# Complex ISNMF: a Phase-Aware Model for Monaural Audio Source Separation

Paul Magron, Tuomas Virtanen, *Senior Member, IEEE*

*Abstract*—This paper introduces a phase-aware probabilistic model for audio source separation. Classical source models in the short-time Fourier transform domain use circularly-symmetric Gaussian or Poisson random variables. This is equivalent to assuming that the phase of each source is uniformly distributed, which is not suitable for exploiting the underlying structure of the phase. Drawing on preliminary works, we introduce here a Bayesian anisotropic Gaussian source model in which the phase is no longer uniform. Such a model permits us to favor a phase value that originates from a signal model through a Markov chain prior structure. The variance of the latent variables are structured with nonnegative matrix factorization (NMF). The resulting model is called complex Itakura-Saito NMF (ISNMF) since it generalizes the ISNMF model to the case of non-isotropic variables. It combines the advantages of ISNMF, which uses a distortion measure adapted to audio and yields a set of estimates which preserve the overall energy of the mixture, and of complex NMF, which enables one to account for some phase constraints. We derive a generalized expectation-maximization algorithm to estimate the model parameters. Experiments conducted on a musical source separation task in a semi-informed setting show that the proposed approach outperforms state-of-the-art phase-aware separation techniques.

*Index Terms*—Nonnegative matrix factorization (NMF), complex NMF, anisotropic Gaussian model, Itakura-Saito divergence, Bayesian inference, phase recovery, audio source separation.

## I. INTRODUCTION

THE goal of audio source separation [1] is to extract underlying *sources* that add up to form an observable audio *mixture*. In this paper, we address the problem of *monaural* source separation, which means that the observed audio signal has been recorded through a single microphone.

To tackle this issue, many techniques act on a time-frequency (TF) representation of the data, such as the short-time Fourier transform (STFT), since the structure of audio signals is more prominent in that domain. In particular, nonnegative matrix factorization (NMF) [2] techniques have shown successful for audio source separation [3], [4]. NMF is a rank-reduction method used for obtaining part-based decompositions of nonnegative data. The NMF problem is expressed as follows: given a matrix $\mathbf{V}$ of dimensions $F \times T$ with nonnegative entries, find a factorization $\mathbf{V} \approx \mathbf{WH}$ where $\mathbf{W}$ and $\mathbf{H}$ are nonnegative matrices of dimensions $F \times K$ and $K \times T$ respectively. To reduce the dimensionality of the data, the rank $K$ is generally chosen so that $K(F + T) \ll FT$. In audio applications $\mathbf{V}$ is usually a magnitude or power

spectrogram, and one can interpret $\mathbf{W}$ as a dictionary of spectral templates and $\mathbf{H}$ as a matrix of temporal activations.

Such a factorization is generally obtained by minimizing a cost function that penalizes the error between $\mathbf{V}$ and $\mathbf{WH}$. Popular choices are the Euclidean distance or Kullback-Leibler (KL) [2] and Itakura-Saito (IS) divergences [4]. NMF may often be framed in a probabilistic framework, where the cost function appears as the negative log-likelihood of the data [4]–[7], and where the model structures the dispersion parameter of the underlying probability distribution rather than its observed realizations. For instance, in additive Gaussian mixtures [8] where the NMF models the variance of the sources, maximum likelihood estimation is equivalent to an NMF with IS divergence (ISNMF) of the power spectrogram [4].

Once the NMF model has been estimated, the complex-valued STFTs are retrieved by means of a Wiener-like filter [9]. This soft-masking of the complex-valued mixture's STFT assigns the phase of the original mixture to each extracted source. However, even if this filter yields quite satisfactory sounding estimates in practice [3], [4], it has been pointed out [10] that when sources overlap in the TF domain, it is responsible for residual interference and artifacts in the separated signals. This is a consequence of assuming that the phase is uniformly distributed [11], and therefore of not exploiting its underlying structure.

To alleviate this issue, the complex NMF (CNMF) model [12] has been proposed. It consists in directly decomposing the complex-valued mixture's STFT into a sum of rank-1 components whose magnitudes are structured by means of an NMF. This model allows for jointly estimating the magnitude and the phase of each source. It is estimated by minimizing the Euclidean distance between the model and the data, to which can be added some regularization terms, such as a sparsity penalty [12]. It was later improved by means of adding a *consistency* constraint [13], that is, to account for the redundancy of the STFT which introduces some dependencies between adjacent TF bins [14], [15].

Alternatively, improved recovery can be achieved by using phase constraints that originate from a signal model. For instance, the model of sums of sinusoids [16] leads to explicit constraints between the phases of adjacent TF bins [17], [18]. Such an approach has been exploited in speech enhancement [19], [20], audio restoration [21] and for a time-stretching application in the phase vocoder algorithm [22]. It has also been incorporated into some phase-constrained CNMF models for audio source separation [23]–[25]. Those developments have shown promising results in terms of interference rejection, though they suffer from two drawbacks. Firstly, the

CNMF model is estimated by minimizing a Euclidean distance, which does not properly characterize the properties of audio (such as its large dynamic range), where alternative divergences (such as KL or IS) are preferred [26]. Secondly, the set of estimated sources does not preserve the overall energy of the mixture, which leads to artifacts in the separated signals.

Drawing on those observations, we proposed in a preliminary work [27] to model the sources with anisotropic Gaussian (AG) variables, i.e., where the phase is no longer uniform. In such a model, one can promote a phase value which is obtained by exploiting the sinusoidal model. Estimation in a minimum mean square error sense results in an anisotropic Wiener filter, which optimally combines the mixture phase and the underlying phase model. We further introduced in [28] a general Bayesian framework in which both magnitudes and phases were modeled as random variables, and the sinusoidal model was promoted through a Markov chain prior structure on the phase location parameter. However, in those preliminary approaches, the variance parameters were left unconstrained and therefore either assumed known or estimated beforehand.

In this paper, we introduce a Bayesian AG model that overcomes the limitations of those approaches. We structure the variance parameters of the sources by means of an NMF model, so we can jointly estimate the magnitudes and the phases in a unified framework. This model, called *complex ISNMF*, combines the benefits of both ISNMF and CNMF:

1) It is phase-aware;
2) The set of estimators is *conservative*, i.e., their sum is equal to the observed mixture;
3) The estimation is based on the minimization of an IS-like divergence, which is appropriate for audio [29].

In order to infer the parameters of the model, we derive a generalized expectation-maximization (EM) algorithm. This model is applied to a musical source separation task in a semi-informed setting. It outperforms both the traditional phase-unaware ISNMF and the phase-constrained CNMF model [25]. This demonstrates the usefulness of such a phase-aware Bayesian AG model to perform the joint estimation of magnitudes and phases for audio source separation.

The rest of this paper is organized as follows. Section II introduces the complex ISNMF model. Section III details the inference procedure. Section IV experimentally validates the potential of this method. Finally, Section V draws some concluding remarks.

## II. COMPLEX ISNMF

Let $\mathbf{X} \in \mathbb{C}^{F \times T}$ be the STFT of a single-channel audio signal, where $F$ and $T$ are the numbers of frequency channels and time frames. $\mathbf{X}$ is the linear and instantaneous mixture of $J$ sources $\mathbf{S}_j \in \mathbb{C}^{F \times T}$, such that for all TF bins $ft$,

$$x_{ft} = \sum_{j=1}^{J} s_{j,ft}. \tag{1}$$

Since all TF bins are treated similarly, we remove the indices $ft$ when appropriate for more clarity.



Fig. 1. Density of the VM distribution.

### A. Modeling magnitude and phase

Let us consider a complex-valued random variable $s = re^{i\phi}$ whose magnitude and phase are assumed independent and denoted $r$ and $\phi$. Drawing on [28], we propose to model $r$ as a Rayleigh random variable $\mathcal{R}(v)$, which is the distribution of the modulus of a circularly-symmetric complex normal distribution with variance $v$. Besides, as in [27], we consider that the phase should be distributed around some favored value $\mu$ and that the relative importance of this value should be adjusted by means of a concentration parameter $\kappa \in [0, +\infty[$: the higher $\kappa$, the more favored $\mu$.

Several non-uniform periodic distributions exist (such as the wrapped Gaussian [30] or wrapped Cauchy distributions) but the von Mises (VM) [31] distribution comes as a natural candidate [32], [33], since its density is easily expressed by:

$$p(\phi|\mu, \kappa) = \frac{e^{\kappa \cos(\phi - \mu)}}{2\pi I_0(\kappa)}, \tag{2}$$

where $I_n$ is the modified Bessel function of the first kind of order $n$ [34], $\mu \in [0; 2\pi[$ is a location parameter and $\kappa \in [0; +\infty[$ is a concentration parameter. In particular, if $\kappa = 0$, the VM distribution becomes uniform. Contrarily, if $\kappa \to +\infty$, it becomes equivalent to a Dirac delta function centered at $\mu$. It is illustrated in Fig. 1.

This methodology results in a model called Rayleigh + von Mises (RVM), in which one can promote some favored phase values (see Section II-C). Such an approach has been originally used in [32], [33] for a speech enhancement application in a speech plus noise model. However, in the present case, since we consider any number of sources $J$, the RVM model is no longer tractable because the density of the mixture does not admit a closed-form expression. Therefore it is not suitable for source separation, where we aim to estimate the model parameters.

Nonetheless, we can compute the moments of $s = re^{i\phi}$ which will be used later in this work. If $\phi \sim \mathcal{VM}(\mu, \kappa)$, the $n$-th circular moment is, $\forall n \in \mathbb{Z}$ (*cf.* [31]):

$$\mathbb{E}(e^{in\phi}) = \frac{I_{|n|}(\kappa)}{I_0(\kappa)} e^{in\mu}. \tag{3}$$

Besides, if magnitude $r \sim \mathcal{R}(v)$, we have:

$$\mathbb{E}(r) = \sqrt{\frac{\pi}{4}v} \text{ and } \mathbb{E}(r^2) = v. \tag{4}$$

This lead to the expression of the mean of $s$:

$$m = \mathbb{E}(re^{i\phi}) = \mathbb{E}(r)\mathbb{E}(e^{i\phi}) = \lambda\sqrt{v}e^{i\mu}, \qquad (5)$$

and its variance $\gamma = \mathbb{E}(|s-m|^2)$:

$$\gamma = \mathbb{E}(|re^{i\phi}|^2) - |m|^2 = (1 - \lambda^2)v, \qquad (6)$$

and the relation term $c = \mathbb{E}((s-m)^2)$, which measures the joint variability of a variable and its complex conjugate:

$$c = \mathbb{E}(r^2)\mathbb{E}(e^{i2\phi}) - m^2 = \rho v e^{i2\mu}, \qquad (7)$$

where

$$\lambda = \frac{\sqrt{\pi}}{2}\frac{I_1(\kappa)}{I_0(\kappa)} \text{ and } \rho = \frac{I_2(\kappa)}{I_0(\kappa)} - \lambda^2. \qquad (8)$$

This relation term $c$ is not commonly introduced in statistical models of audio signals in the TF domain because it is usually assumed to be null [35]. Indeed, most models [4], [9], [36] assume the second-order circularity (or *isotropy*) of the variables, that is, with the same distribution in the complex plane regardless of the orientation. Since this is equivalent to assuming that the phase is uniformly distributed, we propose instead to explicitly consider this relation term as non-zero in general: it enables us to promote the non-circularity of the variable, and therefore the non-uniformity of the phase.

### B. Anisotropic Gaussian sources

To alleviate the non-tractability issue of the RVM model, we propose to approximate it by a Gaussian model[1] in which the moments of the variables are the same ones as in the original RVM model. This approach enables us to keep the phase dependencies in a model which is fully tractable.

Therefore, we assume that each source $s_j$ follows a complex normal distribution: $s_j \sim \mathcal{N}(m_j, \Gamma_j)$, where $m_j = \mathbb{E}(s_j) \in \mathbb{C}$ is the mean of $s_j$ and $\Gamma_j$ is its covariance matrix:

$$\Gamma_j = \begin{pmatrix} \gamma_j & c_j \\ \bar{c}_j & \gamma_j \end{pmatrix}, \qquad (9)$$

where $\gamma_j = \mathbb{E}(|s_j - m_j|^2) \in \mathbb{R}_+$ and $c_j = \mathbb{E}((s_j - m_j)^2) \in \mathbb{C}$ are the variance and relation term of $s_j$, and $\bar{z}$ denotes the complex conjugate of $z$. The density of such a distribution is:

$$p(x|m, \Gamma) = \frac{1}{\pi\sqrt{|\Gamma|}}e^{-\frac{1}{2}(\underline{x}-\underline{m})^{\mathsf{H}}\Gamma^{-1}(\underline{x}-\underline{m})}, \qquad (10)$$

where $\underline{x} = \begin{pmatrix} x & \bar{x} \end{pmatrix}^{\mathsf{T}}$, and where $^{\mathsf{T}}$ and $^{\mathsf{H}}$ denote the transpose and conjugate transpose.

Many previous studies model the sources as circularly-symmetric (or *isotropic*) variables [4], [38] (i.e., such that $m_j = c_j = 0$), which is equivalent to assuming that the phase of each source is uniformly distributed. The keystone of our approach is that, in order to promote a favored phase value, the moments are the same ones as in the original RVM model. Therefore, we use the expressions given by (5), (6) and (7) to estimate the moments which are then used to design the Gaussian model, as illustrated in Fig. 2. The main characteristic of this model is that the relation terms $c_j$ are

[1]This strategy is reminiscent of [37], where the mixture model was a sum of random variables with phase priors.



Fig. 2. Design of the AG model. We first model the magnitudes and phases as Rayleigh and von Mises random variables. The moments in this model are then used to define the equivalent AG model.



Fig. 3. 2-D histograms of 10000 samples generated from the RVM model (left) and AG model (right), with $v = 1$, $\mu = \pi/3$ and $\kappa = 50$. The intersection between the dashed lines represents the mean of the samples.

non-zero in general, which conveys the property of *anisotropy* of the corresponding Gaussian distribution: this is why we refer to it as the anisotropic Gaussian (AG) model.

The additive property of the Gaussian distribution family then implies that $x \sim \mathcal{N}(m_x, \Gamma_x)$ with:

$$m_x = \sum_j m_j, \ \gamma_x = \sum_j \gamma_j, \ c_x = \sum_j c_j, \ \Gamma_x = \sum_j \Gamma_j. \quad (11)$$

**Remark**: If $\kappa = 0$, then $\lambda = \rho = 0$ and consequently $m = c = 0$ and $\gamma = v$: the RVM and AG models are then equivalent since they both become isotropic Gaussian. Contrarily, for important values of $\kappa$, the models still remain quite alike, as illustrated in Fig. 3 for $\kappa = 50$.

### C. Phase model

The non-uniformity of the phase is taken into account in the AG model through the location parameter $\mu$. However, in order to obtain good quality phase estimates, this model can benefit from incorporating some prior knowledge about the phase, for instance by accounting for its structure in time or frequency. We propose to exploit some information about the phase by exploiting the sinusoidal model, which is widely used for representing audio signals [19], [23]. Each source in the time domain is modeled as a sum of sinusoids. Let us assume that there is at most one sinusoid (whose normalized frequency is denoted $\nu_{j,ft}$) per frequency channel. It can be shown [21] that the phase $\mu_j$ follows the unwrapping equation:

$$\mu_{j,ft} \approx \mu_{j,ft-1} + 2\pi l \nu_{j,ft}, \qquad (12)$$

where $l$ is the hop size of the STFT. As in [28], we propose to enforce this property by means of a Markov chain prior structure. We have, for each source:

$$p(\mu_j) = \prod_{f=0}^{F-1} p(\mu_{j,f0}) \prod_{t=1}^{T-1} p(\mu_{j,ft}|\mu_{j,ft-1}). \qquad (13)$$

We then propose the following choice, for $t > 0$:

$$\mu_{j,ft}|\mu_{j,ft-1} \sim \mathcal{VM}(\mu_{j,ft-1} + 2\pi l\nu_{j,ft}, \tau), \qquad (14)$$

and the initial distribution in each frequency channel $p(\mu_{j,f0})$ is Jeffrey's non-informative prior. In this way, we enforce the phase location parameter to approximately follow the sinusoidal model (12). The parameter $\tau \in \mathbb{R}_+$ adjusts the relative importance of this prior. Once again, we choose a VM distribution for modeling the phase location parameter, since it is a natural candidate for accounting for the periodicity of this variable. However, unlike previously, we do not need here to approximate this distribution: since the prior (14) applies independently to each source, it is straightforward to explicitly obtain the log-prior:

$$\log(p(\boldsymbol{\mu})) \stackrel{c}{=} \tau \sum_{j,f,t} \Re\left(e^{i\mu_{j,ft}} e^{-i\mu_{j,ft-1} - 2i\pi l\nu_{j,ft}}\right), \qquad (15)$$

where $\stackrel{c}{=}$ denotes equality up to an additive constant and $\Re$ is the real part. The model therefore depends on two concentration parameters that have a different role: $\kappa$ quantifies the non-uniformity of the phase in the AG model (i.e., how concentrated about a location parameter the phase is), while $\tau$ quantifies how close to the sinusoidal model this location parameter is.

### D. Complex ISNMF

For practical separation applications, it is necessary to constrain the variance parameters of the sources $\mathbf{V}_j$. We propose to structure it by means of an NMF model:

$$\mathbf{V}_j = \mathbf{W}_j\mathbf{H}_j, \qquad (16)$$

where $\mathbf{W}_j$ and $\mathbf{H}_j$ are nonnegative-valued matrices of dimensions $F \times K_j$ and $K_j \times T$ respectively. Therefore, the moments in the AG model become:

$$
\begin{aligned}
m_{j,ft} &= \lambda\sqrt{[\mathbf{W}_j\mathbf{H}_j]_{ft}}\ e^{i\mu_{j,ft}}, \\
\gamma_{j,ft} &= (1-\lambda^2)[\mathbf{W}_j\mathbf{H}_j]_{ft}, \\
c_{j,ft} &= \rho[\mathbf{W}_j\mathbf{H}_j]_{ft}\ e^{i2\mu_{j,ft}},
\end{aligned}
\qquad (17)
$$

where $[\mathbf{W}_j\mathbf{H}_j]_{ft}$ denotes the $(f,t)$-th entry of the matrix $\mathbf{W}_j\mathbf{H}_j$. In particular, if $\kappa = 0$, then $m_j = c_j = 0$ and $\gamma_j = \mathbf{W}_j\mathbf{H}_j$: the model becomes equivalent to ISNMF. Thus, since the proposed model generalizes ISNMF while allowing us to account for some phase constraint, we call it *complex ISNMF*. The whole model is represented as a Bayesian network in Fig. 4



Fig. 4. Bayesian network corresponding to the complex ISNMF model. Latent (resp. observed) variables are represented with empty (resp. shaded) ellipses. The sub-graph contained in each rectangle is repeated according to the index ($k$ or $j$) indicated in the bottom-right corner of the rectangle. The vertical dashed lines mark the limits between successive time frames.

### E. Relation to other models

The AG model along with the NMF variance structure results in a phase-aware extension of ISNMF, as pointed out in Section II-D. However, other models can be seen as particular cases of this general framework. Indeed, in Section II-B we approximated the RVM model with an AG model by equating their moments. As illustrated in Fig. 2, we chose to equate all the moments (mean, variance and relation term), but other approaches are possible.

Firstly, it is possible to set the mean and relation term to 0, in which case the sources follow a circularly-symmetric Gaussian distribution: $s_j \sim \mathcal{N}(0, \gamma_j I)$, where $I$ is the identity matrix. Along with an NMF variance, this results in the ISNMF model [4]. This is therefore another way of seeing the proposed AG model as an extension of ISNMF.

Alternatively, one can only preserve the mean information from the RVM model, and set the covariance matrix to be diagonal with a constant variance $\sigma$: $s_j \sim \mathcal{N}(m_j, \sigma I)$. This is the underlying statistical model from CNMF [12]. Therefore, this AG framework bridges the gap between ISNMF and CNMF since it generalizes both of them in a unified model.

Finally, other approximations are possible. For instance, one can only preserve the second-order statistics from the RVM model and set the mean value at 0 ($s_j \sim \mathcal{N}(0, \Gamma_j)$). Instead, one can set the relation terms at 0 and keep the phase dependencies only through the mean ($s_j \sim \mathcal{N}(m_j, \gamma_j I)$). This leads to alternative versions of Complex ISNMF that simplify the estimation of the NMF parameters (*cf.* Section III-C) or the phase parameters (*cf.* Section III-D). Those will be discussed in the corresponding sections. However, in order to keep the scope of this paper broad enough, we will infer the model in the general case described in Section II-D.

## III. INFERENCE

The model parameters $\Theta = \{\{\mathbf{W}_j\}_j, \{\mathbf{H}_j\}_j, \{\mu_j\}_j\}$ are estimated in a maximum a posteriori sense, which consists in

maximizing the log-posterior distribution:

$$\mathcal{C}_{\text{MAP}}(\Theta) = \log p(\mathbf{X}|\Theta) + \log p(\Theta), \tag{18}$$

where $p(\mathbf{X}|\Theta)$ is the likelihood of the data and $p(\Theta)$ the priors on the parameters. In this work, we only exploit the Markov prior information about the phase, therefore $\log p(\Theta)$ is given by (15). However, this framework is very general and it could be possible to further enforce some desirable property such as harmonicity [39] through priors on the columns of $\mathbf{W}_j$ or temporal continuity [3] through priors on the rows of $\mathbf{H}_j$.

### A. EM framework

Since the direct maximization of the criterion (18) is more involved than in classical isotropic models [4], we propose to adopt an EM [40] strategy which consists in maximizing a lower bound of the log-posterior distribution, given by:

$$\mathcal{Q}^{\text{MAP}}(\Theta, \Theta^{(i-1)}) = \mathcal{Q}^{\text{ML}}(\Theta, \Theta^{(i-1)}) + \log p(\Theta), \tag{19}$$

where $i$ is a step index, $\Theta^{(i-1)}$ contains the current set of estimated parameters (i.e., the parameters estimated at the previous step $i-1$) and $\mathcal{Q}^{\text{ML}}$ is the conditional expectation of the complete-data log-likelihood:

$$\mathcal{Q}^{\text{ML}}(\Theta, \Theta^{(i-1)}) = \int p(\mathbf{Z}|\mathbf{X}; \Theta^{(i-1)}) \log p(\mathbf{X}, \mathbf{Z}; \Theta) d\mathbf{Z}, \tag{20}$$

where $\mathbf{Z}$ denotes a set of latent (hidden) variables. Due to the mixing constraint (1), we use, as in [38], [41], a reduced set of $J' = J - 1$ free variables $\mathbf{Z} = \mathbf{S} = \{\mathbf{s}_{ft}\}_{ft}$, where we note $\mathbf{s}_{ft} = [s_{1,ft}, ..., s_{J',ft}]^{\mathsf{T}}$. Therefore, $s_{J,ft} = x_{ft} - \sum_{j=1}^{J'} s_{j,ft}$.

The EM algorithm consists in alternatively computing the functional $\mathcal{Q}^{\text{MAP}}$ given the current set of parameters $\Theta^{(i-1)}$ (E-step) and maximizing it with respect to $\Theta$ (M-step). This is proven [40] to increase the value of the criterion (18). However, when the maximization of $\mathcal{Q}^{\text{MAP}}$ is too involved, it may be preferable to solely increase its value at the M-step. This has also been proved [40] to lead to a local maximum of (18), and the corresponding procedure is called *generalized* EM. This is the approach we are adopting hereafter.

### B. E-step

Since all $\{s_{j,ft}\}_{j=1}^{J'}$ are independent Gaussian variables, $\mathbf{s}_{ft}$ is a Gaussian vector. It can be shown [35] that $\mathbf{S}|\mathbf{X}$ follows a multivariate complex normal distribution $\mathcal{N}(\mathbf{m}'_{ft}, \boldsymbol{\Xi}_{ft})$. The posterior means of the sources are given by anisotropic Wiener filtering [27]:

$$\underline{m}'_{j,ft} = \underline{m}_{j,ft}^{(i-1)} + \Gamma_{j,ft}^{(i-1)} \left(\Gamma_{x,ft}^{(i-1)}\right)^{-1} (\underline{x}_{ft} - \underline{m}_{x,ft}^{(i-1)}). \tag{21}$$

Note that, given the mixing constraint (1), this expression is also valid for the last source for which $j = J$. The posterior

covariance matrix $\boldsymbol{\Xi}_{ft}$ is given by [41]:

$$\boldsymbol{\Xi}_{ft} = \begin{pmatrix} \Gamma_{1,ft}^{(i-1)} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \Gamma_{J',ft}^{(i-1)} \end{pmatrix}$$
$$- \begin{pmatrix} \Gamma_{1,ft}^{(i-1)} \\ \vdots \\ \Gamma_{J',ft}^{(i-1)} \end{pmatrix} \left(\Gamma_{x,ft}^{(i-1)}\right)^{-1} \begin{pmatrix} \Gamma_{1,ft}^{(i-1)} \\ \vdots \\ \Gamma_{J',ft}^{(i-1)} \end{pmatrix}^{\mathsf{T}}. \tag{22}$$

In particular, the diagonal blocks in the posterior covariance matrix provide the posterior covariance for each source:

$$\Gamma'_{j,ft} = \Gamma_{j,ft}^{(i-1)} - \Gamma_{j,ft}^{(i-1)} \left(\Gamma_{x,ft}^{(i-1)}\right)^{-1} \Gamma_{j,ft}^{(i-1)}. \tag{23}$$

Thanks to (21) and (23), we can compute the posterior mean, variance and relation term of the sources, respectively, denoted by $m'_j$, $\gamma'_j$ and $c'_j$. The computation of (20) is detailed in the appendix and results in:

$$\mathcal{Q}^{\text{ML}}(\Theta, \Theta^{(i-1)}) \stackrel{c}{=} -\sum_{f,t} \sum_{j=1}^{J} \log(\sqrt{|\Gamma_{j,ft}|})$$
$$+ \frac{1}{|\Gamma_{j,ft}|} \left(\gamma_{j,ft}(|m'_{j,ft} - m_{j,ft}|^2 + \gamma'_{j,ft})\right) \tag{24}$$
$$- \frac{1}{|\Gamma_{j,ft}|} \left(\Re(\bar{c}_{j,ft}((m'_{j,ft} - m_{j,ft})^2 + c'_{j,ft}))\right),$$

where $|\Gamma_{j,ft}| = \gamma_{j,f,t}^2 - |c_{j,ft}|^2$ is the determinant of $\Gamma_{j,ft}$.

### C. M-step: NMF parameters

*1) NMF functional:* Let us first rewrite $\mathcal{Q}^{\text{ML}}$ by removing the terms that do not depend on the NMF parameters. Using (24) and (17), we have:

$$\mathcal{Q}^{\text{ML}}(\Theta|\Theta^{(i-1)}) \stackrel{c}{=} -\sum_{j=1}^{J} \sum_{f,t} \log([\mathbf{W}_j \mathbf{H}_j]_{ft}) + \frac{p_{j,ft}}{[\mathbf{W}_j \mathbf{H}_j]_{ft}}$$
$$- \frac{q_{j,ft}}{\sqrt{[\mathbf{W}_j \mathbf{H}_j]_{ft}}}, \tag{25}$$

with:

$$p = \frac{(1-\lambda^2)\left(\gamma' + |m'|^2\right) - \rho\Re\left(e^{-2\mathrm{i}\mu}(c' + m'^2)\right)}{(1-\lambda^2)^2 - \rho^2}, \tag{26}$$

and:

$$q = \frac{2\lambda}{1 - \lambda^2 + \rho} \Re\left(e^{-\mathrm{i}\mu} m'\right), \tag{27}$$

where we removed the indices $j, ft$ for brevity. This highlights two novel quantities $p$ and $q$ on which $\mathcal{Q}^{\text{ML}}$ depends. First, from the derivation conducted in the appendix we remark that:

$$\frac{p_{j,ft}}{[\mathbf{W}_j \mathbf{H}_j]_{ft}} = \mathbb{E}_{\mathbf{S}|\mathbf{X};\Theta^{(i-1)}} \left(\underline{s}_{j,ft}^{\mathsf{H}} \Gamma_{j,ft}^{-1} \underline{s}_{j,ft}\right). \tag{28}$$

In particular, when $\kappa = 0$, $p_{j,ft} = \gamma'_{j,ft} + |m'_{j,ft}|^2$, which is the posterior power of $s_{j,ft}$. Therefore, in the general case, we call the quantity $p$ in (28) the *phase-corrected posterior power* of the sources. Note that since $\Gamma$ is positive-definite, $p$ is necessarily nonnegative. This quantity is interesting because it accounts for the phase while being nonnegative: therefore,

estimating the NMF model from this quantity leads to a phase-aware decomposition of the data.

On the other hand, the physical meaning of the quantity $q$ is not fully clear. In particular, it has the same sign as $\Re\left(e^{-i\mu}m'\right)$, that is, the same sign as $\cos(\mu - \angle m')$. Accounting for the mixture's phase when computing the posterior mean (21) leads to a deviation of $\angle m'$ from the location parameter $\mu$. However, our intuition is that the posterior mean angle will stay relatively close to the location parameter $\mu$. If this angle difference remains relatively small (that is, $|\mu - \angle m'| < \pi/2$), then its cosine (and consequently $q$) is nonnegative. Then, $q$ has the dimension of a magnitude, and can therefore be seen as a *phase-corrected posterior magnitude*. Even though we were not able to formally demonstrate that this intuition holds, we observed experimentally that $q$ was always nonnegative. Therefore, we will assume in what follows that $q$ is nonnegative, and we leave to future work a more in-depth analysis of those quantities.

*2) Majorize-minimization approach:* Since $\mathcal{Q}^{\text{MAP}}$ is equal to $\mathcal{Q}^{\text{ML}}$ up to the log-prior on the phase, which does not depend on the NMF parameters, the problem then becomes that of minimizing the following function, for all sources $j$:

$$\mathcal{H}(\Theta) = \sum_{f,t}\log(\sum_k w_{fk}h_{kt}) + \frac{p_{ft}}{\sum_k w_{fk}h_{kt}} - \frac{q_{ft}}{\sqrt{\sum_k w_{fk}h_{kt}}}. \quad (29)$$

To do so, we propose to adopt a majorize-minimization approach [42]. The core idea of this strategy is to find an auxiliary function $\mathcal{G}$ which majorizes $\mathcal{H}$:

$$\forall(\Theta, \widetilde{\Theta}), \ \mathcal{H}(\Theta) \leq \mathcal{G}(\Theta, \widetilde{\Theta}), \text{ and } \mathcal{H}(\widetilde{\Theta}) = \mathcal{G}(\widetilde{\Theta}, \widetilde{\Theta}). \quad (30)$$

Given some current parameter $\widetilde{\Theta}$, minimizing $\mathcal{G}(\Theta, \widetilde{\Theta})$ with respect to $\Theta$ provides an update on $\Theta$. This approach guarantees that the cost function $\mathcal{H}$ is non-increasing over iterations.

Let us derive the update on $\mathbf{W}_j$. We introduce auxiliary parameters $\widetilde{w}_{fk}$ and we denote $\widetilde{v}_{ft} = \sum_k \widetilde{w}_{fk}h_{kt}$. In a similar fashion as in [43]–[45], we decompose the function $\mathcal{H}$ into its convex and concave parts.

Since $p$ is nonnegative, the term in (29) involving $p$ is convex. Therefore it is majorized by using the Jensen inequality:

$$\frac{p_{ft}}{\sum_k w_{fk}h_{kt}} \leq \sum_k \frac{\widetilde{w}_{fk}^2}{w_{fk}}\frac{p_{ft}h_{kt}}{\widetilde{v}_{ft}^2}. \quad (31)$$

Besides, since we assumed that $q$ is negative, the term in (29) involving $q$ is concave, so it is majorized by its tangent:

$$-\frac{q_{ft}}{\sqrt{\sum_k w_{fk}h_{kt}}} \leq \sum_k \frac{w_{fk}h_{kt}q_{ft}}{\widetilde{v}_{ft}^{3/2}}. \quad (32)$$

Finally, the first term in (29) is majorized as in [44]:

$$\log(\sum_k w_{fk}h_{kt}) \leq \sum_k \frac{w_{fk}h_{kt}}{\widetilde{v}_{ft}}. \quad (33)$$

Combining (31), (32) and (33) results into the following auxiliary function for $\mathcal{H}$:

$$\mathcal{G}(\Theta, \widetilde{\Theta}) = \sum_{f,k}\frac{\widetilde{w}_{fk}^2}{w_{fk}}\sum_t\frac{p_{ft}h_{kt}}{\widetilde{v}_{ft}^2} + w_{fk}\sum_t h_{kt}(\frac{1}{\widetilde{v}_{ft}} + \frac{q_{ft}}{\widetilde{v}_{ft}^{3/2}}). \quad (34)$$

*3) Update rules:* Setting the derivative of $\mathcal{G}$ with respect to $w_{fk}$ at zero and solving leads to the following update:

$$w_{fk} = \widetilde{w}_{fk}\sqrt{\frac{\displaystyle\sum_t \frac{p_{ft}h_{kt}}{\widetilde{v}_{ft}^2}}{\displaystyle\sum_t h_{kt}\left(\frac{1}{\widetilde{v}_{ft}} + \frac{q_{ft}}{\widetilde{v}_{ft}^{3/2}}\right)}}. \quad (35)$$

We can rewrite this update rule onto matrix form as:

$$\mathbf{W}_j \leftarrow \mathbf{W}_j \odot \left(\frac{(\mathbf{P}_j \odot \mathbf{V}_j^{\odot-2})\mathbf{H}_j^{\mathsf{T}}}{(\mathbf{V}_j^{\odot-1} + \mathbf{Q}_j \odot \mathbf{V}_j^{\odot-3/2})\mathbf{H}_j^{\mathsf{T}}}\right)^{\odot 1/2}, \quad (36)$$

where $\odot$, $^{\odot}$ and the fraction bar denote element-wise matrix multiplication, power and division respectively, and where $\mathbf{P}_j$ and $\mathbf{Q}_j$ are the matrices whose entries are the $p_{j,ft}$ and $q_{j,ft}$ defined in (26) and (27). By applying exactly the same methodology, we obtain the update on $\mathbf{H}$:

$$\mathbf{H}_j \leftarrow \mathbf{H}_j \odot \left(\frac{\mathbf{W}_j^{\mathsf{T}}(\mathbf{P}_j \odot \mathbf{V}_j^{\odot-2})}{\mathbf{W}_j^{\mathsf{T}}(\mathbf{V}_j^{\odot-1} + \mathbf{Q}_j \odot \mathbf{V}_j^{\odot-3/2})}\right)^{\odot 1/2}. \quad (37)$$

*4) Relation to other approaches:* We remark that if $\kappa = 0$, then $\lambda = \rho = 0$: therefore, $q_{j,ft} = 0$ and $p_{j,ft}$ becomes the posterior power of $s_{j,ft}$, as mentioned in Section III-C1. Then, we recognize in (25) the IS divergence between $\mathbf{P}_j$ and $\mathbf{W}_j\mathbf{H}_j$, as in the EM algorithm for ISNMF [46]. Consequently, the updates rules (36) and (37) are similar to those obtained in such a scenario [46], up to an additional power $1/2$, which is common when applying the majorize-minimization methodology for estimating ISNMF [44].

Besides, one can consider an alternative AG model as described in Section II-E. If one considers that the sources are centered ($s_j \sim \mathcal{N}(0, \Gamma_j)$), then $\mathbf{Q}_j = 0$: we recognize in (25) the IS divergence between the NMF model and the phase-corrected posterior power. The derivation of the update rules is then easier than in the general case, since it eliminates the need for the majorize-minimization method: one can apply the commonly-used heuristic method described in [2] to obtain alternative multiplicative update rules. This approach is described in more details in [47].

*D. M-step: phase parameters*

Let us now derive the updates on the phase parameters. We rewrite the functional (24) by removing the terms that do not depend on the phase parameters, which leads to:

$$\mathcal{Q}^{\text{ML}}(\Theta|\Theta^{(i-1)}) \stackrel{c}{=} \sum_{j=1}^J \sum_{f,t} \Re\left(\alpha_{j,ft}e^{-2i\mu_{j,ft}} + \beta_{j,ft}e^{-i\mu_{j,ft}}\right), \quad (38)$$

with:

$$\alpha_{j,ft} = \frac{\rho}{((1-\lambda^2)^2 - \rho^2)[\mathbf{W}_j\mathbf{H}_j]_{ft}}(c'_{j,ft} + m'^2_{j,ft}), \quad (39)$$

and:

$$\beta_{j,ft} = \frac{2\lambda(1-\lambda^2-\rho)}{((1-\lambda^2)^2 - \rho^2)\sqrt{[\mathbf{W}_j\mathbf{H}_j]_{ft}}}m'_{j,ft}. \quad (40)$$

Therefore, adding the log-prior over the phase parameters (15) leads to maximizing the following functionals:

$$g_{j,ft}(\mu_{j,ft}) = \Re\left(\alpha_{j,ft}e^{-2i\mu_{j,ft}} + \tilde{\beta}_{j,ft}e^{-i\mu_{j,ft}}\right), \quad (41)$$

with respect to $\mu_{j,ft}$, and where:

$$\tilde{\beta}_{j,ft} = \beta_{j,ft} + \tau\left(e^{i\mu_{j,ft-1}+2i\pi l\nu_{j,ft}} + e^{i\mu_{j,ft+1}-2i\pi l\nu_{j,ft+1}}\right). \quad (42)$$

Let us remove the indexes $j, ft$ in what follows for more clarity. We then seek to maximize:

$$g(\mu) = \Re\left(\alpha e^{-2i\mu} + \tilde{\beta}e^{-i\mu}\right) \quad (43)$$

$$= |\alpha|\cos(2\mu - \angle\alpha) + |\tilde{\beta}|\cos(\mu - \angle\tilde{\beta}), \quad (44)$$

which leads to finding the roots of:

$$g'(\mu) = -2|\alpha|\sin(2\mu - \angle\alpha) - |\tilde{\beta}|\sin(\mu - \angle\tilde{\beta}). \quad (45)$$

Unfortunately, it is not straightforward to write the solutions of this problem in closed-form. Besides, it requires further operations to determine which root maximizes $g$, leading to a quite computationally intensive procedure. Instead, drawing on [28], since we experimentally observed that $|\alpha| << |\beta|$, we propose to approximate (44) by:

$$\tilde{g}(\mu) = \Re\left(\tilde{\beta}e^{-i\mu}\right) = |\tilde{\beta}|\cos(\mu - \angle\tilde{\beta}), \quad (46)$$

which is easily maximized by $\mu = \angle\tilde{\beta}$. This update depends on the values of the phase parameter in frames $t-1$ and $t+1$, so it has to be applied sequentially over time frames (which is common when using Markov chain priors such as in [39]).

To assess the validity of this update scheme, we applied both procedures (maximization of the exact functional (44) and its approximation (46)) on the learning dataset used in the experimental evaluation (see Section IV-A). The average relative difference between the phases obtained with those two approaches was of approximately $10^{-5}$. Consequently, we propose to use the approximate update scheme, since it yields very similar estimates while being significantly faster than performing the exact maximization.

Finally, if one consider an alternative AG model with null relation terms (*cf.* Section II-E), then $\alpha = 0$, which eliminates the need for this simplifying assumption. It also modifies the values of $\beta$, $p$ and $q$, therefore leading to a different procedure, which will be investigated in future work.

### E. Full procedure

The EM procedure is summarized in Algorithm 1. The phase location parameters $\mu_j$ are initialized by assigning the mixture phase to each source. The initialization of the NMF matrices is discussed in Sections IV-A2 and IV-B.

The frequencies $\nu$ are provided as inputs of the algorithm. We estimate them by means of a quadratic interpolated FFT (QIFFT) [48] on the log-spectra of the initial variance estimates $\mathbf{V}_j$. This estimation is performed locally (at each time frame) in order to account for slow variations of the frequencies. The frequency range is then decomposed into *regions of influence* [21] to ensure that the phase in a given channel is unwrapped with the appropriate frequency.

---

**Algorithm 1:** EM algorithm for complex ISNMF

**1 Inputs**: Mixture $\mathbf{X} \in \mathbb{C}^{F \times T}$,
**2** Phase parameters $\kappa$ and $\tau \in \mathbb{R}_+$,
**3** Initial NMF matrices $\forall j$, $\mathbf{W}_j \in \mathbb{R}_+^{F \times K_j}$, $\mathbf{H}_j \in \mathbb{R}_+^{K_j \times T}$,
**4** Initial phases $\forall j$, $\mu_j \in [0, 2\pi[^{F \times T}$,
**5** Normalized frequencies $\forall j$, $\nu_j \in \mathbb{R}^{\times F \times T}$.
**6 Anisotropy parameters**:
**7** Compute $\lambda$ and $\rho$ with (8).
**8 while** stopping criterion not reached **do**
**9**       % E-step
**10**       Update $m$, $\gamma$ and $c$ with (17),
**11**       Update $m_x$, $\gamma_x$ and $c_x$ with (11),
**12**       Update $m'$ with (21),
**13**       Update $\gamma'$ and $c'$ with (23),
**14**       % M-step: NMF
**15**       Update $p$ with (26) and $q$ with (27).
**16**       $\forall j$, Update $\mathbf{W}_j$ with (36) and $\mathbf{H}_j$ with (37),
**17**       Normalize $\mathbf{W}$ and $\mathbf{H}$.
**18**       % M-step: phase
**19**       Update $\beta$ with (40).
**20**       **for** $t = 1$ **to** $T - 2$ **do**
**21**          $\forall(j, f)$, update $\tilde{\beta}_{j,ft}$ with (42),
**22**          $\mu_{j,ft} = \angle\tilde{\beta}_{j,ft}$.
**23**       **end**
**24 end**
**25** Update $m$, $\gamma$ and $c$ with (17),
**26** Update $m_x$, $\gamma_x$ and $c_x$ with (11),
**27** Update $m'$ with (21).
**28 Outputs**: $m' \in \mathbb{C}^{J \times F \times T}$.

---

This algorithm includes a normalization step after updating $\mathbf{W}_j$ and $\mathbf{H}_j$, which eliminates trivial scale indeterminacies and avoids numerical instabilities. We impose a unitary $\ell_2$-norm on each column of $\mathbf{W}_j$ and scale $\mathbf{H}_j$ accordingly, so that the cost function is not affected.

Finally, one final E-step is performed after looping in order to estimate the sources with the most up-to-date parameters.

## IV. EXPERIMENTAL EVALUATION

In this section, we experimentally assess the potential of the proposed complex ISNMF model for a task of monaural musical source separation. Sound excerpts can be found on the companion website for this paper [49]. In the spirit of reproducible research, the code of this experimental study is available online[2].

### A. Protocol

*1) Dataset:* We consider 100 music song excerpts from the DSD100 database, a semi-professionally mixed set of music songs used for the SiSEC 2016 campaign [50]. Each excerpt is 10 seconds long and is made up of $J = 4$ sources: bass, drum, vocals and other. The database is split into two subsets of 50 songs: a learning set, on

---

[2]https://github.com/magronp/complex-isnmf

which the meta-parameters of the algorithms are tuned and the initialization strategies are investigated, and a test set, on which the separation benchmark is performed. The signals are sampled at $44100$ Hz and the STFT is computed with a $92$ ms long Hann window and $75$ % overlap. The resulting STFTs are therefore matrices of dimensions $2049 \times 433$.

*2) Separation scenario:* In *coding-based informed source separation* [51], we assume some side-information can be computed from the isolated sources (the *encoding* stage) and then used to perform separation (the *decoding* stage). A common approach consists of computing a nonnegative matrix or tensor factorization [52]–[54] on the isolated source spectrograms and then using the corresponding decomposition to estimate a Wiener filter at the decoding stage. Here, we consider a *semi-informed* scenario, in which the dictionaries $\mathbf{W}_j$ are estimated on the isolated sources and the activation matrices $\mathbf{H}_j$ computed from the mixture. This setting is less restrictive than a fully-informed setting since we only transmit the dictionaries instead of both NMF matrices. Note than since we use a learning dataset for tuning some parameters, this setting is actually supervised semi-informed, but we refer to it as semi-informed for brevity.

Dictionaries are learned with $200$ iterations of ISNMF applied to each isolated spectrogram, using multiplicative update rules [4], random initial matrices and a rank of factorization $K_j = 50$, which corresponds to an $8$-fold compression ratio. The dictionaries are then fixed at the separation stage, since we experimentally observed that it leads to better results than further updating them on the mixture.

*3) Comparison references:* As baselines, we test the consistent anisotropic Wiener (CAW) filter [41] which combines the consistent [38] and anisotropic [27] Wiener filters, and we also consider the phase-constrained CNMF [23]–[25]. In order to make the comparison fair, we implemented a version of CNMF known as CNMF with intra-source additivity [55]: it consists in modeling the phase $\phi_j$ of each source instead of the phase of each NMF component, as in the classical CNMF model [12]. This significantly reduces the number of parameters of the model, thus it lowers both the memory and computation time required for the estimation of the model, at the cost of a moderate drop in terms of separation quality [55].

Source separation quality is measured with the signal-to-distortion, signal-to-interference, and signal-to-artifact ratios (SDR, SIR, and SAR) [56] expressed in dB, where only a rescaling (not a refiltering) of the reference is allowed.

### B. Initialization strategy

We briefly investigate here on the best strategy for initializing the complex ISNMF algorithm at the separation stage, once the dictionaries are learned. A first approach is to provide a warm start to the algorithm thanks to $50$ iterations of ISNMF computed on the mixture, whose activation matrix is randomly initialized. Besides, it is necessary to have a first estimate of the variances in order to compute the frequencies, which are needed as inputs of Algorithm 1. On top of that initialization, we run $150$ iterations of complex ISNMF. Alternatively, we run $200$ iterations of complex ISNMF on top of a random



Fig. 5. SDR over iterations for an ISNMF (left) and random (right) initialization.



Fig. 6. Influence of the phase parameters $\kappa$ and $\tau$ on the source separation quality (SDR and SAR are similar). The range is limited to $[0, 1]$ and $[0, 5]$ for $\kappa$ and $\tau$ respectively for clarity purpose, since the performance decreases outside of these ranges.

initialization (though we still use the frequencies as computed before), so the total number of iterations is the same in both scenarios.

We present the SDR over iterations in Fig. 5 (results are averaged over the learning set) for $\kappa = \tau = 0.5$: similar conclusions can be drawn from other values of the parameters and from the SIR and SAR. We observe that initializing complex ISNMF with ISNMF provides better results than a random initialization. Consequently, in the following experiments, we will retain this ISNMF-initialization strategy in order to bootstrap the complex ISNMF algorithm, which will use $100$ iterations.

### C. Phase parameters influence

We run the different methods on the $50$ songs that form the learning set in order to learn the optimal phase parameters.

*1) Complex ISNMF:* The results presented in Fig 6 show that for non-null values of the phase parameters, the proposed approach can outperform a phase-unaware approach (for which $\kappa = \tau = 0$) according to the SDR, SIR and SAR. We found that $\kappa = 0.5$ and $\tau = 5$ provides a quite good compromise between the different indicators.

*2) Phase-constrained CNMF:* This method depends on a weight parameter $\sigma_u$ which promotes the sinusoidal model phase constraint. The separation work flow is the same as for complex ISNMF, except we use here an NMF with Euclidean distance [2] for both dictionary learning and initialization on the mixture. Indeed, since CNMF is based on the Euclidean distance, learning IS-based dictionaries would not be consistent with the distortion metric in CNMF. The value $\sigma_u = 10^{-2}$ appears as the best candidate, since the SDR is slightly reduced ($-0.2$ dB) compared to the unconstrained baseline (for which

| | Bass | | | Drums | | | Other | | | Vocals | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SDR | SIR | SAR | SDR | SIR | SAR | SDR | SIR | SAR | SDR | SIR | SAR |
| Wiener | 2.6 | 7.9 | 4.4 | 4.7 | 17.4 | 5.1 | 3.7 | 12.9 | 4.4 | 7.6 | 18.1 | 8.1 |
| AW | 2.6 | 8.1 | 4.3 | 4.4 | **18.5** | 4.7 | 3.6 | **13.1** | 4.2 | 7.5 | **18.9** | 7.9 |
| CAW | 2.8 | 8.1 | **4.5** | 4.8 | 17.6 | 5.1 | 3.8 | 12.9 | 4.4 | 7.0 | 16.7 | 7.5 |
| CNMF | 2.3 | 6.9 | 4.5 | 3.7 | 12.8 | 4.4 | 2.6 | 10.1 | 3.7 | 5.9 | 15.7 | 6.5 |
| Complex ISNMF | **3.0** | **10.1** | 4.1 | **5.4** | 15.9 | **5.9** | **3.8** | 12.4 | **4.6** | **7.7** | 18.4 | **8.2** |

| | SDR | SIR | SAR |
|---|---|---|---|
| Wiener | $4.7 \pm 1.6$ | $14.1 \pm 2.9$ | $5.5 \pm 1.5$ |
| AW | $4.5 \pm 1.7$ | **$14.6 \pm 2.8$** | $5.3 \pm 1.5$ |
| CAW | $4.6 \pm 2.0$ | $13.8 \pm 2.7$ | $5.4 \pm 2.0$ |
| CNMF | $3.6 \pm 1.7$ | $11.4 \pm 2.3$ | $4.8 \pm 1.6$ |
| Complex ISNMF | **$5.0 \pm 1.7$** | $14.2 \pm 2.8$ | **$5.7 \pm 1.6$** |

$\sigma_u = 0$), but it allows for more interference reduction ($+1.4$ dB in SIR). Values of $\sigma_u$ greater than $10^{-2}$ still increase the SIR, but at the cost of a significant drop in SDR.

*3) Wiener filters:* CAW [41] depends on two parameters $\kappa$ and $\delta$ which respectively promote anisotropy and consistency. We first estimate the variances with 150 iterations of ISNMF on the mixture, and then we apply the filter. We propose the following sets of values:

- For $\kappa = 1$ and $\delta = 0$, the SIR is improved by $+0.6$ dB at the cost of a slight decrease in SDR ($-0.1$ dB) compared to the baseline Wiener filtering (for which $\kappa = \delta = 0$). We simply refer to it as AW since the consistency weight is null in this setting.
- For $\kappa = 0.1$ and $\delta = 10^{-3}$, the SIR is very slightly reduced compared to the baseline ($-0.02$ dB) while the SDR is increased by $0.05$ dB. We refer to it as CAW.

One may chose other values for the parameters in order to have the best possible SDR (or SIR/SAR), but the proposed settings yield an overall compromise which does not excessively favor one indicator over the others.

### D. Results of the benchmark

We now consider the 50 songs that form the test set and run the compared methods. The results for each instrumental source are presented in Table I, and the results averaged over instruments are presented in Table II.

We observe that the proposed complex ISNMF approach yields the best results in terms of SDR and SAR for all instruments and among all the compared techniques, except for the bass track in terms of SAR. It also outperforms the phase-unaware Wiener filtering and the phase-constrained CNMF in terms of average SIR. This demonstrates the interest of exploiting some phase information in a probabilistic model to overcome the limitations of those baseline approaches, as stressed in the introduction of this paper.

The complex ISNMF estimates contain slightly more interference than the AW estimates (a $0.4$ dB difference in SIR on average), but less artifacts (a $0.4$ dB difference in SAR on average), which leads to a greater SDR. Therefore, it is overall preferable to employ this method than our preliminary approaches [27], [41] to perform a joint estimation of magnitude and phase.

Let us note that the metrics do not vary much from one technique to another. Indeed, the main difference between them is the phase recovery technique, which has less impact on the SDR, SIR and SAR than differences in terms of magnitude estimation strategy.

An informal perceptual evaluation is consistent with those results (sounds excerpts are available at [49]). In particular, CNMF introduces smearing artifacts in the separated sources, and the bass and drum tracks estimated with the Wiener filters are strongly corrupted by musical noise. In comparison, the proposed complex ISNMF method yields bass estimates which contain fewer artifacts and interference, and drums estimates with neater attacks.

### E. Fitting the data

Finally, we investigate on the capability of the AG model to represent audio data, that is to say, to assess that the mixture variables $x_{ft}$ are well-represented by AG distributions. To do so, we need to normalize the variables $x_{ft}$ so that all TF entries become identically distributed, which allows us to compute their histogram, and therefore to compare their empirical and theoretical densities. Since $x_{ft} \sim \mathcal{N}(m_{x,ft}, \Gamma_{x,ft})$, it can be shown that:

$$y_{ft} = (\underline{x}_{ft} - \underline{m}_{x,ft})^{\mathsf{H}} \Gamma_{x,ft}^{-1} (\underline{x}_{ft} - \underline{m}_{x,ft}) \qquad (47)$$

follows a chi-squared distribution with 2 degrees of freedom [35]. Then, once the model is estimated, we compute the normalized variable $\mathbf{Y}$ from the mixture $\mathbf{X}$ according to (47), and all the entries of $\mathbf{Y}$ are expected to be identically chi-squared distributed. Finally, even if there are some dependencies between the $x_{ft}$ because of the NMF and phase models, they are conditionally independent given the model parameters, which are estimated beforehand in order to compute the $y_{ft}$ with (47). The resulting variables $y_{ft}$ are then independent and identically distributed, thus it becomes possible to plot their histogram.

The setting is the same as in the previous experiments, but we set $\tau$ at 0 and we initialize Algorithm 1 with the true phase values for $\mu_j$. Indeed, a fitting error can be due to a mismatch between the model and the observed data, but also to an estimation error. In this way, we only investigate on the accuracy of the model to represent the data, not on the phase

Fig. 7. Empirical densities of the normalized data for several values of $\kappa$ (solid lines) and reference chi-squared density (dashed line).

estimation itself. The complex ISNMF algorithm is run on one song (similar results are obtained for the other songs) for several values of $\kappa$. The results are presented in Fig. 7.

We observe that small values of $\kappa$ lead to empirical densities that approach the theoretical one from above for small values of $x$ and from below for greater values of $x$. For greater values of $\kappa$, this trend is inverted. In particular, the value $\kappa = 0.5$ leads to a good fit on average, which may explain why this value leads to the best results in terms of separation quality (see Section IV-C).

Overall, a better fit can be obtained with non-null values of $\kappa$, which demonstrates the interest of AG distributions over isotropic variables to represent audio data in the STFT domain.

## V. CONCLUSION

In this paper, we introduced complex ISNMF, a probabilistic model based on the AG distribution. It consists of modeling the sources with anisotropic random variables, which makes it possible to enforce some desirable phase properties, while classical circularly-symmetric variables do not allow one to favor a phase model. Therefore, it combines the advantages of ISNMF and CNMF, that is, using a distortion metric well adapted to audio and phase-awareness. We experimentally showed that it outperforms those two approaches, and thus appears as a good candidate for phase-aware audio source separation in semi-informed settings. This model is also suitable for supervised applications where some training material is available, but then it is required to account for the potential mismatch between training and test materials [57], [58].

An interesting direction for future work is the investigation of alternative phase-aware probabilistic models, in order to extend CNMF to other beta-divergences, as first attempted in [59]. Alternatively, one can exploit the family of multivariate stable distributions [60] with an anisotropic shape matrix in order to combine phase-awareness and robust magnitude modeling [61]. Finally, we could incorporate deep neural networks in this Bayesian framework for estimating the variances

instead of using an NMF model, as it was done in a multichannel scenario with isotropic Gaussian variables [62]. Indeed, deep learning methods have shown remarkably good results for musical source separation [63], but there is still some room for improvement, notably in terms of phase recovery, since those methods usually exploit a phase-unaware Wiener-like mask to estimate the complex-valued sources.

## APPENDIX

In this appendix, we detail the E-step of the proposed algorithm, which consists in computing the functional given by (20), which we recall hereafter:

$$\mathcal{Q}^{\mathrm{ML}}(\Theta, \Theta^{(i-1)}) = \int p(\mathbf{S}|\mathbf{X}; \Theta^{(i-1)}) \log p(\mathbf{X}, \mathbf{S}; \Theta) d\mathbf{S}.$$

The complete data log-likelihood is given by:

$$\log p(\mathbf{X}, \mathbf{S}; \Theta) = \sum_{f,t} \log p(x_{ft}|\mathbf{s}_{ft}; \Theta) + \sum_{j=1}^{J'} \log p(s_{j,ft}; \Theta)$$

$$\overset{c}{=} -\frac{1}{2} \sum_{f,t} \log(|\Gamma_{J,ft}|) + B_{ft} + \sum_{j=1}^{J'} \log(|\Gamma_{j,ft}|) + A_{j,ft},$$

where:

$$A_{j,ft} = (\underline{s}_{j,ft} - \underline{m}_{j,ft})^{\mathsf{H}} \Gamma_{j,ft}^{-1} (\underline{s}_{j,ft} - \underline{m}_{j,ft}),$$

and

$$B_{ft} = (\underline{x}_{ft} - \underline{m}_{J,ft} - \sum_{j=1}^{J'} \underline{s}_{j,ft})^{\mathsf{H}} \Gamma_{J,ft}^{-1} (\underline{x}_{ft} - \underline{m}_{J,ft} - \sum_{j=1}^{J'} \underline{s}_{j,ft}).$$

Therefore, (20) rewrites:

$$\mathcal{Q}^{\mathrm{ML}}(\Theta, \Theta^{(i-1)}) \overset{c}{=} -\frac{1}{2} \sum_{f,t} \sum_{j=1}^{J} \log(|\Gamma_{j,ft}|)$$

$$+ \sum_{f,t} \sum_{j=1}^{J'} \mathbb{E}_{\mathbf{S}|\mathbf{X};\Theta^{(i-1)}} (A_{j,ft}) + \mathbb{E}_{\mathbf{S}|\mathbf{X};\Theta^{(i-1)}} (B_{ft}). \quad (48)$$

Firstly, let us compute the expectation $\mathbb{E}_{\mathbf{S}|\mathbf{X};\Theta^{(i-1)}} (A_{j,ft})$. We remove the indices $j, ft$ and the subscript $\mathbf{S}|\mathbf{X}; \Theta^{(i-1)}$ for clarity. We have, thanks to the trace identity:

$$\mathbb{E}(A) = \mathbb{E} \left( (\underline{s} - \underline{m})^{\mathsf{H}} \Gamma^{-1} (\underline{s} - \underline{m}) \right)$$
$$= (\underline{m}' - \underline{m})^{\mathsf{H}} \Gamma^{-1} (\underline{m}' - \underline{m}) + \mathrm{Tr}(\Gamma^{-1} \Gamma').$$

Besides,

$$\mathrm{Tr}(\Gamma^{-1} \Gamma') = \frac{1}{|\Gamma|} (\gamma \gamma' - \Re(\bar{c}c')),$$

then:

$$\mathbb{E}(A) = \frac{2}{|\Gamma|} \left( \gamma(|m' - m|^2 + \gamma') - \Re(\bar{c}((m' - m)^2 + c')) \right).$$

Now, let us compute $\mathbb{E}(B)$. We use, once again, the trace identity, which leads to:

$$\mathbb{E}(B) = \mathbb{E} \left( (\underline{x} - \underline{m}_J - \sum_{j=1}^{J'} \underline{s}_j)^{\mathsf{H}} \Gamma_J^{-1} (\underline{x} - \underline{m}_J - \sum_{j=1}^{J'} \underline{s}_j) \right)$$

$$= (\underline{x} - \underline{m}_J - \sum_{j=1}^{J'} \underline{m}'_j)^{\mathsf{H}} \Gamma^{-1} (\underline{x} - \underline{m}_J - \sum_{j=1}^{J'} \underline{m}'_j) + \mathrm{Tr}(\Gamma_J^{-1} \Gamma'_J).$$

Thanks to the conservative property of the anisotropic Wiener filtering (21), we have $\sum_{j=1}^{J} \underline{m}'_j = \underline{x} - \underline{m}'_J$, so:

$$\mathbb{E}(B) = (\underline{m}'_J - \underline{m}_J)^{\mathsf{H}} \Gamma_J^{-1} (\underline{m}'_J - \underline{m}_J) + \mathrm{Tr}(\Gamma_J^{-1} \Gamma'_J).$$

Then, $\mathbb{E}(B)$ is similar to $\mathbb{E}(A)$, but applied to the last source $J$. Finally, incorporating the expressions of $\mathbb{E}(A)$ and $\mathbb{E}(B)$ into (48) leads to the expression of $\mathcal{Q}^{\mathrm{ML}}$:

$$\mathcal{Q}^{\mathrm{ML}}(\Theta, \Theta^{(i-1)}) \stackrel{c}{=} -\sum_{f,t} \sum_{j=1}^{J} \log(\sqrt{|\Gamma_{j,ft}|})$$
$$+ \frac{1}{|\Gamma_{j,ft}|} \left( \gamma_{j,ft}(|m'_{j,ft} - m_{j,ft}|^2 + \gamma'_{j,ft}) \right)$$
$$- \frac{1}{|\Gamma_{j,ft}|} \left( \Re(\bar{c}_{j,ft}((m'_{j,ft} - m_{j,ft})^2 + c'_{j,ft})) \right).$$

## REFERENCES

[1] P. Comon and C. Jutten, *Handbook of blind source separation: independent component analysis and applications*. Academic press, 2010.

[2] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, no. 6755, pp. 788–791, 1999.

[3] T. Virtanen, "Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 3, pp. 1066–1074, March 2007.

[4] C. Févotte, N. Bertin, and J.-L. Durrieu, "Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis," *Neural computation*, vol. 21, no. 3, pp. 793–830, March 2009.

[5] T. Virtanen, A. T. Cemgil, and S. Godsill, "Bayesian extensions to non-negative matrix factorisation for audio signal modelling," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2008, pp. 1825–1828.

[6] A. Liutkus, D. Fitzgerald, and R. Badeau, "Cauchy nonnegative matrix factorization," in *Proc. of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, October 2015, pp. 1–5.

[7] U. Simsekli, A. Liutkus, and A. T. Cemgil, "Alpha-stable matrix factorization," *IEEE Signal Processing Letters*, vol. 22, no. 12, pp. 2289–2293, December 2015.

[8] C. Fevotte and J. F. Cardoso, "Maximum likelihood approach for blind audio source separation using time-frequency Gaussian source models," in *Proc. of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, October 2005, pp. 78–81.

[9] A. Liutkus and R. Badeau, "Generalized Wiener filtering with fractional power spectrograms," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2015, pp. 266–270.

[10] P. Magron, R. Badeau, and B. David, "Phase recovery in NMF for audio source separation: an insightful benchmark," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2015, pp. 81–85.

[11] R. M. Parry and I. Essa, "Incorporating phase information for source separation via spectrogram factorization," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2007, pp. II–661II–664.

[12] H. Kameoka, N. Ono, K. Kashino, and S. Sagayama, "Complex NMF: A new sparse representation for acoustic signals," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2009, p. 34373440.

[13] J. Le Roux, H. Kameoka, E. Vincent, N. Ono, K. Kashino, and S. Sagayama, "Complex NMF under spectrogram consistency constraints," in *Proc. of Acoustical Society of Japan Autumn Meeting*, September 2009.

[14] D. Griffin and J. S. Lim, "Signal estimation from modified short-time Fourier transform," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 32, no. 2, pp. 236–243, April 1984.

[15] J. Le Roux, N. Ono, and S. Sagayama, "Explicit consistency constraints for STFT spectrograms and their application to phase reconstruction," in *Proc. of ISCA Workshop on Statistical and Perceptual Audition (SAPA)*, September 2008, pp. 23–28.

[16] R. J. McAuley and T. F. Quatieri, "Speech analysis/Synthesis based on a sinusoidal representation," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 34, no. 4, pp. 744–754, August 1986.

[17] M. Krawczyk and T. Gerkmann, "STFT phase improvement for single channel speech enhancement," in *Proc. of International Workshop on Acoustic Signal Enhancement (IWAENC)*, September 2012, pp. 1–4.

[18] P. Magron, R. Badeau, and B. David, "Model-based STFT phase recovery for audio source separation," *IEEE/ACM Transactions on Audio, Speech and Language Processing*, vol. 26, no. 6, pp. 1095–1105, June 2018.

[19] M. Krawczyk and T. Gerkmann, "STFT phase reconstruction in voiced speech for an improved single-channel speech enhancement," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 12, pp. 1931–1940, December 2014.

[20] P. Mowlaee and J. Kulmer, "Harmonic phase estimation in single-channel speech enhancement using phase decomposition and SNR information," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 9, pp. 1521–1532, September 2015.

[21] P. Magron, R. Badeau, and B. David, "Phase reconstruction of spectrograms with linear unwrapping: application to audio signal restoration," in *Proc. of European Signal Processing Conference (EUSIPCO)*, August 2015, pp. 1–5.

[22] J. Laroche and M. Dolson, "Improved phase vocoder time-scale modification of audio," *IEEE Transactions on Speech and Audio Processing*, vol. 7, no. 3, pp. 323–332, May 1999.

[23] J. Bronson and P. Depalle, "Phase constrained complex NMF: Separating overlapping partials in mixtures of harmonic musical sources," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2014, pp. 7475–7479.

[24] F. J. Rodriguez-Serrano, S. Ewert, P. Vera-Candeas, and M. Sandler, "A score-informed shift invariant extension of complex matrix factorisation for improving the separation of overlapped partials in music recordings," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2016, pp. 61–65.

[25] P. Magron, R. Badeau, and B. David, "Complex NMF under phase constraints based on signal modeling: application to audio source separation," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2016, pp. 46–50.

[26] R. Gray, A. Buzo, A. Gray, and Y. Matsuyama, "Distortion measures for speech processing," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 28, no. 4, pp. 367–376, August 1980.

[27] P. Magron, R. Badeau, and B. David, "Phase-dependent anisotropic Gaussian model for audio source separation," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2017, pp. 513–535.

[28] P. Magron and T. Virtanen, "Bayesian anisotropic Gaussian model for audio source separation," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2018, pp. 166 – 170.

[29] B. King, C. Févotte, and P. Smaragdis, "Optimal cost function and magnitude power for NMF-based speech separation and music interpolation," in *Proc. of IEEE International Workshop on Machine Learning for Signal Processing (MLSP)*, September 2012, pp. 1–6.

[30] Y. Agiomyrgiannakis and Y. Stylianou, "Wrapped Gaussian mixture models for modeling and high-rate quantization of phase data of speech," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 4, pp. 775–786, May 2009.

[31] K. V. Mardia and P. J. Zemroch, "Algorithm AS 86: The von Mises distribution function," *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, vol. 24, no. 2, pp. 268–272, 1975.

[32] T. Gerkmann, "MMSE-optimal enhancement of complex speech coefficients with uncertain prior knowledge of the clean speech phase," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2014, pp. 4478–4482.

[33] ——, "Bayesian estimation of clean speech spectral coefficients given a priori knowledge of the phase," *IEEE Transactions on Signal Processing*, vol. 62, no. 16, pp. 4199–4208, August 2014.

[34] G. N. Watson, *A treatise on the theory of Bessel functions*. Cambridge university press, 1995.

[35] B. Picinbono, "Second-order complex random vectors and normal distributions," *IEEE Transactions on Signal Processing*, vol. 44, no. 10, pp. 2637–2640, October 1996.

[36] A. Liutkus, C. Rohlfing, and A. Deleforge, "Audio source separation with magnitude priors: the BEADS model," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2018, pp. 56 – 60.

[37] P. Beckmann, "Statistical distribution of the amplitude and phase of a multiply scattered field," *Journal of Research of the National Bureau of Standards*, vol. 66D, no. 3, pp. 231–240, May-June 1962.

[38] J. Le Roux and E. Vincent, "Consistent Wiener filtering for audio source separation," *IEEE Signal Processing Letters*, vol. 20, no. 3, pp. 217–220, March 2013.

[39] N. Bertin, R. Badeau, and E. Vincent, "Enforcing harmonicity and smoothness in Bayesian non-negative matrix factorization applied to polyphonic music transcription," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 18, no. 3, pp. 538–549, March 2010.

[40] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the royal statistical society. Series B (methodological)*, vol. 39, no. 1, pp. 1–38, 1977.

[41] P. Magron, J. Le Roux, and T. Virtanen, "Consistent anisotropic Wiener filtering for audio source separation," in *Proc. of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, October 2017, pp. 269–273.

[42] D. R. Hunter and K. Lange, "A tutorial on MM algorithms," *The American Statistician*, vol. 58, no. 1, pp. 30–37, 2004.

[43] C. Févotte and J. Idier, "Algorithms for nonnegative matrix factorization with the beta-divergence," *Neural Computation*, vol. 23, no. 9, pp. 2421–2456, September 2011.

[44] C. Févotte, "Majorization-minimization algorithm for smooth Itakura-Saito nonnegative matrix factorization," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2011, pp. 1980–1983.

[45] A. Lefèvre, F. Bach, and C. Févotte, "Itakura-Saito nonnegative matrix factorization with group sparsity," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2011, pp. 21–24.

[46] P. Magron and T. Virtanen, "Expectation-maximization algorithms for Itakura-Saito nonnegative matrix factorization," in *Proc. of Interspeech*, September 2018, pp. 856–860.

[47] ——, "Towards complex nonnegative matrix factorization with the beta-divergence," in *Proc. of the International Workshop on Acoustic Signal Enhancement (iWAENC)*, September 2018.

[48] M. Abe and J. O. Smith, "Design criteria for simple sinusoidal parameter estimation based on quadratic interpolation of FFT magnitude peaks," in *Audio Engineering Society Convention 117*, May 2004.

[49] http://www.cs.tut.fi/~magron/demos/demo_CISNMF.html.

[50] A. Liutkus, F.-R. Stöter, Z. Rafii, D. Kitamura, B. Rivet, N. Ito, N. Ono, and J. Fontecave, "The 2016 signal separation evaluation campaign," in *Proc. of International Conference on Latent Variable Analysis and Signal Separation (LVA/ICA)*, February 2017, pp. 323–332.

[51] A. Liutkus, J. Pinel, R. Badeau, L. Girin, and G. Richard, "Informed source separation through spectrogram coding and data embedding," *Signal Processing*, vol. 92, no. 8, pp. 1937–1949, 2012.

[52] C. Rohlfing, J. M. Becker, and M. Wien, "NMF-based informed source separation," in *Prof. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2016, pp. 474–478.

[53] C. Rohlfing, J. E. Cohen, and A. Liutkus, "Very low bitrate spatial audio coding with dimensionality reduction," in *Prof. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2017, pp. 741–745.

[54] A. Ozerov, A. Liutkus, R. Badeau, and G. Richard, "Coding-based informed source separation: Nonnegative tensor factorization approach," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 8, pp. 1699–1712, Aug 2013.

[55] B. J. King, "New methods of complex matrix factorization for single-channel source separation and analysis," Ph.D. dissertation, University of Washington, 2012.

[56] E. Vincent, R. Gribonval, and C. Févotte, "Performance measurement in blind audio source separation," *IEEE Transactions on Speech and Audio Processing*, vol. 14, no. 4, pp. 1462–1469, July 2006.

[57] T. Virtanen and A. T. Cemgil, "Mixtures of gamma priors for non-negative matrix factorization based speech separation," in *Proc. of International Conference on Latent Variable Analysis and Signal Separation (LVA/ICA)*, March 2009, pp. 646–653.

[58] D. Kitamura, H. Saruwatari, K. Shikano, K. Kondo, and Y. Takahashi, "Music signal separation by supervised nonnegative matrix factorization with basis deformation," in *Proc. of International Conference on Digital Signal Processing (DSP)*, July 2013, pp. 1–6.

[59] H. Kameoka, H. Kagami, and M. Yukawa, "Complex NMF with the generalized Kullback-Leibler divergence," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2017, pp. 56–60.

[60] S. Leglaive, U. Simsekli, A. Liutkus, R. Badeau, and G. Richard, "Alpha-stable multichannel audio source separation," in *Proc. of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2017, pp. 576–580.

[61] P. Magron, R. Badeau, and A. Liutkus, "Lévy NMF for robust nonnegative source separation," in *Proc. of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, October 2017, pp. 259–263.

[62] A. A. Nugraha, A. Liutkus, and E. Vincent, "Multichannel audio source separation with deep neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 9, pp. 1652–1664, September 2016.

[63] N. Takahashi and Y. Mitsufuji, "Multi-scale multi-band DenseNets for audio source separation," in *Proc. of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, October 2017, pp. 21–25.