# BLIND SOURCE SEPARATION OF MANY SIGNALS IN THE FREQUENCY DOMAIN

*Ryo Mukai      Hiroshi Sawada      Shoko Araki      Shoji Makino*

NTT Communication Science Laboratories, NTT Corporation
2–4 Hikaridai, Seika-cho, Soraku-gun, Kyoto 619–0237, Japan

{ryo,sawada,shoko,maki}@cslab.kecl.ntt.co.jp

## ABSTRACT

This paper describes the frequency-domain blind source separation (BSS) of convolutively mixed acoustic signals using independent component analysis (ICA). The most critical issue related to frequency domain BSS is the permutation problem. This paper presents two methods for solving this problem. Both methods are based on the clustering of information derived from a separation matrix obtained by ICA. The first method is based on direction of arrival (DOA) clustering. This approach is intuitive and easy to understand. The second method is based on normalized basis vector clustering. This method is less intuitive than the DOA based method, but it has several advantages. First, it does not need sensor array geometry information. Secondly, it can fully utilize the information contained in the separation matrix, since the clustering is performed in high-dimensional space. Experimental results show that our methods realize BSS in various situations such as the separation of many speech signals located in a 3-dimensional space, and the extraction of primary sound sources surrounded by many background interferences.

## 1. INTRODUCTION

Blind source separation (BSS) [1] is a technique for estimating individual source signals from their mixtures observed by sensors. The BSS of audio signals has a wide range of applications including speech enhancement. Independent component analysis (ICA) [2] is one of the main statistical methods used for BSS and it achieves separation by using the non-Gaussianity and independence of source signals. In most realistic applications, the number of source signals is large, and the signals are mixed in a convolutive manner with reverberations. This makes the problem difficult.

There are two major approaches to solving the convolutive BSS problem. The first is the time domain approach, where ICA is applied directly to the convolutive mixture model [3, 4, 5]. This approach achieves good separation if the calculation successfully converges to a correct solution, however, it incurs considerable computational cost. Thus it is difficult to obtain a solution in a practical time especially when the number of source signals is large.

The other approach is frequency domain BSS, where ICA is applied to multiple instantaneous mixtures in the frequency domain [6, 7, 8, 9]. This approach takes much less computation time than time domain BSS. However, it poses another problem in that we need to align the output signal order for every frequency bin so that a separated signal in the time domain contains frequency components from one source signal. This problem is known as the permutation problem. We have been studying frequency domain BSS [10, 11, 12, 13, 14, 15], and have developed effective methods for solving the permutation problem. Our methods realize BSS in various situations such as the separation of many speech signals located in a 3-dimensional space, and the extraction of primary sound sources surrounded by many background interferences. In



**Fig. 1**. Flow of frequency domain BSS

this paper, we describe key techniques and demonstrate the effectiveness of our methods experimentally.

## 2. FREQUENCY DOMAIN BSS

When $N$ source signals are $\mathrm{s}_1(t), ..., \mathrm{s}_N(t)$ and the signals observed by $M$ sensors are $\mathrm{x}_1(t), ..., \mathrm{x}_M(t)$, the mixing model can be described by the following equation:

$$\mathrm{x}_j(t) = \sum_{i=1}^{N} \sum_l \mathrm{h}_{ji}(l)\mathrm{s}_i(t-l), \qquad (1)$$

where $\mathrm{h}_{ji}(l)$ is the impulse response from source $i$ to sensor $j$.

Figure 1 shows the flow of frequency domain BSS. First, time-domain observed signals $\mathrm{x}_j(t)$ sampled at frequency $f_s$ are converted into frequency-domain time-series signals $x_j(f, \tau)$ with an $L$-point short-time Fourier transform (STFT):

$$x_j(f, \tau) = \sum_{r=-L/2}^{L/2-1} \mathrm{x}_j(\tau+r)\,\mathrm{win}(r)\,e^{-j2\pi fr}, \qquad (2)$$

where $f \in \{0, \frac{1}{L}f_s, \ldots, \frac{L-1}{L}f_s\}$ is a frequency, $\mathrm{win}(r)$ is a window that tapers smoothly to zero at each end, and $\tau$ is a new index representing time. The convolutive mixtures (1) can be approximated as instantaneous mixtures at each frequency:

$$x_j(f, \tau) \approx \sum_{k=1}^{N} h_{jk}(f)s_k(f, \tau), \qquad (3)$$

where $h_{jk}(f)$ is the frequency response from source $k$ to sensor $j$, and $s_k(f, \tau)$ is a frequency-domain time-series signal of $\mathrm{s}_k(t)$ obtained by an operation similar to (2). A vector notation of (3) is given as:

$$\mathbf{x}(f, \tau) = \sum_{k=1}^{N} \mathbf{h}_k(f)s_k(f, \tau), \qquad (4)$$

where $\mathbf{x} = [x_1, \ldots, x_M]^T$ is an observation vector and $\mathbf{h}_k = [h_{1k}, \ldots, h_{Mk}]^T$ is the vector of the frequency responses from source $s_k$ to all sensors. The first step of the BSS is to obtain frequency components $\mathbf{h}_k(f)s_k(f, \tau)$ for each source signal $k$ from the observation vector $\mathbf{x}(f, \tau)$. To extract the components, we apply ICA to the observation vectors, and we have:

$$\mathbf{y}(f, \tau) = \mathbf{W}(f)\,\mathbf{x}(f, \tau), \qquad (5)$$

where $\mathbf{W}(f)$ is an $N \times M$ separation matrix and $\mathbf{y}(f, \tau) = [y_1(f, \tau), \ldots, y_M(f, \tau)]^T$ is a vector of independent components. When the number of source signals $N$ is known and $N < M$, we can apply Principle Component Analysis (PCA) to the observation vector to reduce its dimensions, otherwise we assume $N = M$ and $\mathbf{W}(f)$ is a square matrix. The ICA algorithm for complex-valued

**Fig. 2**. DOA (source location) clustering (Secs.3.2 and 4.1)



**Fig. 3**. Normalized basis vector clustering (Secs.3.3 and 4.2)



**Fig. 4**. Far-field model

signals is detailed in [10].

ICA maximizes the non-Gaussianity of the output signals $y_i(f,\tau)$, therefore, when the source signals are non-Gaussian and mutually independent, the separation is achieved in each frequency bin. However, the ICA solution suffers permutation and scaling ambiguities. Before constructing output signals in the time domain, we have to align the permutation so that each channel contains frequency components from one source signal. Section 3 details methods for solving the permutation problem. With regard to the scaling problem, there is a simple and reasonable solution that uses one element of the basis vector obtained by (7), which is given in the next section:

$$y_i(f,\tau) \leftarrow a_{Ji}(f)y_i(f,\tau), \qquad (6)$$

where $J$ is a reference sensor. This solution is equivalent to the minimal distortion principle (MDP) [3] or the projection back method [7]. By using this solution, the output signal $\mathrm{y}_i$ becomes an estimation of the reverberant version of source $\mathrm{s}_i$ measured at sensor $J$. After the operations for solving the permutation and scaling problems, time-domain output signals $\mathrm{y}_i(t)$ are obtained by using an inverse STFT (ISTFT):

$$\mathrm{y}_i(\tau+r) = \frac{1}{L \cdot \mathrm{win}(r)} \sum_{f \in \{0, \frac{1}{L}f_s, \ldots, \frac{L-1}{L}f_s\}} y_i(f,\tau)\, e^{j2\pi fr}.$$

## 3. SOLVING PERMUTATION PROBLEM

### 3.1. Basis vector

The inverse (or pseudo inverse when $N < M$) of the separation matrix $\mathbf{W}$ provides useful information for solving the permutation problem. For simplicity, we assume $N = M$ in the following discussions. Let $\mathbf{a}_1, ..., \mathbf{a}_M$ be the column vectors of $\mathbf{W}^{-1}$:

$$[\mathbf{a}_1, \cdots, \mathbf{a}_M] \triangleq \mathbf{W}^{-1}, \ \mathbf{a}_i = [a_{1i}, \ldots, a_{Mi}]^T. \qquad (7)$$

We call $\mathbf{a}_1, ..., \mathbf{a}_M$ basis vectors, because the observation vector $\mathbf{x}(f,\tau)$ is represented by a linear combination of these vectors:

$$\mathbf{x}(f,\tau) = \sum_{i=1}^{M} \mathbf{a}_i(f)y_i(f,\tau), \qquad (8)$$

which is given by multiplying both sides of (5) by $\mathbf{W}^{-1}$. This equation is very important for frequency domain BSS. If a separation matrix $\mathbf{W}(f)$ is successfully calculated by ICA, there exist $i$ and $k$ such that $\mathbf{a}_i(f)y_i(f,\tau)$ corresponds to $\mathbf{h}_k(f)s_k(f,\tau)$ in (4). Determining the correspondences between $i$ and $k$ for all $f$ is equivalent to solving the permutation problem.

The following subsections describe two approaches for solving the permutation problem. Both methods utilize information derived from the basis vectors. The first method is based on the clustering of estimated directions of arrival (DOA) or source locations (Fig. 2). This approach is intuitive and easy to understand, however it needs the assumptions that the number of source signals $N$ is known and $N \le M$. The second method is based on the clustering of normalized basis vectors (Fig. 3). This method is

designed to work even when $N$ is unknown and $N > M$. In such a case, we extract $M$ primary signals instead of separating all $N$ signals.

### 3.2. Clustering estimated DOA

With a far-field model, a frequency response from source $k$ to sensor $j$ can be approximated as:

$$h_{jk}(f) \approx e^{j2\pi fc^{-1}\mathbf{p}_j^T\mathbf{q}_k}, \qquad (9)$$

where $c$ is the wave propagation speed, $\mathbf{p}_j$ is the location of sensor $j$, and $\mathbf{q}_k$ represents a unit vector that points to the direction of source $k$. By taking the ratio for a sensor pair $j$ and $j'$ with this model, we have

$$h_{jk}(f)/h_{j'k}(f) \approx e^{j2\pi fc^{-1}(\mathbf{p}_j-\mathbf{p}_{j'})^T\mathbf{q}_k} \qquad (10)$$

$$= e^{j2\pi fc^{-1}\|\mathbf{p}_j-\mathbf{p}_{j'}\|\cos\theta_k^{jj'}}, \qquad (11)$$

where $\theta_k^{jj'}$ is the direction of source $k$ relative to the sensor pair $j$ and $j'$. In this way, DOA has two kinds of representations; the *absolute* DOA $\mathbf{q}_k$, which is determined in a coordinate system and the *relative* DOA $\theta_k^{jj'}$ which is determined relative to a microphone axis. When we adopt a near-field model, that includes the attenuation of the wave, we can estimate range information in addition to DOA. The details are given in [12, 13].

When the ICA solution is successfully calculated, and assuming correspondences between terms in (4) and (8), the ratio of elements $a_{ji}(f)$ and $a_{j'i}(f)$ in a basis vector $\mathbf{a}_i(f)$ can be expressed as follows:

$$\frac{a_{ji}(f)}{a_{j'i}(f)} = \frac{a_{ji}y_i}{a_{j'i}y_i} \approx \frac{h_{jk}s_k}{h_{j'k}s_k} = \frac{h_{jk}(f)}{h_{j'k}(f)}. \qquad (12)$$

Here, indexes $i$ and $k$ may be different. This represents the permutation ambiguity. By using the arguments of (12) and (11), we can estimate a relative DOA:

$$\hat{\theta}_i^{jj'}(f) = \arccos \frac{\arg[a_{ji}(f)/a_{j'i}(f)]}{2\pi fc^{-1}\|\mathbf{p}_j - \mathbf{p}_{j'}\|}. \qquad (13)$$

Absolute DOA $\mathbf{q}_k$ is estimated by using multiple sensor pairs. By using the arguments of (12) and (10), we have:

$$2\pi fc^{-1}(\mathbf{p}_j-\mathbf{p}_{j'})^T\mathbf{q}_k \approx \arg[a_{ji}(f)/a_{j'i}(f)]. \qquad (14)$$

When we consider (14) for $u$ sensor pairs $(j_1, j_1'), \ldots, (j_u, j_u')$, we have a simultaneous equation

$$2\pi fc^{-1}\mathbf{V}\,\mathbf{q}_k = \mathbf{r}_i(f), \qquad (15)$$

where

$$\mathbf{V} \triangleq [\mathbf{p}_{j_1}-\mathbf{p}_{j_1'}, \ldots, \mathbf{p}_{j_u}-\mathbf{p}_{j_u'}]^T,$$

$$\mathbf{r}_i(f) \triangleq [\arg(a_{j_1 i}/a_{j_1' i}), \ldots, \arg(a_{j_u i}/a_{j_u' i})]^T.$$

In a practical situation, (15) seldom has the exact solution because of an estimation error in $\mathbf{r}_i$. Therefore, we use the Moore-Penrose pseudo-inverse $\mathbf{V}^+ \triangleq (\mathbf{V}^T\mathbf{V})^{-1}\mathbf{V}^T$, and we obtain an approximately optimal solution:

$$\hat{\mathbf{q}}_i(f) = \frac{\mathbf{V}^+\mathbf{r}_i(f)}{2\pi fc^{-1}}, \ \hat{\mathbf{q}}_i(f) \leftarrow \frac{\hat{\mathbf{q}}_i(f)}{\|\hat{\mathbf{q}}_i(f)\|}. \qquad (16)$$

We can group separated frequency components $y_i(f,\tau)$ according

to the clustering result of their estimated DOAs. Section 4.1 shows experimental results obtained using this method.

### 3.3. Clustering normalized basis vector

This section describes a method for solving the permutation problem by clustering normalized basis vectors $\bar{\mathbf{a}}_i(f)$, which are calculated by eliminating frequency dependency from basis vectors $\mathbf{a}_i(f)$. This method is less intuitive than the DOA based method described above, but it has several advantages. First, it does not need sensor array geometry information $\mathbf{p}_1, ..., \mathbf{p}_M$. Secondly, it can fully utilize the information contained in the basis vectors, since clustering is performed in $M$-dimensional complex-valued space $\mathbf{C}^M$, while the DOA based method uses the information reduced onto a sphere $\mathbf{S}^2$ in 3-dimensional space.

Elements of the normalized basis vector $\bar{\mathbf{a}}_i(f) = [\bar{a}_{1i}(f), \dots, \bar{a}_{Mi}(f)]^T$ are calculated by the following formula:

$$\bar{a}_{ji}(f) = |a_{ji}(f)| \exp\left[ \jmath \frac{\arg[a_{ji}(f)/a_{Ji}(f)]}{4fc^{-1}d_{\max}} \right], \quad (17)$$

where $J$ is the index of a reference sensor, and $d_{\max}$ is the maximum distance between the sensor $J$ and a sensor $\forall j \in \{1, \dots, M\}$. This equation eliminates the frequency dependency in $\mathbf{a}_i(f)$. Then, we normalize the vector length to 1 to eliminate the scaling ambiguity.

$$\bar{\mathbf{a}}_i(f) \leftarrow \bar{\mathbf{a}}_i(f) / \|\bar{\mathbf{a}}_i(f)\| \quad (18)$$

According to these operations, the normalized basis vector $\bar{\mathbf{a}}_i(f)$ becomes independent of the frequency $f$, but it depends on the source direction $\mathbf{q}_k$ and the (unknown) sensor locations $\mathbf{p}_1, \dots, \mathbf{p}_M$. Actually, we can confirm this by using a far-field model (9) and equations (14), (17) and (18)

$$\bar{a}_{ji}(f) \approx \frac{1}{\sqrt{M}} \exp\left[ \jmath \frac{\pi}{2} \frac{(\mathbf{p}_j - \mathbf{p}_J)^T \mathbf{q}_k}{d_{\max}} \right].$$

When we use a near-field model, we can also prove that $\bar{a}_{ji}(f)$ is independent of the frequency $f$ [14].

After normalizing all the basis vectors, we employ a clustering algorithm to find clusters $C_1, \dots, C_M$ formed by normalized vectors $\bar{\mathbf{a}}_i(f)$. The centroid $\mathbf{c}_k$ of a cluster $C_k$ is calculated by

$$\mathbf{c}_k \leftarrow \sum_{\bar{\mathbf{a}} \in C_k} \bar{\mathbf{a}} / |C_k|, \quad \mathbf{c}_k \leftarrow \mathbf{c}_k / \|\mathbf{c}_k\|, \quad (19)$$

where $|C_k|$ is the number of vectors in $C_k$. The clustering criterion is to minimize the total sum $\mathcal{J}$ of the squared distances between cluster members and their centroid

$$\mathcal{J} = \sum_{k=1}^M \mathcal{J}_k, \quad \mathcal{J}_k = \sum_{\bar{\mathbf{a}} \in C_k} \|\bar{\mathbf{a}} - \mathbf{c}_k\|^2. \quad (20)$$

This minimization can be achieved by using an ordinary clustering method such as the $k$-means algorithm [16].

Then, to align the permutation ambiguities, we renumber the indexes of the basis vectors by

$$\mathbf{a}_k(f) \leftarrow \mathbf{a}_{\Pi_f(k)}(f), \quad (21)$$

where $\Pi_f : \{1, \dots, M\} \rightarrow \{1, \dots, M\}$ is a one-to-one mapping decided for each frequency $f$ by

$$\Pi_f = \operatorname{argmin}_\Pi \sum_{k=1}^M \|\bar{\mathbf{a}}_{\Pi(k)}(f) - \mathbf{c}_k\|^2. \quad (22)$$

We also renumber independent components $y_1(f, \tau), \dots, y_M(f, \tau)$ accordingly.

Section 4.2 shows the experimental BSS results obtained using this method with multiple target signals surrounded by many background interferences.

## 4. EXPERIMENTS

### 4.1. Separation of six 3-D located sources

We carried out experiments in an ordinary office and evaluated the Signal to Interference Ratio (SIR) performance. We used eight microphones located at the vertex of a 4 cm cube and six signals distributed in a three-dimensional space (Fig. 5). We calculated the



**Fig. 5**. Six 3-D located source signals and eight microphones located at the vertex of a 4 cm cube



**Fig. 6**. Clustering result of estimated DOAs

**Table 1**. Experimental results (dB)

|            | SIR$_1$ | SIR$_2$ | SIR$_3$ | SIR$_4$ | SIR$_5$ | SIR$_6$ | ave. |
|------------|---------|---------|---------|---------|---------|---------|------|
| Input SIR  | −11.6   | −9.0    | −9.0    | −6.6    | −6.9    | −2.5    | −7.6 |
| Output SIR | 7.6     | 12.2    | 16.4    | 14.4    | 13.6    | 13.7    | 13.0 |
| improvement | 19.2   | 21.2    | 25.4    | 21.0    | 20.5    | 16.2    | 20.6 |

separation filter by using live recorded mixtures, and evaluated the SIRs by using individually activated source signals. The separation system $\mathbf{W}$ was calculated by using a complex-valued version of FastICA [17] and further improved by using InfoMax [18] combined with the natural gradient whose nonlinear function is based on the polar coordinate [10]. We applied the $k$-means algorithm to the estimated DOAs obtained by (16). Figure 6 shows a clustering result. As reverberations disrupt the presupposed far-field model, the estimated DOAs are scattering around the centroids. Nevertheless, this information is sufficient for solving the permutation problem. The results of SIR evaluation are shown in Table 1. We obtained good separation in spite of the very low input SIR. The average SIR improvement was more than 20 dB.

### 4.2. Extraction of primary sources in ambient noise

Next we designed experiments to evaluate the effectiveness of the basis vector clustering. We measured impulse responses $h_{jk}(l)$ under the conditions shown in Fig. 7. We used four 3-dimensionally arranged microphones ($M = 4$), and nine loudspeakers. Three of the loudspeakers were located near the microphones and used as primary target signals to be separated. The remaining six loudspeakers were located far from the microphones to provide ambient interferences. These speakers simulate a cocktail party situation. Mixtures were made at the microphones by convolving the impulse responses and 6-second English and Japanese speeches sampled at 8 kHz. The system knew only the maximum distance (4 cm) between the reference microphone (Mic. 1) and the others.

As the number of microphones was $M = 4$, we calculated a

**Fig. 7**. Experimental conditions



Squared distances from corresponding centroid

**Fig. 8**. Result of normalized basis vector clustering (1 target)



Squared distances from corresponding centroid

**Fig. 9**. Result of normalized basis vector clustering (3 targets)

$4 \times 4$ separation matrix $\mathbf{W}(f)$ for each frequency bin. The permutation was aligned according to the information obtained by the normalized basis vector clustering described in Sec. 3.3. Figures 8 and 9 show the clustering results. There are four clusters corresponding to four output channels. The normalized basis vectors are $M$-dimensional complex-valued, and so it is difficult to visualize the clustering results. These figures show the squared distance between the cluster members and each corresponding centroid.

The variances of the clusters $\mathcal{J}_k/|C_k|$ can be used to distinguish primary target signals and background noises. We can infer that clusters with small variances correspond to target signals near the microphones. Figure 8 shows the results we obtained when we activated only one of the three target speakers. We can see that one cluster has a small variance and it corresponds to a primary source signal. Figure 9 shows the results when all three target speakers were activated. There are three clusters with small variances that correspond to primary sources. In this way, we can estimate the number of primary sounds.

We evaluated the separation performance with three targets and six background interferences. Experiments were conducted with ten combinations of nine speeches. Table 2 shows the average SIR. Even under such difficult conditions, our system succeeded in enhancing and separating the target sources. The separation performance can be improved further by employing post-processing using time-frequency masking [14].

**Table 2**. Average SIR (dB)

| Target position | a120 | b120 | c170 |
|---|---|---|---|
| InputSIR$_i$ | −3.9 | −3.6 | −5.9 |
| OutputSIR$_i$ | 8.6 | 10.0 | 8.6 |
| improvement | 12.5 | 13.6 | 14.5 |

## 5. CONCLUSION

We have described frequency domain BSS for a large number of source signals and demonstrated its effectiveness in two different situations. The key technique is a method for solving the permutation problem. We can align the permutation efficiently according to the clustering results of estimated DOAs or normalized basis vectors. The frequency normalization method described in this paper can be used for underdetermined ($N > M$) BSS, which does not use ICA [15].

## 6. REFERENCES

[1] S. Haykin, Ed., *Unsupervised Adaptive Filtering*, John Wiley & Sons, 2000.

[2] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent Component Analysis*, John Wiley & Sons, 2001.

[3] K. Matsuoka and S. Nakashima, "Minimal distortion principle for blind source separation," in *Proc. ICA 2001*, 2001, pp. 722–727.

[4] S. C. Douglas and X. Sun, "Convolutive blind separation of speech mixtures using the natural gradient," *Speech Communication*, vol. 39, pp. 65–78, 2003.

[5] S. C. Douglas, H. Sawada, and S. Makino, "A spatio-temporal FastICA algorithm for separating convolutive mixtures," in *Proc. ICASSP 2005*, 2005, vol. V, pp. 165–168.

[6] P. Smaragdis, "Blind separation of convolved mixtures in the frequency domain," *Neurocomputing*, vol. 22, pp. 21–34, 1998.

[7] N. Murata, S. Ikeda, and A. Ziehe, "An approach to blind source separation based on temporal structure of speech signals," *Neurocomputing*, vol. 41, no. 1-4, pp. 1–24, 2001.

[8] F. Asano, S. Ikeda, M. Ogawa, H. Asoh, and N. Kitawaki, "Combined approach of array processing and independent component analysis for blind separation of acoustic signals," *IEEE Trans. Speech Audio Processing*, vol. 11, no. 3, pp. 204–215, May 2003.

[9] 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.

[10] H. Sawada, R. Mukai, S. Araki, and S. Makino, "Polar coordinate based nonlinear function for frequency-domain blind source separation," *IEICE Trans. Fundamentals*, vol. E86-A, no. 3, pp. 590–596, Mar. 2003.

[11] H. Sawada, R. Mukai, S. Araki, and S. Makino, "A robust and precise method for solving the permutation problem of frequency-domain blind source separation," *IEEE Trans. Speech Audio Processing*, vol. 12, no. 5, pp. 530–538, 2004.

[12] R. Mukai, H. Sawada, S. Araki, and S. Makino, "Frequency domain blind source separation for many speech signals," in *Proc. ICA2004 (LNCS 3195)*, pp. 461–469. Springer-Verlag, 2004.

[13] H. Sawada, R. Mukai, S. Araki, and S. Makino, "Frequency-domain blind source separation," in *Speech Enhancement*, J. Benesty, S. Makino, and J. Chen, Eds., chapter 13, pp. 299–327. Springer, 2005.

[14] H. Sawada, S. Araki, R. Mukai, and S. Makino, "Blind extraction of a dominant source from mixtures of many sources using ICA and time-frequency masking," in *Proc. ISCAS 2005*, 2005, pp. 5882–5885.

[15] S. Araki, H. Sawada, R. Mukai, and S. Makino, "A novel blind source separation method with observation vector clustering," in *Proc. IWAENC 2005*, 2005, pp. 117–120.

[16] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*, Wiley Interscience, 2nd edition, 2000.

[17] E. Bingham and A. Hyvärinen, "A fast fixed-point algorithm for independent component analysis of complex valued signals," *International Journal of Neural Systems*, vol. 10, no. 1, pp. 1–8, Feb. 2000.

[18] 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.