



**SPEAKER IDENTIFICATION
 BASED ON MULTIPULSE EXCITATION AND LPC VOCAL-TRACT MODEL**

Seiichiro HANGAI and Kazuhiro MIYAUCHI

Department of Electrical Engineering
 SCIENCE UNIVERSITY OF TOKYO
 1-3 Kagurazaka Shinjuku-ku Tokyo 162 JAPAN

ABSTRACT

In automatic speaker identification, the reduction of dimensions of template is a key to realize a quick identification and to save storage. In this study, we extract some glottal wave parameters which does not seem to be susceptible to mimicry and combine them with some LPC vocal tract parameters to make a smaller sized template. In the extraction of feature parameters of glottal wave, the multipulse excitation model is adopted under modifying the pulse search algorithm.

According to experimental results obtained from 30 speakers' 5 vowels, 94% identification rate with a template of 10 feature parameters (8 vocal tract and 2 glottal parameters) and 99% with a template of 36 feature parameters (14 vocal tract and 22 glottal parameters) are obtained respectively. In comparison with the number of parameters to get same identification rate only with a formants' template, the number of feature parameters is reduced by 5 in case of getting 94% identification rate and 7 in case of getting 99% identification rate.

1. INTRODUCTION

As feature parameters of template which represents speaker's characteristics, long time power spectrum[1], linear prediction coefficients[2], pitch contours[3] and Cepstral analysis [4] are reported. To get high identification rate with these parameters, a large number of parameters is required. In automatic speaker identification within large speaker population, the reduction of dimensions of template is required to realize quick response and to save storage. This means that parameters in the template should be orthogonal each other. Furthermore, parameters should be stable over time and not be susceptible to mimicry.

In the investigation of speech synthesis, it is reported that the multipulse glottal model[5] makes synthesized voice natural and gives individuality. This means that the locations and sizes of pulses involves individuality. This glottal parameter also seems to be orthogonal to the vocal tract parameter and be susceptible to mimicry.

In this paper, we describe the pulse search algorithm which is modified to match our aim, i.e., searching pulses not in a frame but in a pitch, and the reduction of dimensions of template. And, we also show the relation between identification rate and the number of parameters based on multipulse excitation and LPC vocal-tract

model.

2. PROCEDURE FOR EXTRACTING FEATURE PARAMETERS

2.1 EXTRACTION OF VOCAL TRACT PARAMETERS

Fig.1 shows a flow of extracting individual feature parameters from speech signal. In order to get vocal tract parameters, speech signal is pre-emphasized and weighted by Hamming window. The emphasized and weighted signal is framed in 25.6ms and analyzed to get 12 LPC coefficients. When the signal is vowel, we can regard the transfer function $H(z)$ of vocal tract as an all pole model and write $H(z)$ as

$$H(z) = \frac{G}{1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_{12} z^{-12}} \quad (1)$$

where G is constant and a_i is the i -th LPC coefficient. From this equation, the formant

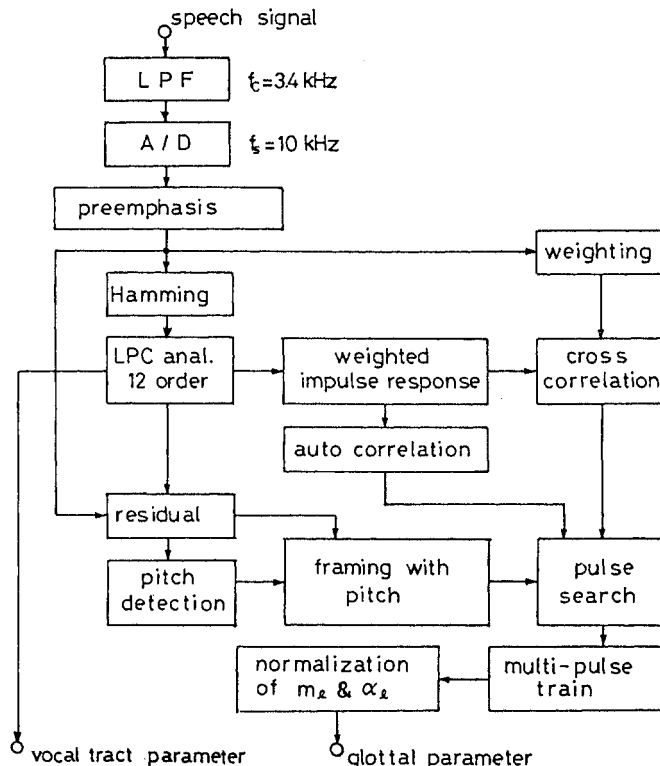


Fig.1 flow of extracting feature parameters

frequencies F_j and bandwidths B_j are given by complex roots as follows,

$$F_j = \frac{1}{2\pi T_s} \tan^{-1} [\text{Im}(z_j) / \text{Re}(z_j)] \quad (2.1)$$

$$B_j = \frac{1}{\pi T_s} \ln |z_j| \quad (2.2)$$

where z_j is one of complex roots and T_s is sampling period. We define 5 pairs of F_j and B_j as the vocal tract parameters.

2.2 EXTRACTION OF GLOTTAL WAVE PARAMETERS

To get glottal wave parameters, we adopt the multipulse excitation glottal model[5]. In this model, pulse train is determined to minimize squared error between original speech and synthesized speech by A-b-S method. This method, however, requires much calculation time. So, we use correlation technique to search pulses successively[6]. As shown in Fig.1, the emphasized speech and the LPC filter's impulse response which is derived from eq.(1) are perceptually weighted respectively. By obtaining cross correlation function of the weighted speech and the weighted response, the amplitude α_l of l -th pulse is successively calculated as follows ($0 \leq l \leq L-1$, L : number of pulses),

$$\alpha_l(m_l) = \frac{\sum_{n=1}^{256} s_w(n) h_w(n-m_l) - \sum_{k=1}^{l-1} \alpha_k \sum_{n=1}^{256} h_w(n-m_k) h_w(n-m_l)}{\sum_{n=1}^{256} h_w(n-m_l) h_w(n-m_l)} \quad (3)$$

where m_l is location of pulse α_l , $s_w(n)$ is the weighted speech and $h_w(n)$ is the weighted impulse response. In the original pulse search algorithm, pulses are determined in one frame. However, this is not suitable for our application, because we do not require natural sounding speech but need the information of locations and sizes of pulses per pitch. Therefore, we modify the algorithm to search pulses in one pitch period.

Fig.2 shows the procedure to get 10 pulses including main pulse (pitch pulse) per one pitch and the waveform, when the speech is vowel /a/. After getting 10 pulses, locations and sizes of pulses are normalized by those of main pulse and plotted on the polar-plane as shown in Fig.3 with cause of treating the pulse at rT_p ($0 \leq r \leq 1$, T_p : pitch period) and the pulse at $(1-r)T_p$ as the same pulse.

We define 9 pairs of abscissa x_l and ordinate y_l of plots of subpulses in Fig.3 as the glottal wave parameters ($l=1,2,\dots,9$).

3. EVALUATION OF FEATURE PARAMETERS

In order to effectively reduce the amount of feature parameter, we measure the F-ratio, i.e., inter-speaker variance to intra-speaker variance ratio for each parameter[7].

Table 1 shows the F-ratio of vocal tract and glottal wave parameters of 5 vowels (/a/, /i/, /u/, /e/, /o/), when 30 persons' 20 times utterances are used as sample. Parameters are sorted in order of F-ratio and 10 significant parameters in each vowel are listed. From Table 1(a), it is found

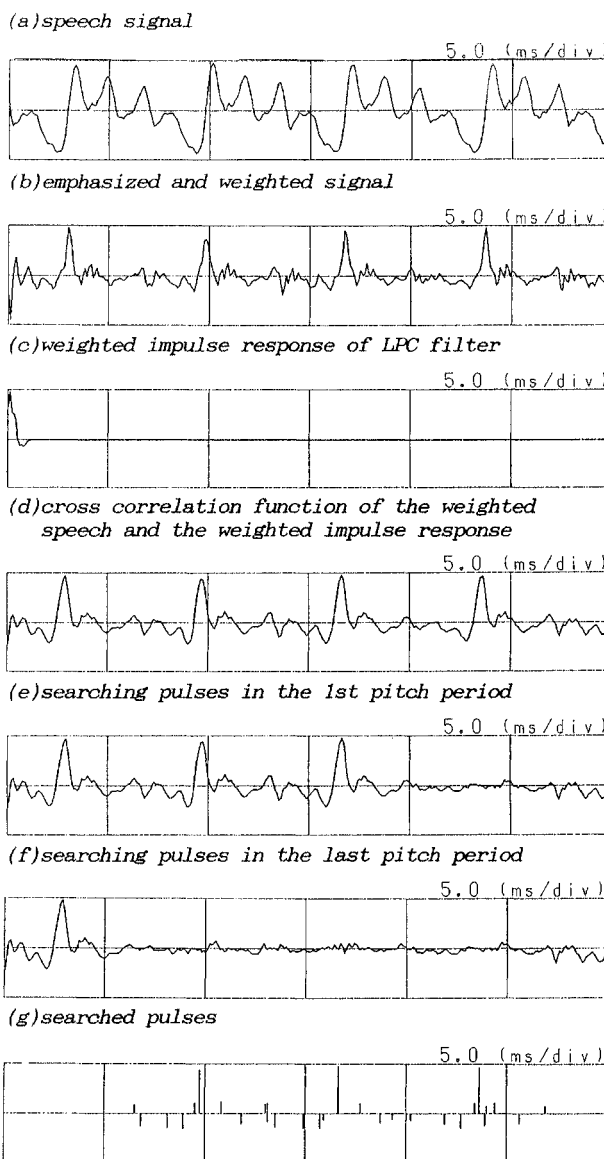


Fig.2 procedure to get 10 pulses per one pitch

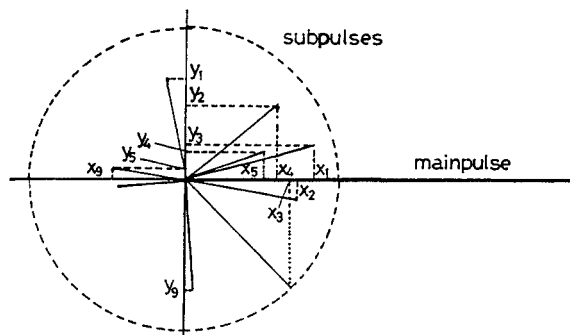


Fig.3 9 sub-pulses plotted on polar-plane

Table 1 F-ratio of vocal tract and glottal wave parameters of 5 vowels

(a) vowel /a/		(b) vowel /i/	
paramtrs.	F-ratio	paramtrs.	F-ratio
F ₁	7.2	F ₁	7.6
F ₂	3.8	B ₂	2.5
x ₈	3.3	F ₃	2.3
x ₇	3.2	y ₁	1.8
x ₉	2.6	x ₁	1.7
x ₆	2.4	F ₄	1.6
x ₅	1.9	x ₂	1.6
B ₃	1.9	x ₃	1.5
x ₁	1.6	x ₄	1.1
x ₄	1.6	y ₂	1.1

(c) vowel /u/		(d) vowel /e/	
paramtrs.	F-ratio	paramtrs.	F-ratio
F ₃	4.4	x ₇	2.8
B ₁	3.7	x ₈	2.7
y ₉	3.5	x ₆	2.6
F ₂	3.4	x ₉	2.3
x ₉	2.8	x ₅	2.1
F ₁	2.2	F ₁	2.0
x ₁	2.1	y ₁	1.6
y ₁	2.1	x ₄	1.5
x ₇	1.8	F ₄	1.4
x ₈	1.7	y ₅	1.3

(e) vowel /o/	
paramtrs.	F-ratio
F ₂	5.9
F ₁	4.4
x ₉	3.0
x ₈	2.3
y ₉	2.2
F ₃	1.8
x ₁	1.6
x ₇	1.6
B ₁	1.6
B ₃	1.5

that vocal tract parameters F₁, F₂ and B₃ of vowel /a/ have high F-ratio but that other 7 parameters are glottal wave parameters. From Table 1(b) to Table 1(e), 4 vocal tract and 6 glottal wave parameters in case of vowel /i/, 4 vocal tract and 6 glottal wave parameters in case of vowel /u/, 2 vocal tract and 8 glottal wave parameters in case of vowel /e/ and 5 vocal tract and 5 glottal wave parameters in case of vowel /o/ is significant parameters. The fact that many glottal wave parameters rank relatively high means the availability of the parameters in speaker identification.

4. SPEAKER IDENTIFICATION

By combining glottal wave parameters and vocal tract parameters, we attempt to reduce dimensions of template for speaker identification. The reference template is made from 5 times utterance of each vowel. In the identification, every Mahalanobis' distance between a template obtained from unknown speaker's 5 vowels and 30 reference templates is calculated and an unknown speaker is identified as the speaker who has a most similar template in the Mahalanobis space.

Table 2 Feature parameters in reference template

num. of paramtrs	feature paramtrs (combined)	feature paramtrs (vocal tract only)
1	/a/ F ₁	/a/ F ₁
2	/i/ F ₁	/i/ F ₁
3	/o/ F ₂	/o/ F ₂
4	/o/ F ₁	/o/ F ₁
5	/u/ F ₃	/u/ F ₃
6	/a/ F ₂	/a/ F ₂
7	/u/ B ₁	/u/ B ₁
8	/u/ y ₉	/u/ F ₂
9	/u/ F ₂	/i/ B ₂
10	/a/ x ₈	/i/ F ₃
11	/a/ x ₇	/u/ F ₁
12	/o/ x ₉	/e/ F ₁
13	/u/ x ₉	/a/ B ₃
14	/e/ x ₇	/o/ F ₃
15	/e/ x ₈	/i/ F ₄
16	/e/ x ₆	/o/ B ₁
17	/a/ x ₉	/u/ B ₄
18	/i/ B ₂	/u/ F ₅
19	/a/ x ₆	/o/ B ₃
20	/i/ F ₃	/e/ F ₄
21	/o/ x ₈	/u/ F ₄
22	/e/ x ₉	/e/ F ₅
23	/o/ y ₉	/a/ F ₄
24	/u/ F ₁	/e/ B ₁
25	/e/ x ₅	/o/ B ₄
26	/u/ x ₁	/e/ F ₃
27	/u/ y ₁	/i/ F ₂
28	/e/ F ₁	/a/ F ₃
29	/a/ x ₅	/e/ F ₂
30	/a/ B ₃	/u/ B ₂
31	/i/ y ₁	/e/ B ₃
32	/o/ F ₃	/i/ B ₃
33	/u/ x ₇	/u/ B ₃
34	/i/ x ₁	/e/ B ₄
35	/u/ x ₆	/e/ B ₂
36	/u/ x ₆	/a/ B ₄
37	/o/ x ₁	/i/ B ₄
38	/e/ y ₁	/u/ B ₅
39	/o/ x ₇	/o/ F ₄
40	/i/ F ₄	/a/ B ₁
41	/a/ x ₁	/o/ F ₅
42	/i/ x ₂	/o/ B ₅
43	/o/ B ₁	/e/ B ₅
44	/u/ x ₂	/o/ B ₂
45	/a/ x ₄	/a/ F ₅
46	/e/ x ₄	/a/ B ₅
47	/i/ x ₃	/i/ B ₅
48	/a/ y ₉	/a/ B ₂
49	/u/ B ₄	/i/ B ₁
50	/u/ F ₅	/i/ F ₅

Table 2 shows what kind of feature parameters are selected in the reference template as number of parameters is increased. The left table consists of glottal wave and vocal tract parameters and the right table consists of vocal tract parameters only. The table shows that no difference is found in both templates when the number of parameters is 7 or less. According to identification test, the identification rate is 16% by 1 parameter, 35% by 2 parameters, 49% by 3 parameters, 67% by 4 parameters, 77% by 5 parameters, 81% by 6 parameters and 83% by 7 parameters.

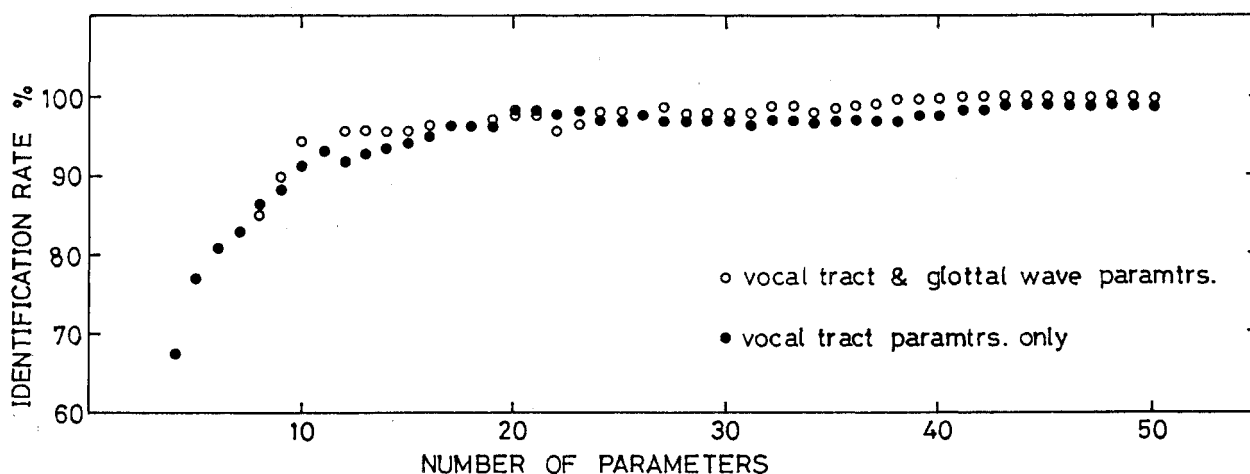


Fig.4 Relation between identification rate and number of parameters

Fig.4 shows the relation between the identification rate and number of parameters with the combined parameters and only with the vocal tract parameters. From the figure, the identification rate using combined parameters exceeds the identification rate using only vocal tract parameters by 0.5% to 2% when number of parameters is 24 or more. When number of parameters is between 17 and 23, however, the identification rate using only vocal tract parameters is better. This is caused by sorting parameters not in optimum order but in F-ratio order.

It is also found, from Table 2 and Fig.4, that 94% identification rate with a template of 10 feature parameters (8 vocal tract and 2 glottal parameters), 96% with a template of 16 feature parameters (8 vocal tract and 8 glottal parameters) and 99% with a template of 36 feature parameters (14 vocal tract and 22 glottal parameters) are obtained. In comparison with the number of parameters to get same identification rate only with formants' template, the number of feature parameters is reduced by 5 in case of getting 94% identification rate, 1 in case of getting 96% identification rate and 7 in case of getting 99% identification rate.

5. CONCLUSION

For the purpose of making a smaller sized template, we extract some glottal wave parameters which does not seem to be susceptible to mimicry and combine them with some LPC vocal tract parameters. As glottal parameters, abscissa and ordinate of subpulses plotted on a polar plane are used.

According to experimental results obtained

from 30 speakers' 5 vowels, glottal parameters show relatively high F-ratio. In comparison with the number of parameters to get same identification rate only with a formants' template, the number of feature parameters is reduced by 5 in case of getting 94% identification rate, 1 in case of getting 96% identification rate and 7 in case of getting 99% identification rate.

Further improvement of identification rate might be done by adaptive subpulse treatment.

The authors wish to thank Masayuki Abe, a graduate student of the Science University of Tokyo, for his technical assistance.

REFERENCES

- [1]S.Furui et al, "Talker recognition by longtime averaged speech spectrum," Trans. on IECEJ, 55-A, No.10, pp549-pp556, 1972
- [2]M.R.Sambur, "Speaker recognition using orthogonal linear prediction," IEEE Trans. on ASSP ASSP-24, No.4, pp.283-pp289, 1976
- [3]B.S.Atal, "Automatic speaker recognition based on pitch contours," J. Acoust.Soc.Amer., vol.52, pp.1687-1697, Dec., 1972.
- [4]S.Furui, "Cepstral Analysis Technique for Automatic Speaker Verification," IEEE Trans. on ASSP, ASSP-29, No.2, pp254-pp272, 1981
- [5]B.S.Atal et al, "A New Model of LPC Excitation for Producing Natural-Sounding Speech at Low Bit Rates," Proc. ICASSP, 82, pp614-617, 1982
- [6]K.Ozawa et al, "Speech Coding Based on Multi-pulse Excitation Method," Tech. Rep of IEICEJ, CS82-161, pp115-pp122, 1982
- [7]S.Furui, "Digital Speech Processing, Synthesis, and Recognition," pp302, Dekker, 1989