

A STUDY ON THE SPEECH SYNTHESIS-BY-RULE SYSTEM APPLIED MULTIBAND EXCITATION SIGNAL

Younjeong Kyung*, Geesoon Kim**, Hwangsoo Lee*, Yanghee Lee**

* Department of Information and Communication Engineering, Korea Advanced Institute of Science and Technology, P.O.Box 201, Cheongryang, Seoul, 130-650, KOREA

** Department of Computer Science, Dong-Duck Women's Univ. Wolgokdong 23-1, Sungbukgu, Seoul, 136-714, KOREA

ABSTRACT

In this paper, we design and implement the Korean speech synthesis by rule system. This system is applied the multiband excitation signal on voiced sounds. The multiband excitation signal is obtained by mixing impulse spectrum and white noise spectrum. We find that the quality of synthesized speech is improved using this application. Also, we classify the voiced sounds by cepstral euclidian distance measure for reducing overhead memory. The representative excitation signal of the same group's voiced sounds is used as excitation signal on synthesis. This method does not affect the quality of synthesized speech. As the result of experiment, this method eliminates the "buzziness" of synthesized speech and reduces the spectral distortion of synthesized speech.

1. INTRODUCTION

Generally, the speech wave production mechanism can be divided into two stages_[1]. These stages consist of sound source production, and articulation by vocal tract. Speech information processing techniques are based on the linear separable equivalent circuit model of the speech production mechanism. In the basic vocoder scheme, speech signals are generated from a simple source-filter model as illustrated in figure 1.

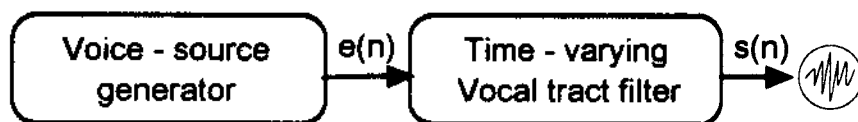


Fig.1 Basic speech production mechanism. The $e(n)$ corresponds to excitation signal and the $s(n)$ corresponds to synthesized speech signal.

The filter characteristics in this model are obtained from an analysis procedure and approximate the short-time power spectral envelope of the input speech. The source signal therefore has a flat spectral envelope and

contributes only fine spectral detail and phase information. In typical excited vocoders the source signal is either a periodic impulse train for generation voiced sounds or a random sequence for generation unvoiced sounds^[1]. This simple excitation signal causes the poor quality of synthetic speech such as the buzziness of voiced sounds and so forth. In this paper, we use the multiband excitation signal which is proposed in speech analysis/synthesis part on voiced sounds of TTS. Also we classify the voiced sounds by cepstral distance for reducing overhead that is required for multiband excitation signal.

In section 2, we briefly summarize our TTS system. In section 3, our grouping result is described. In section 4, our experimentation result and estimation is discussed. Section 5 concludes the paper with a discussion.

2. IMPROVE THE QUALITY OF SYNTHESIZED SPEECH BY MULTIBAND EXCITATION SIGNAL

Typically, in order to synthesize speech, the excitation signal consists of a periodic impulse train in voiced regions or random noise in unvoiced regions. This excitation signal is then filtered using the estimated system parameters. Even though vocoders based on this above-mentioned speech synthesis system have been quite successful in synthesizing intelligible speech, the vocoders have not been successful in synthesizing high quality speech. The example of poor quality of synthesized speech is 'buzziness'. Many researchers have noted that buzziness of synthetic speech could be caused by the use of a purely periodic excitation for segments of natural speech that have turbulent noise-like spectral regions.

Fig.2 shows the spectrum of voiced original speech. We can see the harmonics and noise-like spectrum component in this figure. We ought not to disregard this noise-like component for high quality synthesized speech^[2].

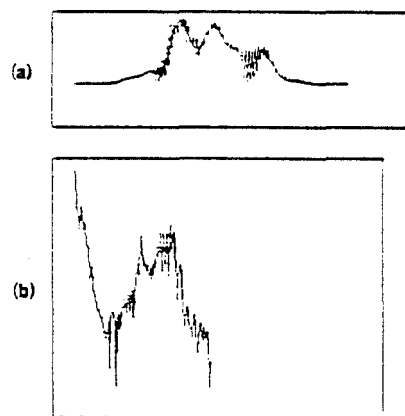


Fig. 2 This figure shows the original data : (a)Waveform, (b) spectrum

Therefore we use the multiband excitation signal. Fig.3 describes the procedure to obtain the multiband excitation signal. This processing follows the proposed method in speech analysis/synthesis fields. And fig.4 shows the multiband excitation signal which is used in our system. The general spectrum

of synthesized speech on voiced, $|Sp'(\omega)|$ is obtained by $|Sp'(\omega)| = |Ip(\omega)| \cdot |Se(\omega)|$ where $|Ip(\omega)|$ is impulse spectrum and $|Se(\omega)|$ is spectral envelope. But we have the spectrum of synthesized speech on voiced, $|Sp'(\omega)| = |(e)| \cdot |Se(\omega)|$ where (e) is obtained by mixing (c) and (d) in accordance with (b) - V/UV information of (a).

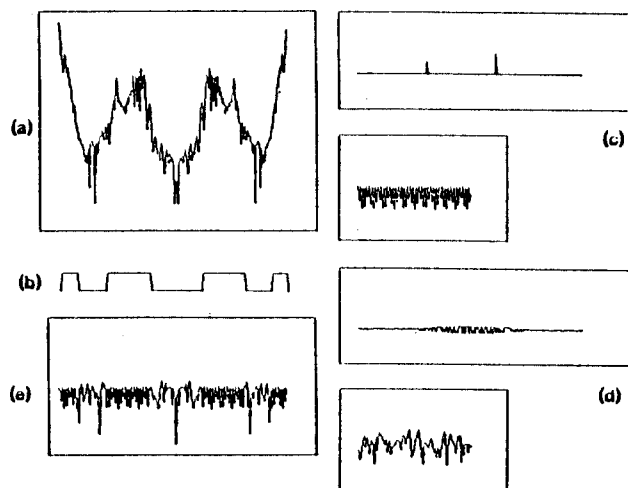


Fig. 3 The procedure of obtaining the mixed excitation signal (a)spectrum (b)V/UV info. (c)impulse spectrum (d)noise spectrum (e)mixed spectrum

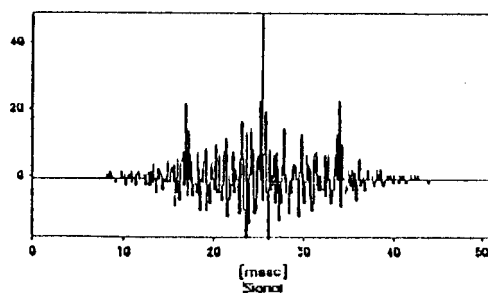


Fig.4 Multiband Excitation Signal

The summary of our system is as follows :

- A/D Conversion : 5KHz LPF, 10KHz Sampling, 12bit Quantization
- Analysis Parameter : Improved Cepstrum
- Analysis(Vocal Tract Approximation)
 - Frame period : 10ms for male, 5ms for female
 - Window function :
 - W - 25.6ms Blackman window for short-term predictor
 - Wp - 40ms Blackman window for long-term predictor
 - Cepstrum order : 30 for male, 25 for female
 - Spectral envelope : by Improved Cepstral method_[3]
 - V/UV determinant : by Spectrum Envelope parameter
 - Pitch detection : by Cepstrum peak picking_[4]
- Synthesis filter : LMA(Log Magnitude Approximation) filter_[5]
- Excitation signal :
 - White noise for unvoiced
 - Mixed signal(Impulse + Noise) for voiced

3. GROUPING OF VOICED SOUNDS FOR THE REPRESENTATIVE EXCITATION SIGNAL

Preliminary results indicate that high quality reproduction can be obtained with this speech synthesis system for both clean and noisy speech without the "buzziness".

The above-stated system has many multiband excitation signals correspondent to each voiced sound. In this case, large memory size is required for excitation signals. Thus we classify the voiced sounds by level of spectral distortion to reduce overhead memory.

Spectral distortion measure is cepstral distance measure. It is defined as follows :

$$Cd = 10 \ln 10 \sqrt{2 \sum_{i=1}^M (c^{v1}[i] - c^{v2}[i])^2}$$

$v1 \neq v2$: voiced sound

It is known as good measure of spectral distortion. Table 1 display the cepstral euclidian distance between two voiced sounds. Table 2 shows the groups of voiced sounds. The classification threshold value is about 2.5dB. Experimentally, the fact that the representative excitation signal of the same group's voiced sounds is used as excitation signal on synthesis, does not affect the quality of synthesized speech.

Table 1. The Cepstral distance of the voiced sounds

	a	o	u	i	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	ɛ	
a	-	2.35015, 59614, 22719, 79215, 11315, 88215, 88013, 17714, 40617, 12414, 14116, 69914, 2131																		
o		-	1.76011, 32717, 28612, 95613, 68713, 85712, 43212, 15514, 17414, 59114, 18912, 7841																	
u			-	1.74118, 02914, 28814, 95515, 34413, 03514, 06114, 32115, 75714, 70713, 7271																
i				-	7.46014, 11914, 50914, 77411, 82712, 52213, 36413, 36014, 85012, 9151															
ɛ					-	3.99313, 32113, 92417, 00617, 26415, 18118, 68312, 286110, 521														
ɛ						-	0.18310, 17113, 14512, 29412, 54714, 43010, 91716, 6901													
ɛ							-	0.09913, 57713, 10412, 31514, 93510, 62117, 5791												
ɛ								-	3.49213, 06713, 21114, 83410, 96617, 9621											
ɛ									-	2.24812, 96812, 25314, 35412, 9041										
ɛ										-	1.49113, 16013, 64412, 6571									
ɛ											-	3.82912, 61913, 5751								
ɛ												-	5.73814, 9341							
ɛ													-	7.7491						
ɛ														-						

Table 2. The Groups of the voiced sounds

Group 1	j, ε, φ, e
Group 2	u, o, ɔ
Group 3	ɪ, ɨ, w
Group 4	n
Group 5	m
Group 6	a
Group 7	i

4. IMPLEMENTATION AND ESTIMATION

Fig.5 shows the waveform of TTS using the representative excitation signal of each group. (a) is the original, (b) is the synthesized speech by impulse train, (c) is the synthesized speech by mixed excitation signal. We can find the (c) is better than (b) in this figure. The method of using the multiband excitation signal is so useful for variation of the pitch in TTS. Fig.6 shows the waveform of synthesized speech which is obtained by multiplication of original pitch value by 1.5. (a) is the synthesized speech by impulse train, (b) is by multiband excitation signal. In the case of (a), as pitch period is longer, synthesized speech's quality is poorer. But in the case of (b), although pitch period is longer, synthesized speech's quality is not especially distorted.

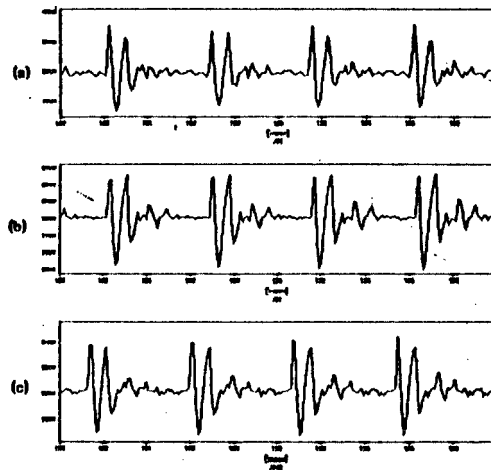


Fig. 5 Synthesized speech waveform (a)original data (b)by impulse (c)by multiband excitation signal

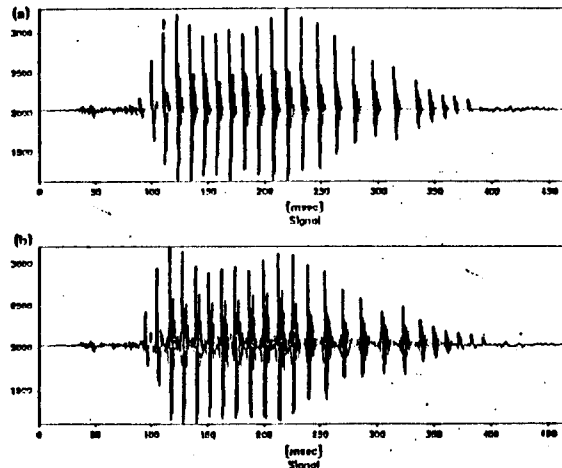


Fig. 6 Synthesized speech waveform is obtained by variation of the pitch (a)by impulse (b)by multiband excitation signal

Fig.7 is the example of TTS synthesized by the mixed excitation signal. (a) is the example of TTS by impulse, (b) is by the mixed excitation signal. It is the waveform of /annjε̃hasim/- a part of synthetic speech /annjε̃hasimnik'a/. According to the observations on synthetic waveform and the result of subjective test, we conclude that the method using the multiband excitation signal is better than that of using the impulse train.

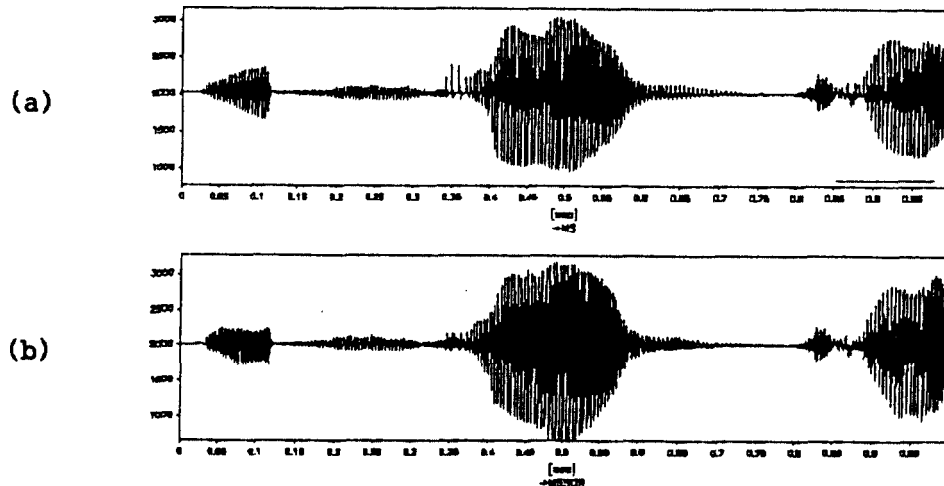


Fig. 7 The example of TTS synthesized speech (a)by impulse (b)by multiband excitation signal

5. CONCLUSION

In this paper, we inquire the problems of the common TTS and propose the solution. The solution is the application of the multiband excitation signal to TTS. As the result of experiment, it is found that the application improves the synthetic speech quality, eliminates the "buzziness". In addition, it reduces the spectral distortion that is caused by the variation of pitch period in TTS. Also, we classify the voiced sounds to reduce overhead memory by cepstral euclidian distance measure.

REFERENCES

- [1] S. Furui, *Digital Speech Processing, Synthesis, and Recognition*, MARCEL DEKKER, 168-169 (1992).
- [2] Denial W. Griffin, et al., "A High Quality 9.6kbps Speech Coding System", ICASSP, 125-128 (1986).
- [3] S. Inai and Y. Abe, "Spectral envelope extraction by improved cepstral method.", Trans IECE Vol.62-A No.4, 217-223 (1979).
- [4] A. Michael Noll, "Cepstrum Pitch Determination", J.Acoust.Soc.Amer., 179-195 (1967).
- [5] S. Inai, et al., "Log Magnitude Approximation (LMA)filter", Trans. IECE Vol.J63-A No.12, 886-893 (1981).