

## SIMPLIFIED SYSTEM FOR ISOLATED WORD RECOGNITION

R. GUBRYNOWICZ, K. MARASEK, W. MIKIEL, W. WIĘZŁAK

Institute of Fundamental Technological Research, Polish Academy of Sciences  
(00-049 Warszawa, ul. Świętokrzyska 21)

This paper presents a general-purpose system for recognition of a limited set of words uttered in isolation. Such a system is intended for voice control of robot's movements. In order to minimize the number of operations performed during the recognition process and to limit the memory requirements frequency analysis of the signal was performed in adequately selected bands. Output signals from filters undergo detection and through an A/D converter are introduced into a computer where they undergo further processing logarithmic conversion and linear time standarization, among others. This leads to a reduction of the number range in further calculations. The DTW algorithm was used in the recognition process, while templates of individual words are introduced once, in principle separately for individual operators. The developed system speaker-dependent, in principle was verified experimentally for various vocabularies (containing 20 to 60 words) uttered by 11 voices (including 1 female voice). The average recognition accuracy for a 60 word vocabulary exceeded 98% for individual voices, while in a case of recognition without system accomodation to given voice the average error of recognition increased by about 10%.

W pracy przedstawiono uniwersalny system rozpoznawania ograniczonego zbioru wyrazów wymawianych w izolacji. System ten jest docelowo przeznaczony do sterowania ruchem robota za pomocą głosu. W celu zminimalizowania liczby operacji wykonywanych podczas procesu rozpoznawania oraz zmniejszania zajętości pamięci zastosowano analizę częstotliwościową sygnału w odpowiednio dobranych pasmach. Sygnały wyjściowe filtrów są poddawane detekcji i wprowadzone poprzez przetwornik A/C do komputera, gdzie następuje ich dalsza obróbka, m.in. konwersja logarytmiczna danych oraz liniowa normalizacja czasowa. Uzyskano dzięki temu znaczną redukcję zakresu liczb, którymi operuje się przy dalszych obliczeniach. W procesie rozpoznawania zastosowano algorytm DTW, przy czym wzorce poszczególnych wyrazów są w zasadzie wprowadzane jednokrotnie, oddzielnie dla poszczególnych operatorów. Opracowany system, który w zasadzie jest typu speaker-dependent został zweryfikowany doświadczalnie dla różnych słowników zawierających od 20 do 60 wyrazów wypowiedzianych przez 11 głosów (w tym 1 żeński). Średnia dokładność rozpoznawania dla słownika 60-wyrazowego dla poszczególnych głosów wyniosła ponad 98%, zaś w przypadku rozpoznawania bez dostosowywania systemu na zadany głos, średni błąd rozpoznawania wzrósł o ok. 10%.

## 1. Introduction

The problem of man-machine communication with the use of speech is at present an issue of great interest in many scientific research centres all over the world. Until the mid 70-ies researches were performed in university centres and, were out of interest of the industry in general. However, the practical application of the results of these researches has been taken into consideration more and more in the second half of this decade. It could be expected that the progress in the realization of automatic speech recognition systems will make their effective operation, also in industry, possible (PELTIN [10]).

The advantages of the use of the so-called acoustic input for data input and control by speech are indisputable. Speech is the most natural and quickest, at the same time, form of communication between people. Table 1 presents the rate of information transmission from a man using different means of communication:

**Table 1.** Average information transmission rate from man

Speech (10 sounds/s)	— 50 bit/s
Keyboard (60 words/min)	— 30 bit/s
Morse code (12 words/min)	— 6 bit/s
Buttons, keys	— < 6 bit/s

It is worth noticing that the values given for non speaking communication are rather maximal. For example, a lower transmission rate, of about 15–20 bit/s can be expected for an average user of a keyboard.

The freedom of movement of the operator who can perform additional manual operations on his work-stand, as well as the possibility of keeping eye contact with the examined elements, surfaces etc. even when registering control information, are further advantages of an acoustic input. The application of an acoustic input for information feeding greatly increases productivity, especially on stands of visual control of product quality (PELTIN [10]).

Two fundamental kinds of approaches can be distinguished in the problem of automatic speech recognition. The first one is realized on the basis of formalized knowledge concerning phonetic-acoustic, phonological, lexical, syntactic structures and semantics of given language. Speech recognition in such a system is based on sets of rules used for speech signal transformation into sequences of symbols representing definite notions of linguistic units with given structure. This approach is used mainly in systems for continuous speech recognition — very complex and generally constructed for the investigation of foundations of unlimited recognition of continuous speech.

The other approach consists in the recognition of a speech signal with methods of pattern recognition, with the application of numerical classification procedures of measurement results or vectors of features, created on their basis. In this case, in

general, classification process is realized on the basis of the division of the space of features into subspace ranges — individually for every class. Boundaries of these ranges and classification rules are defined on the grounds of geometric, topologic or probabilistic criteria.

Recognition models of the first type are very complicated and require the entire linguistic and phonetic knowledge of a given language to be presented in formalized form. In this case the optimization of recognition algorithms is an extremely difficult task, while for models based on the analysis of similarity between the recognised and reference pattern it can be relatively easily formulated formally, since the fact that individual models can be derived automatically from data observation during training is applied. This makes it possible to construct relatively easily a practical recognition system, but with limited effectiveness. It is a result of an assumption that ideal patterns, generated by a fairly simple model, exist. Only in such a case the problem of parameter estimation can be given in analytic and optimized form. Such limitations and simplifications of the description form are possible to accept in definite situations, e.g. when a relatively small number of words spoken in isolation is applied, as it was assumed in this paper.

## **2. Background of the recognition system of isolated words**

Considering the application of the designed system it was accepted that the vocabulary for operator-robot communication, except for 10 digits, has to contain words for the control of its motion, that is such words as “forwards”, “backwards”, “right”, “left”, etc. At the same time the replacement, full or partial, of the communication vocabulary was assumed and the limitation resulting from fine tuning of the system to the voice of a given operator was permitted.

It is also essential for the system to work in a laboratory without additional acoustic adaptations of the room and attenuation of outer noise sources. Because of the condition of vocabulary replacement for operator-robot communication, it was necessary to accept a system with an isolated training stage, which would be repeated every time the operator used a new collection of words. The condition of resistance to noise disturbances imposed the necessity of applying a mouth microphone and developing a specific procedure for endpoints detection. The recognition time of one word, which should not exceed 1 s for a computer with a clock of 4.7 MHz, significantly influenced the applied solution. It radically influenced the elaborated method of speech signal processing and the form of the recognition algorithm.

## **3. Speech signal processing**

While elaborating the isolated word recognition system it was accepted that the speech signal corresponding to individual words is sufficiently defined in the frequency domain in the form of time-dependent power spectra. It is furthermore

assumed that changes of the spectrum within a spoken word are sufficient for its identification and more detailed acousto-phonetic analysis, which makes possible, the detection and recognition of smaller linguistic units such as speech sounds, syllables, etc. is not necessary. Such an approach has a certain advantage — it is not necessary to know the phonetical and phonological structure of a given language, and errors of speech signal segmentation are avoided.

The requirement for operation rate and low cost of the recognition system, made it practically impossible to apply more refined spectrum analysis methods, such as FFT or LPC. Thus, a hybrid solution, with a considerable part of processed signal realized by an analogue system was decided upon. To this end a set of 7 band filters with detection circuits on the output was applied. Filters has bands with different widths (from 1/3 to 2 octaves), while their mid-band frequencies were so chosen, that it was possible to differentiate 6 Polish vowels and to divide (unvoiced fricatives) into two classes. The applied preliminary universal for most voices in principle filtration and detection led to a considerable data reduction, with a drastic decrease of their transmission rate to the computer (about 8.4 kb/s), at the same time.

Yet the application of spectrum analysis makes necessary to apply effective amplitude normalization, either on the analogue side, or after digital conversion. In the first case an amplitude compression system is required. But it increases the system's sensitivity to external noises. Whereas for normalization on the digital side, attenuation of components in higher frequency bands occurs frequently. Moreover, this type of normalization causes an increase of processing time. For these reasons an analogue solution was chosen.

Systems with band filtration are still popular (e.g. CROCHIERE, FLANAGAN [4]), especially in systems with practical application, for which the requirement of simplicity and operational reliability is as important as the low price of the recognition system. However it should be noted that, in general, signal filtration is not oriented as in the described system and, in general basis, frequencies are distributed as in vocoder systems, in linear or mel scale. The selection of filter bands adequately to the structure of chosen Polish language speech sounds made considerable reduction of the number of filters possible, and what follows — further reduction of the data size with an increase of reliability at the same time.

## 4. Description of system's operation

### 4.1. Speech signal parametrization

Figure 1 presents a general block diagram of the system. After compression, filtration and detection the speech signal is fed into a computer through a multi-channel, standard, 12-bit  $A/D$  converter. Used ranges of numbers are reduced by logarithmic matrix conversion of output signals. Figure 2 a presents examples of the

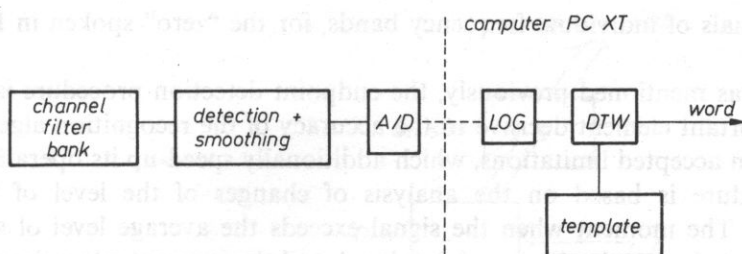


FIG. 1. Block diagram of isolated word recognition system

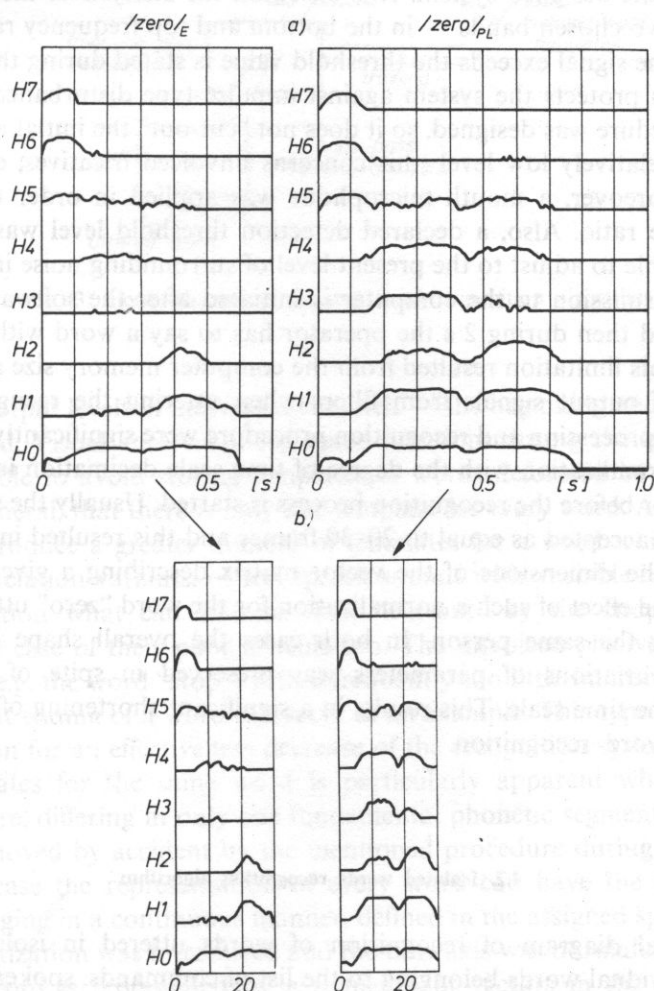


FIG. 2. Examples of output signals from filters before (2a) and after linear time normalization (2b), obtained for the word "zero" said in English and in Polish



output signals of individual frequency bands, for the "zero" spoken in English and Polish.

As it was mentioned previously, the endpoint detection procedure is one of the more important element decisive to the accuracy of the recognition algorithm. This results from accepted limitations, which additionally speed-up its operation. Usually this procedure is based on the analysis of changes of the level of the signal's amplitude. The moment when the signal exceeds the average level of surrounding noise is treated as the beginning of the signal, and the moment when the signal's level becomes lower than the noise level and after which the silence segment of not less than 0.3–0.4 s occurs is accepted as the end of the signal.

A somewhat more complex criterion for the beginning of the signal was elaborated in the designed system. It is based on the analysis of the sum of output signals from two chosen bands — in the bottom and top frequency range. While the fact whether the signal exceeds the threshold value is stated during the first 50 ms of the signal, this protects the system against impulse type disturbances. At the same time this procedure was designed, so it does not "cut-out" the initial and final speech sounds with relatively low level (this concerns unvoiced fricatives, especially) from the signal. Moreover, a mouth microphone was applied in order to improve the signal to noise ratio. Also, a declared detection threshold level was applied, what makes it possible to adjust to the present level of surrounding noise in a given room.

Data transmission to the computer is initiated after the software or hardware interruption and then during 2 s the operator has to say a word with duration time under 1.2 s. This limitation resulted from the computer memory size assigned for the registration of output signals from filters when uttering the recognized word.

The signal processing and recognition procedure were significantly speeded up by linear time normalization with the degree of time scale decimation initially declared by the operator before the recognition process is started. Usually the standard length of a word was accepted as equal to 20–30 frames and this resulted in a reduction by 2–4 times of the dimensions of the vector matrix describing a given word. Figure 2 b presents the effect of such a normalization for the word "zero" uttered in English and Polish by the same person. In both cases the overall shape of the function representing variations of parameters was preserved in spite of an over twice reduction of the time scale. This results in a significant shortening of the calculation time during word recognition.

#### 4.2. Isolated words recognition algorithm

The general diagram of recognition of words uttered in isolation is shown in Fig. 3. Individual words belonging to the list of commands, spoken to control the robot's motions, are uttered — once, in general — in turn by the operator during the training stage. After the word is uttered, the monitor automatically displays

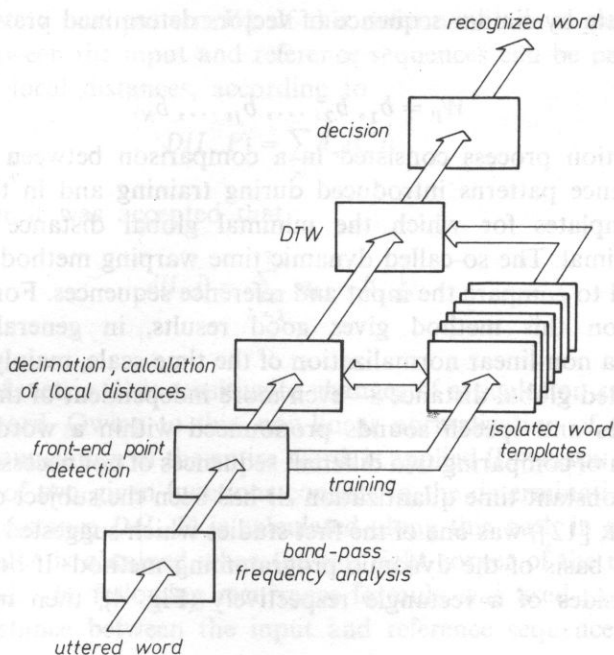


FIG. 3. Block diagram of isolated word recognition algorithm

variations of output levels in individual channels — the beginning and end of spoken word is marked. Visual analysis of registered functions during the training process makes it possible to avoid storing templates of words deformed by chance. It was accepted, in general, that there is only one template for every word. All the same it is possible to introduce a greater number of templates for a word which has various correct pronunciation variants, or has speech sounds with variable level in the final or initial position what can lead to their “cut-out” by the endpoint detection procedure in a case of their poor articulation. The unvoiced plosive /p/ spoken as a final sound (e.g. the word “stop”) with a frequently too little intensity of plosion to form an explicit ending of a word can serve as an example. This type of error can be the main reason for an effectiveness decrease of the recognition system. The need of various templates for the same word is particularly apparent when words with similar structure, differing in only one fundamental phonetic segment which may be sometimes removed by accident by the mentioned procedure during learning stage. In a general case the representation of every word can have the form of vector functions changing in a continuous manner, defined in the assigned space of features. Because discretization was introduced and the time axis was normalized to  $N$  points, hence every word is represented by a sequence of vectors in the following form

$$W_T = a_1, a_2, \dots, a_j, \dots, a_N.$$

while it's template by similar sequence of vectors determined previously during the learning process

$$W_p = b_1, b_2, \dots, b_j, \dots, b_N.$$

The recognition process consisted in a comparison between input word and succeeding reference patterns introduced during training and in the choice of this word from templates for which the minimal global distance from the input sequence is minimal. The so-called dynamic time warping method DTW (Levinson [7]) was applied to compare the input and reference sequences. For cases of isolated word recognition this method gives good results, in general, owing to the introduction of a non-linear normalization of the time-scale, mainly. This makes the value of calculated global distance — even more independent of the variable rate of uttering syllables and speech sounds pronounced within a word.

The problem of comparing two different sequences of not necessarily equal length at maintained constant time quantization  $\Delta t$  has been the subject of much research. Paper (VINTSYUK [12]) was one of the first studies which suggested a solution to this problem on the basis of the dynamic programming method. If both sequences are drawn on the sides of a rectangle respectively (Fig. 4), then minimal distances

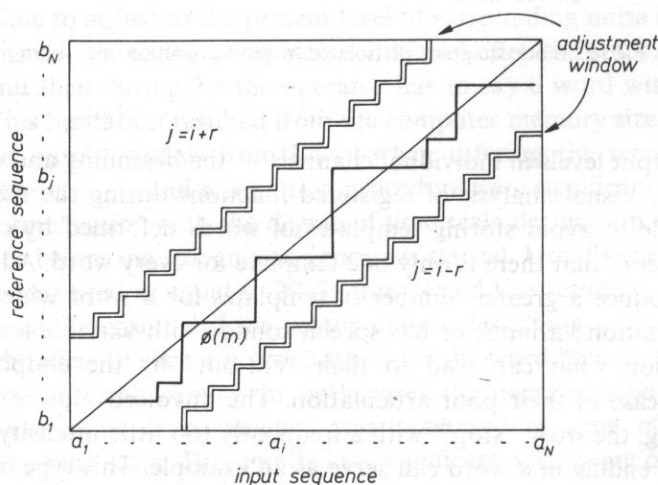


FIG. 4. Determination of cumulative distance between input word and reference sequence with the DTW method

between succeeding segments lie near the rectangle's diagonal when the input sequence and reference sequence are very similar. In a case of perfect compatibility, the line joining minimal local distances  $d(i, k)$  coincides exactly with the diagonal. A different rate of uttering the input word and the reference word leads to a deformation of this line to form  $\phi(m)$  which is the function of deformation of



pattern  $W_I$  adequately into pattern  $W_p$ . If this deformation is neglected, then the global distance between the input and reference sequences can be presented in the form of a sum of local distances, according to

$$D(I, P) = \sum_m d_m(i, j)$$

while in this paper it was accepted that

$$d(i, j) = \sum_{k=1}^7 |a_{i,k} - b_{j,k}|,$$

for  $1 \leq i, j \leq N$ .

Still this type of metric (1) is sensitive to changes of articulation rate of syllables uttered within a word. Owing to this, non-linear normalization of the time scale, including the pronunciation of the entire word, is applied (LEVINSON [7]). Then the similarity analysis of two given functions consists in the determination of the path  $\phi(m)$ . Cumulative distance  $D(I, P)$  is calculated along this path in such a manner that its minimal value is obtained when the top right corner of the matrix of local distances is reached. The following recurrence formula was used to determine the minimal global distance between the input and reference sequence

$$D(i, j) = d(i, j) + \min \{D(i-1, j), D(i-1, j-1), D(i, j-1)\}. \quad (3)$$

When calculating the minimal distance a "corridor" is put along the diagonal  $i = j$  in order to prevent too large deviations of function  $\phi(m)$  (Fig. 4), according to formula

$$|i - j| \leq r. \quad (4)$$

Because the minimal global distance for a given input word can be calculated for various templates along different optimal paths, an additional normalization is introduced with respect to the number of steps  $k$

$$D_n(W_I, W_p) = \frac{1}{k} \times D(s, u), \quad (5)$$

where  $N \leftarrow k \leftarrow 2 \times N$  and  $N - r \leftarrow \delta, u \leftarrow N$ .

The acceleration of the template search is an important problem, as well as the problem of avoiding the inspection of the entire collection of reference patterns every time a template closest to the spoken word is chosen. At the present moment the searching time can be considerably shortened by including two additional factors when comparing the input word with succeeding templates. These factors are: word length and maximal acceptable cumulative distance value for taking a decision whether the input word belongs to the same class as the given template. The elimination of distance calculation in cases when lengths of the input word and template were positively different, was possible owing to durations analysis. The critical difference of 50% was accepted. When it was exceeded the distance was not calculated and the given template was rejected. A more accurate comparison of

lengths of words expressed in the number of syllables, is planned in the future. At that time the entire vocabulary will be preliminary divided into subsets assembling words with the same number of syllables.

The introduction of the maximal critical cumulative distance for taking a positive decision not only shortens the time necessary to review the template collection, but also reduces the number of incorrect recognitions, especially in a case of a word not belonging to the vocabulary.

Methods of speeding up the lexical retrieval and choice of the template subcollection close to the input word are the object of special interest of designers of recognition systems, because the recognition accuracy and size of the word vocabulary depend on the efficiency of developed searching algorithms. The vocabulary size is decisive to the recognition time. An interesting method of searching the so-called "Geometrical Search" was presented recently (FARAGO, GORDOS, LUGOSI [5]). The elaborated algorithm, taking advantage of preliminary calculated distances between individual templates, makes it possible to shorten the search time by one order of magnitude for a vocabulary of 200–300 words at almost unchanged recognition accuracy.

## 5. Experimental material

The developed system was practically evaluated on the basis of various vocabularies uttered by 10 persons. Experiments were performed at various time intervals, reaching 3–4 months. The largest vocabulary, containing 60 words, was uttered 15 times by 8 persons. A total of 5000 statements constituted the testing material. The 60-word vocabulary was chosen to make it possible to control by speech data transmission and processing in a computer. At the same time research included another vocabulary consisting of 24 words and containing instructions for a computer system for robot's motion control.

## 6. Results of recognition

It is worth mentioning at the very beginning that in systems applying the DTW method, recognition accuracy depends on several factors which change the spectrum pattern of an analysed word in an irregular manner. In the first place we should mention: instability of individual characteristics of the operator's voice, considerable information reduction of the signal parametric representation frequently making the discrimination of certain classes of speech sounds impossible and variable phonetic-acoustic structure of individual words included in the accepted vocabulary for communication.

Fluctuations of the level of the input signal are one of the significant elements of

the signal's instability. In this case the application of a logarithmic scale of voltage conversion is insufficient. Therefore, when band-pass frequency analysis is applied an amplitude compressor has to be used. In the case under consideration the compressor converts fluctuations of the input level from about 20 dB to the range of 3 dB.

The stability of individual characteristics and phonetic-acoustic structure greatly depends on the operator's training and knowledge of microphone technique. An operator with big enough practice achieved recognition accuracy of 98–99%, while less trained persons reached about 92% on the average. A recognition accuracy decrease is due to the fact that results of signal parametrization and its acoustic pattern can also depend on the instability of individual characteristics.

The reduction of the parametric description of a speech signal considerably speeds-up the distance calculation. Yet, acoustic patterns of phonetic-acoustic classes with similar physical structure differ only slightly. This leads to difficulties with recognition of words with similar pronunciation. Table 2 contains analysis results of

**Table 2.** Influence of decimation degree (description reduction) on average distances from input word to correct template and closest incorrect template, and on recognition accuracy

Average distances $\backslash$ $N$	10	15	20	25	30
From correct template (a)	97.8	93.3	86.67	80.72	78.38
From nearest incorrect (b)	142.77	145.38	144.26	140.92	138.5
Distance between templates (a) i (b)	44.97	52.08	57.59	60.2	60.12
Recognition score	0.93	0.928	0.98	0.983	0.995

the influence of the number of points  $N$  for one voice, with duration of uttered word subjected to linear normalization. The table presents average distance changes between the input word and correct template and between the input word and nearest incorrect one. The data given above for a 60-word vocabulary proves that 20–25 is the sufficient number of points. This means that decimation of the signal occurs every 15–20 ms for one-syllable words, while every 40–50 ms on the average for 3 and 4 syllable words. It is characteristic that an increase of the number of points is accompanied not only by an increase of recognition accuracy, but also, what is

logical, by an increase of average distances between the correct and closest incorrect template. This indicator makes possible to closely observe the influence of the reduction of the parametric discription. The value of  $N = 20$  was accepted as optimal on the basis of recognition results. This number is two times smaller from the one given usually in literature (Myers et al. [9]), while the optimal width of the window does not differ from the one generally applied in isolated word recognition systems ( $r = 5$ ).

Table 3 presents experimental results of recognition in the speaker – indepen-

**Table 3.** Average distances from input word to correct template and closest incorrect template, and recognition accuracy obtained for different operator's voice in a speaker-independent mode  $N = 20$

Voice	Average distances	From correct template	From closest incorrect template	Between templates	Recognition score
		(a)	(b)	(a) i (b)	
M1		122.17	166.47	44.3	0.9
M2		124.25	161.67	37.42	0.88
M3		137.28	167.75	30.47	0.85
M4		131.35	170.73	39.38	0.92
M5		117.8	157.17	39.37	0.93
M6		109.93	154.78	44.85	0.95
M7		136.66	166.7	30.04	0.84
M8		124.0	155.82	31.82	0.88
M9		150.33	181.05	30.7	0.84
M10		136.05	176.48	40.43	0.9
F1		133.45	165.23	31.82	0.92

dent system, i.e. when reference words are entered by one operator, while recognized words were uttered by other persons. This research included 10 male voices and 1 female voice. A drop of recognition accuracy by about 10–12% on the average occurred in accordance with data given in literature. It can be avoided by creating a greater number of templates for individual words. On the basis of experimental data it was stated (Atal, Rabiner [3]) that in order to obtain a speaker – independent system, 12 templates for one word have to be formed.

Also the structure of the vocabulary itself, or to be precise the similarity of the phonetic-acoustic structure of words included in it, have a significant influence on the recognition accuracy. The relatively low recognition accuracy of about 75% on the average for the words input “wejście” (vejɛtɕe) and output “wyjście” (vijɛtɕe), extremely dependent on the operators voice, can be an example. Globally the number of recognition errors for 6 words (10% of words included in the vocabulary) with similar structure was equal to 50% of the total number of errors. Therefore, the

communication vocabulary in the presented system should not contain words with too similar or too variable acoustic structure.

To conclude it should be emphasized that reference vocabularies were entered only once by the operator and on this account recognition errors were caused by various, although permissible, pronunciations of the same word. Of course it is advisable to increase the number of reference sequences for a given word in such cases.

## 7. Conclusions

Achieved results confirmed the practical usefulness of the presented system, in spite of the application of a simplified spectrum analysis. The introduction of amplitude normalization before and logarithmic conversion after  $A/D$  conversion increased the system's resistance to random level fluctuations of a speech signal and made it possible to apply a 8-bit converter with relatively low sampling frequency (100 Hz). Logarithmic values were taken from tables, not calculated currently, to avoid further extension of calculation. It was possible to shorten the recognition time without a significant decrease of recognition accuracy by considerable reduction of the parametric description. Yet, in the present form the system is speaker-dependent, although a recognition accuracy reduction in a speaker-independent version can be avoided by an adequate choice of vocabulary.

This work was supported by the Problem CPBP 02.13.

## References

- [1] J. ACKENHAUSEN, S. S. ALI, D. BISHOP, L. F. ROSA, R. THORKILDSEN, *Single board general — purpose speech recognition system*, AT and T Technical Journal, **65**, 5, 48–59 (1986).
- [2] J. ALLEN, *A perspective on man — machine communication by speech*, Proc. IEEE, **73**, 11, 1541–1500 (1985).
- [3] B. S. ATAL, R. R. RABINER, *Speech research directions*, AT and T Technical Journal, **65**, 5, 75–85 (1986).
- [4] R. E. CROCHIERE, J. L. FLANAGAN, *Speech processing: An evolving technology*, AT and T Technical Journal, **65**, 5, 2–11 (1986).
- [5] A. FARAGO, S. GORDOS, G. LUGOSI, *Methods for decreasing the response time in isolated word speech recognition*, Proc. Speech Research, 89, Budapest, 255–258 (1989).
- [6] L. L. LAMEL, L. R. RABINER, A. E. ROSENBERG, J. G. WILPON, *An improved endpoint detector for isolated word recognition*, IEEE Trans. on ASSP-29, **4**, 777–785 (1981).
- [7] S. E. LEVINSON, *A unified theory of composite pattern analysis for automatic speech recognition in: Computer Speech Processing* [ed.] F. Fallaido, W. A. Woods, Prentice-Hall, Englewood Cliffs 1985, 243–275.



