

MICROPHONEMES AS FUNDAMENTAL SEGMENTS OF SPEECH WAVE
PRIMARY SEGMENTATION - AUTOMATIC SEARCHING FOR MICROPHONEMES

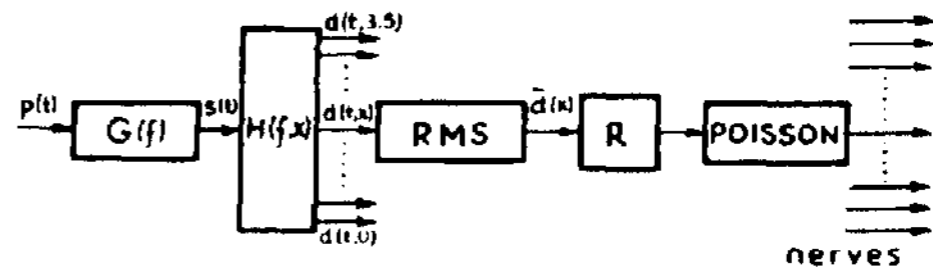
Mgr inż. Andrzej Dziurnikowski
Institute of Organization, Management and Control Sciences, Warsaw, Poland

Abstract

This paper concerns the stage of acoustic analysis in speech recognition of Speech Understanding System SUSY; subsystem AKORD, Acoustic analysis in this system is based on dynamic spectral analysis of the speech wave. It corresponds with the function of the hearing organ expressed by an analog model of human ear. Acoustic analysis in the AKORD system uses an FFT algorithm but is not strictly based on accepting constant intervals of 10-20 ms; analysis does not cause inaccuracy at the level of parametrization because of a considerable dispersion of parameters of particular segments. A conclusion was made in the AKORD system that a properly conducted process of segmentation will help in solving some problems of speech wave analysis, particularly recognition and time compression. This paper attempts to choose the primary segment (microphonem) in the most optimum way and to indicate the algorithm of automatic segmentation based on microphonemes as the dynamic segments.

Introduction

System AKORD was designed for research work concerning the analysis and synthesis of speech waves by means of a ZAM-41 computer (1,2,3) - during the initial phase of research that system was helpful in defining the most essential problems of speech wave analysis- It also made it possible to inculcate and test some algorithm models of speech analysis. However, before accepting a concrete algorithm model of speech analysis it was necessary to realize different functions of the hearing organ. Those functions are expressed by an analog model of the human ear (4). Generally, the track of a given sound in the ear may be presented as follows: the external ear, the ear-drum, the internal ear and its bones (malleus, incus, stirrup). The movements of the stirrup bone cause translocation of cochlea septum : -its resilience is not constant through its whole length (4,5). That vibration is subsequently transmitted through the nerve cells on the septum to ~30 000 auditory nerve fibrils. Figure 1 shows analog model of human ear presented by P.Kolers (4). Although this model is only approximate, a comparison of the characteristics obtained experimentally showed a considerable convergence (4). We accept the conclusion, with no deeper consideration of the ear model presented here, that the internal ear is the centre of sound spectral analysis $S(f)$. Its results are made average



$P(t)$ - external auditory duct pressure, $G(f)$ - linear block transforming pressure $P(t)$ into voluminal translocation $s(t)$ of the stirrup, $H(f,x)$ - this function connects translocation $s(t)$ with translocation of septum $d(t,x)$ at the point x [cm] from the stirrup, RMS - realises the mean-square function of input signal $d(t,x)$ for short duration (~ 10 ras).

Figure 1

in time for the period T_s (~ 10 ms). The model of the auditory organ presented here is a basis in this paper, although its imperfection is realised. But it is not a novelty. The majority of studies concerning the analysis of speech waves have been based on the spectral analysis of the examined signals (6). Some of them used the dynamic spectral analysis, yet they were strictly based on accepting constant intervals of 10-20 ms. Although the acceptance of that constant interval made it possible to realise some functions of analysis more simply, it caused some inaccuracy at the level of parametrization due to a considerable dispersion of parameters of particular segments (12). This fact resulted in substantial difficulties connected with the segment extraction for recognition. It seems relatively easy to find a relevant set of distinctive phoneme features and so the segment corresponding with phonemes would be the most suitable for the recognition (8). The number of segments determining particular classes also points to the same conclusion; in the case of phonemes there would be no more than 40. Yet, difficulties connected with phoneme segmentation as well as those resulting from the lack of univocal dependence between the parameters of some concrete realization of the stochastic process and the process itself, make it too difficult to work out a highly efficient automatic system of speech recognition. Some authors consider it simply impossible to develop studies based on phoneme analysis and recognition (7,8) pointing to the practical difficulties of segment extraction. It is possible to conclude that a solution of this problem will become a basis for the development of some methods of speech analysis by means of particular speech sounds of a given natural language.

The process of segmentation when properly conducted will help in solving some other problems of speech wave analysis e.g. time compression, coding, code ciphering and recognition. Many scientists have noted the role of segmentation in the process of speech wave analysis (6,7,8,9). Yet, they were too strict in their choice of segments derived from the analog model. Naturally, therefore, they did not obtain the type of results which could be acquired through a segmentation based on the segments chosen in the most optimum manner. This paper attempts to indicate the algorithm of automatic segmentation based on microphonemes as the dynamic segments.

Microphoneme as the Fundamental Segment of Voiced Speech-Waves

We are concerned with voiced speech, since only for that class of speech is possible to extract fundamental segments in a simple way, keeping in mind the purpose this extraction could serve. It is also possible to resolve the stochastic signal into the class of stationary signals (as considered in frequency domain) by means of some simple methods, for examples, Fourier's transformation. It corresponds with the model of the ear presented above. If we consider the objective function in speech wave analysis the following requirements arise:

- for the purposes of information recognition (its content) it is necessary to accept segment S_k^l , $l = 1, 2, 3, \dots$, making it possible to realize the set of relations m_i ; for the correct classification, while considering the subset of parameters $t_{Ri} \equiv t$ defined for a given segment. It may be presented as the following condition:

$$1 \quad m_i: t_{k,i}(S_k^l) \xrightarrow{t} t_{Ri}; \quad l=1,2,\dots$$

meaning that there is a relation m_i which makes it possible to transfer the set of parameters defined for any "l" realization of the class K_i , $i=1,2,\dots,m$ ("l" segment S_k) into the set of parameters t_{Ri} relevant to class K_i .

- for the purposes of compressing and coding a speech wave, it is necessary for the extracted segments S_k^l to satisfy condition No 1. Considering the time compression, the best solution (if possible) would be to extract such segment S_k so that they will be realization of the quasi-periodic waveforms within a given class K_i . That quasi-periodicity should be conceived as the periodicity which is discrete from the point of view of segments (time quanta of real signals). In other words the segment S_k would constitute a string of the segments S_k^j

$$2. \quad S_k = \langle S_k^j \rangle; \quad j = 1, 2, \dots$$

In the time domain t (continuous) it would be possible to present it in the following way:

$$3 \quad \bar{S}_k(t) = \prod(t_k - t) \cdot \prod_{i=1}^n \left\{ S_k(t) \cdot \prod_{r=0}^i \left(\sum_{\tau=0}^i T(\tau) - t \right) \cdot \prod_{r=0}^{i-1} T(\tau) \right\}$$

where: \bar{S}_k means the segment S_k defined for a period of time t_k (conceived as a set of

samples of a signal in the interval of $[0, t]$, n is a function of time t and it is calculated according to the following condition

$$4 \quad \sum_{j=1}^n T(j) \leq t$$

It is illustrated in Figure 2.

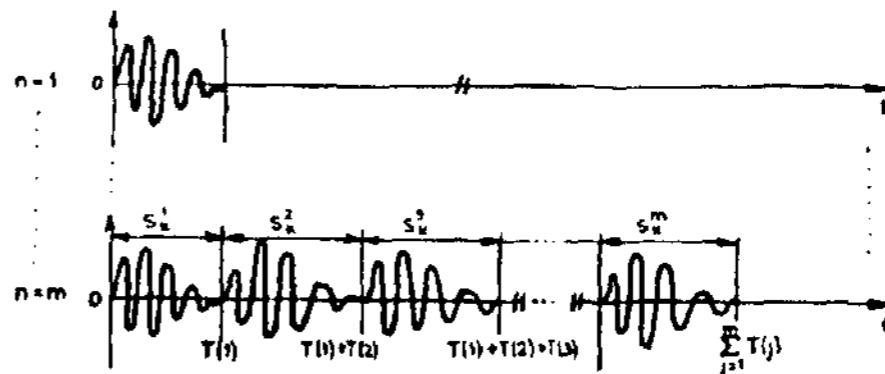


Figure 2

Condition No 1 determines the time for the segment S_k ; it is t_k . Considering the function of compression and coding, i_k should be the maximum length of time of the segment S_k satisfying condition No 1. Fundamental segments S_k^j forming quasi-periodicity of the segment S_k will be called primary segments. Each of them is of different duration $T(j)$ and they form quasi-periodicity within the segment S_k : - see formula No 3. With regard to the function of speech wave time compression, primary segments should satisfy the following condition:

$$5 \quad \bar{S}_k^{T(j)} \equiv \min_t S_k$$

This condition states that a number of primary segments n should be maximum.

- for the purposes of speaker identification it is necessary to indicate such segment S_k^l ; $l=1,2,3,\dots$, for the subset of parameters afterwards defined (identifying the speaker) to satisfy condition No 1 in consideration of a given class of speakers K_p and also:

$$6 \quad \bar{S}_k^l \equiv t_{IP}^l; \quad l=1,2,\dots$$

This condition indicates that the subset t_{IP} of parameters identifying the speaker specified in consideration of a freely accepted segment S_i (free in the sense of time) is the equivalent of the subset of parameters t_{IP} . These parameters are defined by means of one of the possible methods and the subset is calculated with some optional duration of speech for a given speaker, belonging to the class K_p .

- in the process of reducing the noise the most essential parameters are: the speech spectrum and particularly the average spectrum of a speech wave. Let us denote the subset of those parameters: t_s . With this approach it is necessary to define the segment S_0 such as:

$$7 \quad \bar{S}_0^l \equiv t_s^l; \quad l=1,2,\dots$$

This condition means that it is necessary to indicate such a segment S_0 for the calculated parameters of the average speech spectrum, to fix univocally the average spectrum parameters calculated for

a speech wave of any duration and at any moment. In this case it is necessary to indicate a segment which would satisfy condition No 7 and which would at the same time be the most favourable considering the minimal time interval. This approach shortens the process of indicating the parameters ts considerably. With regard to these considerations, it is necessary to state that in the process of speaker identification and noise reduction, the choice of the best segment should be carried out considering its length of time. The shorter the segment S_i or S_o , the larger the possibilities of realizing both processes. Therefore, it will be possible to indicate the minimal segment of a speech wave and through that to identify the speaker. The segment chosen in this way must satisfy condition No 6 and 7, corresponding with condition No 1 in that the set of parameters is indicated for a given speaker p . But for the purposes of compression and codification it is necessary for the segment SK to be maximum and also to satisfy condition No 1 as in the case of recognition. In this case the aim is also to find the segment $SK(j)$ which would satisfy conditions No 1,3 and 5. We must search for a primary segment satisfying the following condition :

$$8 \quad \int \overline{S_n^{T(j)}} \equiv \int \overline{S_n^j} ; j=1,2,\dots$$

within a given class K_i , $i=1,2,\dots,m$. Because it is well known that the acceptance of the segment S_n fixes the number of classes K_i (Ct^fCS*), the minimization of the number of classes K_i is the additional condition that should be taken into consideration while choosing the segment SR . If we accept syllables and words as the segment S_n , the number of classes in which it would be necessary to include the analysed segments would be tens and hundreds of thousands. Therefore it seems that it is necessary to accept a phoneme as the segment S_n or a primary segment which is the equivalent of phoneme considering the condition No 1. The number of classes with a segment so defined is $m \approx 40$ for the Polish language. With this number of classes a given phoneme as a segment S_n^* is the longest segment and it would fix the segment SK within the class K_i . Yet in spite of its usefulness for the purposes of recognition and compression, it would not be the most favourable segment because the condition 5 would not be satisfied. And only finding the primary segment that would meet all the requirements would make it possible to state the existence of the most favourable choice of the speech wave segment, in terms of recognition of the content and time compression of a speech wave.

The number of classes according to condition No 8 would be retained. The additional effect of such an approach on speech wave analysis would be the obtaining of a dynamic analysis of parameters defined within the primary segment SK . Concerning the two remaining goals, speaker identification and reduction of noise, we may accept the primary segments SK as fundamental

only if it is possible to determine the parameters ϵ_r and ϵ_a^* on the basis of any parameter combination estimated within a given segment SK . In the case of speaker identification it is advisable to estimate parameters for the dynamically analysed speech wave. Speed and changes in the larynx vibrations are some of the dynamic parameters characterizing the speaker and his articulation : let us denote the length of time of the larynx tone : TK . It is a function of time for a given voiced speech wave. Therefore it is advisable to accept the primary segment SK ; its duration is $T(j) - TK$. It is different from the majority of algorithms assuming $T(j) = \text{const}$. This approach would also give the phase accordance of the analysed segments (11). The primary segment is based on pitch period equal to TK corresponding to the analog ear model ($TK \approx 10 \text{ ms}$) and is called a microphoneme in this paper. Our studies have entirely confirmed the pertinence of the accepted reasoning and also the fact that a microphoneme satisfies all the above conditions (12,13).

Primary Segmentation Automatic Searching for Microphonemes

Since a voiced speech wave results from the vibration of the vocal cords and is therefore only approximately periodic and since it depends on the individual characteristics of the speaker, it is not possible to assign definitely the frequency which could be considered fundamental for all segments SK of a speech wave. In other words, the length of time of microphonemes is : $T(IW \text{ const})$. The purpose of primary segmentation is the extraction of microphoneme sequence from the continuous speech wave of a given speaker by means of an algorithm realized as a computer programme. The solution of the problem concerning the extraction of quasi-periodic segments of a speech wave is complicated by the fact that not only quasi-periodicity is troublesome but also because the speech waves vary in amplitude and shape. In assigning microphonemes the algorithm uses some of the methods and remarks included in Reddy's algorithm : Pitch Period Determination (10). It also completes Reddy's algorithm with the elements connected with the acceptance of the microphonemes as a primary segment and with other parameters appropriate for Polish speech as well as the accepted method of their representation (1). The algorithm FPD was used as an element to indicate the location neighbourhood of the expected end of a given microphoneme. The algorithm extracts microphonemes and indicates the places of expected microphoneme ends of the analysed speech wave and also gives a proportional indication of the voiced signal content in the analysed speech wave*

As we have already assumed, we shall be searching for the microphonemes only in that part of the speech wave which is voiced according to condition No 3. Therefore the first step of the present segmentation

is to find and indicate voiced segments of the analysed speech wave. Because the whole procedure of primary segmentation operates on the time signal of the speech wave, the first step makes use of amplitude and frequency criteria. At this stage these criteria make it possible, as well as in the PPD algorithm, to distinguish three fundamental kinds of signals: silence, noise and a quasi-periodic signal. The algorithm has been adapted to the processing of the input signal which is first quantized. It assumes the frequency of sampling: $f_p = 12$ kHz, which makes it possible to represent the primary speech wave quite well. The analog-digital converter working in the AKORD system performs signal quantization on 256 levels, enabling the notation in the form of numerical data ranging from $\epsilon - 128$, $+127$. This fact has its reflection in the values of coefficients indicating the amplitude thresholds. Initially the following assumptions have been accepted as in Reddy pro:

1. In any examined speech wave (with the same speaker) the length of neighbouring microphonemes cannot differ more than 20%.
2. When the frequency of the larynx tone is contained within 70-450 Hz the frequency test $f_p = 12$ kHz will produce a length of microphoneme within the bounds of 26-171 samples.
3. The examined speech wave will be also accessible in the possible process of correction. This assumption is easy to realize in the case of computer processing.

The succeeding speech wave segments of $D = 256$ samples (21,3 ms) undergo the amplitude criterion. The accepted length D results from the fact, that in the segment S_l longest microphoneme s_k should be contained ($\max(T(j) \cdot f_p) < D$). The segments S_l , $l=1,2,\dots$ are accepted as silence segments if the maximum amplitude within a given segment does not exceed the threshold value $\delta_1 = 8$, δ_1 -times. In other words a variable α^l

$$\alpha^l = \# \{x_i : x_i \gg \delta_1 \wedge x_i \in S_l^l\}$$

(where $x_i \in S_l^l$ means the samples of the segment l ; its length is D whereas $\#\{ \}$ means the cardinality of the set) assigns some segments S_l accepted as silence in the case of the condition $\alpha^l < \delta$. In this algorithm $\delta = 5$ was accepted. At the next stage only those segments S_l are accepted for which $\alpha^l \geq \delta$. All those segments (of the length D) having $\alpha^l \geq \delta$ where $\delta_2 = 43$ are accepted as the segments of quasi-periodic character. All the remaining segments undergo the frequency criterion. Initially they are included among the class of high frequency signals or among the class of voiced signals or "conditional noise". Frequency criterion is based on the estimation of zero-crossing parameters. The segments S_l for which the number of zero-crossings is larger than $\delta_3 = 96$ (and that with $f_p = 12$ kHz equals the frequency exceeding 2,5 kHz) are considered high fre-

quency noise segments, whereas the segments for which the number of zero-crossings exceeds $\delta_4 = 50$ ($C = 1,2$ kHz) and their maximum amplitude is $\max\{x_i\} < \delta_2$; $x_i \in S_l^l$ are accepted as "conditional noise".

"Conditional noise" may be accepted as a voiced segment or as high frequency noise. If there is a group of "conditional noise" segments in the close neighbourhood of the segment accepted as high frequency noise, the whole group is accepted as high frequency noise. Otherwise the "conditional noise" segments are considered and included among the class of voiced segments. If we denote the segment S_l accepted as silence - 0, voiced segment - 1, "conditional noise" - 2, and high frequency noise - 3, the signal represented by the following sequence of the appropriate classes of the segments S_l

0001110001112221112233112233111

is classified as

00011100011111111113333113333111

During the next stage of the algorithm we make use of the results of the previous stage. Only those segments which have been accepted as voiced are exposed to processing. It is necessary to search for indexes t for which the momentary values correspond with the significant amplitude of a given microphonemes. As well as PR) algorithm the following terms are introduced - local maximum and local minimum, absolute maximum, significant maximum and minimum and also significant extreme. Some modifications concerning the semantic content of the above terms have been introduced and also the algorithm of their indication has changed.

If we denote the vector representing the input signal as X then:

$x_i \in X$ is a local maximum if

$$(x_i > x_{i-1}) \wedge (x_i > x_{i+1}) \wedge (x_i > 0)$$

and $x_i \in X$ is a local minimum if

$$(x_i < x_{i-1}) \wedge (x_i < x_{i+1}) \wedge (x_i < 0)$$

We shall denote them \max_x , \min_x respectively*. On the basis of local maxima the absolute maximum, the significant maxima and significant minima are indicated, we denote them $\max^A x$, $\min^A x$, A . The absolute maximum A is calculated for the succeeding segments of the length $D_1 = 8D$ (C it makes about 170.6 ms). Here, it is necessary to point out that a division of a signal into segments of the length D_1 as dictated by the limitations posed by the ZAH-41 is different from the principles given in the algorithm PPD and aims at using as few memory cells as possible and achieving compatibility with the system AKORD (1,3). Accepting the segments of the length D_1 greatly influences the process of dynamic analysis of the microphoneme periods. It is also extremely significant in the process of correction, during the final stage of algorithm. The indication of the period of the larynx tone for such segments allows independence from fairly essential changes of the length of the larynx tone owing to some changes in articulation, contrary to the indication of the expected period of the larynx tone estimated continually through the whole signal X .

It is particularly important for long signals considering the initially accepted assumption No 1. The indication of segments of duration D_1 influences the first stage of the algorithm in which 8 segments S are indicated successively for each step. Within every segment $S_{D_1}^k$, $k=1,2,\dots$ the absolute maximum A_k is assigned.

$$A_k = \max \{ \max_{i \in X} x \}_k$$

where $\{ \max_{i \in X} x \}_k$ means the set of local maxima indicated within the segment $S_{D_1}^k$. Having A_k , we indicate the significant maxima and minima within the sets of local maxima and local minima. The significant maximum must satisfy the following conditions:

$$W1. \left(\max_{2n} x > 0.9A_k \vee \left\{ \begin{array}{l} \max_{2n} x > \frac{\max_{2n-2} x + \max_{2n-1} x}{2} ; n > 2 \\ \max_{2n} x > 0.6 \cdot A_k ; n < 2 \end{array} \right. \right) \wedge$$

$$\wedge G_2 \leq \rho(\max_{2n} x, \max_{2n-1} x) \leq G_3 \wedge \rho(\max_{2n} x, \max_{2n-1} x) = \min \rho$$

where $\rho(\max_{2n} x, \max_{2n-1} x)$ indicates the distance between the previous and the current significant maximum. That distance should be minimum within the defined interval. If the condition W1 is not satisfied the first "essential" local maximum is chosen according to the condition:

$$W2. \max_{2n} x \wedge \max_{2n-1} x \wedge \rho(\max_{2n} x, \max_{2n-1} x) \leq G_3 \wedge$$

$$\wedge G_2 \leq \rho(\max_{2n} x, \max_{2n-1} x) \leq G_3 \wedge \rho(\max_{2n} x, \max_{2n-1} x) = \min \rho$$

If the significant maximum has been indicated considering the condition W1 and if that condition is satisfied:

$$\rho(\max_{2n} x, \max_{2n-1} x) > 2 \cdot G_2$$

then we make an attempt at finding the additional significant maximum $\max_{2n}^d x$ according to the condition W2 within the interval satisfying the following inequality:

$$G_2 \leq \rho(\max_{2n}^d x, \max_{2n-1} x) \leq \rho(\max_{2n} x, \max_{2n-1} x) - G_2$$

This algorithm is different from the one presented by Roddy, and thus the significant extremes will have different meanings than in the algorithm PPD. Yet owing to this difference, it is possible to eliminate errors of the type of amplitude fluctuation caused by low frequency ($T \approx 2T_k$). The algorithm of searching for the significant minimum is analogous to the search for the significant maximum. Speech wave analysis has shown that because the envelope of the larynx tone is of a rather differentiated character, not all significant maxima univocally assign the places of extremes of the microphoneriac envelope. It can be observed that whenever there exists a significant maximum that we accept as the extreme of the microphoneriac envelope, there also exists a corresponding significant minimum within a small neighbourhood. This can be explained by the occurrence of the greatest turbulence in the speech wave at the moment when the vocal tract is excited by the release of the pressure resulting from the vibration of the vocal cords, the moment in time when the microphoneriac begins. Therefore, as the "significant extremes" only those

significant maxima are chosen which have the corresponding significant minima in their neighbourhood (4.5 ms - 42 samples as 1/4 of the longest expected microphoneriac) . Thus obtained, microphoneriac markers undergo the logical processing. During the stage of the logical processing we make use of the already presented assumption that microphoneriacs are searched for in the quasi-periodic signal and that the following inequality is satisfied:

$$0.8 \leq \frac{T(j)}{T(j-1)} \leq 1.2$$

This assumption makes it possible to eliminate the errors of the second stage, e.g. finding an extreme in a place where it should not be or not finding it in the place where it should be. In order to eliminate such errors an expected value on the average period of microphoneriacs T_{sp} is calculated. It is calculated as a discrete statistical distribution of distance among the selected extremes within the segment $S_{D_1}^k$, consistent with the earlier accepted settlements. We shall denote it T_{sp}^k . Thus calculated, T_{sp}^k is used for the stage of correction on the same analogous basis as it was done by Reddy (10). The acceptance of T_{sp}^k calculated along the fixed length D_1 of the segment $S_{D_1}^k$ makes it possible to eliminate the errors, resulting from intonation changes of longer statements. Thus the correction's efficiency has been considerably increased while retaining all of Reddy's algorithm and the coefficients accepted by him.

The process described above concerning the logical processing is iterated for the same segments according to need. The retrieved markers indicate the significant places within the microphoneriac and those regions make a basis for finding the ends of microphoneriacs. The author will consider the zero-crossing point or the point where the signal assumes a value which is the closest to zero, the end or beginning of the microphoneriac. The end or beginning of microphoneriac is indicated in the region where the envelope of the signal assumes a minimum. Such a choice for the place of the microphoneriac end reflects the physical aspect connected with the excitement of the vocal cords and also it corresponds with the results of the spectral analysis. The spectral analysis concerning microphoneriacs showed great parameter stability. The agreement of phase as well as the great relevance of microphoneriac and phoneme characteristics, belonging to the same speaker and to the same class of speech sounds, proves the correct segmentation. In accordance with the phenomenon of vocal cords vibration we shall search for the place of the microphoneriac end (or the beginning) before the maximum of turbulence assigned by the significant extremes, i.e. search for the ends of the microphoneriac before the significant extremes or before the proceeding local maxima. We accept a principle that the end of the microphoneriac is calculated directly before the local maximum which is the furthest on the left and its value is not less than 1/2 of the

value of the succeeding significant extreme. Let there be a condition where the distance between the indicated end of the microphoneme and the succeeding extreme should not be larger than $0.5 T_{sp}$. Let us denote this distance as E^j ; where j is a succeeding number of the microphoneme. It is used for indicating the average distance between the expected end of the microphoneme and the succeeding extreme.

$$E_{sr} = \frac{\sum_{j=1}^m E^j}{m}; m = \max j$$

This value is used for the correction of erroneously calculated ends of microphonemes. Such errors may appear in the situation shown in Figure 3.

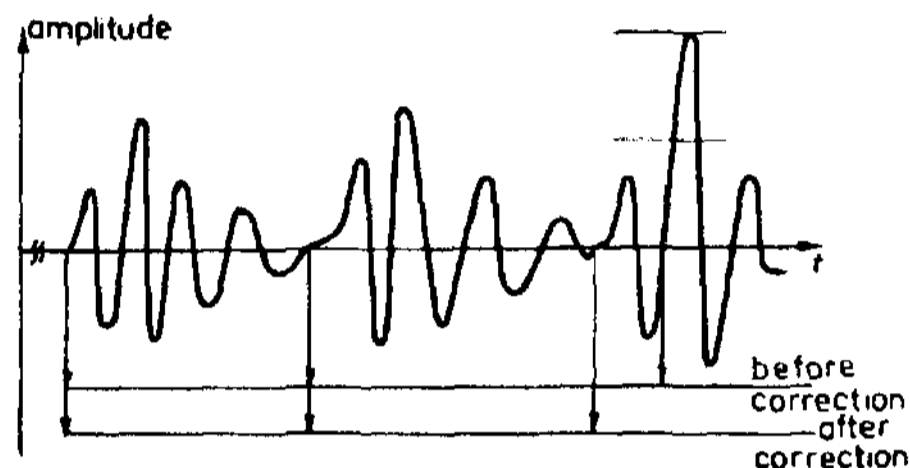


Figure 3

All the ends of microphonemes which satisfy the condition:

$$W3. |E_{sr} - E^j| < \frac{1}{16} T_{sp}$$

are considered correct.

In the other cases, a process of correction is undertaken which may lead to indicating a new and better end of the microphoneme. The algorithm of correction aims at finding a new end of the microphoneme, satisfying condition W3. This phase ends the stage of correction for a given microphoneme or the stage of indication of the microphoneme end, which is the least distant from the expected end, marked by E_{sr} . The ends of microphonemes are calculated during the stage of correction according to the analogous principle as before but the values of succeeding local maxima need not exceed half of the value of the respective extreme. Because the process of correction may be repeated, there is a condition indicated, concerning the character of this process. The process of correction should be "unilaterally convergent", by which we mean that a new end of the microphoneme will be considered correct, and the stage of correction in which it was calculated will be considered essential, if it is placed closer to the expected end than the previous incorrect one and also, if it is on the same side, in relation to the expected end, as the end considered incorrect. In the case of "oscillation" of the correction process towards the expected end, we accept the end of the microphoneme indicated during the previous stage of correction as correct. As a consequence of a possible correction we obtain the markers of micropho-

nemes considered the results of the algorithm of primary segmentation.

Summing up, the algorithm of primary segmentation is divided into four stages:

- Stage I - The indication of a quasi-periodic signal and its extraction from an input signal.
- Stage II - The calculation of the significant extremes of a voiced signal.
- Stage III - The correction of the extremes calculated in stage II
- Stage IV - The calculation of the ends of microphonemes
 - part I - An indication of the expected ends
 - part II - The correction; an indication of the resultant ends of microphonemes.

Conclusions

The procedure described in this paper was used for the analysis of 10 single words* with satisfactory results. In the majority of cases this algorithm found more microphonemes than a man could have done on his own. In particular this cases concerned word endings. It was impossible to carry on more exacting research on longer signals because of the ZAM-41 computer's slowness. Due to this difficulty, the above described algorithm is processed in Fortran on the H-6030 computer. This new programme will thus make it possible to test a larger number of longer utterances. Therefore, by working with more and longer utterances, it will be possible to indicate if this algorithm is better, in any cases, than those previously used for frequency analysis of speech sounds.

References

1. A. Dziurnikowski, P. Zdrojek: Organizacja współpracy urządzenia analogowo-cyfrowego AKORD z maszyną cyfrową ZAM-41
Papers of Symposium: Metody bezpośredniego wprowadzania i wyprowadzania informacji tekstowej i obrazowej w systemach informatycznych, Jabłonna 1973.
2. L. Macherzyński: Określenie widma przebiegu wprowadzanego do EMC przez urządzenie analogowo-cyfrowe AKORD i wyprowadzanie przebiegów o zadanym widmie
Papers of Symposium: Metody bezpośredniego wprowadzania i wyprowadzania informacji tekstowej i obrazowej w systemach informatycznych, Jabłonna 1973.
3. A. Dziurnikowski: System AKORD, SUP-1 CAM-TRA
Internal report of W.U.C.C. 1972.
4. P. Kolars, Murray e Den: Recognizing Patterns- Studies in living and Automatic Systems
MIT Press, Cambridge, Mass., 1968

* This words are: ala, rojal, rzeka, wąż, jodła, nora, jęba, masz, foka, ryba.

5. S. Zwicker, R.Feldkeller : Das Ohr als Nachrichtenempfänger
Stuttgart ,1967
6. Pierre Vicens : Aspects of Speech Recognition by Computer
(a dissertation) Com. Scien. Dep.
Stanford University, 1969
7. E.A. Sapozkow : 3ygnal mowy w telekomunikacji i cybernetyce
.arszawa, VNT, 1966
8. R. Tadeusiewicz : Sioniczna koncepcja mowy
Papers of Symposium: Zietody Bezposredniego wprowadzania informacji tekstowej i obrazowej w systemach informatycznych, Jablonna, 1973.
9. J.L. Flanagan : Speech Analysis, Synthesis and Perception
Berlin, Springer-Verio^, 1965.
10. D.R. Reddy : Pitch period determination of speech sound
Communication of the ACM, 1967, vol.10
11. G.Modena, C. Scagliola, E. Vivalda : Influence of phase handling on short time speech spectra
The 8 International Congres on Acoustic, London, July 1974, Vol. 1. p. 238. Contributed Papers
12. A.Dziurnikowski : Segmentacja pierwotna sygnaiu mowy i ekstrakcja mikrofonemow
Papers of Symposium : Komputerowe systemy przetwarzania danych doswiadczalnych, Kazimierz, 1974
13. A.Dziurnikowski : System rozumienia mowy SUSY
Papers of Symposium : Metody rozpoznawania, klasyfikacji i wyszukiwania informacji, Jadwisin 1975.