

A design of adaptive beamformer based on average speech spectrum for noisy speech recognition

Yuka Okada¹, Takanobu Nishiura^{2,3}, Satoshi Nakamura³, Takeshi Yamada⁴ and Kiyohiro Shikano¹

¹Graduate School of Information Science, Nara Institute of Science and Technology, 8916-5 Takayama, Ikoma, 630-0101 Japan

²Faculty of Systems Engineering, Wakayama University, 930 Sakaedani, Wakayama, 640-8510 Japan

³ATR Spoken Language Translation Research Laboratories, 2-2-2 Hikaridai Seika-cho, Soraku-gun, Kyoto, 619-0288 Japan

⁴Institute of Information Sciences and Electronics, University of Tsukuba, 1-1-1, Tennoudai, Tsukuba, 305-8573 Japan

(Received 17 April 2002, Accepted for publication 25 June 2002)

Keywords: Microphone array, Adaptive beamformer, Average speech spectrum, Noisy speech recognition, Real environments
PACS number: 43.72.Ew

1. Introduction

It is very important for the natural interfaces of machines like self-moving robots to capture and recognize distant-talking speech with high accuracy. However, background noise and room reverberations seriously degrade the sound capture quality in the real acoustic environments. A microphone array is an ideal candidate for capturing distant-talking speech. With a microphone array, a desired speech signal can be acquired selectively by steering the directivity. Accordingly, super-high directivity is necessary to reduce noise signals.

To form directivity, delay-and-sum beamformers [1,2] and adaptive beamformers [3,4] have been proposed as conventional beamformers. A delay-and-sum beamformer forms super-high directivity to the desired signal, and an adaptive beamformer forms null directivity to the noise signal. However, delay-and-sum beamformers have two serious drawbacks: the performance is not good enough to capture the desired signal without a sufficient number of transducers, and performance degrades in highly reverberant rooms. On the other hand, adaptive beamformers can form null directivity with a small number of transducers. Furthermore, they can form sharper directivity than delay-and-sum beamformer. Consequently, adaptive beamformers are often used for the front-end processing of ASR (Automatic Speech Recognition) [5].

AMNOR (Adaptive Microphone-array for NOise Reduction) [4] is an adaptive beamformer proposed by Kaneda et al. in 1986. AMNOR is an effective beamformer for capturing and recognizing desired distant signals in noisy environments. Also, it can be easily designed with an adaptive filter for noise reduction in real environments, because it only allows small distortion for capturing the desired distant signal.

However, if we knew the spectrum characteristics of desired distant signals when designing the adaptive filter of AMNOR, its performance could be further improved. The conventional AMNOR is designed to suppress the spectrum distortion of the desired distant signal on all frequency bands, but in many cases, the purpose of signal capture is limited to

speech capture. Therefore, in this paper we regard speech as the desired distant signal and design AMNOR by using the speech spectrum for distant-talking speech capture and recognition.

2. AMNOR (Adaptive Microphone-array for NOise Reduction)

Figure 1 shows a block diagram of the adaptive beamformer. In Fig. 1, $S(\omega)$ is the Fourier transform of the desired signal and $Y(\omega)$ is the Fourier transform of the output signal. $G_m(\omega)$ is the acoustic transfer function from the desired sound source to the m -th microphone element and $H_m(\omega)$ is the frequency response of the m -th filter. The frequency response $F(\omega)$ of the adaptive beamformer to the desired signal is represented as

$$F(\omega) = \sum_{m=1}^M G_m(\omega)H_m(\omega), \quad (1)$$

where M is the number of microphone elements. The concept of the adaptive beamformer is to minimize the output noise energy while constraining $F(\omega)$ to the desired frequency response. AMNOR [4] has the constraint shown in Eq. (2):

$$D = \int |1 - F(\omega)|^2 d\omega \leq \hat{D}. \quad (2)$$

This constraint attains maximum noise reduction while allowing a small distortion D in the frequency response to

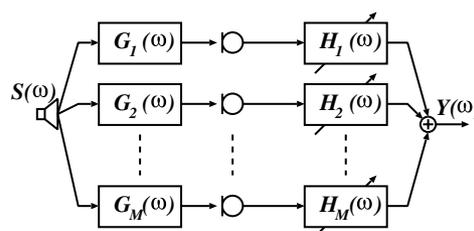


Fig. 1 Block diagram of adaptive beamformer.

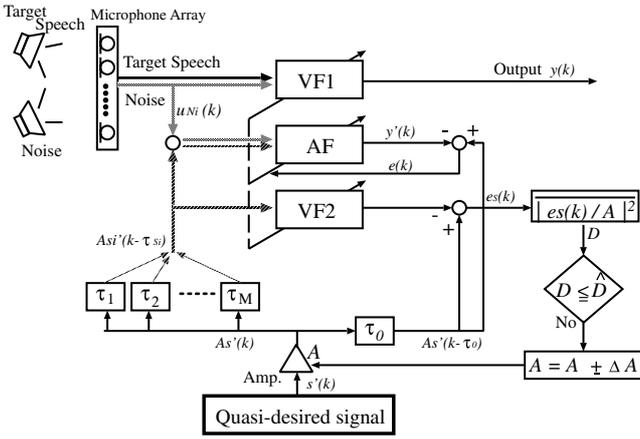


Fig. 2 Overview of AMNOR.

the desired signal. In this paper, we focus on suitable control of the admissible distortion \hat{D} in the frequency response for noisy speech recognition. Figure 2 shows a general overview of AMNOR. In Fig. 2, each VF1, AF, and VF2 is a FIR filter with M -input and 1-output. AF is the adaptive filter, and VF1 and VF2 are variable filters that have the same filter coefficients as AF. A quasi-desired signal $s'(k)$ is indispensable for designing the adaptive filter of AMNOR because AMNOR attains maximum noise reduction with a quasi-desired signal and an environmental noise signal from the environment. The quasi-desired signal $s'(k)$ derives $As'_i(k - \tau_{si})$ from amplifier and time delay τ_{si} , $i = 1, \dots, M$, which is calculated subject to the known desired sound source's DOA (Direction Of Arrival). This situation assumes the simulation where signal $As'(k)$ arrives from the desired sound source with known DOA to the microphone array. In addition, the microphone only captures the noise signal $u_{Ni}(k)$, $i = 1, \dots, M$ (not including the desired signal), and it is inputted in the adaptive filter AF after adding it to quasi-desired signal $As'_i(k - \tau_{si})$. AF controls the filter coefficients based on $e(k)$ as the following Eq. (3).

$$e(k) = As'(k - \tau_0) - y'(k), \quad (3)$$

where τ_0 is the constant delay for cause and effect. $es(k)$ is calculated by using VF2 after designing the filter coefficients by AF, and current distortion D is derived from Eq. (4).

$$D = \overline{|es(k)/A|^2}. \quad (4)$$

By comparing current distortion D and admissible distortion \hat{D} , amplitude A is renewed with the amplifier until $D \leq \hat{D}$. In the above algorithm, AMNOR attains higher noise reduction performance in real acoustic environments.

3. Suitable design of AMNOR based on average speech spectrum

The conventional AMNOR uses a white Gaussian signal that has flat frequency characteristics as a quasi-desired signal in order to suppress the spectrum distortion of the desired signal on all frequency bands. But in many cases, the purpose of signal capture is limited to speech capture. Therefore, if we knew the spectrum characteristics of desired distant signals in

advance, it may be possible to improve the performance of AMNOR by designing a suitable adaptive filter for the environment. In this paper, we regard speech as the desired distant signal and design AMNOR by using the speech spectrum for distant talking speech capture and recognition. First, we calculate the average speech spectrum weight by Eq. (5).

$$W_{sp}(\omega) = \frac{1}{L \cdot N} \sum_{l=1}^L \sum_{n=1}^N SP_l(\omega; n), \quad (5)$$

where L represents the number of speech (words), N represents the number of frames, $SP_l(\omega; n)$ represents the Fourier transform of speech signed $sp_l(t)$, and $W_{sp}(\omega)$ represents the average speech spectrum weight. The quasi-desired signal based on the average speech spectrum is derived from weighting the white Gaussian spectrum with the average speech spectrum weight $W_{sp}(\omega)$. Figure 3 shows the spectrum of white Gaussian as the quasi-desired spectrum for the conventional AMNOR and the spectrum of average speech weighted as quasi-desired spectrum for the proposed AMNOR. Compared with the spectra in Fig. 3, the average speech weighted spectrum is enhanced at lower frequencies. We attempted to improve the ASR performance by using the average speech spectrum weighted quasi-desired signal for AMNOR, and this modified system was named S-AMNOR.

In addition, we also investigated whether the average speech spectrum weight is normalized to keep the energy ratio equivalent between vowels and consonants on each frame when estimating $W_{sp}(\omega)$ in Eq. (5). We further consider a new spectrum weight defined by Eq. (6). This weight is capable to balance the occurrence of vowel and consonant frames.

$$W_{sp}(\omega) = \frac{1}{2} \left(\frac{1}{L_c} \sum_{l_c=1}^{L_c} \frac{1}{N_{l_c}} \sum_{n=1}^{N_{l_c}} SP_{l_c}(\omega; n) + \frac{1}{L_v} \sum_{l_v=1}^{L_v} \frac{1}{N_{l_v}} \sum_{n=1}^{N_{l_v}} SP_{l_v}(\omega; n) \right), \quad (6)$$

where L_v represents the number of vowels, L_c represents the number of consonants, N_{l_v} represents the number of vowel frame on each speech (word), and N_{l_c} represents the number of consonant frame on each speech (word). The system using this modified W_{sp} was named Normalized S-AMNOR.

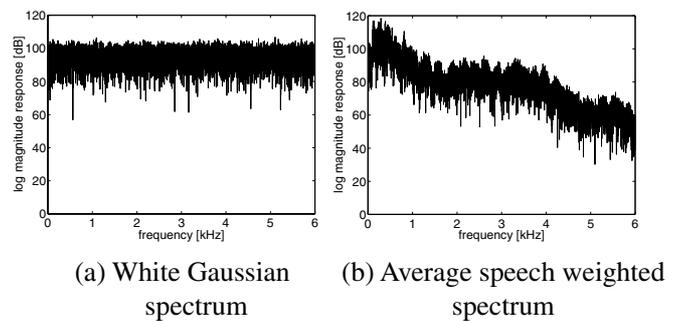


Fig. 3 Spectrum of quasi-desired signal.

4. Evaluation experiments

4.1. Experimental conditions

We evaluated the ASR performance in a real acoustic room. Figure 4 shows the experimental environment, and Table 1 shows the experimental condition. The desired distant signal arrives from the front direction (90 degrees), and the noise signal arrives from the right and left directions (40 degrees and 120 degrees, respectively). The distance between the sound source and the microphone array is two meters. In this situation, the ASR performance was evaluated by variations in the admissible distortion \hat{D} as Eq. (2). ASR

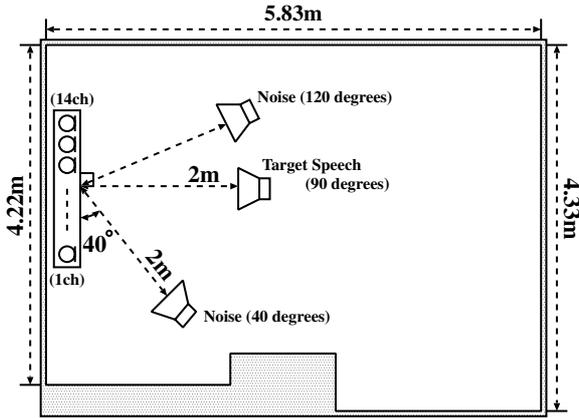


Fig. 4 Experimental environment.

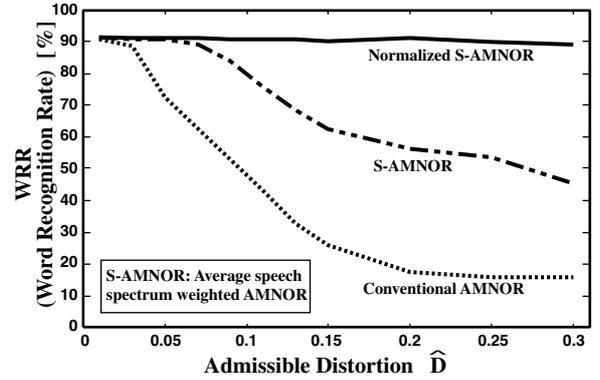
Table 1 Experimental conditions.

Recording conditions	
Reverberation time	$T_{[60]} = 180$ ms
Microphone array	Linear type 14 transducers, 2.83 cm spacing
Sampling frequency	12 kHz (Quantization: 16 bit)
Experimental conditions for ASR	
Frame length	32 ms (Frame interval: 8 ms)
HMM	Gaussian mixture density (3 states)
Feature vector	MFCC (16 orders, 4 mixtures), Δ MFCC (16 orders, 4 mixtures), Δ power (1 order, 2 mixtures)
Average speech spectrum weight	
Speech DB	ATR speech DB SetA [6] and ASJ continuous speech corpus [7]
Speech (L)	2,620 words \times 4 subjects and 150 sentences \times 64 subjects
Consonants (L_c, N_c)	L_c : 28,156 phonemes, N_c : 94,315 frames ($= \sum N_{l_c}$)
Vowels (L_v, N_v)	L_v : 23,440 phonemes, N_v : 107,085 frames ($= \sum N_{l_v}$)
Test data (Open)	
Desired speech signal	Speech: 216 words \times 2 subjects (1 female and 1 male)
Noise signal	Female speech, male speech or white Gaussian noise
SNR	3 dB

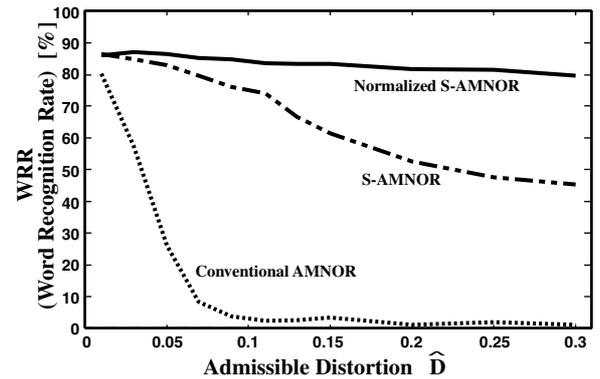
performance was also evaluated by the WRR (Word Recognition Rate).

4.2. Experimental results for ASR performance

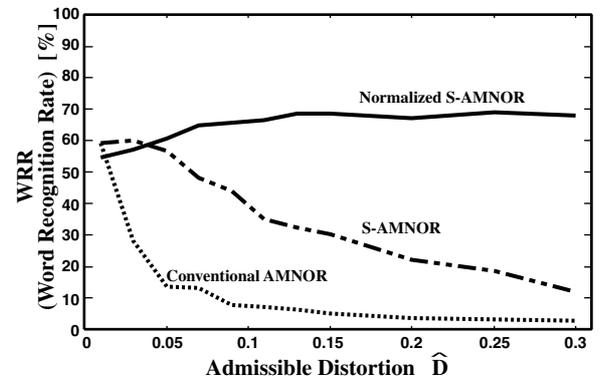
Figure 5 shows the ASR performance in the evaluation environment. In this experiment, the sound source position is known for designing the adaptation filter. In Fig. 5, (a) shows the results in an environment of one desired speech [90 degrees DOA], (b) shows the results in an environment of one desired speech [90 degrees DOA] and one noise (female



(a) Environment of one desired speech.



(b) Environment of one desired speech and one noise (female speech).



(c) Environment of one desired speech and two noises (female speech and white Gaussian noise).

Fig. 5 ASR performance.

speech [40 degrees DOA]), and (c) shows the results in an environment of one desired speech [90 degrees DOA] and two noises (female speech [40 degrees DOA] and white Gaussian signal [120 degrees DOA]).

As a result of our evaluation experiments, we could confirm that the average speech spectrum weighted AMNOR (S-AMNOR) provides higher ASR performance than the conventional AMNOR. The effectiveness of S-AMNOR is also confirmed with large admissible distortion \hat{D} . It also showed the same tendency in the environment of one desired speech [90 degrees DOA] and one noise (male speech [40 degrees DOA]) and in that of one desired speech [90 degrees DOA] and one noise (white Gaussian noise [40 degrees DOA]). In addition, we could confirm that the normalized speech spectrum weighted AMNOR (Normalized S-AMNOR) is more effective than the basic S-AMNOR. This is because the adaptive filter of Normalized S-AMNOR has a more greatly optimized energy balance between vowels and consonants than that of S-AMNOR.

Figure 5(c) shows that the Normalized S-AMNOR with large admissible distortion \hat{D} may provide higher ASR performance than that with small admissible distortion \hat{D} in noisy environment. Therefore, we could guess that the spectrum distortion of the desired signal may be suppressed even with large admissible distortion \hat{D} , if we can acquire the optimum quasi-desired signal in advance. As a result, ASR performance improves with large admissible distortion \hat{D} because noise signal can be vastly reduced by Normalized S-AMNOR with large admissible distortion \hat{D} .

Next, we show the maximum ASR performance with the optimum admissible distortion \hat{D} in Fig. 6, which we manually selected. In Fig. 6, we confirm that if we estimate the optimum admissible distortion \hat{D} in advance, ASR performance improves by 5–10% with the normalized speech spectrum weight in a noisy environment.

4.3. Experimental results for the spectrum characteristics of designed adaptive filter

Figure 7 shows the spectrum characteristics of adaptive filters. We investigated performance with the signal as a flat spectrum characteristic and the designed adaptive filters. In Fig. 7, (a) shows the input spectrum with a white Gaussian

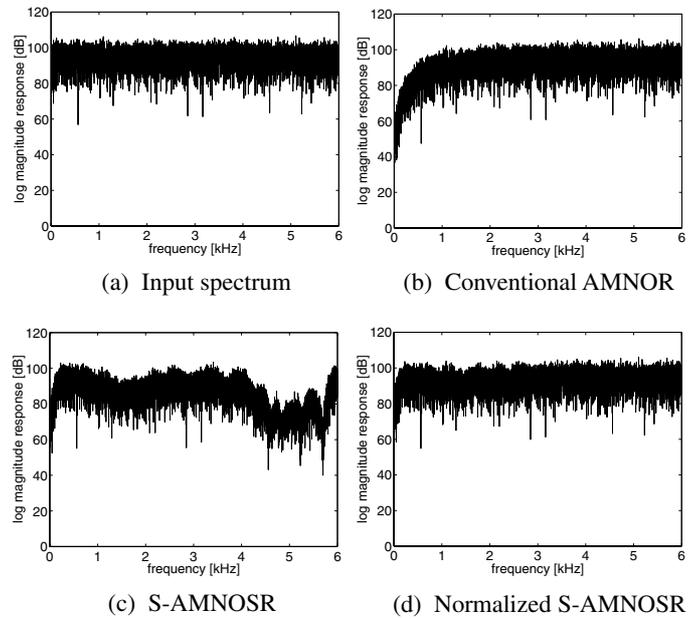


Fig. 7 Spectrum characteristics of adaptive filter based on average speech spectrum weight when admissible distortion $\hat{D} = 0.1$.

signal as the flat spectrum characteristic, (b) shows the output spectrum with an adaptive filter based on the conventional AMNOR, (c) shows the output spectrum with an adaptive filter based on the average speech spectrum weight (S-AMNOR), and (d) shows the output spectrum with an adaptive filter based on the normalized average speech spectrum weight (Normalized S-AMNOR). By comparing the results from Figs. 7(b) and 7(d), we could confirm that the adaptive filter based on the normalized average speech spectrum weight (Normalized S-AMNOR) shows almost no distortion on any frequency band although the adaptive filter based on the white Gaussian (Conventional AMNOR) shows severe distortion in the lower frequency bands which are indispensable for speech recognition. In addition, we could also confirm, by comparing the results using Figs. (c), and (d), that normalization of the average speech spectrum (Normalized S-AMNOR) improves the signal capturing performance. In the above evaluation experiments, we confirmed that AMNOR based on the normalized average speech spectrum weight (Normalized S-AMNOR) is more effective than the conventional AMNOR for noisy speech recognition.

5. Conclusions

In this paper, we proposed a method to improve ASR performance by AMNOR (Adaptive Microphone-array for NOise Reduction) with the average speech spectrum weight in noisy environments. As a result of evaluation experiments in real acoustic environments, we confirmed that ASR performance is improved by using the normalized average speech spectrum weighted AMNOR (Normalized S-AMNOR). In the future, we will improve ASR performance by integrating the proposed AMNOR with talker localization [8] and automatically estimating the optimum admissible distortion \hat{D} for ASR in noisy environments.

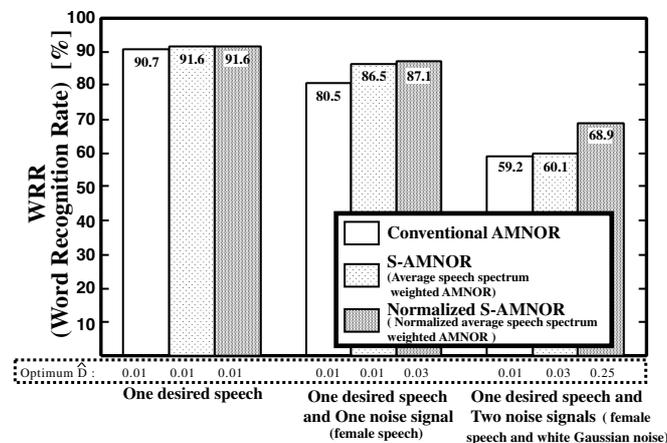


Fig. 6 ASR performance with optimum admissible distortion \hat{D} .

Acknowledgement

This research was partially supported by The Ministry of Education, Culture, Sports, Science and Technology of Japan under Grant-in-Aid No. 14780288.

References

- [1] J. L. Flanagan, J. D. Johnston, R. Zahn and G. W. Elko, "Computer-steered microphone arrays for sound transduction in large rooms," *J. Acoust. Soc. Am.*, **78**, 1508–1518 (1985).
- [2] S. U. Pillai, *Array Signal Processing* (Springer-Verlag, New York, 1989).
- [3] L. J. Griffiths and C. W. Jim, "An alternative approach to linearly constrained adaptive beam-forming," *IEEE Trans. AP*, **AP-30**, 27–34 (1982).
- [4] Y. Kaneda and J. Ohga, "Adaptive microphone-array system for noise reduction," *IEEE Trans. Acoust. Speech Signal Process.*, **ASSP-34**, 1391–1400 (1986).
- [5] M. Omologo, M. Matassoni, P. Svaizer and D. Giuliani, "Microphone array based speech recognition with different talker-array position," *Proc. ICASSP 97*, pp. 227–230 (1997).
- [6] K. Takeda, Y. Sagisaka and S. Katagiri, "Acoustic-phonetic labels in a Japanese speech database," *Proc. Eur. Conf. Speech Technology*, Vol. 2, pp. 13–16 (1987).
- [7] T. Kobayashi, S. Itahashi and T. Takezawa, "ASJ continuous speech corpus for research," *J. Acoust. Soc. Jpn. (J)*, **48**, 888–893 (1992).
- [8] T. Nishiura, S. Nakamura and K. Shikano, "Statistical sound source identification in real acoustic environment for robust speech recognition using a microphone array," *Proc. EURO-SPEECH 2001*, pp. 2611–2614 (2001).