

Blind Source Separation for Convolved Signals Based on Properties of Acoustic Transfer Function in Real Environments

Takaaki Ishibashi

Abstract—Frequency domain independent component analysis has a scaling indeterminacy and a permutation problem. The scaling indeterminacy can be solved by use of a decomposed spectrum. For the permutation problem, we have proposed the rules in terms of gain ratio and phase difference derived from the decomposed spectra and the source's coarse directions.

The present paper experimentally clarifies that the gain ratio and the phase difference work effectively in a real environment but their performance depends on frequency bands, a microphone-space and a source-microphone distance. From these facts it is seen that it is difficult to attain a perfect solution for the permutation problem in a real environment only by either the gain ratio or the phase difference.

For the perfect solution, this paper gives a solution to the problems in a real environment. The proposed method is simple, the amount of calculation is small. And the method has high correction performance without depending on the frequency bands and distances from source signals to microphones. Furthermore, it can be applied under the real environment. From several experiments in a real room, it clarifies that the proposed method has been verified.

Keywords—blind source separation, frequency domain independent component analysis, permutation correction, scale adjustment, target extraction.

I. INTRODUCTION

MANY noise reduction methods using ICA (independent component analysis) [1-5] have been proposed. ICA can separate unknown sources from their mixtures without information on the transfer functions, provided that the sources are statistically independent. For the instantaneous mixtures, the original sources can be completely recovered in the time domain except for indeterminacy of scaling and permutation.

In a real environment, the signals observed at microphones are not instantaneous mixtures but are convolved version of the sound sources. On account of this, there have been reported many trials to separate the convolved mixtures in the frequency domain.

However, in the frequency domain ICA, the indeterminacy of scaling and permutation appears at every frequency. In order to recover the sources properly, these indeterminacies must be essentially solved before making an inverse transformation from the frequency to the time domain.

In order to solve the problem, several methods have been proposed. Scaling indeterminacy, i.e. amplitude and phase indeterminacy can be solved by use of decomposed spectrum [5].

T. Ishibashi is with the Department of Information, Communication and Electronic Engineering, Kumamoto National College of Technology, Kumamoto, Japan e-mail: ishishashi@kncet.ac.jp.

We have derived that the decomposed spectrum is uniquely expressed as a product of a source spectrum and a transfer function from a source to a microphone [6].

For the permutation problem, there have been tried a method using similarities between separate spectra [5], a method taking advantage of directivity of array-microphones [7]. The authors have proposed some rules in terms of the gain ratio and the phase difference derived from the decomposed spectra and the source's coarse directions [6]. These correction methods have problems; its calculation is not simple and it does not function well at such frequencies that the sound component is very small to be susceptible to ambient noises.

In the present paper, we newly define the gain and phase differences between two acoustic transfer functions from a single source to two microphones. These differences are easily calculated in terms of the decomposed spectra and are also useful quantities to resolve the permutation problem. Their characteristics are fully examined in a real environment. It is clarified that the characteristics depend on frequency bands, a microphone-space and a source-microphone distance, and that these facts are also true for permutation correction rules.

For the perfect solution of the permutation problem, the paper gives a solution to the problems in a real environment. The proposed method is simple, the amount of calculation is small. And the method has high correction performance without depending on the frequency bands and distances from source signals to microphones. Furthermore, it can be applied under the real environment. From several experiments in a real room, it clarifies that the proposed method has been verified.

II. BLIND SOURCE SEPARATION

A. Independent Component Analysis in Frequency Domain

The source signals $s_i(t)$ ($i = 1, 2$) are denoted by

$$\mathbf{s}(t) = [s_1(t), s_2(t)]^T \quad (1)$$

under the assumption that each component $s_i(t)$ is statistically independent of each other. The observed mixtures $x_j(t)$ ($j = 1, 2$) are represented by

$$\mathbf{x}(t) = [x_1(t), x_2(t)]^T. \quad (2)$$

Consider the case that two statistically independent sound sources are observed by two microphones as

$$\mathbf{x}(t) = G(t) * \mathbf{s}(t) \quad (3)$$

where $*$ denotes the convolutional operator and $G(t)$ a matrix whose elements are transfer functions from the sources to the microphones.

The mixtures are transformed into the short time spectra by the discrete Fourier transform:

$$x_j(\omega, k) = \sum_t e^{-\sqrt{-1}\omega t} x_j(t) w(t - k\tau) \quad (4)$$

where $\omega (= 0, 2\pi/M, \dots, 2\pi(M-1)/M)$ denotes a frequency bin, M the number of samples in a frame, k the frame number, τ the frame shift, $w(t)$ a window function.

In the frequency domain, the spectra $\mathbf{x}(\omega, k)$ of the convolved mixtures are expressed as a product of the source spectra $\mathbf{s}(\omega, k)$ and the frequency transfer function $G(\omega)$.

$$\mathbf{x}(\omega, k) = G(\omega)\mathbf{s}(\omega, k) \quad (5)$$

In order to estimate the unknown sources, the separating matrix $H(\omega)$ is estimated by frequency domain ICA algorithm [1-5]. The separated spectra $\mathbf{u}(\omega, k) = [u_1(\omega, k), u_2(\omega, k)]^T$ can be obtained as

$$\mathbf{u}(\omega, k) = H(\omega)\mathbf{x}(\omega, k). \quad (6)$$

It seems to be possible that the separated signal $u_i(t)$ ($\hat{i} = 1, 2$) for the sound source $s_i(t)$ can be obtained from the separated spectra $u_i(\omega, k)$ by simply applying the inverse transform to time domain as

$$u_i(t) = \frac{1}{2\pi} \frac{1}{W(t)} \sum_k \sum_\omega e^{\sqrt{-1}\omega(t-k\tau)} u_i(\omega, k) \quad (7)$$

where $W(t) = \sum_k w(t - k\tau)$.

However, the separated spectra $u_i(\omega, k)$ have the problem with permutation and scaling indeterminacy. The order \hat{i} of the separated spectrum $u_i(\omega, k)$ is not necessarily consistent with i of the source spectrum $s_i(\omega, k)$, and the scale of $u_i(\omega, k)$ is not consistent with that of $s_i(\omega, k)$. Therefore, the indeterminacy of permutation, amplitude and phase must be settled to get a meaningful signal $u_i(t)$ before inversely transforming $u_i(\omega, k)$ from the frequency domain to the time domain.

B. A Solution for Scaling Indeterminacy

In order to solve the scaling indeterminacy, a decomposed spectrum $\mathbf{v}_i(\omega, k) = [v_{i1}(\omega, k), v_{i2}(\omega, k)]^T$ ($\hat{i} = 1, 2$) is introduced follows [5].

$$\mathbf{v}_1(\omega, k) = H^{-1}(\omega)[u_1(\omega, k), 0]^T \quad (8)$$

$$\mathbf{v}_2(\omega, k) = H^{-1}(\omega)[0, u_2(\omega, k)]^T \quad (9)$$

We have derived that the decomposed spectrum $v_{i\hat{j}}(\omega, k)$ is uniquely determined as a product of the source spectrum $s_i(\omega, k)$ and the transfer function $g_{j\hat{i}}(\omega)$, although their combination differ depending on whether permutation occur or not [6]. This fact means that the scaling factor of the decomposed spectrum is the transfer function itself and the decomposed spectrum has no scaling indeterminacy. Table I shows the correspondence between the decomposed spectrum and the sources.

TABLE I
DECOMPOSED SPECTRA IN PERMUTATION AND THOSE IN NO-PERMUTATION

		$g_{11}s_1$	$g_{21}s_1$	$g_{12}s_2$	$g_{22}s_2$
No-permutation	($\hat{i} = i$)	v_{11}	v_{12}	v_{21}	v_{22}
Permutation	($\hat{i} \neq i$)	v_{21}	v_{22}	v_{11}	v_{12}

TABLE II
DECISION RULE FOR PERMUTATION

No-permutation	($\hat{i} = i$)	$ r_1(\omega) > r_2(\omega) $	$\angle r_1(\omega) > \angle r_2(\omega)$
Permutation	($\hat{i} \neq i$)	$ r_1(\omega) < r_2(\omega) $	$\angle r_1(\omega) < \angle r_2(\omega)$

C. A Solution for Permutation

It assumes that one source is closer to the first microphone and another source is closer to the second microphone. From this assumption, the gain and the phase inequalities on the transfer function are derived.

$$|g_{ii}(\omega)| > |g_{ji}(\omega)| \quad (10)$$

$$\angle g_{ii}(\omega) > \angle g_{ji}(\omega) \quad (11)$$

It is found from Table I that of the source $s_i(\omega, k)$ there exist two candidate estimates $v_{i,j=1}(\omega, k) = g_{j=1,i}(\omega)s_i(\omega, k)$ and $v_{i,j=2}(\omega, k) = g_{j=2,i}(\omega)s_i(\omega, k)$. Here, we adopt $v_{i,j=i}(\omega, k) = g_{j=i,i}(\omega)s_i(\omega, k)$ as the estimate of $s_i(\omega, k)$ because the decomposed spectra observed at the nearer microphone ($j = i$) are larger in power and are less disturbed by ambient noise than those observed at the other microphone ($j \neq i$).

The ratio $r_i(\omega)$ between two decomposed spectra is defined.

$$r_i(\omega) = \frac{1}{K} \sum_{k=0}^{K-1} \frac{v_{i1}(\omega, k)}{v_{i2}(\omega, k)} \quad (12)$$

This ratio leads to basic two permutation decision rules in terms of $|r_i(\omega)|$ and $\angle r_i(\omega)$ as shown Table II. From these decision rules, further, two type of permutation correction methods are derived: one is based on the gain inequality, and the other the phase inequality.

[Rule1 : gain-based method]

$$y_i(\omega, k) = \begin{cases} v_{i=i,j=i}(\omega, k) & \text{if } |r_1(\omega)| > |r_2(\omega)| \\ v_{i \neq i,j=i}(\omega, k) & \text{if } |r_1(\omega)| < |r_2(\omega)| \end{cases} \quad (13)$$

[Rule2 : phase-based method]

$$y_i(\omega, k) = \begin{cases} v_{i=i,j=i}(\omega, k) & \text{if } \angle r_1(\omega) > \angle r_2(\omega) \\ v_{i \neq i,j=i}(\omega, k) & \text{if } \angle r_1(\omega) < \angle r_2(\omega) \end{cases} \quad (14)$$

The amount of calculation in either of these correction methods is much less than that in the method based on the power of decomposed spectra [6].

After the above permutation correction, the source spectra $y_i(\omega, k)$ is given by

$$y_i(\omega, k) = g_{ii}(\omega)s_i(\omega, k). \quad (15)$$

This implies that $y_i(\omega, k)$ can be an estimate of $s_i(\omega, k)$.

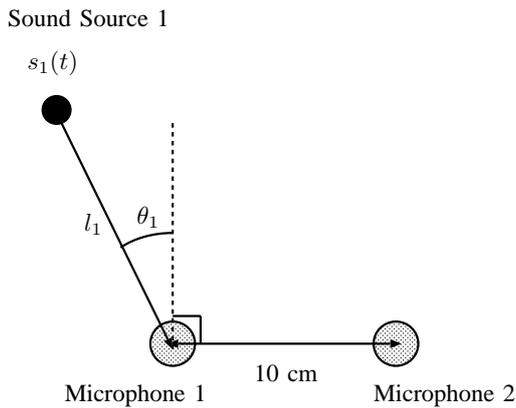


Fig. 1. Placement of one sound source and two microphones

TABLE III
RATIO WHERE EQS. (10) AND (11) HOLDS IN REAL ENVIRONMENT

Source location	10 cm	30 cm
Gain inequality	99.88 %	98.50 %
Phase inequality	97.84 %	99.64 %

III. PROPERTIES OF ACOUSTIC TRANSFER FUNCTION

A. Frequency properties of acoustic transfer function ratio

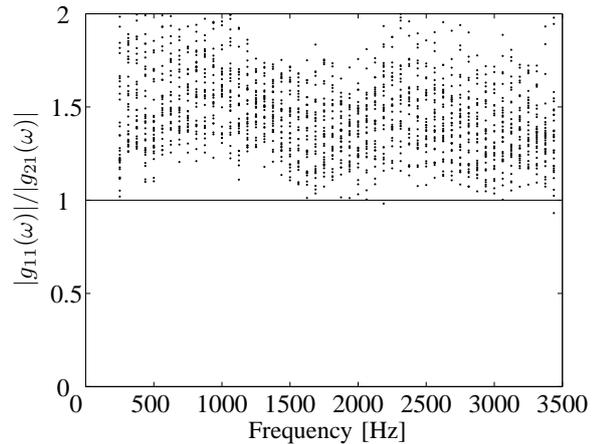
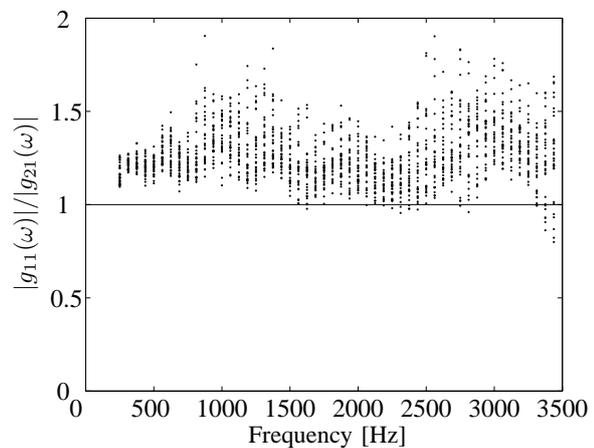
Under the condition that only $s_1(\omega, k)$ is active while $s_2(\omega, k) = 0$, the mixtures are described by $x_1(\omega, k) = g_{11}(\omega, k)s_1(\omega, k)$ and $x_2(\omega, k) = g_{21}(\omega, k)s_1(\omega, k)$. The ratio of mixtures $r_x(\omega, k) = x_1(\omega, k)/x_2(\omega, k)$ is found to be equal to $g_{11}(\omega)/g_{21}(\omega)$, i.e., the ratio of the transfer functions from source 1 to microphone 1 and microphone 2. Therefore it can be measured by the ratio $r_x(\omega, k)$ whether the inequalities in Eqs. (10) and (11) hold. The quantities $|r_x(\omega)|$ and $\angle r_x(\omega)$ respectively mean a gain ratio and a phase difference between $g_{11}(\omega)$ and $g_{21}(\omega)$.

Several experiments on the gain ratio and the phase difference were carried out in a room (L7.3 m \times W6.5 m \times H2.9 m). The room reverberation time was 500 ms and the ambient noise was 48 dB. Fig.1 shows the configuration of a source and two microphones. The source distance l_1 from microphone 1 was set at 10 cm and 30 cm with a fixed angle $\theta_1 = 10^\circ$. The two microphones with 10 cm spacing were condenser microphones with band width 200-5000 Hz.

At every l_1 , 32 data sets (8 humans \times 4 words [8]) were observed. The mixture signals were sampled at a rate of 8000 Hz with 16 bit resolution. The sampled data were windowed by the Hamming window with a frame length 16 ms and a frame period 8 ms.

Table III shows the ratios where the inequalities of Eqs. (10) and (11) hold. From these results, the inequalities can be confirmed to hold almost surely.

Figures 2 and 3 are the plots of the gain ratio $|r_x(\omega)|$ at $l_1=10$ and 30 cm. In these plots, the exceptions to $r_x(\omega) > 1$ are found at middle and high frequencies. The reason for this can be considered as follows: in general a speech sound energy

Fig. 2. Frequency properties of the gain ratio at $l_1 = 10$ cmFig. 3. Frequency properties of the gain ratio at $l_1 = 30$ cm

is large at lower frequencies but small at high frequencies, and its energy decreases with its propagation proceeding. Thus measurement of gain ratio becomes susceptible to ambient noises at high frequencies. The number of exceptions is larger in Fig.3 than in Fig.2. This can be considered to be because the magnitude of gain ratio approaches to unity with the source distance l_1 becoming longer and is susceptible to ambient noises; in addition to this, the influence of reverberation should be taken into consideration when the source distance l_1 becoming longer.

The exceptions to $\angle r_x(\omega) > 0$ are found at low frequencies in Figs. 4 and 5, respectively, the plot of the phase difference $\angle r_x(\omega)$ at $l_1=10$ and 30 cm. The reason for this can be considered as follows: The phase difference is calculated by dividing the time difference of sound arrival to the two microphones by the wave length. Thus $\angle r_x(\omega)$ for the same time difference takes a smaller value in low frequencies than in high frequencies. Thus the phase difference becomes indiscernible in lower frequencies and is susceptible to ambient noises.

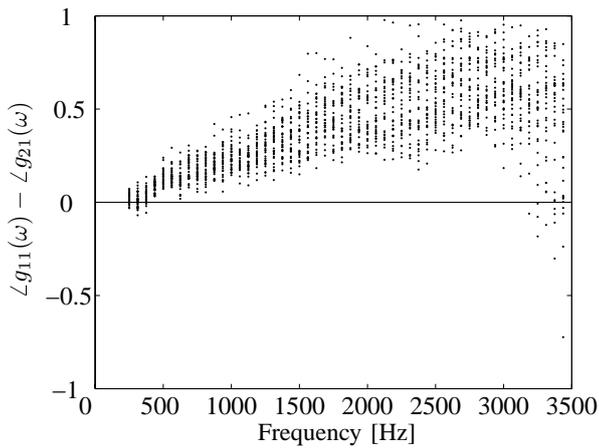


Fig. 4. Frequency properties of the phase difference at $l_1 = 10$ cm

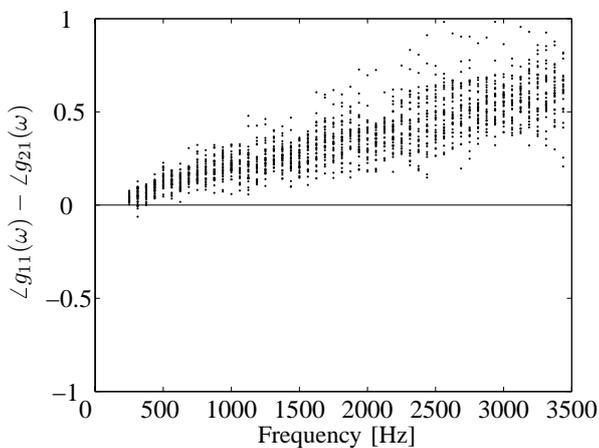


Fig. 5. Frequency properties of the phase difference at $l_1 = 30$ cm

The exceptions at high frequencies over about 3200 Hz in Fig.4 are ascribed to the spatial aliasing. Under the assumption that the maximum frequency is 3400 Hz and the sound velocity is 340 m/s, the aliasing is derived to occur if the difference between the two distances from the sound source to two microphones is longer than 5 cm. The difference is 5.32 cm when $l_1 = 10$ cm and 3.23 cm when $l_1 = 30$ cm. Thus in the placement of Fig.1, the aliasing occurs at frequency over 3195.5 Hz when $l_1 = 10$ cm while never occurs when $l_1 = 30$ cm.

From these results, it is concluded that the gain inequality of Eq.(10) holds rather in low frequencies while the phase inequality Eq.(11) holds rather in high frequencies unless the spacial aliasing occurs.

B. Properties of gain-based and phase-based correction methods

In order to evaluate the frequency characteristics of the the gain-based and phase-based methods, several experiments

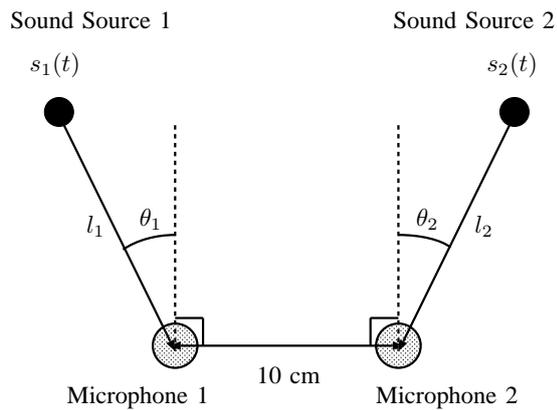


Fig. 6. Placement of two sound sources and two microphones

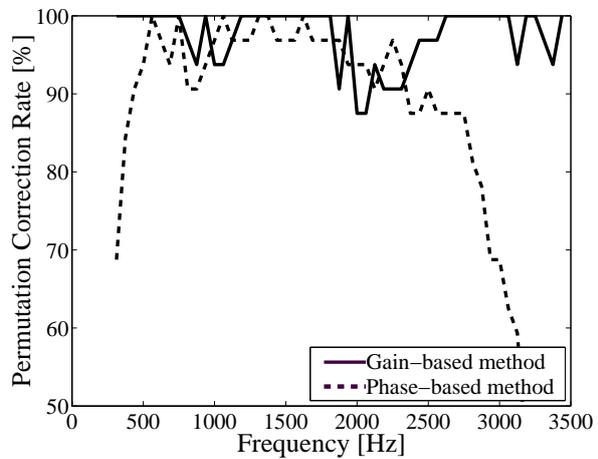
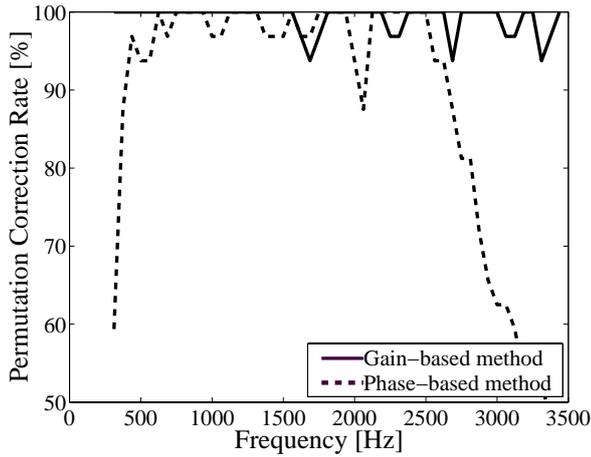
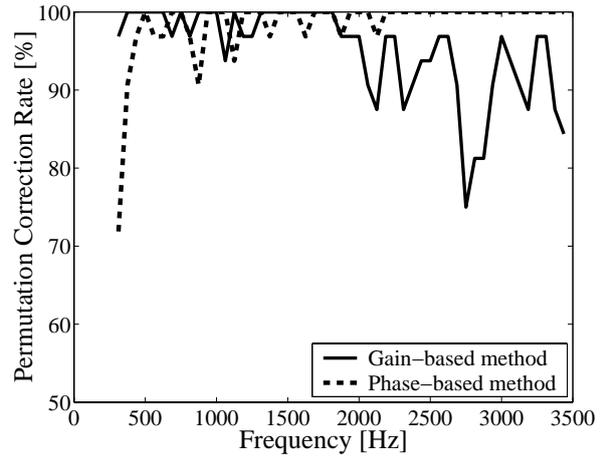
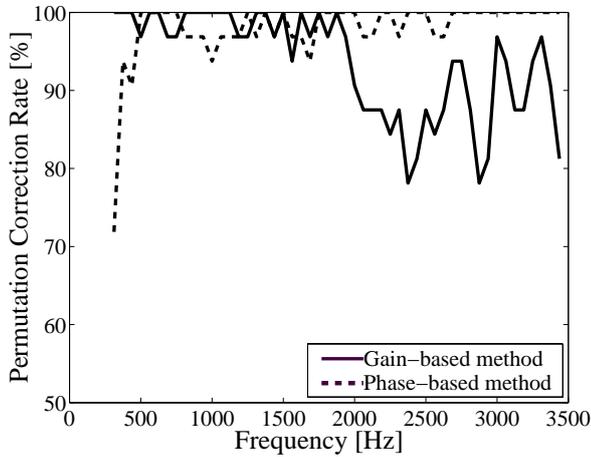


Fig. 7. Correction rates (l_1, l_2) = (10, 30)

were carried out. As shown in Fig.6, source 1 distance l_1 from microphone 1 was set at l_1 cm with $\theta_1 = 10^\circ$ and source 2 distance l_1 from microphone 2 was set at l_2 cm with $\theta_2 = 10^\circ$. The other conditions are the same as in the above.

Over a loudspeaker located at source 2, we play a roaring train noise recorded at a station premises [9]; departure and arrival announcements and passengers' conversations were also included in the noise. The noise level was in average 82.1 dB at a distance of $l_2=30$ cm and 76.3 db at a distance of $l_2=60$ cm from the loudspeaker. Under these conditions, 4 words were uttered by 4 male and 4 female speakers at the position of (l_1 cm, l_2 cm)=(10,30), (10,60), (30,60) and (30,100). In total, 128 data sets were obtained.

The experimental results are shown in Figs. 7-10. In the figures, the horizontal axis denotes the frequency and the vertical axis the permutation correction rates. A comparison of the results in Figs. 7 and 8 show that both the gain-based and the phase-based methods give higher correction rates with the noise source distance l_2 to the microphones becoming longer. This fact is almost true in Figs. 9 and 10.

Fig. 8. Correction rates $(l_1, l_2) = (10, 60)$ Fig. 10. Correction rates $(l_1, l_2) = (30, 100)$ Fig. 9. Correction rates $(l_1, l_2) = (30, 60)$

From all the figures it is found that the correction rates by the gain-based method is higher in low frequencies than in high frequencies. This can be considered to be because a speech sound concentrates its energy in low frequencies and thus because the sound components in low frequencies are not susceptible to ambient noise as compared with those in high frequencies.

In low frequencies, the phase-based method does not give poor results compared with the gain-based one and in particular does not function normally at frequencies less than 500 Hz. The reason for this can be considered as follows. The phase difference is calculated by dividing the time difference of sound arrival to the two microphones by the wave length. Even if the time difference is the same, the phase difference is evaluated smaller in low frequencies than in high frequencies. Thus the phase difference is more susceptible to ambient noise in lower frequencies than in high frequencies.

From a comparison of Figs. 7 and 8 to Figs. 9 and 10, it is found that at frequencies over 1800 Hz the gain-based method

is more robust at $l_1=10$ cm than at $l_1=30$ cm. This can be ascribed to a decrease in sound energy with the propagation distance. As seen in Figs. 7 and 8, the correction rates due to the phase-based method drop suddenly at high frequencies over 2500 Hz in the case of $l_1=10$ cm. This is ascribed to the special aliasing. In the case of $l_1=30$ cm, the aliasing never occurs. So the phase-based method gives almost perfect results as shown in Figs. 9 and 10.

C. An Integrated Permutation Correction Method

From the above facts, it is concluded that the gain-based and the phase-based methods should be used according to frequency bands; because their frequency characteristics depend on the microphone-space and the source distance from the microphones. In this section, we propose a new permutation correction method based on the frequency characteristics.

In the case where the spacial aliasing doesn't occur, the gain-based method as Eq.(13) is functional on the low frequency bands, the phase-based method as Eq.(14) works well on the high frequency bands. Therefore, an integrated permutation correction method is given by

$$y_i(\omega, k) = \begin{cases} v_{i=i, j=i}(\omega, k) & \text{if } |r_1(\omega)| > |r_2(\omega)|, f < f_0 \\ v_{i \neq i, j=i}(\omega, k) & \text{if } |r_1(\omega)| < |r_2(\omega)|, f < f_0 \\ v_{i=i, j=i}(\omega, k) & \text{if } \angle r_1(\omega) > \angle r_2(\omega), f_0 \leq f \\ v_{i \neq i, j=i}(\omega, k) & \text{if } \angle r_1(\omega) < \angle r_2(\omega), f_0 \leq f \end{cases} \quad (16)$$

where f_0 denotes the boundary frequency. In this experimental condition, f_0 only has to set the frequency from about 1000 Hz to 1500 Hz.

D. Experimental Results in Real Room

In order to verify our proposals, several experiments were carried out in same laboratory room as **III-B**. f_0 was set to 1800 Hz. Table IV shows the permutation correction rates and estimated SN ratio [10]. The power-based method is corrected about 90% but the correction rate depends on the distance from the sources to the microphones. The proposed method

TABLE IV
PERMUTATION CORRECTION RATES AND SNR AGAINST STATION PREMISE
NOISES ROARED FROM A LOUD SPEAKER %

(l_{11}, l_{22})	Correction rate %		SNR dB	
	(30,60)	(30,100)	(30,60)	(30,100)
FastICA	88.61	87.62	6.31	7.11
Similarity	93.87	95.20	15.17	15.45
Power	88.60	90.75	16.42	16.88
Rule 1	94.04	95.51	17.30	17.66
Rule 2	98.13	98.58	17.78	17.84
Proposed	99.98	99.99	18.20	18.29

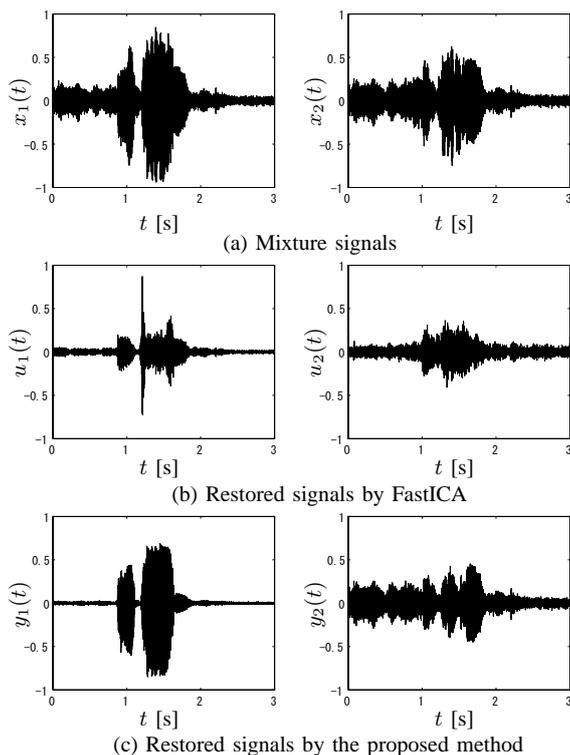


Fig. 11. Experimental results when a female speaker uttered ($l_{11}=10$ cm) under station premise noises from a loud speaker ($l_{22}=30$ cm): (a) Mixture signals recorded by microphone 1 and 2 respectively. (b) Restored signals by FastICA. (c) Restored signals by the proposed method.

can solve almost perfectly in the all conditions. Furthermore, the permutation in the time domain does not occur using by the proposed method. In other words, the separated signal $y_1(t)$ is estimated the source $s_1(t)$ and $y_2(t)$ is estimated $s_2(t)$. It is clarified that the proposed method can not only the permutation correction but also the target extraction.

Fig.11 shows the source estimates $y_i(t)(i = 1, 2)$ recovered in the time domain. From the waveforms, it is found that there exist less cross talk in Fig.11(c) by the proposal while in Fig.11(b) by FastICA there exist much cross talk because of the unresolved permutation described above. It is also found that the waveforms in Fig.11(c) are properly scaled while those in Fig.11(b) not. These facts were also confirmed by articulation tests.

IV. CONCLUSIONS

From the experiments in a real environment, it is found that the gain and the phase of the acoustic transfer function are dependent on three factors, namely, the frequency bands, the microphone-space and the source-microphone distance. In a real application of blind source separation, the distances from sources to microphones are unknown. For better performance, therefore, we should design the permutation correction, especially taking the frequency bands and the microphone-space into consideration. This paper propose a new permutation correction method for frequency domain ICA based on properties of acoustic transfer functions. The proposed method has been verified from several experiments in a real room.

ACKNOWLEDGMENT

This research has been supported by the Kayamori Foundation of Informational Science Advancement and the Tateishi Science and Technology Foundation.

REFERENCES

- [1] A. J. Bell and T. J. Sejnowski: An information maximization approach to blind separation and blind deconvolution, *Neural Computation*, Vol. 7, No. 6, pp. 1129-1159, 1995.
- [2] A. Cichocki and S. Amari: Adaptive blind signal and image processing, Learning algorithm and applications, *John Wiley & Sons, Ltd*, 2002.
- [3] T. W. Lee, M. Girolami and T. J. Sejnowski: Independent component analysis using an extended informax algorithm for mixed subgaussian and supergaussian sources, *Neural Computation*, Vol. 11, No. 2, pp. 417-441, 1999.
- [4] A. Hyvärinen, J. Karhunen and E. Oja, "Independent component analysis," John Wiley & Sons, Ltd, 2001.
- [5] N. Murata, S. Ikeda and A. Ziehe: An approach to blind source separation based on temporal structure of speech signals, *Neurocomputing*, Vol. 41, Issue 1-4, pp. 1-24, 2001.
- [6] H. Gotanda, K. Nobu, T. Koya, K. Kaneda, T. Ishibashi and N. Haratani: Permutation correction and speech extraction based on split spectrum through FastICA, *4th International Symposium on Independent Component Analysis and Blind Signal Separation (ICA2003)*, pp. 379-384, 2003.
- [7] H. Saruwatari, S. Kurita, K. Takeda, F. Itakura, T. Nishikawa and K. Shikano: Blind source separation combining independent component analysis and beamforming, *EURASIP Journal on Applied Signal Processing*, Vol.2003, No.11, pp.1135-1146, 2003.
- [8] Acoustical Society of Japan: ASJ continuous speech corpus japanese newspaper article sentences, *JNAS Vols.1-16*, 1997.
- [9] NTT Advanced Technology Corporation: Ambient noise database for telephony 1996, 1996.
- [10] T. Koya, N. Iwasaki, T. Ishibashi, G. Hirano, H. Shiratsuchi and Hiromu Gotanda, "SN ratio estimation and speech segment detection of extracted signals through Independent Component Analysis," *Journal of Advanced Computational Intelligence and Intelligent Informatics*, Vol. 14, No. 4, pp.364-374, 2010.



Takaaki Ishibashi received the Ph.D. degree in Engineering from Kyushu Institute of Technology, in 2007.

He is Associate Professor of the Department of Information, Communication and Electronic Engineering, Kumamoto National College of Technology. And his research area of interest is Digital Signal Processing, and Human Interface.

Dr. Ishibashi is a member of the Institute of Electronics, Information and Communication Engineers (IEICE), the Institute of Systems, Control and Information Engineers (ISCIE), The Society of Instrument and Control Engineers (SICE) and the Astronomical Society of Japan (ASJ).