

RESEARCH

Open Access

# Consistent independent low-rank matrix analysis for determined blind source separation



Daichi Kitamura<sup>1†\*</sup>  and Kohei Yatabe<sup>2†</sup>

## Abstract

Independent low-rank matrix analysis (ILRMA) is the state-of-the-art algorithm for blind source separation (BSS) in the determined situation (the number of microphones is greater than or equal to that of source signals). ILRMA achieves a great separation performance by modeling the power spectrograms of the source signals via the nonnegative matrix factorization (NMF). Such a highly developed source model can solve the permutation problem of the frequency-domain BSS to a large extent, which is the reason for the excellence of ILRMA. In this paper, we further improve the separation performance of ILRMA by additionally considering the general structure of spectrograms, which is called *consistency*, and hence, we call the proposed method *Consistent ILRMA*. Since a spectrogram is calculated by an overlapping window (and a window function induces spectral smearing called main- and side-lobes), the time-frequency bins depend on each other. In other words, the time-frequency components are related to each other via the uncertainty principle. Such co-occurrence among the spectral components can function as an assistant for solving the permutation problem, which has been demonstrated by a recent study. On the basis of these facts, we propose an algorithm for realizing Consistent ILRMA by slightly modifying the original algorithm. Its performance was extensively evaluated through experiments performed with various window lengths and shift lengths. The results indicated several tendencies of the original and proposed ILRMA that include some topics not fully discussed in the literature. For example, the proposed Consistent ILRMA tends to outperform the original ILRMA when the window length is sufficiently long compared to the reverberation time of the mixing system.

**Keywords:** Audio source separation, Convolutional mixture, Demixing filter estimation, Phase-aware signal processing, Spectrogram consistency

## 1 Introduction

Blind source separation (BSS) is a technique for separating individual sources from an observed mixture without knowing how they were mixed. BSS for multichannel audio signals observed by multiple microphones has been particularly studied [1–13]. The BSS problem can be divided into two situations: underdetermined (the number of microphones is less than the number of sources) and (over-)determined (the number of microphones is

greater than or equal to the number of sources) cases. This paper focuses on the determined BSS problem, as high-quality separation can be achieved compared with the underdetermined BSS methods.

Independent component analysis (ICA) is the most popular and successful algorithm for solving the determined BSS problem [1]. It estimates a demixing matrix (the inverse system of the mixing process) by assuming statistical independence between the sources. For a mixture of audio signals, ICA is usually applied in the time-frequency domain via the short-time Fourier transform (STFT) because the sources are mixed up by convolution. This strategy is called frequency-domain ICA (FDICA)

\*Correspondence: [d-kitamura@ieee.org](mailto:d-kitamura@ieee.org)

<sup>†</sup>Daichi Kitamura and Kohei Yatabe contributed equally to this work.

<sup>1</sup>National Institute of Technology, Kagawa College, 355 Chokushi, Takamatsu, Kagawa, 761-8058, Japan

Full list of author information is available at the end of the article

[2] and independently applies ICA to the complex-valued signals in each frequency. Then, the estimated frequency-wise demixing matrices must be aligned over all frequencies so that the frequency components of the same source are grouped together. Such alignment of the frequency components is called a *permutation problem* [3–6], and a complete solution to it has not been established. Therefore, a great deal of research has tackled this problem.

To avoid the permutation misalignment as much as possible, various sophisticated source models have been proposed. Independent vector analysis (IVA) [7–10] is one of the most successful methods in the early stage of the development. It assumes higher-order dependencies (co-occurrence among the frequency components) of each source by utilizing a spherical generative model of the source frequency vector. This assumption enables IVA to simultaneously estimate the frequency-wise demixing matrices and solve the permutation problem to a large extent using only one objective function. It has been further developed by improving its source model. One natural and powerful extension of IVA is independent low-rank matrix analysis (ILRMA) [11, 12], which integrates the source model of nonnegative matrix factorization (NMF) [14, 15] based on the Itakura–Saito divergence (IS-NMF) [16] into IVA. This extension has greatly improved the performance of separation by taking the low-rank time-frequency structure (co-occurrence among the time-frequency bins) of the source signals into account. ILRMA has achieved the state-of-the-art performance and been further developed by several researchers [17–29]. In this respect, ILRMA can be considered the new standard of the determined BSS algorithms. However, the separation performance of IVA and ILRMA is still inferior compared to the ideal performance of ICA-based frequency-domain BSS. In [30], the performances of IVA and ILRMA were compared with that of FDICA with perfect permutation alignment using reference sources (ideal permutation solver), and it was confirmed that there is still a noticeable room for improvement of ILRMA-based BSS. In fact, IVA and ILRMA often encounter the block permutation problem, that is, group-wise permutation misalignment of components between sources [31].

The *consistency* of a spectrogram is another promising approach for solving the permutation problem. A recent study has shown that STFT can provide some effective information related to the co-occurrence among the time-frequency bins [32]. Since an overlapping window is utilized in STFT, the time-frequency bins are related to each other based on the overlapping segments. The frequency components within a segment are also related to each other because of the spectral smearing called main- and side-lobes of the window. In other words, the time-frequency components are not independent but related to

each other via the *uncertainty principle* of time-frequency representation. Such relations have been well-studied in phase-aware signal processing [33–43] by the name of spectrogram consistency [44–47]. In the previous study [32], the spectrogram consistency was imposed on BSS to help the algorithm solve the permutation problem. This is an approach very different from the conventional studies of determined BSS because it utilizes the general property of STFT *independent of the source model* (in contrast to the abovementioned methods that focused on modeling of the source signals without considering the property of STFT). As the spectrogram consistency can be incorporated with any source model, its combination with the state-of-the-art algorithm should achieve a high separation performance.

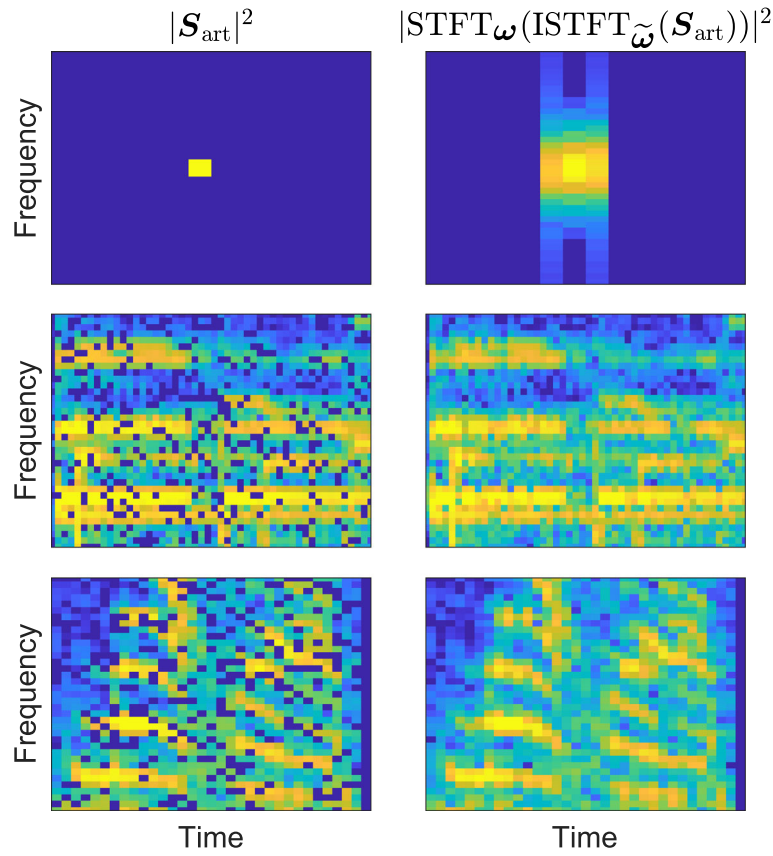
However, the paper that proposed the combination of consistency and determined BSS [32] only showed the potential of consistency in an experiment using FDICA and IVA. The paper claimed that it was a first step of incorporating the spectrogram consistency with determined BSS, and no advanced method was tested. In particular, ILRMA was not considered because its algorithm is far more complicated than that derived in [32], and thus, it is not clear whether (and how much) the spectrogram consistency might improve the state-of-the-art BSS algorithm.

In this paper, we propose a new variant of ILRMA called *Consistent ILRMA* that considers the spectrogram consistency within the algorithm of ILRMA. The combination of IS-NMF and spectral smoothing of the inverse STFT (see Figs. 1 and 2 in Section 2.3) achieves the source modeling for a complex spectrogram. In particular, the spectral smearing in the frequency direction ties the adjacent frequency bins together, and this effect of spectrogram consistency helps ILRMA to solve the permutation problem. Since consistency is a concept depending on the parameters related to a window function, we extensively tested the separation performance of Consistent ILRMA through experiments with various window lengths and shift lengths. The results clarified several tendencies of the conventional and proposed methods, including that the proposed method outperforms the original ILRMA when the window length is sufficiently long compared to the reverberation time of the mixing system.

## 2 Permutation problem of frequency-domain BSS and spectrogram consistency

### 2.1 Formulation of frequency-domain BSS

Let the  $l$ th sample of a time-domain signal be denoted as  $x[l]$ , and  $N$  source signals be observed by  $M$  microphones. Then, the  $l$ th samples of the multichannel source, observed, and separated signals are respectively denoted as:



**Fig. 1** Inconsistent power spectrograms  $|S_{\text{art}}|^2$  (left column) and their consistent version (right column) obtained by applying inverse STFT and STFT. The top-left spectrogram is artificially produced with random phase. The middle-left and the bottom-left spectrograms are music and speech signals with random dropout. Enforcing spectrogram consistency can be viewed as a smoothing process of the inconsistent spectrogram along both time and frequency axes

$$s[l] = [s_1[l], s_2[l], \dots, s_n[l], \dots, s_N[l]]^T \in \mathbb{R}^N, \quad (1)$$

$$x[l] = [x_1[l], x_2[l], \dots, x_m[l], \dots, x_M[l]]^T \in \mathbb{R}^M, \quad (2)$$

$$y[l] = [y_1[l], y_2[l], \dots, y_n[l], \dots, y_N[l]]^T \in \mathbb{R}^N, \quad (3)$$

where  $n = 1, \dots, N$ ,  $m = 1, \dots, M$ , and  $l = 1, \dots, L$  are the indexes of sources, microphones (channels), and discrete time, respectively, and  $\cdot^T$  denotes the transpose. BSS aims at recovering the source signal  $s$  from the observed signal  $x$ , i.e., making  $y$  as close to  $s$  as possible.

In the frequency-domain BSS, those signals are handled in the time-frequency domain via STFT. Let the window length and shifting step of STFT be denoted as  $Q$  and  $\tau$ , respectively. Then, the  $j$ th segment of a signal  $z[l]$  is defined as:

$$\begin{aligned} z^{[j]} &= [z[(j-1)\tau+1], z[(j-1)\tau+2], \dots, z[(j-1)\tau+Q]]^T, \\ &= [z^{[j]}[1], z^{[j]}[2], \dots, z^{[j]}[q], \dots, z^{[j]}[Q]]^T \in \mathbb{R}^Q, \end{aligned} \quad (4)$$

where  $j = 1, \dots, J$  and  $q = 1, \dots, Q$  are the indexes of the segments and in-segment samples, respectively, and

the number of segments is given by  $J = L/\tau$  with some zero-padding for adjusting the signal length  $L$  if necessary. STFT of a signal  $z = [z[1], z[2], \dots, z[L]]^T \in \mathbb{R}^L$  is denoted by:

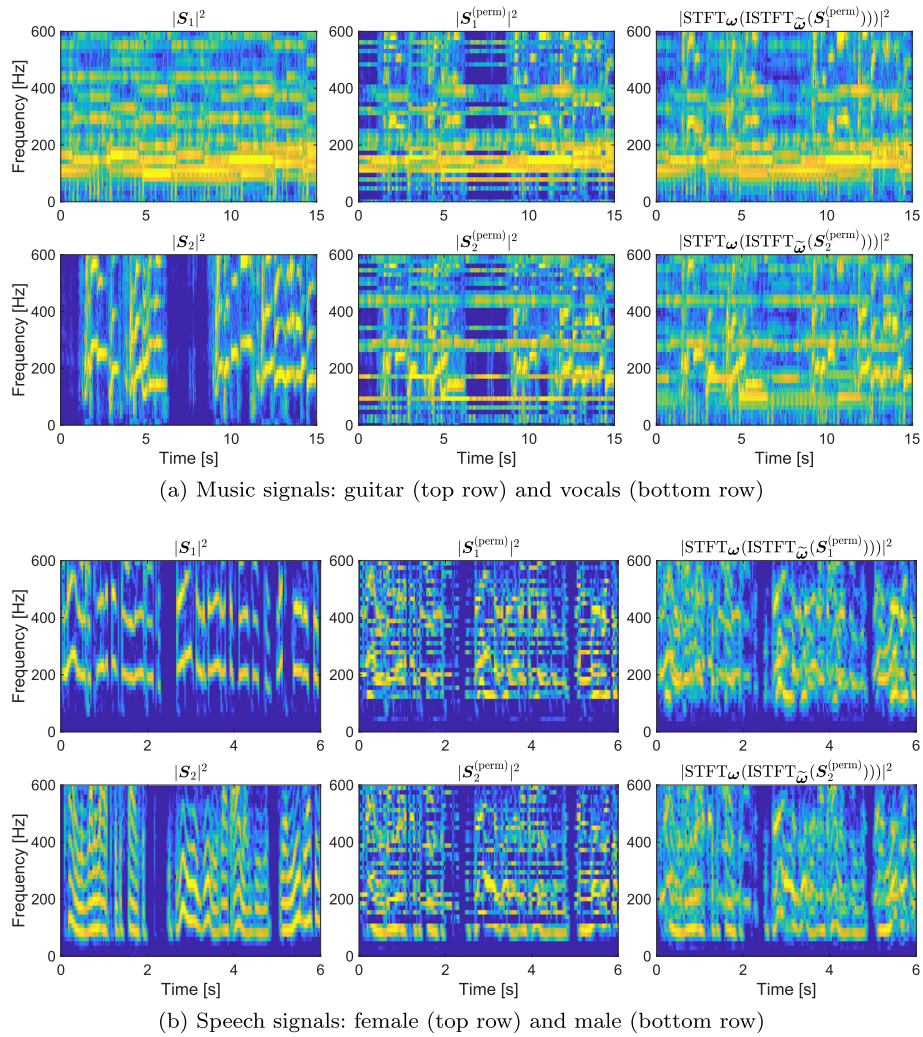
$$Z = \text{STFT}_\omega(z) \in \mathbb{C}^{I \times J}, \quad (5)$$

where the  $(i, j)$ th bin of the spectrogram  $Z$  is given as:

$$z_{ij} = \sum_{q=1}^Q \omega[q] z^{[j]}[q] e^{-i2\pi(q-1)(i-1)/F}, \quad (6)$$

$i = 1, \dots, I$  is the index of frequency bins,  $F$  is an integer satisfying  $\lfloor F/2 \rfloor + 1 = I$ ,  $\lfloor \cdot \rfloor$  is the floor function,  $i$  denotes the imaginary unit, and  $\omega$  is an analysis window. The inverse STFT with a synthesis window  $\tilde{\omega}$  is also defined in the usual way and denoted as  $\text{ISTFT}_{\tilde{\omega}}(\cdot)$ . In this paper, we assume that the window pair satisfies the following perfect reconstruction condition:

$$z = \text{ISTFT}_{\tilde{\omega}}(\text{STFT}_\omega(z)) \quad \forall z \in \mathbb{R}^L. \quad (7)$$



**Fig. 2** Smoothing effect of spectrogram consistency applied to permutation misaligned signals: **a** music and **b** speech. The left column shows the original source signals  $|S_n|^2$ , and the center column shows their randomly permuted versions, which simulates the permutation problem and is denoted as  $S_n^{(\text{perm})}$ . The right column shows the consistent versions of  $S_n^{(\text{perm})}$ . The smoothing effect mixes up the signals

By applying STFT, the  $(i, j)$ th bin of the spectrograms of source, observed, and separated signals can be written as:

$$\mathbf{s}_{ij} = [s_{ij1}, s_{ij2}, \dots, s_{ijn}, \dots, s_{ijN}]^T \in \mathbb{C}^N, \quad (8)$$

$$\mathbf{x}_{ij} = [x_{ij1}, x_{ij2}, \dots, x_{ijm}, \dots, x_{ijM}]^T \in \mathbb{C}^M, \quad (9)$$

$$\mathbf{y}_{ij} = [y_{ij1}, y_{ij2}, \dots, y_{ijn}, \dots, y_{ijN}]^T \in \mathbb{C}^N. \quad (10)$$

We also denote the spectrograms corresponding to the  $n$ th or  $m$ th signals in (8)–(10) as  $S_n \in \mathbb{C}^{I \times J}$ ,  $X_m \in \mathbb{C}^{I \times J}$ , and  $Y_n \in \mathbb{C}^{I \times J}$ , whose elements are  $s_{ijn}$ ,  $x_{ijm}$ , and  $y_{ijn}$ , respectively. In the ordinary frequency-domain BSS, an instantaneous mixing process for each frequency bin is assumed:

$$\mathbf{x}_{ij} = \mathbf{A}_i \mathbf{s}_{ij}, \quad (11)$$

where  $\mathbf{A}_i \in \mathbb{C}^{M \times N}$  is a frequency-wise mixing matrix. The mixture model (11) is approximately valid when the reverberation time is sufficiently shorter than the length of the analysis window used in STFT [48].

Hereafter, we consider the determined case, i.e.,  $M = N$ . In this case, BSS can be achieved by estimating the inverse of  $\mathbf{A}_i$  for all frequency bins. By denoting an approximate inverse as  $\mathbf{W}_i \approx \mathbf{A}_i^{-1}$ , the separation process can be written as:

$$\mathbf{y}_{ij} = \mathbf{W}_i \mathbf{x}_{ij}, \quad (12)$$

where  $\mathbf{W}_i = [\mathbf{w}_{i1}, \mathbf{w}_{i2}, \dots, \mathbf{w}_{iN}]^H \in \mathbb{C}^{N \times M}$  is a frequency-wise demixing matrix and  $\cdot^H$  denotes the Hermitian transpose. The aim of a determined BSS algorithm is to find



the demixing matrices for all frequency bins so that the separated signals approximate the source signals.

## 2.2 Permutation problem in determined BSS

In practice, the scale and permutation of the separated signals are unknown because the information of the mixing process is missing. That is, when the separation is correctly performed by some demixing matrix  $\mathbf{W}_i$  as in (12), the following signal is also a solution to the BSS problem:

$$\hat{\mathbf{y}}_{ij} = \hat{\mathbf{W}}_i \mathbf{x}_{ij} \quad \left( \hat{\mathbf{W}}_i = \mathbf{D}_i \mathbf{P}_i \mathbf{W}_i \right), \quad (13)$$

where  $\mathbf{D}_i \in \mathbb{C}^{N \times N}$  and  $\mathbf{P}_i \in \{0, 1\}^{N \times N}$  are arbitrary diagonal and permutation matrices, respectively. While the signal scale can easily be recovered by applying the back projection [49], the permutation of the estimated signals  $\hat{\mathbf{y}}_{ij}$  must be aligned for all frequency bins, i.e.,  $\mathbf{P}_i$  must be the same for all  $i$ . This alignment of the permutation of estimated signals is the permutation problem, which is the main obstacle of the frequency-domain determined BSS.

In FDICA, a permutation solver (realignment process of  $\mathbf{P}_i$ ) is utilized as a post-processing applied to the frequency-wise separated signals  $\hat{\mathbf{y}}_{ij}$  [4–6]. In recent frequency-domain BSS methods, an additional assumption on sources (or source model) is introduced to circumvent the permutation problem. For example, IVA assumes simultaneous co-occurrence of all frequency components in the same source, and ILRMA assumes a low-rank structure of the power spectrogram  $\mathbf{Y}_n$ . Other source models have also been proposed for improving the separation performance [50–52]. These source models can avoid the permutation problem to some extent during the estimation of  $\hat{\mathbf{W}}_i$ . Recent developments of determined BSS have been achieved via the quest to find a better source model that represents the source signals more precisely.

## 2.3 Solving permutation problem by spectrogram consistency

A recent paper reported another approach for solving the permutation problem based on the general property of STFT called spectrogram consistency [32]. The consistency is a fundamental property of a spectrogram. Since any time-frequency representation has a theoretical limitation called the uncertainty principle, the time-frequency bins of a spectrogram are not independent but related to each other. The inverse STFT always modifies the spectrogram  $\mathbf{Z}_n$  that violates this kind of inter-time-frequency relation so that the relation is recovered. That is, a spectrogram  $\mathbf{Z}_n$  properly retains the inter-time-frequency relation if and only if

$$\mathcal{E}(\mathbf{Z}_n) = \mathbf{Z}_n - \text{STFT}_\omega(\text{ISTFT}_{\tilde{\omega}}(\mathbf{Z}_n)) \quad (14)$$

is zero, i.e.,  $\|\mathcal{E}(\mathbf{Z}_n)\| = 0$  for a norm  $\|\cdot\|$ . Such spectrogram  $\mathbf{Z}_n$  satisfying  $\|\mathcal{E}(\mathbf{Z}_n)\| = 0$  is said to be *consistent*.

Figure 1 demonstrates the effect of spectrogram consistency, where  $\mathcal{S}_{\text{art}} \in \mathbb{C}^{I \times J}$  is an artificially produced complex-valued spectrogram and  $|\mathcal{S}_{\text{art}}|^2$  is its power spectrogram. The notation  $|\cdot|^2$  for a matrix input represents the element-wise squared absolute value. By applying  $\text{STFT}_\omega(\text{ISTFT}_{\tilde{\omega}}(\cdot))$ , the inconsistent spectrogram  $\mathcal{S}_{\text{art}}$  shown in the left column of Fig. 1 is converted into the corresponding consistent spectrogram, which is a smoothed version of  $\mathcal{S}_{\text{art}}$ , as shown in the right column. This smoothing process occurs because the main- and side-lobes of the window function (and the overlap-add process) spread the energy of a time-frequency bin.

Since the inverse STFT is a process of recovering the consistency (the inter-time-frequency relation), it has the capability of aligning the frequency components. This is also demonstrated in Fig. 2. As a simulation of the permutation problem, the frequency bins in  $\mathcal{S}_1$  and  $\mathcal{S}_2$  were randomly shuffled to obtain the spectrogram with permutation misalignment,  $\mathcal{S}_n^{(\text{perm})}$  (the center column in the figure), which is a typical output signal of FDICA. Note that these misaligned spectrograms are perfectly separated for each frequency because each time-frequency bin contains only one of the two sources. By enforcing spectrogram consistency, the smoothing process spreads the time-frequency components as shown in the right column of Fig. 2. In other words, the inverse STFT mixes up the separated signals if the frequency-wise permutation is not aligned correctly. Therefore, enforcing consistency within a BSS algorithm by applying  $\text{STFT}_\omega(\text{ISTFT}_{\tilde{\omega}}(\cdot))$  can improve the separation performance to some extent [32].

## 3 Proposed method

By incorporating spectrogram consistency into ILRMA, we propose a novel BSS method named *Consistent ILRMA*. In this section, after stating our motivation and contributions, we first review the standard ILRMA introduced in [11, 12] and then propose the consistent version of ILRMA with an algorithm that achieves Consistent ILRMA and is openly available on the web.

### 3.1 Motivations and contributions

The previous paper [32] only reported that the performances of traditional BSS algorithms, FDICA and IVA, were improved by enforcing consistency during the estimation of the demixing matrix  $\mathbf{W}_i$ . In addition, no detailed experimental analysis related to STFT parameters was provided, even though the parameters of window functions in the STFT and inverse STFT directly affect the smoothing effect of spectrogram consistency.

The spectrogram consistency is a general property of STFT, and therefore, it can be combined with any source model for determined BSS. Its combination with state-of-the-art models, including ILRMA, is of great interest

because the current mainstream algorithm for determined audio source separation is centered on ILRMA, which is based on an NMF-based richer time-frequency source model. Indeed, many recent papers are based on the framework of ILRMA [17–29]. Even though combining ILRMA with the spectrogram consistency should be able to exceed the limit of existing BSS algorithms, no such method has been investigated in the literature.

In this paper, we propose a new BSS algorithm that combines ILRMA and spectrogram consistency. Our first contribution is an algorithm that achieves Consistent ILRMA by inserting  $\text{STFT}_\omega(\text{ISTFT}_\omega(\cdot))$  into the iterative optimization algorithm of ILRMA. The second contribution is to apply a scale-aligning process called *iterative back projection* within the iterative algorithm. This process enhances the separation performance when it is combined with spectrogram consistency. The third contribution is an experimental finding that spectrogram consistency can work properly with the iterative back projection. We found that both Consistent IVA and Consistent ILRMA require iterative back projection to achieve a good performance. Our fourth contribution is to provide the massive experimental results for several window functions, window lengths, shift lengths, reverberation times, and source types. We also provide discussions for clarifying the tendency of ILRMA with spectrogram consistency.

### 3.2 Standard ILRMA [12]

The original ILRMA [12] was derived from the following generative model of the spectrograms of the separated signals:

$$Y_n \sim p(Y_n) = \prod_{ij} \mathcal{N}_c(0, r_{ijn}) = \prod_{ij} \frac{1}{\pi r_{ijn}} \exp\left(-\frac{|y_{ijn}|^2}{r_{ijn}}\right), \quad (15)$$

where  $\mathcal{N}_c(\mu, r)$  is the circularly symmetric complex Gaussian distribution with mean  $\mu$  and variance  $r$ . In this model, the source component  $y_{ijn}$  is assumed to obey a zero-mean and isotropic distribution, i.e., the phase of  $y_{ijn}$  is generated from the uniform distribution in the range  $[0, 2\pi)$  and the real and imaginary parts of  $y_{ijn}$  are mutually independent. The validity of this assumption is shown in the [Appendix](#). The variance  $r_{ijn}$  can be viewed as an expectation value of  $|y_{ijn}|^2$ . This variance  $r_{ijn}$  as a two-dimensional array indexed by  $(i, j)$  is denoted as  $\mathbf{R}_n \in \mathbb{R}_{>0}^{I \times J}$ , which is called the variance spectrogram corresponding to the  $n$ th source. In ILRMA, the variance matrix  $\mathbf{R}_n$  is modeled using the rank- $K$  NMF, as:

$$\mathbf{R}_n = \mathbf{T}_n \mathbf{V}_n, \quad (16)$$

where  $\mathbf{T}_n \in \mathbb{R}_{>0}^{I \times K}$  and  $\mathbf{V}_n \in \mathbb{R}_{>0}^{K \times J}$  are the basis and activation matrices in NMF. The basis vectors in  $\mathbf{T}_n$ , which

represent spectral patterns of the  $n$ th source signal, are indexed by  $k = 1, \dots, K$ . As in FDICA, statistical independence between the source signals is also assumed in ILRMA:

$$p(Y_1, Y_2, \dots, Y_N) = \prod_n p(Y_n). \quad (17)$$

ILRMA estimates the demixing matrix  $\mathbf{W}_i$  so that the power spectrograms of the separated signals  $|Y_n|^2$  have a low-rank structure that can be well-approximated by  $\mathbf{T}_n \mathbf{V}_n$  with small  $K$ . This BSS principle of ILRMA is illustrated in Fig. 3. When the low-rank source model can appropriately fit to the power spectrograms of the original source signals  $|S_n|^2$ , ILRMA provides an excellent separation performance without explicitly solving the permutation problem afterward.

The demixing matrix  $\mathbf{W}_i$  and the nonnegative matrices  $\mathbf{T}_n$  and  $\mathbf{V}_n$  can be obtained through maximum likelihood estimation. The negative log-likelihood to be minimized, denoted by  $\mathcal{L}$ , is given as [12]:

$$\begin{aligned} \mathcal{L} &= -\log p(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M), \\ &= -\sum_{ij} \log |\det \mathbf{W}_i|^2 - \log p(Y_1, Y_2, \dots, Y_N), \\ &\stackrel{c}{=} -2J \sum_i |\det \mathbf{W}_i| + \sum_{i,j,n} \left( \frac{|\mathbf{w}_{in}^H \mathbf{x}_{ij}|^2}{\sum_k t_{ikn} v_{kjn}} + \log \sum_k t_{ikn} v_{kjn} \right), \end{aligned} \quad (18)$$

where  $\stackrel{c}{=}$  denotes equality up to constant factors, and  $t_{ikn} > 0$  and  $v_{kjn} > 0$  are the elements of  $\mathbf{T}_n$  and  $\mathbf{V}_n$ , respectively. The minimization of (18) can be performed by iterating the following update rules for the spatial model parameters,

$$\mathbf{u}_{in} \leftarrow \frac{1}{J} \sum_j \frac{1}{\sum_k t_{ikn} v_{kjn}} \mathbf{x}_{ij} \mathbf{x}_{ij}^H, \quad (19)$$

$$\mathbf{w}_{in} \leftarrow (\mathbf{W}_i \mathbf{u}_{in})^{-1} \mathbf{e}_n, \quad (20)$$

$$\mathbf{w}_{in} \leftarrow \mathbf{w}_{in} (\mathbf{w}_{in}^H \mathbf{u}_{in} \mathbf{w}_{in})^{-\frac{1}{2}}, \quad (21)$$

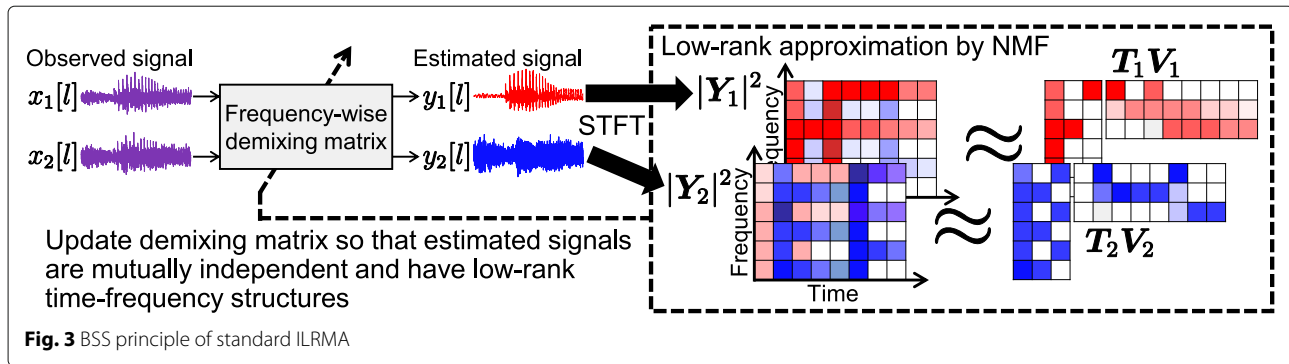
$$y_{ijn} \leftarrow \mathbf{w}_{in}^H \mathbf{x}_{ij}, \quad (22)$$

and for the source model parameters,

$$t_{ikn} \leftarrow t_{ikn} \sqrt{\frac{\sum_j |y_{ijn}|^2 (\sum_{k'} t_{ik'n} v_{k'jn})^{-2} v_{kjn}}{\sum_j (\sum_{k'} t_{ik'n} v_{k'jn})^{-1} v_{kjn}}}, \quad (23)$$

$$v_{kjn} \leftarrow v_{kjn} \sqrt{\frac{\sum_i |y_{ijn}|^2 (\sum_{k'} t_{ik'n} v_{k'jn})^{-2} t_{ikn}}{\sum_i (\sum_{k'} t_{ik'n} v_{k'jn})^{-1} t_{ikn}}}, \quad (24)$$

where  $\mathbf{e}_n \in \{0, 1\}^N$  is the unit vector with the  $n$ th element equal to unity. Update rules (19)–(24) ensure the monotonic non-increase of the negative log-likelihood function



$\mathcal{L}$ . After iterative calculations of updates (19)–(24), the separated signal can be obtained by (12).

Equation 22 is equivalent to beamforming [53] to  $x_{ij}$  with the beamformer coefficients  $w_{in}$ . Thus, FDICA, IVA, and ILRMA can be interpreted as an adaptive estimation process of beamforming coefficients without having to know the geometry of microphones and sources [54]. For this reason, the estimated signal  $Y_n$  obtained by (22) is a complex-valued spectrogram, and we do not need to recover its phase components using, for example, Griffin–Lim algorithm-based techniques [37–40, 43, 55–59]. Both the amplitude and phase components of each source are recovered by the complex-valued linear separation filter  $w_{in}$ .

### 3.3 Proposed Consistent ILRMA

To further improve the separation performance of the standard ILRMA, we introduce the spectrogram consistency into the parameter update procedure. In the proposed Consistent ILRMA, the following combination of forward and inverse STFT is performed at the beginning of each iteration of parameter updates:

$$Y_n \leftarrow \text{STFT}_\omega(\text{ISTFT}_{\tilde{\omega}}(Y_n)). \quad (25)$$

This procedure is the projection of the spectrogram of a separated signal  $Y_n$  onto the set of consistent spectrograms [32]. That is,  $\text{STFT}_\omega(\text{ISTFT}_{\tilde{\omega}}(Y_n))$  performs nothing if  $Y_n$  is consistent, but otherwise, it smooths the complex spectrogram  $Y_n$ , by going through the time domain, so that the uncertainty principle is satisfied.

In Consistent ILRMA, the calculation of (25) is performed in each iteration of parameter updates based on (19)–(24). Enforcing the spectrogram consistency for the temporary separated signal  $Y_n$  in each iteration guides the parameters  $W_i$ ,  $T_n$ , and  $V_n$  to better solutions, which results in higher separation performance compared to that of conventional ILRMA.

Note that this simple update (25) may increase the value of the negative log-likelihood function (18), and therefore, the monotonicity of the algorithm is no longer guaranteed. However, we will see later in the experiments that

the value of the negative log-likelihood function stably decreases as in the standard ILRMA. The amount of the inconsistent component (14) also settles down to some specific value after several iterations.

### 3.4 Iterative back projection

Since frequency-domain BSS cannot determine the scales of estimated signals (represented by  $D_i$  in (13)), the spectrogram of a separated signal  $Y_n$  after an iteration is inconsistent due to the scale irregularity. To take full advantage of the projection enforcing spectrogram consistency in (25), we also propose applying the following back projection at the end of each iteration so that the frequency-wise scales are aligned.

In determined BSS, the back projection is a standard procedure for recovering the frequency-wise scales. It can be written as [49]:

$$\tilde{y}_{ijn} = W_i^{-1}(\mathbf{e}_n \circ \mathbf{y}_{ij}) = y_{ijn} \lambda_{in}, \quad (26)$$

where  $\tilde{\mathbf{y}}_{ijn} = [\tilde{y}_{ijn1}, \tilde{y}_{ijn2}, \dots, \tilde{y}_{ijnM}]^T \in \mathbb{C}^M$  is the  $(i, j)$ th bin of the scale-fitted spectrogram of the  $n$ th separated signal,  $\lambda_{in} = [\lambda_{in1}, \lambda_{in2}, \dots, \lambda_{inM}]^T \in \mathbb{C}^M$  is a coefficient vector of back projection for the  $n$ th signal at the  $i$ th frequency, and  $\circ$  denotes the element-wise multiplication. In the proposed method, this update (26) is performed at the end of each iteration so that the projection (25) at the beginning of the next iteration properly smooths the spectrograms without the effect of scale indeterminacy.

One side effect of this back projection is that the value of the negative log-likelihood function (18) is also changed due to the scale modification. In IVA, this problem cannot be avoided because the only parameter in IVA is the demixing matrix  $W_i$ . However, in ILRMA, since both the demixing matrix  $W_i$  and the source model parameter  $T_n V_n$  can determine the scale of estimated signal  $Y_n$ , the likelihood variation can be avoided by appropriately adjusting  $w_{in}$  and  $T_n$  after the back projection. To prevent the likelihood variation, the following updates are required after performing (26):

**Algorithm 1** Consistent ILRMA**Input:**  $\{\mathbf{x}_{ij}\}_{i=1,j=1}^{I,J}$ , maxlter**Output:**  $\{\mathbf{y}_{ij}\}_{i=1,j=1}^{I,J}$ 

- 1: Initialize  $\{\mathbf{T}_n\}_{n=1}^N, \{\mathbf{V}_n\}_{n=1}^N, \{\mathbf{W}_i\}_{i=1}^I$
- 2: **for** iter = 1, 2,  $\dots$ , maxlter **do**
- 3:   Ensure consistency by calculating (25)  $\forall n$
- 4:   Update source model by calculating (23) and (24)  $\forall i, j, k, n$
- 5:   Update spatial model by calculating (19)–(22)  $\forall i, j, n$
- 6:   Apply back projection by calculating (26)  $\forall i, j, n$
- 7:   Update parameters by calculating (27)–(29)  $\forall i, j, k, n$
- 8: **end for**

$$\mathbf{w}_{in} \leftarrow \mathbf{w}_{in} \lambda_{inm_{\text{ref}}}, \quad (27)$$

$$y_{ijn} \leftarrow \mathbf{w}_{in}^H \mathbf{x}_{ij}, \quad (28)$$

$$t_{ikn} \leftarrow t_{ikn} |\lambda_{inm_{\text{ref}}}|^2, \quad (29)$$

where  $m_{\text{ref}}$  is the index of the reference channel utilized in the back projection.

The overall algorithm of the proposed Consistent ILRMA is summarized in Algorithm 1. The iterative loop for the parameter optimization appears in the second to eighth lines. The spectrogram consistency of the temporary separated signal  $\mathbf{Y}_n$  is ensured in the third line, and the iterative back projection is applied in the sixth and seventh lines. Note that an algorithm for the conventional ILRMA can be obtained by performing only the fourth and fifth lines (i.e., ignoring the third, sixth, and seventh lines). A Python code of the conventional ILRMA is openly available online (<https://pyroomacoustics.readthedocs.io/en/pypi-release/pyroomacoustics.bss.ilrma.html>), and therefore, the proposed Consistent ILRMA with Python can be easily implemented by slightly modifying the codes. A MATLAB code of Consistent ILRMA is also available online (<https://github.com/d-kitamura/ILRMA/blob/master/consistentILRMA.m>).

## 4 Experiments

In this section, we conducted two experiments using synthesized and real-recorded mixtures. The synthesized mixtures were produced by convoluting the impulse responses to dry audio sources, while the real-recorded mixtures were actually recorded by using a microphone array in an ordinary room with ambient noise.

### 4.1 BSS of synthesized mixtures

#### 4.1.1 Conditions

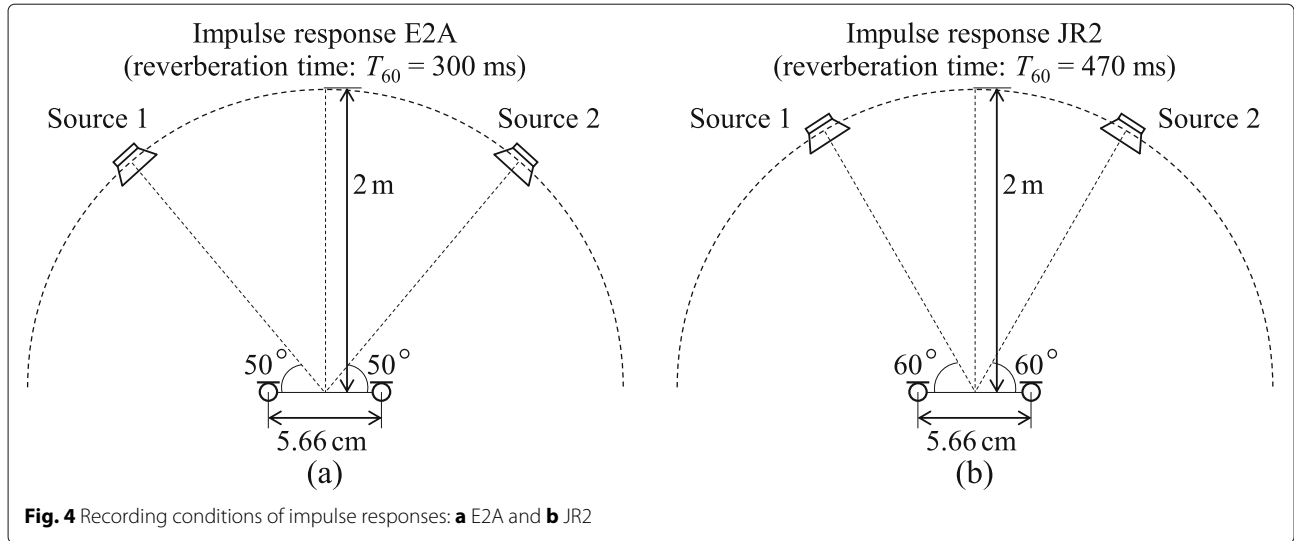
We conducted determined BSS experiments using synthesized music and speech mixtures with two sources and two microphones ( $N = M = 2$ ). The dry sources of music and speech signals, listed in Table 1, were respectively obtained from professionally produced music and underdetermined separation tasks provided as a part of SiSEC2011 [60]. They were convoluted with the impulse response E2A ( $T_{60} = 300$  ms) or JR2 ( $T_{60} = 470$  ms), obtained from the RWCP database [61], to simulate the multichannel observation signals. The recording conditions of these impulse responses are shown in Fig. 4.

In this experiment, we compared the performance of six methods: three conventional and three proposed. The conventional methods were the standard IVA [10], Consistent IVA [32], and standard ILRMA [11]. The proposed methods were Consistent IVA with iterative back projection (Consistent IVA+BP), Consistent ILRMA, and Consistent ILRMA with iterative back projection (Consistent ILRMA+BP). For all methods, the initial demixing matrix was set to an identity matrix. For the ILRMA-based

**Table 1** Music and speech dry sources obtained from SiSEC2011

Signal	Data name	Source (1/2)
Music 1	bearlin-roads	acoustic_guit_main/vocals
Music 2	bearlin-roads	piano/acoustic_guit_main
Music 3	bearlin-roads	piano/vocals
Music 4	another_dreamer-the_ones_we_love	guitar/vocals
Music 5	another_dreamer-the_ones_we_love	drums/guitar
Music 6	fort_minor-remember_the_name	violins_synth/vocals
Music 7	fort_minor-remember_the_name	vocals/drums
Music 8	tamy-que_pena_tanto_faz	guitar/vocals
Music 9	ultimate_nz_tour	guitar/synth
Music 10	ultimate_nz_tour	drums/vocals
Speech 1	dev1_female4	src_1/src_2
Speech 2	dev1_female4	src_1/src_4
Speech 3	dev1_female4	src_2/src_3
Speech 4	dev1_female4	src_2/src_4
Speech 5	dev1_female4	src_3/src_4
Speech 6	dev1_male4	src_1/src_2
Speech 7	dev1_male4	src_1/src_4
Speech 8	dev1_male4	src_2/src_3
Speech 9	dev1_male4	src_2/src_4
Speech 10	dev1_male4	src_3/src_4





methods, the nonnegative matrices  $T_n$  and  $V_n$  were initialized using uniformly distributed random values in the range (0, 1). Five trials were performed for each condition using different pseudorandom seeds. The number of bases for each source,  $K$ , was set to 10 for music mixtures and 2 for speech mixtures, where it was experimentally confirmed that these conditions provide the best performance for the conventional ILRMA [11]. To satisfy the perfect reconstruction condition (7), the inverse STFT was implemented by the canonical dual of the analysis window. For both Consistent IVA+BP and Consistent ILRMA+BP, the iterative back projection was applied, where the reference channel was set to  $m_{\text{ref}} = 1$ . Since the property of spectrogram consistency depends on the window length, shift length, and type of window function, various combinations of them were tested. The experimental conditions are summarized in Table 2.

For quantitative evaluation of the separation performance, we measured the source-to-distortion ratio (SDR), source-to-interference ratio (SIR), and source-to-artifact ratio (SAR). In a noiseless situation, SDR, SIR, and SAR are defined as follows [62]:

$$\text{SDR} = 10 \log_{10} \frac{\sum_l |s_t[l]|^2}{\sum_l |e_i[l] + e_a[l]|^2}, \quad (30)$$

$$\text{SIR} = 10 \log_{10} \frac{\sum_l |s_t[l]|^2}{\sum_l |e_i[l]|^2}, \quad (31)$$

$$\text{SAR} = 10 \log_{10} \frac{\sum_l |s_t[l] + e_i[l]|^2}{\sum_l |e_a[l]|^2}, \quad (32)$$

where  $s_t[l]$ ,  $e_i[l]$ , and  $e_a[l]$  are the  $l$ th sample of target signal, interference, and artificial components of the estimated signal, respectively, in the time domain. SIR

and SAR are used to quantify the amount of interference rejection and the absence of artificial distortion of the estimated signal, respectively. SDR is used to quantify the overall separation performance, as SDR is in good agreement with both SIR and SAR for determined BSS.

In this experiment, the energy of sources was not adjusted, i.e., the energy ratio of sources (source-to-source ratio) was automatically determined by the initial volume of the dry sources and the level of the impulse responses. That is, the source-to-source ratio of each mixture signal is different from the others. To equally evaluate the performances of different mixtures, we calculated SDR improvement ( $\Delta\text{SDR}$ ) and SIR improvement ( $\Delta\text{SIR}$ ) defined as:

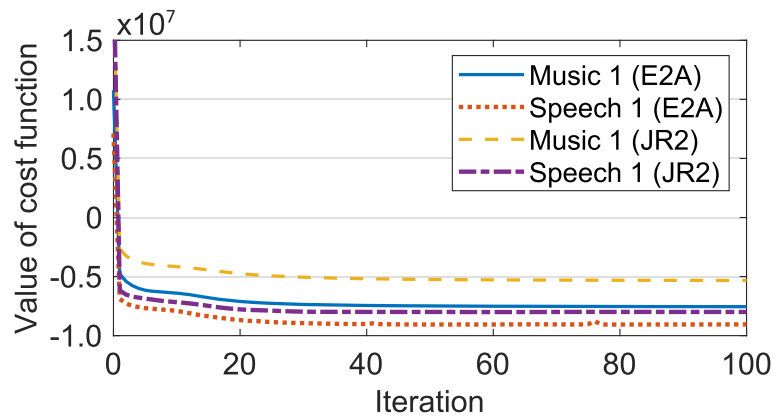
$$\Delta\text{SDR} = \text{SDR}_{\text{sep}} - \text{SDR}_{\text{input}}, \quad (33)$$

$$\Delta\text{SIR} = \text{SIR}_{\text{sep}} - \text{SIR}_{\text{input}}, \quad (34)$$

where  $\text{SDR}_{\text{sep}}$  and  $\text{SIR}_{\text{sep}}$  are the SDR and SIR of the separated signal, and  $\text{SDR}_{\text{input}}$  and  $\text{SIR}_{\text{input}}$  are the SDR and SIR of the initial mixture signal input to the BSS methods. Note that SAR improvement cannot be defined because its value of the signal without artificial processing cannot be defined ( $\text{SAR}_{\text{input}} = \infty$ ).

**Table 2** Experimental conditions

Window function	Hann/Hamming/Blackman window
Window length	64, 128, 256, 512, 768, 1024 ms
Window shift length	1/16, 1/8, 1/4, 1/2 of window length
Number of bases $K$ for each source in ILRMA	10 for music signals and 2 for speech signals
Number of iterations	100



**Fig. 5** Values of negative log-likelihood function (18) of Consistent ILRMA+BP (window length, 256 ms; shift length, 32 ms)

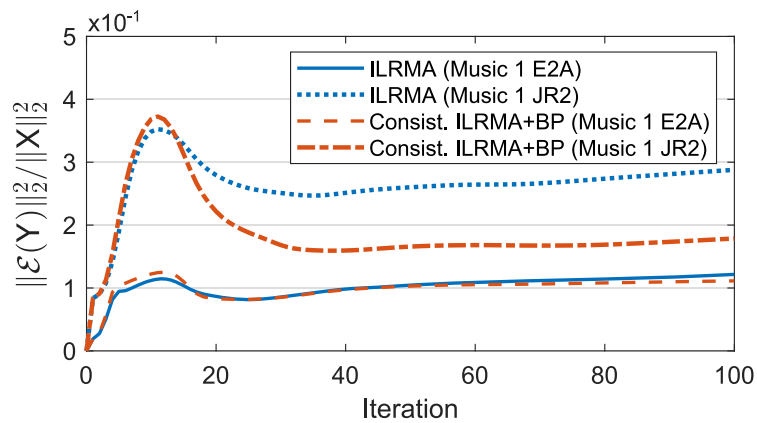
#### 4.1.2 Results and discussions

Figure 5 shows examples of the value of the negative log-likelihood function (18) of Consistent ILRMA+BP. Although the algorithmic convergence of the proposed method has not been theoretically justified because of the additional projection (25), we experimentally confirmed a smooth decrease of the cost function. We also confirmed that such behavior was common for the other experimental conditions and mixtures. This result indicates that the additional procedure in the proposed method does not have a harmful effect on the behavior of the overall algorithm.

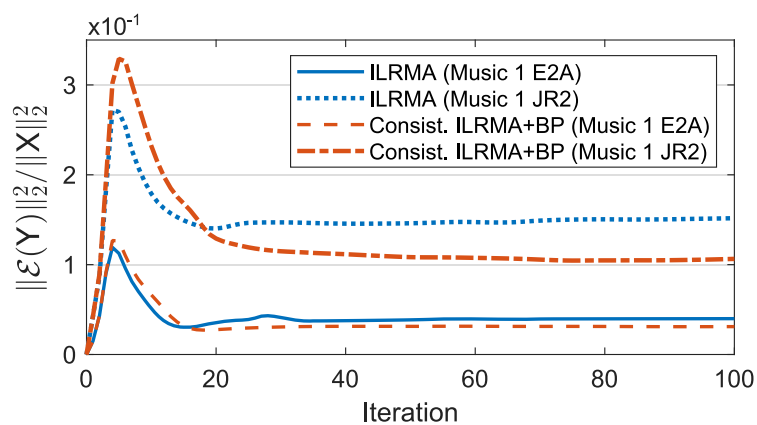
Figures 6 and 7 show examples of the energy of the inconsistent components (14) of standard ILRMA and Consistent ILRMA+BP. The energy was normalized by that of the initial spectrograms in order to align the vertical axis. Note that the energy of inconsistency components is not directly related to the degree of permutation misalignment or the separation performance. These figures are shown to confirm whether the proposed algorithm can properly reduce the degree of inconsistency. These values are completely zero when the separated spectrograms are consistent, and hence, those at the 0th iteration (the leftmost values) are zero because no processing is performed at that point. By iterating the algorithms, this energy rapidly increased because the demixing matrix for each frequency independently tried to process and separate the signals. However, the normalized energy tended toward some specific values after several iterations. We confirmed that the converged values of Consistent ILRMA+BP were always lower than those of standard ILRMA. This result indicates that Consistent ILRMA+BP reduces the amount of the inconsistent components and tries to make the separated spectrogram more consistent. In addition, similar to Fig. 5, the algorithmic stability of Consistent ILRMA+BP can be confirmed from Figs. 6 and 7.

Figures 8 and 9 summarize the SDR improvements for the music mixtures and speech mixtures, respectively. The window function was the Hann window. Each box contains 50 results (i.e., 5 pseudorandom seeds  $\times$  10 mixtures in Table 1), where  $\Delta$ SDRs of the two separated sources in each mixture were averaged. The central lines of the box plots indicate the median, and the bottom and top edges of the box indicate the 25th and 75th percentiles, respectively. Each row corresponds to the same window length, while each column corresponds to the same shift length. As we conducted the experiment for six window lengths, four shift lengths, and two impulse responses, each figure consists of  $6 \times 4 \times 2$  subfigures. In each subfigure, six boxes are shown to illustrate the results of (1) IVA, (2) Consistent IVA, (3) Consistent IVA+BP, (4) ILRMA, (5) Consistent ILRMA, and (6) Consistent ILRMA+BP. Since the tendency of the results was the same as Figs. 8 and 9, we provide the SDR improvements for the other windows (Hamming and Blackman) in the Appendix. The SIR improvement and SAR are also given in the Appendix.

Since IVA and ILRMA assume the instantaneous mixing model (11) for each frequency in the time-frequency domain, the window length should be long relative to the reverberation time to achieve accurate separation. At the same time, too long a window degrades the separation performance of IVA and ILRMA, as discussed in [30]. This is because capturing the source activity and spectral patterns becomes difficult for IVA and ILRMA as the time resolution of the spectrograms becomes low due to a long window. The robustness of IVA and ILRMA is also deteriorated by a long window because the effective number of time segments is decreased. This trade-off of the separation performance caused by window length in STFT can be easily confirmed from the results for both music (Fig. 8) and speech (Fig. 9) mixtures, which is consistent with the results in [30]. As shown in the figures, the performance was poor for the shorter windows ( $\leq 128$  ms), and



(a) Window length: 256 ms, shift length: 32 ms



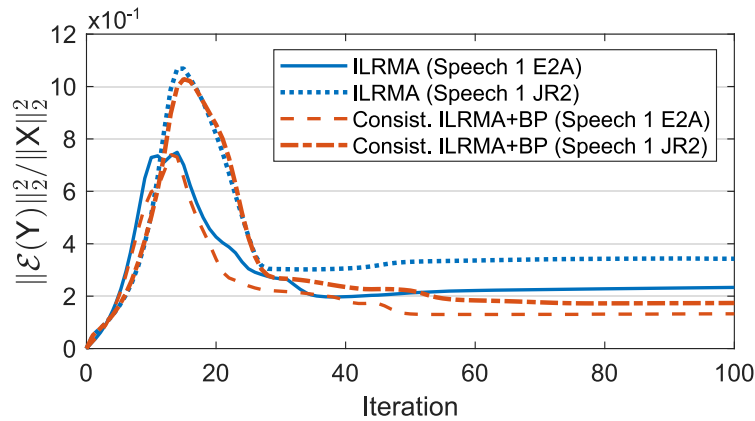
(b) Window length: 1024 ms, shift length: 512 ms

**Fig. 6** Examples of normalized energy of inconsistent components ( $\|\mathcal{E}(Y)\|_2^2/\|X\|_2^2$ ) of ILRMA and Consistent ILRMA+BP for music 1: **a** 256-ms-long window and 32-ms shifting and **b** 1024-ms-long window and 512-ms shifting, where  $X = [X_1, X_2]$ ,  $Y = [Y_1, Y_2]$ , and  $\mathcal{E}(\cdot)$  is in (14)

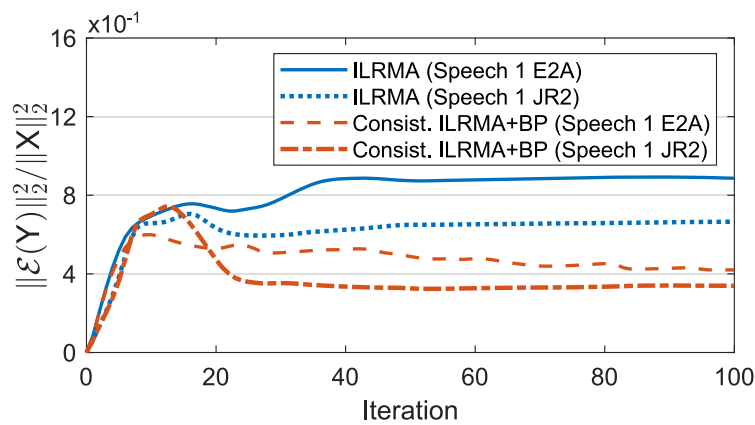
the performance for the longer windows ( $\geq 768$  ms) was more varied than that of the shorter ones. The window length best suited for these conditions (combinations of source signals and impulse responses) seems to be around 256 ms or 512 ms. While the maximum achievable performance becomes higher as the window length becomes longer due to the mixing model (11), these results indicate that the source modeling becomes difficult for both IVA and ILRMA when the window length is too long. This trade-off should be important for discussing the results further.

By comparing the performances of the conventional (IVA, Consistent IVA, and ILRMA) and proposed (Consistent IVA+BP, Consistent ILRMA, and Consistent ILRMA+BP) methods, we can see that the proposed methods tend to outperform the conventional ones. Some comparisons are made as follows:

- **Conventional and proposed IVAs.** The proposed Consistent IVA+BP performed better than the conventional IVAs (IVA and Consistent IVA) in Figs. 8b and 9b when the window length was sufficiently long ( $\geq 256$  ms). In those cases, the conventional Consistent IVA resulted in a worse performance than IVA, which indicates that just using spectrogram consistency cannot improve the performance of IVA. This demonstrates the importance of the iterative back projection when spectrogram consistency is considered within determined BSS.
- **Conventional and proposed ILRMAs.** The proposed Consistent ILRMA without BP performed comparably to the conventional ILRMA. In Figs. 8a and 9a, Consistent ILRMA performed better than ILRMA when the window length was long ( $\geq 768$  ms). In contrast, in Figs. 8b and 9b, Consistent



(a) Window length: 256 ms, shift length: 32 ms



(b) Window length: 1024 ms, shift length: 512 ms

**Fig. 7** Examples of normalized energy of inconsistent components ( $\|\mathcal{E}(Y)\|_2^2/\|X\|_2^2$ ) of ILRMA and Consistent ILRMA+BP for speech 1: **a** 256-ms-long window and 32-ms shifting and **b** 1024-ms-long window and 512-ms shifting, where  $X = [X_1, X_2]$ ,  $Y = [Y_1, Y_2]$ , and  $\mathcal{E}(\cdot)$  is in (14)

tent ILRMA performed worse than ILRMA. This is presumably because the scale ambiguity prevented the spectrogram consistency from working properly. By incorporating iterative back projection into Consistent ILRMA, the proposed Consistent ILRMA+BP performed better than the conventional ILRMA. In the best situation (the top-left subfigure of Fig. 8), Consistent ILRMA+BP performed 8 dB better than ILRMA by bringing out the potential of spectrogram consistency in determined BSS.

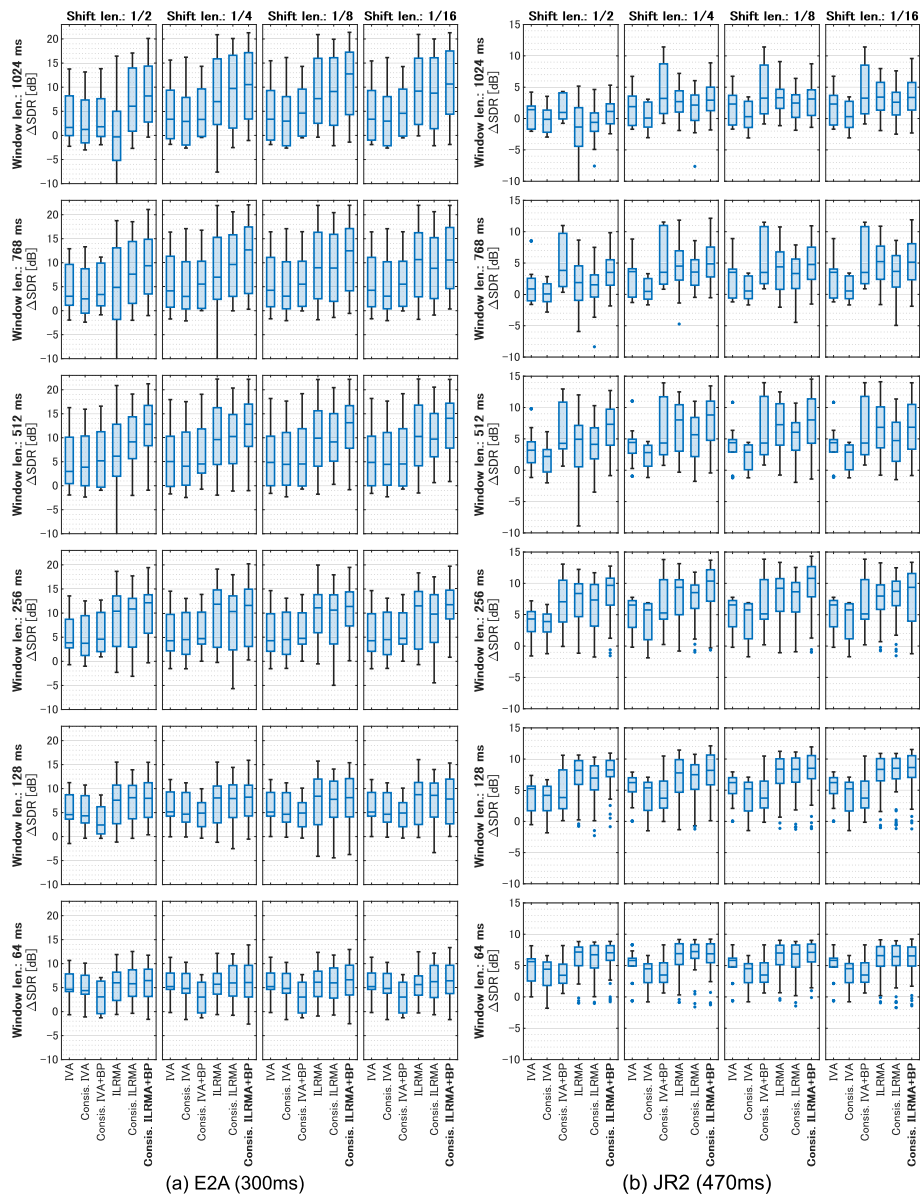
To further explain the experimental results, some notable tendencies are summarized as follows:

- **Short window.** When the window length was short (64 ms), all methods performed similarly in terms of  $\Delta$ SDR. This is because the achievable performance was already limited by the window length that

was shorter than the reverberation time. This result contradicted our expectation before performing the experiment. Since enforcing the consistency spreads the frequency components based on the main-lobe of the window function, we expected that the ability to solve the permutation problem would be higher when the window length was shorter because of the wider main-lobe. In reality, we found that the spectrogram consistency could assist IVA and ILRMA except for the cases where the window length was short ( $\leq 128$  ms in this experiment) compared to the reverberation time.

- **Large window shift.** When the shift length was 1/2 of the window length, the performance of ILRMA significantly dropped compared to smaller shift lengths (1/4, 1/8, and 1/16), especially when the window length was long (e.g., 1024 ms). This is pre-

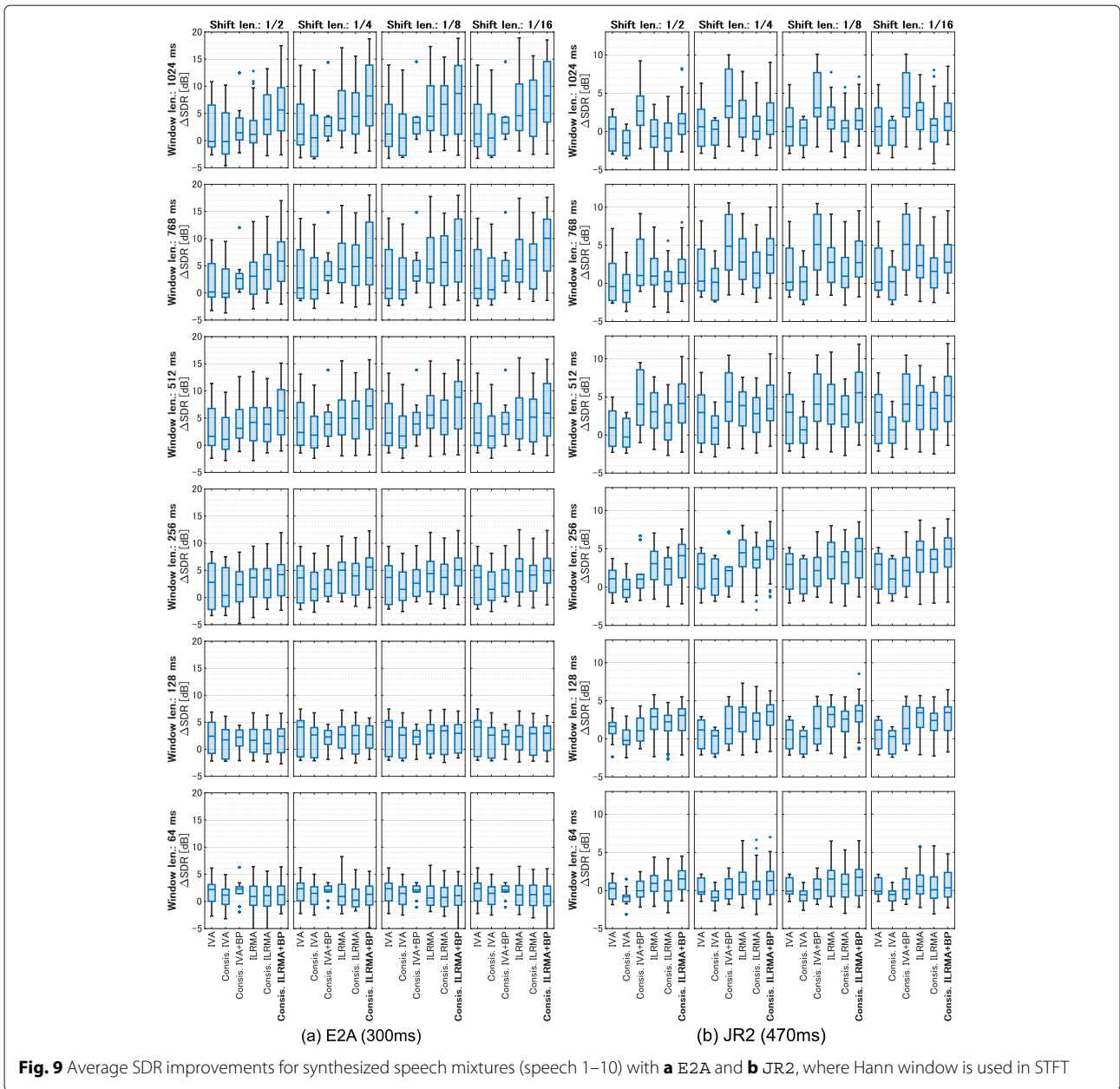




**Fig. 8** Average SDR improvements for synthesized music mixtures (music 1–10) with **a** E2A and **b** JR2, where Hann window is used in STFT

sumably because the number of time segments was small, i.e., NMF in ILRMA failed to model the source signals from the given amount of data. In addition, for a larger window, distinguishing spectral patterns of the sources became difficult for ILRMA due to the time-directional blurring effect caused by a longer window. Such performance degradation was alleviated for Consistent ILRMA+BP. This might be because the smoothing process of the inverse STFT provides some additional information for the source modeling from the adjacent bins.

- Length of boxes.** When the length of the box of ILRMA was long, as in Figs. 8a and 9a, Consistent ILRMA+BP was able to improve the performance. Conversely, when the length of the box of ILRMA was short, as in Figs. 8b and 9b, Consistent ILRMA+BP was only able to slightly improve the performance. Note that the vertical axes are different. This result indicates that the achievable performance decided by the mixing model (11) limits the improvement obtained by spectrogram consistency. Since consistency is the characteristic of a spectro-



gram, it cannot manage the mixing process. The demixing-filter update of ILRMA, which is the same for the conventional and proposed methods, manages the mixing process. Hence, when the mixing model has a mismatch with the observed condition, there is less room for spectrogram consistency to improve the performance.

- Improvement by consistency.** The proposed method tended to achieve a good performance when the conventional ILRMA also worked well, e.g., Figs. 8a and 9a. This tendency indicates that the spectro-

gram consistency effectively promotes the separation when the estimated source  $Y_n$  accurately approaches the original source  $S_n$  during the optimization, as  $S_n$  is naturally a consistent spectrogram. This is the reason we feel that the consistency can be an assistant of the frequency-domain BSS. An important aspect is that the source model (e.g., NMF in ILRMA) actually informs the separation cue, and the spectrogram consistency enhances the separation performance when the source modeling functions correctly.

## 4.2 BSS of real-recorded mixtures

### 4.2.1 Conditions

Next, we evaluated the conventional and proposed methods using live-recorded music and speech mixtures obtained from underdetermined separation tasks in SiSEC2011 [60], where only two sources were mixed to make the BSS problem determined ( $M = N = 2$ ). The signals used in this experiment are listed in Table 3. The reverberation time of these signals was 250 ms, and the microphone spacing was 1 m (see [60]). Since these source signals were actually recorded using a microphone array in an ordinary room with ambient noise, the observed signals are more realistic compared to those in Section 4.1.

For simplicity, in this experiment, we used STFT with a fixed condition, the 512-ms-long Hann window with 1/4 shifting. The experimental conditions other than the window were the same as those in Section 4.1.1.

### 4.2.2 Results and discussion

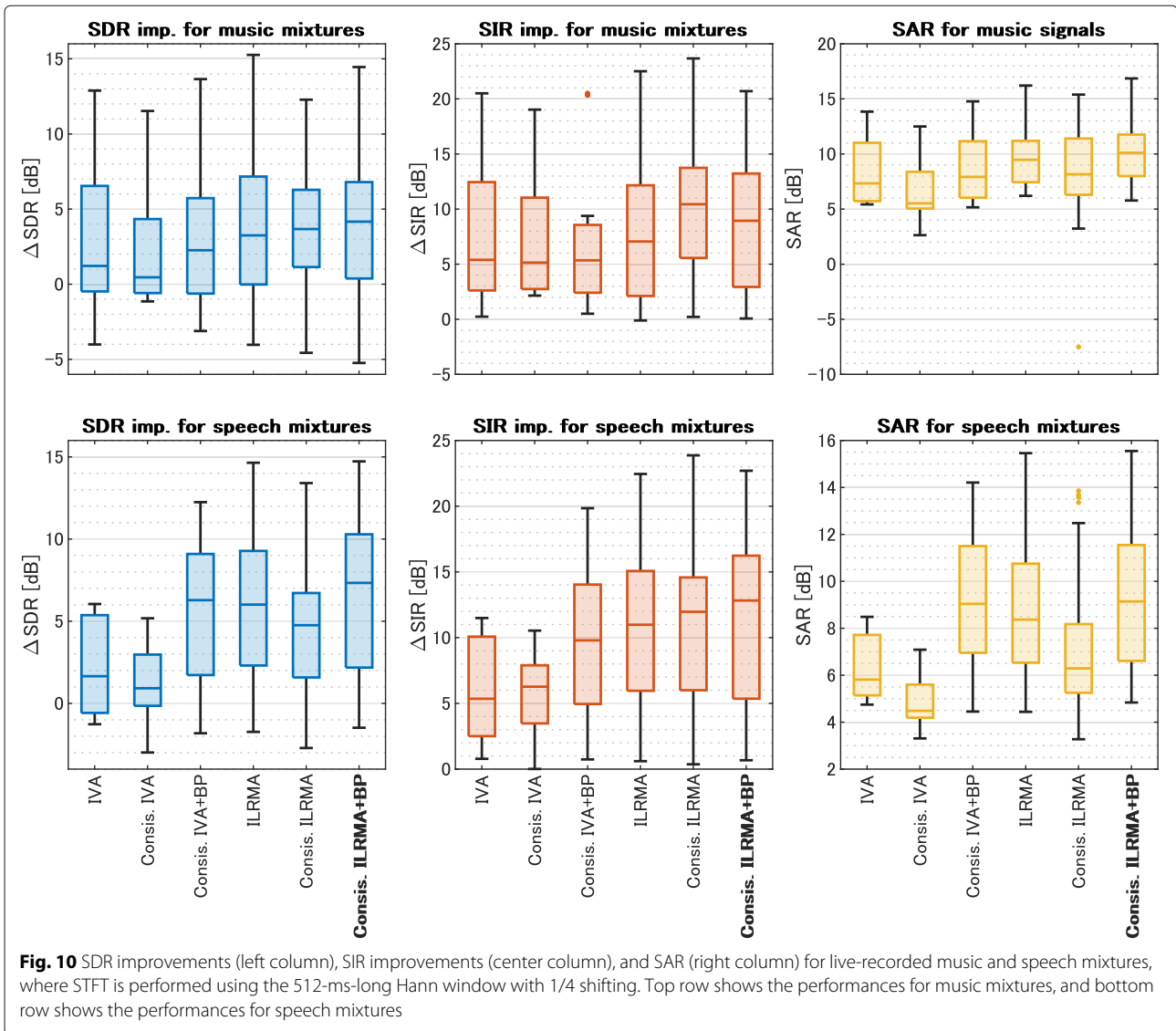
Figure 10 shows the results of live-recorded music and speech mixtures. The absolute scores were lower than those for the synthesized mixtures discussed in Section 4.1.2 due to the existence of ambient noise. Still, we can confirm the improvements of the proposed Consistent IVA+BP and Consistent ILRMA+BP compared to the conventional IVA and ILRMA, respectively, for both the music (upper row) and speech (lower row) mixtures. In particular, Consistent IVA+BP improved more than 4 dB over IVA in terms of the median of the  $\Delta$ SDR of speech mixtures. Consistent ILRMA+BP achieved the highest performance in terms of the median of the SDR improvement for both music and speech mixtures. These results confirm that the combination of spectrogram consistency and iterative back projection can assist the separation of determined BSS for a more realistic situation.

## 5 Conclusion

In this paper, we have proposed a new variant of the state-of-the-art determined BSS algorithm called Consistent ILRMA. It utilizes the smoothing effect of the inverse STFT in order to assist the separation and enhance the performance. Experimental results showed that the proposed method can improve the separation performance when the window length is sufficiently large ( $\geq 256$  ms in the experimental condition of this paper). These results demonstrate the potential of considering spectrogram consistency within the state-of-the-art determined BSS algorithm. In addition, we experimentally confirmed the importance of iterative back projection for considering spectrogram consistency within determined BSS. It should be possible to construct a new source model in consideration of the spectrogram consistency, which can pave the way for the next direction of research on determined BSS.

**Table 3** Live-recorded music and speech signals obtained from SiSEC2011

Signal	Source (1/2)
Music 1	dev1_nodrums_liverec_250ms_1m_sim_1/dev1_nodrums_liverec_250ms_1m_sim_2
Music 2	dev1_nodrums_liverec_250ms_1m_sim_1/dev1_nodrums_liverec_250ms_1m_sim_3
Music 3	dev1_nodrums_liverec_250ms_1m_sim_2/dev1_nodrums_liverec_250ms_1m_sim_3
Music 4	dev1_wdrums_liverec_250ms_1m_sim_1/dev1_nodrums_liverec_250ms_1m_sim_2
Music 5	dev1_wdrums_liverec_250ms_1m_sim_1/dev1_nodrums_liverec_250ms_1m_sim_3
Music 6	dev1_wdrums_liverec_250ms_1m_sim_2/dev1_nodrums_liverec_250ms_1m_sim_3
Music 7	dev2_nodrums_liverec_250ms_1m_sim_1/dev1_nodrums_liverec_250ms_1m_sim_2
Music 8	dev2_nodrums_liverec_250ms_1m_sim_1/dev1_nodrums_liverec_250ms_1m_sim_3
Music 9	dev2_nodrums_liverec_250ms_1m_sim_2/dev1_nodrums_liverec_250ms_1m_sim_3
Music 10	dev2_wdrums_liverec_250ms_1m_sim_1/dev1_nodrums_liverec_250ms_1m_sim_2
Music 11	dev2_wdrums_liverec_250ms_1m_sim_1/dev1_nodrums_liverec_250ms_1m_sim_3
Music 12	dev2_wdrums_liverec_250ms_1m_sim_2/dev1_nodrums_liverec_250ms_1m_sim_3
Speech 1	dev1_female4_liverec_250ms_1m_sim_1/dev1_female4_liverec_250ms_1m_sim_2
Speech 2	dev1_female4_liverec_250ms_1m_sim_3/dev1_female4_liverec_250ms_1m_sim_4
Speech 3	dev1_male4_liverec_250ms_1m_sim_1/dev1_male4_liverec_250ms_1m_sim_2
Speech 4	dev1_male4_liverec_250ms_1m_sim_3/dev1_male4_liverec_250ms_1m_sim_4
Speech 5	dev1_female4_liverec_250ms_1m_sim_1/dev1_male4_liverec_250ms_1m_sim_2
Speech 6	dev2_female4_liverec_250ms_1m_sim_3/dev1_male4_liverec_250ms_1m_sim_4
Speech 7	dev2_female4_liverec_250ms_1m_sim_1/dev1_female4_liverec_250ms_1m_sim_2
Speech 8	dev2_female4_liverec_250ms_1m_sim_3/dev1_female4_liverec_250ms_1m_sim_4
Speech 9	dev2_male4_liverec_250ms_1m_sim_1/dev1_male4_liverec_250ms_1m_sim_2
Speech 10	dev2_male4_liverec_250ms_1m_sim_3/dev1_male4_liverec_250ms_1m_sim_4
Speech 11	dev2_male4_liverec_250ms_1m_sim_1/dev1_female4_liverec_250ms_1m_sim_2
Speech 12	dev2_male4_liverec_250ms_1m_sim_3/dev1_female4_liverec_250ms_1m_sim_4



## Appendix

### Independence between real and imaginary parts of spectrogram

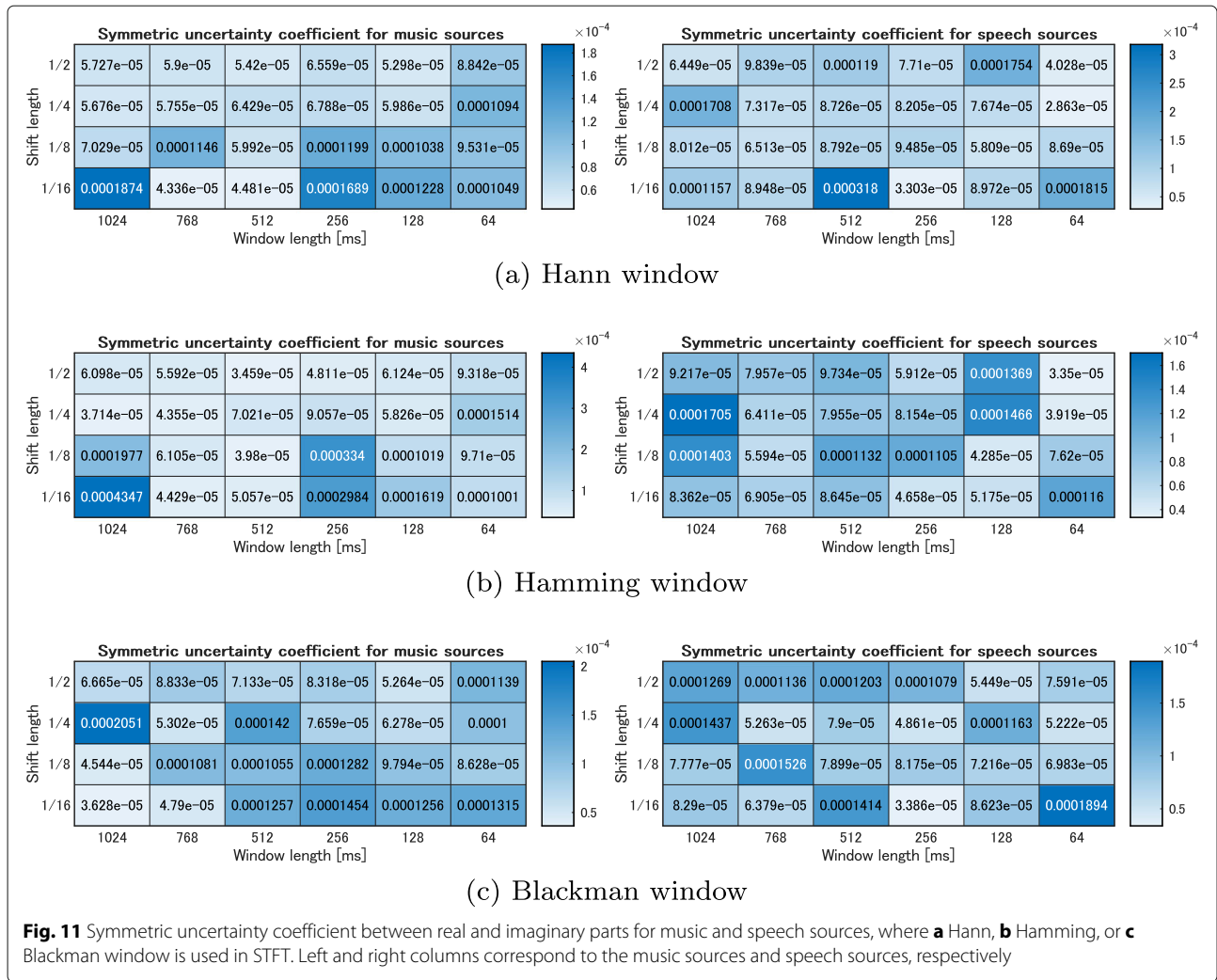
The source generative model (15) assumes that the real and imaginary parts of a source in the time-frequency domain are mutually independent because the generative model has a zero-mean and circularly symmetric shape in the complex plane. The independence between real and imaginary parts or amplitude and phase has been investigated, but its validity may depend on the parameters of STFT. Independence can be measured by a symmetric uncertainty coefficient [63–65]:

$$C(q_1, q_2) = 2 \frac{H(q_1) + H(q_2) - H(q_1, q_2)}{H(q_1) + H(q_2)}, \quad (35)$$

where  $q_1$  and  $q_2$  are random variables,  $H(q_1)$  and  $H(q_2)$  are their entropy, and  $H(q_1, q_2)$  is the joint entropy of  $q_1$  and  $q_2$ . Since the numerator of (35) corresponds to the mutual information of  $q_1$  and  $q_2$ , the symmetric uncertainty coefficient can be interpreted as normalized mutual information. When  $q_1$  and  $q_2$  are mutually independent, (35) becomes zero. In contrast, when  $q_1$  and  $q_2$  are completely dependent, (35) becomes one.

We calculated the symmetric uncertainty coefficient (35) between the real and imaginary parts of a time-frequency bin obtained by applying STFT to music or speech sources. Let  $s$  be a complex-valued time-frequency bin of a source (the indexes of frequency and time are omitted here). The independence between the real and imaginary parts can be measured by  $C(\text{Re}(s), \text{Im}(s))$ ,





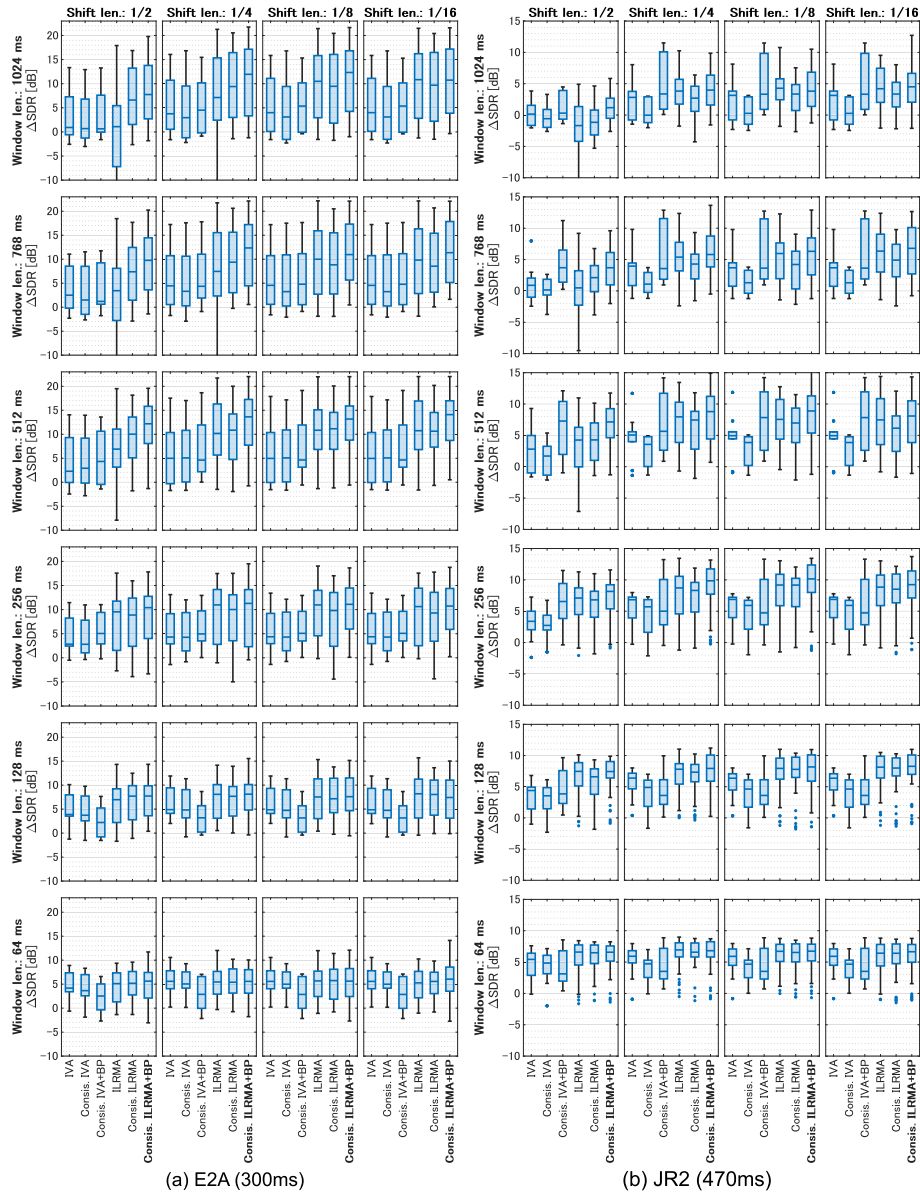
where  $\text{Re}(\cdot)$  and  $\text{Im}(\cdot)$  return the real and imaginary parts of an input complex value, respectively. Here,  $H(\text{Re}(s))$ ,  $H(\text{Im}(s))$ , and  $H(\text{Re}(s), \text{Im}(s))$  were approximately obtained by calculating the histograms of  $\text{Re}(s)$  and  $\text{Im}(s)$ . The number of bins in the histograms was set to 10,000. We used the dry sources listed in Table 1: 15 music (instrumental) and eight speech sources. The parameters of STFT were the same as those in Section 4.

Figure 11 shows the symmetric uncertainty coefficients averaged over all bins and sources. Their values  $C(\text{Re}(s), \text{Im}(s))$  were almost zero for all STFT conditions and source types (music or speech), and thus, the assumption of independence between real and imaginary parts is valid for music and speech sources. This fact leads

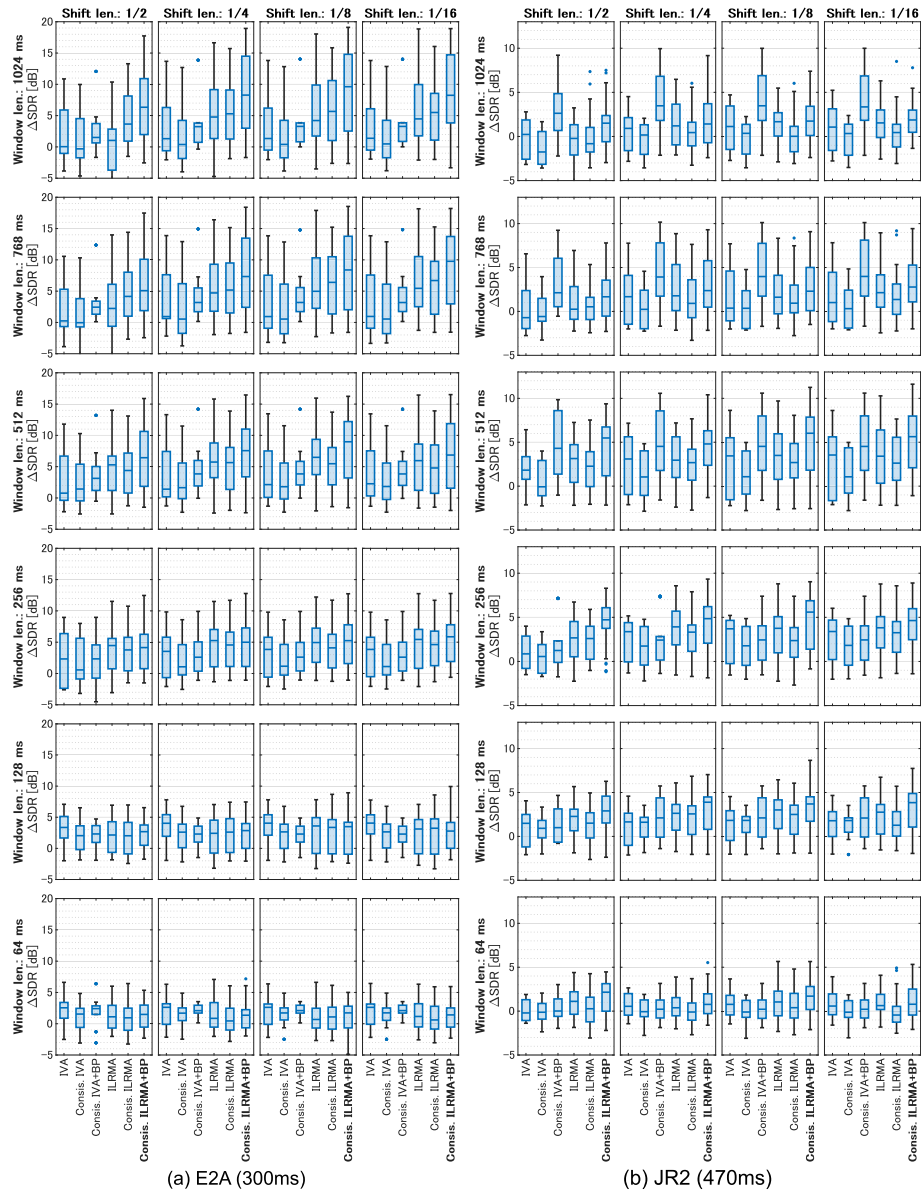
to the generative model assumed in ILRMA. Note that those symmetric uncertainty coefficients validated the independence of real and imaginary parts at each time-frequency bin. That is, the inter-bin relation is not considered here. The proposed method captures such inter-bin relations imposed by the spectrogram consistency, which is not apparent in these bin-wise assessments of independence.

**Additional experimental results for synthesized mixtures**  
 Figures 12–15, 16–21, and 22–27 show the SDR improvements, SIR improvements, and SAR, respectively, for synthesized music and speech mixtures. These figures correspond to the results and discussions in Section 4.1.2.



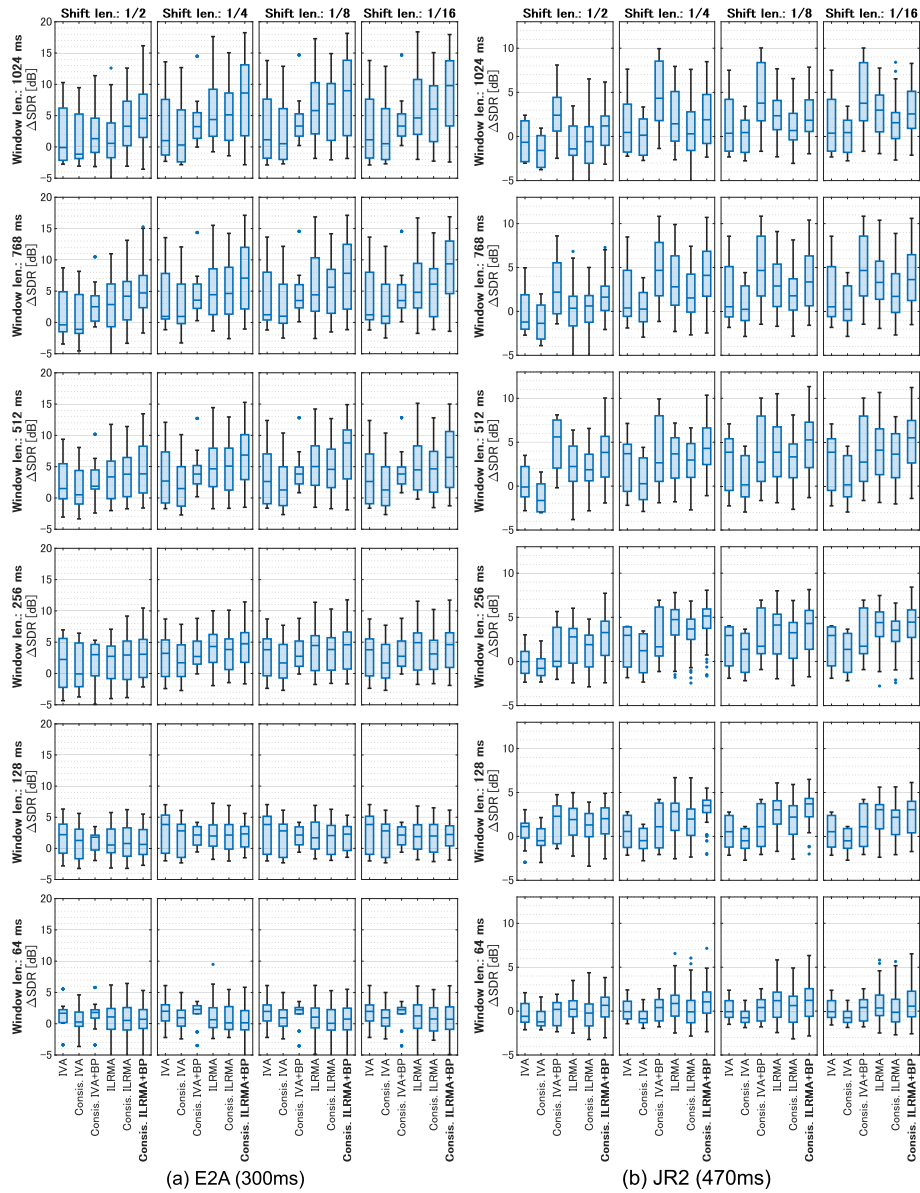


**Fig. 13** Average SDR improvements for synthesized music mixtures (music 1–10) with **a** E2A and **b** JR2, where Blackman window is used in STFT

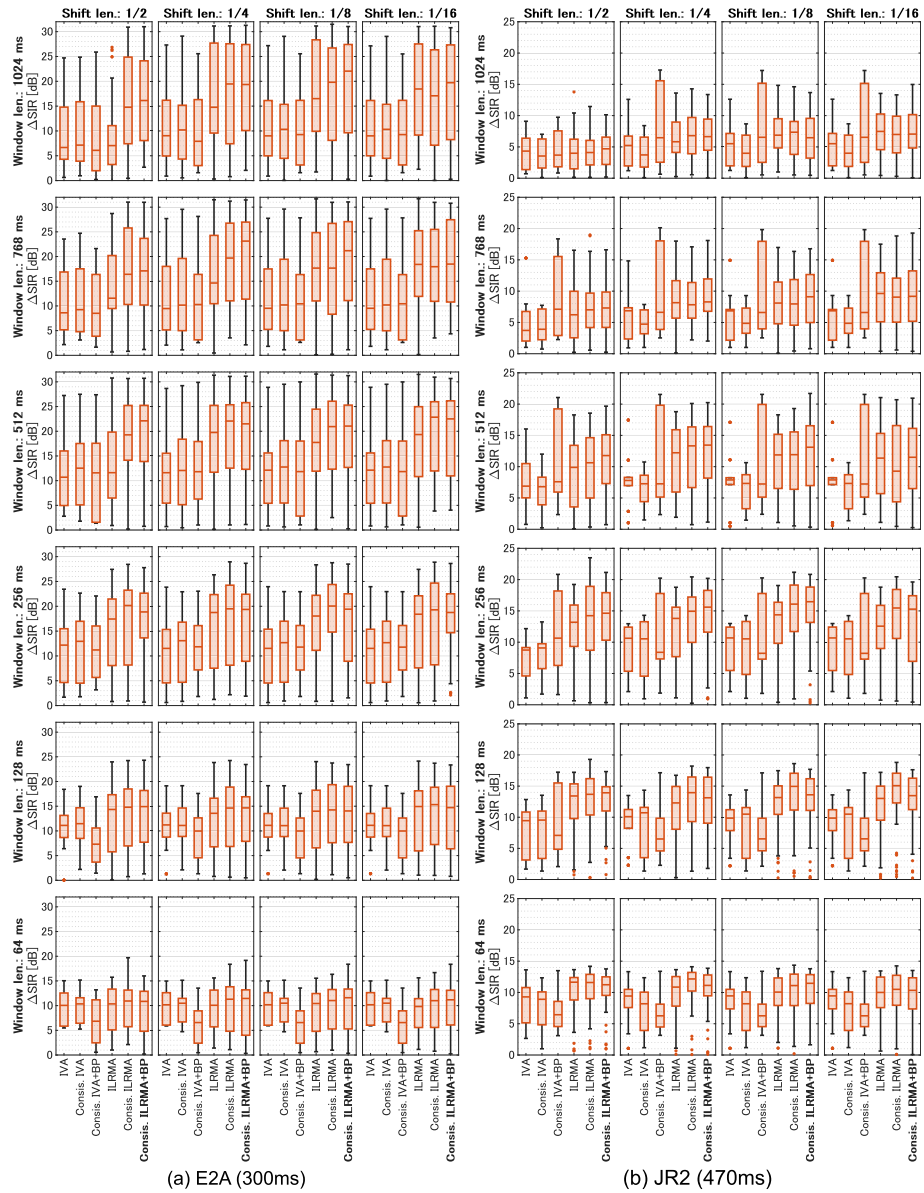


**Fig. 14** Average SDR improvements for synthesized speech mixtures (speech 1–10) with **a** E2A and **b** JR2, where Hamming window is used in STFT

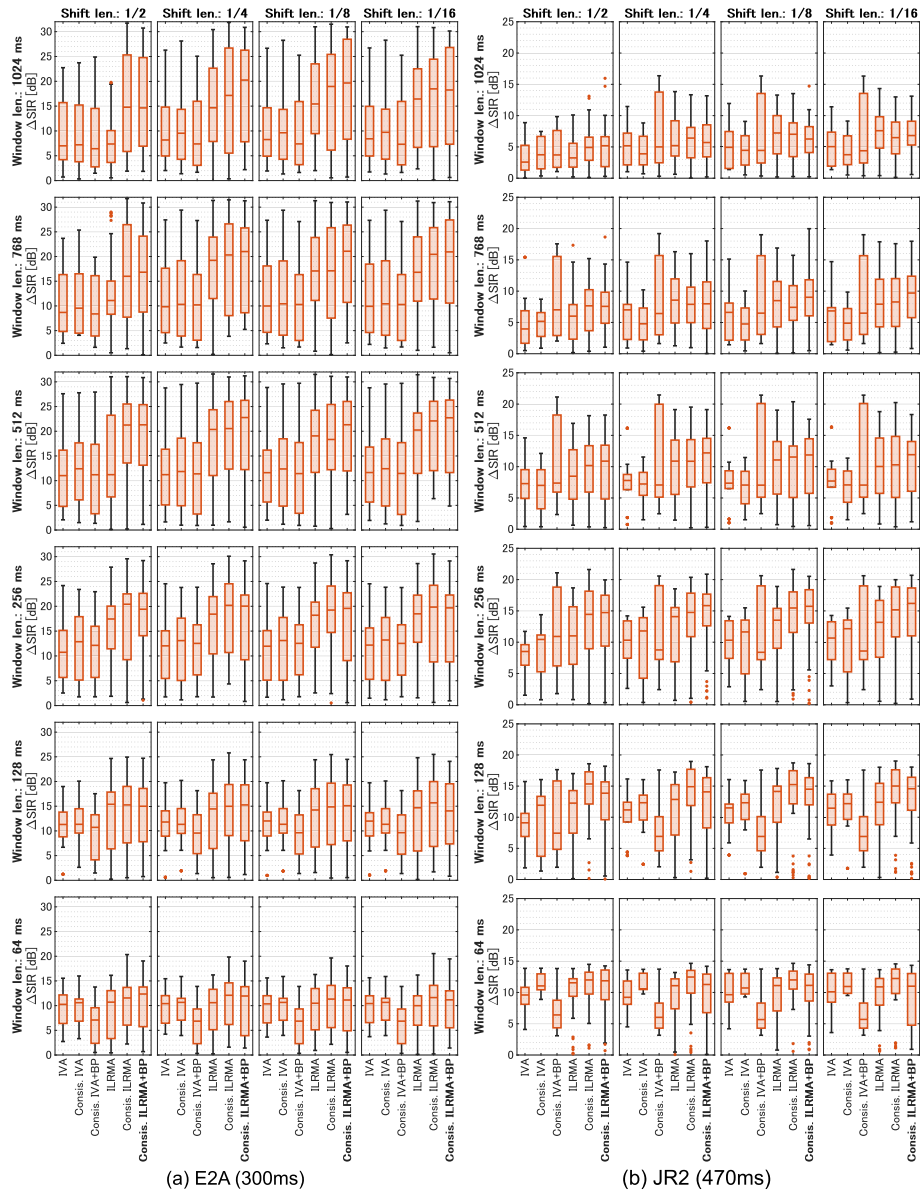




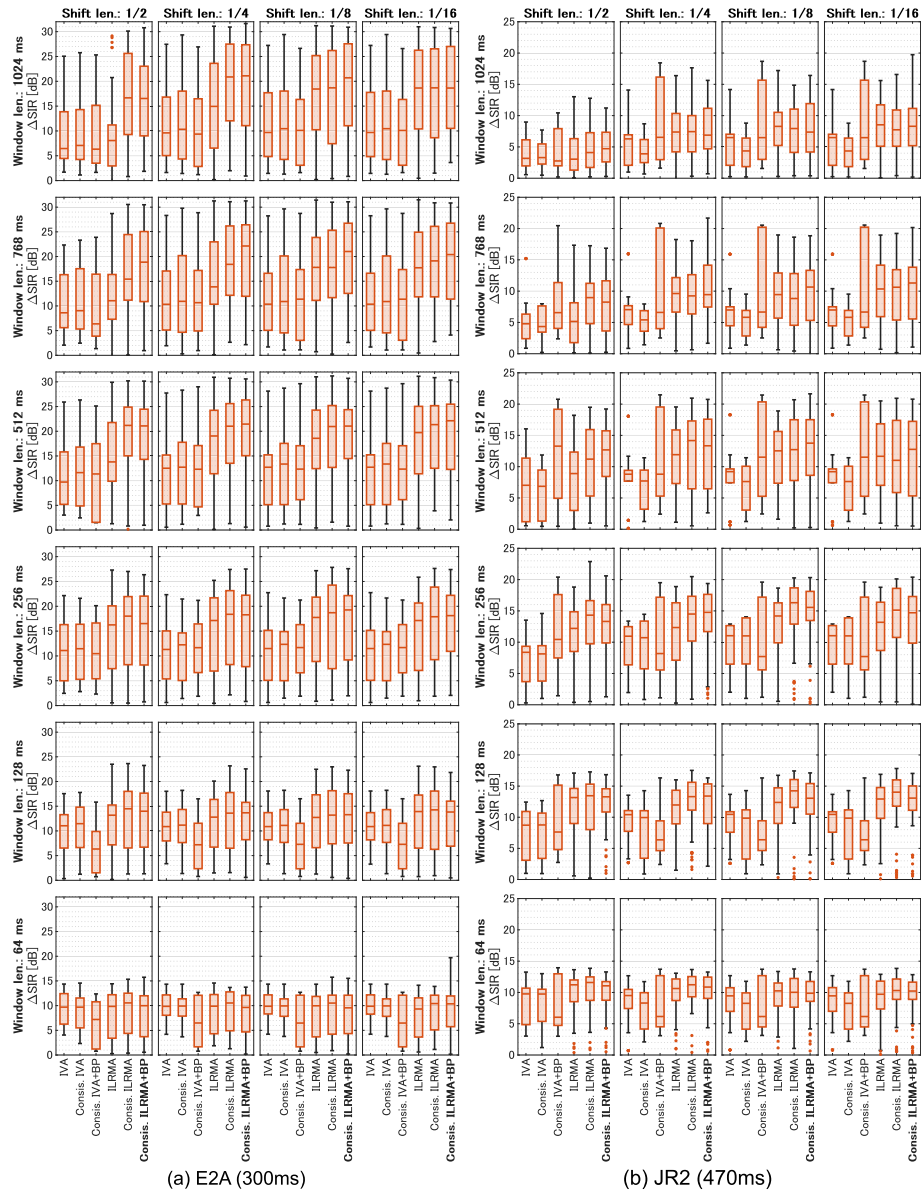
**Fig. 15** Average SDR improvements for synthesized speech mixtures (speech 1–10) with **a** E2A and **b** JR2, where Blackman window is used in STFT



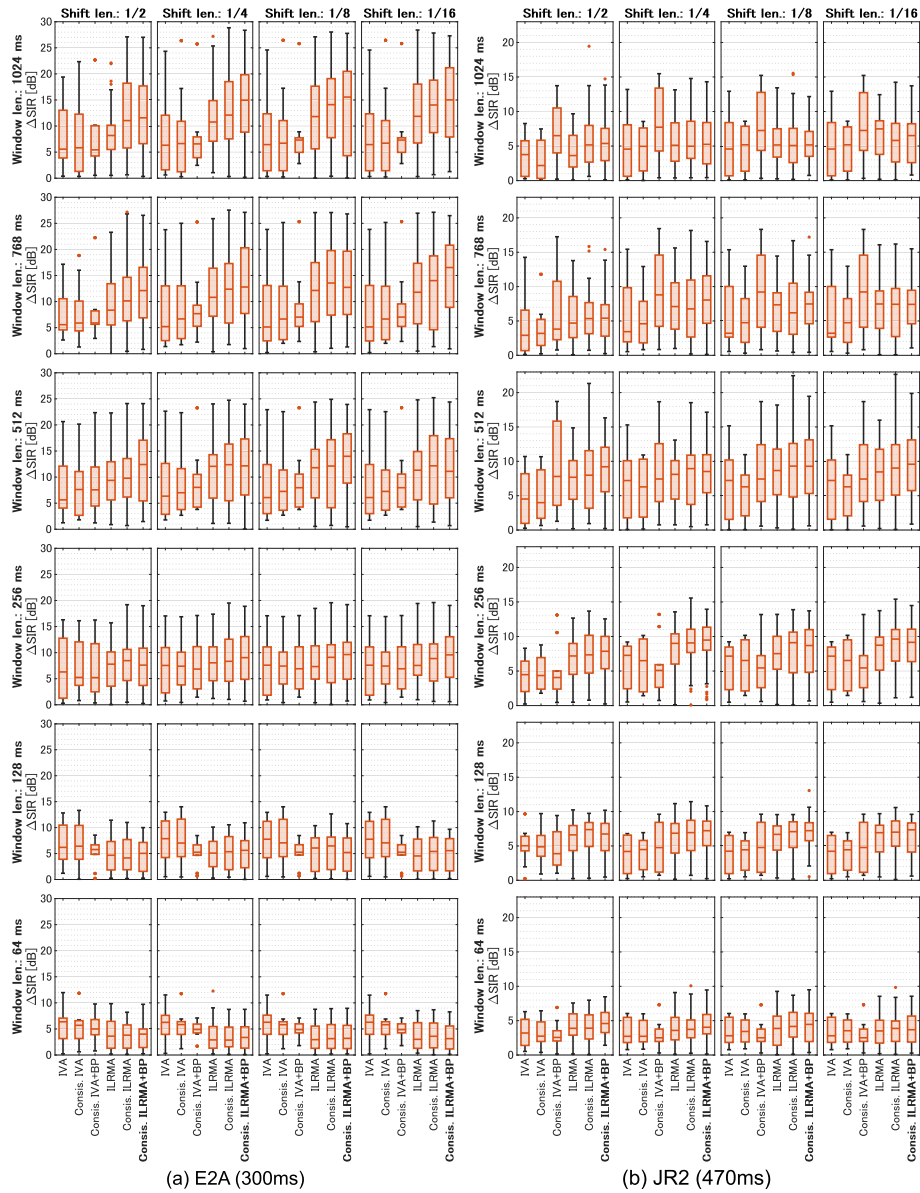
**Fig. 16** Average SIR improvements for synthesized music mixtures (music 1–10) with **a** E2A and **b** JR2, where Hann window is used in STFT



**Fig. 17** Average SIR improvements for synthesized music mixtures (music 1–10) with **a** E2A and **b** JR2, where Hamming window is used in STFT

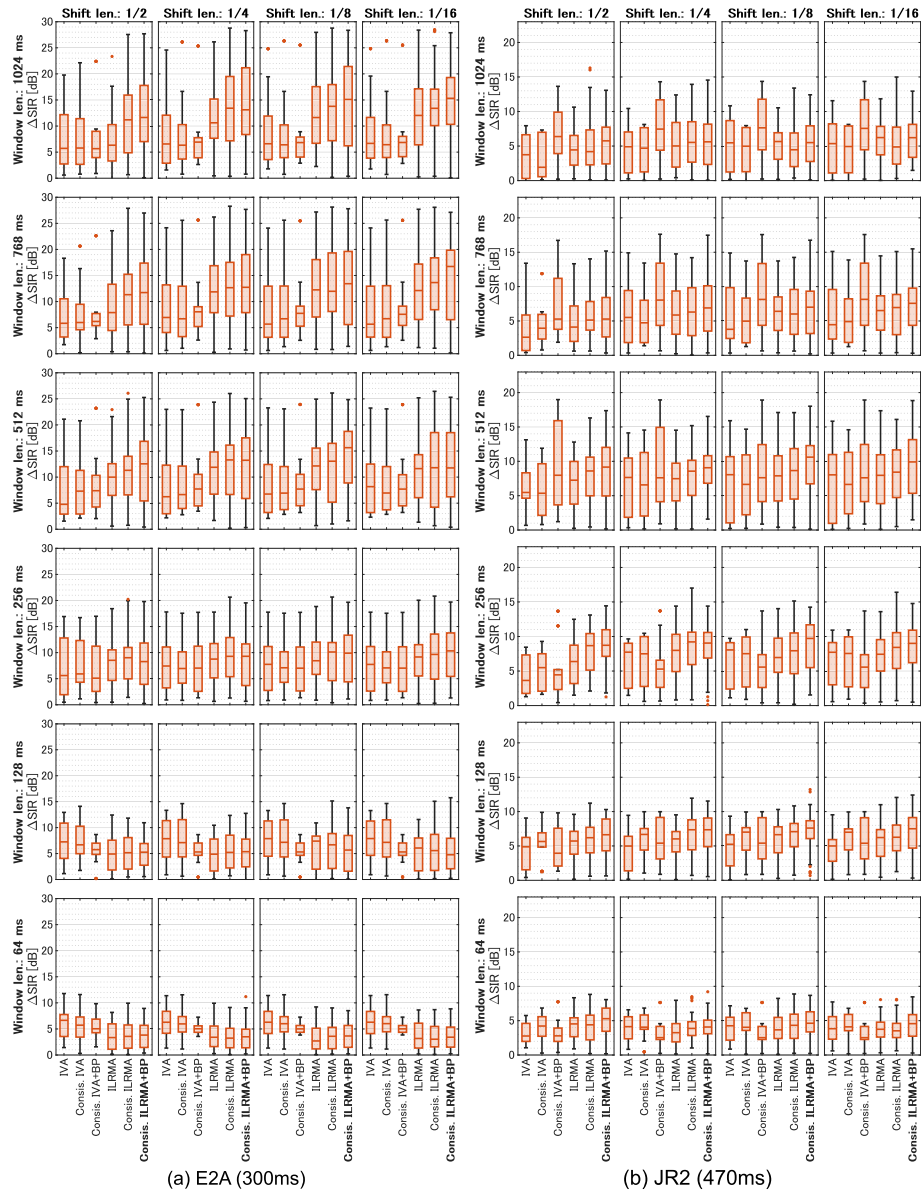


**Fig. 18** Average SIR improvements for synthesized music mixtures (music 1–10) with **a** E2A and **b** JR2, where Blackman window is used in STFT

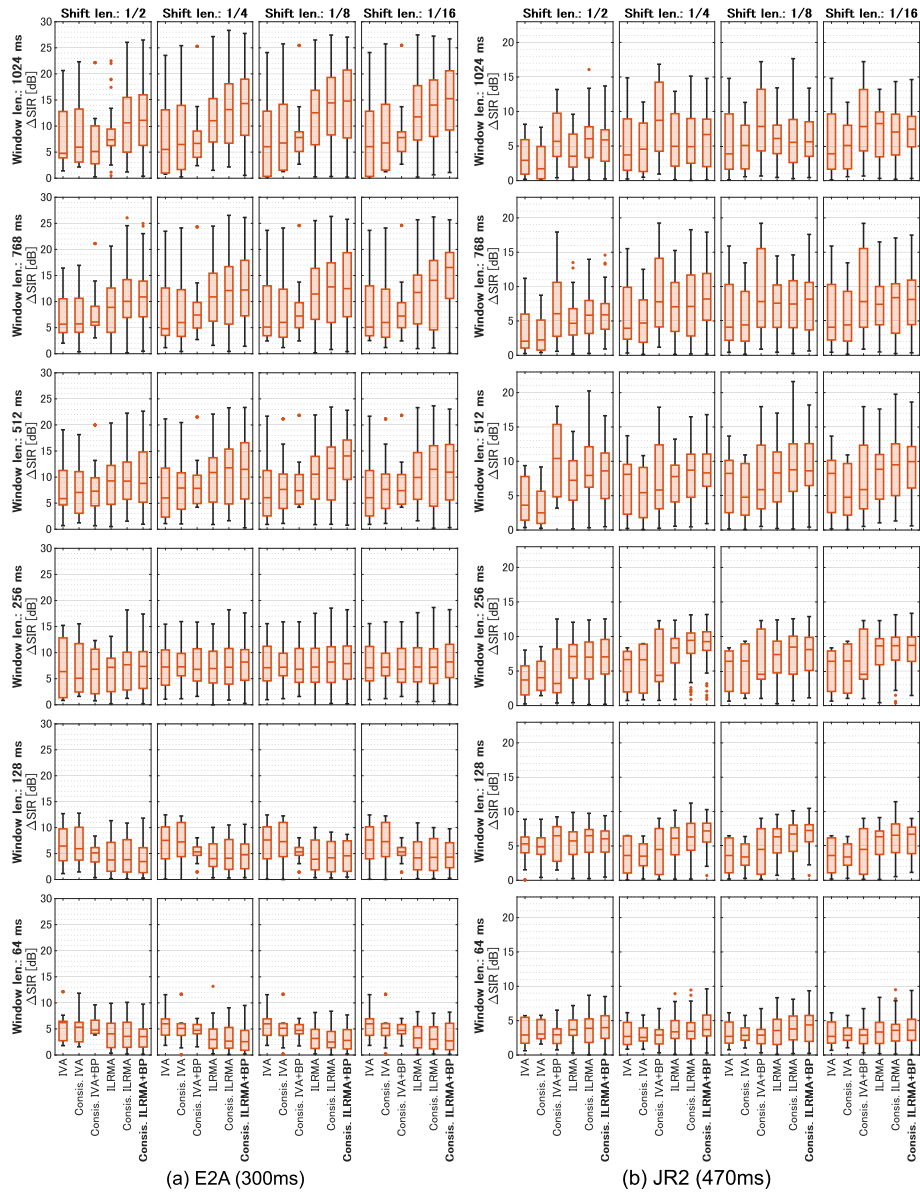


**Fig. 19** Average SIR improvements for synthesized speech mixtures (speech 1–10) with **a** E2A and **b** JR2, where Hann window is used in STFT

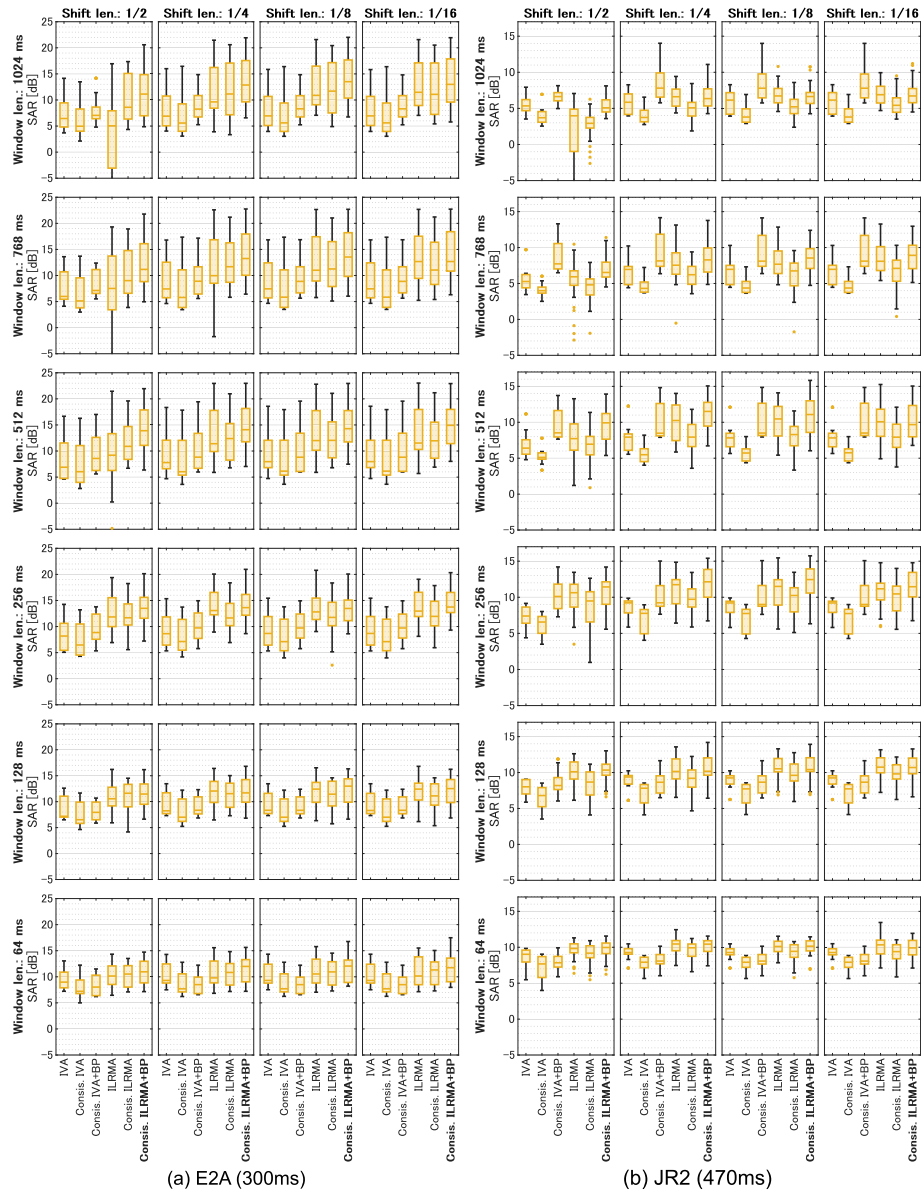




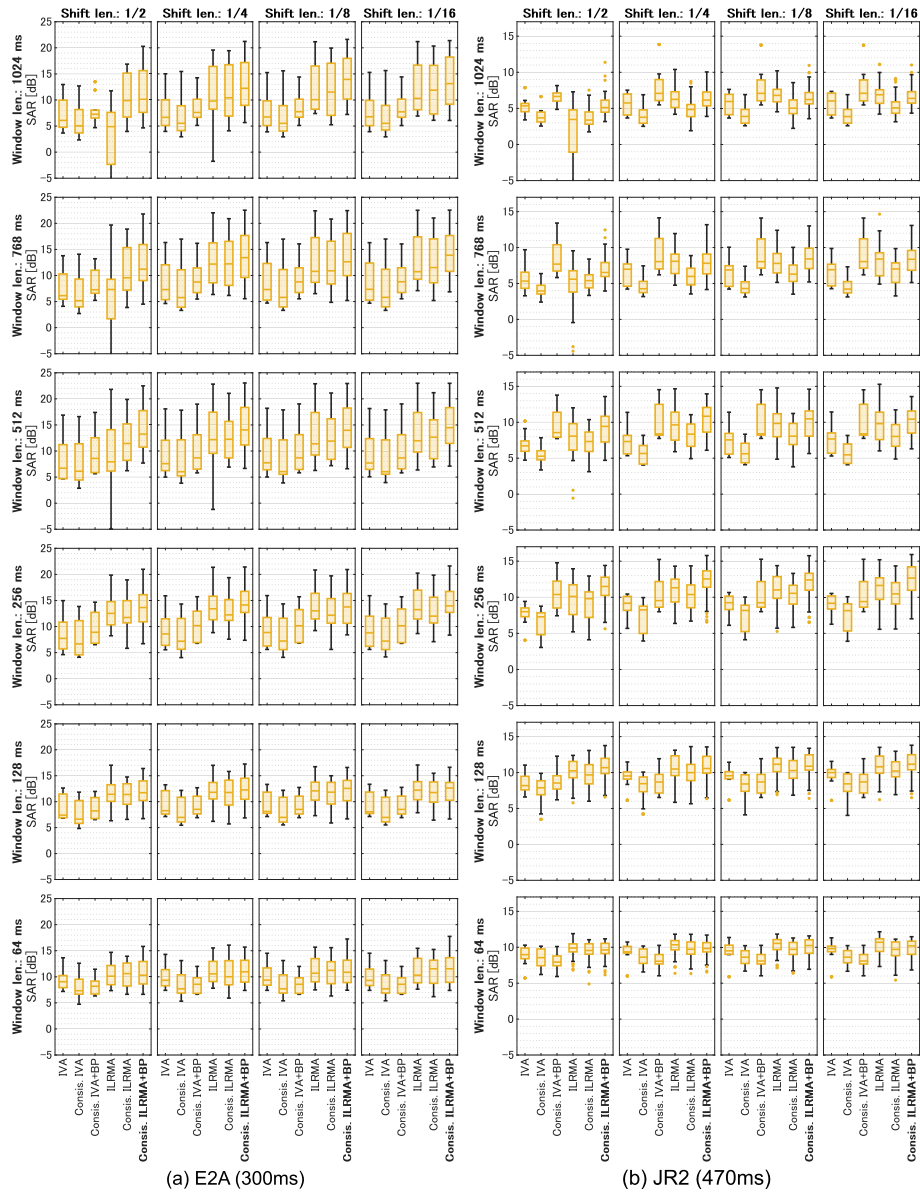
**Fig. 20** Average SIR improvements for synthesized speech mixtures (speech 1–10) with **a** E2A and **b** JR2, where Hamming window is used in STFT



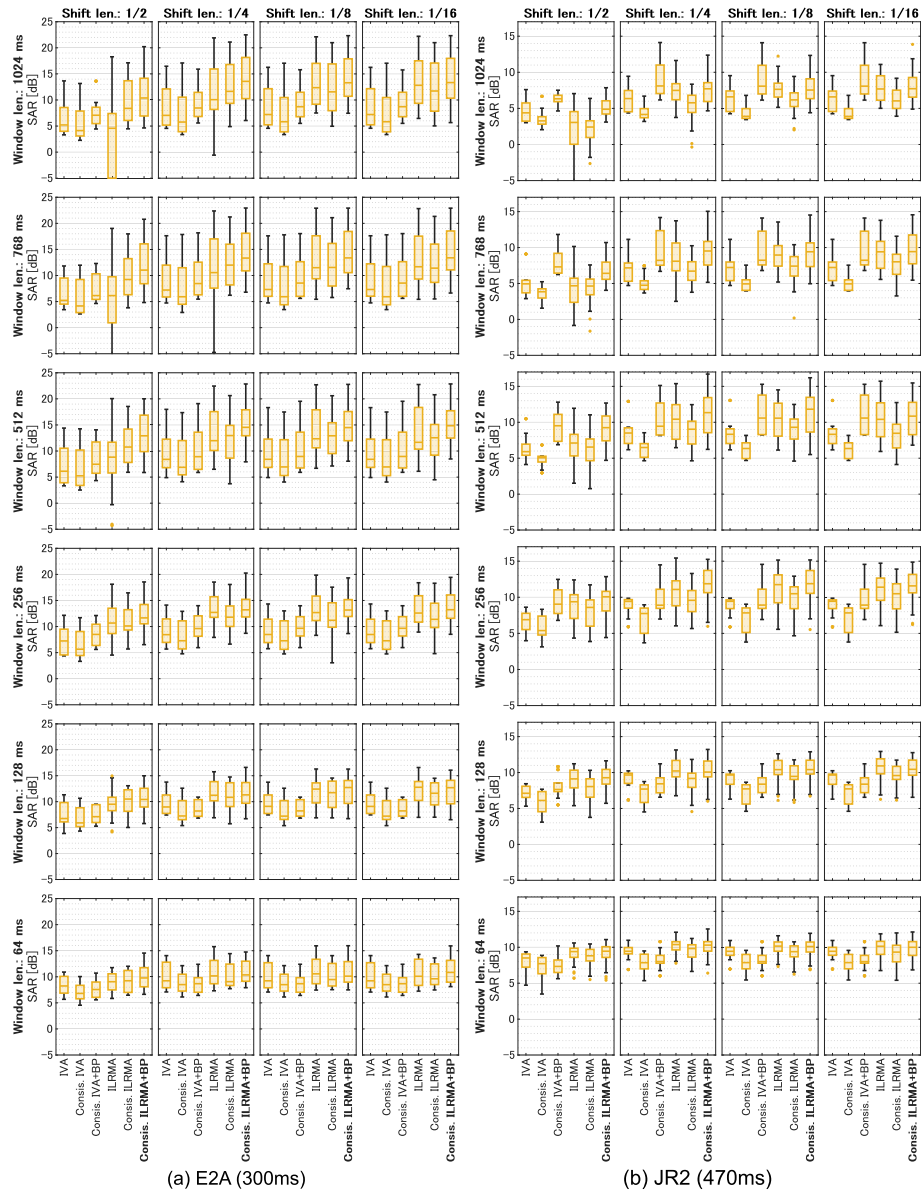
**Fig. 21** Average SIR improvements for synthesized speech mixtures (speech 1–10) with **a** E2A and **b** JR2, where Blackman window is used in STFT



**Fig. 22** Average SAR for synthesized music mixtures (music 1–10) with **a** E2A and **b** JR2, where Hann window is used in STFT

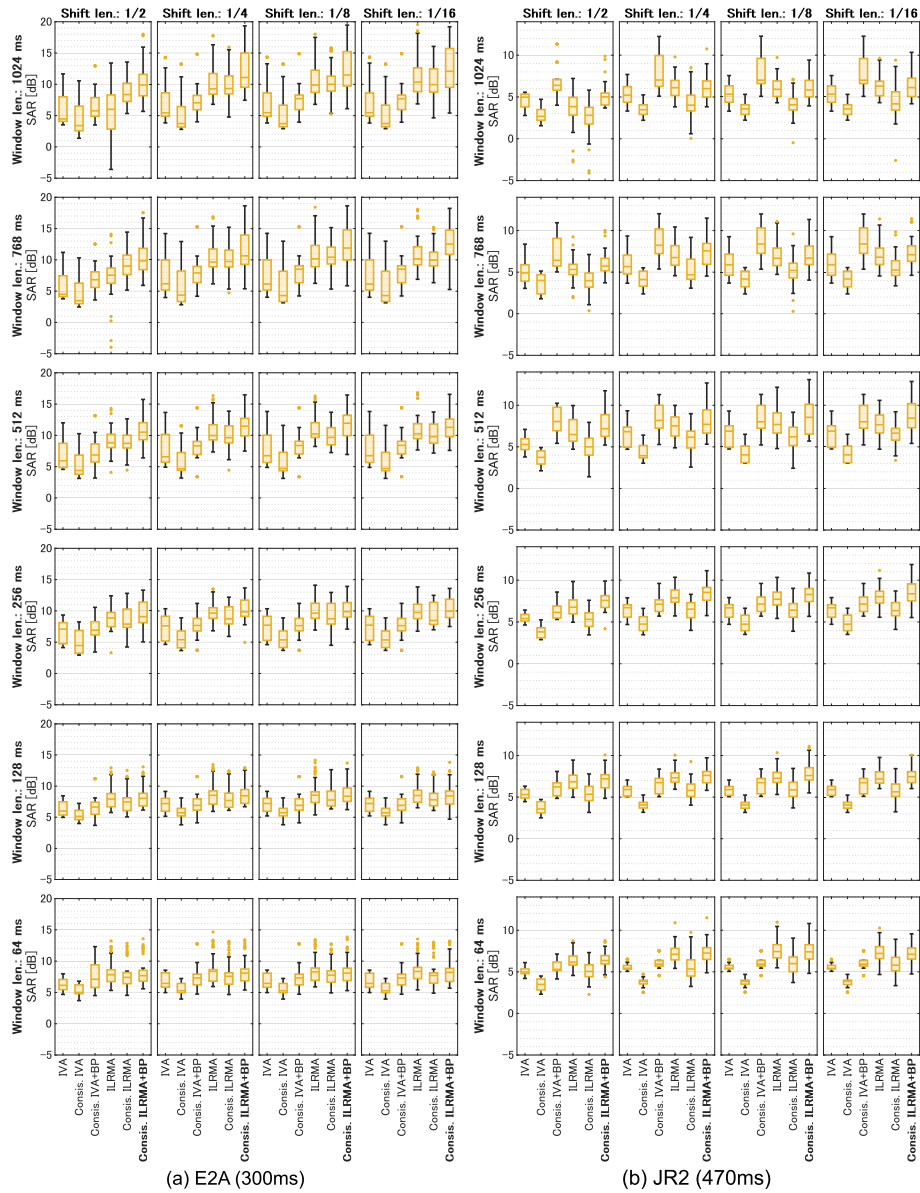


**Fig. 23** Average SAR for synthesized music mixtures (music 1–10) with **a** E2A and **b** JR2, where Hamming window is used in STFT



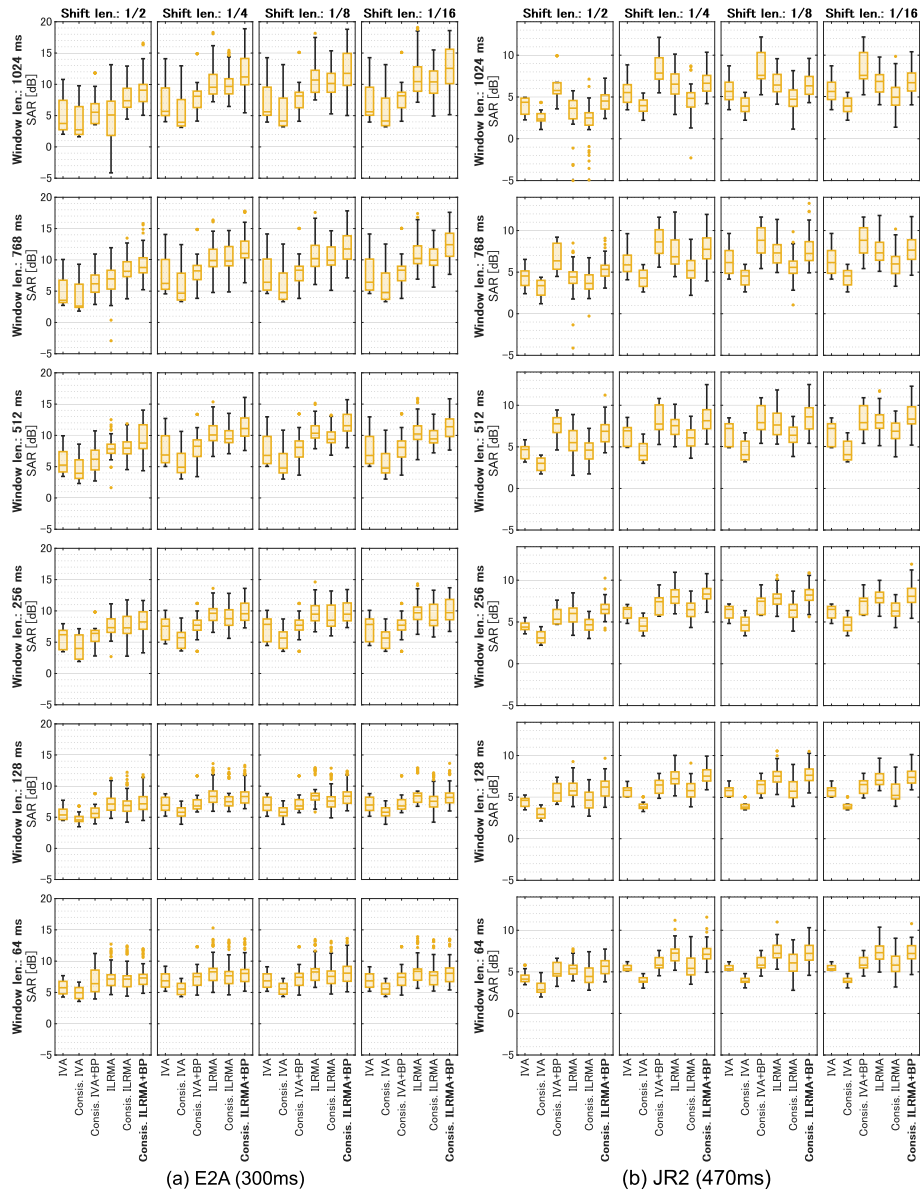
**Fig. 24** Average SAR for synthesized music mixtures (music 1–10) with **a** E2A and **b** JR2, where Blackman window is used in STFT





**Fig. 25** Average SAR for synthesized speech mixtures (speech 1–10) with **a** E2A and **b** JR2, where Hann window is used in STFT





**Fig. 27** Average SAR for synthesized speech mixtures (speech 1–10) with **a** E2A and **b** JR2, where Blackman window is used in STFT

### Abbreviations

BSS: Blind source separation; ICA: Independent component analysis; STFT: Short-time Fourier transform; FDICA: Frequency-domain independent component analysis; IVA: Independent vector analysis; ILRMA: Independent low-rank matrix analysis; NMF: Nonnegative matrix factorization; IS-NMF: Nonnegative matrix factorization based on the Itakura–Saito divergence; SDR: Source-to-distortion ratio; SIR: Source-to-interference ratio; SAR: Source-to-artifact ratio

### Acknowledgements

The authors would like to thank Nao Toshima for his support on the experiment. Also, the authors would like to thank the anonymous reviewers for their valuable comments and suggestions that helped improve the quality of this manuscript.

### Authors' contributions

DK derived the algorithm, performed the experiment, drafted the manuscript for initial submission, and revised the manuscript. KY proposed the main idea, gave advice, mainly wrote the manuscript for initial submission, and corrected the draft of revised manuscript. The authors read and approved the final manuscript.

### Funding

This work was partially supported by JSPS Grants-in-Aid for Scientific Research 19K20306 and 19H01116.

### Availability of data and materials

The datasets used for the experiments in this paper are openly available: SiSEC 2011 (<http://sisec2011.wiki.irisa.fr/>) and RWCP-SSD (<http://research.nii.ac.jp/src/en/RWCP-SSD.html>). Our MATLAB implementation of the proposed method is also openly available at the following site: <https://github.com/d-kitamura/ILRMA/blob/master/consistentLRMA.m>

### Competing interests

The authors declare that they have no competing interests.

### Author details

<sup>1</sup>National Institute of Technology, Kagawa College, 355 Chokushi, Takamatsu, Kagawa, 761-8058, Japan. <sup>2</sup>Waseda University, 3-4-1 Okubo, Shinjuku-ku, Tokyo, 169-8555, Japan.

Received: 2 July 2020 Accepted: 29 October 2020

Published online: 16 November 2020

### References

- P. Comon, Independent component analysis, a new concept? *Signal Process.* **36**(3), 287–314 (1994)
- P. Smaragdis, Blind separation of convolved mixtures in the frequency domain. *Neurocomputing.* **22**, 21–34 (1998)
- S. Kurita, H. Saruwatari, S. Kajita, K. Takeda, F. Itakura, in *Proc. ICASSP*. Evaluation of blind signal separation method using directivity pattern under reverberant conditions, vol. 5 (IEEE, 2000), pp. 3140–3143
- N. Murata, S. Ikeda, A. Ziehe, An approach to blind source separation based on temporal structure of speech signals. *Neurocomputing.* **41**(1–4), 1–24 (2001)
- H. Saruwatari, T. Kawamura, T. Nishikawa, A. Lee, K. Shikano, Blind source separation based on a fast-convergence algorithm combining ICA and beamforming. *IEEE Trans. ASLP.* **14**(2), 666–678 (2006)
- H. Sawada, R. Mukai, S. Araki, S. Makino, A robust and precise method for solving the permutation problem of frequency-domain blind source separation. *IEEE Trans. SAP.* **12**(5), 530–538 (2004)
- A. Hiroe, in *Proc. ICA*. Solution of permutation problem in frequency domain ICA using multivariate probability density functions (Springer, Berlin, Heidelberg, 2006), pp. 601–608
- T. Kim, T. Eltoft, T.-W. Lee, in *Proc. ICA*. Independent vector analysis: an extension of ICA to multivariate components (Springer, Berlin, Heidelberg, 2006), pp. 165–172
- T. Kim, H.T. Attias, S.-Y. Lee, T.-W. Lee, Blind source separation exploiting higher-order frequency dependencies. *IEEE Trans. ASLP.* **15**(1), 70–79 (2007)
- N. Ono, in *Proc. WASPAA*. Stable and fast update rules for independent vector analysis based on auxiliary function technique (IEEE, 2011), pp. 189–192
- D. Kitamura, N. Ono, H. Sawada, H. Kameoka, H. Saruwatari, Determined blind source separation unifying independent vector analysis and nonnegative matrix factorization. *IEEE/ACM Trans. ASLP.* **24**(9), 1626–1641 (2016)
- D. Kitamura, N. Ono, H. Sawada, H. Kameoka, H. Saruwatari, in *Audio Source Separation*, ed. by S. Makino. Determined blind source separation with independent low-rank matrix analysis (Springer, Cham, 2018), pp. 125–155
- T. Tachikawa, K. Yatabe, Y. Oikawa, in *Proc. IWAENC*. Underdetermined source separation with simultaneous DOA estimation without initial value dependency (IEEE, 2018), pp. 161–165
- D.D. Lee, H.S. Seung, Learning the parts of objects by non-negative matrix factorization. *Nature.* **401**(6755), 788–791 (1999)
- D.D. Lee, H.S. Seung, in *Proc. NIPS*. Algorithms for non-negative matrix factorization, (2000), pp. 556–562
- C. Févotte, N. Bertin, J.-L. Durrieu, Nonnegative matrix factorization with the Itakura–Saito divergence. With application to music analysis. *Neural Comput.* **21**(3), 793–830 (2009)
- Y. Mitsui, D. Kitamura, S. Takamichi, N. Ono, H. Saruwatari, in *Proc. ICASSP*. Blind source separation based on independent low-rank matrix analysis with sparse regularization for time-series activity (IEEE, 2017), pp. 21–25
- H. Kagami, H. Kameoka, M. Yukawa, in *Proc. ICASSP*. Joint separation and dereverberation of reverberant mixtures with determined multichannel non-negative matrix factorization (IEEE, 2018), pp. 31–35
- R. Ikeshita, Y. Kawaguchi, in *Proc. ICASSP*. Independent low-rank matrix analysis based on multivariate complex exponential power distribution (IEEE, 2018), pp. 741–745
- D. Kitamura, S. Mogami, Y. Mitsui, N. Takamune, H. Saruwatari, N. Ono, Y. Takahashi, K. Kondo, Generalized independent low-rank matrix analysis using heavy-tailed distributions for blind source separation. *EURASIP J. Adv. Signal Process.* **2018**, 28 (2018)
- K. Yoshii, K. Kitamura, Y. Bando, E. Nakamura, T. Kawahara, in *EUSIPCO*. Independent low-rank tensor analysis for audio source separation (IEEE, 2018), pp. 1657–1661
- R. Ikeshita, in *EUSIPCO*. Independent positive semidefinite tensor analysis in blind source separation (IEEE, 2018), pp. 1652–1656
- R. Ikeshita, N. Ito, T. Nakatani, H. Sawada, in *WASPAA*. Independent low-rank matrix analysis with decorrelation learning (IEEE, 2019), pp. 288–292
- N. Makishima, S. Mogami, N. Takamune, D. Kitamura, H. Sumino, S. Takamichi, H. Saruwatari, N. Ono, Independent deeply learned matrix analysis for determined audio source separation. *IEEE/ACM Trans. ASLP.* **27**(10), 1601–1615 (2019)
- K. Sekiguchi, Y. Bando, A.A. Nugraha, K. Yoshii, T. Kawahara, Semi-supervised multichannel speech enhancement with a deep speech prior. *IEEE/ACM Trans. ASLP.* **27**(12), 2197–2212 (2019)
- S. Mogami, N. Takamune, D. Kitamura, H. Saruwatari, Y. Takahashi, K. Kondo, N. Ono, Independent low-rank matrix analysis based on time-variant sub-Gaussian source model for determined blind source separation. *IEEE/ACM Trans. ASLP.* **28**, 503–518 (2019)
- Y. Takahashi, D. Kitahara, K. Matsuura, A. Hirabayashi, in *Proc. ICASSP*. Determined source separation using the sparsity of impulse responses (IEEE, 2020), pp. 686–690
- M. Togami, in *Proc. ICASSP*. Multi-channel speech source separation and dereverberation with sequential integration of determined and underdetermined models (IEEE, 2020), pp. 231–235
- S. Kanoga, T. Hoshino, H. Asoh, Independent low-rank matrix analysis-based automatic artifact reduction technique applied to three BCI paradigms. *Front. Hum. Neurosci.* **14**, 17 (2020)
- D. Kitamura, N. Ono, H. Saruwatari, in *Proc. EUSIPCO*. Experimental analysis of optimal window length for independent low-rank matrix analysis, (2017), pp. 1210–1214
- Y. Liang, S.M. Naqvi, J. Chambers, Overcoming block permutation problem in frequency domain blind source separation when using AuxIVA algorithm. *Electron. Lett.* **48**(8), 460–462 (2012)
- K. Yatabe, Consistent ICA: determined BSS meets spectrogram consistency. *IEEE Signal Process. Lett.* **27**, 870–874 (2020)
- T. Gerkmann, M. Krawczyk-Becker, J. Le Roux, Phase processing for single-channel speech enhancement: history and recent advances. *IEEE Signal Process. Mag.* **32**(2), 55–66 (2015)

34. P. Mowlaee, R. Saeidi, Y. Stylianou, Advances in phase-aware signal processing in speech communication. *Speech Commun.* **81**, 1–29 (2016)
35. P. Mowlaee, J. Kulmer, J. Stahl, F. Mayer, *Single channel phase-aware signal processing in speech communication: theory and practice*. (Wiley, 2016)
36. K. Yatabe, Y. Oikawa, in *Proc. ICASSP*. Phase corrected total variation for audio signals (IEEE, 2018), pp. 656–660
37. K. Yatabe, Y. Masuyama, Y. Oikawa, in *Proc. IWAENC*. Rectified linear unit can assist Griffin–Lim phase recovery (IEEE, 2018), pp. 555–559
38. Y. Masuyama, K. Yatabe, Y. Oikawa, in *Proc. IWAENC*. Model-based phase recovery of spectrograms via optimization on Riemannian manifolds (IEEE, 2018), pp. 126–130
39. Y. Masuyama, K. Yatabe, Y. Oikawa, Griffin–Lim like phase recovery via alternating direction method of multipliers. *IEEE Signal Process. Lett.* **26**(1), 184–188 (2019)
40. Y. Masuyama, K. Yatabe, Y. Koizumi, Y. Oikawa, N. Harada, in *Proc. ICASSP*. Deep Griffin–Lim iteration (IEEE, 2019), pp. 61–65
41. Y. Masuyama, K. Yatabe, Y. Oikawa, in *Proc. ICASSP*. Phase-aware harmonic/percussive source separation via convex optimization (IEEE, 2019), pp. 985–989
42. Y. Masuyama, K. Yatabe, Y. Oikawa, in *Proc. ICASSP*. Low-rankness of complex-valued spectrogram and its application to phase-aware audio processing (IEEE, 2019), pp. 855–859
43. Y. Masuyama, K. Yatabe, Y. Koizumi, Y. Oikawa, N. Harada, in *Proc. ICASSP*. Phase reconstruction based on recurrent phase unwrapping with deep neural networks (IEEE, 2020), pp. 826–830
44. J.L. Roux, H. Kameoka, N. Ono, S. Sagayama, in *Proc. DAFX*. Fast signal reconstruction from magnitude STFT spectrogram based on spectrogram consistency, (2010)
45. J. Le Roux, E. Vincent, Consistent Wiener filtering for audio source separation. *IEEE Signal Process. Lett.* **20**(3), 217–220 (2013)
46. N. Perraudin, P. Balazs, P.L. Søndergaard, in *Proc. WASPAA*. A fast Griffin–Lim algorithm (IEEE, 2013), pp. 1–4
47. K. Yatabe, Y. Masuyama, T. Kusano, Y. Oikawa, Representation of complex spectrogram via phase conversion. *Acoust. Sci. Tech.* **40**(3), 170–177 (2019)
48. M. Kowalski, E. Vincent, R. Gribonval, Beyond the narrowband approximation: wideband convex methods for under-determined reverberant audio source separation. *IEEE Trans. ASLP.* **18**(7), 1818–1829 (2010)
49. K. Matsuoka, S. Nakashima, in *Proc. ICA*. Minimal distortion principle for blind source separation, (2001), pp. 722–727
50. K. Yatabe, D. Kitamura, in *Proc. ICASSP*. Determined blind source separation via proximal splitting algorithm (IEEE, 2018), pp. 776–780
51. K. Yatabe, D. Kitamura, in *Proc. ICASSP*. Time-frequency-masking-based determined BSS with application to sparse IVA (IEEE, 2019), pp. 715–719
52. K. Yatabe, D. Kitamura, Determined BSS based on time-frequency masking and its application to harmonic vector analysis. *arXiv:2004.14091* (2020)
53. M. Brandstein, D. Ward, *Microphone arrays: signal processing techniques and applications*. (Springer Science & Business Media, 2013)
54. S. Araki, S. Makino, Y. Hinamoto, R. Mukai, T. Nishikawa, H. Saruwatari, Equivalence between frequency-domain blind source separation and frequency-domain adaptive beamforming for convolutive mixtures. *EURASIP J. Adv. Signal Process.* **2003**(11), 1157–1166 (2003)
55. D. Griffin, J. Lim, Signal estimation from modified short-time Fourier transform. *IEEE Trans. Acoust. Speech Signal Process.* **32**(2), 236–243 (1984)
56. D. Gunawan, D. Sen, Iterative phase estimation for the synthesis of separated sources from single-channel mixtures. *IEEE Signal Process. Lett.* **17**(5), 421–424 (2010)
57. N. Sturmel, L. Daudet, L. Girin, in *Proc. DAFX*. Phase-based informed source separation of music, (2012)
58. M. Watanabe, P. Mowlaee, in *Proc. INTERSPEECH*. Iterative sinusoidal-based partial phase reconstruction in single-channel source separation, (2013)
59. F. Mayer, D. Williamson, P. Mowlaee, D.L. Wang, Impact of phase estimation on single-channel speech separation based on time-frequency masking. *J. Acoust. Soc. Am.* **141**, 4668–4679 (2017)
60. S. Araki, F. Nesta, E. Vincent, Z. Koldovsky, G. Nolte, A. Ziehe, A. Benichoux, in *Proc. LVA/ICA*. The 2011 signal separation evaluation campaign (SISEC2011): -Audio source separation, (2012), pp. 414–422
61. S. Nakamura, K. Hiyané, F. Asano, T. Nishiura, T. Yamada, in *Proc. LREC*. Acoustical sound database in real environments for sound scene understanding and hands-free speech recognition, (2000), pp. 965–968
62. E. Vincent, R. Gribonval, C. Févotte, Performance measurement in blind audio source separation. *IEEE Trans. ASLP.* **14**(4), 1462–1469 (2006)
63. W.H. Press, S.A. Teukolsky, W.T. Vetterling, B.P. Flannery, *Numerical Recipes in C: The Art of Scientific Computing*. (Cambridge University Press, New York, 1992)
64. I. Andrianakis, P. White, Speech spectral amplitude estimators using optimally shaped gamma and chi priors. *Speech Comm.* **51**(1), 1–14 (2009)
65. P. Mowlaee, J. Stahl, Single-channel speech enhancement with correlated spectral components: limits-potential. *Speech Comm.* **121**, 58–69 (2020)

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Submit your manuscript to a SpringerOpen<sup>®</sup> journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

---

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)

---