

# Multichannel Source Separation Using Time-Deconvolutive CNMF

Thadeu Luiz Barbosa Dias, Wallace Alves Martins, and Luiz Wagner Pereira Biscainho

**Abstract**—This paper addresses the separation of audio sources from convolutive mixtures captured by a microphone array. We approach the problem using complex-valued non-negative matrix factorization (CNMF), and extend previous works by tailoring advanced (single-channel) NMF models, such as the deconvolutive NMF, to the multichannel factorization setup. Further, a sparsity-promoting scheme is proposed so that the underlying estimated parameters better fit the time-frequency properties inherent in some audio sources. The proposed parameter estimation framework is compatible with previous related works, and can be thought of as a step toward a more general method. We evaluate the resulting separation accuracy using a simulated acoustic scenario, and the tests confirm that the proposed algorithm provides superior separation quality when compared to a state-of-the-art benchmark. Finally, an analysis of the effects of the introduced regularization term shows that the solution is in fact steered toward a sparser representation.

**Index Terms**—Blind source separation, convolutive mixture, NMF, deconvolutive NMF

## I. INTRODUCTION

Blind source separation (BSS) is an extensively researched topic with a wide variety of applications [2]. A celebrated example is the use of independent component analysis (ICA) [3, 4, 5] to separate muscular activity interference from brain activity in encephalographic scans [6]. Another interesting example is the use of BSS in speech enhancement for hearing aid devices [7].

Among traditional techniques for source separation, non-negative matrix factorization (NMF) [8], a single-channel method, has been successfully employed in the literature [9]. NMF factorizes an input matrix with non-negative entries into two lower-rank matrices with non-negative entries, and is able to extract the most significant components that explain the observed data, i.e., a model and a set of parameters that produce a satisfactory estimate of the data. In comparison with other rank-reducing methods such as singular-value decomposition (SVD), the fact of dealing only with non-negative quantities is a distinctive feature of NMF: this is suitable for parameters

that are non-negative by nature as is the case of magnitude or power spectra, and prevents mutual cancelling of components by destructive interference, since the model is strictly additive.

In the source separation scenario, the resulting NMF factors from a mixture spectrogram can be thought of as a set of spectral signatures and temporal activation patterns [10]. It is expected that subsets of the extracted signatures explain each source, and a typical challenge is how to assign which components correspond to each source. Usually, this assignment relies on some other prior information.

A more powerful way to perform source separation is to exploit spatial information, as in the case of multichannel processing methods. The complex-valued NMF (CNMF) [11] is a development in this direction, with the introduction of Hermitian positive semidefinite matrices as data points, derived from the measured signals' complex-valued spectrograms. Building on the CNMF model, an alternative factorization, with geometric constraints, is proposed in [12], providing spatially-coherent estimation, thus enhancing the separation quality of the method.

In [11, 12, 13], the factorization models were based on the standard NMF; by construction, the standard NMF model does not take into account the temporal sequence of samples, that is, a random shuffle in the time frames is irrelevant from the decomposition viewpoint. This approach, therefore, disregards useful continuity structures that can be observed in some data, such as musical samples, and often, the emission signatures cannot be efficiently represented by single NMF components. In order to address this limitation, we propose using a deconvolutive NMF (NMFD) scheme [14]. The extended NMFD signatures have a user-defined span of time frames, opening up the possibility for a more concise representation of some musical emissions, and eventually enhancing the separation quality.

In this paper we show that the deconvolutive NMF model can be tailored to the CNMF framework with good results, and that we may regularize the related cost function toward a sparse solution. We start with the problem statement in Section II, give a brief introduction to the NMF models in Section III, describe the transformations of the original data into the CNMF data points in Section IV, and the constrained construction of channel matrices as well as the application of the deconvolutive model in Section V. We derive the estimation framework for a Euclidean cost function in Section VI, and present the numerical results in Section VII.

A note on notation: In this paper, we denote scalars by regular lower-case letters, e.g.  $x$ , vector variables are lower-case bold letters, e.g.  $\mathbf{x}$ , matrices are upper-case bold, e.g.  $\mathbf{X}$ , and

This work is an extended version of our previous work [1], presented at the XXXVII Brazilian Symposium on Telecommunications and Signal Processing (SBrt-2019).

Mr. Thadeu L.B. Dias is with the Electrical Engineering Program (PEE/COPPE) of the Federal University of Rio de Janeiro (UFRJ); Prof. Wallace A. Martins is with the Department of Electronics and Computer Engineering (DEL/Poli) & PEE/COPPE, UFRJ (on leave), and with the Interdisciplinary Centre for Security, Reliability and Trust (SnT) of the University of Luxembourg; Prof. Luiz W.P. Biscainho is with DEL/POLI & PEE/COPPE, UFRJ. E-mails: {thadeu.dias, wallace.martins, wagner}@smt.ufrj.br. This work was partially supported by CAPES (88882.331631/2019-01), CNPq (PQ 306331/2017-9), and FAPERJ.

Digital Object Identifier: 10.14209/jcis.2020.11

higher-order tensors are calligraphic upper-case bold, e.g.  $\mathcal{X}$ . An indexed element from a higher-order tensor corresponds to a slice of the original tensor, e.g.  $[\mathcal{B}_k]_{it} = [\mathcal{B}]_{ikt}$ . Conversely, when some index is omitted and the variable gets promoted to a higher-order form, we refer to the entire collection of elements with the same name, that is,  $\mathcal{W} = \{\mathbf{W}_{io} \forall i, o\}$ .

## II. PROBLEM STATEMENT

We state our problem as that of separating meaningful acoustic sources captured by a microphone array with known geometry placed in a reverberant environment. Meaningful in this case means a source that is spatially concentrated, and emissions from different positions are considered as being from different sources. We consider that the propagation media is linear and time-invariant, and effects such as reverberation and multipath-propagation are then modeled by the unknown source-sensor channel impulse responses.

Under such considerations and assuming an array with  $M$  microphones, the acquired signals can then be described by the following model:

$$x_m(t) = \sum_{q=1}^Q (h_{qm} * s_q)(t), \quad (1)$$

where  $Q$  is the true number of sources,  $x_m(t)$  is the  $m^{\text{th}}$  sensor measurement,  $h_{qm}(t)$  is the impulse response relative to the channel between the source-sensor pair  $(q, m)$ ,  $s_q(t)$  is the true emission of source  $q$ , and  $*$  denotes the convolution operation. We would then like to be able to reconstruct estimates for the individual source images  $\hat{y}_{qm}(t) = (h_{qm} * s_q)(t)$ , as captured by the array.

We now describe the basic NMF formulations, and the steps toward integrating Eq. (1) into the non-negative framework.

## III. NMF BASICS

Non-negative matrix factorization is, at its core, the factorization of a matrix of non-negative entries into two reduced-rank matrices with non-negative entries.

In its most basic formulation [15], we have a data matrix  $\mathbf{X} \in \mathbb{R}_+^{I \times L}$ , and want to find a  $K$ -rank approximation of the original matrix according to a suitable criterion, subject to a non-negative constraint on the components; the approximation is given by matrices  $\mathbf{B} \in \mathbb{R}_+^{I \times K}$  and  $\mathbf{G} \in \mathbb{R}_+^{K \times L}$  such that

$$\mathbf{X} \approx \mathbf{B}\mathbf{G} = \sum_{k=1}^K \mathbf{b}_k \mathbf{g}_k^T, \quad (2)$$

as illustrated in Fig. 1.

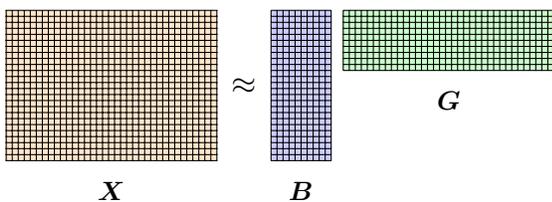


Fig. 1. Simple NMF model.

Usually, the approximation criterion is described by some cost function  $\mathcal{L}(\mathbf{X}, \mathbf{B}\mathbf{G})$ , and the factors are obtained by jointly minimizing the overall cost:

$$\min_{\mathbf{B}, \mathbf{G}} \mathcal{L}(\mathbf{X}, \mathbf{B}\mathbf{G}) \quad (3)$$

subject to non-negative constraints.

A useful interpretation of the extracted factors is the following: if the rows of  $\mathbf{X}$  are features, and each column of  $\mathbf{X}$  represents the set of features from a given observation; then, after a successful factorization, the columns of  $\mathbf{B}$  can be thought of as recurrent signatures, forming a basis for the observations, and the rows of  $\mathbf{G}$  as coefficients denoting the activation of each corresponding signature [9]. The presence of a given signature could then be related to some hidden factor of the underlying data; for instance, in a spectrogram factorization, an extracted signature can be linked to a particular emission from a musical instrument.

Often, however, the problem is so underdetermined that simply relying on the fit yielded by the minimization of  $\mathcal{L}$  is not satisfactory. A common solution is then to restrict the model, imposing extra constraints on the obtained factors, such as sparsity, smoothness, or orthogonality [16, 17, 18].

One particular limitation of the simple NMF model is the fact that the extracted patterns are static, that is, the shape of the estimate provided by each component  $\mathbf{b}_k$  is equal across all observations. Taking temporal continuity into consideration, such that the sequence of observations (columns of  $\mathbf{X}$ ) are in fact a time sequence, a more powerful model can be designed, allowing the patterns to be longer in length than a single observation. The deconvolutive-NMF is a generalization that provides such features by extending the basis vectors  $\mathbf{b}_k$  to matrices  $\mathbf{B}_k \in \mathbb{R}_+^{I \times T}$ , for some user-chosen  $T$ ; the corresponding activation vector  $\mathbf{g}_k$  is then time-shifted by  $t$  and applied to each subcomponent  $\mathbf{b}_{kt}$ :

$$\mathbf{X} \approx \sum_{k=1}^K \sum_{t=1}^T \mathbf{b}_{kt} \mathbf{g}_k^T; \quad (4)$$

the dataflow for the deconvolutive NMF is illustrated in Fig. 2.

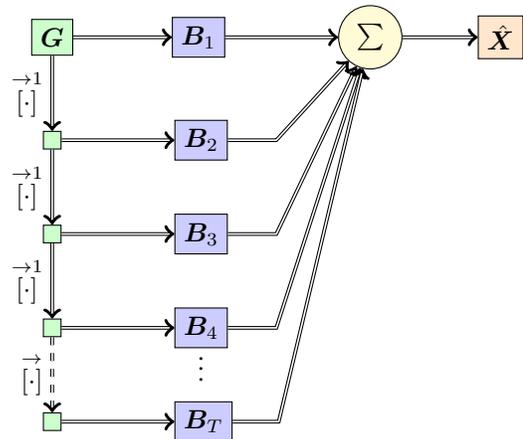


Fig. 2. Deconvolutive NMF model generating the approximation  $\hat{\mathbf{X}}$ . Here, the overall basis tensor  $\mathcal{B}$  is sliced across the deconvolutive mode, indexed by  $t$ .

The right-shift operator  $\overset{\rightarrow}{[\cdot]}$  is such that for a given vector

$$\overset{\rightarrow}{[g_1, g_2, \dots, g_L]} = [0, g_1, \dots, g_{L-1}];$$

for  $t = 0$ , the operator reduces to identity, and with  $t > 1$ , it is equivalent to  $t$  successive applications of the operator with unit shift. The left-shift operator is defined analogously, and the extension to matrices and tensors follows naturally, shifting over some specified dimension.

The main idea behind deconvolutive-NMF in musical analysis is that longer musical note profiles can be efficiently modeled by the framework, leading to an increase in separation accuracy.

#### IV. SIGNAL REPRESENTATION

We now describe the steps toward a phase-aware non-negative multichannel representation, aimed at the decomposition in non-negative factors. The first thing to do is to extend the basic data matrix  $\mathbf{X}$ , introduced in Section III, to the multichannel case. One way to do this is by associating each entry of this data matrix with an  $M \times M$  matrix modeling relations between pairs of sensors; but then we also have to extend the notion of non-negativeness. When dealing with real numbers, being non-negative means belonging to the convex cone<sup>1</sup>  $\mathbb{R}_+$ . Positive-semidefiniteness is the matrix counterpart of the scalar non-negativeness property; indeed, positive-semidefiniteness is preserved under conical combination of positive-semidefinite (PSD) matrices.<sup>2</sup>

Translating the relationship of Eq. (1) to the short-time Fourier transform (STFT) domain, each (complex) frequency-time point measurement is

$$x_{ilm} = \sum_{q=1}^Q h_{iqm} s_{ilq}, \quad (5)$$

where  $i$  denotes the frequency bin,  $l$  is an index for the time frame,  $h_{iqm}$  is the frequency response at bin  $i$  of the channel relative to the source-sensor pair  $(q, m)$ , and  $s_{ilq}$  is the STFT of the emission of source  $q$  at frequency-time point  $(i, l)$ .

In order to represent the overall measurements as Hermitian PSD matrices, Sawada et al. [11] take the outer product of the vector  $\mathbf{x}_{il}$ —formed by the measurements across all sensors at a single frequency-time point—with itself, obtaining the matrices

$$\mathbf{P}_{il} = \mathbf{x}_{il} \mathbf{x}_{il}^H = \sum_{q=1}^Q \sum_{q'=1}^Q \mathbf{h}_{iq} \mathbf{h}_{iq'}^H s_{ilq} \bar{s}_{ilq'}, \quad (6)$$

where  $\mathbf{v}^H$  is the Hermitian of  $\mathbf{v}$  and  $\bar{z}$  is the complex conjugate of  $z$ . Additionally, assuming that different sources  $q$  and  $q'$  are orthogonal, and then transposing the correlation property to the

<sup>1</sup>A subset  $\mathcal{C}$  of a given vector space is a cone if  $\alpha \mathbf{x} \in \mathcal{C}$  for any  $\mathbf{x} \in \mathcal{C}$  and any  $\alpha \geq 0$ . It is a convex cone when it is closed for conical combinations, i.e., when  $\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 \in \mathcal{C}$  for any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C}$  and any  $\alpha_1, \alpha_2 \geq 0$ .

<sup>2</sup>The space of PSD matrices,  $\mathbb{S}_+$ , is itself a cone, that is, any element  $\mathbf{S} \in \mathbb{S}_+$  can be written as a conical combination of some basis matrices.

STFT coefficients, a useful approximation can be obtained by neglecting the crossed terms:

$$\mathbf{P}_{il} \approx \sum_{q=1}^Q \mathbf{h}_{iq} \mathbf{h}_{iq}^H |s_{ilq}|^2. \quad (7)$$

Lastly, in order to factor a magnitude instead of a power spectrum, the STFT coefficients are mapped through a magnitude square-root function applied elementwise, as proposed in [11]:  $\phi(z) = \frac{z}{\sqrt{|z|}}$ . In this new context, equation (6) is replaced by the corresponding expression

$$\mathbf{X}_{il} = \phi(\mathbf{x}_{il}) \phi(\mathbf{x}_{il})^H \approx \sum_{q=1}^Q \boldsymbol{\eta}_{iq} \boldsymbol{\eta}_{iq}^H |s_{ilq}|, \quad (8)$$

where  $\boldsymbol{\eta}_{iq}$  is a factor similar to  $\mathbf{h}_{iq}$ , allowing to define  $\mathbf{H}_{iq} = \boldsymbol{\eta}_{iq} \boldsymbol{\eta}_{iq}^H$  as a matrix that encodes the phase properties of source  $q$  at bin  $i$ . The entries of  $\mathbf{H}_{iq}$  encode the phase difference between the responses of each channel pair. By the outer product construction,  $\mathbf{H}_{iq}$  preserves phase information without actually modeling the absolute phase of the measurements. Expression (8) then motivates a joint factorization of the sources' magnitude spectra and spatial-property matrices.

#### V. FACTORIZATION MODEL

We now describe the steps toward the factorization of the multichannel model in (8) into non-negative factors. The standard CNMF [11] simply attaches to each NMF component some phase information in the form of a set of PSD matrices, as will be described in Subsection V-A. A limitation of this approach is that no structure is imposed on the family of PSD matrices associated with each component; we intend to separate sources based on spatial cues, and one could expect some form of coherence in terms of how the spatial information is encoded within a single component, as shown in [12]. In Subsection V-B, prior knowledge about the sensor array geometry is incorporated into the process in the form of beamforming kernels, from which the spatial matrices associated with each component are derived from. Finally, in Subsection V-C, the complete factorization model is glued together, defining some useful additional constraints, and in Subsection V-D the process of recovering the source image spectrograms from the CNMF parameters is described.

The main idea of factorization model is to explore the compressibility of the magnitude representation in order to find  $K$  components that best explain the measurements. In the context of CNMF [11], we seek to explain the measured data points  $\mathbf{X}_{il}$  as non-negative combinations of positive semidefinite matrices. We assign to each NMF component a family of spatial-property matrices, and cluster components based on their spatial properties when reconstructing the sources. The phase-magnitude model for the measured data points can be written as

$$\mathbf{X}_{il} \approx \sum_{k=1}^K \mathbf{H}_{ik} \tilde{s}_{ilk}, \quad (9)$$

where  $\mathbf{H}_{ik}$  encodes spatial properties for a component and  $\tilde{s}_{ilk}$  is a magnitude estimate that shall be computed through

the NMF framework. In the following, we detail how the parameters are obtained.

### A. Magnitude activation model

In the single channel case, the standard NMF finds a low-rank approximation to some data matrix  $\mathbf{X} \in \mathbb{R}_+^{I \times L}$  as the product of two smaller non-negative matrices  $\mathbf{B} \in \mathbb{R}_+^{I \times K}$  and  $\mathbf{G} \in \mathbb{R}_+^{K \times L}$ , where, usually,  $K \ll \text{Rank}(\mathbf{X})$ . If the input data matrix consists of a magnitude spectrogram, with bins as rows and time frames as columns, a useful interpretation of the extracted matrices arises: the columns  $\mathbf{b}_k \in \mathbb{R}_+^I$  of  $\mathbf{B}$  are spectral signatures present in the measurements, and the rows  $\mathbf{g}_k \in \mathbb{R}_+^L$  of  $\mathbf{G}$  are activation patterns for such signatures across time frames. The application of the simple NMF to the CNMF model would lead to the magnitude estimate  $\check{s}_{ilk} = b_{ik}g_{kl}$ . We propose the estimation of  $\check{s}_{ilk}$  through a deconvolutive NMF model as presented in Section III. Using the deconvolutive model, the estimate of the instantaneous magnitudes due to the  $k^{\text{th}}$  component is

$$\check{s}_{ilk} = \sum_{t=1}^T b_{itk} [\mathbf{g}_k]_l^{\overrightarrow{t-1}}. \quad (10)$$

In fact, the standard form of the CNMF with the deconvolutive model can be written as

$$\hat{\mathbf{X}}_{il} = \sum_{k=1}^K \mathbf{H}_{ik} \sum_{t=1}^T b_{itk} [\mathbf{g}_k]_l^{\overrightarrow{t-1}}, \quad (11)$$

and this model shares the single-channel deconvolutive NMF properties of being able to efficiently extract spectral patterns that vary with time. Considering that continuous emissions occur in many audio signals, this model is appropriate when moving toward a more powerful separation model.

### B. Spatial covariance model

We apply the direction-of-arrival (DoA) based factorization method introduced by Nikunen and Virtanen [12] to the channel matrices  $\mathbf{H}_{ik}$ . An issue with the unconstrained estimation of matrices  $\mathbf{H}_{ik}$  is that there is no guarantee that the set of matrices  $\mathcal{H}_k$  (all matrices  $\mathbf{H}_{ik}$  with fixed  $k$ ) actually encodes a single coherent single-input multiple-output channel between some component and the sensor array. Instead, the set  $\mathcal{H}_k$  is constructed as a non-negative linear combination of geometrically-defined beamforming kernel matrices  $\mathbf{W}_{io}$ .

Consider the scheme depicted in Fig. 3: for a sufficiently far emission source somewhere along the direction of  $\mathbf{k}_o$  (a unit-length vector), such that the wavefronts can be considered planar, the relative time delay between sensor acquisitions can be defined in terms of the sensor array geometry, wave propagation velocity, and incidence direction. In fact, the difference in propagation length can be calculated as the inner product  $\langle \mathbf{p}_{m'} - \mathbf{p}_m, \mathbf{k}_o \rangle$  so that the relative time-difference of arrival (TDoA) is simply

$$\tau_{mm'}(\mathbf{k}_o) = \frac{\langle \mathbf{p}_{m'} - \mathbf{p}_m, \mathbf{k}_o \rangle}{c}, \quad (12)$$

where  $c$  denotes the wave propagation velocity.

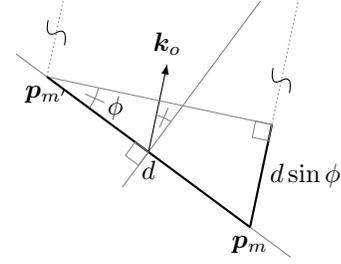


Fig. 3. TDoA as function of array geometry and wave incidence direction.

It is straightforward to find the frequency-dependent phase lag using Fourier transform properties, and from the non-normalized STFT bin frequencies, the per-bin phase lag can be calculated as

$$\theta_{mm'}(f_i, \mathbf{k}_o) = -2\pi f_i \tau_{mm'}(\mathbf{k}_o). \quad (13)$$

The idea for modeling  $\mathcal{H}$  is to sample  $O$  directions from the unit sphere around the array and form the beamforming kernels  $\mathbf{W}_{io}$  for all STFT bins for each DoA sample. The beamforming kernels are  $M \times M$  Hermitian matrices containing the relative phase shifts (for a set frequency  $f_i$  and direction  $\mathbf{k}_o$ ) as complex factors:

$$\mathbf{W}_{io} = \begin{bmatrix} 1 & e^{j\theta_{12}(f_i, \mathbf{k}_o)} & \dots & e^{j\theta_{1M}(f_i, \mathbf{k}_o)} \\ e^{j\theta_{21}(f_i, \mathbf{k}_o)} & 1 & \dots & e^{j\theta_{2M}(f_i, \mathbf{k}_o)} \\ \vdots & \vdots & \ddots & \vdots \\ e^{j\theta_{M1}(f_i, \mathbf{k}_o)} & e^{j\theta_{M2}(f_i, \mathbf{k}_o)} & \dots & 1 \end{bmatrix}.$$

Finally, the channel matrices  $\mathbf{H}_{ik}$  can be described in terms of the DoA kernels as conic combinations

$$\mathbf{H}_{ik} = \sum_{o=1}^O z_{ko} \mathbf{W}_{io}, \quad (14)$$

where the factors  $z_{ko} \in \mathbb{R}_+$  are shared across all frequencies for a given component, making this a spatially coherent factorization. Additionally, this model allows the spatial properties encoded by the vectors  $\mathbf{z}_k$  to be clustered, since it is expected that components with similar spatial signatures belong to the same source.

### C. Complete model

The complete model for the measured covariance matrices in terms of deconvolutive CNMF parameters can be written as

$$\hat{\mathbf{X}}_{il} = \sum_{k=1}^K \sum_{o=1}^O z_{ko} \mathbf{W}_{io} \sum_{t=1}^T b_{ikt} [\mathbf{g}_k]_l^{\overrightarrow{t-1}}. \quad (15)$$

This factorization has a possible scaling ambiguity, so additional constraints are introduced, namely  $\sum_o z_{ko}^2 = 1$ ,  $\sum_l g_{kl}^2 = 1$ , and  $\|\mathbf{W}_{io}\|_F = 1$ , where  $\|M\|_F = \sqrt{\text{tr}(M^H M)}$  denotes the Frobenius norm of a matrix  $M$ . In order to find the best estimates for the parameters, a statistical model can be employed, and an estimation framework can be built—see Section VI.

#### D. Source reconstruction

Given the CNMF parameter estimates, the per-source spectral images can be reconstructed through a Wiener filter

$$\mathbf{y}_{ilq} = \mathbf{x}_{il} \frac{\sum_{k,o} \beta_{qk} z_{ko} \sum_t b_{ikt} [\mathbf{g}_k]_l^{t-1}}{\sum_{q,k,o} \beta_{qk} z_{ko} \sum_t b_{ikt} [\mathbf{g}_k]_l^{t-1}}, \quad (16)$$

where  $\beta_{qk}$  are learned membership coefficients relating a given component  $k$  to source  $q$ . The membership coefficients can be obtained through a regular clustering algorithm, such as k-means, c-means, or even NMF (considering that the spatial factors  $z_{ko}$  are also non-negative). The overall effect is to multiply the input spectrograms by a mask of ratios of estimated spectral magnitudes, preserving the original phase. Finally, the discrete-time version of  $\hat{y}_{qm}(t)$  can be retrieved from the inverse-STFT of the coefficients.

### VI. PARAMETER ESTIMATION

Following previous works [11, 12], the generative model adopted for the entries of the data matrices  $\mathbf{X}_{il}$  is that of independent complex Gaussian variables with unit variance, which allows to recast the likelihood estimate as a squared error minimization problem. In fact, the likelihood function for the parameters, considering the overall measurements, can be written as

$$\mathcal{L}(\mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}) \propto \prod_{i=1}^I \prod_{l=1}^L \exp\left(-\|\mathbf{X}_{il} - \hat{\mathbf{X}}_{il}\|_{\mathbb{F}}^2\right), \quad (17)$$

leading to the alternative cost function

$$\ell(\mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}) = \sum_{i=1}^I \sum_{l=1}^L \|\mathbf{X}_{il} - \hat{\mathbf{X}}_{il}\|_{\mathbb{F}}^2, \quad (18)$$

related to the colog-likelihood. It can be useful to embed some prior on the parameters. This knowledge can be directly related to a regularization factor, steering the algorithm toward a solution with some desirable properties. In this paper we consider the generative model for the spectral signatures  $b_{ikt}$  as a one-sided exponential distribution with some scaling factor  $\alpha_{\mathbf{B}}$ , leading to the regularized likelihood<sup>3</sup>

$$\mathcal{L}_{\mathbf{R}} \propto \exp(-2\alpha_{\mathbf{B}} \|\mathbf{B}\|_1) \mathcal{L}(\mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}) \quad (19)$$

and a regularized cost function

$$\ell_{\mathbf{R}}(\mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}) = 2\alpha_{\mathbf{B}} \|\mathbf{B}\|_1 + \sum_{i=1}^I \sum_{l=1}^L \|\mathbf{X}_{il} - \hat{\mathbf{X}}_{il}\|_{\mathbb{F}}^2, \quad (20)$$

where the tensor  $\ell_1$ -norm is defined as  $\sum_{i,k,t} |b_{ikt}|$ . This is inspired by a LASSO [19] regression, aiming to enforce the selectivity of signatures  $\mathbf{B}_k$  and, consequently, produce a sparser representation, suitable to pitched audio signals.

<sup>3</sup>We extract a factor 2 from  $\alpha_{\mathbf{B}}$  for convenience.

#### A. Minimization procedure

To obtain the maximum a posteriori estimate, one must minimize the regularized cost (20). The convoluted interdependence among the parameters  $\mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}$  in this cost function hinders its direct (gradient-based) minimization, calling for alternative procedures. A possible alternative would be to minimize an auxiliary function, related to the actual cost function, that depends on an extra set of parameters while being, in some sense, smoother. The individual gradients w.r.t. the CNMF parameters of the auxiliary function can be straightforwardly derived, and the extra set of parameters can be chosen in a way that explicit computation is avoided. In this subsection, we shall describe a convenient auxiliary function, along with the specific choice of additional parameters as well as the update rules for the non-negative tensors. At last, the update for the kernel matrices, which requires additional projection and normalization steps, is described.

While (20) is non-convex relative to the CNMF parameters, it is individually convex over  $\mathbf{Z}, \mathbf{W}, \mathbf{B}$ , and  $\mathbf{G}$ . Thus, a block relaxation minimization scheme [20] may be employed with good results. We consider an auxiliary function to (20), namely:

$$\ell_{\mathbf{R}}^+ = 2\alpha_{\mathbf{B}} \|\mathbf{B}\|_1 + \sum_{i,l,k,o,t} \frac{1}{r_{ilkot}} \|\mathbf{S}_{ilkot} - z_{ko} \mathbf{W}_{io} b_{ikt} [\mathbf{g}_k]_l^{t-1}\|_{\mathbb{F}}^2, \quad (21)$$

where  $r_{ilkot}$  are any positive variables satisfying  $\sum_{k,o,t} r_{ilkot} = 1$ , and  $\mathbf{S}_{ilkot}$  are Hermitian matrices satisfying  $\sum_{k,o,t} \mathbf{S}_{ilkot} = \mathbf{X}_{il}$ , for which we derive the following proposition:

**Proposition 1.** For all  $\mathbf{S}_{ilkot} \in \mathbb{S}$  and  $r_{ilkot} \in \mathbb{R}_+$  such that  $\sum_{k,o,t} \mathbf{S}_{ilkot} = \mathbf{X}_{il}$  and  $\sum_{k,o,t} r_{ilkot} = 1$ , then

$$\ell_{\mathbf{R}}^+(\mathcal{S}, \mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}) \geq \ell_{\mathbf{R}}(\mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}) \quad \text{and} \quad (22)$$

$$\min_{\mathcal{S}} \ell_{\mathbf{R}}^+(\mathcal{S}, \mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}) = \ell_{\mathbf{R}}(\mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}). \quad (23)$$

Furthermore,

$$\begin{aligned} \mathbf{S}_{ilkot}^* &\triangleq \arg \min_{\mathbf{S}_{ilkot}} \ell_{\mathbf{R}}^+(\mathcal{S}, \mathbf{Z}, \mathbf{W}, \mathbf{B}, \mathbf{G}) \\ &= z_{ko} \mathbf{W}_{io} b_{ikt} [\mathbf{g}_k]_l^{t-1} - r_{ilkot} \mathbf{E}_{il}. \end{aligned} \quad (24)$$

*Proof.* See Appendix A.  $\square$

Through the majorizer conditions (22) and (23), the individual minimization of (21) across the CNMF variables with  $\mathcal{S}$  set as the optimal  $\mathcal{S}^*$  is guaranteed to be non-increasing. With the auxiliary definition

$$\hat{\mathbf{x}}_{il} = \sum_{k,o,t} z_{ko} b_{ikt} [\mathbf{g}_k]_l^{t-1}, \quad (25)$$

a useful way to define  $r_{ilkot}$  arises, namely

$$r_{ilkot} = \frac{z_{ko} b_{ikt} [\mathbf{g}_k]_l^{t-1}}{\hat{\mathbf{x}}_{il}}. \quad (26)$$

Although this definition is not strictly positive, it is safe to ignore the zeroed values, since

$$\lim_{r_{ilkot} \rightarrow 0} \frac{1}{r_{ilkot}} \|\mathbf{S}_{ilkot}^* - z_{ko} \mathbf{W}_{io} b_{ikt} [\mathbf{g}_k]_l\|_F^2 = 0,$$

and this definition allows for implicit computation of the  $\mathbf{S}_{ilkot}^*$  factors.

Replacing  $\mathbf{S}^*$  and  $r_{ilkot}$  with their definitions, the multiplicative rules for non-negative factors can be obtained (see Appendix B):

$$z_{ko} \leftarrow z_{ko} \left[ \frac{\sum_{i,l,t} (\hat{x}_{il} + \text{tr}(\mathbf{E}_{il} \mathbf{W}_{io})) b_{ikt} [\mathbf{g}_k]_l}{\sum_{i,l,t} \hat{x}_{il} b_{ikt} [\mathbf{g}_k]_l} \right], \quad (27)$$

$$b_{ikt} \leftarrow b_{ikt} \left[ \frac{\sum_{l,o} (\hat{x}_{il} + \text{tr}(\mathbf{E}_{il} \mathbf{W}_{io})) z_{ko} [\mathbf{g}_k]_l}{\alpha_{\mathbf{B}} + \sum_{l,o} \hat{x}_{il} z_{ko} [\mathbf{g}_k]_l} \right], \quad (28)$$

$$g_{kl} \leftarrow g_{kl} \left[ \frac{\sum_{i,o,t} ([\hat{\mathbf{x}}_i]_l + \text{tr}([\mathbf{E}_i]_l \mathbf{W}_{io})) z_{ko} b_{ikt}}{\sum_{i,o,t} [\hat{\mathbf{x}}_i]_l z_{ko} b_{ikt}} \right]. \quad (29)$$

The update process for the kernel matrices is slightly different, since only magnitude optimizations are allowed, and the positive semidefinite constraint must be accounted for. What follows is an optimization scheme similar to a projected gradient algorithm: the possibly unfeasible point that minimizes the cost function is calculated as (see Appendix C)

$$\hat{\mathbf{W}}_{io} \leftarrow \frac{\sum_{l,k,t} z_{ko} b_{ikt} [\mathbf{g}_k]_l [\hat{x}_{il} \mathbf{W}_{io} + \mathbf{E}_{il}]}{\sum_{l,k,t} \hat{x}_{il} z_{ko} b_{ikt} [\mathbf{g}_k]_l}; \quad (30)$$

this point is projected onto the positive semidefinite cone by rectification of its eigenvalues:

$$\mathbf{V}_{io} \mathbf{\Lambda}_{io} \mathbf{V}_{io}^H \leftarrow \hat{\mathbf{W}}_{io} \quad (31)$$

$$\hat{\mathbf{W}}_{io}^+ \leftarrow \mathbf{V}_{io} \mathbf{\Lambda}_{io}^+ \mathbf{V}_{io}^H; \quad (32)$$

finally, only the entries' magnitudes are updated, as the true update is obtained as

$$\mathbf{W}_{io} \leftarrow \text{abs}(\hat{\mathbf{W}}_{io}^+) \odot \text{sign}(\mathbf{W}_{io}), \quad (33)$$

where  $\text{abs}(\cdot)$  and  $\text{sign}(\cdot)$  both operate elementwise on their arguments, and  $\odot$  denotes the matrix Hadamard product. Concerning the scaling factors, after each update,  $z_k$  and  $g_k$  are normalized to unity, while the reciprocal correction factor is applied to  $\mathbf{B}_k$ , that is:

$$v_k \leftarrow \|z_k\|_2 : z_k \leftarrow \frac{z_k}{v_k} \quad \mathbf{B}_k \leftarrow v_k \mathbf{B}_k, \quad \text{and}$$

$$v_k \leftarrow \|g_k\|_2 : g_k \leftarrow \frac{g_k}{v_k} \quad \mathbf{B}_k \leftarrow v_k \mathbf{B}_k.$$

Similarly,  $\mathbf{W}_{io}$  is rescaled to unity Frobenius norm, but no rescaling of other parameters is needed:

$$\mathbf{W}_{io} \leftarrow \frac{\mathbf{W}_{io}}{\|\mathbf{W}_{io}\|_F}.$$

## VII. NUMERICAL RESULTS

The program developed to assess the accuracy of the proposals was coded in Python, using TensorFlow v1.14, and executed on an Intel Xeon Gold 5120. We attempt the separation of two sound sources positioned 90 degrees apart, from four mixtures synchronously captured by omnidirectional microphones placed approximately 8 cm apart from each other, in the form of a tetrahedron. The audio tracks are two musical samples about 20 s long, sampled at 22.05 kHz. A closed room of dimensions 5 m  $\times$  4 m  $\times$  3 m (length, width, height) with  $\text{RT}_{60} \approx 450$  ms was simulated using CATT-Acoustic v9.0c [21], with the microphone array at the center and the sources on the horizontal plane 1.5 m away from it, as depicted in Fig. 4. The DoA vectors were randomly sampled from the unit sphere in a way to approximately maximize the cosine distance between the closest vectors, with  $O = 110$  directions. An illustration of the DoA coordinates is shown in Fig. 5.

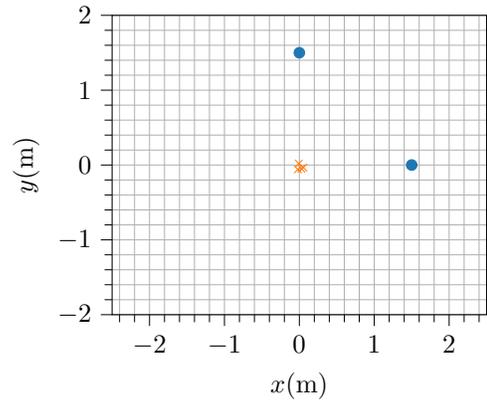


Fig. 4. Top-down view of the source-sensors arrangement inside the room. The array is indicated by orange cross markers, and the two sources by blue dots. The reference axes are aligned with the walls of the room.

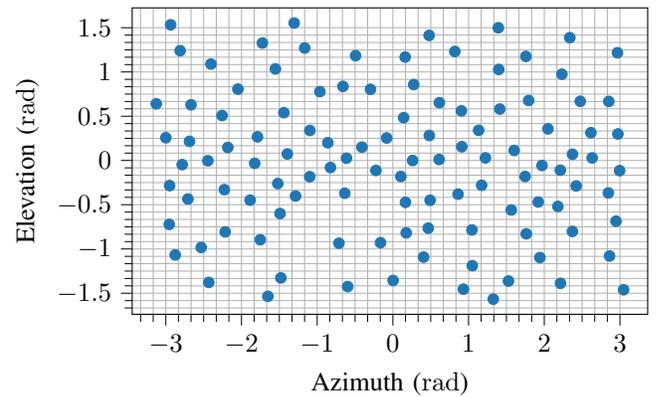


Fig. 5. Coordinates of the DoA samples on the sphere.

Square-rooted Hanning windows were used for STFT analysis and synthesis, with 50% overlap. Frame length was chosen as 1024 samples, corresponding to approximately 46 ms,  $I = 513$  bins, and  $L = 978$  frames. We considered three scenarios: a baseline run, with  $(K = 60, T = 1, \alpha_{\mathbf{B}} = 0)$ , equivalent to [12]; a deconvolutive run, with  $(K = 60, T = 5, \alpha_{\mathbf{B}} = 0)$  testing the proposed extension; and a sparsity-promoting run,

with  $(K = 60, T = 5, \alpha_{\mathcal{B}} = 0.5)$ . The algorithm was run for 500 iterations in all cases.

We performed weighted k-means on the vectors  $\mathbf{z}_k$ , with weights corresponding to the component energy  $\|\mathbf{B}_k\|_F^2$ . The source-component membership coefficients  $\beta_{qk}$  were set to 1 or 0 based on the obtained clustering.

First, we analyze the directional sensitivity of the method. With k-means clustering, it is possible to retrieve the centroids corresponding to the spatial signatures for each source; based on the spherical coordinates of the DoA vectors and the obtained centroid coefficients, a bivariate spherical spline interpolation was used to create a visualization of the spatial properties for each source. The results are depicted in Figs. 6 and 7. The true positions are depicted with blue dots.

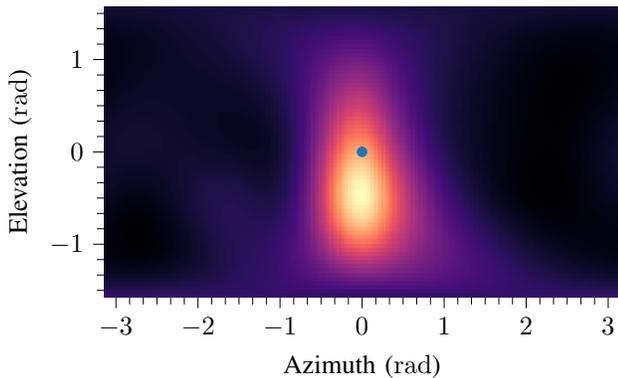


Fig. 6. Visualization of spatial signature for the first source. A main lobe is clearly discernible, aligned with the source’s true direction (blue dot).

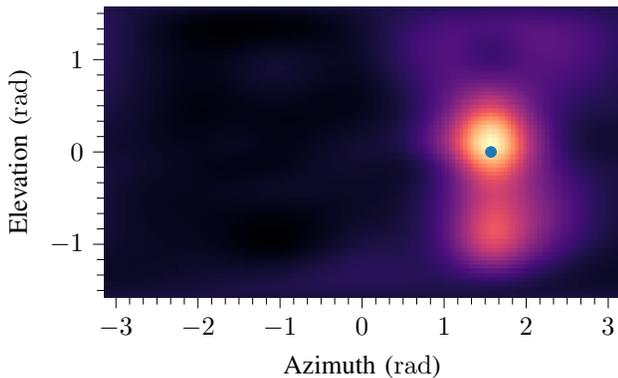


Fig. 7. Visualization of spatial signature for the second source. A main lobe is clearly discernible, closely matching the source’s true direction (blue dot).

The centroids’ peak activations closely match the true directions for each source, such that a rough estimate of the sources directions can be obtained from the method, although the estimate is likely to deteriorate in heavily reverberant environments.

The separation quality [22, 23] was measured using the `mir_eval` suite [24, 25]. From the suite’s provided image evaluation function, we obtain the source-to-distortion ratio (SDR) and source-to-interference ratio (SIR) for each source, and the results are depicted in Fig. 8. SDR encompasses several types of distortions into a single metric, being a robust form of evaluation for BSS techniques; SIR, on the other hand

measures only the overall crosstalk between estimated sources, being another useful metric for separation evaluation.

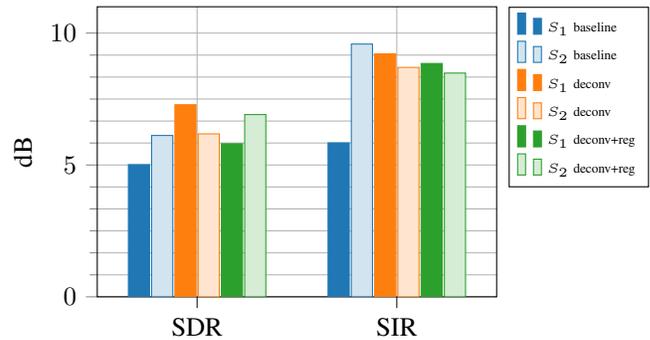


Fig. 8. Separation metrics for the three tests. Both deconvolutive and regularized tests obtained superior averages than the baseline, while the regularized factorization was slightly worse than the non-regtest.

Due to the stochastic nature of the fit, an unbalance between sources can be noted in all scenarios, but a slight advantage can be noted in the deconvolutive case. The regularization slightly degraded the overall measures in comparison to the unregularized deconvolutive test, although they are still superior (on average) to those of the baseline case.

We would like to investigate the sparsity-inducing properties on the extracted signatures  $\mathcal{B}$ . While the factors obtained by fitting the data using the Euclidean cost are often sparse, they are not the sparsest possible representation. Due to the nature of the range of the values assumed by  $b_{ikt}$ , we plotted the histogram of  $\log_{10}(b_{ikt})$  using 50 bins as shown in Fig. 9. A considerable increase in frequency of the lowest bin can be observed: the count on the lowest bin for the unregularized run was 84991 hits, while the regularized test had 145982 hits, a 71% increase. It is then clear that the regularizing component of  $\mathcal{L}_R$  steered the factor search toward a solution with higher selectivity.

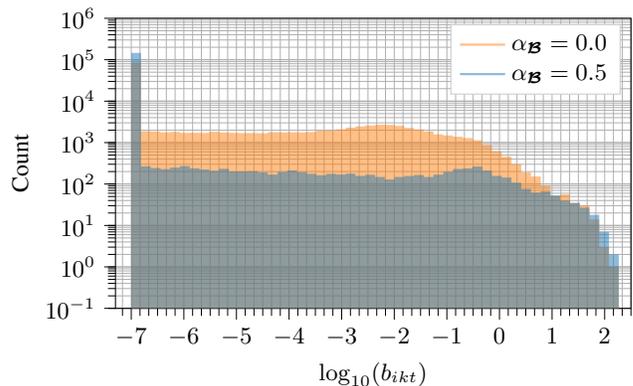


Fig. 9. Histogram of  $\mathcal{B}$  coefficients. Due to numerical issues, coefficients lower than some  $\epsilon = 10^{-7}$  were clipped to  $\epsilon$ .

In general terms, one desires sparse signatures if the data have a sparse nature; this property is usually manifested in tonal emissions, where the emitted energy is well localized in the frequency domain. In this sense, the imposed regularization is expected to enhance the precision of the algorithm for tonal

emissions, while some degradation can occur in percussive or other atonal emissions. In our test scenarios, the regularization did result in some slight degradation in separation, suggesting that the tonal assumption is not strictly valid—which is not surprising in real-life composed signals. Nevertheless,  $\ell_1$  penalty parameters are ubiquitous in regular NMF algorithms [26, 27], so that the added model flexibility is welcome.

### VIII. CONCLUSIONS

We proposed an extended version [1] of the CNMF algorithm [11], leveraging the efficient representation from the deconvolutive NMF model. We also provided the user with control over the distribution of the extracted signatures through regularization, which steers the method toward a sparse solution. Our proposed method is a generalization of the baseline algorithm, with added flexibility, able to efficiently factorize signals with complex spectral signature.

Although testing was not exhaustive, due to the large number of user-chosen parameters and different possible configurations, our simulations corroborated the method's capabilities in the separation task.

Future works include the application of these ideas using different generative models for the matrices  $\mathbf{X}_{il}$ , as the Gaussian assumption and the associated Frobenius norm method are not entirely adequate for distance measures on the PSD cone; however, the fact that, by construction,  $\mathbf{X}_{il}$  is rank-1 and lives on the boundary of  $\mathbb{S}_+$  imposes challenges on the application of traditional measures, such as log det divergences [28], which are only defined in the interior of the cone. Also, application of other NMF frameworks, such as multi-layer NMF, could improve the accuracy and consistency of the method. Proper initialization is also a critical point in traditional NMF methods [29], and the design of a smart initialization routine is likely to have a considerable effect on convergence time and accuracy.

#### APPENDIX A

##### PROOF OF PROPOSITION 1

First, the Lagrangian function for Eq. (21) w.r.t.  $\mathbf{S}_{ilkot}$  restriction is computed:

$$\psi_{il} = \sum_{k,o,t} \frac{\|\mathbf{S}_{ilkot} - z_{ko} \mathbf{W}_{io} b_{ikt} [\mathbf{g}_k]_l\|_F^2}{r_{ilkot}} - \text{tr}(\mathbf{M}_{il} (\mathbf{X}_{il} - \sum_{k,o,t} \mathbf{S}_{ilkot})^H), \quad (34)$$

where the regularization term was omitted for brevity. The partial derivative w.r.t.  $\mathbf{S}_{ilkot}$  is given by

$$\frac{\partial \psi_{il}}{\partial \mathbf{S}_{ilkot}} = \frac{2(\mathbf{S}_{ilkot} - z_{ko} \mathbf{W}_{io} b_{ikt} [\mathbf{g}_k]_l)}{r_{ilkot}} - \mathbf{M}_{il}. \quad (35)$$

Summing over  $k$ ,  $o$ , and  $t$  after equating to  $\mathbf{0}$  results in

$$\sum_{k,o,t} r_{ilkot} \mathbf{M}_{il} = 2 \left( \sum_{k,o,t} \mathbf{S}_{ilkot} - \sum_{k,o,t} z_{ko} \mathbf{W}_{io} b_{ikt} [\mathbf{g}_k]_l \right) \quad (36)$$

$$\mathbf{M}_{il} = 2(\mathbf{X}_{il} - \hat{\mathbf{X}}_{il}) = 2\mathbf{E}_{il}.$$

Finally, backsolving Eq. (35) for  $\mathbf{S}_{ilkot}$  yields

$$\mathbf{S}_{ilkot}^* = z_{ko} \mathbf{W}_{io} b_{ikt} [\mathbf{g}_k]_l - r_{ilkot} \mathbf{E}_{il}; \quad (37)$$

replacing  $\mathbf{S}_{ilkot}^*$  in Eq. (21) leads directly to the equality in Eq. (23), and the fact that Eq. (21) is quadratic in  $\mathbf{S}_{ilkot}$ , with a single global minimum, is sufficient to satisfy inequation (22).

#### APPENDIX B

##### DERIVATION OF EQUATIONS (27), (28), AND (29)

The partial derivatives of  $\ell_R^+$  w.r.t. the non-negative parameters are:

$$\frac{\partial \ell_R^+}{\partial z_{ko}} = \sum_{i,l,t} \frac{-2}{r_{ilkot}} \left[ -b_{ikt} [\mathbf{g}_k]_l \text{tr}(\mathbf{S}_{ilkot} \mathbf{W}_{io}) + z_{ko} b_{ikt}^2 [\mathbf{g}_k]_l^2 \right], \quad (38)$$

$$\frac{\partial \ell_R^+}{\partial b_{ikt}} = 2\alpha_{\mathbf{B}} + \sum_{l,o} \frac{-2}{r_{ilkot}} \left[ -z_{ko} [\mathbf{g}_k]_l \text{tr}(\mathbf{S}_{ilkot} \mathbf{W}_{io}) + z_{ko}^2 b_{ikt} [\mathbf{g}_k]_l^2 \right], \quad (39)$$

$$\frac{\partial \ell_R^+}{\partial g_{kl}} = \sum_{i,o,t} \frac{-2}{[r_{ilkot}]_l} \left[ -z_{ko} b_{ikt} \text{tr} \left( [\mathbf{S}_{ilkot}]_l \mathbf{W}_{io} \right) + z_{ko}^2 b_{ikt}^2 g_{kl} \right], \quad (40)$$

(note that the left-shift operator appears in (40)), where the fact that  $\text{tr}(\mathbf{W}_{io} \mathbf{W}_{io}^H) = 1$  was used.

Applying the definition of  $r_{ilkot}$  in terms of  $z_{ko}$ ,  $b_{ikt}$ ,  $g_{kl}$ , and  $\hat{x}_{il}$ , expressions of type  $\frac{b_{ikt} [\mathbf{g}_k]_l}{r_{ilkot}}$  are replaced with  $\frac{\hat{x}_{il}}{z_{ko}}$ , and those of type  $\frac{z_{ko} b_{ikt}^2 [\mathbf{g}_k]_l^2}{r_{ilkot}}$  with  $\hat{x}_{il} b_{ikt} [\mathbf{g}_k]_l$ ; solving for the respective variables yields:

$$z_{ko}^* = \frac{\sum_{i,l,t} \hat{x}_{il} \text{tr}(\mathbf{S}_{ilkot} \mathbf{W}_{io})}{\sum_{i,l,t} \hat{x}_{il} b_{ikt} [\mathbf{g}_k]_l}, \quad (41)$$

$$b_{ikt}^* = \frac{\sum_{l,o} \hat{x}_{il} \text{tr}(\mathbf{S}_{ilkot} \mathbf{W}_{io})}{\alpha_{\mathbf{B}} + \sum_{l,o} \hat{x}_{il} z_{ko} [\mathbf{g}_k]_l}, \quad (42)$$

$$g_{kl}^* = \frac{\sum_{i,o,t} [\hat{x}_i]_l \text{tr} \left( [\mathbf{S}_{ilkot}]_l \mathbf{W}_{io} \right)}{\sum_{i,o,t} [\hat{x}_i]_l z_{ko} b_{ikt}}. \quad (43)$$

Finally, replacing  $\mathbf{S}$  with  $\mathbf{S}^*$  leads to the multiplicative updates.

APPENDIX C  
DERIVATION OF EQUATION (30)

The gradient of  $\ell_R^+$  w.r.t.  $\mathbf{W}_{io}$  is

$$\frac{\partial \ell_R^+}{\partial \mathbf{W}_{io}} = \sum_{l,k,t} \frac{2}{r_{ilkot}} \left[ -z_{ko} b_{ikt} [\vec{\mathbf{g}}_k]_l \mathbf{S}_{ilkot} + z_{ko}^2 \mathbf{W}_{io} b_{ikt}^2 [\vec{\mathbf{g}}_k]_l^2 \right]. \quad (44)$$

Once more, applying the choice of  $r_{ilkot}$ , and solving for  $\frac{\partial \ell_R^+}{\partial \mathbf{W}_{io}} = \mathbf{0}$  yields

$$\hat{\mathbf{W}}_{io} = \frac{\sum_{l,k,t} \hat{x}_{il} \mathbf{S}_{ilkot}}{\sum_{l,k,t} \hat{x}_{il} z_{ko} b_{ikt} [\vec{\mathbf{g}}_k]_l}. \quad (45)$$

Replacing  $\mathbf{S}$  with  $\mathbf{S}^*$  leads to the multiplicative update.

REFERENCES

[1] T. L. B. Dias, L. W. P. Biscainho, and W. A. Martins, "Time-deconvolutive CNMF for multichannel blind source separation," in *Anais do XXXVII Simpósio Brasileiro de Telecomunicações e Processamento de Sinais (SBrT)*. Petrópolis, Brazil: SBrT, Sep. 2019, pp. 1–5, doi: 10.14209/sbrt.2019.1570557364.

[2] M. Pal, R. Roy, J. Basu, and M. S. Bepari, "Blind source separation: A review and analysis," in *2013 International Conference Oriental COCODA held jointly with 2013 Conference on Asian Spoken Language Research and Evaluation (O-COCODA/CASLRE)*. Gurgaon, India: IEEE, Nov. 2013, pp. 1–5, doi: 10.1109/ICSDA.2013.6709849.

[3] A. Hyvärinen, J. Karhunen, and E. Oja, "What is independent component analysis?" in *Independent Component Analysis*. John Wiley & Sons, Ltd, 2002, ch. 7, pp. 145–164, doi: 10.1002/0471221317.ch7.

[4] —, "ICA by maximization of nongaussianity," in *Independent Component Analysis*. John Wiley & Sons, Ltd, 2002, ch. 8, pp. 165–202, doi: 10.1002/0471221317.ch8.

[5] —, "ICA by maximum likelihood estimation," in *Independent Component Analysis*. John Wiley & Sons, Ltd, 2002, ch. 9, pp. 203–219, doi: 10.1002/0471221317.ch9.

[6] M. Congedo, C. Gouy-Pailler, and C. Jutten, "On the blind source separation of human electroencephalogram by approximate joint diagonalization of second order statistics." *Clinical Neurophysiology*, vol. 119, no. 12, pp. 2677–2686, Dec. 2008, doi: 10.1016/j.clinph.2008.09.007.

[7] K. Reindl, Y. Zheng, and W. Kellermann, "Speech enhancement for binaural hearing aids based on blind source separation," in *4th International Symposium on Communications, Control and Signal Processing (IS-CCSP)*. Limassol, Cyprus: IEEE, Mar. 2010, pp. 1–6, doi: 10.1109/ISCCSP.2010.5463317.

[8] 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, Oct. 1999, doi: 10.1038/44565.

[9] N. Gillis, "The why and how of nonnegative matrix factorization," in *Regularization, Optimization, Kernels, and Support Vector Machines*, ser. Machine Learning and Pattern Recognition Series, J. A. K. Suykens, M. Signoretto, and A. Argyriou, Eds. Boca Raton, USA: Chapman & Hall/CRC, 2014, ch. 12, pp. 257–291, doi: 10.5555/2700548.

[10] P. Smaragdis and J. C. Brown, "Non-negative matrix factorization for polyphonic music transcription," in *2003 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. New Paltz, USA: IEEE, Oct. 2003, pp. 177–180, doi: 10.1109/AS-PAA.2003.1285860.

[11] H. Sawada, H. Kameoka, S. Araki, and N. Ueda, "New formulations and efficient algorithms for multichannel NMF," in *2011 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. New Palz, USA: IEEE, Oct. 2011, pp. 153–156, doi: 10.1109/ASPAA.2011.6082275.

[12] J. Nikunen and T. Virtanen, "Direction of arrival based spatial covariance model for blind sound source separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 3, pp. 727–739, Mar. 2014, doi: 10.1109/TASLP.2014.2303576.

[13] J. Nikunen and T. Virtanen, "Multichannel audio separation by direction of arrival based spatial covariance model and non-negative matrix factorization," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Florence, Italy, May 2014, pp. 6677–6681, doi: 10.1109/ICASSP.2014.6854892.

[14] P. Smaragdis, "Non-negative matrix factor deconvolution; extraction of multiple sound sources from monophonic inputs," in *Independent Component Analysis and Blind Signal Separation*, C. G. Puntonet and A. Prieto, Eds. Granada, Spain: Springer, Sep. 2004, pp. 494–499, doi: 10.1007/978-3-540-30110-3\_63.

[15] A. Cichocki, R. Zdunek, A. H. Phan, and S.-I. Amari, "Introduction – problem statements and models," in *Nonnegative Matrix and Tensor Factorizations*. John Wiley & Sons, Ltd, 2009, ch. 1, pp. 1–80, doi: 10.1002/9780470747278.ch1.

[16] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006, doi: 10.5555/1162264.

[17] 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, Mar. 2010, doi: 10.1109/TASL.2010.2041381.

[18] J. Yoo and S. Choi, "Nonnegative matrix factorization with orthogonality constraints," *Journal of Computer Science and Engineering*, vol. 4, no. 2, pp. 97–109, Jun. 2010, doi: 10.5626/JCSE.2010.4.2.097.

[19] R. Tibshirani, "Regression shrinkage and selection via the LASSO," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996, doi: 10.1111/j.2517-6161.1996.tb02080.x.

- [20] J. de Leeuw, "Block-relaxation algorithms in statistics," in *Information Systems and Data Analysis*, H.-H. Bock, W. Lenski, and M. M. Richter, Eds. Berlin, Heidelberg: Springer, 1994, pp. 308–324, doi: 10.1007/978-3-642-46808-7\_28.
- [21] "CATT-Acoustic," (Visited on 10-Feb-2019). [Online]. Available: <http://catt.se>
- [22] E. Vincent, R. Gribonval, and C. Fevotte, "Performance measurement in blind audio source separation," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 4, pp. 1462–1469, Jul. 2006, doi: 10.1109/TSA.2005.858005.
- [23] E. Vincent, H. Sawada, P. Bofill, S. Makino, and J. P. Rosca, "First stereo audio source separation evaluation campaign: Data, algorithms and results," in *Independent Component Analysis and Signal Separation*, M. E. Davies, C. J. James, S. A. Abdallah, and M. D. Plumbley, Eds. London, UK: Springer, Sep. 2007, pp. 552–559, doi: 10.1007/978-3-540-74494-8\_69.
- [24] C. Raffel, B. McFee, E. J. Humphrey, J. Salamon, O. Nieto, D. Liang, D. P. Ellis, and C. C. Raffel, "mir\_eval: A transparent implementation of common MIR metrics," in *15th International Society for Music Information Retrieval Conference*. Taipei, Taiwan: ISMIR, Oct. 2014, pp. 367–372, doi: 10.1.1.722.1267.
- [25] C. Raffel, "mir\_eval documentation," [http://craffel.github.io/mir\\_eval/](http://craffel.github.io/mir_eval/), [Online; accessed February 20, 2019].
- [26] A. Cichocki, R. Zdunek, A. H. Phan, and S.-I. Amari, "Multiplicative iterative algorithms for NMF with sparsity constraints," in *Nonnegative Matrix and Tensor Factorizations*. John Wiley & Sons, Ltd, 2009, ch. 3, pp. 131–202, doi: 10.1002/9780470747278.ch3.
- [27] —, "Alternating least squares and related algorithms for NMF and SCA problems," in *Nonnegative Matrix and Tensor Factorizations*. John Wiley & Sons, Ltd, 2009, ch. 4, pp. 203–266, doi: 10.1002/9780470747278.ch4.
- [28] A. Cichocki, S. Cruces, and S. ichi Amari, "Log-determinant divergences revisited: Alpha-beta and gamma log-det divergences," *Entropy*, vol. 17, no. 5, pp. 2988–3034, May 2014, doi: 10.3390/e17052988.
- [29] C. Boutsidis and E. Gallopoulos, "SVD based initialization: A head start for nonnegative matrix factorization," *Pattern Recognition*, vol. 41, no. 4, pp. 1350–1362, Apr. 2008, doi: 10.1016/j.patcog.2007.09.010.



**Thadeu Luiz Barbosa Dias** was born in Rio de Janeiro, Brazil, in 1994. B.Sc. degree in electronics and computer engineering from the Universidade Federal do Rio de Janeiro (POLI/UFRJ) in 2020, currently pursuing his M.Sc. degree in electrical engineering at COPPE (PEE/COPPE/UFRJ). Mr. Dias's research interests include digital signal processing, image processing, machine learning and information geometry.



**Wallace Alves Martins** received the Electronics Engineer degree from Universidade Federal do Rio de Janeiro (UFRJ, Brazil) in 2007, the M.Sc. and D.Sc. degrees in Electrical Engineering also from UFRJ in 2009 and 2011, respectively. He was a Research Visitor at University of Notre Dame (USA, 2008), at Université de Lille 1 (France, 2016), and at Universidad de Alcalá (Spain, 2018). From 2010 to 2013 he was an Associate Professor of Centro Federal de Educação Tecnológica Celso Suckow da Fonseca (CEFET/RJ, Brazil). Since 2013 he has been

with the Department of Electronics and Computer Engineering (DEL/Poli) and Electrical Engineering Program (PEE/COPPE) at UFRJ, where he is presently an Associate Professor (on leave). He is currently a Research Associate working with the Interdisciplinary Centre for Security, Reliability and Trust (SnT) at Université du Luxembourg. His research interests are in the fields of digital signal processing, especially adaptive signal processing and graph signal processing, as well as digital communications, with focus on equalization and precoding for wireless communications. Dr. Martins has authored more than 60 technical papers in refereed international journals and conferences, and 1 scientific book. Also, he received the Best Student Paper Award from EURASIP at EUSIPCO-2009, Glasgow, Scotland, and the 2011 Best Brazilian D.Sc. Dissertation Award from Capes. He is currently a member of IEEE (Institute of Electrical and Electronics Engineers) and SBrT (Brazilian Telecommunications Society).



**Luiz Wagner Pereira Biscaíno** was born in Rio de Janeiro, Brazil, in 1962. He received the diploma of electronic engineer (magna cum laude) from the EE (now Poli) at Universidade Federal do Rio de Janeiro (UFRJ), Brazil, in 1985, and the M.Sc. and D.Sc. degrees in electrical engineering from the COPPE at UFRJ in 1990 and 2000, respectively. Having worked in the telecommunication industry between 1985 and 1993, Dr. Biscaíno is now Associate Professor at DEL/Poli and PEE/COPPE, at UFRJ. His research area is digital signal processing, particularly

audio processing. He is currently a member of IEEE (Institute of Electrical and Electronics Engineers), AES (Audio Engineering Society), SBrT (Brazilian Telecommunications Society), and SBC (Brazilian Computer Society).