ji}) - \exp(-\alpha\tau_{ij})]|$$

and, as can easily be verified,

$$|\alpha^{-1}[\exp(-\alpha\tau_{ji}) - \exp(-\alpha\tau_{ij})]| \leqq |\tau_{ji} - \tau_{ij}|.$$

Therefore, for $i \neq j$,

$$|m_{ij}| \leqq a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji}) + a_{ij}\tfrac{1}{2}|\tau_{ji} - \tau_{ij}| = a_{ij}\max(\tau_{ij}, \tau_{ji}).$$

However, by the hypothesis of the sublemma,

$$c_i > \sum_{j \neq i} a_{ij}\max(\tau_{ij}, \tau_{ji})$$

for all $i$, and thus $C + (i\omega)^{-1}\tilde{A}_S$ is a strongly-dominant hermitian matrix. This means that $C + (i\omega)^{-1}\tilde{A}_S$ is positive definite, and hence that $\alpha$ in equation (23) is negative, which is a contradiction. $\square$

*Sublemma 3:* *If* $b_1, b_2, \cdots, b_n$ *are positive constants such that*

$$b_i > \sum_{j \neq i} a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji})$$

*for all* $i$, *then there exists a positive constant* $\sigma$ *with the property that if* $c_i \geqq b_i$ *for all* $i$, *and if* $\triangle(s) = 0$ *with* $s = \alpha + i\omega$ *and* $\alpha \geqq 0$ *and* $\omega \neq 0$, *then* $\alpha > \sigma$.

*Proof:* With $s = \alpha + i\omega$ and $\alpha \geqq 0$ and $\omega \neq 0$, and proceeding as in the proof of Sublemma 2, with any $C$ such that $c_i \geqq b_i$ for all $i$ we have

for some nonzero vector $x$

$$2\alpha + \frac{x^*Cx}{\|x\|^2} + \frac{x^*Mx}{\|x\|^2} = 0$$

in which $M$ is hermitian with elements $m_{ii} = 0$ for all $i$ and $m_{ij}$ ($i \neq j$) such that

$$| m_{ij} | \leq a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji}) + \tfrac{1}{2}a_{ij} \, | \, \omega^{-1}[\exp(-\alpha\tau_{ji}) - \exp(-\alpha\tau_{ij})] \, |.$$

By Sublemma 1,

$$\omega^2 \geq \alpha^2 + \alpha b$$

in which $b = \min_i b_i$. Thus $| \omega |^{-1} \leq (\alpha^2 + \alpha b)^{-\frac{1}{2}}$. Since (for all $\alpha \geq 0$)

$$| \exp(-\alpha\tau_{ji}) - \exp(-\alpha\tau_{ij}) | \leq \alpha \, | \, \tau_{ji} - \tau_{ij} \, |,$$

we have (for all $i \neq j$)

$$| m_{ij} | \leq a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji}) + \tfrac{1}{2}a_{ij}\alpha(\alpha^2 + \alpha b)^{-\frac{1}{2}} \, | \, \tau_{ji} - \tau_{ij} \, |.$$

Let

$$\epsilon = \min_i \left\{ b_i - \sum_{j \neq i} a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji}) \right\}.$$

Of course $\epsilon > 0$. Since $b > 0$, it is clear that there exists a positive constant $\sigma$ such that (for all $i \neq j$)

$$(2n)^{-1} \epsilon \geq \tfrac{1}{2}a_{ij}\alpha(\alpha^2 + \alpha b)^{-\frac{1}{2}} \, | \, \tau_{ji} - \tau_{ij} \, |$$

for all $0 \leq \alpha \leq \sigma$. It follows that $(C + M)$ is positive definite for all $\alpha$ such that $0 \leq \alpha \leq \sigma$. This means that if $\alpha \geq 0$ satisfies equation (23), then $\alpha > \sigma$. $\square$

*Sublemma 4:*   *There exists a positive constant $k$, independent of $C$, such that if $\triangle(s) = 0$ with $s = \alpha + i\omega$ and $\alpha \geq 0$, then*

$$| s | < \max_i c_i + k.$$

*Proof:*   The elements of $\tilde{A}$ are uniformly bounded on $\mathrm{Re}\,[s] \geq 0$. Let $k$ be any positive constant greater than

$$\sup \left\{ \left| \frac{x^*Ax}{\|x\|^2} \right|^{\frac{1}{2}} : \mathrm{Re}\,[s] \geq 0, \quad x \neq \theta \right\}.$$

If $\triangle(s) = 0$ with $s = \alpha + i\omega$ and $\alpha \geq 0$, then for some $x \neq \theta$

$$s^2 + s\frac{x^*Cx}{\|x\|^2} + \frac{x^*\tilde{A}x}{\|x\|^2} = 0$$

and

$$2s = -\frac{x^*Cx}{\|x\|^2} \pm \left[ \left(\frac{x^*Cx}{\|x\|^2}\right)^2 - 4\frac{x^*\tilde{A}x}{\|x\|^2}\right]^{\frac{1}{2}}.$$

Thus, since

$$\left| \left[\left(\frac{x^*Cx}{\|x\|^2}\right)^2 - 4\frac{x^*Ax}{\|x\|^2}\right]^{\frac{1}{2}} \right| \leqq \frac{x^*Cx}{\|x\|^2} + 2\left|\frac{x^*\tilde{A}x}{\|x\|^2}\right|^{\frac{1}{2}},$$

it follows that

$$|s| \leqq \frac{x^*Cx}{\|x\|^2} + \left|\frac{x^*\tilde{A}x}{\|x\|^2}\right|^{\frac{1}{2}}$$

$$\leqq \max_i c_i + k. \quad \square$$

By Sublemmas 1 and 3, if

$$c_i > \sum_{j \neq i} a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji}) \tag{24}$$

for all $i$, then $\triangle(s)$ has no zeros on the positive-real axis of the $s$ plane, and $\triangle(i\omega) \neq 0$ for all real $\omega \neq 0$. We are now in a position to show that if condition (24) is satisfied for all $i$, then $\triangle(s) \neq 0$ for Re $[s] > 0$.

Let $\underline{c}_i$ for $i = 1, 2, \cdots, n$ be any set of positive constants such that

$$\underline{c}_i > \sum_{j \neq i} a_{ij}\tfrac{1}{2}(\tau_{ij} + \tau_{ji})$$

for all $i$, and let $\bar{c}_i$ for $i = 1, 2, \cdots, n$ be any set of positive constants such that $\bar{c}_i \geqq \underline{c}_i$ and

$$\bar{c}_i > \sum_{j \neq i} a_{ij} \max (\tau_{ij}, \tau_{ji})$$

for all $i$. Let $c_i \, \varepsilon \, [\underline{c}_i, \bar{c}_i]$ for all $i$. Then, by Sublemmas 3 and 4 there exists a "half-moon-shaped" finite region $\mathfrak{R}$ in the strict right-half plane bounded by a line Re $[s] = \sigma$ and a semicircular arc $|s| \leqq K$ such that if $c_i \, \varepsilon \, [\underline{c}_i, \bar{c}_i]$ for all $i$ and $\triangle(s) = 0$ with Re $[s] > 0$, then $s$ lies *inside* $\mathfrak{R}$. By Sublemma 2, for $c_i = \bar{c}_i$ for all $i$, $\triangle(s) \neq 0$ for Re $[s] > 0$. On the boundary of $\mathfrak{R}$, $\triangle(s) \neq 0$ for all $c_i \, \varepsilon \, [\underline{c}_i, \bar{c}_i]$.

For each $i$, we can write $\triangle(s)$ as $A_i(s) + c_iB_i(s)$ in which $A_i$ and $B_i$ are entire functions which are independent of $c_i$. For $c_i = \bar{c}_i$ for $i = 2, 3, \cdots, n$

$$A_1(s) + c_1B_1(s) \neq 0$$

on the boundary of $\mathfrak{R}$ for all $c_1 \, \varepsilon \, [\underline{c}_1, \bar{c}_1]$. Thus $\triangle(s)$ has no zeros inside

$\mathcal{R}$ for $c_1 = \underline{c}_1$ and $c_i = \bar{c}_i$ for $i = 2, 3, \cdots, n$.* Similarly, for $c_1 = \underline{c}_1$ and $c_i = \bar{c}_i$ for $i = 3, 4, \cdots, n$

$$A_2(s) + c_2 B_2(s) \neq 0$$

on the boundary of $\mathcal{R}$ for all $c_2 \; \varepsilon \; [\underline{c}_2, \bar{c}_2]$. Therefore $\triangle(s)$ has no zeros inside $\mathcal{R}$ for $c_1 = \underline{c}_1$, $c_2 = \underline{c}_2$ and $c_i = \bar{c}_i$ for $i = 3, 4, \cdots, n$. By continuing in this manner we find that $\triangle(s)$ has no zeros inside $\mathcal{R}$ when $c_i = \underline{c}_i$ for all $i$, which shows that if condition (24) is satisfied for all $i$, then $\triangle(s) \neq 0$ for Re $[s] > 0$.

The following argument completes the proof of Lemma 2.

If $\triangle(s) = 0$ with $s = -\alpha + i\omega$, then for some $x \neq \theta$

$$\alpha^2 - \omega^2 - \alpha \frac{x^* C x}{\|x\|^2} + \frac{x^* \tilde{A}_H x}{\|x\|^2} = 0.$$

Since $\triangle(s)$ is an entire function of $s$, in any circle in the $s$ plane $\triangle(s)$ has at most a finite number of zeros. Thus, either there exists a positive constant $\alpha_0$ such that if $\triangle(s) = 0$ with $s = -\alpha + i\omega$ and $\alpha > 0$, then $\alpha \geqq \alpha_0$, or there exists an infinite sequence $s_1, s_2, \cdots$, such that $\triangle(s_n) = 0$ for all $n$ and, with $s_n = -\alpha_n + i\omega_n$, $\omega_n \to \infty$ and $\alpha_n \to 0+$ as $n \to \infty$. In the later case, since for $n = 1, 2, \cdots$

$$\omega_n^2 = \alpha_n^2 - \alpha_n \frac{x_n^* C x_n}{\|x_n\|^2} + \frac{x_n^* \tilde{A}_H x_n}{\|x_n\|^2}$$

with $x_n \neq \theta$ an associate of $s_n$, we would have

$$\alpha_n^2 - \alpha_n \frac{x_n^* C x}{\|x\|^2} + \frac{x_n^* A_H x_n}{\|x\|^2} \to \infty$$

as $\alpha_n \to 0+$ which is impossible since for any real $\beta < 0$

$$\frac{x_n^* \tilde{A}_H x_n}{\|x_n\|^2}$$

is bounded on the strip $\beta \leqq$ Re $[s] \leqq 0$ uniformly in $x_n$. $\square$

*Proof of Lemma 3:*   We have

$$\ddot{F} + C\dot{F} + \underline{A}F = \theta, \qquad t \geqq 0. \tag{25}$$

---

\* Here we use the following known result.[8] Let $\mathcal{R}$ be a closed region in the $s$ plane, the boundary of which consists of a finite number of regular arcs; let the functions $f(s)$ and $h(s)$ be regular on $\mathcal{R}$. Assume that for no value of the real parameter $c$, running through the interval $a \leq c \leq b$, does the function $f(s) + ch(s)$ become equal to 0 on the boundary of $\mathcal{R}$. Then the number $N(c)$ of the zeros of $f(s) + ch(s)$ inside $\mathcal{R}$ is independent of $c$ for $a \leq c \leq b$.

Using the Bellman–Gronwall Lemma (see p. 31 of Ref. 7), we can prove that there are constants $k_1$, $k_2 > 0$ such that for all $i$

$$| F_i(t) | \leq k_1 \exp (k_2 t), \qquad t \geq 0.$$

Thus $F(\cdot)$ possesses a Laplace transform

$$\int_0^\infty F(t)e^{-st} \, dt \triangleq \hat{F}(s)$$

which is well defined for all $s$ such that Re $[s] > k_2$. Therefore

$$s^2\hat{F}(s) + sC\hat{F}(s) + \tilde{A}\hat{F}(s) = sV_1 + V_2 + W(s)$$

in which the constant vectors $V_1$ and $V_2$ take into account the values of $\dot{F}(0)$ and $F(0)$, and for all $i$

$$[W(s)]_i = \sum_{j \neq i} a_{ij} \exp (-s\tau_{ij}) \int_{-\tau_{ij}}^0 F_j(t)e^{-st} \, dt.$$

Thus

$$\hat{F}(s) = [s^2 1_n + sC + \tilde{A}]^{-1}[sV_1 + V_2 + W(s)].$$

Of course

$$F(t) = (2\pi i)^{-1} \int_{(l+)} \hat{F}(s)e^{st} \, ds, \qquad t \geq 0$$

in which $(l+)$ is some line parallel to and to the right of the imaginary axis of the $s$ plane. Let $(l-)$ denote the line indicated in Fig. 1. Then



Fig. 1 — Relation between $(l-)$ and $(l+)$.

(see Fig. 1)

$$\int_{(l+)} = \lim_{\beta\to\infty} \int_{\alpha_1-i\beta}^{\alpha_2+i\beta} ,$$

and for $t \geqq 0$

$$\int_{\alpha_1-i\beta}^{\alpha_1+i\beta} + \int_{\alpha_1+i\beta}^{-\alpha_2+i\beta} + \int_{-\alpha_2+i\beta}^{-\alpha_2-i\beta} + \int_{-\alpha_2-i\beta}^{\alpha_1-i\beta} = K$$

in which $K$ is a constant $n$-vector (since the contour contains a single pole at the origin). Since the integrand $\to \theta$ uniformly in $-\alpha_2 \leqq$ Re $[s] \leqq \alpha_1$ as $|$ Im $[s] | \to \infty$, the integrals over the horizontal pieces $\to \theta$ as $\beta \to \infty$, and therefore for $t \geqq 0$

$$\lim_{\beta\to\infty} \int_{\alpha_1-i\beta}^{\alpha_1+i\beta} = K - \lim_{\beta\to\infty} \int_{-\alpha_2+i\beta}^{-\alpha_2-i\beta} .$$

We now show that

$$\lim_{\beta\to\infty} \int_{-\alpha_2+i\beta}^{-\alpha_2-i\beta} \to \theta$$

*exponentially* as $t \to \infty$.

On $(l-)$ each element of $W(s)$ is bounded. Thus on $(l-)$ each element of

$$Q(s) \triangleq [1_n s^2 + sC + \tilde{A}]^{-1}[V_2 + W(s)]$$

(we consider the part involving $V_1$ later) is bounded and of order at most $| s |^{-2}$ as $| s | \to \infty$. Therefore

$$\int_{(l-)} Q(s)e^{st} \, ds = \exp(-\alpha_2 t) \lim_{\beta\to\infty} \int_{i\beta}^{-i\beta} Q(-\alpha_2 + i\omega)e^{i\omega t} \, d(i\omega)$$

But $\int_{i\beta}^{-i\beta}$ is uniformly bounded on the $t$-set $[0, \infty)$, and hence

$$\int_{(l-)} Q(s)e^{st} \, ds$$

approaches zero exponentially.

Consider now the integral

$$\int_{(l-)} [1_n s^2 + sC + \tilde{A}]^{-1} sV_1 e^{st} \, ds.$$

We cannot directly apply the same argument as for the integration of $Q$ because here it need not be true that on $(l-)$ each component of the

integrand is of order at most $|s|^{-2}$ as $|s| \to \infty$. However, if on $(l-)$ some component $R_i(s)$ of $[1_n s^2 + sC + \tilde{A}]^{-1} s V_1$ is not of order at most $|s|^{-2}$ as $|s| \to \infty$, then it can be written (to within a constant multiplier) as

$$R_i(s) = \frac{s^{2n-1} \Sigma^{(2n-1)} + s^{2n-2} \Sigma^{(2n-2)} + \cdots + \Sigma^{(0)}}{s^{2n} + \xi_D(s)}$$

in which $s^{-2n} \xi_D(s) \to 0$ as $|s| \to \infty$ on $(l-)$, and $\Sigma^{(2n-1)}$, $\Sigma^{(2n-2)}, \cdots, \Sigma^{(0)}$ denote sums of exponentials (some sums may be constants). Thus the sum of $R_i(s)$ and

$$-\frac{\Sigma^{(2n-1)}}{s + \sigma} \qquad (\sigma > \alpha_0)$$

is of order at most $|s|^{-2}$ as $|s| \to \infty$ on $(l-)$. Since the inverse Laplace transform of

$$\frac{\Sigma^{(2n-1)}}{s + \sigma}$$

vanishes faster than some damped exponential, we see that

$$\int_{(l-)} [1_n s^2 + sC + \tilde{A}]^{-1} s V_1 e^{st} \, ds, \qquad t \geq 0$$

can be written as the sum of two integrals, the components of both of which approach zero at least as fast as some damped exponential. $\square$

*Proof of Lemma 4:* Since $F(\cdot)$ satisfies

$$\ddot{F} + C\dot{F} + \underline{A}F = \theta, \qquad t \geq 0$$

we have

$$\dot{F}(t) = e^{-Ct} \dot{F}(0) - \int_0^t e^{-C(t-\tau)} (\underline{A}F)(\tau) \, d\tau, \qquad t \geq 0.$$

But $F(t) \to K$ as $t \to \infty$. Therefore as $t \to \infty$,

$$-\dot{F}(t) \to \lim_{t\to\infty} \int_0^t e^{-C(t-\tau)} AK \, d\tau = C^{-1} AK.$$

In addition, since

$$F(t) = \int_0^t \dot{F}(t) \, dt + F(0), \qquad t \geq 0$$

we see that $-C^{-1}AK$ [that is, the limit of $\dot{F}(t)$ as $t \to \infty$] must be the zero vector, for otherwise $F(t)$ could not approach a constant vector as

$t \rightarrow \infty$. Thus $AK = \theta$. We note that $A(1, 1, \cdots, 1)^{tr} = \theta$, and that $A$ is of rank $(n - 1)$. Therefore $K = \rho(1, 1, \cdots, 1)^{tr}$ for some constant $\rho$. $\square$

### 3.2 *Proof of Theorem 2*

By Lemma 2 of the proof of Theorem 1, $\triangle(s) \neq 0$ for Re $[s] \geqq 0$ and $s \neq 0$. Let $\triangle(s) = 0$ with $s = -\alpha + i\omega$, $\alpha > 0$, and $\omega \neq 0$. Then for some $x \neq \theta$

$$-2\alpha i\omega + i\omega \frac{x^*Cx}{\|x\|^2} + \frac{x^*\tilde{A}_s x}{\|x\|^2} = 0$$

or

$$-2\alpha + \frac{x^*Cx}{\|x\|^2} + \frac{x^*(i\omega)^{-1}\tilde{A}_s x}{\|x\|^2} = 0.$$

Let $\bar{\tau} = \max_{i \neq j} \{\tau_{ij}\}$. Then $(i\omega)^{-1}\tilde{A}_s = e^{\alpha\bar{\tau}}B_s$ with

$$(B_s)_{ij} = 0, \quad \text{for all } i = j$$

$$= a_{ij}\tau_{ij} \exp\left[\alpha(\tau_{ij} - \bar{\tau})\right]\frac{\sin \omega\tau_{ij}}{\omega\tau_{ij}}, \quad \text{for all } i \neq j.$$

Thus since $| (B_s)_{ij} | \leqq a_{ij}\tau_{ij}$ for all $i \neq j$, all $\alpha > 0$, and all $\omega \neq 0$; and

$$c_i - \sum_{j \neq i} a_{ij}\tau_{ij} \geqq \delta$$

for all $i$, it follows that

$$\frac{x^*Cx}{\|x\|^2} \pm \frac{x^*B_s x}{\|x\|^2} \geqq \delta$$

for all $\alpha > 0$ and all $\omega \neq 0$. Therefore with

$$c = \frac{x^*Cx}{\|x\|^2} \quad \text{and} \quad b = \frac{x^*B_s x}{\|x\|^2},$$

we have

$$-2\alpha + c + be^{\alpha\bar{\tau}} = 0$$

with $| b | \leqq c - \delta$ and $c \leqq \max_i c_i$. But

$$-2\alpha + c \leqq | -2\alpha + c | = | b | e^{\alpha\bar{\tau}} \leqq (c - \delta)e^{\alpha\bar{\tau}}$$

and hence

$$-2\alpha + c \leqq (c - \delta)e^{\alpha\bar{\tau}}, \quad c \leqq \max_i c_i.$$

Therefore $\alpha$ must not be less than $\alpha_0$ , the unique solution of

$$-2\alpha_0 + \max_i c_i = [\max_i c_i - \delta]e^{\alpha_0\bar{\tau}}. \quad \square$$

### 3.3 *Proof of Theorem 3*

We have for all $i$

$$\dot{f}_i = c_i(f_{0i} - f_i) + \sum_{j \neq i} a_{ij}b_{ij} , \qquad t \geq 0$$

with

$$b_{ij}(t) = [f_j(t - \tau_{ij}) - f_i(t)].$$

Observe that

$$\int_0^\tau f_i(t - \tau_{ij}) \, dt$$

$$= \int_{-\tau_{ij}}^{\tau - \tau_{ij}} f_i(u) \, du$$

$$= \int_{-\tau_{ij}}^0 f_i(u) \, du + \int_0^{\tau - \tau_{ij}} [f_i(u) - \rho] \, du + (\tau - \tau_{ij})\rho$$

$$= (\tau - \tau_{ij})\rho + \int_{-\tau_{ij}}^0 f_i(u) \, du + \int_0^\infty [f_i(u) - \rho] \, du$$

$$- \int_{\tau - \tau_{ij}}^\infty [f_i(u) - \rho] \, du.$$

Since

$$f_i = \exp(-c_it)f_i(0) + \int_0^t \exp[-c_i(t - \tau)]$$

$$\cdot [c_if_{0i} + \sum_{j \neq i} a_{ij}b_{ij}(\tau)] \, d\tau, \qquad t \geq 0$$

for all $i$,

$$f_i = \exp(-c_it)f_i(0) + \int_0^t \exp[-c_i(t - \tau)]\Big\{c_if_{0i} + \sum_{j \neq i} a_{ij}b_{ij}(0)$$

$$- \sum_{j \neq i} a_{ij}\tau_{ij}\rho + \sum_{j \neq i} a_{ij} \int_{-\tau_{ij}}^0 f_i(u) \, du + \sum_{j \neq i} a_{ij}(p_j - p_i)$$

$$- \sum_{j \neq i} a_{ij}[q_{ij}(\tau) - r_i(\tau)]\Big\} \, d\tau, \qquad t \geq 0 \tag{26}$$

for all $i$, in which

$$p_i = \int_0^\infty [f_i(u) - \rho]\, du$$

$$q_{ij}(\tau) = \int_{\tau - \tau_{ij}}^\infty [f_i(u) - \rho]\, du$$

$$r_i(\tau) = \int_\tau^\infty [f_i(u) - \rho]\, du.$$

The functions $q_{ij}(\tau)$ and $r_i(\tau)$ approach zero as $\tau \to \infty$. In order that the asymptotic values (as $t \to \infty$) of both sides of equation (26) agree, we must have

$$\rho c_i = c_i f_{0i} + \sum_{j \neq i} a_{ij} b_{ij}(0) - \sum_{j \neq i} a_{ij} \tau_{ij} \rho$$

$$+ \sum_{j \neq i} a_{ij} \int_{-\tau_{ij}}^0 f_i(u)\, du + \sum_{j \neq i} a_{ij}(p_j - p_i)$$

for all $i$. But (using the assumption that $a_{ij} = a_{ji}$ for all $i \neq j$)

$$\sum_i \sum_{j \neq i} a_{ij}(p_j - p_i) = 0.$$

Therefore

$$\rho \sum_i c_i = \sum_i c_i f_{0i} + \sum_i \sum_{j \neq i} a_{ij} b_{ij}(0)$$

$$- \rho \sum_i \sum_{j \neq i} a_{ij} \tau_{ij} + \sum_i \sum_{j \neq i} a_{ij} \int_{-\tau_{ij}}^0 f_i(u)\, du. \quad \square *$$

### 3.4 *Proof of Theorem 4*

The following lemma and Lemmas 1, 3, and 4 of the proof of Theorem 1 prove Theorem 4.

*Lemma 2′:  If*

$$c_i \geq \left(2 \sum_{j \neq i} a_{ij}\right)^{\frac{1}{2}} \qquad \text{for all } i,$$

*then there exists a positive constant $\alpha_0$ such that $\triangle(s) = 0$ and $s \neq 0$ imply that $\mathrm{Re}\,[s] \leq -\alpha_0$.*

*Proof of Lemma 2′:*  With $D_0 = \mathrm{diag}\,\{\sum_{j \neq i} a_{ij}\}$ and $\tilde{A}' = \tilde{A} - D_0$,

---

* The last part of this proof, which involves the observation that the double sum is zero, is similar to an argument used by Brilliant (see the appendix of Ref. 5).

$$\Delta(s) = \det [1_n s^2 + sC + \tilde{A}]$$
$$= \det [1_n s^2 + sC + D_0 + \tilde{A}']$$
$$= \det [1_n s^2 + sC + D_0] \cdot \det [1_n + (1_n s^2 + sC + D_0)^{-1} \tilde{A}'].$$

The diagonal elements of $D_0$ are positive. It is therefore clear that there is a positive constant $\alpha_0'$ such that if $\det [1_n s^2 + sC + D_0] = 0$, then Re $[s] \leq -\alpha_0'$. It is also clear that $\det [1_n + (1_n s^2 + sC + D_0)^{-1} \tilde{A}']$ is not zero for all $s$ such that Re $[s] \geq 0$ and

$$| (s^2 + sc_i + \sum_{j \neq i} a_{ij})^{-1} \sum_{j \neq i} a_{ij} | < 1 \qquad \text{for all } i, \qquad (27)$$

since for those values of $s$ the matrix $[1_n + (1_n s^2 + sC + D_0)^{-1} \tilde{A}']$ is strongly row-sum dominant.

By assumption,

$$c_i \geq (2 \sum_{j \neq i} a_{ij})^{\frac{1}{2}} \qquad \text{for all } i.$$

It is a simple matter to verify that this assumption implies that condition (27) is satisfied for all $s = i\omega$ with $\omega$ real and $\omega \neq 0$. Thus $\Delta(i\omega) \neq 0$ for all $\omega \neq 0$. But for all $i$

$$(s^2 + sc_i + \sum_{j \neq i} a_{ij})^{-1} \sum_{j \neq i} a_{ij}$$

is analytic throughout the closed right-half $s$ plane, and

$$| (s^2 + sc_i + \sum_{j \neq i} a_{ij})^{-1} \sum_{j \neq i} a_{ij} | \leq 1$$

for all $i$ and for all $s = i\omega$. By the Maximum Modulus Theorem, condition (27) is satisfied for all $s$ such that Re $[s] > 0$.[*] Therefore $\Delta(s) \neq 0$ for all $s \neq 0$ such that Re $[s] \geq 0$. Finally, the argument used to prove the last part of Lemma 2 shows that there exists a positive constant $\alpha_0$ such that if $\Delta(s) = 0$ with $s \neq 0$, then Re $[s] \leq -\alpha_0$. $\square$

APPENDIX

*Left-Half-Plane Zeros of det $[I - B(s)]$*

We wish to show here that all of the left-half-plane zeros of $\det [I - B(s)]$, in which $I - B(s)$ is as defined in Ref. 3, lie to the left of some line which is parallel to and lies to the left of the imaginary axis of the complex $s$ plane, provided that each $H_i(s)$, which enters into the definition of $B(s)$, is a meromorphic function of $s$ such that there

---

[*] This type of argument is also used in Ref. 3.

exist positive constants $\sigma_i$ and $K_i$ with the property that $\mid H_i(s) \mid \leqq K_i$ for all $s$ with $-\sigma_i \leqq \text{Re } [s] < 0$.

Assume that what we wish to prove is false. Then, as in the proof of the last part of Lemma 2, there would exist a sequence $\mathcal{S} = \{s_k\}_0^\infty$ such that Re $[s_k] < 0$ for all $k$, Re $[s_k] \to 0$ as $k \to \infty$, Im $[s_k] \to \infty$ as $k \to \infty$, and det $[I - B(s_k)] = 0$ for all $k \geqq 0$. But the complex numbers $\tilde{a}_{ij}$ of Ref. 3 are bounded on $\mathcal{S}$ and $H_i(s_k)[s_k + H_i(s_k)]^{-1} \to 0$ as $k \to \infty$. This means (see Ref. 3) that there is a positive number $k'$ such that the matrix $[I - B(s_k)]$ is strongly dominant for all $k \geqq k'$, which contradicts the assumption that det $[I - B(s_k)] = 0$ for all $k \geqq 0$, and proves that our assertion is true.

REFERENCES

1. Pierce, J. R., "Synchronizing Digital Networks," B.S.T.J., *48*, No. 3 (March 1969), pp. 615–636.
2. Karnaugh, M., "A Model for the Organic Synchronization of Communications Systems," B.S.T.J., *45*, No. 10 (December 1966), pp. 1705–1735.
3. Gersho, A., and Karafin, B. J., "Mutual Synchronization of Geographically Separated Oscillators," B.S.T.J., *45*, No. 10 (December 1966), pp. 1689–1704.
4. Brilliant, M. B., "The Determination of Frequency in Systems of Mutually Synchronized Oscillators," B.S.T.J., *45*, No. 10 (December 1966), pp. 1737–1748.
5. Brilliant, M. B., "Dynamic Response of Systems of Mutually Synchronized Oscillators," B.S.T.J., *46*, No. 2 (February 1967), pp. 319–356.
6. Candy, J. C., and Karnaugh, M., "Organic Synchronization: Design of the Controls and Some Simulation Results," B.S.T.J., *47*, No. 2 (February 1968), pp. 227–259.
7. Bellman, R., and Cooke, K. L., *Differential-Difference Equations*, New York: Academic Press, 1963, p. 171.
8. A. M. Ostrowski, *Solution of Equations and Systems of Equations*, New York: Academic Press, 1966, p. 222.

# Bounds on the Bias of Signal Parameter Estimators

By JACOB ZIV

*Any estimator which is constrained to take values in a finite range is, in general, biased. Many times the bias is unknown; furthermore, in some cases the bias may become the main contributor to the mean square error of an estimator. This paper derives upper and lower bounds on the bias of a finite-range, signal parameter estimator.*

## I. INTRODUCTION AND MAIN RESULTS

### 1.1 *Introduction*

Let the parameter be denoted by $a$ and let $a$ take values in $[-\alpha, \alpha]$. We refer to $2\alpha$ as the *a priori* range (or space) of $a$. We assume that there exists probabilistic mapping from the parameter space to an observation space, that is, a probability law that governs the effect of $a$ on the observation.[1] This probability law will be referred to as the "channel." After observing the "outcome" which is a point in the observation space, we estimate the value of $a$. Let this estimate be denoted by $\hat{a}$. Clearly, $\hat{a}$ is a random variable.

We assume, throughout this paper, that $\hat{a}$ takes values in $[-A, A]$. Let the bias be defined

$$b(a) \triangleq E_a[\hat{a} - a] = \int (\hat{a} - a) \, dp(\hat{a} \mid a) \tag{1}$$

where $p(\hat{a} \mid a)$ is the probability distribution function of $\hat{a}$ given $a$. Assume that we are now told that the true value of the parameter $a$ is either $a_1$ or $-a_1$ with equal probabilities. Let $H_{a_1}$ be the hypothesis that $a = a_1$ and let $H_{-a_1}$ be the hypothesis that $a = -a_1$. The minimum probability of error is (dropping the subscript 1 from $a_1$):

$$P_e\{a, -a\} \triangleq \text{Min} \, \{\tfrac{1}{2}[Pr \, \{a \mid -a\} + Pr \, \{-a \mid a\}]\}$$

where $\Pr\{a \mid -a\}$ is the probability that the decision will be $a$, given that $-a$ is transmitted, and where the minimization is carried over all possible decision rules. (A decision rule is a mapping from the observation space to the set $\{-a; a\}$.) Then we show in Section II that

$$\tfrac{1}{2}[b(a) - b(-a)] \leqq -AP_e\{-a; a\} + (A - a); \quad a \geqq 0; \quad (2\text{a})$$

$$\tfrac{1}{2}[b(a) - b(-a)] \geqq AP_e\{-a; a\} - (A + a); \quad a \geqq 0. \quad (2\text{b})$$

By equation (2a) we have

$$\tfrac{1}{2}[\mid b(-a) \mid + \mid b(a) \mid] \geqq AP_e\{-a; a\} - (A - a); \quad a \geqq 0.$$

Hence, an estimator must have a nonzero bias if

$$\frac{\mid a \mid}{A} > 1 - P_e\{-a; a\}. \quad (3)$$

The bounds of equation (2) are the main result of this paper.

If we assume that for any $a$, $b(a) = -b(-a)$ we have by equation (2) that

$$b(a) \leqq -AP_e\{-a; a\} + (A - a) \quad (4\text{a})$$

and

$$b(a) \geqq AP_e(-a; a) - (A + a). \quad (4\text{b})$$

These bounds are sketched in Fig. 1.

If, in addition, we assume a symmetry around $a$ in the sense that

$$\int_{-A+\mid a \mid \leqq \hat{a}-a \leqq A-\mid a \mid} (\hat{a} - a) \, dp(\hat{a} \mid a) = 0. \quad (5\text{a})$$

We show (see Section II) that in this case

$$b(a) \geqq -aP_e\{-a; a\}; \quad -A \leqq a \leqq -\frac{A}{2} \quad (5\text{b})$$

$$(A - a)P_e(-a; a) \geqq b(a) \geqq 0; \quad -\frac{A}{2} \leqq a \leqq 0 \quad (5\text{c})$$

and

$$b(a) \leqq -aP_e\{-a; a\}; \quad A \geqq a \geqq \frac{A}{2} \quad (5\text{d})$$

$$-(A + a)P_e(-a; a) \leqq b(a) \leqq 0; \quad \frac{A}{2} \geqq a \geqq 0. \quad (5\text{e})$$

The bias $b(a)$ is unknown, in general. However, the probability

Fig. 1 — Bounds on the bias of an estimator.

$P_e\{a; -a\}$ is known for many important cases. The bounds derived here depend on the channel probability law through $P_e(-a; a)$ only, and therefore are easy to compute in many cases.

### 1.2 *Sharpness*

Section III shows that, for one special case, $b(a)$ is given by

$$b(a) = -2AP_e\{-a; a\} + (A - a); \quad a \geqq 0, \qquad (6a)$$

$$b(a) = 2AP_e\{-a; a\} - (A + a); \quad a \leqq 0. \qquad (6b)$$

Section III also shows that, for another special case, $b(a)$ is given by

$$b(a) = 2AP_e\{-a; a\} - (A + a); \quad a \geqq 0, \qquad (6c)$$

$$b(a) = -2AP_e\{-a; a\} + (A - a); \quad a < 0. \qquad (6d)$$

The comparison of equation (6) with the bounds of equations (4) and (2) indicates the degree of sharpness of these bounds (see Fig. 1).

### 1.3 *Examples*

Let the received message be a sample function of the random process

$$r(t) = s(t - a) + n(t), \qquad (7a)$$

where $a$ is an unknown parameter constrained to take values in $(-\alpha; \alpha)$. The term $n(t)$ is assumed to be white gaussian noise with (double sided) spectral density $N_o$.

Let

$$\int_{-\infty}^{\infty} s^2(t) \, dt = E, \tag{7b}$$

$$\rho(2a) = \frac{1}{E} \int_{-\infty}^{\infty} s(t - a)s(t + a) \, dt, \tag{7c}$$

$$q = ([1 - \rho(2a)]E/2N_o)^{\frac{1}{2}}. \tag{7d}$$

Hence, in this case,[2]

$$P_e(-a; a) = (2\pi)^{-\frac{1}{2}} \int_q^{\infty} \exp(-x^2/2) \, dx$$

$$\geqq (2\pi)^{-\frac{1}{2}} \int_{(E/N_o)^{\frac{1}{2}}}^{\infty} \exp(-x^2/2) \, dx. \tag{8a}$$

Hence, by equation (3)

$$\mid b(a) \mid > 0 \quad \text{for any} \quad \frac{\mid a \mid}{A} > 1 - (2\pi)^{-\frac{1}{2}} \int_{(E/N_o)^{\frac{1}{2}}}^{\infty} \exp(-x^2/2) \, dx. \tag{8b}$$

Furthermore, if the channel is that of equation (7) and if $\hat{a}$ is a maximum likelihood estimator, then it follows from equation (7) that the conditions of equation (5a) are satisfied, since the maximum likelihood procedure for estimating $a$ is to evaluate

$$\lambda(a^*) = \int_{-\infty}^{\infty} r(t)s(t - a^*) \, dt$$

$$= \int_{-\infty}^{\infty} s(t - a)s(t - a^*) \, dt + \int_{-\infty}^{\infty} n(t)s(t - a^*) \, dt$$

and to set $\hat{a}$ to the value of $a^*(-A \leq a^* \leq A)$ for which $\lambda(a^*)$ is maximum. Hence, the statistics of $\lambda(a_1^*)$ are the same as those of $\lambda(a_2^*)$ if $\frac{1}{2}(a_1^* + a_2^*) = a$; also, $b(a) = -b(-a)$. Therefore by equations (5) and (8a),

$$b(a) \geqq -a(2\pi)^{-\frac{1}{2}} \int_{(E/N_o)^{\frac{1}{2}}}^{\infty} \exp(-x^2/2) \, dx; \quad -A \leq a \leq -\frac{A}{2} \tag{9a}$$

$$b(a) \geqq 0; \quad -\frac{A}{2} \leq a \leq 0 \tag{9b}$$

$$b(a) \leqq -a(2\pi)^{-\frac{1}{2}} \int_{(E/N_o)^{\frac{1}{2}}}^{\infty} \exp\left(-x^2/2\right) dx; \qquad \frac{A}{2} \leqq a \leqq A \qquad (9\text{c})$$

$$b(a) \leqq 0; \qquad 0 < a \leqq \frac{A}{2}. \qquad (9\text{d})$$

## II. DERIVATION OF THE BOUNDS[3]

$$b(a) \triangleq \int (\hat{a} - a)\, dp(\hat{a} \mid a)$$

$$= \int_{\hat{a}>0} (\hat{a} - a)\, dp(\hat{a} \mid a) + \int_{\hat{a}\leqq 0} (\hat{a} - a)\, dp(\hat{a} \mid a). \qquad (10)$$

Now,

$$\int_{\hat{a}>0} (\hat{a} - a)\, dp(\hat{a} \mid a) \geqq -a \Pr\{\hat{a} > 0 \mid a\} \qquad (11\text{a})$$

$$\int_{\hat{a}>0} (\hat{a} - a)\, dp(\hat{a} \mid a) \leqq (A - a) \Pr\{\hat{a} > 0 \mid a\} \qquad (11\text{b})$$

$$\int_{\hat{a}\leqq 0} (\hat{a} - a)\, dp(\hat{a} \mid a) \geqq -(A + a) \Pr\{\hat{a} \leqq 0 \mid a\} \qquad (11\text{c})$$

$$\int_{\hat{a}\leqq 0} (\hat{a} - a)\, dp(\hat{a} \mid a) \leqq -a \Pr\{\hat{a} \leqq 0 \mid a\}. \qquad (11\text{d})$$

Also, we have that

$$\Pr\{\hat{a} > 0 \mid a\} = 1 - \Pr\{\hat{a} \leqq 0 \mid a\}. \qquad (12)$$

Inserting equations (11) and (12) into equation (10), we have

$$b(a) \geqq A \Pr\{\hat{a} > 0 \mid a\} - (A + a) \qquad (13)$$

$$b(a) \leqq -A \Pr\{\hat{a} \leqq 0 \mid a\} + (A - a). \qquad (14)$$

Consider the following detection problem. Assume that $a$ and $-a\,(a > 0)$ are used as two signals for equiprobable binary signalling; decide on $a$ if $\hat{a} > 0$ and decide on $-a$ if $\hat{a} \leqq 0$. The probability of error associated with this detection procedure is given by

$$P_a \triangleq \tfrac{1}{2} \Pr\{\hat{a} > 0 \mid -a\} + \tfrac{1}{2} \Pr\{\hat{a} \leqq 0 \mid a\}; \qquad a > 0. \qquad (15)$$

The error probability $P_a$ is lower bounded by $P_e(-a; a)$ which is the probability of error that is associated with the optimal binary detection scheme for this detection problem. In the same way $P_a$ is upper bounded by $1 - P_e(-a; a)$.

Hence,

$$1 - P_e(-a; a) \geqq \tfrac{1}{2}[\text{Pr } \{\hat{a} \geqq 0 \mid -a\}$$
$$+ \text{Pr } \{\hat{a} \leqq 0 \mid a\}] \geqq P_e(-a; a); \qquad a \geqq 0. \qquad (16)$$

By equations (13), (14), and (16) we get equations (2a) and (2b).

Now, if

$$\int_{-A + |a| \leqq \hat{a} - a \leqq A - |a|} (\hat{a} - a) \, dp(\hat{a} \mid a) = 0$$

then

$$b(a) = \int_{2a + A}^{A} (\hat{a} - a) \, dp(\hat{a} \mid a) \geqq 0; \qquad -A \leqq a \leqq 0. \qquad (17a)$$

Hence

$$b(a) > -a \, \text{Pr } [\hat{a} > 0 \mid a]; \qquad -A \leqq a \leqq -\frac{A}{2} \qquad (17b)$$

$$b(a) < (A - a) \, \text{Pr } [\hat{a} > 0 \mid a]; \qquad -\frac{A}{2} \leqq a \leqq 0. \qquad (17c)$$

Also

$$b(a) = \int_{-A}^{2a - A} (\hat{a} - a) \, dp(\hat{a} \mid a) \leqq 0; \qquad A \geqq a \geqq 0. \qquad (17d)$$

Hence

$$b(a) < -a \, \text{Pr } [\hat{a} \leqq 0 \mid a]; \qquad \frac{A}{2} \leqq a \leqq A \qquad (17e)$$

$$b(a) > -(A + a) \, \text{Pr } [\hat{a} \leqq 0 \mid a]; \qquad 0 \leqq a \leqq \frac{A}{2}. \qquad (17f)$$

Equation (5) follows from equations (17) and (11).

### III. THE SHARPNESS OF THE BOUNDS

In order to check the sharpness of the bounds on $b(a)$, let us discuss the following example.

Let $\hat{a}$ be some estimation of the parameter $a$.
Let $\hat{\hat{a}}$ be defined as

$$\hat{\hat{a}} = A \qquad \text{if} \qquad \hat{a} > 0$$
$$\hat{\hat{a}} = -A \qquad \text{if} \qquad \hat{a} \leqq 0.$$

Now, regard $\hat{\hat{a}}$ as an estimation of $a$. The bias of $\hat{\hat{a}}$ is given by

$$b(a) = (A - a) \Pr [\hat{a} > 0 \mid a] + (-A - a)[1 - \Pr (\hat{a} > 0 \mid a]$$

$$= 2A \Pr \{\hat{a} > 0 \mid a\} - (A + a); \qquad a \leqq 0 \tag{18}$$

and also

$$b(a) = -2A \Pr \{\hat{a} \leqq 0 \mid a\} + (A - a); \qquad a > 0. \tag{19}$$

Compare equations (18) and (19) with equations (13) and (14).

In the special case where $\hat{a}$ is a maximum likelihood estimator and the channel is the one given by equation (7), we have that

$$\Pr \{\hat{a} \leqq 0 \mid a\} = \Pr \{\hat{a} > 0 \mid - a\}$$

$$= P_e\{-a; a\}; \qquad a \geqq 0. \tag{20}$$

Inserting equation (20) into equations (18) and (19) yields equation (6a) and (6b). By making $\hat{\hat{a}} = -A$ if $\hat{\hat{a}} > 0$ and $\hat{\hat{a}} = A$ if $\hat{a} \leqq 0$ we get equations (6c) and (6d) in a similar way.

## IV. APPLICATIONS

### 4.1 *Postdetection Integration*

Assume that one makes $n$ independent, equally distributed, estimations of $a$: $\hat{a}_1, \hat{a}_2, \hat{a}_3, \cdots, \hat{a}_i, \cdots, \hat{a}_n$, and let

$$\tilde{a} = \frac{1}{n} \sum_{i=1}^{n} \hat{a}_i \; ;$$

$\tilde{a}$ is sometimes called the "postdetection estimation of $a$". Such an estimator appears in many applications: radar range estimation, postdetection diversity combiners in communication systems, and so on. Now

$$\epsilon_a^2 \triangleq E\{(\tilde{a} - a)^2 \mid a\} = E[(\tilde{a} - b(a) - a)^2 \mid a] + b^2(a)$$

$$= \frac{1}{n} \sigma_a^2 + b^2(a)$$

where $\sigma_a^2$ is the variance of $\hat{a}_i$ (for any $i$), given $a$. Clearly, if the estimator $\hat{a}$ is unbiased, the mean square error that is associated with $\tilde{a}$ can be made arbitrarily small by making $n$ large enough. However, if $\hat{a}$ is biased, then, for any $n$, the mean square error is lower bounded by

$$\epsilon_a^2 \geqq b^2(a) \geqq b_L^2(a) \tag{21}$$

where $b_L(a)$ is the lower bound on $\mid b(a) \mid$ given by equations (2), (4), or (5).

*Example:* Let the channel be given by equation (7) and let $\hat{a}$ be

a maximum likelihood estimator; then by equations (9a), (9c), and (21)

$$\epsilon_a^2 \geq a^2 \frac{1}{2\pi} \left[ \int_{(E/N_o)^{\frac{1}{2}}}^{\infty} \exp\left(-x^2/2\right) dx \right]^2 ; \qquad A \geq |a| \geq \frac{A}{2}.$$

Assume that the *a priori* range of $a$, is smaller than $(-A, A)$. Then

$$\epsilon_a^2 \geq \frac{a^2}{2\pi} \left[ \int_{(E/N_o)^{\frac{1}{2}}}^{\infty} \exp\left(-x^2/2\right) dx \right]^2 ; \qquad \frac{A}{2} \leq |a| \leq \alpha.$$

Now, let

$$\hat{\epsilon}^2 \triangleq \max_a \epsilon_a^2 .$$

Then, unless $A \geq 2\alpha$ (that is, unless the range of $\hat{a}$ is at least twice as large as the *a priori* range of $a$), we have that

$$\hat{\epsilon}^2 \geq \frac{\alpha^2}{2\pi} \left[ \int_{(E/N_o)^{\frac{1}{2}}}^{\infty} \exp\left(-x^2/2\right) dx \right]^2$$

even if $n \to \infty$.

### 4.2 *Predetection Integration*

Let the channel be the one given by equation (7). Assume that the estimation is based on $n$ repeated measurements; namely, the received signal is given by

$$r(t) = n(t) + \sum_{i=0}^{n-1} s(t - a - i2A).$$

In this case, an estimation is being made only after observing the complete received signal ("predetection integration"). It then follows that for a maximum likelihood estimation of $a$

$$b^2(a) \geq a^2 \frac{1}{2\pi} \left[ \int_{(nE/N_o)^{\frac{1}{2}}}^{\infty} \exp\left(-x^2/2\right) dx \right]^2 ; \qquad \frac{A}{2} \leq |a| \leq A$$

which is the same as for single measurement except for $E$ being replaced by $nE$. In this case, unlike the previous case, the lower bound vanishes as $n$ tends to infinity.

REFERENCES

1. Van Trees, H. L., "Detection, Estimation and Modulation Theory," Part 1, New York: John Wiley and Sons, 1968, p. 53.
2. Viterbi, A., "Principles of Coherent communication," New York: McGraw-Hill, 1966, Chapter 7.
3. Ziv, J., and Zakai, M., "Some Lower Bounds on Signal Parameter Estimation," IEEE Trans. Information Theory, *IT-15*, No. 3 (May 1969), pp. 386–391.

# Solutions of Fokker-Planck Equation with Applications in Nonlinear Random Vibration

## By S. C. LIU

*In the course of analyzing the dynamic behavior of mechanical systems subjected to random excitations, we investigate the associated Fokker–Planck equation. We also discuss the relationship between the characteristics of the random excitation and the nonlinear intensity coefficients governed by the physical properties of the system. This relationship leads to some simplified methods for solving the response probability density of certain nonlinear systems. We present general solutions to a class of multidimensional problems with desirable constraints. The random motion of a single-mode mechanical oscillator with various nonlinear stiffnesses and a charged particle moving in an electromagnetic field are examples. Cosine-power and sech-power distributions are found to be associated with the steady state response of a tangent stiffness system and a hyperbolic tangent stiffness system, respectively. When the influencing magnetic vector potential M is irrotational, the stationary probability for the moving particle in the 6-dimensional response phase-space is statistically independent.*

## I. INTRODUCTION

Although the theory of stochastic processes has found wide applications in information and communication sciences for many years, only recent advances in rocket propulsion and aerospace industries have made random vibration problems subjects of growing importance in mechanical and civil engineering. These problems involve structural responses due to random loadings and are in general nonlinear resulting from large motions.[1] Such nonlinear random transformation problems often encountered in practice are generally memory-dependent; that is, the equations of motion are described by nonlinear differential equations.[2-4] Under the Markov and Smoluchowski assumptions, it has been shown that the probability density function

of a large class of random processes satisfies equations of the Fokker-Planck (F-P) type.[5-9] Recently Pawula showed that generalized Fokker-Planck equations can be derived for many cases with both these assumptions removed.[7] Many interesting problems with their governing equations of the Fokker-Planck type have been investigated by various researchers. Rosenbluth, and others, studied the Fokker-Planck equation for the distribution function for gases with an inverse-square particle interaction force;[10] van Kampen used an Fokker-Planck equations to describe the thermal fluctuations in linear and nonlinear systems;[11] Ariaratnam found the steady state response distribution for a class of nonlinear two-mode mechanical oscillators by applying certain constraints to decouple the governing Fokker-Planck equation;[12] and Hempstead and Lax used Fourier transform techniques to eliminate the phase variable from the Fokker-Planck equation in the polar coordinates for a rotating-wave Van der pol oscillator and found the phase and amplitude spectra.[13]

The Fokker-Planck equation, satisfied by the random response probability density function of a dynamic system, is a parabolic partial differential equation which generally is rather difficult to solve. Although approximate results may be obtained by using the perturbation and equivalent linearization tehcniques (for which brief accounts are given in Appendixes A and B), the formal solution yielded by the appropriate Fokker-Planck equation is still the most sought one.[14-17] In this paper, we investigate the relationship between the random input characteristics and the intensity coefficients of the response process as governed by the physical properties of the system. Based on this relationship, we establish theorems concerning classes of potential-type and uncoupled-type solutions to the multidimensional stationary Fokker-Planck equations. Then, we present simple methods, based on these theorems, for solving such classes of random transformations and we describe the required natural restraints which justify applying these methods on physical grounds.

As examples we analyze two separate cases, a simple mechanical oscillator with various nonlinear spring resistances and a charged particle moving in an electromagnetic field, which we also consider to be subject to random excitation.

## II. MARKOV PROCESSES AND RANDOMLY EXCITED DYNAMIC SYSTEMS

A stochastic process which has no aftereffect is called a Markov process. If $\mathbf{y}(t) = [y_1(t), y_2(y), \cdots, y_N(t)]$ is such a process (where

bold type indicates a vector) we can write

$$p(\mathbf{y}_1 \mid \mathbf{y}_2, \cdots, \mathbf{y}_N) = p(\mathbf{y}_1 \mid \mathbf{y}_2), \qquad N > 2 \tag{1}$$

and

$$p(\mathbf{y}_1, \mathbf{y}_2, \cdots, \mathbf{y}_N) = p(\mathbf{y}_1 \mid \mathbf{y}_2)p(\mathbf{y}_2 \mid \mathbf{y}_3) \cdots p(\mathbf{y}_{N-1} \mid \mathbf{y}_N)p(\mathbf{y}_N) \tag{2}$$

where $p(\mid)$ and $p(\;)$ represent the conditional probability density and the joint probability density, respectively; and $\mathbf{y}_1 = \mathbf{y}(t_1), \cdots, \mathbf{y}_N = \mathbf{y}(t_N)$.

In examining the motion of a dynamic system under purely random disturbance, it is found that the phase point $\mathbf{y}(t_2)$ of the system at time $t_2$ depends only on the phase-point position $\mathbf{y}(t_1)$ at the previous time $t_1$. Therefore, the trajectory of the phase-point of a system under purely random disturbance is described by a Markov process $\mathbf{y}(t) = [y_1(t), y_2(t), \cdots, y_N(t)]$ in the phase space. The components $y_i$, $i = 1, 3, 5, \cdots, N - 1$ represent the generalized coordinates of the system and components $y_{i+1}$ represent the first time derivatives of $y_i$. These $N$ variables completely defined the dynamic state of a viboratory system in the $N$-dimensional phase-space.

If such a Markov process $\mathbf{y}(t) = [y_1(t), y_2(t), \cdots, y_N(t)]$ satisfies the following conditions

(i) the Smoluchowski equation

$$p(\mathbf{y}_1 \mid \mathbf{y}_2, \triangle t) = \int d\mathbf{y}_3 \, p(\mathbf{y}_1 \mid \mathbf{y}_3, t + \triangle t)p(\mathbf{y}_3 \mid \mathbf{y}_2, t) \tag{3}$$

holds for every $\mathbf{y}_1$, $\mathbf{y}_2$ and $\mathbf{y}_3$ defined in the $N$-dimensional phase space,

(ii) the higher order intensity coefficients vanish, that is,

$$K_s(\mathbf{y}) = \lim_{\triangle t \to 0} \frac{1}{\triangle t} \langle (\mathbf{y}_{\triangle t} - \mathbf{y})^s \rangle = 0 \quad \text{for} \quad s \geqq 3 \tag{4}$$

and the first and second intensity coefficients

$$K_1(\mathbf{y}) = \lim_{\triangle t \to 0} \frac{1}{\triangle t} \langle y_{i,\triangle t} - y_i \rangle = A_i(\mathbf{y}, t) \tag{5}$$

$$K_2(\mathbf{y}) = \lim_{\triangle t \to 0} \frac{1}{\triangle t} \langle (y_{i,\triangle t} - y_i)(y_{j,\triangle t} - y_j) \rangle = B_{ij}(\mathbf{y}, t) \tag{6}$$

exist, where the symbol $\langle\;\rangle$ indicates ensemble average,*

(iii) $A_i$, $B_{ij}$ are continuous and bounded,

(iv) $\partial A_i / \partial y_i$ exist for every $y_i$, and are continuous and bounded,

---

* A Markov process which satisfies condition (ii) is said to be continuous.

(v) $\partial B_{ij}/\partial y_i$ and $\partial^2 B_{ij}/\partial y_i \partial y_j$ exist for every $y_i$ and $y_j$ , and are continuous and bounded, and

(vi) $B_{ij}(\mathbf{y}, t)$ form a positive definite matrix,

then the conditional probability $p(\mathbf{y}_0 \mid \mathbf{y}, t)$ of the process $y(t)$ satisfies the Fokker–Planck equation

$$\frac{\partial p}{\partial t} = - \sum_{i=1}^{N} \frac{\partial}{\partial y_i} [A_i p] + \frac{1}{2} \sum_{i,j=1}^{N} \frac{\partial^2}{\partial y_i \, \partial y_j} [B_{ij} p] \tag{7}$$

and the initial condition

$$p(\mathbf{y}_0 \mid \mathbf{y}, t) = \prod_{i=1}^{N} \delta(y_i - y_{i0}) \quad \text{as} \quad t \to 0, \tag{8}$$

where $\mathbf{y}_0$ is the initial state of $\mathbf{y}$ and $\delta$ represents the Dirac delta function.

Consider a class of $n$-degree-of-freedom nonlinear discrete dynamic systems whose motions are defined by the following system of differential equations

$$\ddot{x}_i + a_i \dot{x}_i [1 + \epsilon_i D_i(x_1 , x_2 , \cdots , x_n , \dot{x}_1 , \dot{x}_2 , \cdots , \dot{x}_n)]$$
$$+ b_i x_i [1 + \mu_i S_i(x_1 , x_2 , \cdots , \dot{x}_n , \dot{x}_1 , \cdots , \dot{x}_n)] = \eta_i(t)$$
$$i = 1, 2, \cdots , n \tag{9}$$

where $a_i$ and $b_i$ are linear damping and stiffness coefficients, respectively; $\epsilon_i$ and $\mu_i$ are nonlinear parameters; $D_i$ and $S_i$ are nonlinear functions; and $\eta_i(t)$ are excitations and are, in general, random processes.

It is always convenient to transform the motion in $n$-dimensional generalized coordinates space into a $2n$-dimensional phase-space, that is, let $y_i = x_i$ and $y_{i+1} = \dot{x}_i$ . Equation (9) can be written in a system of $2n$ first order differential equations as

$$\dot{y}_i = f_i(y_1 , y_2 , \cdots , y_{2n}) + \eta_i(t) \quad i = 1, 2, \cdots , 2n. \tag{10}$$

Assuming $\langle \eta_i(t) \rangle = 0$ and $\langle \eta_i(t_1) \eta_j(t_2) \rangle = S_{ij} \delta_{ij}(t_1 - t_2)$ for constant $D_{ij}$ and applying equations (4) through (6) to equation (9), we obtain

$$A_i(\mathbf{y}, t) = y_{i+1} ,$$

$$A_{i+1}(\mathbf{y}, t) = -a_i y_{i+1} [1 + \epsilon_i D_i(y_1 , \cdots , y_{2n})]$$
$$- b_i y_i [1 + \mu_i S_i(y_1 , \cdots , y_{2n})],$$

$$B_{ii}(\mathbf{y}, t) = 0,$$

$$B_{ij}(\mathbf{y}, t) = 0,$$

and

$$B_{i+1,i+1}(\mathbf{y}, t) = \lim_{\Delta t \to 0} \frac{1}{\Delta t} \left\langle \left\{ \int_t^{t+\Delta t} \eta_i(\tau) \, d\tau \right. \right.$$

$$\left. \left. - a_i y_{i+1}[1 + \epsilon_i D_i] \, \Delta t - b_i y_i[1 + \mu_i S_i] \, \Delta t \right\}^2 \right\rangle$$

$$= \lim_{\Delta t \to 0} \frac{1}{\Delta t} \left\langle \int_t^{t+\Delta t} \int_t^{t+\Delta t} \eta_i(\tau_1) \eta_i(\tau_2) \, d\tau_1 \, d\tau_2 + O(\Delta t^2) \right\rangle$$

$$= S_{ii} \, .$$

From the above results we notice that the coefficients $A_{i+1}$ are determined by the specific nonlinear damping and nonlinear spring functions while coefficients $B_{ij}$ depend only upon the statistical properties of the random excitation functions. Therefore, we can conclude that for the random response components $y_i(t)$ in equation (10) if the limits of their first and second increments

$$\lim_{\Delta t \to 0} \frac{1}{\Delta t} \langle \Delta y_i(t) \rangle \quad \text{and} \quad \lim_{\Delta t \to 0} \frac{1}{\Delta t} \langle \Delta y_i(t) \, \Delta y_j(t) \rangle$$

exist, then the probability density function of the many random functions $y_1, \cdots, y_{2n}$ satisfies the Fokker-Planck equation (7).

Equation (7) can be written as

$$\dot{p} = - \sum_{i=1}^{2n} \frac{\partial G_i(\mathbf{y})}{\partial y_i} \tag{11}$$

where

$$G_i(\mathbf{y}) = A_i(\mathbf{y})p - \frac{1}{2} \sum_{j=1}^{2n} \frac{\partial}{\partial y_j} [B_{ij}(\mathbf{y})p] \tag{12}$$

describe the components of a probability current vector

$$\mathbf{G} = (G_1, G_2, \cdots, G_{2n}).$$

Because the general solution to the above multidimensional parabolic partial differential equation with arbitrary boundary condition is difficult and impracticable to find, no such attempt is made in this study. However, with certain constraints on the properties of the system as well as on the characteristics of the random input, equation (11) may be reduced to simpler and more readily solvable forms. Such forms are considered in Section III.

III. STATIONARY SOLUTION OF JOINT PROBABILITY DISTRIBUTION

In the limiting case of equation (11) when $t \to \infty$ the transition probability density $p(\mathbf{y}_0 \mid \mathbf{y})$ tends to a stationary joint probability density $p_{st}(\mathbf{y}) = p_{st}(y_1, y_2, \cdots, y_{2n})$ which is independent of the initial conditions and approaches a stationary (steady state) value as the time of passage is sufficiently large. Under this stationary situation, equation (11) becomes

$$\sum_{i=1}^{2n} \frac{\partial G_i(\mathbf{y})}{\partial y_i} = 0 \tag{13}$$

or

$$\overline{\nabla}_{y_i} \cdot \mathbf{G} = 0, \qquad i = 1, 2, \cdots, 2n. \tag{14}$$

Therefore $\mathbf{G}$ may be regarded as being incompressible and there are no sources or sinks in the region $R$. Notice that since rotational probability flows can occur even for cases of zero boundary conditions, that is, the probability current $\mathbf{G}(G_1, G_2, \cdots, G_{2n})$ satisfies

$$G_i(\mathbf{y}) = 0, \qquad i = 1, 2, \cdots, 2n \tag{15}$$

on the boundary of the region $R$ under consideration, $\mathbf{G}$ may not vanish within $R$. In a special case, however, discussed in Section 3.1, the current vector $\mathbf{G}$ vanishes in the whole region $R$.

Two situations which readily yield steady state solutions to equation (10) are investigated in Section 3.1.

### 3.1 *Potential Distribution*

Under the assumption that the probability current vector $\mathbf{G}$ vanishes everywhere in $R$, a solution of the potential form $\exp(-U)$ can be constructed where $U$ is a Liapunov type potential function of the system.

If the equations of motion of a dynamic system satisfy the following conditions:

(*i*)  $G_i = $ for all $i$,

(*ii*)  the matrix $[B_{ij}]$ is not singular,

(*iii*)  $\dfrac{\partial}{\partial y_\alpha} \sum_i D_{\beta i} \left( \sum_j \dfrac{\partial B_{ij}}{\partial y_j} - 2A_i \right)$

$$= \frac{\partial}{\partial y_\beta} \sum_i D_{\alpha i} \left( \sum_j \frac{\partial B_{ij}}{\partial y_j} - 2A_i \right), \qquad \alpha, \beta = 1, 2, \cdots, 2n$$

and

(*iv*) there is no probability flow through the boundary of $R$,

then the solution to the steady state Fokker-Planck equation (13) is

$$p_{st}(\mathbf{y}) = C \exp \left\{ - \int_{a_1, a_2, \cdots, a_{2n}}^{y_i, y_2, \cdots, y_{2n}} \sum_{\alpha} \left[ \sum_{i,j} D_{\alpha i} \frac{\partial B_{ij}}{\partial y_j} - 2 \sum_i D_{\alpha i} A_i \right] dy_\alpha \right\}$$

(16)

where $[D_{ai}] = [B_{ij}]^{-1}$, that is,

$$\sum_i D_{\alpha i} B_{ij} = \delta_{\alpha j} ,$$

and $C$ is the normalization constant determined by

$$\int_{2n\text{-fold}} \cdots \int p_{st} \, dy_1 , \cdots , dy_{2n} = 1.$$

(17)

Equation (16) can be easily verified by a direct substitution of $p(\mathbf{y}) = \exp(-U)$ into equation (13) to solve for $U$. A special case of interest is when the matrix $[B_{ij}]$ is isotropic, that is, when

$$[B_{ij}] = B(\mathbf{y})[\delta_{ij}].$$

(18)

The solution for such a case, which can be obtained by substituting equation (18) into (16), is

$$p_{st}(\mathbf{y}) = C \exp(-U)$$

(19)

$$U = \log \frac{1}{B(\mathbf{y})} - 2 \int_{a_1, \cdots, a_{2n}}^{y_i, \cdots, y_{2n}} \sum_i^{2n} \frac{A_i}{B(\mathbf{y})} \, dy_i$$

where $C$ is the usual normalization factor.

### 3.2 *Uncoupled Distribution*

There are cases when all generalized response variables of a system in the $2n$ phase-space coordinates are independent of one another. There are two approaches by which the solutions can be achieved.

### 3.2.1 *Forward Approach*

The governing Fokker–Planck equation may be reduced to a form in which the final probability distribution function clearly shows a statistical independent character. For this type of motion it is sometimes possible to find appropriate partial operators which, when linearly operated on functions of the type $g_i(y_i)p + h_i(y_i)(\partial p/\partial y_i)$, generate an equivalent nonlinear partial differential equation. If equation (13)

can be put in the form

$$\sum_{i=1}^{2n} L_i \left[ g_i(y_i)p + h_i(y_i) \frac{\partial p}{\partial y_i} \right] = 0 \tag{20}$$

where the coefficients $L_i$ are arbitrary first order partial differential operators, and if there exists a $p(\mathbf{y})$ satisfying each

$$g_i(y_i)p + h_i(y_i) \frac{\partial p}{\partial y_i} = 0, \qquad (i = 1, 2, \cdots, 2n),$$

then by Gray's uniqueness theorem such $p(\mathbf{y})$ is the unique solution of equation (20).[18] Such a solution is

$$p_{st}(\mathbf{y}) = C \prod_{i=1}^{2n} \exp \left[ - \int_0^{y_i} \frac{g_i(\lambda_i)}{h_i(\lambda_i)} \, d\lambda_i \right] \tag{21}$$

and $C$ is the normalization factor.

Notice that previous investigators such as Ariaratnam and Klein had their problems satisfying equation (20) and therefore obtained their solutions in a form similar to equation (21).

### 3.2.2 *Backward Approach*

Sometimes it is more convenient to work backward. By this procedure, the statistically independent property is assumed in order to derive the desirable solution for the Fokker-Planck equation under investigation. This method gives a close insight into the physical properties of the system and enables the natural boundary conditions to be deduced. These deduced boundary conditions provide the necessary constraints for randomly excited systems which have independent distributions in their response variables.

If equation (13) satisfies the conditions                                    (22)

(*i*) $[B_{ij}]$ is a constant diagonal matrix for $i, j = 1, 2, \cdots, 2n$, and,

(*ii*) the first order intensity coefficients $A_i$ are functions of $y_i$ and $y_{i+1}$ only, and there are no cross terms in $A_i$, that is,

$$\left. \begin{array}{l} A_i = A_{i,i}(y_i) + A_{i,i+1}(y_{i+1}) \\ A_{i+1} = A_{i+1,i}(y_i) + A_{i+1,i+1}(y_{i+1}) \end{array} \right\}$$

$$\text{for} \quad i = 1, 3, \cdots, 2n - 1; \tag{23}$$

then by setting $p_{st}(y) = \prod_{i=1}^{2n} p_i(y_i)$ in equation (13), we obtain

$$\frac{1}{2} \sum_{i=1}^{2n} B_{ii} \frac{p(\mathbf{y})}{p_i(y_i)} \frac{\partial^2}{\partial y_i^2} p_i(y_i)$$

$$= \sum_{i=1}^{2n} \left[ \frac{p(\mathbf{y})}{p_i(y_i)} (A_{i,i} + A_{i,i+1}) \frac{\partial p_i(y_i)}{\partial y_i} + p(\mathbf{y}) \frac{\partial A_{i,i}(y_i)}{\partial y_i} \right].$$

Dividing the above equation by $p(\mathbf{y})$, we obtain

$$\frac{1}{2} \sum_{i=1}^{2n} B_{ii} \frac{p_i''}{p_i} = \sum_{i=1}^{2n} \left[ (A_{i,i} + A_{i,i+1}) \frac{p_i'}{p_i} + A_{i,i}' \right]$$

where the $'$ denotes $\partial/\partial y_i$.

The above equation can be reorganized as

$$\sum_{i=1}^{2n} \left[ \frac{B_{ii}}{2} \frac{p_i''}{p_i} - \frac{p_i'}{p_i} A_{i,i} - A_{i,i}' \right]$$

$$= \sum_{i=1,3,\cdots}^{2n-1} \left[ \frac{p_i'}{p_i} A_{i,i+1} + \frac{p_{i+1}'}{p_{i+1}} A_{i+1,i} \right].$$

A sufficient solution $p(\mathbf{y})$ for the above equation requires that it satisfies the following relations

(i)    $\dfrac{B_{ii}}{2} \dfrac{p_i''}{p_i} - \dfrac{p_i'}{p_i} A_{i,i} - A_{i,i}' = 0 \quad \text{for} \quad i = 1, 2, \cdots, 2n$      (24)

and

(ii)    $\dfrac{p_i'}{p_i} \dfrac{1}{A_{i+1,i}} = -\dfrac{p_{i+1}'}{p_{i+1}} \dfrac{1}{A_{i,i+1}} = E = \text{constant}$

$$\text{for} \quad i = 1, 3, \cdots, 2n - 1. \quad (25)$$

Equation (24) can be reduced to

$$\frac{B_{ii}}{2} \frac{d}{dy_i} (p_i') = \frac{d}{dy_i} (p_i A_{i,i}).$$

Integration with respect to $y_i$ yields

$$\frac{B_{ii}}{2} p_i' = p_i A_{i,i} + c_1 .$$

When compared with equation (25), $c_1$ vanishes and the following condition must hold

$$\frac{2A_{i,i}}{A_{i+1,i} B_{i,i}} = -\frac{2A_{i+1,i+1}}{A_{i,i+1} B_{i+1,i+1}} = E, \quad (i = 1, 3, \cdots, 2n - 1). \quad (26)$$

Equation (26) is the natural restraint under which the backward method can be applied. From the above analysis and the uniqueness theorem stated, we can claim. If the steady state Fokker–Planck equation satisfies conditions (22), (23), and (26), and if the phase-space coordinates $y_1$, $y_2$, $\cdots$, $y_{2n}$ are statistically independent, then the

unique solution $p_{st}(\mathbf{y})$ of equation (13) is

$$p_{st}(\mathbf{y}) = C \prod_{i=1,3,\cdots}^{2n-1} \exp\left[\int_0^{y_i} EA_{i+1,i}(\lambda_i) \, d\lambda_i \right.$$
$$\left. - \int_0^{y_{i+1}} EA_{i,i+1}(\lambda_{i+1}) \, d\lambda_{i+1}\right] \qquad (27)$$

where $C$ is the normalized constant.

## IV. EXAMPLES

### 4.1 *Randomly Excited Nonlinear Simple Mechanical Oscillator*

When subjected to dynamic loadings $\eta(t)$ the equation of motion for a single-mode oscillator with nonlinear spring $F(x)$ is

$$\ddot{x} + \beta\dot{x} + F(x) = \eta(t) \qquad (28)$$

where the excitation $\eta(t)$ is a gaussian stationary process with

$$\left.\begin{aligned}
\langle\eta(t)\rangle &= 0 \\
\langle\eta(t_1)\eta(t_2)\rangle &= 2\pi S_o \, \delta(t_1 - t_2)
\end{aligned}\right\} \qquad (29)$$

in which $S_o$ is a constant power spectral density.

Letting $y_1 = x$ and $y_2 = \dot{x} = \dot{y}_1$, the intensity coefficients $A_i$ and $B_{ij}$ can be found by using equations (4) and (5) as follows

$$A_1 = A_{11} + A_{12} = y_2, \quad \text{hence} \quad A_{11} = 0, \quad A_{12} = y_2;$$

$$A_2 = A_{21} + A_{22} = -F(y_1) - \beta y_2,$$

$$\text{hence} \quad A_{21} = -F(y_1), \quad A_{22} = -\beta y_2;$$

$$B_{22} = 2\pi S_o.$$

Therefore the governing Fokker-Planck equation is

$$\pi S_o \frac{\partial^2 p}{\partial y_2^2} - \frac{\partial}{\partial y_1}(y_2 p) + \frac{\partial}{\partial y_2}\{[F(y_1) + \beta y_2]p\} = 0. \qquad (30)$$

Forward approach:

Equation 30 can be written in the form of equation (20) as

$$\frac{\partial}{\partial y_2}\left[F(y_1)p + \frac{\pi S_o}{\beta}\frac{\partial p}{\partial y_1}\right] + \left[\beta\frac{\partial}{\partial y_2} - \frac{\partial}{\partial y_1}\right]\left[y_2 p + \frac{\pi S_o}{\beta}\frac{\partial p}{\partial y_2}\right] = 0,$$

from which we see that

$$L_1 = \frac{\partial}{\partial y_2}, \qquad L_2 = \frac{\partial}{\partial y_2} - \frac{\partial}{\partial y_1},$$

$$g_1(y_1) = F(y_1), \qquad h_1(y_1) = \frac{\pi S_o}{\beta},$$

and

$$g_2(y_2) = y_2, \qquad h_2(y_2) = \frac{\pi S_o}{\beta}.$$

Substitution of $g_1$, $g_2$, $h_1$ and $h_2$ into equation (21) gives the following steady state solution for the system described by equation (28)

$$p_{st}(\mathbf{y}) = C \exp\left\{ -\int_0^{y_1} \frac{\beta F(\zeta_1)}{\pi S_o} \, d\zeta_1 - \int_0^{y_2} \frac{\beta \zeta_2}{\pi S_o} \, d\zeta_2 \right\}$$

$$= C \exp\left\{ -\frac{\beta}{\pi S_o} \left[ \int_0^{y_1} F(\zeta_1) \, d\zeta_1 + \frac{y_2^2}{2} \right] \right\}. \tag{31}$$

Backward approach:

In this two-dimensional case, it can be shown that condition (26) is satisfied:

$$E = \frac{2A_{11}}{B_{11}A_{21}} = \frac{-2A_{22}}{B_{22}A_{12}} = \frac{-2(-\beta y_2)}{2\pi S_o y_2} = \frac{\beta}{\pi S_o}.$$

Therefore, according equations (22) through (27) the stationary solution can be written:

$$p_{st}(y_1, y_2) = C \exp\left[ \int_0^{y_1} E A_{21}(\zeta_1) \, d\zeta_1 - \int_0^{y_2} E A_{12}(\zeta_2) \, d\zeta_2 \right]$$

$$= C \exp\left\{ -\frac{\beta}{\pi S_o} \left[ \int_0^{y_1} F(\zeta_1) \, d\zeta_1 + \frac{y_2^2}{2} \right] \right\}$$

which is the same as equation (31).

### 4.1.1 Cubic Stiffness Characteristics

Let us consider various nonlinear spring resistance functions $F(x)$. The group of cubic stiffness characteristics is the classical case concerned with the hardening spring type of nonlinearity, which is represented by

$$F(x) = k_o x + \epsilon x^3, \tag{32}$$

where $k_o$ is the initial linear stiffness and $\epsilon$ is the nonlinear coefficient.

Substituting equation (32) into (31), we obtain

$$p_{st}(x, \dot{x}) = C \exp\left[ -\frac{\beta}{\pi S_o} \left( \frac{k_o x^2}{2} + \frac{\epsilon x^4}{4} + \frac{\dot{x}^2}{2} \right) \right]. \qquad (33)$$

Notice from equation (33) that the marginal probability density disitrbutions for $x$ and $\dot{x}$ are statistically independent and that $p(\dot{x})$ follows the gaussian law. As $\epsilon \rightarrow 0$, the system becomes linear and $p(x)$ approaches gaussian.

### 4.1.2 Tangent Stiffness Characteristics[19]

The spring function $F(x)$ in this case is shown in Fig. 1 and is represented by

$$F(x) = \left( \frac{2k_o d}{\pi m} \right) \tan \left( \frac{\pi x}{2d} \right), \qquad -d < x < d \qquad (34)$$

in which $m$ and $d$ are constants. Notice that $d$ is the threshold of the oscillator and $m$ can be regarded as the mass of the oscillator. Again it follows from equation (31) that

$$p_{st}(x, \dot{x}) = p(\dot{x})p(x)$$

$$= C\left[ \exp \left( -\frac{\dot{x}^2}{2\sigma_o^2 \omega_o^2} \right) \exp \left( \frac{4d^2}{\pi^2 \sigma_o^2} \ln \cos \frac{\pi x}{2d} \right) \right], \qquad (35)$$
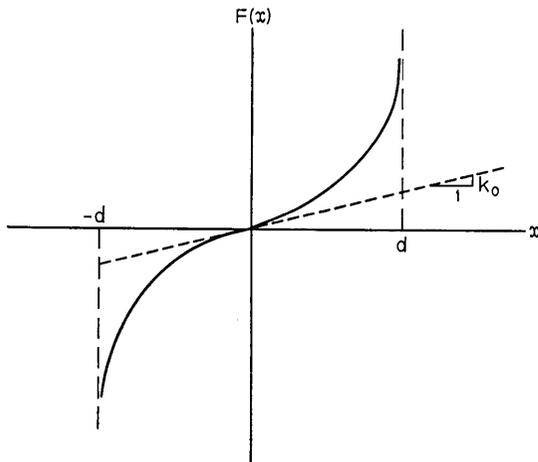


Fig. 1 — Tangent stiffness characteristics.

where $\omega_o^2 = k_o/m$ and $\sigma_o^2 = \pi S_o/2\beta\omega_o^2$ is the corresponding linear mean square response [that is, if $F(x) = k_o x$].

The cosine-power distribution $p(x)$ is shown in Fig. 2 for various values of $n = 4d^2/\pi^2\sigma_o^2$. It is interesting that for fixed $\sigma_o^2$, $p(x)$ approaches the gaussian distribution as $d \rightarrow \infty$ and approaches the uniform distribution as $d \rightarrow 0$.

### 4.1.3 Hyperbolic Tangent Stiffness Characteristics

Figure 3 shows a hyperbolic tangent stiffness model; the spring resistance function $F(x)$ is given by

$$F(x) = \frac{k_o}{mb}\tanh\,(bx), \qquad b,\,k_o\,,\,m > 0. \tag{36}$$

Notice that the resistance $F(x)$ developed during the vibration is bounded between $k_o/bm$ and $-k_o/bm$. Therefore $k_o/bm$ may be regarded as the yield level and $1/bm$ the corresponding yielding response; equation (36) can be used to model the familiar elastic-perfect-plastic behavior observed in many physical realities. The joint probability density function for this can be obtained from equation (31),

$$p_{st}(x,\,\dot{x}) = p(\dot{x})p(x) = C\,\exp\left[-\frac{\dot{x}^2}{2\omega_o^2\sigma_o^2} - \frac{1}{b^2\sigma_o^2}\ln\,\cosh\,(bx)\right]. \tag{37}$$

There are some interesting features about the limiting behavior of the marginal distribution $p(x)$ which is of sech-power type

$$p(x) = C_1(b)\,\mathrm{sech}^{(1/\sigma_o^2 b^2)}\,bx. \tag{38}$$

The sech-power term in equation (38) belongs to a class of distribution function $f_n(x)$ of a monotone decreasing sequence, integrable on $[-\infty,\infty]$ such that $f_n(x) \leqq 1$ for all $x$. It can be shown that

$$\lim_{n\to\infty}\left[f_n(x)\,\bigg/\,\int_{-\infty}^{\infty}f_n(x)\,dx\right] = 0;[20]\ \text{and for}\ \lim_{b\to 0}p(x) = \frac{1}{(2\pi)^{\frac{1}{2}}\,\sigma_o}\exp\left(-\frac{x^2}{2\sigma_o^2}\right),$$

the limit of the distribution function in equation (38) at zero, that is, $\lim_{b\to 0}\,p(x)$ converges positively to a normal distribution with zero mean and variance $\sigma_o^2$. The sech-power distribution $p(x)$ is shown in Fig. 4 for various values of $m' = 1/\sigma_o^2 b^2$.

### 4.2 Random Vibration of a Charged Particle Moving in an Electromagnetic Field

As a second example, we consider a particle of mass $m$ carrying charge $q$ subjected to a random loading $\mathbf{n}(t) = \eta_1(t)\hat{\imath} + \eta_2(t)\hat{\jmath} + \eta_3(t)\hat{k}$
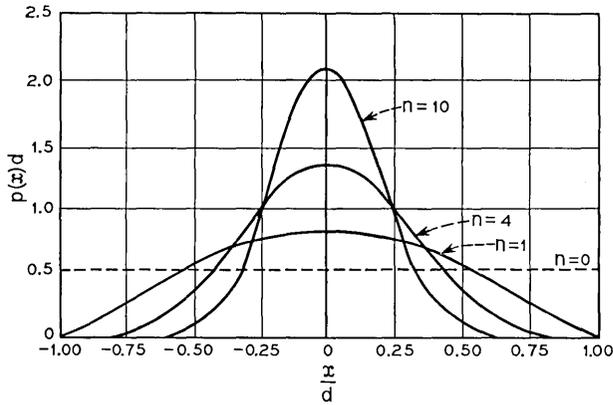
Fig. 2 — Cosine-power probability density distribution function.

where $\hat{\imath}, \hat{\jmath}, \hat{k}$ are base vectors in the Cartesian coordinates. The particle, moving in an electromagnetic field $\mathbf{M} = M_1(x, y, z)\hat{\imath} + M_2(x, y, z)\hat{\jmath} + M_3(x, y, z)\hat{k}$, is also subjected to a friction force $F_f = -\nabla_v \mathfrak{F}$ where $\mathfrak{F} = \frac{1}{2}(\lambda_1 v_x^2 + \lambda_2 v_y^2 + \lambda_3 v_z^2)$ is an energy dissipation function. The complete force on the particle is therefore,

$$\mathbf{F} = q\left\{-\nabla\phi - \frac{1}{c}\frac{\partial \mathbf{M}}{\partial t} + \frac{1}{c}(\mathbf{v} \times \nabla \times \mathbf{M})\right\} \tag{39}$$

where $\phi$ is the scalar potential and $c$ is a constant.

If we introduce a velocity-dependent potential $w$, such that

$$w = q\phi - \frac{q}{c}\mathbf{M}\cdot\mathbf{v} \tag{40}$$

then equation (39) can be written as

$$\mathbf{F} = -\nabla_d w + \frac{d}{dt}\nabla_v w. \tag{41}$$

Therefore the Lagrangian function $L$ can be expressed in terms of $w$ as

$$L = T - q\phi + \frac{q}{c}\mathbf{M}\cdot\mathbf{v} \tag{42}$$

where $T$ represents the kinetic energy of the particle, and the equation of motion of the charged particle can be derived from the Lagrange equation

$$\frac{d}{dt}\left(\frac{\partial L}{\partial \dot{q}_i}\right) - \frac{\partial L}{\partial q_i} + \frac{\partial \mathfrak{F}}{\partial \dot{q}_i} = \eta_i(t) \tag{43}$$
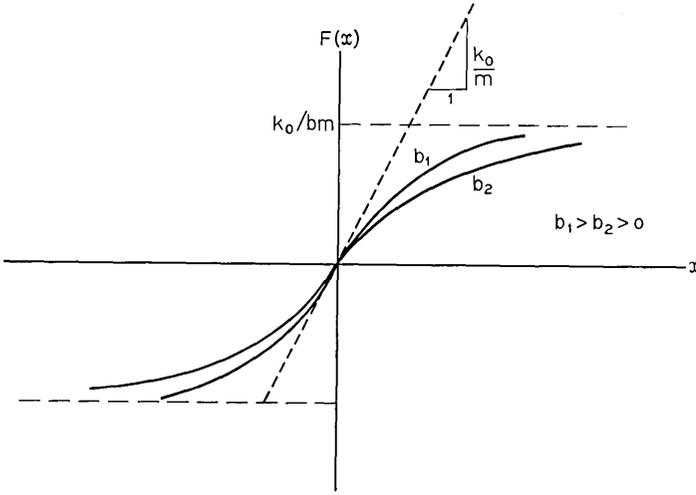
Fig. 3 — Hyperbolic tangent stiffness characteristics.

where $q_i$ represent the generalized coordinates of the motion. Using equations (42) and (43), one obtains

$$m\ddot{u}_1 + q\frac{\partial \phi}{\partial u_1} + \frac{q}{c}\left(\frac{\partial M_1}{\partial t} + \frac{\partial M_1}{\partial u_3}u_4\right.$$

$$\left. + \frac{\partial M_1}{\partial u_5}u_6 - \frac{\partial M_2}{\partial u_1}u_4 - \frac{\partial M_3}{\partial u_1}u_6\right) + \lambda_1 u_2 = \eta_1$$
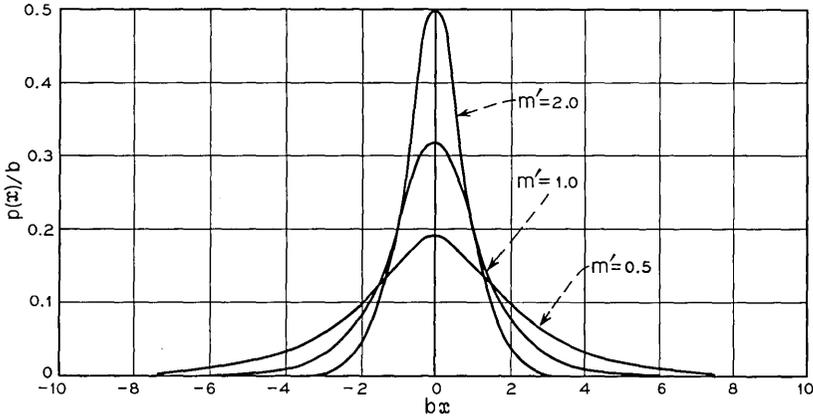


Fig. 4 — Sech-power probability density distribution function.

$$m\ddot{u}_3 + q\frac{\partial \phi}{\partial u_3} + \frac{q}{c}\left(\frac{\partial M_2}{\partial t} + \frac{\partial M_2}{\partial u_1}u_2\right.$$

$$\left. + \frac{\partial M_2}{\partial u_5}u_6 - \frac{\partial M_1}{\partial u_3}u_2 - \frac{\partial M_3}{\partial u_3}u_6\right) + \lambda_2 u_4 = \eta_2 \tag{44}$$

$$m\ddot{u}_5 + q\frac{\partial \phi}{\partial u_5} + \frac{q}{c}\left(\frac{\partial M_5}{\partial t} + \frac{\partial M_3}{\partial u_3}u_4\right.$$

$$\left. + \frac{\partial M_3}{\partial u_1}u_2 - \frac{\partial M_1}{\partial u_5}u_2 - \frac{\partial M_2}{\partial u_5}u_4\right) + \lambda_3 u_6 = \eta_3$$

where in equation (44), transformations $u_1 = x$, $u_2 = \dot{x} = \dot{u}_1$, $u_3 = y$, $u_4 = \dot{y} = \dot{u}_3$, $u_5 = z$, and $u_6 = \dot{z} = \dot{u}_5$ have been made.

Define

$$\left.\begin{aligned} D_1 &= (d_{13} - d_{21})u_4 + (d_{15} - d_{31})u_6\\ D_2 &= (d_{21} - d_{13})u_2 + (d_{25} - d_{33})u_6\\ D_3 &= (d_{31} - d_{15})u_2 + (d_{33} - d_{25})u_4 \end{aligned}\right\} \tag{45}$$

where $d_{\alpha\beta} = d_{\alpha\beta}(u_1, u_3, u_5) = \partial M_\alpha/\partial u_\beta$, $\alpha = 1, 2$, or $3$ and $\beta = 1, 3$, or $5$, are independent of the velocity variables.

The following properties are assumed for the forcing function $\mathbf{n}(t)$

$$\langle \eta_\alpha \rangle = 0 \qquad \text{for} \quad \alpha = 1, 2, \text{ or } 3 \tag{46}$$

$$\begin{aligned} \langle \eta_\alpha(t_1)\eta_\alpha(t_2)\rangle &= 0 &&\text{for}\quad \alpha \neq \beta\\ &= 2\pi S_\alpha(t_1 - t_2) &&\text{for}\quad \alpha = \beta \end{aligned}\Bigg\} \cdot \tag{47}$$

From equations (44) through (47), the six-dimensional Fokker–Planck equation governing the transition probability density $p(u_1, u_2, \cdots, u_6)$, derived by the standard technique, is

$$\dot{p}(u_1, \cdots, u_6)$$

$$= -\frac{\partial}{\partial u_1}(u_2 p) - \frac{\partial}{\partial u_2}\left\{\left[-\frac{\lambda_1}{m}u_2 - \frac{q}{m}\frac{\partial \phi}{\partial u_1} - \frac{q}{mc}\left(\frac{\partial M_1}{\partial t} + D_1\right)\right]p\right\}$$

$$- \frac{\partial}{\partial u_3}(u_4 p) - \frac{\partial}{\partial u_4}\left\{\left[-\frac{\lambda_2}{m}u_4 - \frac{q}{m}\frac{\partial \phi}{\partial u_3} - \frac{q}{mc}\left(\frac{\partial M_2}{\partial t} + D_2\right)\right]p\right\}$$

$$- \frac{\partial}{\partial u_5}(u_6 p) - \frac{\partial}{\partial u_6}\left\{\left[-\frac{\lambda_3}{m}u_6 - \frac{q}{m}\frac{\partial \phi}{\partial u_5} - \frac{q}{mc}\left(\frac{\partial M_3}{\partial t} + D_3\right)\right]p\right\}$$

$$+ \frac{\pi}{m^2}\left[S_1\frac{\partial^2 p}{\partial u_2^2} + S_2\frac{\partial^2 p}{\partial u_4^2} + S_3\frac{\partial^2 p}{\partial u_6^2}\right]. \tag{48}$$

The initial condition for equation (48) is

$$p(u_1 , u_2 , \cdots , u_6) \mid_{t_o} = \prod_{i=1}^{6} \delta[u_i(t) - u_i(t_o)]. \tag{49}$$

At $t \to \infty$, the motion becomes stationary; therefore $\dot{p} = 0$ and $p = p_{st}$ if $\partial \mathbf{M}/\partial t = 0$. Imposing these conditions on equation (48) and making use of equation (45), the corresponding steady state equation for (48) is:

$$\left(\frac{\lambda_1}{m^2} \frac{\partial}{\partial u_2} - \frac{1}{m} \frac{\partial}{\partial u_1}\right)\left(mu_2 p + \frac{\pi S_1}{\lambda_1} \frac{\partial p}{\partial u_2}\right) + \frac{1}{m} \frac{\partial}{\partial u_2} \left(q \frac{\partial \phi}{\partial u_1} p + \frac{\pi S_1}{\lambda_1} \frac{\partial p}{\partial u_1}\right)$$

$$+ \left(\frac{\lambda_2}{m^2} \frac{\partial}{\partial u_4} - \frac{1}{m} \frac{\partial}{\partial u_3}\right)\left(mu_4 p + \frac{\pi S_2}{\lambda_2} \frac{\partial p}{\partial u_4}\right) + \frac{1}{m} \frac{\partial}{\partial u_4} \left(q \frac{\partial \phi}{\partial u_3} p + \frac{\pi S_2}{\lambda_2} \frac{\partial p}{\partial u_3}\right)$$

$$+ \left(\frac{\lambda_3}{m^2} \frac{\partial}{\partial u_6} - \frac{1}{m} \frac{\partial}{\partial u_5}\right)\left(mu_6 p + \frac{\pi S_3}{\lambda_3} \frac{\partial p}{\partial u_6}\right) + \frac{1}{m} \frac{\partial}{\partial u_6} \left(q \frac{\partial \phi}{\partial u_5} p + \frac{\pi S_3}{\lambda_3} \frac{\partial p:}{\partial u_5}\right)$$

$$+ q\left(D_1 \frac{\partial p}{\partial u_2} + D_2 \frac{\partial p}{\partial u_4} + D_3 \frac{\partial p}{\partial u_6}\right) = 0. \tag{50}$$

Notice that in the above equation $p = p_{st}$. Terms in the last parentheses of equation (50) vanish if

$$(i) \quad D_1 = D_2 = D_3 = 0 \tag{51a}$$

or

$$(ii) \quad u_2 : u_4 : u_6 = (d_{25} - d_{33}) : (d_{31} - d_{15}) : (d_{13} - d_{21}) \tag{51b}$$

or

$$(iii) \quad \nabla \times \mathbf{M} = 0 \quad \text{or} \quad \mathbf{M} = G \tag{51c}$$

where $G$ is a scalar potential function of $x$, $y$ and $z$.

If any one of these conditions is satisfied, equation (50) is of the same form as equation (20) and its solution can be found immediately by using equation (21), that is:

$$p_{st}(u_1 , \cdots , u_6)$$

$$= C \exp \left[- \frac{m\lambda_1}{\pi S_1} \left(\frac{u_2^2}{2}\right) - \frac{q\lambda_1}{\pi S_1} \int_0^{u_1} \frac{\partial \phi}{\partial \zeta_1} d\zeta_1 - \frac{m\lambda_2}{\pi S_2} \left(\frac{u_4^2}{2}\right)\right.$$

$$\left.- \frac{q\lambda_2}{\pi S_2} \int_0^{u_3} \frac{\partial \phi}{\partial \zeta_3} d\zeta_3 - \frac{m\lambda_3}{\pi S_3} \left(\frac{u_6^2}{2}\right) - \frac{q\lambda_3}{\pi S_3} \int_0^{u_5} \frac{\partial \phi}{\pi \zeta_5} d\zeta_5\right], \tag{52}$$

where $C$ is the normalization factor determined by

$$\int \cdots \int p_{st} \, du_1 , \cdots , du_6 = 1.$$

If damping factors in three directions are identical and the random loading $\mathbf{n}$ is uniform, that is, if

$$\mathfrak{F} = \tfrac{1}{2}\lambda(v_x^2 + v_y^2 + v_z^2) \tag{53}$$

and

$$S_1 = S_2 = S_3 = S, \tag{54}$$

then equation (52) becomes

$$p_{st} = C \exp\left[\frac{-\lambda}{\pi S}(T + q\phi)\right]. \tag{55}$$

Applying this result to single-mode conservative oscillators with potential $V(x)$ and subjected to gaussian white random loadings, the stationary response probability density is

$$p_{st} = C \exp\left(-\frac{\lambda}{\pi S} H\right) \tag{56}$$

where $H = T + V$ is the Hamiltonian function of the system.

It is interesting that if the magnetic vector potential is irrotational, the steady state response probability density of a charged particle under white noise type random disturbances is statistically independent. The solution can be immediately written down in terms of a quadrature. Extension of this result to a conservative dynamic system shows that stationary probability solution is of the form $p_{st} = C \exp[-f(\gamma, H)]$, where $H$ is the Hamiltonian function of the system and $\gamma$ is a coefficient depending on the random input characteristics and the energy dissipation mechanism of the system.

APPENDIX A

*Perturbation Technique*

The perturbation method is based on a series expansion in powers of the nonlinearity coefficient $\epsilon$. This method is valid for small values of $\epsilon$ only. For example, consider a single-mode oscillator with non-linear damping and stiffness, when subjected to random force $\eta(t)$, the equation of motion is

$$\ddot{x} + 2\beta\dot{x}\left(1 + \epsilon \sum_n b_n \dot{x}^n\right) + \omega_o^2 x\left(1 + \epsilon \sum_n a_n x^n\right) = \eta(t). \tag{57}$$

A series solution can be assumed[*]

$$x(t) = x_o + \epsilon x_1 + \epsilon^2 x_2 + \cdots + \epsilon^n x_n.\tag{58}$$

From equations (57) and (58) it follows that

$$\ddot{x}_o + 2\beta\dot{x}_o + \omega_o^2 x_o = \eta(t)\tag{59}$$

$$\ddot{x}_1 + 2\beta\dot{x}_1 + \omega_o^2 x_1 = -\sum_n (2\beta b_n \dot{x}_o^{n+1} + \omega_o^2 a_n x_o^{n+1}).\tag{60}$$

Therefore the correlation function is

$$R_{xx}(\tau) = R_{x_o x_o}(\tau) + \epsilon R_{x_o x_1}(\tau) + \epsilon R_{x_1 x_o}(\tau)\tag{61}$$

where

$$R_{x_y} = \langle x(t)y(t + \tau)\rangle.\tag{62}$$

Equations (59) and (60) are linear and their solutions can be readily obtained

$$x_o(t) = \int_0^\infty h(\tau)\eta(t - \tau)\, d\tau,\tag{63}$$

and

$$x_1(t) = -\int_0^\infty h(\tau) \sum_n [2\beta b_n \dot{x}_o^{n+1}(t - \tau) + \omega_o^2 a_n x_o^{n+1}(t - \tau)]\, d\tau\tag{64}$$

where $h(\tau) = [e^{-\beta\tau}/\omega_o(1 - \beta^2)^{\frac{1}{2}}] \sin (1 - \beta^2)^{\frac{1}{2}}\omega_o\tau$ is the transfer function for system described by equation (59).

From equations (63) and (64), the nonlinear response moments can be found. Considering only the first-order perturbation,

$$R_{x_o x_1}(t) = -\int_0^\infty h(\tau) \sum_n [2\beta b_n \langle x_o(t)\dot{x}_o^{n+1}(t - \tau)\rangle$$

$$+ \omega_o^2 a_n \langle x_o(t)x_o^{n+1}(t - \tau)\rangle]\, d\tau,$$

which can be evaluated explicitly if $x_o(t)$ is gaussian.

APPENDIX B

*Equivalent Linearization Technique*

Consider the equation of motion

$$\ddot{x} + \beta\dot{x} + \omega_o^2 x + \epsilon g(x, \dot{x}, t) = \eta(t),\tag{65}$$

---

[*] Under fairly general conditions it can be shown that this series solution is convergent.

for a nonlinear function $g$, dependent on both $x$ and $\dot{x}$. The following equation is said to be equivalent to equation (65) in the sense that the mean square deficiency is minimized:

$$\ddot{x} + \beta_e \dot{x} + \omega_e^2 x + e(x, \dot{x}, t) = \eta(t). \tag{66}$$

The deficiency $e(x, \dot{x}, t) = (\beta - \beta_e)\dot{x} + (\omega_o^2 - \omega_e^2)x + \epsilon g(x, \dot{x}, t)$, in which $\beta_e$ and $\omega_e$ are equivalent damping and frequency is determined by,

$$\frac{\partial \langle e^2 \rangle_{\text{av}}}{\partial \beta_e} = 0 \quad \text{and} \quad \frac{\partial \langle e^2 \rangle_{\text{av}}}{\partial \omega_e^2} = 0, \tag{67}$$

where the $\langle \ \rangle_{\text{av}}$ indicates time average, that is,

$$\langle e^2 \rangle_{\text{av}} = \lim_{T \to \infty} \frac{1}{2T} \int_{-T}^{T} e^2(x, \dot{x}, t) \, dt. \tag{68}$$

Using equations (67) and (68), we obtain

$$\beta_e = \beta + \epsilon \frac{\langle \dot{x}g \rangle_{\text{av}}}{\langle \dot{x}^2 \rangle_{\text{av}}} \tag{69}$$

$$\omega_e^2 = \omega_o^2 + \epsilon \frac{\langle xg \rangle_{\text{av}}}{\langle x^2 \rangle_{\text{av}}} \tag{70}$$

From equations (69) and (70) and by neglecting the deficiency term $e$, equation (66) can be solved by using standard linear theory. If the system is nonhereditary, the time averages in equations (67) through (70) are replaced by ensemble averages. $\beta_e$ and $\omega_e^2$ for this situation can be solved by a prior assumption for the probability density function of $x(t)$.

REFERENCES

1. Lyon, R. H., "Empirical Evidence for Nonlinearity and Directions for Future Work," J. Acoust. Soc. Amer., *35*, No. 11 (November 1963), pp. 1712–1721.
2. Deutsch, R., "Nonlinear Transformations of Random Processes," Englewood Cliffs, New Jersey: Prentice-Hall, Inc., 1962.
3. Stratonovich, R. L., "Topics in the Theory of Random Noise," vol. 1, New York: Gordon and Breach, 1963.
4. Kuznetsov, P. I., Stratonovich, R. L., and Tikhonov, V. I., "Nonlinear Transformations of Stochastic Processes," New York: Pergamon Press Ltd., 1965.
5. Wang, M. C., and Uhlenbeck, G. G., "On the Theory of Brownian Motion II," Rev. Mod. Phys., *17*, Nos. 2 and 3 (April–July 1945), pp. 323–342.
6. Fine, T., "Partial Differential Equations for Densities of Random Processes," Technical Report No. 445, Craft Laboratory, Harvard University, Cambridge, Massachusetts, May 4, 1964.
7. Pawula, R. F., "Generalizations and Extensions of the Fokker-Planck-Kolmogorov Equation," IEEE Trans. on Inform. Theory, *IT-13*, No. 1 (January 1967), pp. 33–41.
8. Kushner, H. J., "On the Differential Equations Satisfied by Conditional Probability Densities of Markov Processes, with Applications," J. SIAM Control Series A, *2*, No. 1 (1962), pp. 106–119.

9. Kushner, H. J., "On the Dynamical Equations of Conditional Probability Density Functions, with Applications to Optimal Stochastic Control Theory," J. Math. Analy. Appl., 8 (1964), pp. 332–344.
10. Rosenbluth, M. N., MacDonald, W. M., and Judd, D. L., "Fokker-Planck Equation for an Inverse-Square Force," Phys. Rev., 107, No. 1 (July 1957), pp. 1–6.
11. Van Kampen, N. G., "Thermal Fluctuations in Nonlinear Systems," J. Math. Phys., 4, No. 2 (February 1963), pp. 190–194.
12. Ariaratnam, S. T., "Random Vibration of Nonlinear Suspensions," J. Mechanical Eng. Sci., 2, No. 3 (1960), pp. 195–201.
13. Hempstead, R. D., and Lax, M., "Classical Noise VI-Noise in Self-Sustained Oscillators near Threshold," Phys. Rev., 161, No. 2 (September 1967), pp. 350–366.
14. Crandall, S. H., "Perturbation Technique for Random Vibration of Nonlinear Systems," J. Acoust. Soc. Amer., 35, No. 11 (November 1961), pp. 1700–1705.
15. Khabbaz, G. R., "Power Spectral Density of the Response of a Nonlinear System to Random Excitation," J. Acoust. Soc. Amer., 38, No. 5 (November 1965), pp. 847–850.
16. Caughey, T. K., "Equivalent Linearization Techniques," J. Acoust. Soc. Amer., 35, No. 11 (November 1963), pp. 1706–1711.
17. Booton, R. C., "The Analysis of Nonlinear Control Systems with Random Inputs," Proc. Symp. Nonlinear Circuit Anal., Polytechnic Inst., Brooklyn, New York, 2 (1953).
18. Gray, A. H., "Uniqueness of Steady-State Solutions to the Fokker-Planck Equation," J. Math. Phys., 6, No. 4 (April 1965), pp. 644–647.
19. Klein, G. H., "Random Excitation of a Nonlinear System with Tangent Elasticity Characteristics," J. Acoust. Soc. Amer., 36, No. 11 (November 1964), pp. 2095–2105.
20. Liu, S. C., unpublished work.

# Contributors to This Issue

CLEO D. ANDERSON, B.S.E.E. 1960, University of Idaho; M.E.E., 1962, New York University; Bell Telephone Laboratories 1960—. Mr. Anderson has been mainly concerned with system analysis of submarine cable systems. He is now supervisor of the High Frequency Radio Group. Member, IEEE, Sigma Tau, Phi Kappa Phi, and Eta Kappa Nu.

WILLIAM F. BODTMANN, Monmouth College, 1957-61; Bell Telephone Laboratories 1941—. Mr. Bodtmann has been engaged in research on long- and short-haul microwave radio systems, frequency feedback receivers, and FM multiplex systems. He is working with communication systems operating at millimeter wavelengths.

SOO YOUNG CHAI, B.S.E.E., 1961, and M.S.E.E., 1962, Ohio State University; Ph.D., 1966, University of California at Berkeley; acting assistant professor of electrical engineering, University of California at Berkeley, 1966-1967; Bell Telephone Laboratories, 1967—. Mr. Chai was originally engaged in opto-electronics research with emphasis on optical devices using a color-selective spatial low-pass filter. Now he is doing exploratory development of a color video telephone system. Member, Eta Kappa Nu, Tau Beta Pi, Sigma Xi.

HERBERT YU-PANG CHANG, B.S., 1960; M.S., 1962, and Ph.D. (E.E.), 1964, University of Illinois; Bell Telephone Laboratories, 1964—. Mr. Chang has been engaged in the exploratory studies of fault diagnosis techniques for electronic switching systems, including the development of fault dictionary techniques for No. 1 ESS and exploratory development of a digital fault simulator for large processors. He is currently involved in studies of the techniques for the design of self-diagnosable digital machines and the reliability and maintainability studies for digital systems. Member, Sigma Xi, Tau Beta Pi, Eta Kappa Nu, Pi Mu Epsilon, IEEE, Association for Computing Machinery.

ARTHUR B. CRAWFORD, B.S. in E.E., 1928, Ohio State University; Bell Telephone Laboratories, 1928—. Mr. Crawford has specialized in radio research in the utrashort wave and microwave regions. He has been concerned with measurement techniques, propagation, and

2053

antenna studies. He designed the horn-reflector antenna used at Craw-ford Hill in the Project Echo and Project *Telstar*® communication satellite experiments and for radio astronomy studies. As Head of the Radio Techniques Research Department, he is in charge of a group concerned with antennas for short-hop microwave systems and satellite communications, radio astronomy, and certain devices for use in co-herent optics. Fellow, IEEE; member, Pi Mu Epsilon, Eta Kappa Nu, Tau Beta Pi, Sigma Xi.

JAMES C. DALY, B.S., 1960, University of Connecticut; M.E.E., 1962, Ph.D. 1967, Rensselaer Polytechnic Institute; member of the faculty of the Electrical Engineering Department at Rensselaer, 1962–1966; Bell Telephone Laboratories, 1966–1969. At Rensselaer, he did research on bulk semiconductor microwave interactions. At Bell Laboratories he has been concerned with optical wave guidance. Currently on leave of absence from Bell Telephone Laboratories, Mr. Daly is a visiting member of the faculty of the Electrical Engineering Department of the University of Rhode Island. Member, IEEE, Tau Beta Pi, Eta Kappa Nu, Sigma Xi.

CORRADO DRAGONE, Laurea in E.E., 1961, Padua University (Italy) ; Libera Docenza, 1968, Ministero della Pubblica Istruzione (Italy) ; Bell Telephone Laboratories, 1961—. Mr. Dragone has been engaged in experimental and theoretical work on microwave antennas and solid-state power sources. He is currently involved in solid-state radio systems experiments.

JOHN D. GABBE, B.A. 1950, New York University; M.S. 1951, Uni-versity of Illinois; Ph.D., 1957, New York University; Bell Telephone Laboratories, 1956—. Mr. Gabbe was first associated with the *Picture-phone*® visual telephone project, then with studies of the earth's magnetosphere. At present, he is engaged in research concerning the methodology of data analysis. Member, American Physical Society.

MRS. ANNE E. FREENY, B.A., 1957, University of Connecticut; M.S., 1959, Cornell University; Bell Telephone Laboratories, 1959—. Mrs. Freeny has worked in data analysis, concentrating primarily on the development of programs which would apply new statistical techniques to various large bodies of data. She has also worked on the organiza-tion of the results of the analyses. Member, Phi Beta Kappa, Phi Kappa Phi; associate member, Sigma Xi.

BERNARD GLANCE, Dipl. Ing., Ecole Spéciale de Mécanique et Electricité, 1958, Dipl. Ing., 1960, Ecole Superieure d'Electricité, Paris, France); C.S.F., Research Center of Corbeville, Orsay, France, 1960–1966; Dipl. Docteur (Ing.), Sorbonne, Paris, 1964; Bell Telephone Laboratories, 1968—. At C.S.F., Mr. Glance had been engaged in research on microwave tubes. At S.F.D. Laboratories, he had worked on high power microwave amplifiers. Mr. Glance is presently working on microwave solid-state integrated circuits.

F. E. GUILFOYLE, Newark College of Engineering; Bell Telephone Laboratories, 1955—. His work with the Radio Research Group has been concerned with components for the Echo and Telstar communication experiments, and laser communication measurements and techniques.

HARRY E. KELLER, Bell Telephone Laboratories, 1942—. Mr. Keller's work has included microwave branching filters, microwave radio systems, multiplex for radio systems, frequency modulation with feedback receiver for Echo and Holmdel *Telstar®* communication satellite receiver. His most recent project is the rainfall measuring network data recording system.

L. U. KIBLER, B.S., 1950, U. S. Coast Guard Academy; M.S.E.E., 1956, Massachusetts Institute of Technology; Ph.D., 1968, Polytechnic Institute of Brooklyn; Bell Telephone Laboratories, 1956—. Mr. Kibler has been concerned with experimental research in parametric amplifiers, tunnel diodes, lasers, microwave photo diodes, and Schottky-barrier diode converters. He participated in the design and operation of the receivers for the Echo and *Telstar®* projects. He is engaged in millimeter wave antenna investigations. Member, IEEE, Eta Kappa Nu, Sigma Xi.

KANEYUKI KUROKAWA, B.S., 1951, and Ph.D., (Engineering), 1958, University of Tokyo; Assistant Professor, University of Tokyo, 1957–1963; Bell Telephone Laboratories, 1963—. Mr. Kurokawa has been concerned mainly with microwave solid-state circuits, including amplifiers, oscillators, and switches. He supervises a group responsible for the exploratory development of solid-state functional devices and circuits. Member, Institute of Electronics and Communications Engineers of Japan, IEEE.

S. C. LIU, B.S. in C.E., 1960, National Taiwan University; M.S., 1964, and Ph.D., 1967, University of California at Berkeley; Bell Telephone Laboratories, 1967—. Mr. Liu has been doing research in applied mechanics, structural dynamics, random vibrations, and earthquake engineering. Member, American Society of Civil Engineers, Seismological Society of America.

DEAN E. McCUMBER, B.E., 1952, and M.E., 1955, Yale University; A.M., 1956, and Ph.D., 1960, Harvard University; Engineering Division Officer, U.S.S. Tripoli, 1952–54; National Science Foundation Postdoctoral Fellow, 1959–61; Bell Telephone Laboratories, 1961—. Since joining Bell Laboratories, Mr. McCumber has been concerned with the physical theory of optical impurities in solids, of lasers, of electron bulk-effect devices, and, currently, of superconductor tunneling and weak-link devices. He is Head of the Crystal Electronics Research Department. Fellow, American Physical Society; Member, Tau Beta Pi.

THOMAS L. OSBORNE, B.S.E.E., 1961, M.S.E.E., 1963, Auburn University; Bell Telephone Laboratories, 1963—. Mr. Osborne has been involved in research on solid-state microwave radio systems and associated circuits. Member, IEEE, Sigma Xi, Phi Kappa Phi, Tau Beta Pi, Eta Kappa Nu, Pi Mu Epsilon.

HARRY RUDIN, JR., B.E., 1958, M. Eng., 1960, and D. Eng., 1964, Yale University; Bell Telephone Laboratories, 1964–1968; IBM Research Laboratory (Zurich, Switzerland), 1968—. Mr. Rudin was Instructor in Electrical Engineering at Yale University from 1961 until 1964. At Bell Telephone Laboratories he worked in the data communication area, concentrating on automatic equalization. At IBM he is a full-time consultant working in the general area of computer-related communications. He is a former executive of the IEEE Connecticut Section and is a member of the Yale Engineering Association executive board.

CLYDE L. RUTHROFF, B.S.E.E., 1950, and M.A., 1952, University of Nebraska; Bell Telephone Laboratories, 1952—. Mr. Ruthroff has published contributions on the subjects of FM distortion theory, broadband transformers, FM limiters, threshold extension by feedback, and microwave radio systems for satellite and terrestrial use.

He is interested in the extension of radio communication into the millimeter and optical wavelengths. Member, A.A.A.S., I.E.E.E., Sigma Xi.

IRWIN W. SANDBERG, B.E.E., 1955, M.E.E., 1956, and D.E.E., 1958, Polytechnic Institute of Brooklyn; Bell Telephone Laboratories, 1958—. Mr. Sandberg has been concerned with analysis of military systems, synthesis and analysis of active and time-varying networks, studies of properties of nonlinear systems, and some problems in communication theory and numerical analysis. He is head of the Systems Theory Research Department. Member, IEEE, Eta Kappa Nu, Sigma Xi, Tau Beta Pi.

M. V. SCHNEIDER, M.S., 1956, and Ph.D., 1959, Swiss Federal Institute of Technology, Zurich, Switzerland; Bell Telephone Laboratories, 1962—. Mr. Schneider has been engaged in experimental work on thin-film solid-state devices, optical detectors, and microwave integrated circuits. Member, IEEE, American Vacuum Society.

R. A. SEMPLAK, B.S. (Physics), 1961, Monmouth College; Bell Telephone Laboratories, 1955—. Mr. Semplak has been engaged in research on microwave antennas and propagation. He participated in Project Echo and Project *Telstar*® Communication satellite experiments. He currently is concerned with the attenuation effects of rain on propagation at 18.5 and 30.9 GHz. Member, Sigma Xi, American Association for the Advancement of Science.

WILLIAM W. SNELL, JR., 1951, Williamsport Technical Institute; Bell Telephone Laboratories, 1955—. His early work for the radio research department centered around waveguide components in the 4, 6, and 11 GHz common carrier bands. This included ferrite devices, microwave mixers, and polarization couplers. He later participated in the Shotput experiment, a suborbital proving test for Project Echo. During Project Echo he designed, built, and patented several components of the Crawford Hill receiving terminal. More recently he has been concerned with high order varactor frequency multipliers and varactor diode fabrication. He is presently interested in making hybrid integrated circuits for RF systems above 10 GHz.

LEROY C. TILLOTSON, B.S.E.E., 1938, University of Idaho; M.S.E.E., 1940, University of Missouri; D.Sc. (Hon.), 1966, University of

Idaho; Bell Telephone Laboratories, 1941—. Mr. Tillotson initially worked on the design of filters and networks. In 1954 he was appointed head of a group working on radio relay systems. Beginning in July 1958, he spent more than a year on a leave of absence with the Institute for Defense Analysis, Washington D. C., where he was concerned with communication satellites and other space-related activities. He was appointed Director of Radio Research in 1963.

RICHARD H. TURRIN, B.S.E.E., 1956, Newark College of Engineering; M.E.E., 1960, New York University; Bell Telephone Laboratories, 1956—. Mr. Turrin has been involved in antenna research and development using microwaves and millimeter waves. He has also been involved in propagation studies related to atmospheric and environmental effects and is concerned with satellite antenna problems. Member, IEEE, Eta Kappa Nu, Tau Beta Pi.

JACOB ZIV, B.Sc. 1954, Engineering Diploma 1955, and M.Sc. 1957, Technion-Israel Institute of Technology, Haifa, Israel; D.Sc., 1962, Massachusetts Institute of Technology; Scientific Department, Israel Ministry of Defence, 1955–1959, 1962–1968; Bell Telephone Laboratories, 1968—. (On leave of absence from Israel Ministry of Defence.) Mr. Ziv has been engaged in research in information theory and statistical communication theory. Member, IEEE.

CONTENTS

(Continued from front cover)