E. H. Rothauser

# The Integrated Vocoder and its Application in Computer Systems

**Abstract:** This paper reviews the conceptual features and the applications of the integrated vocoder. In a comparison with the structural properties of the channel vocoder and other experimental types, the integrated vocoder is shown to be a technically simple solution to the well-known problems of pitch description and vocoder signal transmission. The design reduces the complexity of the over-all speech processing equipment and features a saving in device hardware. The potential of this vocoder design for time division multiplex transmission systems is discussed and an application of the new vocoder concept is shown for computer input-output equipment.

## 1. Introduction

### 1.1 Background

For digital transmission of speech and for processing of speech signals with digital computers some kind of analog-digital conversion, e.g. code modulation, is necessary. The use of vocoder techniques enhances coding efficiency and permits reduction, by a factor of 10 or more, of channel capacity necessary for the transmission of speech signals. The fundamental concepts of the vocoder were introduced by H. Dudley[1] in 1939, but for many years afterwards only a few papers reported work on vocoder development. Among these was the work of Halsey and Swaffield[2] at the British Post Office. For the last ten years there has been a rapidly increasing interest in vocoder techniques, which have been found to be important not only for speech transmission itself, but also for speech processing with digital computers.[3,4] The reduction of the required channel capacity facilitates computer input/output operations and permits a larger quantity of speech signals to be stored and processed in a given computer.

The present paper reviews the basic structure of some typical vocoders and summarizes some of the results of an effort which began at the University of Technology in Vienna, Austria, in 1954 and continued[5-10] at the IBM Laboratory in Vienna since 1961. The new vocoder structure described here will be designated the "integrated vocoder."

### 1.2 Difficulties of vocoder applications

Since the first vocoder was built by Dudley nearly thirty years ago, other groups have also built experimental vocoder systems, nearly all of which differ in structural details. Many have exhibited excellent performance under laboratory conditions. Experience shows that channel vocoders allow speech transmission over channels with a capacity as low as 2400 bits/sec. At this rate the speech quality of a good channel vocoder can still be sufficient for many purposes. Speech intelligibility will be high, but difficulties may be encountered with the naturalness of the synthetic voice sounds and in identifying speakers. For good "telephone quality" a channel capacity of the order of 10 kbits/sec seems still to be required. In contrast to the success of the vocoder in laboratory tests, there are two key problems which restrict widespread vocoder applications in the field: pitch detection and multiplexing.

### 1.21 Pitch detection

The first of these problems is pitch detection and pitch measurement for distorted speech signals. The voiced-unvoiced discrimination and the pitch measurement are relatively easy for high-fidelity speech input. But in most field applications the input speech is not much better than normal telephone quality; particularly, it is frequency-limited at the lower end and does not contain the fundamental of voice sounds. Furthermore, the signal-to-noise ratio is small, due not only to line noise and ambient noise in the environment of the speaker, but also to special properties of the carbon microphones ordinarily used in telephone handsets. The described effects, together with the natural inherent pitch variations of speech signals, make a sure voiced-unvoiced discrimination very difficult, if not impossible, and lead to errors in the pitch measurements.

As a consequence, the resulting synthesized output speech signal will be of low quality. While its intelligibility can still be sufficient, because speech signals even with a
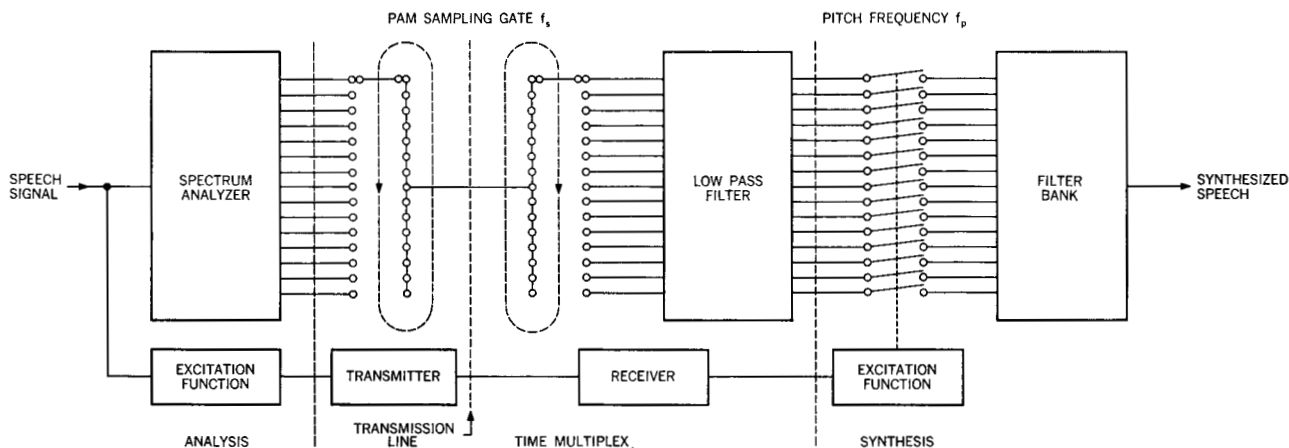
**455**

**Figure 1** PAM transmission of vocoder signals.

constant artificial pitch can be highly intelligible, naturalness and discriminability among speakers may suffer considerably.[4,7]

### 1.22 Multiplexing

The second problem stems from the fact that the vocoder analyzer describes the input speech signal in terms of a relatively large set of signals. One of these signals carries the pitch information, and a set of ten or more signals is required to transmit a description of the spectral envelope of the speech signal. Since transmission of all these signals over physically separated channels seems to be out of question for all conceivable applications, frequency multiplexing or time multiplexing transmission systems will practically always be necessary.

Frequency multiplexing techniques yield a compound signal which requires less bandwidth than a normal telephone signal would need. A given telephone transmission channel might thus handle more simultaneous conversations. The main difficulty with this approach is cost. It is not easy to devise a low-cost carrier frequency system for 20 channels having individual bandwidths of 25 Hz without wasting bandwidth with double-sideband transmission or with large safety gaps between adjacent channels. The present state of the art seems to indicate that a vocoder combined with a frequency multiplexing transmission system is economically feasible only for very expensive channels, e.g., transoceanic lines.

Time multiplexing techniques lend themselves to a quantized, digital description of the speech signal, e.g., the PAM system shown in Fig. 1 may be easily modified to a PCM system with fully quantized signal transmission. But applications for vocoder systems with this type of signal transmission have also been very limited because of the high costs of the over-all system and the lack of sufficiently large digital networks. In special cases, e.g. mili-

tary applications with security problems, transmission over telephone networks of time-multiplexed and digitized vocoder signals can be justified. Additional interface equipment is then required which transmits the digital signals over the telephone line in a representation that is matched to the typical properties of such a line.

### ❧ 1.3 Proposed solutions

### 1.31 The voice-excited vocoder

The obvious solution to overcome one part of the pitch problem described in Section 1.21 is to look for a way of specifying the excitation function that does not require voiced-unvoiced decisions. This approach leads to two different solutions: the well-known voice-excited vocoder[11] and the integrated vocoder to be described in this paper. The two solutions differ not only in their way of handling the excitation signal but also in their optimum way of multiplexing the transmission signals. These specific features will be stressed in the following, while other important problems, like "spectral flattening" of the excitation signal, are neglected here because they lead in all vocoder types to much more similar solutions.

The voice-excited vocoder is most suitable for a frequency multiplexing transmission system. Instead of making voiced-unvoiced decisions and pitch measurements, a sufficiently large part of the low-frequency components of the speech signal is transmitted directly to the synthesizer. This is the so-called "base-band." By means of a special nonlinear device which increases the number of zero-crossings of the base-band signal, the excitation synthesizer spreads the original spectrum of the base-band over the whole spectrum range of the speech signal. This band-spread signal is used for the excitation of the different spectrum channels.

The base-band signal itself may also be utilized directly

as one component of the output speech signal. For moderate bandwidth reductions the speech quality of a base-band vocoder can be excellent, e.g., if the base-band goes up to 2 kHz and only speech components from 2 to 10 kHz are encoded, thus allowing transmission of speech components up to 10 kHz over a normal telephone channel with a passband from 300 to 3300 Hz. A frequency-multiplexing transmission system must be used for this special purpose.

### 1.32 The integrated vocoder

As its name implies, the integrated vocoder integrates a channel vocoder with its accompanying multiplex transmission system. This becomes possible through a special approach to the problem of pitch detection and measurement. As in the voice-excited vocoder no discrimination is necessary between voiced and unvoiced sounds. In contrast to the conventional form of the voice-excited vocoder, the integrated vocoder does not use continuous transmission of analog signals but instead utilizes pulse techniques. In the following section this approach will be explained and is shown to lead to significant reductions in the over-all hardware system.

## 2. The integrated vocoder

### • The special excitation concept

The integrated vocoder is a special type of channel vocoder. The introduction of the concept starts with a description of the excitation method used in this system. While in the conventional excitation scheme only voiced sounds are assumed to have a fundamental frequency $f_p$, in the integrated vocoder the unvoiced sounds are also ascribed to have an "ostensible fundamental" frequency. The fundamental of voiced sounds is given by the vibration frequency of the vocal cords of the speaker; the "ostensible fundamental" of unvoiced sounds may arbitrarily be defined with the restriction that it has to allow for the derivation of an adequate excitation signal. The integrated vocoder requires the excitation signal to be defined as a pulse train during both the voiced and the unvoiced portions of the signal. During voiced portions the distances between successive pulses are random within certain time limits. Thus no discrimination is necessary between voiced and unvoiced sounds as in other vocoder systems. Hybrid forms of voiced and unvoiced sounds can be excited; the distances between successive pulses will then fluctuate around a mean value corresponding to the pitch period. The use of pulse excitation for voiced and unvoiced sounds is not new in itself.[5] But here the essential idea in the chosen definition of the excitation signal is the required minimum distance between excitation pulses. In the integrated vocoder, it will be shown that it is necessary to fix this minimum distance corresponding to the highest anticipated pitch frequency
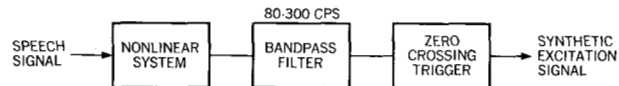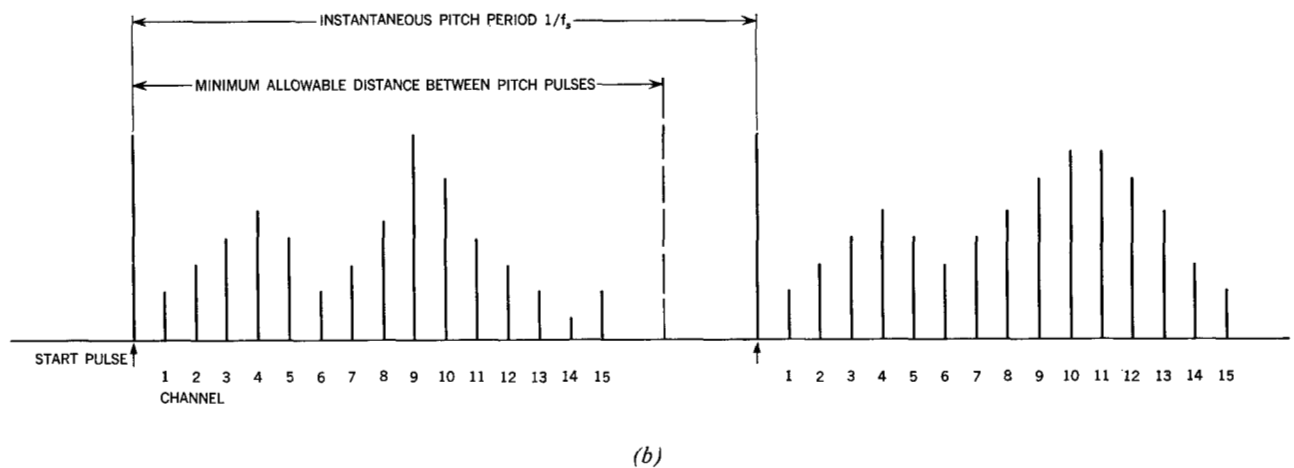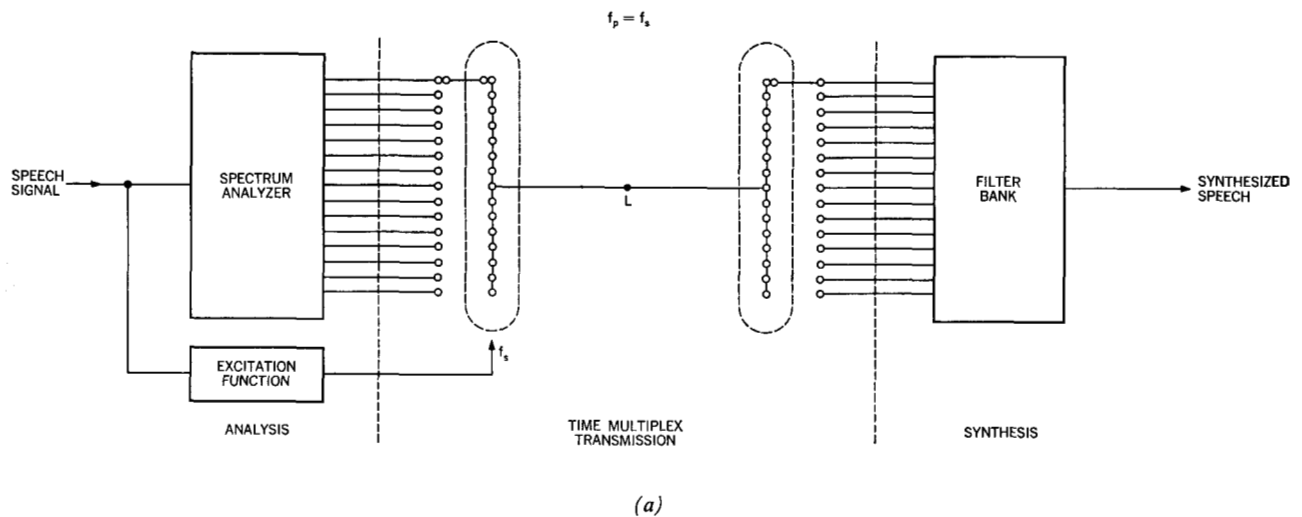
80-300 CPS

**Figure 2** A simple excitation generator.

of the potential users of the system. Experience shows that it is indeed possible to maintain minimum distances between excitation pulses for most male speakers corresponding to about 150 Hz without getting into serious problems with the excitation of unvoiced sounds.

The generation of the excitation pulse train may be performed in several ways. At the present state of the art a compromise between cost and complexity versus performance seems to be necessary which demands a very high price for increases in performance, as soon as the quality of the input speech signals is not ideally high.[12] In accordance with the aim of this paper to present a vocoder system which attempts to be a "minimum hardware" system, the simplest device possible will be described for the generation of the excitation pulse train (Fig. 2).* The input to this device is the speech signal which may have been restricted to the telephone band. The speech signal is first distorted by a nonlinear system. Then the spectrum of the distorted speech is restricted to the range of possible fundamental frequencies by a bandpass filter. At each positive (or negative) zero crossing of the output of this bandpass filter, an excitation pulse is generated. There are two reasons for the nonlinear distortion. During voiced sounds the fundamental frequency has to be emphasized in relation to its harmonics or even regenerated, if due to a restriction of the speech signal to the telephone band the original fundamental has been eliminated. During unvoiced sounds, spectral energy from higher frequency bands is shifted to the range of possible fundamental frequencies, thus yielding a noise output level of the bandpass filter sufficiently high to insure adequate pulse generation. The upper limiting frequency of the bandpass filter ensures a minimum distance between successive excitation pulses. This will be shown to be imperative for the integrated vocoder described in the next section.

---

* The extremely simple pitch extractor shown here is a good example of the fact that for a given speech quality a certain hardware approach may be advantageous but is completely inadequate for a next higher level of speech quality. The pitch extractor described solves a problem which was very annoying with the conventional pitch measurement circuits. The described solution is relatively insensitive to spurious harmonics in the pitch band. They may create additional pulses, but as the pulse train is used directly for excitation, even a relatively weak fundamental is preserved in the pulse positions. Thus the subjective pitch of the synthesized speech signal will not jump occasionally to the first or even the second harmonic, as happened with some conventional pitch extractors. Quite the contrary of this adverse effect occurs: The additional pulses will help to reduce a certain roughness of the vowels by giving additional perceptive cues for the correct harmonic structure of the original vowels. But all this is true only up to a speech quality level which is slightly below "good telephone quality." If a higher speech quality is required and is also possible with regard to the other components of the vocoder system, the described pitch extractor has to be replaced by more refined approaches.[13, 14]

$f_p = f_s$

SPEECH SIGNAL → SPECTRUM ANALYZER

EXCITATION FUNCTION

$f_s$

ANALYSIS

L

TIME MULTIPLEX TRANSMISSION

FILTER BANK → SYNTHESIZED SPEECH

SYNTHESIS

(a)

INSTANTANEOUS PITCH PERIOD $1/f_s$

MINIMUM ALLOWABLE DISTANCE BETWEEN PITCH PULSES

START PULSE

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
CHANNEL

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

(b)

**Figure 3** The integrated vocoder. (a) Block diagram. (b) Typical signal pattern on the transmission line [L in diagram (a)] for a 15-channel vocoder.

● *2.2 Basic structures*

*2.21 The integrated vocoder with PAM signal transmission*

A block diagram of a vocoder combined with a conventional PAM-multiplex system has been shown in Fig. 1. The step from this conventional scheme to the integrated vocoder[6] as shown in Fig. 3 becomes possible if one considers the limited sensitivity of the human ear to time delays between the different spectrum channels[7] of a vocoder. Thus, instead of parallel operation, the modulator switches in the synthesizer may be driven serially without any noticeable additional distortions. The transition from the conventional scheme of Fig. 1 to the integrated vocoder of Fig. 3 is made by using the variable instantaneous fundamental frequency $f_p$ as the sampling rate $f_s$ of the PAM system. Then the demodulation of the transmitted pulses is no longer necessary. Additionally,

the sampling switch at the transmitter side performs the function of the modulator switches of the synthesizer. The receiver now needs only a distribution switch for serial-to-parallel conversion and a set of synthesizer bandpass filters.

Figure 3(b) shows a typical signal pattern on the transmission line. It makes clear that a special type of pulse amplitude modulation is being used. Instead of having a constant time frame for all pulses, only the distances within a sampling group are fixed. The distances between corresponding pulses of successive sampling groups are continuously variable and describe the instantaneous pitch frequency. Two reasons for a minimum distance between pitch pulses, as indicated in Sec. 2.1, become evident here. In order to keep the message rate on the transmission line low, the number of pitch pulses for unvoiced sounds has to be kept as low as possible. Furthermore, the distance between pitch pulses must not be shorter than a sampling cycle.

**458**

## 2.22 The integrated vocoder with PCM signal transmission

In many practical applications a PAM system, as it is described in the preceding section, will form only an intermediate stage in a superposed digital transmission system. The system shown in Fig. 3 can easily be modified for a special type of PCM transmission. Only one coder and one decoder would be necessary. As in the corresponding PAM system, coding would be done on a start-stop basis triggered by the excitation pulses.

Experience shows that eight amplitude quantization levels are sufficient for the encoding of the PAM signal. Thus for digital transmission three times as many pulses are necessary to describe a sampling cycle. As a consequence of the larger size of the resulting pulse groups additional marker pulses within these groups may be necessary to ensure synchronization between transmitter and receiver.

Additional quantization in the time domain requires a sufficiently high quantization frequency in order to avoid pitch distortions. It will allow for the time multiplexing of several speech signals over one channel. If pitch is digitally encoded, there is no more need for pitch-synchronous signal transmission, and the integrated vocoder then tends to become a normal channel vocoder. An interesting exception to this rule is the connection of an integrated vocoder to a digital computer, as is described in the next section.

## 2.23 The integrated vocoder as a special I/O device

In speech processing studies with a computer a valid approach is to reduce the number of specialized, peripheral hardware attachments to the computer as far as possible. This policy not only cuts development time for these attachments, but also increases the reliability of the over-all system. If a vocoder is to be utilized as an I/O device, the concept of the integrated vocoder offers considerable hardware savings on the output side of the system as compared to the conventional vocoder.[9] The outlines of a system combining the integrated vocoder with a computer will be described by a specific example. For speech processing work in our Vienna laboratory we connected an integrated vocoder to an IBM 1401.

The speech input part of this experimental system is shown in Fig. 4. Samples of the 20 spectrum-channel signals are taken at a rate controlled by the computer. The samples are quantized to eight levels. The difference of amplitude between adjacent levels is 5 dB. The IBM 1401 computer with I/O attachment can read six bits in parallel at a rate up to 87 kHz. Thus with two analog-to-digital converters in parallel and a distribution switch with ten positions, it is possible to read one sample of each of the 20 channels in an interval of 115 $\mu$s.
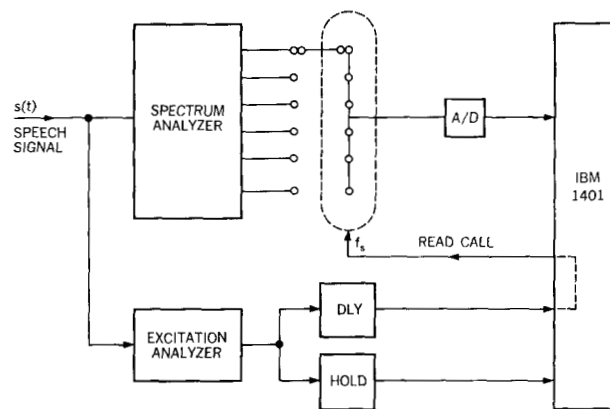
At the beginning of a recording, spectrum data as well



**Figure 4** Speech input to an IBM 1401.

as pitch data are read into the 16 K core storage of the computer, which can accommodate about 8 sec of connected speech. As soon as the core store is full, reading is stopped and the data are transferred to tape. A special feature of the voice input system is the HOLD circuit. It locates distances between excitation pulses that are longer than a given time.

The computer uses these breaks to clear the core store. As breaks of sufficient duration occur rather often in normal speech, it is thus possible to record any length of connected speech.

The data thus stored are further processed. In any case spectrum samples that do not differ significantly from the preceding ones are eliminated. The reduced data are stored on tape and may be printed out for visual inspection, e.g., in a form as shown in Fig. 5. The printout shows cross sections of the 20 spectrum signals in the left column and portrays the pitch contour in the middle part. The right column specifies time. Such printouts not only form a convenient basis for handmade corrections or processing of the stored speech data, but they also allow for observation of the effects of programmed processing steps.

The speech output part of the system is shown in Fig. 6. The computer provides spectrum data in the form of six-bit words. As in the operation of a spectrum analyzer an electronic switch distributes the data to the corresponding synthesizer channels. The distribution switch is actuated by a sequence of ten WRITE calls which follow one another at intervals of 11.5 $\mu$s. The output rate for the spectrum data is controlled by the pitch information. The distances between two successive groups of WRITE calls are equal to the distances between successive excitation pulses. The digital spectrum data are converted separately for each channel to analog signals and are then fed directly to the respective synthesizer bandpass filters.

As previously mentioned, the proposed encoding scheme requires a higher channel capacity for transmission of the vocoder signals than the conventional technique. This is

**459**

| ADDRESS | AGGREGATE FUNCTION | EXCITATION FUNCTION | PAUSE SECONDS |
|---|---|---|---|
| 005001000 | 00000 00000 00000 00000 | 914 | .010 |
| 005001013 | | 788 | .019 |
| 005001016 | | 777 | .028 |
| 005001019 | | 807 | .037 |
| 005001022 | | 809 | .047 |
| 005001025 | | 736 | .055 |
| 005001028 | | 498 | .061 |
| 005001031 | | 794 | .070 |
| 005001034 | | 786 | .079 |
| 005001037 | | 807 | .088 |
| 005001040 | | 629 | .095 |
| 005001043 | | 652 | .103 |
| 005001046 | | 749 | .112 |
| 005001049 | | 807 | .121 |
| 005001052 | 00000 00000 02000 00000 | 757 | .130 |
| 005001065 | | 609 | .137 |
| 005001068 | 222 | | .139 |
| 005001071 | | 426 | .144 |
| 005001074 | 00000 00000 04222 10000 | 800 | .153 |
| 005001087 | | 881 | .163 |
| 005001090 | | 836 | .173 |
| 005001093 | 00100 00002 23222 00000 | 481 | .178 |
| 005001106 | | 646 | .186 |
| 005001109 | | 519 | .192 |
| 005001112 | 67642 22223 33101 00000 | 717 | .200 |
| 005001125 | | 706 | .208 |
| 005001128 | | 688 | .216 |
| 005001131 | | 683 | .224 |
| 005001134 | 67775 44556 45323 20000 | 680 | .232 |
| 005001147 | | 673 | .240 |
| 005001150 | | 676 | .247 |
| 005001153 | 67765 44465 45323 10000 | 679 | .255 |
| 005001166 | | 695 | .263 |
| 005001169 | | 725 | .271 |
| 005001172 | 76611 00131 00000 00000 | 687 | .279 |
| 005001185 | | 665 | .287 |
| 005001188 | | 675 | .295 |
| 005001191 | 64310 01220 00000 00000 | 664 | .302 |
| 005001204 | | 689 | .310 |
| 005001207 | | 710 | .319 |
| 005001210 | 76422 24200 C000C 00000 | 701 | .327 |
| 005001223 | | 693 | .335 |
| 005001226 | | 689 | .342 |
| 005001229 | 67622 45200 CC00C 00000 | 674 | .350 |
| 005001242 | | 656 | .358 |
| 005001245 | | 664 | .365 |
| 005002000 | | 663 | .373 |
| 005002003 | 67743 34542 24123 10000 | 646 | .380 |
| 005002016 | | 638 | .388 |
| 005002019 | | 635 | .395 |
| 005002022 | 67752 21344 56335 20000 | 638 | .402 |
| 005002035 | | 641 | .410 |
| 005002038 | | 650 | .417 |
| 005002041 | 67741 21113 56445 20000 | 653 | .425 |
| 005002054 | | 665 | .432 |
| 005002057 | | 683 | .440 |
| 005002060 | | 705 | .448 |
| 005002063 | 75520 00C00 23202 00000 | 694 | .456 |

**Figure 5** Specimen printout of reduced spectral data.

not true for the storage of speech signals according to such a scheme in a computer. Since spectral samples might be repeated for several excitation cycles, the necessary storage space can easily be reduced to about 2000 bits/sec of stored speech.

### 3. Integrated vocoder vs conventional vocoder

The integrated vocoder of Fig. 3 may be favorably compared with the conventional vocoder of Fig. 1 using PAM transmission, when the number of hardware elements for both systems is considered. As compared to Fig. 1, in the integrated vocoder the PAM transmission channel and the vocoder have one set of modulators in common. There are no lowpass filters and a second set of modulators is unnecessary. Transmission of pitch information does not require a separate channel. The same hardware savings apply, if the integrated vocoder is used as special I/O equipment for a computer.

There is no difference in speech quality between an integrated vocoder and its conventional equivalent. Experience shows that switching of a channel vocoder between serial and parallel pulse excitation, corresponding to the operation of the integrated and conventional vocoder, does not lead to perceptible changes of the speech output signal.

In order to estimate the efficiency of the integrated vocoder in encoding speech signals (Table 1) the following experimental conditions are assumed:

(1) Only male speech with a pitch pulse below 150 Hz, and restricted to a 4 kHz bandwidth, will be transmitted.

(2) Since 16 channels seems to be an average operating condition for conventional vocoders, this number is adopted as a standard in the following examples.

(3) The data rates are specified for each of the following hypothetical systems on the assumption that they will provide the approximate equivalent of "telephone quality." This assumption is difficult to verify because of the lack of a suitable procedure for the objective measurement of

E. H. ROTHAUSER

speech quality. Since we are concerned only with the order of magnitude of the specified data rates, a small variation in speech quality or proposed data rates does not materially affect the comparative analysis in Table 1.

The table was compiled under the above restrictions and cites conventional transmission systems with typical data rates.

The table below, while not intended to offer precise data, shows that the required channel capacity for the integrated vocoder lies between that of the conventional channel vocoder and the respective normal PAM or PCM system. Even when the transmission of female voices with pitches up to 300 Hz is considered, the integrated vocoder needs only half the channel capacity of the equivalent normal PCM system. If the input speech signal is not limited to 4 kHz, thus allowing for higher output speech quality, the integrated vocoder and all other vocoder systems will gain more than another factor of 2 when they are compared to normal PAM or PCM.



**Figure 6** Speech output from an IBM 1401.

**Table 1** Encoding of speech signals: integrated vocoder compared with some other systems

a) *Conventional vocoder, conventional PAM*
Spectrum samples required, including, e.g.,
    four markers $(16 + 4) \times 30$ Hz bandwith = 600 Hz
Required pitch sampling frequency
    $(2 \times 150$ Hz) = 300 Hz
Required pulse frequency $\cong 1$ kHz

b) *Integrated vocoder with PAM (Start-Stop)*
Sixteen spectrum samples plus four
    markers in one group = 20
Maximum pulse frequency
    (150 Hz $\times$ 20 samples) = 3 kHz

c) *Normal PAM (no vocoder)*
Required pulse frequency $(2 \times 4$ kHz) = 8 kHz

d) *Conventional vocoder, conventional PCM*
Spectrum samples (as specified in (a))
    $600 \times 3$ bits = 1800 Hz
pitch (as specified in (a))
    $300 \times 7$ bits = 2100 Hz
Required pulse frequency $\cong 4$ kHz

e) *Integrated vocoder with PCM (Start-Stop)*
Spectrum samples (as in (b)) $20 \times 3$ bits = 60
Maximum pulse rate (150 Hz $\times$ 60 samples) = 9 kHz

f) *Normal PCM (no vocoder)*
Required pulse frequency:
    $2 \times 4$ kHz $\times$ 5 bits = 40 kHz

g) *Integrated vocoder connected to digital computer*
Storage space necessary in computer memory for intermediate storing of speech signals (after deletion of spectral samples which do not differ sufficiently from preceding spectral pattern) = 2 kbits/s

**References**

1. H. W. Dudley, "The Vocoder," *Bell Lab. Rec.* 18, No. 4, 122 (1939); also, H. W. Dudley, "Remaking Speech," *J. Acoust. Soc. Amer.* 11, 169 (1939).
2. R. J. Halsey and J. Swaffield, "Analysis-Synthesis Telephony with Special Reference to the Vocoder," *J. I. E. E.* 95, 391 (1948).
3. J. L. Flanagan, *Speech Analysis, Synthesis and Perception,* Springer Verlag, Berlin, 1965.
4. M. R. Schroeder, "Vocoders, Analysis and Synthesis of Speech," *Proc. IEEE* 54, 720–34 (1966).
5. Ch. Schwiedernoch, "Entwicklung eines Vocoders," Doctor's Thesis, 1955, Inst. for Telecommunications, Technical University, Vienna.
6. E. H. Rothauser, "Ein Impulsverfahren zur Sprachübertragung nach dem Vocoderprinzip," Doctor's Thesis, 1960, Institute for Telecommunications, Technical University, Vienna.
7. E. H. Rothauser, "Dependence of Speech Quality on Transmitted Information Data in a Band Compression System." *Information Processing,* 354 (1962).
8. E. Paulus and E. H. Rothauser, "Normalization of Spectral Information in a Pulse Excited Channel Vocoder System." Abstract in *J. Acoust. Soc. Amer.* 36, 2002 (1964).
9. E. Paulus and E. H. Rothauser, "A Pulse Excited Vocoder System." Abstract in *J. Acoust. Soc. Amer.* 36, 2002 (1964).
10. E. H. Rothauser, E. Paulus and M. Fleck "Regeneration of Formant Coded Speech Using a Channel Vocoder Synthesizer," *Proceedings of the 5th Congress on Acoustics,* Liege, Belgium, 1965.
11. D. Schroeder, "A Vocoder for Transmitting 10 kc/s Speech over a 3.5 kc/s Channel," *Acoustica* 10, No. 1, 35 (1960).
12. A. M. Noll, "Short-Time Spectrum and 'Cepstrum' Techniques for Vocal Pitch Detection," *J. Acoust. Soc. Amer.* 36, 296–302 (1964).
13. Knauft, Lamparter and Spruth, "Some New Methods for Digital Encoding of Voice Signals and for Voice Code Transmission," *IBM Journal* 10, 244 (1966).
14. J. Tierney, B. Gold, V. Sferrino, J. A. Dumanian, and E. Aho, "Channel Vocoder with Digital Pitch Extractor," *J. Acoust. Soc. Amer.* 36, 1901–05 (1964).