

# The Nature of Speech and Its Interpretation<sup>1</sup>

By HARVEY FLETCHER

## INTRODUCTION

**V**ARIOUS phases of this subject have received serious study by phoneticians, otologists, and physicists. On account of its universal interest, it has received attention from men in many branches of science. In spite of the large amount of time devoted to the subject, the progress in understanding its fundamental aspects has been rather slow. At the present time the physical properties which differentiate the various fundamental speech sounds are understood in only a very fragmentary way. Some very interesting and painstaking work has been done on the physical analysis of vowel sounds, but the results to date are far from conclusive. Although several theories have been advanced to explain the way in which the ear interprets sound waves, they are still in the controversial stage.

The material which is presented here is the result of an investigation which has been carried on in the Research Laboratories of the American Telephone and Telegraph Company and Western Electric Company during the past few years.

To make a quantitative study of speech and hearing it is necessary to obtain the speech sounds at varying degrees of loudness and with definitely known amounts of distortion. The main reason why so few real results have been obtained in the investigation of speech sounds is due to the fact that it is extremely difficult to change the volume and distortion of these sounds by acoustic means. Due to recent developments in the electrical transmission of speech it is possible to produce the equivalent of these changes by electrical means. For this purpose a telephone system was constructed which reproduced speech with practically no distortion. It was arranged so that by means of distortionless attenuators the volume of reproduced speech could be varied through a very wide range, and so that by introducing various kinds of electrical apparatus the transmitted speech wave could be distorted in definitely known ways.

A method was developed for measuring quantitatively the ability of the ear to interpret the transmitted speech sounds under different conditions of distortion and loudness. By choosing these conditions properly, considerable information was gained concerning both speech and hearing. This indirect method of attack has a distinct advantage

<sup>1</sup> Presented at a meeting of the Electrical Section of the Franklin Institute held Thursday, March 30, 1922. Reprinted from the Journal of the Franklin Institute for June 1922.

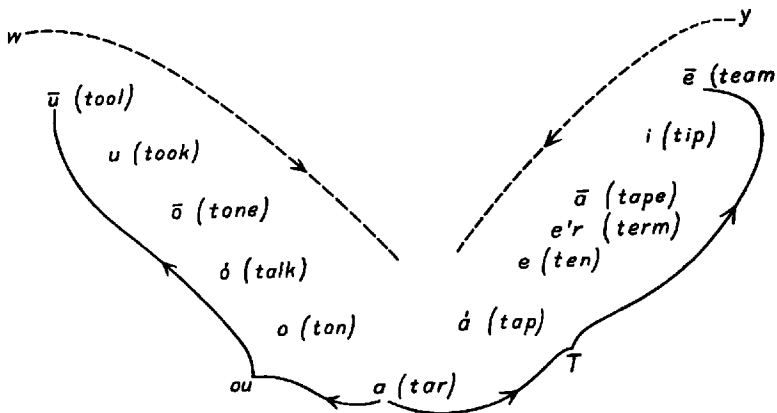
for engineering purposes, in that it measures directly the thing of most interest, namely, the degrading effect upon telephone conversation of introducing electrical distortion into the transmission circuit. However, the application of the results is not limited to this particular field.

METHOD OF MEASURING THE QUALITY OF SPEECH

Briefly stated the method consists in pronouncing detached speech sounds into the transmitting end of the system and having observers write the sounds which they hear at the receiving end. The comparison of the called sounds with those observed shows the number and kinds of errors which are made. The per cent of the total sounds spoken which are correctly received is called the articulation of the system.

TABLE I.  
Classification of the Speech Sounds.

Pure Vowels



Combinational and Transitional Vowels

w - y - ou - i - h

Semi-vowels

l - r

Stop Consonants

Voiced

b  
d  
j  
g

Unvoiced

p  
t  
ch  
k

Nasalized

m  
n  
—  
ng

Formation of Stop

lip against lip  
tongue against teeth  
tongue against hard palate  
tongue against soft palate

Fricative Consonants

Voiced

v  
z  
th (then)  
zh (azure)

Unvoiced

f  
s  
th (thin)  
sh

Formation of Air Outlet

lip to teeth  
teeth to teeth  
tongue to teeth  
tongue to hard palate

In order to understand the construction of the articulation lists and also to interpret the results of this investigation, I desire to give here a brief classification of the speech sounds, which is based upon the position of the various speech organs when the sounds are being produced. It is shown in the accompanying table (Table I).

The pure vowels are arranged in the vowel triangle, which is familiar to phoneticians. Starting with the sound  $\bar{u}$  the lips are rounded and there is formed a single resonant cavity in the front part of the mouth. Passing along the left side of the triangle from  $\bar{u}$  to  $\bar{e}$  the mouth is gradually opened with the tongue lowered to form the successive vowels. Going along the right side of the triangle from  $\bar{a}$  to  $\bar{e}$ , the tongue is gradually raised to the front part of the mouth forming two resonant chambers in the mouth cavity. An infinite number of different shadings of these vowels may be produced by placing the mouth in the various intermediate positions, but the ones which are shown were chosen as being the most distinct.

The sounds  $w$ ,  $y$ ,  $ou$ ,  $\bar{i}$  and  $h$  are classed as combinational and transitional vowels. As the mouth is placed in the position to say  $\bar{u}$  and then suddenly changed so as to form any other vowel in the triangle, the result obtained is signified in writing by placing the letter  $w$  before the vowel. In a similar way we get the effect usually designated by  $y$  if the position of the vowel suddenly changes from  $\bar{e}$  to any other vowel. An infinite variety of diphthongs can be formed by changing the position of the mouth necessary to form one vowel to that to form another without interrupting the voice. The most distinct and principal ones used in our language are formed by passing from the sound  $\bar{a}$  to either extreme corner of the triangle and are known as  $ou$  and  $\bar{i}$ . When a vowel commences a syllable it is formed by suddenly opening the glottis, permitting the air, which has been held in the lungs, to escape into the mouth, which is formed for the proper vowel. If the glottis remains open and the vowel is started by the sudden contraction of the lungs, we have the effect which is represented in writing by placing an  $h$  before the vowel. The sounds  $l$  and  $r$  are called semi-vowels because the voice train is partially interrupted, although the sound can be continued. The stop and fricative consonants are classified in a manner which is familiar to phoneticians.

It will be noticed that the markings are not those used in the international phonetic alphabet which were entirely too complicated for practical use. Only the bar and accent stroke are used. These can be written quickly and with little chance of error.

In order to pronounce these speech sounds properly, they must

be combined into syllables. For the purpose of this investigation they were combined into mono-syllables of the simple types consonant-vowel, vowel-consonant, and consonant-vowel-consonant.

To eliminate memory effects every possible combination of the sounds into these types of syllables was used unless there was a good reason for excluding it. The complete list contained 8700 syllables. For convenience of testing these syllables were divided into groups of fifty. Each group contained the same kind and number of syllable forms and an equal number of each of the fundamental vowel and consonant sounds.

TABLE II.  
*Speech-sound Testing List. List No. 160*

	Speech-sound	Key-word		Speech-sound	Key-word
1	ha	ho(t)	26	gōb	go + b
2	hā	hay	27	shōl	shoal
3	wā	wa(g)	28	ros	rus(t)
4	wi	wi(th)	29	jod	ju(g) + d
5	vou	vow	30	bok	buck
6	ār	air	31	zīk	z + (d)ike
7	ez	e(bb) + z	32	bīch	buy + ch
8	ūsh	you + sh	33	kīth	ki(te) + th
9	an	on	34	gīt	gui(de) + t
10	id	(l)id	35	yīf	y + if
11	jouv	jow(l) + v	36	sin	sin
12	moush	mou(nd) + sh	37	tērm	term
13	rour	r + our	38	mērl	m + earl
14	zūth	z + (s)oothe	39	pērv	p + (n)erve
15	hūs	who + s	40	yēt	y + eat
16	chush	ch + (p)ush	41	bēl	b + eel
17	jum	j + (f)oo(t) + m	42	zef	ze(al) + f
18	thup	th + (s)oo(t) + p	43	weng	whe(n) + ng
19	fuch	foo(t) + ch	44	kev	k + ev(er)
20	wōng	wa(l) + ng	45	hāng	hang
21	chōth	cha(lk) + th	46	pāg	p + (r)ag
22	tōj	ta(l) + j	47	yās	y + ace
23	kōg	k + aug(er)	48	dāp	d + ape
24	fōn	(tele)phone	49	yang	ya(cht) + ng
25	dōs	dose	50	lan	l + on

To illustrate the technique of articulation testing a sample list is given in Table II. In the first column the syllable is given in its phonetic form. A key-word showing how each syllable is pronounced is given in the second column. These syllables were written on cards which were shuffled each time before they were used, so that the order in which they were pronounced was entirely haphazard. One hundred and seventy-four similar lists were used in this work. In order to eliminate personal peculiarities, several

callers and observers were used. In Table III are shown the results obtained by an observer when this list was transmitted over a system which eliminated all frequencies above 1250 cycles per second.

TRANSMISSION BRANCH  
ARTICULATION TEST RECORDING SHEET

WORD  
ARTICULATION  
40 %

TITLE OF TEST J20311

CONDITION TESTED Low Pass Filter - 1250~

DATE 2-7-20

Attenuation 5 nepiers down OBSERVER M.A.

TEST No. 11

CALLER H.E.D.

LIST No. 160

No.	OBSERVED	CALLED	ERRORS	No.	OBSERVED	CALLED	ERRORS
1	tan	t'erm	r-a m-n	26	zip	thup	th-z u-i
2	zit	g'it	g-z t-i	27	ko'd	to'j	t-k j-d
3	wa	wa'	a'-a	28	t'ish	chush	ch-t u-i
4	dāp	✓		29	yang	✓	
5	gōb	✓		30	zēt	zūth	ū-ō th-t
6	yis	yif	f-s	31	ref	ras	a-g e-f
7	māl	merl	r-ā	32	jum	✓	
8	thin	sin	s-th	33	jo'g	ko'g	k-j
9	zip	zik	k-p	34	jad	jo'd	o-a h-r
10	jouv	✓		35	tūth	hūs	s-th
11	yāt	yās	s-t	36	id	✓	
12	thou	vou	v-th	37	ha	✓	
13	b'p	b'ch	ch-p	38	fōn	✓	
14	hāng	✓		39	ko'th	cho'th	ch-k
15	mīs	moush	ou-i sh-s	40	rou	✓	
16	dāch	dōs	ō-ā s-ch	41	an	✓	
17	kev	✓		42	bok	✓	
18	tig	pāg	p-t ō-i	43	yēt	✓	
19	kīs	kīth	th-s	44	o'r	a'r	a'-o'
20	hā	✓		45	yōth	ūsh	y inserted ū-ō th-th
21	weng	✓		46	wōng	✓	
22	dāl	bāl	b-d	47	kōv	perv	p-k e'r-ō
23	thich	fuch	f-th u-i	48	zēt	zēf	f-t
24	wif	wi	f inserted	49	lan	✓	
25	ez	✓		50	shēl	✓	

TABLE III.

The correct word is written opposite all of the syllables which were recorded incorrectly. The errors for each of the fundamental sounds were taken from this original sheet and recorded on an analysis



sheet as shown in Table IV, for example it will be noticed that p was recorded as k 24.4 per cent, as p 45 per cent, and as t 22.2 per cent of the times called. On the other hand the sound w was only recorded incorrectly 1 per cent of the times called.

For this system the consonant articulation was 65.8 and the vowel articulation 83.4.

DESCRIPTION OF THE SYSTEM FOR REPRODUCING SPEECH SOUNDS

The telephone system used in this investigation is probably more nearly perfect than any other which has yet been built. Its essential elements are a condenser transmitter to receive the speech waves and transform them into the electrical form, an amplifier for magnifying the intensity of the electrical speech currents, an attenuator for controlling the intensity, an equalizing network, and a receiver for delivering the speech to the ear. A schematic arrangement of the circuit is shown in Fig. 1.

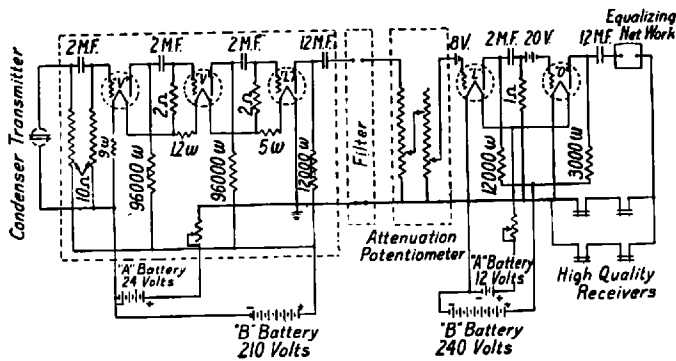


FIG. 1.—High Quality Telephone System

A detailed description of the construction and operation of the condenser transmitter has been given by Crandall and Wentz and published in the *Physical Review*.<sup>1</sup> It is simply an air condenser, one of its plates being a flexible metal diaphragm.

A five-stage vacuum tube amplifier was used. Particular care was taken in coupling the stages together, so that the amplifier was practically free from frequency distortion.

The attenuator consisted of a potentiometer arrangement which could reduce the amplitude of the speech waves to approximately one-millionth of their maximum values.

The equalizing network was an arrangement of resistances, con-

<sup>1</sup> Crandall, *Phys. Rev.*, June, 1918; Wentz, *Phys. Rev.*, July, 1917.

densers and inductance coils having a frequency selectivity which was the complement of that of the rest of the system.

The telephone receiver was a bipolar type having a special construction which was designed to broaden the range of frequency response.

The reproducing efficiency of the system from the mouth of the speaker to the ear of the listener for each frequency is shown in Fig. 2.

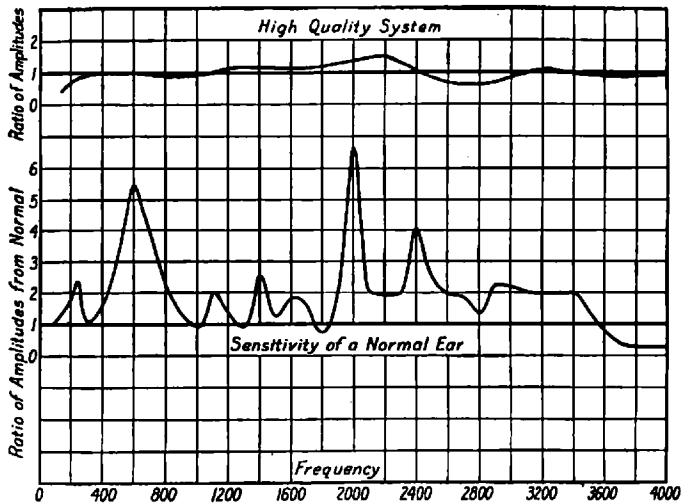


FIG. 2.

The pitch or frequency of the tone is given on the X axis. The ordinates represent amplitude ratios or the number of times the amplitude of the tone reaching the ear was greater than that which entered the transmitter. It will be seen that this high quality system has practically a uniform response for all frequencies throughout the speech range.

In order that its uniformity may be appreciated, a comparison curve is given. This curve shows the deviation in the sensitivity of a typical individual ear from the average sensitivity of a large number of ears. The ordinates represent the ratio of amplitudes at the various pitches which was necessary to bring the tone to the threshold of audibility. It is evident that this deviation is much larger than the departure of the high quality circuit from uniformity.

To show that this particular individual's curve is typical, the curves for both ears of 20 women are given in Fig 3. For convenience these curves are plotted on logarithmic paper. If an arithmetic



scale is used, all of the curves below the mean are crowded together in the small space between zero and one, and all those above the mean are stretched out from one to infinity. By using a logarithmic plot a symmetrical distribution is obtained. The method of obtain-

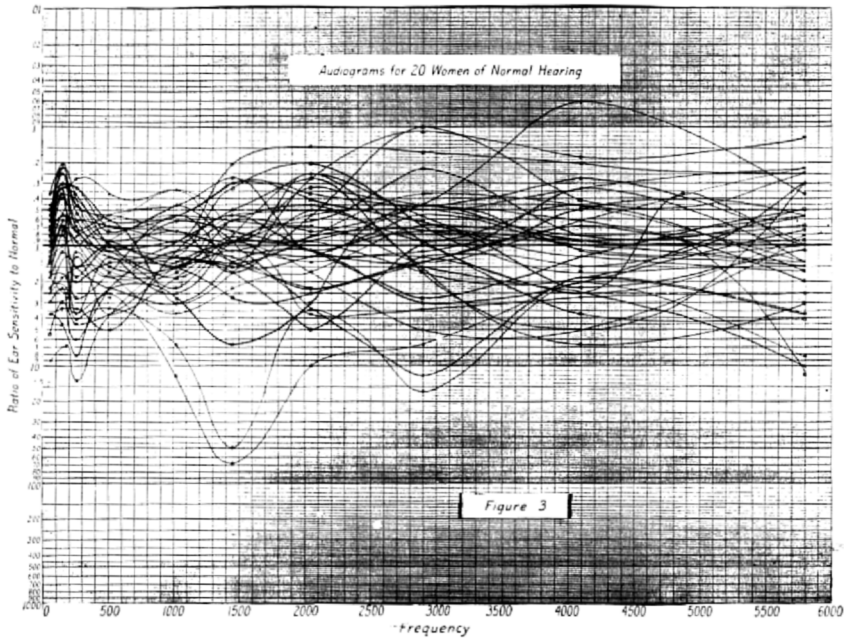


FIG. 3.

ing these ear sensitivity curves was fully described in a recent paper<sup>2</sup> given before the Natural Academy of Sciences.

It is interesting to note that they indicate that each individual has a hearing characteristic which is quite different from other individuals. Consequently speech sounds differently to different persons. Any distortions of the speech sounds will necessarily affect some persons differently from others. It is evident then that in discussing speech and hearing we must deal with statistical averages.

Experimental articulation tests showed that the ear interpreted the speech which was transmitted over this high quality system practically as well as that transmitted through the air. Some may wonder why such good quality is not furnished telephone users in commercial practice: Scientifically speaking, it is possible to furnish such quality, but it is evident that the equipment involved is so com-

<sup>2</sup> Fletcher and Wegel, *Proc. Nat. Acad. Science*, Vol. 8, No. 1, pp. 5-6, Jan., 1922.

plicated that such service would be altogether too costly for commercial use; people could not afford to pay for it.

#### THE RELATION BETWEEN THE VOLUME AND ARTICULATION OF UNDISTORTED SPEECH

Articulation tests were made upon the high quality telephone system described above when it was set to deliver various intensities from the threshold of audibility to very large values. The results shown as syllable articulation values are given by the curve

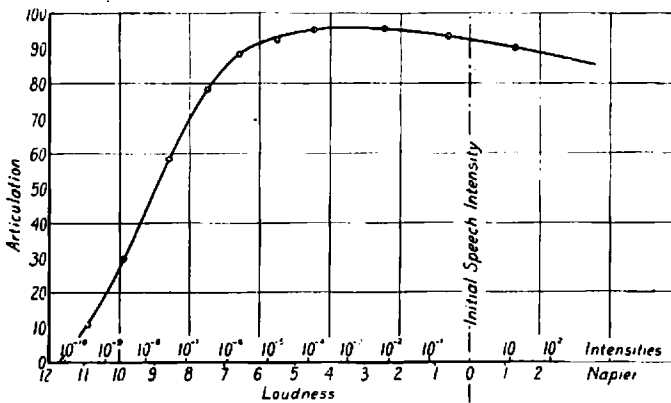


FIG. 4.

in Fig. 4. The abscissas in this curve represent loudness and are expressed as the natural logarithm of the number of times the speech wave amplitude has been decreased from the initial intensity at  $\frac{1}{2}$  inch in front of the mouth of the callers. This unit of loudness has never been given a name, and as a matter of convenience in this work it is called a napier. It will be noticed that when the volume is reduced  $11\frac{1}{2}$  napiers below the initial speech intensity the articulation becomes zero. This point also represents the value at which the speech becomes inaudible and corresponds to approximately  $1/1000$  dynes per square centimetre pressure variation against the ear drum. In energy units it is a reduction of ten billion times below the initial speech intensity. For very loud initial speech this point is shifted about 1 napier. For purposes of comparison the intensity reductions are also indicated on the loudness axis.

At 3 napiers below or at about  $1/1000$  of the initial speech intensity the articulation becomes a maximum. Louder speech than this seems to deaden the nerves so that a person makes a less accurate

interpretation of the received speech. These results were obtained in a room which was especially constructed to exclude outside noise. When noise is present at the receiving station the optimum loudness increases as the noise increases.

The articulation data were analyzed so as to show the errors of each of the fundamental sounds. The curves given in Fig. 5 show the results of this analysis. It will be noticed that the volume at which errors begin to be appreciable is different for the different sounds and is usually higher for the consonants than for the vowels.

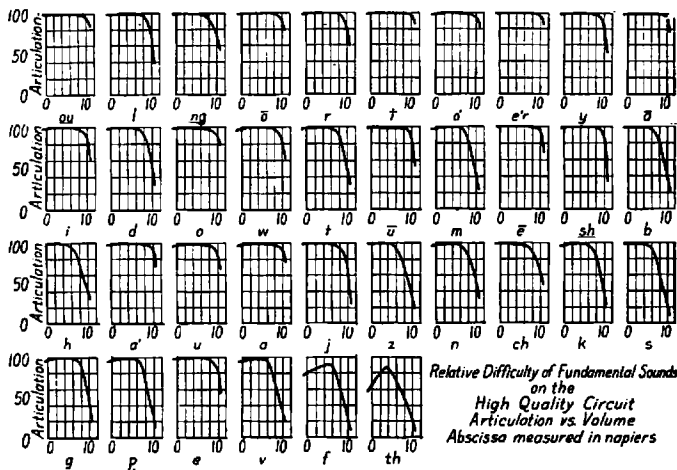


FIG. 5.

Within the precision of the test the intersection point on the X axis was the same for all the sounds, namely at 11.5 napiers.

It will be noticed that the consonants are usually harder to hear than the vowels. However, the speech sounds e and l, r, ng form notable exceptions to this general rule, since the former is among the most difficult, while the latter are among the very easiest speech sounds. The order in which the speech sounds are given here represents their relative difficulty of interpretation when received at average intensities. At all intensities, the sounds th, f and v are the most difficult. Z, h and s become very difficult at weak volumes. The sounds i, ou, er and ó are missed less than 10 per cent of the time, even with "very weak" intensity. At "average" volumes there are only three sounds more difficult than e while at "very weak" volumes there are 23 sounds more difficult. At very weak volumes l, which is the easiest sound at "average" volumes is missed three times as often as e.

We will now pass to a consideration of the effect of distortion upon the articulation of the sounds.

#### DESCRIPTION OF ELECTRICAL FILTERS USED TO PRODUCE DISTORTION

In order to investigate distortion we would like to be able to take the train of speech waves going from the mouth to the ear and operate upon it in various ways such as eliminating frequencies in certain regions without marring or disturbing other frequencies. For ex-

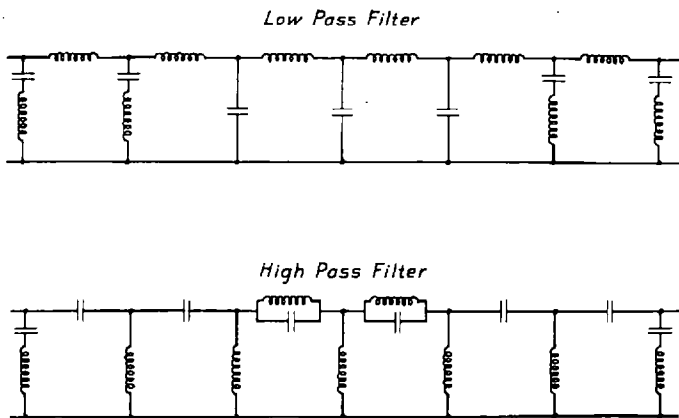


FIG. 6.

ample, if all frequencies above 1000 were eliminated, it would be possible to determine what intelligibility is carried by this range of frequencies.

Fortunately one of the recent electrical inventions is admirably adapted for this purpose, namely, the electrical wave filter invented by Dr. G. A. Campbell. This device was used extensively in this investigation.

The schematic circuit diagrams of the two types of filters which were used are given in Fig. 6.

This arrangement of coils and condensers produces an electrical conductor with the unusual properties that it transmits without appreciable diminution in amplitude any frequency between certain limits and reduces the amplitude of all frequencies outside these limits to less than 1/1000 of their original value. By varying the numerical values of the inductances and capacities this transmitted range can be placed at any desired position. In the arrangement which was used in the investigation these coils and condensers were

housed in two boxes. The switching mechanism was arranged so that by turning a dial the condensers and coils were connected in such a way that the filter transmitted different frequency bands.

In Fig. 7 are shown the transmission properties of the low pass filter when the dial is set to transmit frequencies from 0 to 1500. It is seen that for frequencies below 1400 the amplitudes of the transmitted tones are always greater than .8 of their initial values, while for frequencies above 1500 the amplitudes are decreased to less than .001 of their initial values. These electrical filters were connected into the high quality circuit between the third and fourth stages

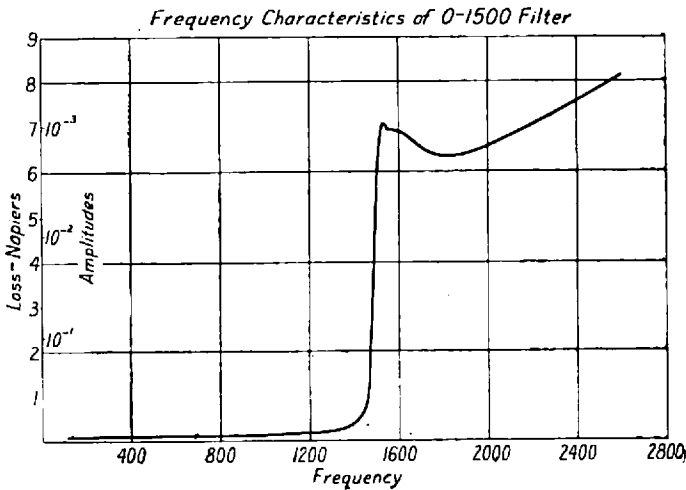


FIG. 7.

of the amplifier as indicated in Fig. 1. This combination formed a system which would pick up a complex sound wave and transmit faithfully to the ear those component frequencies in any desired region and eliminate all other frequencies.

#### RESULTS OF ARTICULATION TESTS WITH FILTER SYSTEMS

Articulation tests were made with these filter systems and the results analyzed as described above. In Fig. 8 the syllable articulation results are shown in graphical form. The ordinates for the solid curves represent the per cent of the articulation syllables called into the system which were correctly recorded at the observing end. The abscissas represent the so-called "cut off" frequency of the filter. For example on the curve labelled "Articulation L" the point (1000, 40) means that a system which transmits only frequencies

below 1000 cycles per second has a syllable articulation of 40 per cent. Similarly on the curve labelled "Articulation H" the point (1000, 86) means that a system which transmits only frequencies above 1000 cycles per second has a syllable articulation of 86 per cent. The dotted curves show the per cent of the total speech energy which is transmitted through the filter systems used in the articulation tests. These curves are derived from the results of Crandall and MacKenzie which were recently published.<sup>3</sup>

It will be seen that although the fundamental cord tones with their first few harmonies carry a large portion of the speech energy,

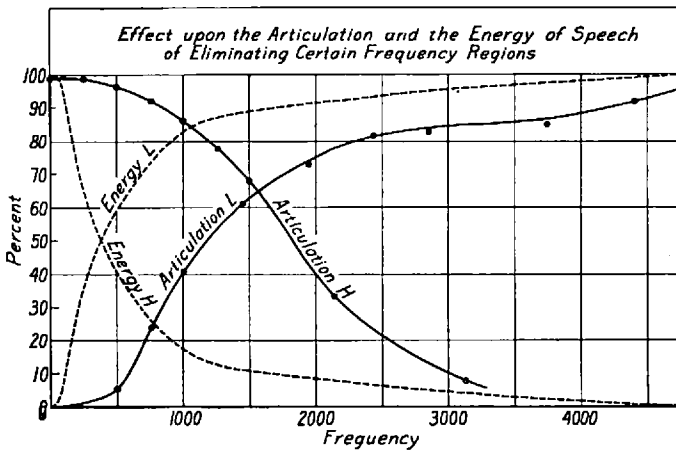


FIG. 8.

they carry practically none of the speech articulation. A filter system which eliminates all frequencies below 500 cycles per second eliminates 60 per cent of the energy in speech; but only reduces the articulation 2 per cent. A system which eliminates frequencies above 1500 cycles per second eliminates only 10 per cent of the speech energy, but reduces the articulation 35 per cent. A system which eliminates all frequencies above 3000 cycles per second has as low a value for the articulation as one which eliminates all frequencies below 1000 cycles per second. This last statement may appear rather astonishing since it is contrary to the popular notion of the relative importance of various voice frequencies from an interpretation standpoint.

The two solid curves intersect on the 1550 cycle abscissa and at 65 per cent articulation, which shows that using only frequencies

<sup>3</sup> See preceding paper.

above or frequencies below 1550 cycles an articulation of 65 per cent will be obtained. The two dotted curves necessarily intersect at 50 per cent.

The curves in Fig. 9 show how the articulation of some of the fundamental speech sounds was affected by eliminating certain frequency regions. The ordinate gives the number of times the sound was written correctly per 100 times called. As in Fig. 8 the

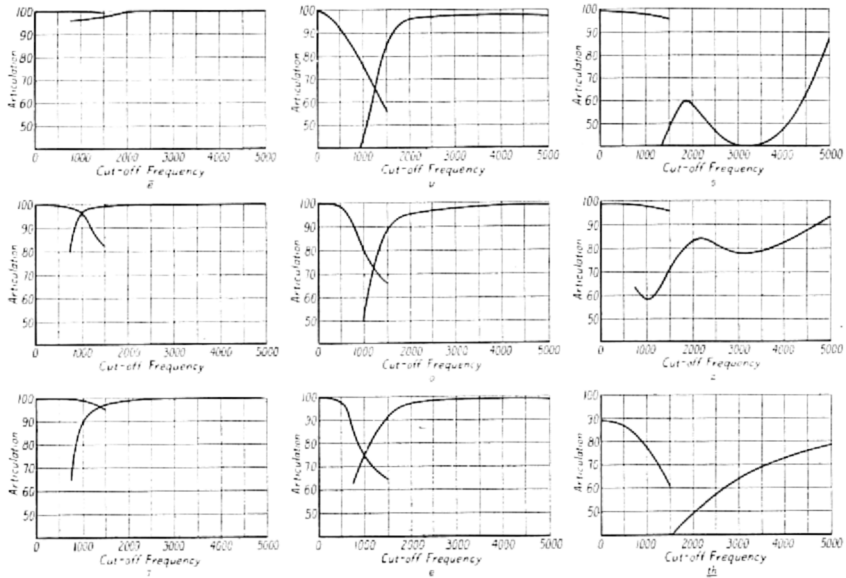


FIG. 9.

left hand curve shows the effect of eliminating all frequencies below and the right hand curve the effect of eliminating all frequencies above the frequency specified by the abscissa.

These nine speech sounds were chosen as representing three important classes. It is seen that the long vowels  $\bar{e}$ ,  $l$  and  $\bar{i}$  can be transmitted with an error of less than 3 per cent when using either half of the range of frequencies. When using either frequencies from 0 to 1700 or from 1700 to infinity  $\bar{e}$  was interpreted correctly 98 per cent of the time. Similarly  $l$  was interpreted correctly 97 per cent of the time when using either the range from 0 to 1000 or 1000 to infinity, and  $\bar{i}$  96 per cent of the time when using either the range from 0 to 1350 or from 1350 to infinity. The short vowels,  $u$ ,  $o$  and  $e$  are seen to have important characteristics carried by frequencies below 1000. More than a 20 per cent error is made on any of these

three sounds when frequencies below 1000 are eliminated. The elimination of frequencies above 2000 produces almost no effect.

The fricative consonants *s*, *z* and *th* are seen to be affected very differently from those in the other two classes. These sounds are very definitely affected when frequencies above 5000 are eliminated. The sounds *s* and *z* are not affected by the elimination frequencies below 1500. It is principally due to these three sounds that the syllable articulation is reduced from 98 per cent to 82 per cent when frequencies above 2500 cycles are eliminated.

A more detailed analysis of the articulation results on all the speech sounds showing the kind as well as the number of errors will be given in a future paper.

### CONCLUSION

In conclusion then we see that the intensity of undistorted speech which is received by the ear can be varied from 100 times greater to one-millionth less than the initial speech intensity without noticeably affecting its interpretation. The intensity must be reduced to one-ten-billionth of that initial speech intensity to reach the threshold of audibility for the average ear. Also it is seen that any apparatus designed to reproduce speech and preserve all of its characteristic qualities must transmit frequencies from 100 to above 5000 cycles with approximately the same efficiency. Although most of the energy in speech is carried by frequencies below 1000, the essential characteristics which determine its interpretation are carried mostly by frequencies above 1000 cycles. In ordinary conversation the sounds *th*, *f* and *v* are the most difficult to hear and are responsible for 50 per cent of the mistakes of interpretation. The characteristics of these sounds are carried principally by the very high frequencies.

It is evident that progress in the knowledge of speech and hearing has a great human interest. It will greatly aid the linguists, the actors, and the medical specialists. It may lead to improved devices which will alleviate the handicaps of deaf and dumb persons. Furthermore this knowledge will be of great importance to the telephone engineer, and since the telephone is so universally used, any improvement in its quality will be for the public good.

These humanitarian and utilitarian motives as well as the pure scientific interest have already attracted a number of scientists to this field. Now that new and powerful tools are available, it is expected that in the near future more will be led to pursue research along those lines.