

Geometrically Constrained Blind Moving Source Extraction based on Constant Separation Vector and Auxiliary Function Technique

Ruifeng Zhang, Tetsuya Ueda and Shoji Makino
Waseda University

2-7 Hibikino, Wakamatsu-ku, Kitakyushu, Fukuoka 808-0135, Japan

Abstract—This paper focuses on the permutation problem in Blind Source Extraction (BSE) by employing beamforming-based Geometrical Constraints (GC). Specifically, we focus on the recently proposed auxiliary function-based Independent Vector Extraction (Aux-IVE) with Constant Separation Vector (CSV) mixing model. Building upon this mixing model, we provide the source signals' spatial information, utilizing GC based on beamforming. We facilitate the extraction of moving target source and improve the extraction performance. Furthermore, we discuss the effect of the number of constraints on the extraction performance. Experimental results confirmed that our proposed method, GC-CSV-Aux-IVE, exhibits superior extraction performance and achieves nearly 100% accuracy in extracting the moving target source.

I. INTRODUCTION

Recently, the applications of Blind Source Separation (BSS) [1]–[3] span across a wide range of fields, contributing to advancements in speech recognition, audio processing, and speech enhancement. BSS algorithms can separate individual speech signals in scenarios where multiple speeches overlap, improving speech recognition accuracy and estimating the speaker's location.

In the field of BSS, Independent Component Analysis (ICA) [4], [5] and Independent Vector Analysis (IVA) [6], [7] have emerged as prominent techniques for separating mixed signals into their constituent sources. Auxiliary function-based Independent Vector Analysis (Aux-IVA) [8], [9] is an extension of IVA that incorporates auxiliary functions to improve the separation performance.

Moreover, in environments with noisy or degraded speech, Blind Source Extraction (BSE) methods [10]–[12] offer the capability to enhance the target speech by extracting it from the mixture, leading to improved speech intelligibility and quality. BSE techniques have emerged as effective methods for extracting specific source signals from mixtures. By utilizing various algorithms and signal processing methods [13], BSE focuses on extracting the Source of Interest (SOI) while suppressing or removing interference from other sources in the mixture.

In recent years, the Constant Separating Vector (CSV) model [14], [15] proposed in the BSE domain has demon-

strated significant advancements in the moving source extraction. Based on the CSV model, [16] proposed CSV-Aux-IVE, which estimates a time-varying mixing vector and a time-invariant filter and utilizes an auxiliary function technique to enhance the extraction performance and the convergence speed of Independent Vector Extraction (IVE).

However, since the permutation problem exists, it is often difficult to determine which source will be extracted in the field of BSE. An important consideration is ensuring that the extracted signal corresponds to SOI. In [14], they addressed this problem by employing the pilot signal which is a signal that is mutually dependent with the corresponding source signal. This method uses prior information about the sound source to resolve the permutation ambiguity.

On the other hand, we propose to solve this ambiguity using Geometric Constraints (GC) [17]. Moreover, we aim to give several constraints for the moving speaker. By employing GC method [18]–[21], we can enhance the accuracy of signal extraction by only requiring approximate location information of the target signal and sensors positions. Consequently, it can be widely applied in scenarios where the signal's moving range is known, such as conferences and speech monitoring. GC method is based on beamforming, which utilizes spatial information and sensor positions to optimize BSS problems. It exploits spatial information to guide the separation matrix to obtain the signal in the desired direction. This method optimizes the constraint weights to ensure that extracted source signal keeps unit energy. Therefore, it is our objective to address the permutation problem and improve the performance of moving target source extraction by applying the GC method.

II. PROBLEM FORMULATION

A. Blind Source Extraction

Considering the situation in BSE, we assume there are M microphones in the room. After the Short Time Fourier Transform (STFT), we obtain microphone signals $x_{f,l}$ and source signals $s_{f,l}$, where f and l represent the frequency bin and time frame index. Our main task is to extract SOI $s_{f,l}$ from M microphones, and $z_{f,l}$ represent the remaining background signals. Generally, we can assume that $s_{f,l}$ corresponds to the

first source of $\mathbf{s}_{f,l}$ without loss of generality. So the relationship of microphone signals $\mathbf{x}_{f,l} = [x_{f,l,1}, \dots, x_{f,l,M}]^T \in \mathbb{C}^M$ and source signals $\mathbf{s}_{f,l} = [s_{f,l}, \mathbf{z}_{f,l}^T]^T \in \mathbb{C}^N$ where $\mathbf{z}_{f,l} = [z_{f,l,1}, \dots, z_{f,l,N-1}]$ can be written as

$$\mathbf{x}_{f,l} = \mathbf{A}_{f,l} \mathbf{s}_{f,l} = \mathbf{A}_{f,l} \begin{bmatrix} s_{f,l} \\ \mathbf{z}_{f,l} \end{bmatrix}, \quad (1)$$

where $\mathbf{A}_{f,l}$ stands for $M \times N$ mixing matrix and $(\cdot)^T$ denotes the transpose. If M is equal to N , the inverse matrix of $\mathbf{A}_{f,l}$ exists and is called separation matrix $\mathbf{W}_{f,l} \in \mathbb{C}^{M \times M}$, so (1) can be written as:

$$\mathbf{s}_{f,l} = \mathbf{W}_{f,l} \mathbf{x}_{f,l} = \begin{bmatrix} \mathbf{w}_f^H \\ \mathbf{B}_{f,l} \end{bmatrix} \mathbf{x}_{f,l}. \quad (2)$$

Here $(\cdot)^H$ denotes Hermit transpose. Then we divide mixing matrix as $\mathbf{A}_{f,l} = [\mathbf{a}_{f,l} \ \mathbf{D}_{f,l}]$ which implies that $\mathbf{a}_{f,l}$ is the first column of $\mathbf{A}_{f,l}$ while $\mathbf{D}_{f,l}$ is the remaining part. Similarly, \mathbf{w}_f^H refers to the first row of the matrix, and $\mathbf{B}_{f,l}$ refers to the rest. Then the mixing and separation model (1) and (2) can be written as:

$$\mathbf{x}_{f,l} = \mathbf{a}_{f,l} s_{f,l} + \mathbf{D}_{f,l} \mathbf{z}_{f,l}, \quad (3)$$

and

$$\begin{bmatrix} s_{f,l} \\ \mathbf{z}_{f,l} \end{bmatrix} = \begin{bmatrix} \mathbf{w}_f^H \mathbf{x}_{f,l} \\ \mathbf{B}_{f,l} \mathbf{x}_{f,l} \end{bmatrix}. \quad (4)$$

B. CSV Mixing Model

In the CSV model, the signal after STFT is divided into T blocks, and the length of each block is defined as L_b . At this time, the index of each block is $t \in \{1, \dots, T\}$. Thus, the microphone signal is divided into several blocks. We consider the signal in each block approximately stationary. Accordingly, the expressions of the corresponding mixing matrix $\mathbf{A}_{f,t}$ in (1) and separation matrix $\mathbf{W}_{f,t}$ in (2) will be rewrite as [12]

$$\mathbf{A}_{f,t} = [\mathbf{a}_{f,t}, \mathbf{D}_{f,t}] = \begin{bmatrix} \gamma_{f,t} & \mathbf{h}_f^H \\ \mathbf{g}_{f,t} & \frac{1}{\gamma_{f,t}} (\mathbf{g}_{f,t} \mathbf{h}_f^H - \mathbf{I}_{M-1}) \end{bmatrix}, \quad (5)$$

$$\mathbf{W}_{f,t} = \begin{bmatrix} \mathbf{w}_f^H \\ \mathbf{B}_{f,t} \end{bmatrix} = \begin{bmatrix} \beta_f^* & \mathbf{h}_f^H \\ \mathbf{g}_{f,t} & -\gamma_{f,t} \mathbf{I}_{M-1} \end{bmatrix}. \quad (6)$$

Here $(\cdot)^*$ denotes the complex conjugate, and \mathbf{I}_M denotes to $M \times M$ identity matrix. In the remaining sections of this paper, we omit the frame index $l = \{1, \dots, L\}$. Given that the mixing matrix and separation matrix are inverse to each other, we can obtain $\mathbf{w}_f^H \mathbf{a}_{f,t} = 1$, which is called distortionless constraint. Similarly, we can get $\mathbf{B}_{f,t} \mathbf{a}_{f,t} = 0$, where $\mathbf{B}_{f,t}$ represents the block matrix. Using the assumption above, we can rewrite (1) and (2) as follows:

$$\mathbf{x}_{f,t} = \begin{bmatrix} \gamma_{f,t} & \mathbf{h}_f^H \\ \mathbf{g}_{f,t} & \frac{1}{\gamma_{f,t}} (\mathbf{g}_{f,t} \mathbf{h}_f^H - \mathbf{I}_{M-1}) \end{bmatrix} \begin{bmatrix} s_{f,t} \\ \mathbf{z}_{f,t} \end{bmatrix}, \quad (7)$$

and

$$\begin{bmatrix} s_{f,t} \\ \mathbf{z}_{f,t} \end{bmatrix} = \begin{bmatrix} \beta_f^* & \mathbf{h}_f^H \\ \mathbf{g}_{f,t} & -\gamma_{f,t} \mathbf{I}_{M-1} \end{bmatrix} \mathbf{x}_{f,t}. \quad (8)$$

We assume that $\mathbf{s}_{f,t}$ and $\mathbf{z}_{f,t}$ are mutually independent across all time frequency bins. Based on this assumption, we can express the joint pdf of them as follows:

$$p(\{\mathbf{s}_t, \mathbf{z}_{f,t}\}_{f,t}) = \prod_t p(\mathbf{s}_t) \prod_{f,t} p(\mathbf{z}_{f,t}), \quad (9)$$

here $p(\mathbf{s}_t)$ denote the joint pdf of the SOI vector component $\mathbf{s}_t = [s_{1,t}, \dots, s_{F,t}] \in \mathbb{C}^F$ and $p(\mathbf{z}_{f,t})$ denote the pdf of $\mathbf{z}_{f,t}$, respectively.

We consider the block-dependent variance of SOI, which can vary from block to block. We introduce pdf, denoted as $p(\mathbf{s}_t)$, that captures this block-dependent variance and accurately represents the characteristics of the source signal. Specifically, the pdf $p(\mathbf{s}_t)$ is defined as follows:

$$p(\mathbf{s}_t) = h \left(\left\{ \begin{bmatrix} s_{f,t} \\ \hat{\sigma}_{f,t} \end{bmatrix} \right\}_f \right) \left(\prod_{f=1}^F \hat{\sigma}_{f,t} \right)^{-2}, \quad (10)$$

where $\hat{\sigma}_{f,t} = \sqrt{\mathbf{w}_f^H \hat{\mathbf{C}}_{f,t} \mathbf{w}_f}$ and $\hat{\mathbf{C}}_{f,t} = \mathbb{E} [\mathbf{x}_{f,t} \mathbf{x}_{f,t}^H]$. $h(\cdot)$ is a pdf corresponding to a normalized non-Gaussian random variable:

$$h \left(\left\{ \begin{bmatrix} s_{f,t} \\ \hat{\sigma}_{f,t} \end{bmatrix} \right\}_f \right) = C \exp(-G_R(r_t)), \quad (11)$$

where C is a coefficient and G_R is a continuous and differentiable function of a real variable r satisfying that $\psi(r) = \frac{G'_R(r)}{r}$ is continuous and monotonically decreasing in $r \geq 0$. The background noise $\mathbf{z}_{f,t}$ is assumed to follow a circular Gaussian distribution with zero mean and a variance matrix $\hat{\mathbf{C}}_{\mathbf{z}_{f,t}} = \mathbb{E} [\mathbf{z}_{f,t} \mathbf{z}_{f,t}^H]$.

After applying the auxiliary function method described in [8], we can obtain

$$\hat{\mathbb{E}} \left[-\log h \left(\left\{ \begin{bmatrix} s_{f,t} \\ \hat{\sigma}_{f,t} \end{bmatrix} \right\}_f \right) \right] \leq \frac{1}{2} \frac{1}{\hat{\sigma}_{f,t}^2} \mathbf{w}_f^H \mathbb{E} [\psi(r_t) \mathbf{x}_{f,t} \mathbf{x}_{f,t}^H] \mathbf{w}_f + R_t, \quad (12)$$

where r_t is an auxiliary variable and R_t depends purely on r_t . Based on the above assumptions, we proceed to apply the auxiliary function method for updating the objective function. Consequently, we obtain the auxiliary function for CSV as follows:

$$\begin{aligned} \mathcal{L}_{\text{Aux}} = & \frac{1}{T} \sum_{t=1}^T \sum_{f=1}^F \left\{ \frac{1}{2} \frac{\mathbf{w}_f^H \mathbf{V}_{f,t} \mathbf{w}_f}{\hat{\sigma}_{f,t}^2} + \log \hat{\sigma}_{f,t}^2 \right. \\ & \left. + \mathbb{E} \left[\mathbf{z}_{f,t}^H \mathbf{C}_{\mathbf{z}_{f,t}}^{-1} \mathbf{z}_{f,t} \right] - (M-2) \log |\gamma_{f,t}|^2 \right\} + R_t, \end{aligned} \quad (13)$$

where

$$\mathbf{V}_{f,t} = \mathbb{E} [\psi(r_t) \mathbf{x}_{f,t} \mathbf{x}_{f,t}^H]. \quad (14)$$

To minimize the auxiliary function with respect to the normal variables, we calculate the partial derivative with respect

to \mathbf{w}_f . In order to minimize the auxiliary function (13), we derivate it and get

$$\frac{\partial \mathcal{L}_{\text{Aux}}}{\partial \mathbf{w}_f^*} = \frac{1}{2T} \sum_{t=1}^T \left\{ \frac{\mathbf{V}_{f,t}}{\hat{\sigma}_{f,t}^2} \mathbf{w}_f - \frac{\mathbf{w}_f^H \mathbf{V}_{f,t} \mathbf{w}_f}{\hat{\sigma}_{f,t}^2} \mathbf{a}_{f,t} \right\}, \quad (15)$$

where we used the same technique as the previous researches [12], [16] to replace the derivative of the third and fourth term in (13) as $\sum_{f=1}^F \mathbf{a}_{f,t}$. Using conventional technique [16], we get an update rule for the separation filter \mathbf{w}_f by taking linearized solution by fixing $\frac{\mathbf{w}_f^H \mathbf{V}_{f,t} \mathbf{w}_f}{\hat{\sigma}_{f,t}^2}$ and $\hat{\sigma}_{f,t}^2$ as constant terms. Then the update rule of \mathbf{w}_f can be written as

$$\mathbf{w}_f = \left(\sum_{t=1}^T \frac{\mathbf{V}_{f,t}}{\hat{\sigma}_{f,t}^2} \right)^{-1} \sum_{t=1}^T \frac{\mathbf{w}_f^H \mathbf{V}_{f,t} \mathbf{w}_f}{\hat{\sigma}_{f,t}^2} \mathbf{a}_{f,t}. \quad (16)$$

III. PROPOSED METHOD

So far, the CSV-Aux-IVE in Section II-B still encounters the permutation problem. In other words, without prior information, it is not determined whether the separate output corresponds to the SOI or Interference (IR). Therefore, we introduce GC [17] to ensure the extraction of the SOI. Then, let us consider GC that restrict the far-field response of the separation filter \mathbf{w}_f at the direction θ , which is described as

$$\mathcal{L}_{\text{GC}} = \lambda_{\text{GC}} \sum_{f=1}^F \sum_{\theta \in \Theta} |\mathbf{w}_f^H \mathbf{d}_{f\theta} - 1|^2, \quad (17)$$

where $\mathbf{d}_{f\theta}$ is the steering vector pointing to the direction θ , λ_{GC} is a parameter weighting the importance of GC in the objective function, and we can preserve the SOI by designing this constraint when θ is the DOA (Direction of Arrival) of the target, which we call Unit Response (UR) constraints. The set Θ represents the collection of angles that will be taken into account.

The target source we aim to extract is moving, while the filter \mathbf{w}_f remains time-invariant. Thus, a challenge lies in how to provide the filter with information regarding the movement of SOI. In this study, we make the assumption that the range of movement for the signals is known. We divide this range into p equally spaced angles and subsequently apply UR constraints to the signals at each angle.

We combine two loss functions (13) and (17) to apply UR constraints to SOI by summing them together to create a composite objective function, denoted as

$$\mathcal{Q}_{\text{CSV+GC}} = \mathcal{L}_{\text{Aux}} + \mathcal{L}_{\text{GC}}. \quad (18)$$

Then we calculate the partial derivative of the objective function. Thus the resulting derivative is as follows:

$$\begin{aligned} \frac{\partial \mathcal{Q}_{\text{CSV+GC}}}{\partial \mathbf{w}_f^*} &= \frac{1}{2T} \sum_{t=1}^T \left\{ \frac{\mathbf{V}_{f,t} \mathbf{w}_f}{\hat{\sigma}_{f,t}^2} - \frac{\mathbf{w}_f^H \mathbf{V}_{f,t} \mathbf{w}_f}{\hat{\sigma}_{f,t}^2} \mathbf{a}_{f,t} \right\} \\ &+ \lambda_{\text{GC}} \sum_{\theta \in \Theta} \mathbf{d}_{f\theta} \mathbf{d}_{f\theta}^H \mathbf{w}_f - \lambda_{\text{GC}} \sum_{\theta \in \Theta} \mathbf{d}_{f\theta}. \end{aligned} \quad (19)$$

Subsequently, we take the linearized solution of \mathbf{w}_f :

$$\mathbf{w}_f = \left(\lambda_{\text{GC}} \sum_{\theta \in \Theta} \mathbf{d}_{f\theta} \mathbf{d}_{f\theta}^H + \frac{1}{2T} \sum_{t=1}^T \frac{\mathbf{V}_{f,t}}{\hat{\sigma}_{f,t}^2} \right)^{-1} \left(\lambda_{\text{GC}} \sum_{\theta \in \Theta} \mathbf{d}_{f\theta} + \frac{1}{2T} \sum_{t=1}^T \frac{\mathbf{w}_f^H \mathbf{V}_{f,t} \mathbf{w}_f}{\hat{\sigma}_{f,t}^2} \mathbf{a}_{f,t} \right). \quad (20)$$

IV. EXPERIMENT

To evaluate the effectiveness of GC-CSV-Aux-IVE, we conducted a speech separation experiments. The evaluation criteria included extraction performance and order of the output signal, which we calculate its average as its accuracy. We compared our proposed method, GC-CSV-Aux-IVE, with the existing CSV-Aux-IVE method without being supervised.

A. Setup

In this experiment, we utilized a subset of the ATR Japanese Speech Database [22], consisting of recordings from three male speakers and three female speakers. We randomly selected two distinct speeches from this dataset to create mixtures. The indoor impulse response is generated through image method [23], with the observation signal generated using signal generator¹. The layout of the simulated room is illustrated in Fig. 1. In total, we generated 100 pairs of original signals and mixed signals. The length of each signal is about 20 seconds. Among these, the first 50 signals had the SOI moving from 90° to 150°, while the remaining 50 signals had the SOI moving from 150° to 90°. Here SOI was set to move at a speed of 3° per second, while the IR is fixed at 45°. We have 6 microphones, 1 moving target source, and 1 fixed interference source.

To make this experiment more convincing, we added a white noise signal to the mixing signal. We adjusted the SNR to input $\text{SNR} = 10 \log \frac{\lambda_1}{\lambda_2}$ as specified value, where λ_1 corresponds to the sample variance of mixing signal and λ_2 corresponds to that of noises. The other experiment conditions are shown in Table I.

In this experiment, we assumed that the movement range of SOI is known. We divided this movement range into p equal parts to set the collection of angles Θ . Then we applied UR constraints to SOI at each angle. Here is an example of how the Θ can be set:

- when $p = 4$, $\Theta = \{90^\circ, 110^\circ, 130^\circ, 150^\circ\}$,
- when $p = 6$, $\Theta = \{90^\circ, 102^\circ, 114^\circ, 126^\circ, 138^\circ, 150^\circ\}$,
- when $p = 16$, $\Theta = \{90^\circ, 94^\circ, 98^\circ, \dots, 142^\circ, 146^\circ, 150^\circ\}$.

For GC-CSV-Aux-IVE, we updated \mathbf{w}_f 50 times while for CSV-Aux-IVE, we updated \mathbf{w}_f 100 times. After applying GC, the convergence speed of the iterations is greatly improved, which results in the reduction of the number of convergence iterations required for convergence, compared to CSV-Aux-IVE. As a result, the convergence can be achieved in significantly fewer iterations.

¹<https://www.audiolabs-erlangen.de/fau/professor/habets/software/signal-generator>

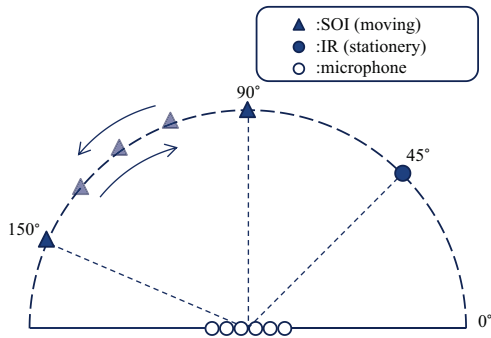


Fig. 1. Configurations of sources and microphones

TABLE I
EXPERIMENTAL CONDITIONS

Sampling rate	16 kHz
STFT window function	Hann
STFT window length	512 samples
STFT shift size	128 samples
Reverberation time	100 ms
Microphone spacing	2 cm
Input SNR	15 dB
L_b	100
λ_{GC}	0.2

TABLE II
AVERAGE SDR, SIR, AND ACCURACY OF EXTRACTING SOI ($M=6$)

method	SDR [dB]	SIR [dB]	Accuracy.
CSV-Aux-IVE	7.89	10.53	40.4%
GC-CSV-Aux-IVE ($p=4$)	11.42	15.73	97.4%
GC-CSV-Aux-IVE ($p=6$)	11.79	16.50	98.2%
GC-CSV-Aux-IVE ($p=16$)	12.36	17.87	99.6%

We exhibited the source extraction performance using the average source-to-distortion ratio (SDR), the source-to-interference ratio (SIR), and the accuracy of target source extraction. These evaluations were calculated by using BSS-eval [24] and setting the dry source as a reference. In order to obtain a more accurate extraction accuracy, we divided the signals into 20 blocks for comparison purposes. Hence we applied the BSS-eval to each block individually and calculated the average of the results from the 20 blocks. And we evaluated the accuracy of the extraction by assessing the output signal order.

B. Result

Table II shows the average SDR, SIR, and accuracy of the output signal. After applying GC, the accuracy of extracting SOI significantly improved from the initial 40% to over 98%. Additionally, the SDR experienced an improvement of approximately 4 dB, while the SIR showed an improvement of around 6 dB. Furthermore, we examined the performance of the extraction method in the different cases of p : $p > M$, $p = M$, and $p < M$ in Table II. Then we evaluated the extraction results under these three scenarios. Interestingly,

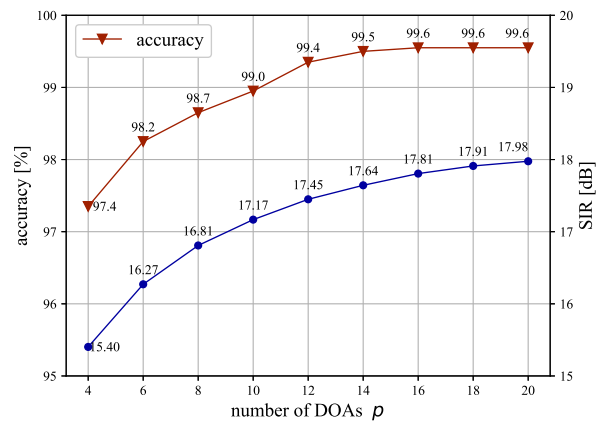


Fig. 2. Effect of increasing p on extraction performance

as the number of UR constraints p was increased, a notable improvement in the extraction performance was observed. The trend of this improvement is depicted in Fig. 2. The increasing trend of both the accuracy and SIR were observed with the increment of p .

V. CONCLUSION

This paper discovered a novel and effective BSE method for moving source. We apply a GC-based approach to address the permutation problem inherent in BSE. Our approach utilizes prior spatial information to aid in the filter updates. We extract the moving SOI by constraining multiple angles within the known range. Then we explore the effect of the changing number of constraint angles on the results. When the number of DOAs p is increased, the extraction effect is also improved. The results also demonstrated that, under similar conditions to the moving target source, GC-CSV-Aux-IVE algorithm achieved an average improvement of approximately 7dB in terms of SIR compared to the original algorithm and achieved an accuracy rate of nearly 100%.

ACKNOWLEDGEMENTS

This work was supported by JSPS KAKENHI Grant Number 23H03423.

REFERENCES

- [1] S. Makino, *Audio source separation*. Springer, 2018.
- [2] S. Y. Low, R. Togneri, and S. Nordholm, "Spatio-temporal processing for distant speech recognition," in *Proc. ICASSP*, vol. 1, 2004, pp. 1–1001.
- [3] P. C. Loizou, *Speech enhancement: theory and practice*. CRC press, 2013.
- [4] A. Hyvärinen and E. Oja, "Independent component analysis: Algorithms and applications," *Neural networks*, vol. 13, no. 4-5, pp. 411–430, 2000.
- [5] H. Sawada, R. Mukai, and S. Makino, "Direction of arrival estimation for multiple source signals using independent component analysis," in *Proc. ISSPA*, vol. 2, 2003, pp. 411–414.
- [6] T. Kim, T. Eltoft, and T.-W. Lee, "Independent vector analysis: An extension of ICA to multivariate components," in *Proc. ICA*, 2006, pp. 165–172.
- [7] N. Ono, "Fast stereo independent vector analysis and its implementation on mobile phone," in *Proc. IWAENC*, 2012, pp. 1–4.

- [8] N. Ono, "Stable and fast update rules for independent vector analysis based on auxiliary function technique," in *Proc. WASPAA*, 2011, pp. 189–192.
- [9] T. Taniguchi, N. Ono, A. Kawamura, and S. Sagayama, "An auxiliary-function approach to online independent vector analysis for real-time blind source separation," in *Proc. HSCMA*, 2014, pp. 107–111.
- [10] R. Scheibler and N. Ono, "Independent vector analysis with more microphones than sources," in *Proc. WASPAA*, 2019, pp. 185–189.
- [11] R. Ikeshita, T. Nakatani, and S. Araki, "Block coordinate descent algorithms for auxiliary-function-based independent vector extraction," *IEEE Trans. SP*, vol. 69, pp. 3252–3267, 2021.
- [12] Z. Koldovský and P. Tichavský, "Gradient algorithms for complex non-Gaussian independent component/vector extraction, question of convergence," *IEEE Trans. SP*, vol. 67, no. 4, pp. 1050–1064, 2018.
- [13] J. Málek, J. Čmejla, and Z. Koldovský, "Blind extraction of target speech source: Three ways of guidance exploiting supervised speaker embeddings," in *Proc. IWAENC*, 2022, pp. 1–5.
- [14] Z. Koldovský, J. Málek, and J. Janský, "Extraction of independent vector component from underdetermined mixtures through block-wise determined modeling," in *Proc. ICASSP*, 2019, pp. 7903–7907.
- [15] V. Kautský, Z. Koldovský, P. Tichavský, and V. Zanzoso, "Cramér-rao bounds for complex-valued independent component extraction: Determined and piecewise determined mixing models," *IEEE Trans. SP*, vol. 68, pp. 5230–5243, 2020.
- [16] J. Janský, Z. Koldovský, J. Málek, T. Kounovský, and J. Čmejla, "Auxiliary function-based algorithm for blind extraction of a moving speaker," *EURASIP JASMP*, vol. 2022, no. 1, pp. 1–16, 2022.
- [17] L. Parra and C. Alvino, "Geometric source separation: merging convolutive source separation with geometric beamforming," *IEEE Trans. SAP*, vol. 10, no. 6, pp. 352–362, 2002.
- [18] H. Saruwatari, T. Kawamura, T. Nishikawa, A. Lee, and K. Shikano, "Blind source separation based on a fast-convergence algorithm combining ICA and beamforming," *IEEE Trans. ASLP*, vol. 14, no. 2, pp. 666–678, 2006.
- [19] L. Li and K. Koishida, "Geometrically constrained independent vector analysis for directional speech enhancement," in *Proc. ICASSP*, 2020, pp. 846–850.
- [20] K. Goto, T. Ueda, L. Li, T. Yamada, and S. Makino, "Geometrically constrained independent vector analysis with auxiliary function approach and iterative source steering," in *Proc. EUSIPCO*, 2022, pp. 757–761.
- [21] A. H. Khan, M. Taseska, and E. A. P. Habets, "A geometrically constrained independent vector analysis algorithm for online source extraction," in *Proc. LVA/ICA*, 2015, pp. 396–403.
- [22] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara, and K. Shikano, "ATR Japanese speech database as a tool of speech recognition and synthesis," *Speech communication*, vol. 9, no. 4, pp. 357–363, 1990.
- [23] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," *The Journal of the Acoustical Society of America*, vol. 65, no. 4, pp. 943–950, 1979.
- [24] E. Vincent, R. Gribonval, and C. Févotte, "Performance measurement in blind audio source separation," *IEEE Trans. ASLP*, vol. 14, no. 4, pp. 1462–1469, 2006.