

FastMVAE2: On improving and accelerating the fast variational autoencoder-based source separation algorithm for determined mixtures

Li Li, *Member, IEEE*, Hirokazu Kameoka, *Senior Member, IEEE*, and Shoji Makino, *Fellow, IEEE*,

Abstract—This paper proposes a new source model and training scheme to improve the accuracy and speed of the multichannel variational autoencoder (MVAE) method. The MVAE method is a recently proposed powerful multichannel source separation method. It consists of pretraining a source model represented by a conditional VAE (CVAE) and then estimating separation matrices along with other unknown parameters so that the log-likelihood is non-decreasing given an observed mixture signal. Although the MVAE method has been shown to provide high source separation performance, one drawback is the computational cost of the backpropagation steps in the separation-matrix estimation algorithm. To overcome this drawback, a method called “FastMVAE” was subsequently proposed, which uses an auxiliary classifier VAE (ACVAE) to train the source model. By using the classifier and encoder trained in this way, the optimal parameters of the source model can be inferred efficiently, albeit approximately, in each step of the algorithm. However, the generalization capability of the trained ACVAE source model was not satisfactory, which led to poor performance in situations with unseen data. To improve the generalization capability, this paper proposes a new model architecture (called the “ChimeraACVAE” model) and a training scheme based on knowledge distillation. The experimental results revealed that the proposed source model trained with the proposed loss function achieved better source separation performance with less computation time than FastMVAE. We also confirmed that our methods were able to separate 18 sources with a reasonably good accuracy.

Index Terms—Multichannel source separation, multichannel variational autoencoder (MVAE), fast algorithm, auxiliary classifier VAE, knowledge distillation

I. INTRODUCTION

BLIND source separation (BSS) is a technique for separating observed signals recorded by a microphone array into individual source signals without prior information about the sources or mixing conditions. This technique has been used in a wide range of applications, including hearing aids, automatic speech recognition (ASR), telecommunications systems, music editing, and music information retrieval.

Acoustic signals are convolved with the impulse responses of acoustic environments and so the signal observed at a

particular position is usually given as the convolutive mixture of nearby source signals. Although it is possible to take a time-domain approach to the BSS problem, it can be computationally expensive since it requires directly estimating and applying demixing filters with thousands of taps. In contrast, the time-frequency-domain approach is advantageous in that the convolution operations can be replaced by multiplications to achieve computationally efficient algorithms, and it allows the flexible use of various models for the time-frequency (TF) representations of source signals. Independent vector analysis (IVA) [2], [3] is an example of the time-frequency-domain approach, which makes it possible to solve frequency-wise source separation and permutation alignment simultaneously by assuming that the magnitudes of the frequency components originating from the same source vary coherently over time. Multichannel nonnegative matrix factorization (MNMF) [4], [5] and independent low-rank matrix analysis (ILRMA) [6]–[8] are other examples, which employ the concept of NMF [9] to model the TF structures of sources. Specifically, they assume that the power spectrum of each source signal can be approximated as the sum of a limited number of basis spectra scaled by time-varying amplitudes. IVA can be understood as a special case of ILRMA where only one flat basis spectrum is used for representing each source. This indicates that ILRMA can capture the TF structure of each source more flexibly than IVA, and this flexibility has been shown to be advantageous in improving the source separation performance [7].

Recently, the success of deep neural network (DNN)-based speech separation methods [10]–[18], including deep clustering (DC) [11], [12] and permutation invariant training (PIT) [13], [14], has proven that DNNs have an excellent ability to capture and learn the structure of spectrograms. The general idea of these methods is to train a network that predicts a TF mask or clean signals given the spectral and spatial features of observed mixture signals. Meanwhile, time-domain methods based on end-to-end training have also been extensively studied and have shown excellent performance [19]–[21]. Some attempts [22], [23] have been made to combine beamforming with the time-domain methods to avoid artifacts introduced by nonlinear processing. Although such an end-to-end approach provides reasonably good separation performance, one drawback is that it suffers from the limitation that the test conditions need to be similar to the training ones, such as the number of speakers and reverberation conditions.

There have also been some attempts to incorporate DNNs into the BSS methods mentioned earlier [24]–[29]. Indepen-

L. Li and H. Kameoka are with NTT Communication Science Laboratories, NTT Corporation, 3-1 Morinosato Wakamiya, Atsugi-shi, Kanagawa 243-0198, Japan. L. Li is currently also with Nagoya University, Furocho, Chikusa-ku, Nagoya, 464-8601, JAPAN, (email: lili-0805@ieee.org, hirokazu.kameoka.uh@hco.ntt.co.jp).

S. Makino is with Waseda University, 2-7 Hibikino, Wakamatsu-ku, Kitakyushu city, Fukuoka, 808-0135, Japan, and University of Tsukuba, 1-1-1 Tennodai, Tsukuba, Ibaraki, 2050821, Japan, (email: s.makino@waseda.jp).

This work was partially supported by JST CREST JPMJCR19A3. A preprint version of this paper has already been made publicly available at [1].

dent deeply low-rank matrix analysis (IDLMA) [25], [30] is one such method, where each DNN is trained using the utterances of a different speaker. After training, the trained DNNs are used to refine the estimated power spectra at each iteration of the source separation algorithm. Namely, each DNN can be seen as a speaker-dependent speech enhancement system. One drawback of IDLMA would be that it can perform poorly in speaker-independent scenarios due to its discriminative training scheme. Within the DNN framework, deep generative models such as variational autoencoders (VAEs) [31], [32], generative adversarial networks (GANs) [33], and normalizing flow (NF) [34] have proven to be powerful in source separation tasks [26]–[29], [35]–[43]. An attempt to employ VAE for semi-supervised single-channel speech enhancement was made in [27] under the name of the “VAE-NMF” method, which uses a VAE to model each single-frame spectrum in an utterance of a target speaker and an NMF model to express a noise spectrogram. Several variants of this method have subsequently been developed, including the incorporation of loudness gain for robust speech modeling [28], the adoption of a noise model based on alpha-stable distribution instead of a complex Gaussian distribution [37], and the extension to multichannel scenarios [36], [38].

Independently, around the same time, we proposed a method called the “multichannel variational autoencoder (MVAE)”. This was the first to incorporate the VAE concept into the multichannel source separation framework, and it has proven to be very successful in supervised determined source separation tasks. Unlike the VAE-NMF methods, the MVAE method uses a conditional VAE (CVAE) with a fully convolutional architecture to model the entire spectrogram of each utterance. The CVAE is trained with the spectrograms of clean speech samples along with the corresponding speaker ID as a conditioning class variable. This is done so that the trained decoder distribution can be used as a generative model of signals produced by all the sources included in a given training set, where the latent space variables and the class variables are the parameters to be estimated from an input mixture signal. The generative model trained in this way is called the *CVAE source model*. At the separation phase, the MVAE algorithm iteratively updates the separation matrix using the iterative projection (IP) method [44] and the underlying parameters of the CVAE source model using a gradient descent method, where the gradients of the latent variables are calculated using backpropagation. The main feature of this optimization algorithm is that the log-likelihood is guaranteed to be non-decreasing if the step size is carefully chosen or if a backtracking line search [45] is applied for the backpropagation algorithm. Furthermore, since the MVAE uses a CVAE to model single source and the demixing matrices are estimated only at separation phase, a trained CVAE source model is principle able to handle arbitrary number of sources and different recording conditions, which is significantly differ from discriminative methods. However, one major drawback of the MVAE method is that the backpropagation required for each iteration makes the optimization algorithm very time-consuming, which can be problematic in practice.

To address this problem, we previously proposed a fast

algorithm called “FastMVAE” [46], which uses an auxiliary classifier VAE (ACVAE) [47] to model the generative distribution of source spectrograms. In this method, the encoder and auxiliary classifier are trained in such a way that they learn to infer the latent space variables and class variables, respectively, given a spectrogram. This allows us to replace the backpropagation steps in the source separation algorithm with the forward propagation of the two networks and thus significantly reduce the computational cost. Furthermore, we showed that FastMVAE can achieve source separation performance comparable to the MVAE method when the training and test conditions are sufficiently close to being consistent. However, when there is mismatch between the training and test conditions, due to, for example, the presence of long reverberation or under speaker-independent conditions, FastMVAE tends to perform worse than the MVAE method. This may be because the encoder and classifier cannot generalize well to inputs that are very different from the training data. To stabilize the parameter inference process under such mismatched conditions, we derived an improved update rule based on the Product-of-Experts (PoE) framework [48]. However, this method requires manual selection of the optimal weights in advance, forcing us to rely on heuristics.

FastMVAE being weak against the mismatch between the training and test conditions may be because the model is structured in such a way that the output of the auxiliary classifier is fed into the encoder and so the error in the classifier output can directly affect the encoder output. One way to avoid this would be to assume a conditional independence between the outputs of the encoder and auxiliary classifier so that they can perform their tasks in parallel. Instead of preparing two separate networks, we propose merging the encoder and classifier into a single multitask network to allow them to share information. We call this new model the “*ChimeraACVAE source model*”.

Another important issue is how to train the above model to have good generalization ability. A number of techniques have been developed with the aim of improving the generalization ability of DNNs. These techniques can be roughly classified into regularization-based [49]–[52], data augmentation-based [53], and training strategy-based methods [54]–[56]. Knowledge distillation (KD), a model compression and acceleration technique that has been rapidly gaining attention in recent years, is typically used to transfer knowledge of a teacher model to a more compact student model. KD has been shown to not only accelerate the inference process through model compression but also provide better generalization ability to the compressed model. In this paper, we propose adopting KD to train the ChimeraACVAE source model. Specifically, we use a pretrained CVAE model as a teacher model and transfer its knowledge to the ChimeraACVAE model by using as a criterion the Kullback-Leibler (KL) divergence between the distributions of the outputs of the encoder and decoder of the CVAE and ChimeraACVAE models.

In summary, the two main contributions of this paper are as follows:

- We propose a new network architecture that replaces the ACVAE source model in FastMVAE, which we call the

“ChimeraACVAE” source model. It merges the encoder and classifier into a single multitask network so that it can handle the tasks of the encoder and classifier simultaneously.

- We propose a loss function based on the KD framework that allows the ChimeraACVAE source model to acquire excellent generalization capability. We show that the model trained in this way can improve source separation performance in both speaker-dependent and speaker-independent conditions.

The rest of this paper is structured as follows. After describing the formulation of the determined multichannel BSS problem and reviewing the original MVAE method in Section II, we describe the ACVAE source model and the FastMVAE method in Section III. In Section IV, we provide technical details of the proposed ChimeraACVAE source model and its training strategy. The effectiveness of the proposed method is demonstrated in Section V by evaluating the source separation performance of speaker-dependent and speaker-independent scenarios. We conclude the article in Section VI.

II. MVAE

A. Problem Formulation

Let us consider a situation where I source signals are captured by I microphones. We use $x_i(f, n)$ and $s_j(f, n)$ to denote the short-time Fourier transform (STFT) coefficients of the signal observed at the i th microphone and j th source signal, where f and n are the frequency and time indices, respectively. If we use

$$\mathbf{x}(f, n) = [x_1(f, n), \dots, x_I(f, n)]^T \in \mathbb{C}^I, \quad (1)$$

$$\mathbf{s}(f, n) = [s_1(f, n), \dots, s_I(f, n)]^T \in \mathbb{C}^I, \quad (2)$$

to denote the vectors containing $x_1(f, n), \dots, x_I(f, n)$ and $s_1(f, n), \dots, s_I(f, n)$, the relationship between the observed signals and source signals can be approximated as

$$\mathbf{s}(f, n) = \mathbf{W}^H(f) \mathbf{x}(f, n), \quad (3)$$

$$\mathbf{W}(f) = [\mathbf{w}_1(f), \dots, \mathbf{w}_I(f)] \in \mathbb{C}^{I \times I}, \quad (4)$$

under a determined mixing condition, where $\mathbf{W}^H(f)$ represents the separation matrix, and $(\cdot)^T$ and $(\cdot)^H$ denote the transpose and Hermitian transpose of a matrix or a vector, respectively. The goal of BSS is to determine $\mathcal{W} = \{\mathbf{W}(f)\}_f$ solely from the observation $\mathcal{X} = \{\mathbf{x}(f, n)\}_{f,n}$. Here, the notation $\{\mathbf{E}_b\}_b$ is used as an abbreviation for $\{\mathbf{E}_b \mid b \in \mathcal{B}\}$, where \mathcal{B} denotes the set of all possible indices.

In the following, we assume that $s_j(f, n)$ independently follows a zero-mean complex proper Gaussian distribution with variance (power spectral density) $v_j(f, n) = \mathbb{E}[|s_j(f, n)|^2]$:

$$p(s_j(f, n)|v_j(f, n)) = \mathcal{N}_{\mathbb{C}}(s_j(f, n)|0, v_j(f, n)), \quad (5)$$

This assumption is often referred to as the local Gaussian model (LGM) [57], [58]. If $s_j(f, n)$ and $s_{j'}(f, n)$ are independent for $\forall j \neq j'$, the density of $\mathbf{s}(f, n)$ becomes

$$p(\mathbf{s}(f, n)|\mathbf{V}(f, n)) = \prod_j p(s_j(f, n)|v_j(f, n))$$

$$= \mathcal{N}_{\mathbb{C}}(\mathbf{s}(f, n)|\mathbf{0}, \mathbf{V}(f, n)), \quad (6)$$

where $\mathbf{V}(f, n) = \text{diag}[v_1(f, n), \dots, v_I(f, n)]$. From (3) and (6), the density of $\mathbf{x}(f, n)$ is obtained as

$$p(\mathbf{x}(f, n)|\mathbf{W}(f), \mathbf{V}(f, n)) = |\mathbf{W}^H(f)|^2 p(\mathbf{s}(f, n) = \mathbf{W}^H(f) \mathbf{x}(f, n)|\mathbf{V}(f, n)), \quad (7)$$

where $|\mathbf{W}^H(f)|^2$ is the Jacobian of the mapping $\mathbf{x}(f, n) \mapsto \mathbf{s}(f, n)$. Therefore, the log-likelihood of $\mathcal{W} = \{\mathbf{W}(f)\}_f$ and $\mathcal{V} = \{v_j(f, n)\}_{f,n,j}$, given $\mathcal{X} = \{\mathbf{x}(f, n)\}_{f,n}$ is expressed as

$$\begin{aligned} & \log p(\mathcal{X}|\mathcal{W}, \mathcal{V}) \\ &= 2N \sum_f \log |\det \mathbf{W}^H(f)| + \sum_j \log p(\mathcal{S}_j|\mathcal{V}_j) \\ &\stackrel{c}{=} 2N \sum_f \log |\det \mathbf{W}^H(f)| \\ &\quad - \sum_{f,n,j} \left(\log v_j(f, n) + \frac{|\mathbf{w}_j^H(f) \mathbf{x}(f, n)|^2}{v_j(f, n)} \right), \quad (8) \end{aligned}$$

where we have used $\stackrel{c}{=}$ to denote equality up to constant terms and a bold italic font to indicate a set consisting of TF elements, namely $\mathcal{S}_j = \{s_j(f, n)\}_{f,n}$ and $\mathcal{V}_j = \{v_j(f, n)\}_{f,n}$. The log-likelihood will be split into F frequency-wise terms if no additional constraint is imposed on $v_j(f, n)$ or $\mathbf{W}(f)$, implying that there is a permutation ambiguity in the separated components for each frequency. Thus, the separated components of different frequency bins that originate from the same source need to be grouped together in order to complete source separation. This process is called permutation alignment [59], [60].

B. CVAE Source Model

Incorporating an appropriate constraint into the power spectrogram $\mathcal{V}_j = \{v_j(f, n)\}_{f,n}$ not only helps eliminate the permutation ambiguity but also provides a clue for estimating \mathcal{W} . In the MVAE method, the complex spectrogram of a single source $\mathcal{S} = \{s(f, n)\}_{f,n}$ is modeled using a CVAE [31] conditioned on a class variable \mathbf{c} . Here, \mathbf{c} is a one-hot vector consisting of C elements that indicates to which class the separated signal belongs. For example, speaker IDs can be used as the class category in multispeaker separation tasks, where the entries of \mathbf{c} will be 1 at the index of a certain speaker and 0 at all other indices.

Since the following applies to all sources, index j will be omitted throughout this paragraph. A CVAE consists of decoder and encoder networks. The decoder network is designed to produce the parameters of the distribution $p_{\theta}^*(\mathcal{S}|\mathbf{z}, \mathbf{c})$ of data \mathcal{S} given a latent space variable \mathbf{z} and a class variable \mathbf{c} . The encoder network is designed to generate the parameters of a conditional distribution $q_{\phi}^*(\mathbf{z}|\mathcal{S}, \mathbf{c})$ that approximates the exact posterior $p_{\theta}^*(\mathbf{z}|\mathcal{S}, \mathbf{c})$. The goal of the CVAE training is to find the weight parameters in the encoder and decoder networks, namely θ and ϕ , such that the encoder distribution $q_{\phi}^*(\mathbf{z}|\mathcal{S}, \mathbf{c})$ becomes consistent with the posterior $p_{\theta}^*(\mathbf{z}|\mathcal{S}, \mathbf{c}) \propto p_{\theta}^*(\mathcal{S}|\mathbf{z}, \mathbf{c})p(\mathbf{z})$. Note that the KL divergence between $q_{\phi}^*(\mathbf{z}|\mathcal{S}, \mathbf{c})$ and $p_{\theta}^*(\mathbf{z}|\mathcal{S}, \mathbf{c})$ is shown

to be equal to the difference between the log marginal likelihood $p_\theta^*(\mathbf{S}|\mathbf{c}) = \int_{\mathbf{z}} p(\mathbf{S}|\mathbf{z}, \mathbf{c}) p(\mathbf{z}) d\mathbf{z}$ and its variational lower bound. Hence, minimizing the KL divergence between $q_\phi^*(\mathbf{z}|\mathbf{S}, \mathbf{c})$ and $p_\theta^*(\mathbf{z}|\mathbf{S}, \mathbf{c})$ amounts to maximizing the following variational lower bound [61]:

$$\mathcal{J} = \mathbb{E}_{(\mathbf{S}, \mathbf{c})} [\mathbb{E}_{\mathbf{z} \sim q_\phi^*(\mathbf{z}|\mathbf{S}, \mathbf{c})} [\log p_\theta^*(\mathbf{S}|\mathbf{z}, \mathbf{c})] - \text{KL}[q_\phi^*(\mathbf{z}|\mathbf{S}, \mathbf{c}) || p(\mathbf{z})]], \quad (9)$$

where we have used $\mathbb{E}_{(\mathbf{S}, \mathbf{c})}[\cdot]$ to denote the sample mean of its argument over the training examples $\{\mathbf{S}_m, \mathbf{c}_m\}_{m=1}^M$, and $\text{KL}[\cdot || \cdot]$ to denote the KL divergence. Although it is difficult to obtain an analytical form of the expectation $\mathbb{E}_{\mathbf{z} \sim q_\phi^*(\mathbf{z}|\mathbf{S}, \mathbf{c})}[\cdot]$ in the first term of \mathcal{J} , we can use a reparameterization trick [61] to obtain a form that allows us to compute the gradient with respect to ϕ using a Monte Carlo approximation. Now, $q_\phi^*(\mathbf{z}|\mathbf{S}, \mathbf{c})$, $p_\theta^*(\mathbf{S}|\mathbf{z}, \mathbf{c})$, and $p(\mathbf{z})$ are distributions that need to be modeled. In the MVAE method, $p(\mathbf{z})$ and $q_\phi^*(\mathbf{z}|\mathbf{S}, \mathbf{c})$ are described as Gaussian distributions as with a regular CVAE:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}), \quad (10)$$

$$q_\phi^*(\mathbf{z}|\mathbf{S}, \mathbf{c}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_\phi^*(\mathbf{S}, \mathbf{c}), \text{diag}(\boldsymbol{\sigma}_\phi^{*2}(\mathbf{S}, \mathbf{c}))), \quad (11)$$

where $\boldsymbol{\mu}_\phi^*(\mathbf{S}, \mathbf{c})$ and $\boldsymbol{\sigma}_\phi^{*2}(\mathbf{S}, \mathbf{c})$ are the encoder network outputs. For stable training, the total energy of each training utterance is normalized to 1. However, the energy of each source in a test mixture does not necessarily equal 1. To fill this gap, a scale factor g is additionally introduced into the decoder distribution as a free parameter to be estimated at test time. Specifically, we use an expression of the decoder distribution with variance scaled by g . Hence, the decoder distribution for the complex spectrogram \mathbf{S} is expressed as

$$p_\theta^*(\mathbf{S}|\mathbf{z}, \mathbf{c}, g) = \prod_{f,n} \mathcal{N}_{\mathbb{C}}(s(f, n)|0, g\sigma_\theta^{*2}(f, n; \mathbf{z}, \mathbf{c})), \quad (12)$$

where $\sigma_\theta^{*2}(f, n; \mathbf{z}, \mathbf{c})$ denotes the (f, n) th element of the decoder network output. (12) is called the CVAE source model.

If we use the above CVAE source model to represent the complex spectrogram of the j th signal in a mixture signal, \mathbf{z}_j , \mathbf{c}_j , and g_j are the unknown parameters to be estimated. Since the CVAE source model is given in the same form as the LGM in (5) if we denote $g_j\sigma_\theta^{*2}(f, n; \mathbf{z}_j, \mathbf{c}_j)$ by $v_j(f, n)$, using this as the generative model for each source gives the log-likelihood in the same form as (8).

C. Optimization Algorithm

The goal of the source separation algorithm in the MVAE method is to maximize the posterior $p(\mathcal{W}, \Psi, \mathcal{G}|\mathcal{X}) \propto p(\mathcal{X}|\mathcal{W}, \Psi, \mathcal{G})p(\mathbf{z})p(\mathbf{c})$ with respect to \mathcal{W} , $\Psi = \{\mathbf{z}_j, \mathbf{c}_j\}_j$, and $\mathcal{G} = \{g_j\}_j$, where \mathbf{z} is assumed to follow $\mathcal{N}(\mathbf{0}, \mathbf{I})$, and $p(\mathbf{c})$ is the empirical distribution of the training examples $\{\mathbf{c}_m\}_m$, expressed as a multinomial distribution. Hence, the objective function is $\log p(\mathcal{X}|\mathcal{W}, \Psi, \mathcal{G}) + \log p(\mathbf{z}) + \log p(\mathbf{c})$. A stationary point of this function can be found by iteratively updating \mathcal{W} , Ψ , and \mathcal{G} so that the function value is guaranteed to be non-decreasing. To update \mathcal{W} , we can use the IP method [44]:

$$\mathbf{w}_j(f) \leftarrow (\mathbf{W}^H(f)\boldsymbol{\Sigma}_j(f))^{-1}\mathbf{e}_j, \quad (13)$$

$$\mathbf{w}_j(f) \leftarrow \frac{\mathbf{w}_j(f)}{\sqrt{\mathbf{w}_j^H(f)\boldsymbol{\Sigma}_j(f)\mathbf{w}_j(f)}}, \quad (14)$$

where $\boldsymbol{\Sigma}_j(f) = \frac{1}{N} \sum_n \mathbf{x}(f, n)\mathbf{x}^H(f, n)/v_j(f, n)$ and \mathbf{e}_j denotes the j th column of an $I \times I$ identity matrix. As for \mathcal{G} , the update rule

$$g_j \leftarrow \frac{1}{FN} \sum_{f,n} \frac{|\mathbf{w}_j^H(f)\mathbf{x}(f, n)|^2}{\sigma_\theta^{*2}(f, n; \mathbf{z}_j, \mathbf{c}_j)} \quad (15)$$

maximizes the objective function with respect to g_j when \mathcal{W} and Ψ are fixed. Under fixed \mathcal{W} and \mathcal{G} , the optimal \mathbf{z}_j and \mathbf{c}_j that maximize the objective function can be found using the gradient descent method. Note that \mathbf{c}_j can be updated under the sum-to-one constraint by inserting an appropriately designed softmax layer that outputs \mathbf{c}_j .

One important feature of VAE in general is its generalization capability, namely the ability to learn the distribution of unseen data. Thanks to this feature, we expect that the CVAE source model trained on speech samples of sufficiently many speakers can generalize somewhat well to the spectrograms of unknown speakers, thus allowing the above algorithm to handle speaker-independent scenarios reasonably well. Another advantage is that it is guaranteed to converge to a stationary point, making it easy to handle in practical use. However, the downside is that the backpropagation algorithm required for each iteration can be computationally expensive.

III. FASTMVAE

A. ACVAE Source Model

The motivation behind the FastMVAE method is to accelerate the process of updating Ψ . Under fixed \mathcal{W} and \mathcal{G} , the objective function of the MVAE method is equal to the sum of $\log p(\mathbf{z}_j, \mathbf{c}_j|\mathbf{S}_j, g_j)$ up to a constant, where \mathbf{S}_j is the set $\{\mathbf{w}_j^H(f)\mathbf{x}(f, n)\}_{f,n}$, namely the complex spectrogram of the signal separated from the observed signal using the current estimate of \mathcal{W} . The idea of the FastMVAE method is to express this posterior as $p(\mathbf{z}_j, \mathbf{c}_j|\mathbf{S}_j, g_j) = p(\mathbf{z}_j|\mathbf{S}_j, \mathbf{c}_j, g_j)p(\mathbf{c}_j|\mathbf{S}_j, g_j)$ and use two trainable networks to approximate these two conditional distributions. Once these networks have been trained, an approximation of the maximum point of the posterior $p(\mathbf{z}_j, \mathbf{c}_j|\mathbf{S}_j, g_j)$ can be obtained by finding the maximum points of the two approximate distributions.

To obtain approximations of the two conditional distributions, the FastMVAE method employs the idea of ACVAE training [47]. ACVAE is a CVAE variant that incorporates the expectation of the mutual information [62]

$$I(\mathbf{c}, \mathbf{S}|\mathbf{z}) = \mathbb{E}_{\mathbf{c} \sim p_D(\mathbf{c}), \mathbf{S} \sim p_\theta(\mathbf{S}|\mathbf{z}, \mathbf{c}), \mathbf{c}' \sim p(\mathbf{c}|\mathbf{S})} [\log p(\mathbf{c}'|\mathbf{S})] + H(\mathbf{c}), \quad (16)$$

into the training criterion with the aim of making the decoder output as correlated as possible with the class variable \mathbf{c} . Here, $p_D(\mathbf{c})$ is the empirical discrete distribution of the samples of \mathbf{c} in the training set and $H(\mathbf{c})$ represents the entropy of \mathbf{c} , which can be regarded as a constant. Since it is difficult to express $I(\mathbf{c}, \mathbf{S}|\mathbf{z})$ in analytical form, rather than using it directly, ACVAE uses its variational lower bound

$$\mathcal{L} = \mathbb{E}_{(\mathbf{S}, \mathbf{c}'), \mathbf{z} \sim q_\phi^*(\mathbf{z}|\mathbf{S}, \mathbf{c}')} [\mathbb{E}_{\mathbf{c} \sim p_\theta^*(\mathbf{S}|\mathbf{z}, \mathbf{c})} [\log \tau_\psi^*(\mathbf{c}|\mathbf{S}, g)]] \quad (17)$$

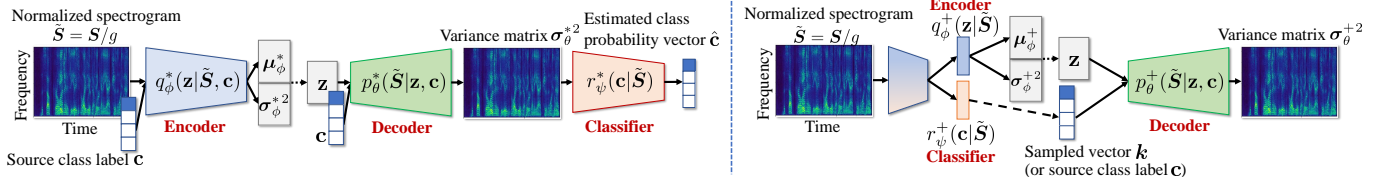


Figure 1: Illustration of the ACVAE model in FastMVAE (left) and the ChimeraACVAE model in FastMVAE2 (right). We use $\tilde{S} = S/g$ to denote a normalized spectrogram and omit $g = 1$ in encoder, decoder, and classifier distributions.

defined using a variational distribution $r_\psi^*(c|S, g) = \text{Mult}(c|\rho_\psi^*(S/g))$ for optimization, where $\mathbb{E}_{(S, c')}[\cdot]$ is equivalent to $\mathbb{E}_{(S, c)}[\cdot]$, $\mathbb{E}_c[\cdot]$ denotes the mean of its argument over all one-hot vectors $c \sim p_D(c)$, which can be approximated by a Monte Carlo approximation, and $\mathbb{E}_{z \sim q_\phi^*(z|S, c)}[\cdot]$ and $\mathbb{E}_{S \sim p_\theta^*(S|z, c')}[\cdot]$ are approached by a Monte Carlo approximation after reparameterization tricks. Here, $\text{Mult}(c|\rho) \propto \prod_i \rho_i^{c_i}$ denotes a multinomial distribution, where $c = [c_1, \dots, c_I]^T$ and $\rho = [\rho_1, \dots, \rho_I]^T$. $\rho_\psi^*(S/g)$ is a neural network that takes S normalized by g as an input and produces a probability vector consisting of C elements that sum to 1. $r_\psi^*(c|S, g)$ is an auxiliary classifier. Since the exact bound is obtained when $r_\psi^*(c|S, g) = p(c|S, g)$, the trained auxiliary classifier $r_\psi^*(c|S, g)$ is expected to be a good approximation of the distribution $p(c|S, g)$ of interest. ACVAE also uses the negative cross-entropy

$$\mathcal{I} = \mathbb{E}_{(S, c)}[\log r_\psi^*(c|S, g)] \quad (18)$$

as the training criterion. Therefore, the entire training criterion to be maximized is given by

$$\mathcal{J} + \lambda_{\mathcal{L}}\mathcal{L} + \lambda_{\mathcal{I}}\mathcal{I}, \quad (19)$$

where $\lambda_{\mathcal{L}}, \lambda_{\mathcal{I}} \geq 0$ denote the regularization weights that weight the importance of the regularization terms. The set of the networks trained in this way using the spectrograms of the training utterances is called the *ACVAE source model*. An illustration of ACVAE is shown on the left of Fig. 1.

B. Optimization Algorithm

After ACVAE training, we achieve $p(z_j, c_j|S_j, g_j) \approx r_\psi^*(c_j|S_j, g_j)q_\phi^*(z_j|S_j, c_j, g_j)$. Since the maximum points of $r_\psi^*(c_j|S_j, g_j)$ and $q_\phi^*(z_j|S_j, c_j, g_j)$ can be found through the forward passes of the auxiliary classifier and encoder, respectively, we can quickly find an approximate solution to $(z_j, c_j) = \arg\max_{z_j, c_j} p(z_j, c_j|S_j, g_j)$ without resorting to gradient descent updates. Specifically, c_j is given as the probability vector produced by the auxiliary classifier network:

$$c_j \leftarrow \rho_\psi^*(S_j/g_j), \quad (20)$$

and z_j is given as the mean of the encoder distribution:

$$z_j \leftarrow \mu_\phi^*(S_j/g_j, c_j). \quad (21)$$

Here, if the j th separated signal corresponds to a speaker unseen in the training set, the elements of (20) can be interpreted as quantities indicating how similar that speaker is to all the speakers in the training set. If the signal of any

Algorithm 1 FastMVAE algorithm w/ PoE

Require: Network parameter θ, ϕ, ψ trained using (19), observed mixture signal $x(f, n)$, iteration number \mathcal{L} , weight parameter α

- 1: randomly initialize \mathcal{W}, Ψ
- 2: **for** $\ell = 1$ to \mathcal{L} **do**
- 3: **for** $j = 1$ to J **do**
- 4: $y_j(f, n) = \mathbf{w}_j^H(f)\mathbf{x}(f, n)$
- 5: (updating source model parameters)
- 6: initialize g_j using (15)
- 7: normalize $\tilde{S}_j = \{y_j(f, n)/g_j\}_{f, n}$
- 8: update c_j using (20)
- 9: update z_j using (22)
- 10: compute $\sigma_j^{*2}(f, n; z_j, c_j, g_j = 1, \theta)$
- 11: update g_j using (15)
- 12: compute $v_j(f, n) = g_j \cdot \sigma_j^{*2}(f, n; z_j, c_j, g_j = 1, \theta)$
- 13: (updating separation matrices)
- 14: **for** $f = 1$ to F **do**
- 15: update $\mathbf{w}_j(f)$ by IP method with (13), (14)
- 16: **end for**
- 17: **end for**
- 18: **end for**

speaker can be assumed to be expressed as a point in the manifold spanned by all the speakers in the training set, our algorithm is expected to be able to handle even mixtures of unknown speakers.

However, our preliminary experiments revealed that directly using the mean of the encoder distribution tends to degrade source separation performance for speakers not included in the training set. To stabilize the inference for unknown speakers, we previously proposed reapplying the prior $p(z_j)$ to the encoder output based on the PoE framework [48] to ensure that z_j will not be updated to an outlier. Namely, the prior $p(z_j)$ is redefined as the product of two distributions with respect to z_j , namely, $\arg\max_{z_j} p(z_j|S_j, c_j, g_j)p(z_j)^\alpha$. Accordingly, the modified update rule of z_j is given as

$$z_j \leftarrow \Sigma_{\phi, j}^{-1}(\Sigma_{\phi, j}^{-1} + \alpha\mathbf{I})^{-1}\mu_\phi^*(S_j/g_j, c_j). \quad (22)$$

Here, α is a parameter that weights the importance of the prior $p(z_j)$ in the inference, and $\Sigma_{\phi, j} = \text{diag}(\sigma_\phi^{*2}(S_j/g_j, c_j))$. Note that (22) reduces to the mean of the encoder distribution when $\alpha = 0$. The algorithm of the FastMVAE method is summarized in **Algorithm 1**.

IV. PROPOSED: FASTMVAE2

While the FastMVAE method can significantly reduce the computation time compared to the MVAE method, its source separation accuracy has been confirmed to be somewhat less than that of the MVAE method [46]. We believe that this is due to the limitations of the generalization capabilities of the encoder and classifier obtained from the ACVAE training. In this paper, we propose introducing a new model architecture and training scheme to overcome these limitations, rather than implementing a heuristic solution at the inference stage.

A. ChimeraACVAE source model

We first describe our motivation and ideas for developing an improved version of the ACVAE source model, which we call the “ChimeraACVAE” source model.

1) *Multitask encoder*: When performing source separation, it is desirable that the speaker identity of each separated signal does not change over time. This is because a change of the identity of each separated signal means a failure in source separation. However, constraining the identity not to change is not an easy task if the decoder is not conditioned on \mathbf{c} (as in a regular VAE), since it will be trained so that \mathbf{z} becomes an entangled mixture of linguistic and speaker-identity information. In contrast, conditioning the decoder on \mathbf{c} is expected to promote disentanglement between \mathbf{z} and \mathbf{c} so that \mathbf{z} represents only the linguistic information and \mathbf{c} represents only the speaker identity. This allows our source separation system to always ensure that the speaker identity of each separated signal is time-invariant. Thus, it is essential for the decoder to remain conditioned on \mathbf{c} , and it is the encoder that we propose to modify. Specifically, we unify the encoder and auxiliary classifier into a single network with two branches that output the parameters of the encoder distribution $q_{\phi}^{+}(\mathbf{z}|\mathbf{S}, g) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\phi}^{+}(\mathbf{S}/g), \text{diag}(\boldsymbol{\sigma}_{\phi}^{+2}(\mathbf{S}/g)))$ and those of the class distribution $r_{\psi}^{+}(\mathbf{c}|\mathbf{S}, g) = \text{Mult}(\mathbf{c}|\boldsymbol{\rho}_{\psi}^{+}(\mathbf{S}/g))$, respectively. Here, the latent variable \mathbf{z} and speaker identity \mathbf{c} are assumed to be conditionally independent. We believe that the main reason for the performance degradation in FastMVAE under the speaker-independent condition is the cascade structure of the classifier and encoder, where errors in the classifier directly affect the outputs of the encoder. The conditional independence assumption in the ChimeraACVAE source model allows us to parallelize the processes by the classifier and encoder and prevent error propagation. Furthermore, the sharing of the layers in the unified encoder network is expected to improve the generalization capability through multitask learning.

2) *Network details*: The original ACVAE source model is designed to include batch normalization layers in its networks. However, since the computation of batch normalization depends on the mini-batch size, the learned parameters may be suboptimal in inference situations where the number of sources differs from the mini-batch size during training. To avoid inconsistencies in computation during training and inference, we replace batch normalization [55] with layer normalization [63]. In addition, we use a sigmoid linear unit (SiLU) [64] instead of a gated linear unit (GLU) [65] to reduce model

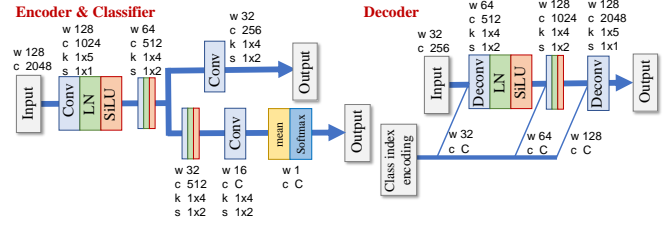


Figure 2: Network architectures of the unified encoder and decoder in the ChimeraACVAE source model. The inputs and outputs are assumed to be vector sequences. A spectrogram is interpreted as a sequence of spectra, with frequency regarded as the channel dimension. “w” denotes the length of the input sequence. “Conv” and “Deconv” denote one-dimensional convolution and deconvolution, respectively, where “c”, “k”, and “s” denote the channel number, kernel size, and stride size, respectively. “LN” and “SiLU” stand for the layer normalization and sigmoid linear unit, respectively. “mean” denotes the operation of averaging the input sequence along the time direction, and “softmax” denotes the operation of applying a softmax function to the input vector. In the decoder, the “class index encoding” \mathbf{c} is concatenated to the input of each deconvolution layer along the channel direction after being repeated along the time direction so that it has the shape compatible with the input.

size. SiLU, also known as the swish activation function, is a self-gated activation function, which can be expressed as

$$\odot_l = (\odot_{l-1} * \mathbb{W}_l + \mathbb{b}_l) \otimes \sigma(\odot_{l-1} * \mathbb{W}_l + \mathbb{b}_l) \quad (23)$$

when applied to a convolution layer. Here, \mathbb{W}_l and \mathbb{b}_l are weight and bias parameters of the l th layer, and \odot_l and \odot_{l-1} denote the output and input of the l th layer, respectively. \otimes denotes element-wise multiplication, and $\sigma(\cdot)$ is the sigmoid function. Both SiLU and GLU are data-driven gates, which control the information passed in the hierarchy. Unlike GLU, where the linear and gate functions are parametrized separately, SiLU uses the same parameters to represent them. This halves the number of parameters in a single layer.

An illustration of the proposed ChimeraACVAE source model is shown on the right in Fig. 1, and the network architecture used to configure the model is shown in Fig. 2. Table I shows the number of the parameters of the CVAE, ACVAE, and ChimeraACVAE models used in the following experiments. Note that the number of parameters depend on the number of speakers in the training dataset. As can be seen from this comparison, the ChimeraACVAE source model with the above modifications has reduced the number of parameters to about 40% of the original ACVAE source model, which is even smaller than that in the CVAE model used in the MVAE method.

B. Training criterion based on KD

Since the latent variable \mathbf{z} no longer depends on \mathbf{c} , we must first rewrite the training loss of ACVAE, i.e., (19), by replacing $q_{\phi}^{*}(\mathbf{z}|\mathbf{S}, \mathbf{c})$ with $q_{\phi}^{+}(\mathbf{z}|\mathbf{S})$. Note that we omit g in this subsection, assuming that g is set to 1 and normalized

Table I: Number of parameters of CVAE, ACVAE, and ChimeraACVAE model used in the experiments.

Model	Number of parameters [M]	
	Spk-dep	Spk-ind
CVAE	10.6	12.5
ACVAE	17.0	18.9
ChimeraACVAE	7.0	7.9

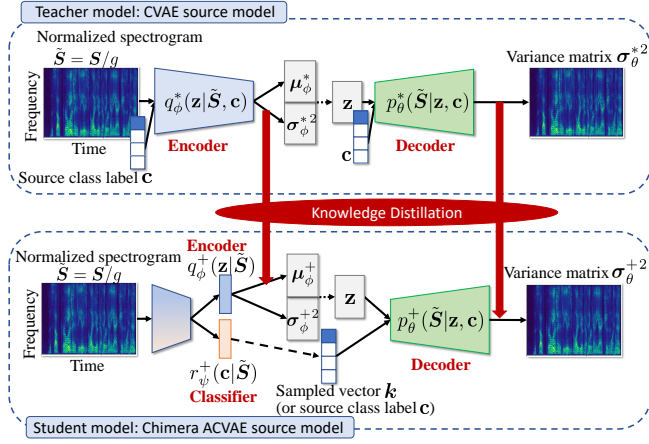


Figure 3: Illustration of the response-based KD from a pre-trained CVAE source model to the ChimeraACVAE source model. We use $\tilde{S} = S/g$ to denote a normalized spectrogram. Note that $g = 1$ is omitted from the expressions of the encoder, decoder, decoder distributions owing to space limitations.

spectrograms are used during training. Thus, the reformulated training criteria are given as

$$\mathcal{J} = \mathbb{E}_{S,c} [\mathbb{E}_{z \sim q_\phi^+(z|S)} [\log p_\theta^+(S|z,c)] - \text{KL}[q_\phi^+(z|S) \| p(z)]], \quad (24)$$

$$\mathcal{L} = \mathbb{E}_{S',z \sim q_\phi^+(z|S')} [\mathbb{E}_{c,S \sim p_\theta^+(S|z,c)} [\log r_\psi^+(c|S)']], \quad (25)$$

$$\mathcal{I} = \mathbb{E}_{S,c} [\log r_\psi^+(c|S)]. \quad (26)$$

Here, $\mathbb{E}_{S'}[\cdot]$ in (25) denote the mean of the arguments over all spectrograms $S' \sim p_D(S)$ in the training dataset. The superscript $+$ is used to distinguish the networks in the ChimeraACVAE model from those in the original ACVAE model superscripted with $*$.

Unlike in the training phase, where the class label c is known and given, in the separation phase, the spectrogram S needs to be constructed using the estimated z and c . Therefore, it is reasonable to simulate this situation in the training phase as well. Namely, we consider not only the reconstruction error defined using the given label c but also the reconstruction error defined using the estimated $c \sim r_\psi^+(c|S)$. Thus, we propose including

$$\mathcal{J}' = \mathbb{E}_{S,z \sim q_\phi^+(z|S), c \sim r_\psi^+(c|S)} [\log p_\theta^+(S|z,c)], \quad (27)$$

$$\mathcal{L}' = \mathbb{E}_{S',z \sim q_\phi^+(z|S'), c \sim r_\psi^+(c|S')} [\mathbb{E}_{S \sim p_\theta^+(S|z,c)} [\log r_\psi^+(c|S)']], \quad (28)$$

in the training objective. Here, it should be noted that both \mathcal{J}' and \mathcal{L}' involve expectations over $c \sim r_\psi^+(c|S')$. However,

there is currently no known reparametrization trick that can be applied to random variables that follow multinomial distributions. Instead, we use the Gumbel-Softmax (GS) distribution as an approximation to the multinomial distribution, which allows the use of the reparameterization trick [66], [67]. The GS distribution of a continuous multivariate variable $\mathbf{k} = [k_1, \dots, k_I]^\top$ is defined as

$$p_{\rho,\tau}(\mathbf{k}) = \Gamma(I)\tau^{I-1} \left(\sum_{i=1}^I \rho_i/k_i^\tau \right)^{-1} \prod_{i=1}^I (\rho_i/k_i^{\tau+1}). \quad (29)$$

This expression is derived analytically as a distribution that is followed by the variables

$$k_i = \frac{\exp((\log \rho_i + g_i)/\tau)}{\sum_{i'=1}^I \exp((\log \rho_{i'} + g_{i'})/\tau)} \quad (i = 1, \dots, I) \quad (30)$$

where g_i , $i = 1, \dots, I$ are Gumbel samples drawn independently and identically from $\text{Gumbel}(0, 1)$, ρ is the class probability vector produced by the classifier, and τ is called the softmax temperature. Here, it is important to note that (29) is shown to become identical to $r_\psi^+(\mathbf{k}|S')$ as τ approaches 0. By replacing $r_\psi^+(\mathbf{k}|S')$ with (29), (27) and (28) can be approximated as

$$\mathcal{J}'_{\text{GS}} = \mathbb{E}_{S,z \sim q_\phi^+(z|S), \mathbf{k} \sim p_{\rho,\tau}(\mathbf{k})} [\log p_\theta^+(S|z,\mathbf{k})], \quad (31)$$

$$\mathcal{L}'_{\text{GS}} = \mathbb{E}_{S',z \sim q_\phi^+(z|S'), \mathbf{k} \sim p_{\rho,\tau}(\mathbf{k}), S \sim p_\theta^+(S|z,\mathbf{k})} [\log r_\psi^+(\mathbf{k}|S)']. \quad (32)$$

Unlike the original expressions, these expressions allow the computations of the derivatives with respect to ψ using the reparameterization trick.

With the reduced number of model parameters, the challenge is how to make the ChimeraACVAE model have a high generalization capability. To this end, we further introduce training criteria derived based on the KD [51] using a pre-trained CVAE model as the teacher model. KD, also known as teacher-student learning, is a technique to transfer the knowledge from a teacher model to a student model, originally proposed for model compression [51] and later shown to improve the generalization capability of the student model [68]. There are three types of knowledge that can be transferred between models: response-based knowledge, feature-based knowledge, and relation-based knowledge. These refer to the knowledge of the last output layer, the knowledge of each output layer, and the knowledge of the relationship between layers, respectively. Since the networks in both the teacher and student models are reasonably shallow, we consider response-based KD to be sufficient, as it requires a minimal increase in training cost.

Specifically, we transfer the knowledge of the distributions of the latent variable $q_\phi^*(z|S,c)$ and the complex spectrograms $p_\theta^*(S|z,c)$ learned by the CVAE model into the ChimeraACVAE model by using these distributions as priors. We use the KL divergences to measure the differences between the distributions estimated by a student model and the pretrained teacher model such that

$$\mathcal{K}_z = \mathbb{E}_{S,c} [\text{KL}[q_\phi^*(z|S,c) \| q_\phi^+(z|S)']], \quad (33)$$

$$\mathcal{K}_S = \mathbb{E}_{S,c,z \sim q_\phi^*(z|S,c), z^+ \sim q_\phi^+(z|S)}$$

$$[\text{KL}[p_{\theta}^*(\mathbf{S}|\mathbf{z}^*, \mathbf{c})||p_{\theta}^+(\mathbf{S}|\mathbf{z}^+, \mathbf{c})]], \quad (34)$$

$$\mathcal{K}'_{\mathbf{S}} = \mathbb{E}_{\mathbf{S}, \mathbf{c}, \mathbf{z}^* \sim q_{\phi}^*(\mathbf{z}|\mathbf{S}, \mathbf{c}), \mathbf{z}^+ \sim q_{\phi}^+(\mathbf{z}|\mathbf{S}), \mathbf{k} \sim p_{\rho, \tau}(\mathbf{k})} [\text{KL}[p_{\theta}^*(\mathbf{S}|\mathbf{z}^*, \mathbf{c})||p_{\theta}^+(\mathbf{S}|\mathbf{z}^+, \mathbf{k})]]. \quad (35)$$

Here, (35) is a criterion that measures the difference between the teacher distribution and decoder distribution computed using the GS distribution. An illustration of KD for training the ChimeraACVAE model is shown in Fig. 3.

The total training criterion of the ChimeraACVAE is a weighted linear combination of the above-mentioned criteria:

$$\mathcal{J} + \lambda_{\mathcal{L}}\mathcal{L} + \lambda_{\mathcal{I}}\mathcal{I} + \lambda_{\mathcal{J}'}\mathcal{J}'_{\text{GS}} + \lambda_{\mathcal{L}'}\mathcal{L}'_{\text{GS}} - \lambda_{\mathcal{K}_{\mathbf{z}}}\mathcal{K}_{\mathbf{z}} - \lambda_{\mathcal{K}_{\mathbf{S}}}\mathcal{K}_{\mathbf{S}} - \lambda_{\mathcal{K}'_{\mathbf{S}}}\mathcal{K}'_{\mathbf{S}}. \quad (36)$$

Here, λ_* denotes a non-negative parameter that weighs the importance of that term.

With the trained ChimeraACVAE source model, we can use the same procedure as **Algorithm 1** to perform source separation. We call it *FastMVAE2* to distinguish it from the method using the ACVAE source model. Note that in FastMVAE2, the PoE-based update rule is no longer required thanks to the improved generalization capability, but of course it can be used in addition.

V. EXPERIMENTAL EVALUATIONS

To evaluate the effectiveness of the proposed training procedure, we compare the source separation performance in speaker-dependent and speaker-independent situations.

A. Datasets

For the speaker-dependent source separation experiment, we used speech utterances of two male speakers (SM1, SM2) and two female speakers (SF1, SF2) excerpted from the Voice Conversion Challenge (VCC) 2018 dataset [69]. The audio files for each speaker were about seven minutes long and manually segmented into 116 short sentences, where 81 and 35 sentences (about five and two minutes long, respectively) served as training and test sets, respectively. We used two-channel mixture signals of two sources as the test data, which were synthesized using simulated room impulse responses (RIRs) generated using the image method [70] and real RIRs measured in an anechoic room (ANE) and an echo room (E2A). The reverberation times (RT_{60}) [71] of the simulated RIRs were set at 78 and 351 ms, which were controlled by setting the reflection coefficient of the walls at 0.20 and 0.80, respectively. For the measured RIRs, we used the data included in the RWCP Sound Scene Database in Real Acoustic Environments [72]. The RT_{60} of ANE and E2A were 173 and 225 ms, respectively. The test data included four pairs of speakers, i.e., SF1+SF2, SF1+SM1, SM1+SM2, and SF2+SM2. For each speaker pair, we generated ten mixture signals. Hence, there were a total of 40 test signals for each reverberation condition, each of which was about four to seven seconds long.

The datasets for the speaker-independent experiment were generated in the same way by using the Wall Street Journal (WSJ0) corpus [73]. All the utterances in WSJ0 folder

si_tr_s (around 25 hours) were used as the training set, which consists of 101 speakers in total. A test set was created by randomly mixing two different speakers selected from the WSJ0 folders si_dt_05 and si_et_05, where the number of speakers was 18. We generated test data using simulated RIRs with $RT_{60} = 78$ ms and $RT_{60} = 351$ ms, where 100 mixture signals were generated under each reverberation condition. All the speech signals were resampled at 16 kHz. The STFT and inverse STFT were calculated by using a Hamming window with a length of 128 ms and half overlap.

B. Experimental settings

We chose ILRMA [7], the MVAE method [26]¹, and the FastMVAE method [46] as the baseline methods for both the speaker-dependent and speaker-independent cases, and IDLMA [30] as another baseline method only for the speaker-dependent scenario. For all the methods, the parameter optimization algorithms were run for 60 iterations, and the separation matrix $\mathbf{W}(f)$ was initialized with an identity matrix.

We set the basis number of ILRMA at 2, which is the optimal setting for speech separation. For IDLMA, we used the same network architecture and training settings as those in [30] except for the optimization algorithm, where we used Adam [74] instead of Adadelta [75]. Note that unlike other methods where speaker information is estimated, IDLMA requires speaker information in order to properly select the corresponding pre-trained network. The network architectures for the CVAE and ACVAE source models were the same as those used in [46], where the encoder consisted of 2 convolutional layers using GLU following a regular convolutional layer, the decoder consisted of 2 deconvolutional layers using GLU following a regular deconvolutional layer, and the classifier consisted of 3 convolutional layers using GLU following a regular convolutional layer. All the GLU layers used batch normalization to stabilize the training. Adam was used to train the networks and estimate \mathbf{z}_j and \mathbf{c}_j in the MVAE method. In the training of ChimeraACVAE, the weight parameters were empirically set with the KD criterion $\mathcal{K}_{\mathbf{z}}$ as 10 and the rest as 1. The temperature τ for the GS distribution was set at 1.

We calculated the source-to-distortions ratio (SDR), source-to-interferences ratio (SIR), and sources-to-artifacts ratio (SAR) [76] to evaluate the source separation performance, and used perceptual evaluation of speech quality (PESQ)² [77] and short-time objective intelligibility (STOI)³ [78] to ascertain the speech quality and intelligibility of the separated waveforms.

C. Multi-speaker separation performance

We first investigated the effectiveness of each training criterion proposed in Subsection IV-B in training the proposed ChimeraACVAE source model. The correspondence between the models and their training criteria are shown in Table

¹Code: <https://github.com/lili-0805/MVAE>

²Code: <https://github.com/vBaiCai/python-pesq>

³Code: <https://github.com/mpariente/pystoi>

Table II: Evaluated models and corresponding training criteria, which are weighted linear combinations of the equations.

Model	Training criterion
ACVAE	(24), (25), (26)
+ estimated_label	(24), (25), (26), (31), (32)
+ KD_z	(24), (25), (26), (31), (32), (33)
+ KD_S	(24), (25), (26), (31), (32), (34), (35)
+ KD_both	(24), (25), (26), (31), (32), (33), (34), (35)
+ KD_z	(24), (25), (26), (33)
+ KD_S	(24), (25), (26), (34)
+ KD_both	(24), (25), (26), (33), (34)

Table III: SDR [dB], SIR [dB], SAR [dB], PESQ, and STOI achieved by using ChimeraACVAE source model trained with different loss functions. Bold font shows the highest scores.

Scenario	Training criteria	SDR	SIR	SAR	PESQ	STOI
Spk-dep	ACVAE	10.74	16.02	13.79	2.45	0.8170
	+ estimated_label	13.29	18.87	15.87	2.64	0.8409
	+ KD_z	15.90	22.23	17.78	2.79	0.8580
	+ KD_S	13.29	18.75	16.04	2.66	0.8378
	+ KD_both	15.40	21.63	17.43	2.77	0.8565
	+ KD_z	9.89	15.48	12.79	2.38	0.8114
	+ KD_S	12.41	17.69	15.23	2.56	0.8253
	+ KD_both	8.05	12.96	11.76	2.24	0.7880
Spk-ind	ACVAE	15.81	22.73	18.60	3.14	0.8855
	+ estimated_label	12.35	18.38	16.01	3.04	0.8634
	+ KD_z	16.89	24.74	18.79	3.17	0.8917
	+ KD_S	15.18	21.95	17.99	3.12	0.8832
	+ KD_both	17.04	24.87	18.85	3.19	0.8945
	+ KD_z	16.16	23.68	18.33	3.14	0.8863
	+ KD_S	15.66	22.65	18.48	3.14	0.8893
	+ KD_both	16.07	23.47	18.39	3.14	0.8892

Table IV: Comparison of SDR [dB], SIR [dB], SAR [dB], PESQ, and STOI among compact FastMVAE, FastMVAE, and FastMVAE2 with the optimal parameter settings. Bold font shows the highest scores.

Scenario	Method	SDR	SIR	SAR	PESQ	STOI
Spk-dep	Compact FastMVAE w/o PoE	7.52	12.71	10.80	2.29	0.7916
	Compact FastMVAE w/ PoE	7.75	12.84	11.06	2.30	0.7933
	FastMVAE w/o PoE [37]	13.78	19.51	16.16	2.03	0.8465
	FastMVAE w/ PoE [37]	13.95	19.54	16.33	2.66	0.8452
	FastMVAE2	15.40	21.63	17.43	2.77	0.8565
Spk-ind	Compact FastMVAE w/o PoE	8.16	12.62	12.47	2.60	0.8119
	Compact FastMVAE w/ PoE	10.55	17.51	12.66	2.78	0.8453
	FastMVAE w/o PoE [37]	10.43	15.41	15.73	2.73	0.8358
	FastMVAE w/ PoE [37]	14.41	21.21	17.35	3.04	0.8776
	FastMVAE2	17.04	24.87	18.85	3.19	0.8945

Table V: Comparison of SDR [dB], SIR [dB], SAR [dB], PESQ, and STOI between FastMVAE2 and baseline methods with the optimal parameter settings. Bold font shows the highest scores.

Scenario	Method	SDR	SIR	SAR	PESQ	STOI
Spk-dep	ILRMA	13.62	19.79	15.83	1.92	0.8570
	IDLMA [46]	14.15	21.11	15.59	1.77	0.8692
	MVAE [46]	17.03	23.75	18.61	2.24	0.8717
	FastMVAE [46]	13.95	19.54	16.33	2.66	0.8452
	FastMVAE2	15.40	21.63	17.43	2.77	0.8565
Spk-ind	ILRMA	14.43	20.98	17.45	2.28	0.8850
	MVAE [46]	17.58	25.13	19.26	2.65	0.8934
	FastMVAE [46]	14.41	21.21	17.35	3.04	0.8776
	FastMVAE2	17.04	24.87	18.85	3.19	0.8945

II. Table III shows the results, which were calculated by averaging over the entire dataset including multiple rever-

Table VI: Lengths [sec] of mixture signals in each case.

Number of sources	Minimum	Maximum	Average
2	5.70	13.86	8.56
3	8.71	13.68	11.47
6	9.04	16.23	12.76
9	9.49	16.33	12.60
12	10.48	15.32	12.77
15	11.75	14.71	13.12
18	11.43	15.83	13.51

beration conditions. The results show that it is effective to further exploit the reconstruction loss and classification loss of the spectrograms reconstructed with the estimated class label in the speaker-dependent scenario, where small amounts of training data were available. Comparing the models trained without KD (1st and 2nd rows) with that trained with KD (3th to 5th rows), we found an improvement in SDR of about 2.6 dB in speaker-dependent situations and more than 1 dB in speaker-independent ones, which confirmed that KD can significantly improve source separation performance. In particular, knowledge transfer of the distribution of the latent variable \mathbf{z} was effective in stabilizing the inference accuracy even for unseen speakers. A further improvement was achieved in the speaker-independent setting by transferring knowledge between distributions of generated complex spectrograms, but no improvement was seen in the speaker-dependent setting.

In Table IV, we show a comparison of source separation performance between the FastMVAE and FastMVAE2 methods. To demonstrate the effectiveness of the proposed training criterion, we trained an ACVAE with the architecture that respectively replaces BN and GLU with LN and SiLU, which is referred to as “compact FastMVAE”. The results of FastMVAE and compact FastMVAE indicate that the replacement of the normalization method and nonlinear activation did not lead to an improvement of the source separation performance. Therefore, the performance improvement by FastMVAE2 can be attributed mainly to the proposed training criterion. The FastMVAE2 method obtained the highest scores in terms of all the criteria. Particularly, FastMVAE2 achieved an SDR improvement of about 6.6 and 2.6 dB from the FastMVAE without and with PoE, respectively. These results indicated that the ChimeraACVAE source model had good generalization to unseen data, which made the FastMVAE2 method able to handle speaker-independent scenario without the heuristic inference method. Table V shows the average scores achieved by each method with their optimal parameter settings. The proposed method significantly outperformed ILRMA and the FastMVAE method, and narrowed the performance gap with the MVAE method.

D. Comparison of computational time in situations with more sources and channels

In this subsection, we investigate the computational time of each method. We conducted speaker-independent experiments with more sources and channels, and compared the computation time of each method for each update iteration and overall processing time.

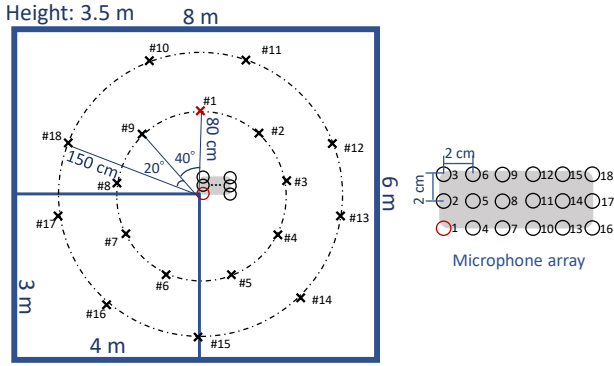


Figure 4: Configuration of sources and microphone array, where red points represent the first microphone and source.

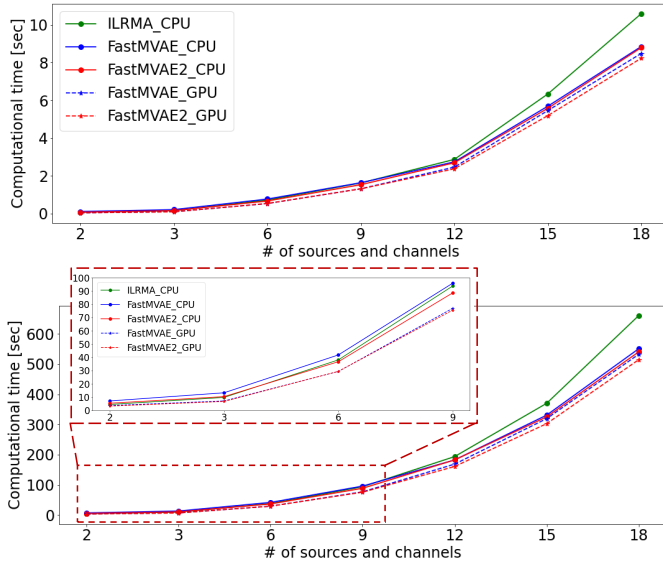


Figure 5: Average inference time [sec] of each iteration (upper) and overall processing (bottom).

As in the above speaker-independent experiment, the simulated RIRs in the $\{2, 3, 6, 9, 12, 15, 18\}$ -channel cases were generated using the image method [70] with the reflection coefficient of the walls set at 0.20. The details of the room configuration and microphone array are shown in Fig. 4. In each case, more sound sources and microphones were added and placed in the order of increasing numbers. Speech utterances were randomly selected from the WSJ0 folders `si_dt_05` and `si_et_05`. We generated 10 samples for each case. The minimum, maximum, and average lengths of the mixture signals are shown in Table VI. The average SDR of the generated mixture signals for each case is shown in the first row of Table VIII. All algorithms were processed using an Intel(R) Xeon(R) Gold 6130 CPU @ 2.10GHz and a Tesla V100 GPU. Other experimental settings were the same as those in the above speaker-independent experiment.

The inference times of ILRMA, FastMVAE, and FastMVAE2 are shown in Fig. 5, and those of MVAE are shown in Table VII as a reference. The fast algorithms performed extremely fast by using a GPU. Comparing the computation

Table VII: Average inference time [sec] of MVAE.

Type	Number of sources and channels						
	2	3	6	9	12	15	18
Each iteration	0.70	1.05	2.65	4.36	9.24	10.43	14.03
Overall processing	43.72	65.11	155.77	266.80	478.08	583.02	872.83

times in the CPU, we found that the FastMVAE2 method achieved runtimes comparable to ILRMA in the 2-source and 3-source cases, and faster than ILRMA in cases with more than 3 sources. This indicates that the proposed method is more efficient in situations with a large number of sources and microphones. The average SDR scores obtained by each method are shown in Table VIII. The proposed FastMVAE2 outperformed ILRMA and the FastMVAE without PoE, and even outperformed the MVAE method in the 2-source case, demonstrating the effectiveness of the proposed ChimeraAC-VAE source model. Note that although the performance of ILRMA was superior to the proposed method in the cases of 3 and 6 sources, this might change with different initialization of the basis and activation matrices of the NMF. On the other hand, the performance of the proposed method is independent of the initialization. We show an example of the magnitude spectrograms of separated signals obtained by ILRMA, MVAE, and FastMVAE2 with their corresponding ground truth signals in Fig. 6⁴. We found that although the MVAE and FastMVAE2 methods suffered from the phenomenon called block permutation [79], [80], in which the permutations in different frequency blocks are inconsistent, the deep generative model-based source models improved the estimation accuracy in the low-frequency band (0-2 kHz), which resulted in a more remarkable SDR improvement compared with ILRMA.

E. Spatialized-WSJ0-2mix benchmark

In this subsection, we evaluate the proposed FastMVAE2 in the spatialized WSJ0-2mix benchmark [18], which is widely used for evaluating the recent DNN-based methods. There are 20,000 (~ 30 h), 5,000 (~ 10 h), and 3000 (~ 5 h) utterances in the training, validation, and test sets. The training and validation mixtures were generated from data in `si_tr_s` folder and the test mixtures were generated from data in `si_dt_05` and `si_et_05` folders. Therefore, the speaker-independent settings mentioned in the above experiments are still valid. RIR used for every utterance was simulated with a random configuration, including room characteristics, speaker locations, and microphone geometry. RT_{60} for the reverberant case was randomly selected from 200 to 600 ms.

We compared FastMVAE2 with (1) oracle ideal binary mask (IBM), (2) oracle ideal ratio mask (IRM), (3) oracle mask-based minimum variance distortionless response (MVDR) beamformer [81], (4) oracle signal-based MVDR, 5) time-domain audio separation network (TasNet) [19], (6) multi-channel TasNet [23], and (7) Beam-TasNet [23]. The oracle IBM and IRM were obtained using the first channel of the spatialized clean sources and applied to the first channel of observed mixture signals. The difference between the oracle

⁴Audio samples are available at <http://www.kecl.ntt.co.jp/people/kameoka.hirokazu/Demos/mvae-ss/index.html>

Table VIII: Comparison of SDR [dB] between FastMVAE2 and baseline methods with the optimal parameter settings in situations with different numbers of sources and channels. Values in parentheses indicate the improvement over unprocessed. Bold font shows the highest scores.

Method	Number of sources and channels						
	2	3	6	9	12	15	18
Unprocessed	0.09	-3.92	-8.13	-10.45	-12.15	-13.03	-13.86
ILRMA	20.89 (20.80)	23.04 (26.96)	7.54 (15.67)	1.61 (12.06)	-0.11 (12.04)	-3.79 (9.24)	-5.92 (7.94)
MVAE	26.63 (26.54)	25.17 (29.09)	11.32 (19.45)	9.26 (19.71)	7.34 (19.49)	5.00 (18.03)	2.58 (16.34)
FastMVAE w/o PoE	15.77 (15.68)	7.59 (11.51)	3.32 (11.45)	4.23 (14.68)	0.69 (12.84)	-0.06 (12.97)	-1.96 (11.90)
FastMVAE2	28.58 (28.49)	21.50 (25.41)	6.53 (14.66)	5.77 (16.23)	4.04 (16.19)	2.90 (15.94)	0.22 (14.08)

mask-based and signal-based MVDR was the signal used for computing spatial covariance matrices, where the former used the multichannel IRM of each source and the later directly used the clean reverberant speech of each source. We investigated window lengths of 128 ms and 512 ms. Settings for TasNet, multichannel TasNet, and Beam-TasNet are available in [23]⁵. One important factor here is the window length. Beam-TasNet used a length of 512 ms to meet the instantaneous mixture model for reverberant signals, while the proposed method used that of 128 ms since the motivation of the FastMVAE2 is to bridge the high performance of MVAE and real-time applications with low latency.

We first show the results of the spatialized anechoic WSJ0-2mix dataset in Table IX. With the anechoic setup, beamforming algorithms achieved higher performance than mask-based methods. From these results, we confirmed that the MVAE and proposed FastMVAE2 achieved even better performance than the oracle mask-based MVDR beamformer, indicating the effectiveness of these two methods when the instantaneous mixture model assumption is satisfied. Next, we show the results of the spatialized reverberant WSJ0-2mix dataset in Table X. The performance of the MVAE and FastMVAE2 degraded significantly due to reverberations. The main reason was the instantaneous mixture model assumption, which was not satisfied anymore with heavy reverberation and short window length. We found that even the performance of oracle MVDRs degraded significantly when the window length became shorter. Two promising approaches can be considered to deal with this problem, including using longer window length and performing separation along with dereverberation [40], [82], [83]. It is straightforward that longer window length helps deal with heavy reverberant conditions, which has also been confirmed from the results of oracle MVDRs with longer window length and Beam-TasNet. However, a longer window length is undesirable and should be avoided in real-time applications because it increases algorithmic latency. Therefore, we consider the second approach, performing separation and dereverberation simultaneously, as one direction of our future works to overcome this limitation of the FastMVAE2 method.

VI. CONCLUSION

In this paper, we proposed an improved ACVAE source model named “ChimeraACVAE” source model for the fast algorithm of the MVAE method, which we call “FastMVAE2”.

⁵We would like to appreciate Dr. Tsubasa Ochiai from NTT Communication Science Laboratories for providing us with the test dataset and evaluation script so that we could compare our methods with results reported in [23].

Table IX: Comparison of SDR [dB] for spatialized anechoic WSJ0-2mix dataset. “1ch” and “2ch” indicate the number of channels used for processing.

Method	window length [ms]	SDR [dB]
Mixture	—	-0.4
Oracle IBM (1ch)	128	13.66
Oracle IRM (1ch)	128	13.55
Oracle mask-based MVDR (2ch)	128	23.26
Oracle signal-based MVDR (2ch)	128	39.68
Oracle mask-based MVDR (2ch)	512	15.66
Oracle signal-based MVDR (2ch)	512	48.02
MVAE (2ch)	128	28.49
FastMVAE2 (2ch)	128	31.31

Table X: Comparison of SDR [dB] for spatialized reverberant WSJ0-2mix dataset. “1ch” and “2ch” indicate the number of channels used for processing.

Method	window length [ms]	SDR [dB]
Mixture	—	0.1
Oracle IBM (1ch)	128	13.41
Oracle IRM (1ch)	128	13.29
Oracle mask-based MVDR (2ch)	128	8.16
Oracle signal-based MVDR (2ch)	128	8.14
Oracle mask-based MVDR (2ch)	512	11.95
Oracle signal-based MVDR (2ch)	512	16.32
TasNet (1ch) [23]	—	11.3
Multichannel TasNet (2ch) [23]	—	12.7
Beam-TasNet (1ch) [23]	512	12.9
Beam-TasNet (2ch) [23]	512	13.8
MVAE (2ch)	128	5.35
FastMVAE2 (2ch)	128	6.02

ChimeraACVAE is a more compact source model that consists of a unified encoder and classifier network and a decoder, which are composed of fully convolutional layers with layer normalization and an SiLU activation function. The KD framework was applied to train the ChimeraACVAE source model to improve the generalization capability to unseen data. The experimental results demonstrated that the FastMVAE2 method achieved significant performance improvement in both speaker-dependent and speaker-independent multispeaker separation tasks, approaching the performance that of the MVAE method. Moreover, the proposed method significantly reduced the model size and improved the computational efficiency, which achieved computational time comparable to ILRMA in cases of two and three sources and lower computational time than ILRMA in cases of more sources.

REFERENCES

- [1] L. Li, H. Kameoka, and S. Makino, “FastMVAE2: On improving and accelerating the fast variational autoencoder-based source separation algorithm for determined mixtures,” arXiv: 2109.13496, 2021.

- [2] T. Kim, T. Eltoft and T.-W. Lee, “Independent vector analysis: An extension of ICA to multivariate components,” in *Proc. ICA*, pp. 165–172, 2006.
- [3] A. Hiroe, “Solution of permutation problem in frequency domain ICA using multivariate probability density functions,” in *Proc. ICA*, pp. 601–608, 2006.
- [4] A. Ozerov and C. Févotte, “Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation,” *IEEE Trans. ASLP*, vol. 18, no. 3, pp. 550–563, 2010.
- [5] H. Sawada, H. Kameoka, S. Araki and N. Ueda, “Multichannel extensions of non-negative matrix factorization with complex-valued data,” *IEEE Trans. ASLP*, vol. 21, no. 5, pp. 971–982, 2013.
- [6] H. Kameoka, T. Yoshioka, M. Hamamura, J. Le. Roux and K. Kashino, “Statistical model of speech signals based on composite autoregressive system with application to blind source separation,” in *Proc. LVA/ICA*, pp. 245–253, 2010.
- [7] D. Kitamura, N. Ono, H. Sawada, H. Kameoka and H. Saruwatari, “Determined blind source separation unifying independent vector analysis and nonnegative matrix factorization,” *IEEE/ACM Trans. ASLP*, vol. 24, no. 9, pp. 1622–1637, 2016.
- [8] D. Kitamura, N. Ono, H. Sawada, H. Kameoka, and H. Saruwatari, “Determined blind source separation with independent low-rank matrix analysis,” in *Audio Source Separation*, S. Makino, Ed., Springer, pp. 125–155, 2018.
- [9] D. D. Lee and H. S. Seung, “Algorithms for non-negative matrix factorization,” in *Adv. NIPS*, pp. 556–562, 2001.
- [10] D. Wang and J. Chen, “Supervised speech separation based on deep learning: An overview,” *IEEE/ACM Trans. ASLP*, vol. 26, no. 10, pp. 1702–1726, 2018.
- [11] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, “Deep clustering: Discriminative embeddings for segmentation and separation,” in *Proc. ICASSP*, pp. 31–35, 2016.
- [12] Y. Isik, J. Le Roux, Z. Chen, S. Watanabe, and J. R. Hershey, “Single-channel multi-speaker separation using deep clustering,” in *Proc. Interspeech*, pp. 545–549, 2016.
- [13] D. Yu, M. Kolbæk, Z. H. Tan, and J. Jensen, “Permutation invariant training of deep models for speaker-independent multi-talker speech separation,” in *Proc. ICASSP*, pp. 241–245, 2017.
- [14] M. Kolbæk, D. Yu, Z. H. Tan, and J. Jensen, “Multitalker speech separation with utterance-level permutation invariant training of deep recurrent neural networks,” *IEEE/ACM Trans. ASLP*, vol. 25, no. 10, pp. 1901–1913, 2017.
- [15] Y. Liu and D. Wang, “Divide and conquer: A deep CASA approach to talker-independent monaural speaker separation,” *IEEE/ACM Trans. ASLP*, vol. 27, no. 12, pp. 2092–2102, 2019.
- [16] J. Le Roux, G. Wichern, S. Watanabe, S. Sarroff, and J. R. Hershey, “Phasebook and friends: Leveraging discrete representations for source separation,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 2, pp. 370–382, 2019.
- [17] M. Delfarah and D. Wang, “Deep learning for talker-dependent reverberant speaker separation: An empirical study,” *IEEE/ACM Trans. ASLP*, vol. 27, no. 11, pp. 1839–1848, 2019.
- [18] Z.-Q. Wang, J. Le Roux, and J. R. Hershey, “Multi-channel deep clustering: Discriminative spectral and spatial embeddings for speaker independent speech separation,” in *Proc. ICASSP*, pp. 1–5, 2018.
- [19] Y. Luo and N. Mesgarani, “Conv-TasNet: Surpassing ideal time–frequency magnitude masking for speech separation,” *IEEE/ACM TASLP*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [20] E. Nachmani, Y. Adi, and L. Wolf, “Voice separation with an unknown number of multiple speakers,” in *Proc. ICML*, pp. 7164–7175, 2020.
- [21] C. Subakan, M. Ravanelli, S. Cornell, M. Bronzi, and J. Zhong, “Attention is all you need in speech separation,” in *Proc. ICASSP*, pp. 21–25, 2021.
- [22] T. Higuchi, K. Kinoshita, M. Delcroix, K. Zmolíková, and T. Nakatani, “Deep clustering-based beamforming for separation with unknown number of sources,” in *Proc. Interspeech*, pp. 1183–1187, 2017.
- [23] T. Ochiai, M. Delcroix, R. Ikeshita, K. Kinoshita, T. Nakatani, and S. Araki, “Beam-TasNet: Time-domain audio separation network meets frequency-domain beamformer,” in *Proc. ICASSP*, pp. 6384–6388, 2020.
- [24] A. A. Nugraha, A. Liutkus and E. Vincent, “Multichannel Audio Source Separation With Deep Neural Networks,” *IEEE/ACM Trans. ASLP*, vol. 24, no. 9, pp. 1652–1664, 2016.
- [25] N. Makishima, S. Mogami, N. Takamune, D. Kitamura, H. Sumino, S. Takamichi, H. Saruwatari, and N. Ono, “Independent deeply learned matrix analysis for determined audio source separation,” *IEEE/ACM Trans. ASLP*, vol. 27, no. 10, pp. 1601–1615, 2019.
- [26] H. Kameoka, L. Li, S. Inoue, and S. Makino, “Supervised determined source separation with multichannel variational autoencoder,” *Neural Computation*, vol. 31, no. 9, pp. 1891–1914, 2019.
- [27] Y. Bando, M. Mimura, K. Itoyama, K. Yoshii and T. Kawahara, “Statistical speech enhancement based on probabilistic integration of variational autoencoder and non-negative matrix factorization,” in *Proc. ICASSP*, pp. 716–720, 2018.
- [28] S. Leglaive, L. Girin and R. Horaud, “A variance modeling framework based on variational autoencoders for speech enhancement,” in *Proc. MLSP*, 2018.
- [29] L. Li, H. Kameoka and S. Makino, “Determined audio source separation with multichannel star generative adversarial network,” in *Proc. MLSP*, 2020.
- [30] S. Mogami, H. Sumino, D. Kitamura, N. Takamune, S. Takamichi, H. Saruwatari, and N. Ono, “Independent deeply learned matrix analysis for multichannel audio source separation,” in *Proc. EUSIPCO*, pp. 1571–1575, 2018.
- [31] D. P. Kingma, S. Mohamed, D. J. Rezende and M. Welling, “Semi-supervised learning with deep generative models,” in *Adv. NIPS*, pp. 3581–3589, 2014.
- [32] K. Sohn, H. Lee and X. Yan, “Learning structured output representation using deep conditional generative models,” in *Adv. NIPS*, pp. 3483–3491, 2015.
- [33] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville and Y. Bengio, “Generative Adversarial Networks,” in *Adv. NIPS*, pp. 2672–2680, 2014.
- [34] D. Rezende and S. Mohamed, “Variational inference with normalizing flows,” in *Proc. PMLR*, pp. 1530–1538, 2015.
- [35] K. Sekiguchi, Y. Bando, K. Yoshii, and T. Kawakara, “Bayesian multi-channel speech enhancement with a deep speech prior,” in *Proc. APSIPA*, pp. 1233–1239, 2018.
- [36] K. Sekiguchi, Y. Bando, A. A. Nugraha, K. Yoshii, and T. Kawahara, “Semi-supervised multichannel speech enhancement with a deep speech prior,” *IEEE/ACM Trans. ASLP*, vol. 27, no. 12, pp. 2197–2212, 2019.
- [37] S. Leglaive, U. Simsekli, A. Liutkus, L. Girin, and R. Horaud, “Speech enhancement with variational autoencoders and alpha-stable distributions,” in *Proc. ICASSP*, pp. 541–545, 2019.
- [38] S. Leglaive, L. Girin, and R. Horaud, “Semi-supervised multichannel speech enhancement with variational autoencoders and non-negative matrix factorization,” in *Proc. ICASSP*, pp. 101–105, 2019.
- [39] S. Seki, H. Kameoka, L. Li, T. Toda, and K. Takeda, “Underdetermined source separation based on generalized multichannel variational autoencoder,” *IEEE Access*, vol. 7, 168104–168115, 2019.
- [40] S. Inoue, H. Kameoka, L. Li, S. Seki, and S. Makino, “Joint separation and dereverberation of reverberant mixtures with multichannel variational autoencoder,” in *Proc. ICASSP*, pp. 96–100, 2019.
- [41] A. A. Nugraha, K. Sekiguchi, M. Fontaine, Y. Bando and K. Yoshii, “Flow-based independent vector analysis for blind source separation,” *IEEE Signal Processing Letters*, vol. 27, pp. 2173–2177, 2020.
- [42] J. Neri, R. Badeau, and P. Depalle, “Unsupervised blind source separation with variational auto-encoders,” in *Proc. EUSIPCO*, pp. 311–315, 2021.
- [43] Y. Bando, K. Sekiguchi, Y. Masuyama, A. A. Nugraha, M. Fontaine, and K. Yoshii, “Neural full-rank spatial covariance analysis for blind source separation,” *IEEE SPL*, vol. 28, pp. 1670–1674, 2021.
- [44] N. Ono, “Stable and fast update rules for independent vector analysis based on auxiliary function technique,” in *Proc. WASPAA*, pp. 189–192, 2011.
- [45] W. Sun and Y. X. Yuan, “Optimization theory and methods: nonlinear programming,” *Springer Science & Business Media*, 2006.
- [46] L. Li, H. Kameoka, S. Inoue and S. Makino, “FastMVAE: A fast optimization algorithm for the multichannel variational autoencoder method,” *IEEE Access*, vol. 8, pp. 228740–228753, 2020.
- [47] H. Kameoka, T. Kaneko, K. Tanaka, and N. Hojo, “ACVAE-VC: Non-parallel voice conversion with auxiliary classifier variational autoencoder,” *IEEE/ACM Trans. ASLP*, vol. 27, no. 9, pp. 1432–1443, 2019.
- [48] G. E. Hinton, “Training products of experts by minimizing contrastive divergence,” *Neural computation*, vol. 14, no. 8, pp. 1771–1800, 2002.
- [49] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [50] A. Krogh and J. A. Hertz “A simple weight decay can improve generalization,” *Proc. NIPS*, pp. 950–957, 1992.
- [51] G. Hinton, O. Vinyals and J. Dean, “Distilling the knowledge in a neural network,” *preprint arXiv:1503.02531*, Mar. 9, 2015.

- [52] Y. Zhang and Q. Yang, "A survey on multi-task learning," *IEEE Trans. Knowledge and Data Engineering*, 2021.
- [53] C. Shorten and T. M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of Big Data*, vol. 6, no. 1, pp. 1–48, 2019.
- [54] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2009.
- [55] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. PMLR*, pp. 448–456, 2015.
- [56] J. L. Ba, J. R. Kiros and G. E. Hinton, "Layer normalization," *preprint arXiv:1607.06450*, Jul 21, 2016.
- [57] C. Févotte and J. F. Cardoso, "Maximum likelihood approach for blind audio source separation using time-frequency Gaussian models," in *Proc. WASPAA*, pp. 78–81, 2005.
- [58] E. Vincent, S. Arberet and R. Gribonval, "Underdetermined instantaneous audio source separation via local Gaussian modeling," in *Proc. ICA*, pp. 775–782, 2009.
- [59] M. Z. Ikram and D. R. Morgan, "A beamforming approach to permutation alignment for multichannel frequency-domain blind speech separation," in *Proc. ICASSP*, pp. 881–884, 2002.
- [60] H. Sawada, R. Mukai, S. Araki, and S. Makino, "A robust and precise method for solving the permutation problem of frequency-domain blind source separation," *IEEE Trans. SAP*, vol. 12, no. 5, pp. 530–538, 2004.
- [61] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *preprint arXiv:1312.6114v10*, May 1, 2014.
- [62] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever and P. Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," in *Adv. NIPS*, pp. 2172–2180, 2016.
- [63] J. L. Ba, J. R. Kiros and G. E. Hinton, "Layer normalization," *preprint arXiv:1607.06450*, Jul 21, 2016.
- [64] P. Ramachandran, B. Zoph, and Q. V. Le, "Swish: a self-gated activation function," *preprint arXiv:1710.05941*, Oct. 16, 2017.
- [65] Y. N. Dauphin, A. Fan, M. Auli and D. Grangier, "Language modeling with gated convolutional networks," in *Proc. ICML*, pp. 933–941, 2017.
- [66] E. Jang, S. Gu, and B. Poole, "Categorical reparametrization with gumbel-softmax," in *Proc. ICLR*, 2017.
- [67] C. J. Maddison, A. Mnih, and Y. W. Teh, "The concrete distribution: A continuous relaxation of discrete random variables," in *Proc. ICLR*, 2017.
- [68] P. Shen, X. Lu, S. Li and H. Kawai, "Feature representation of short utterances based on knowledge distillation for spoken language identification," in *Proc. Interspeech*, pp. 1813–1817, 2018.
- [69] J. Lorenzo-Trueba, J. Yamagishi, T. Toda, D. Saito, F. Villavicencio, T. Kinnunen and Z. Ling, "The voice conversion challenge 2018: Promoting development of parallel and nonparallel methods," *preprint arXiv: 1804.04262*, Apr. 2018.
- [70] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating smallroom acoustics," *J. Acoust. Soc. Am.*, vol. 65, no. 4, pp. 943–950, 1979.
- [71] M. R. Schroeder, "New method of measuring reverberation time," *J. Acoust. Soc. Am.*, vol. 37, no. 3, pp. 409–412, 1965.
- [72] S. Nakamura, K. Hiyane, F. Asano and T. Endo, (1999) "Sound scene data collection in real acoustical environments," *J. Acoust. Soc. Jpn. (E)*, vol. 20, no. 3, pp. 225–231, 1999.
- [73] J. S. Garofolo, et al. CSR-I (WSJ0) Complete LDC93S6A. Web Download. Philadelphia: Linguistic Data Consortium, 1993.
- [74] D. Kingma, J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015.
- [75] M. D. Zeiler, "Adadelta: an adaptive learning rate method," *preprint arXiv: 1212.5701*, Dec. 2012.
- [76] 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.
- [77] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, "Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs," in *Proc. ICASSP*, Cat. No. 01CH37221, vol. 2, pp. 749–752, 2001.
- [78] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "A short-time objective intelligibility measure for time-frequency weighted noisy speech," in *Proc. ICASSP*, pp. 4214–4217, 2010.
- [79] Y. Liang, S. M. Naqvi, and J. Chambers, "Overcoming block permutation problem in frequency domain blind source separation when using AuxIVA algorithm," *Electronics letters*, vol. 48, no. 8, pp. 460–462, 2012.
- [80] Y. Mitsui, N. Takamune, D. Kitamura, H. Saruwatari, Y. Takahashi, and K. Kondo, "Vectorwise coordinate descent algorithm for spatially regularized independent low-rank matrix analysis," in *Proc. ICASSP*, pp. 746–750, 2018.
- [81] H. L. Van Trees, *Optimum array processing: Part IV of detection, estimation, and modulation theory*, John Wiley & Sons, 2004.
- [82] T. Yoshioka, T. Nakatani, M. Miyoshi, and H. G. Okuno, "Blind separation and dereverberation of speech mixtures by joint optimization," *IEEE Trans. ASLP*, vol. 19, no. 1, pp. 69–84, 2011.
- [83] H. Kagami, H. Kameoka, and M. Yukawa, "Joint separation and dereverberation of reverberant mixtures with determined multichannel non-negative matrix factorization," in *Proc. ICASSP*, pp. 31–35, 2018.

PLACE
PHOTO
HERE

the Acoustical Society of Japan and the Signal Processing Society Japan Student Conference Paper Award.

PLACE
PHOTO
HERE

Hirokazu Kameoka received the B.E., M.S. and Ph.D. degrees from the University of Tokyo, Japan, in 2002, 2004, and 2007, respectively. He is currently a Senior Distinguished Researcher at NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation and an Adjunct Associate Professor at the National Institute of Informatics. From 2011 to 2016, he was an Adjunct Associate Professor at the University of Tokyo. His research interests include audio, speech, and music signal processing and machine learning. He has been an associate editor of the IEEE/ACM Transactions on Audio, Speech, and Language Processing since 2015, a Member of IEEE Audio and Acoustic Signal Processing Technical Committee since 2017, and a Member of IEEE Machine Learning for Signal Processing Technical Committee since 2019. He has received 17 awards, including the IEEE Signal Processing Society 2008 SPS Young Author Best Paper Award. He is the author or co-author of about 150 articles in journal papers and peer-reviewed conference proceedings.

Li Li received the B.E. degree from Shanghai University of Finance and Economics, China, in 2014, and the M.S. and Ph.D. degrees from the University of Tsukuba, Japan, in 2018 and 2021, respectively. She is currently a postdoctoral researcher with Nagoya University and an adjunct researcher with NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation. Her research interests include audio and speech signal processing, source separation, and machine learning. She received the Student Presentation Award from the Acoustical Society of Japan and the Signal Processing Society Japan Student Conference Paper Award.

PLACE
PHOTO
HERE

Shoji Makino received his B.E., M.E., and Ph.D. degrees from Tohoku University, Japan, in 1979, 1981, and 1993, respectively. He joined NTT in 1981. He is now a Professor at the University of Tsukuba. His research interests include adaptive filtering technologies, the realization of acoustic echo cancellation, blind source separation of convolutive mixtures of speech, and acoustic signal processing for speech and audio applications.

He received the ICA Unsupervised Learning Pioneer Award in 2006, the IEEE MLSP Competition Award in 2007, the IEEE SPS Best Paper Award in 2014, the Achievement Award for Science and Technology by the Minister of Education, Culture, Sports, Science and Technology in 2015, the Hoko Award of the Hattori Hokokai Foundation in 2018, the Outstanding Contribution Award of the IEICE in 2018, the Technical Achievement Award of the IEICE in 2017 and 1997, the Outstanding Technological Development Award of the ASJ in 1995, and 8 Best Paper Awards. He is the author or co-author of more than 200 articles in journals and conference proceedings and is responsible for more than 150 patents. He was a Keynote Speaker at ICA2007 and a Tutorial Speaker at EMBC2013, Interspeech2011 and ICASSP2007.

He has served on IEEE SPS Board of Governors (2018-20), Technical Directions Board (2013-14), Awards Board (2006-08), Conference Board (2002-04), and Fellow Evaluation Committee (2018-20). He was a member of the IEEE Jack S. Kilby Signal Processing Medal Committee (2015-18) and the James L. Flanagan Speech & Audio Processing Award Committee (2008-11). He was an Associate Editor of the IEEE Transactions on Speech and Audio Processing (2002-05) and an Associate Editor of the EURASIP Journal on Advances in Signal Processing (2005-12). He was a Guest Editor of the Special Issue of the IEEE Signal Processing Magazine (2013-14). He was the Chair of SPS Audio and Acoustic Signal Processing Technical Committee (2013-14) and the Chair of the Blind Signal Processing Technical Committee of the IEEE Circuits and Systems Society (2009-10). He was the General Chair of IWAENC 2018, WASPAA2007, IWAENC2003, the Organizing Chair of ICA2003, and is the designated Plenary Chair of ICASSP2012. Dr. Makino is an IEEE SPS Distinguished Lecturer (2009-10), an IEEE Fellow, an IEICE Fellow, a Board member of the ASJ, and a member of EURASIP.

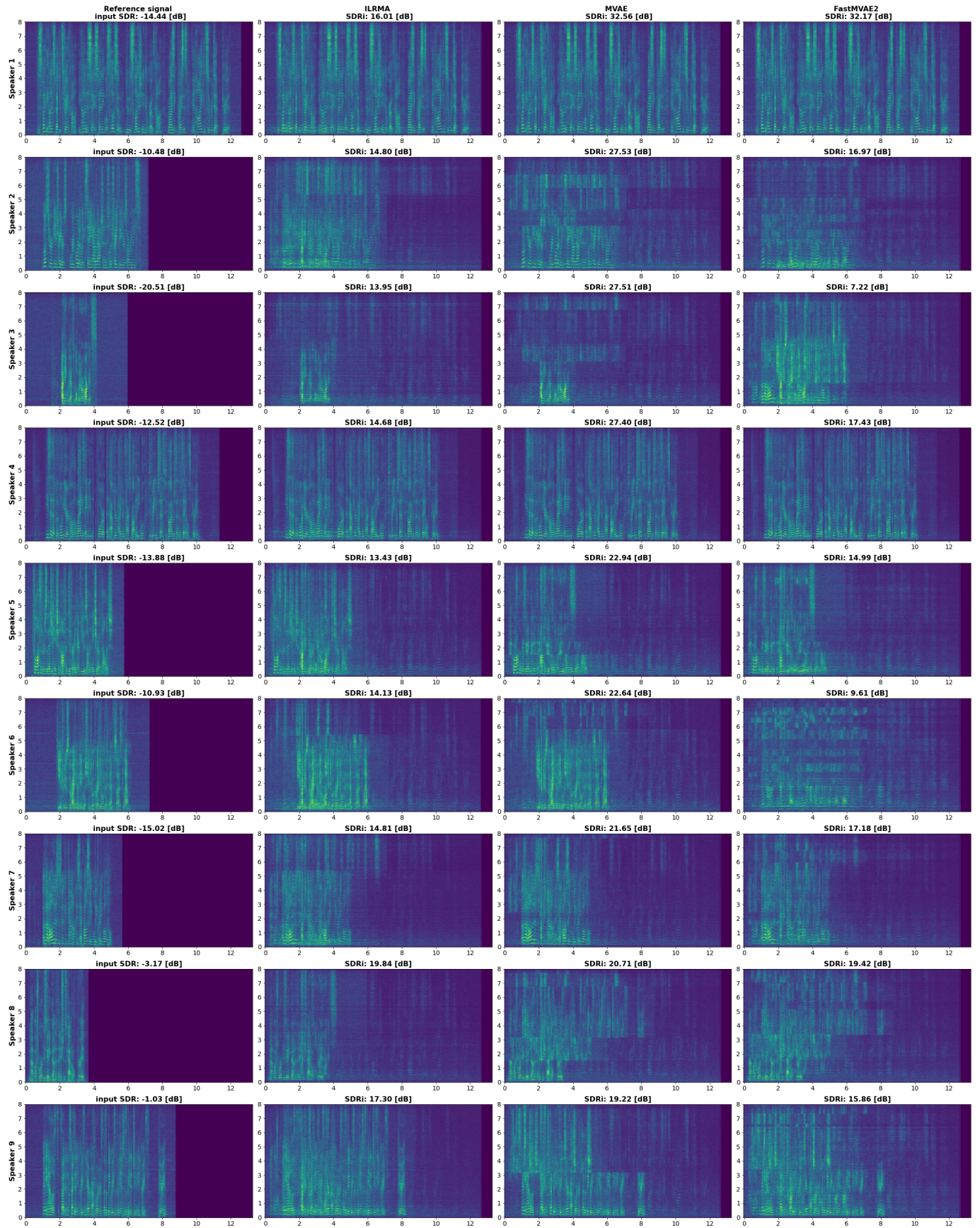


Figure 6: Magnitude spectrograms of ground truth signals (first column) and separated signals obtained by ILRMA (second column), MVAE (third column), and FastMVAE2 (fourth column) in a nine-source case. SDR of input mixture signal with respect to each speaker is shown in the top of figures in first column and SDR improvement achieved by each method is shown in the top of each figure in second to fourth. The x and y axes of each figure denote time [sec] and frequency [kHz], respectively. Audio samples are available at <http://www.kecl.ntt.co.jp/people/kameoka.hirokazu/Demos/mvae-ss/index.html>.