# A Unified Bayesian View on Spatially Informed Source Separation and Extraction based on Independent Vector Analysis

Andreas Brendel, *Student Member, IEEE,* Thomas Haubner, and Walter Kellermann, *Fellow, IEEE*

*Abstract*—Signal separation and extraction are important tasks for devices recording audio signals in real environments which, aside from the desired sources, often contain several interfering sources such as background noise or concurrent speakers. Blind Source Separation (BSS) provides a powerful approach to address such problems. However, BSS algorithms typically treat all sources equally and do not resolve uncertainty regarding the ordering of the separated signals at the output of the algorithm, i.e., the outer permutation problem. This paper addresses this problem by incorporating prior knowledge into the adaptation of the demixing filters, e.g., the position of the sources, in a Bayesian framework. We focus here on methods based on Independent Vector Analysis (IVA) as it elegantly and successfully deals with the internal permutation problem. By including a background model, i.e., a model for sources we are not interested to separate, we enable the algorithm to extract the sources of interest in overdetermined and underdetermined scenarios at a low computational complexity. The proposed framework allows to incorporate prior knowledge about the demixing filters in a generic way and unifies several known and newly proposed algorithms using a Bayesian view. For all algorithmic variants, we provide efficient update rules based on the iterative projection principle. The performance of a large variety of representative algorithmic variants, including very recent algorithms, is compared using measured Room Impulse Responses (RIRs).

*Index Terms*—Source Separation, Independent Vector Analysis, ILRMA, Geometric Constraint, Independent Vector Extraction

## I. INTRODUCTION

SOURCE separation and signal extraction are essential tasks for acoustic signal processing on a variety of devices such as mobile phones, smart home assistants, hearing aids, conference systems etc. For these tasks many algorithms have been proposed in the recent years, e.g., [1], [2] which can roughly be divided into two highly overlapping groups originating from different paradigms: beamforming methods [3] and Blind Source Separation (BSS) [2], [4], [5]. In this paper, we focus on the latter one.

As a first class of BSS algorithms, we consider here algorithms which are based on Independent Component Analysis (ICA) [4], [6], and use the statistical independence of the source signals to derive algorithms capable of separating nongaussian sources. These methods are in general based on

a linear instantaneous mixing and demixing model, which makes them not directly applicable for reverberant enclosures for which the recorded signals are filtered and superimposed versions of the source signals, so that a convolutive mixture model should be applied. As a solution, it has been proposed to apply the ICA algorithm independently in different frequency bins [7]. However, due to the well-known inner permutation problem, i.e., the uncertainty about the assignment of the demixed signals to the output channels in each frequency bin, the ordering of the channels has to be recovered by repair mechanisms [8]. For avoiding the inner permutation problem, Independent Vector Analysis (IVA) [9] has been introduced, which enforces statistical dependence between the frequency bins of the demixed signals. For identifying the demixing system, stable, fast and parameter-free update rules based on the Majorize-Minimize (MM) principle have been proposed in [10].

Another class of algorithms for multichannel source separation is based on Multichannel NMF (MNMF) [11], which is an extension of Nonnegative Matrix Factorization (NMF) [12]. The main idea here is to model the source signal spectrum by a superposition of nonnegative basis vectors. This approach is especially powerful if a distinct spectral structure can be exploited, e.g., for music signals [13] or certain types of noise signals [14].

An approach which synthesizes the ideas of IVA and MNMF has been introduced as Independent Low Rank Matrix Analysis (ILRMA) [15], [16]. ILRMA can either be understood as a special case of MNMF using a rank-1 spatial model or as IVA with a time-varying Gaussian source model [17] whose variance is estimated via NMF. The benefits of this approach are its faster convergence compared to MNMF and the higher separation performance of sources with distinct spectral structure, e.g., music signals. However, if applied blindly, the permutation of the output channels remains arbitrary. Clustering based on the associated identified spatial models is difficult in a static and determined scenario, where the number of sources and sensors is equal. If the sources are moving or the scenario is underdetermined, i.e., there are more sources than sensors, such a clustering-based method is likely to fail.

For signal extraction, a Background (BG) model has been proposed in [18] which leads to the Independent Vector Extraction (IVE) algorithm. Here, one desired source is separated from a set of other sources forming the BG, for which no effort is spent to separate them. The same model has been

used in [19] to derive an MM-based optimization scheme for IVA in overdetermined scenarios. In both cases it is argued that the coupling of the Sources Of Interest (SOI) and the BG is only weakly expressed in the cost function, i.e., the cost function consists of a part only depending on the SOI filters and another part only depending on the BG filters. As a remedy, an orthogonality constraint is imposed on the demixing filters corresponding to SOIs and BG, which yields the update rules for the BG filters. For the selection of the SOI filters, a directional constraint and a supervised adaptation based on a reference signal [20] has been suggested in [21] for IVE. For [19] no such selection strategy exists so far.

Many ways have been proposed to incorporate spatial prior knowledge about the sources into the adaptation of the demixing filters of BSS algorithms to speed up convergence or to ensure the extraction of a desired source [22]. A geometric constraint has also been used in TRIple-N Independent component analysis for CONvolutive mixtures (TRINICON)-based signal extraction [23], [24], [25] and for IVA in [26]. An optimization algorithm for spatially regularized ILRMA based on vector-wise coordinate descent has recently been proposed in [27].

Besides geometric constraints, [28] proposed to use spatial models for the reverberant component of the observed sound signals together with free-field models to obtain a full-rank spatial covariance model. In [29], previously obtained demixing filters are introduced as prior knowledge into BSS.

In this paper, we propose a novel generic Bayesian framework for informed source separation based on IVA. This framework allows to incorporate prior knowledge on the demixing matrices in a generic way and provides fast converging Iterative Projection (IP)-based update rules at a low computational complexity at the same time. Various known and novel algorithmic variants are identified as special cases of the generic framework. Several strategies for incorporating prior knowledge in the Bayesian sense are discussed and exemplified by priors based on a free-field model, which allows to steer spatial ones and nulls. A BG model is introduced, which can also incorporate priors and allows for a significant reduction of computational cost. For the SOIs, several source models are discussed including NMF and fast and stable update rules for all algorithmic variants based on the MM principle are proposed. A new perspective is taken in the derivation of the update rules for the BG filters based on IP. The proposed framework allows the solution of the outer permutation problem of BSS as well as signal extraction and separation in determined and overdetermined scenarios and signal extraction in underdetermined scenarios. This paper is an extension of [30], where we discussed a very specific realization of the generic Bayesian framework presented here.

In the following, scalar variables are typeset as lower-case letters, vectors as bold lower-case letters, matrices as bold upper-case letters and sets as calligraphic upper-case letters. $\mathbf{I}_d$ and $\mathbf{0}_d$ denote a quadratic identity or all-zero matrix, respectively, of dimensions $d \times d$, and $\mathbf{0}_{d_1 \times d_2}$ denotes an all-zero matrix of dimensions $d_1 \times d_2$. $(\cdot)^{\mathrm{H}}$ and $(\cdot)^{\mathrm{T}}$ denote a Hermitian (complex conjugate transpose) and transposed matrix, respectively. Complex-conjugated quantities are marked

| | |
|---|---|
| $\mathbf{I}, \mathbf{0}$ | Identity and all-zero matrix |
| $f, F$ | Frequency bin index and number of frequency bins |
| $k, K$ | Channel index and number of channels |
| $l, L$ | Iteration index and number of iterations |
| $m, M$ | Microphone index and number of microphones |
| $n, N$ | Time block index and number of blocks |
| $\nu, N_{\mathrm{bases}}$ | Basis index and number of bases |
| $\mathbf{W}, \mathbf{w}$ | Demixing matrix and demixing vector |
| $\mathbf{P}$ | Precision matrix of spatial prior |
| $J$ | Cost function |
| $\mathbf{A}, \mathbf{a}$ | Mixing matrix and mixing vector |
| $t, v$ | Basis element and activation of NMF |
| $\mathbf{C}$ | Microphone covariance matrix of frequency bin $f$ |
| $\mathbf{z}$ | BG signal vector |
| $\mathbf{B}, \mathbf{b}$ | BG filter matrix and vector |
| $Q$ | Number of sources |
| $r$ | Estimated demixed signal variance |
| $U(\cdot|\cdot)$ | Upper bound |
| $\mathbf{V}$ | Weighted microphone covariance matrix |
| $\mathbf{x}$ | Microphone signal vector |
| $\mathbf{s}$ | SOI signal vector |
| $\mathbf{q}$ | Source signal vector |
| $\mathbf{y}$ | Demixed signal vector |
| $\mathbf{B}^{M,K}$ | BG filter submatrix |
| $\underline{\mathbf{y}}$ | Broadband demixed signal vector |
| $\underline{\mathbf{z}}$ | Broadband BG signal vector |
| $\mathcal{Y}, \mathcal{X}$ | Set of demixed signals and microphone signals |
| $\mathcal{W}$ | Set of demixing matrices |
| $[N], [F]$ | Index set of time blocks and frequency bins |
| $\mathbf{h}_f(\vartheta)$ | Free-field steering vector for direction $\vartheta$ |

TABLE I
NOTATIONS USED

by $(\cdot)^*$ and the derivative of a function w.r.t. its argument is denoted by $(\cdot)'$. The set $\{1, 2, \ldots, N\}$ is denoted by $[N]$. The notation of important variables is given in Tab. I for later reference.

The remainder of the paper is structured as follows: Sec. II defines the signal model, the probabilistic model for the SOIs and the BG and introduces prior Probability Density Functions (PDFs). The fundamental principle of MM algorithms is described in Sec. III. In the same section, an upper bound for the previously derived cost function is constructed and optimized, and update rules for the demixing filters based on the iterative projection principle are proposed. Experimental results are presented in Sec. IV. The paper is concluded in Sec. V.

## II. MODELS

The following section introduces the underlying source models for SOIs and BG signals, the probabilistic model for the demixing system including prior PDFs which allow to incorporate prior knowledge about the demixing filters.

### A. Signal Model

We consider an acoustic scene in an enclosure comprising $M$ microphones and $Q$ simultaneously active acoustic point sources observed by the microphones as a convolutive mixture. In this contribution, we are interested in separating $K \leq Q$ SOIs out of the observed mixture of $Q$ sources. The remaining sources, if there are any, are associated with the so-called Background (BG) in the following.

With $f \in [F]$ denoting the frequency bin index and $n \in [N]$ the discrete time index, we assume a linear time-invariant mixing model in the Short-Time Fourier Transform (STFT) domain

$$\mathbf{x}_{f,n} = \mathbf{A}_f \mathbf{q}_{f,n}, \tag{1}$$

with the source signal vector

$$\mathbf{q}_{f,n} = [q_{1,f,n}, \ldots, q_{Q,f,n}]^{\mathrm{T}} \in \mathbb{C}^Q, \tag{2}$$

the microphone signal vector

$$\mathbf{x}_{f,n} = [x_{1,f,n}, \ldots, x_{M,f,n}]^{\mathrm{T}} \in \mathbb{C}^M \tag{3}$$

and the mixing matrix containing the acoustic transfer functions at frequency bin $f$ from the source positions to the microphones

$$\mathbf{A}_f \in \mathbb{C}^{M \times Q}. \tag{4}$$

Note that the number of sources $Q$, the number of microphones $M$ and the number of SOIs $K$ can be different in general.

In the following, the demixing model is introduced as illustrated in Fig. 1. The SOIs and the BG signals are obtained by

$$\mathbf{y}_{f,n} = \mathbf{W}_f \mathbf{x}_{f,n} \tag{5}$$

where the demixing matrix applied in frequency bin $f$

$$\mathbf{W}_f = \begin{bmatrix} \mathbf{W}_f^{\mathrm{SOI}} \\ \mathbf{B}_f \end{bmatrix} \in \mathbb{C}^{M \times M} \tag{6}$$

contains two parts: One set of filters extracting the SOIs $\mathbf{s}_{f,n}$

$$\mathbf{W}_f^{\mathrm{SOI}} = \left[\mathbf{w}_f^1, \ldots, \mathbf{w}_f^K\right]^{\mathrm{H}} \in \mathbb{C}^{K \times M}, \tag{7}$$

and another set of filters

$$\mathbf{B}_f = \left[\mathbf{b}_f^1, \ldots, \mathbf{b}_f^{M-K}\right]^{\mathrm{H}} = \left[\mathbf{B}_f^{M,K} \quad -\mathbf{I}_{M-K}\right] \in \mathbb{C}^{M-K \times M} \tag{8}$$

estimating the BG signals $\mathbf{z}_{f,n}$. Note that $\mathbf{B}_f$ is structured according to the model proposed in [18] with the identity matrix $\mathbf{I}_{M-K}$ and a submatrix $\mathbf{B}_f^{M,K}$ capturing the free parameters of $\mathbf{B}_f$, which have to be identified together with the SOI filters $\mathbf{W}_f^{\mathrm{SOI}}$. For a given time frame $n$ and frequency bin $f$, the vector of output signals $\mathbf{y}_{f,n} = \left[\mathbf{s}_{f,n}^{\mathrm{T}}, \mathbf{z}_{f,n}^{\mathrm{T}}\right]^{\mathrm{T}}$ contains the vector of demixed SOIs denoted as

$$\mathbf{s}_{f,n} = \mathbf{W}_f^{\mathrm{SOI}} \mathbf{x}_{f,n} = [s_{1,f,n}, \ldots, s_{K,f,n}]^{\mathrm{T}} \in \mathbb{C}^K, \tag{9}$$

and the vector of BG signals denoted as

$$\mathbf{z}_{f,n} = \mathbf{B}_f \mathbf{x}_{f,n} = [z_{1,f,n}, \ldots, z_{M-K,f,n}]^{\mathrm{T}} \in \mathbb{C}^{M-K}. \tag{10}$$

Note that only if $K < M$ holds, BG signals can be extracted by the assumed $M \times M$ demixing matrix $\mathbf{W}_f$.

For the determined case, i.e., $K = M$, no BG signals are estimated and the demixing matrix separates only the SOIs $\mathbf{W}_f = \mathbf{W}_f^{\mathrm{SOI}}$. Furthermore, we define the broadband signal vector of the $k$th SOI and BG signal at time frame $n$

$$\underline{\mathbf{s}}_{k,n} = [s_{k,1,n}, \ldots, s_{k,F,n}]^{\mathrm{T}}, \quad \underline{\mathbf{z}}_{k,n} = [z_{k,1,n}, \ldots, z_{k,F,n}]^{\mathrm{T}} \in \mathbb{C}^F.$$

With the definitions

$$\underline{\mathbf{s}}_n = \left[\underline{\mathbf{s}}_{1,n}^{\mathrm{T}}, \ldots, \underline{\mathbf{s}}_{K,n}^{\mathrm{T}}\right]^{\mathrm{T}} \in \mathbb{C}^{KF} \tag{11}$$
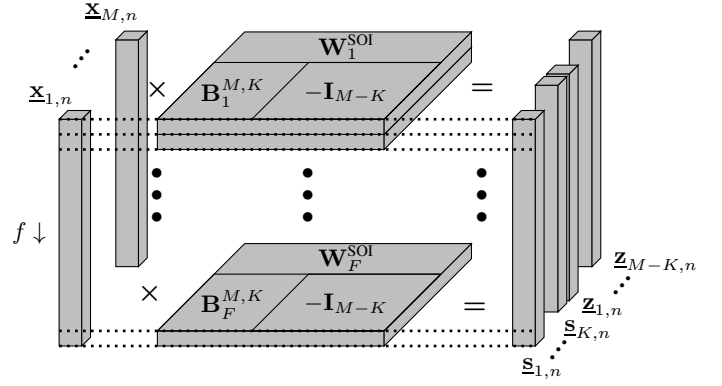


Fig. 1. Illustration of the demixing process. The demixing matrix $\mathbf{W}_f$ is applied in each frequency bin separately to the broadband vectors of microphone signals $\underline{\mathbf{x}}_{k,n}$, $k \in [M]$. The results are the extracted SOIs $\underline{\mathbf{s}}_{k,n}$, $k \in [K]$, and the BG signals $\underline{\mathbf{z}}_{k,n}$, $k \in [M-K]$.

and

$$\underline{\mathbf{z}}_n = \left[\underline{\mathbf{z}}_{1,n}^{\mathrm{T}}, \ldots, \underline{\mathbf{z}}_{M-K,n}^{\mathrm{T}}\right]^{\mathrm{T}} \in \mathbb{C}^{(M-K)F} \tag{12}$$

we can write the signal vector containing all output signals as

$$\underline{\mathbf{y}}_n = \left[\underline{\mathbf{s}}_n^{\mathrm{T}}, \underline{\mathbf{z}}_n^{\mathrm{T}}\right]^{\mathrm{T}} \in \mathbb{C}^{MF}. \tag{13}$$

Note that for the determined case, i.e., $M = K$, $\underline{\mathbf{y}}_n = \underline{\mathbf{s}}_n$ holds.

### B. Probabilistic Model of the Demixing System

For treating the identification of the demixing matrix as a Bayesian estimation problem, we derive the posterior density of the demixing matrices in the following. Before starting the derivation we define the set of all demixing matrices $\mathcal{W} = \left\{\mathbf{W}_f \in \mathbb{C}^{M \times M} | f \in [F]\right\}$, the set of all demixed signal vectors $\mathcal{Y} = \left\{\underline{\mathbf{y}}_n \in \mathbb{C}^{MF} | n \in [N]\right\}$ and the set of all microphone observations $\mathcal{X} = \left\{\mathbf{x}_{f,n} \in \mathbb{C}^M | f \in [F], n \in [N]\right\}$.

Using these definitions, the joint posterior of demixing matrices $\mathcal{W}$ and demixed signals $\mathcal{Y}$ can be written as

$$\begin{aligned} p(\mathcal{W}, \mathcal{Y}|\mathcal{X}) &= p(\mathcal{W}, \mathcal{Y})\frac{p(\mathcal{X}|\mathcal{W}, \mathcal{Y})}{p(\mathcal{X})} \\ &\propto p(\mathcal{W})p(\mathcal{Y}|\mathcal{W})p(\mathcal{X}|\mathcal{W}, \mathcal{Y}). \end{aligned} \tag{14}$$

We choose the following likelihood function for frequency bin $f$ and time step $n$, under the assumption that $\mathbf{W}_f$ is invertible

$$p\left(\mathbf{x}_{f,n}|\mathcal{W}, \mathbf{y}_{f,n}\right) = \delta\left(\mathbf{x}_{f,n} - \mathbf{W}_f^{-1}\mathbf{y}_{f,n}\right), \tag{15}$$

where $\delta(\cdot)$ denotes the Dirac distribution. From (15) a simplistic likelihood for all frequency bins $f \in [F]$ and time steps $n \in [N]$ can be constructed by using an i.i.d. assumption

$$p(\mathcal{X}|\mathcal{W}, \mathcal{Y}) = \prod_{n=1}^{N} \prod_{f=1}^{F} \delta\left(\mathbf{x}_{f,n} - \mathbf{W}_f^{-1}\mathbf{y}_{f,n}\right). \tag{16}$$

Moreover, a simplistic probabilistic model for the sources can be formulated under the assumption of independence between all time frames as

$$p(\mathcal{Y}|\mathcal{W}) = \prod_{n=1}^{N} p\left(\underline{\mathbf{y}}_n\right) = \prod_{n=1}^{N} p\left(\underline{\mathbf{z}}_n\right) \prod_{k=1}^{K} p\left(\underline{\mathbf{s}}_{k,n}\right), \tag{17}$$

where in the rightmost term the realistic assumption of mutual statistical independence of the SOIs and the independence of the SOIs from the BG sources is included. Note that $p\left(\underline{\mathbf{z}}_n\right)$ and $p(\underline{\mathbf{s}}_{k,n})$ are multivariate PDFs capturing all frequency bins. Now, the posterior of the demixing matrices is computed by marginalizing the demixed signals $\mathcal{Y}$ out of the joint posterior (14)

$$p(\mathcal{W}|\mathcal{X}) \propto p(\mathcal{W}) \int p(\mathcal{Y}|\mathcal{W})p(\mathcal{X}|\mathcal{W}, \mathcal{Y})d\underline{\mathbf{y}}_1 \ldots d\underline{\mathbf{y}}_N. \quad (18)$$

Inserting the models (16) and (17) yields

$$p(\mathcal{W}|\mathcal{X}) \propto p(\mathcal{W}) \prod_{n=1}^{N} \int p(\underline{\mathbf{y}}_n) \prod_{f=1}^{F} \delta\left(\mathbf{x}_{f,n} - \mathbf{W}_f^{-1}\mathbf{y}_{f,n}\right) d\underline{\mathbf{y}}_n.$$

Applying the rules for a linear transform of complex random variables [31] to the transform $\underline{\mathbf{y}}_{f,n} = \mathbf{W}_f\underline{\mathbf{x}}_{f,n}$ and using the sifting property of the Dirac distribution yields finally

$$p(\mathcal{W}|\mathcal{X}) \propto p(\mathcal{W}) \prod_{f=1}^{F} |\det \mathbf{W}_f|^{2N} \prod_{n=1}^{N} p\left(\underline{\mathbf{z}}_n\right) \prod_{k=1}^{K} p(\underline{\mathbf{s}}_{k,n}). \quad (19)$$

Optimizing the posterior for the demixing matrices considering the logarithm of (19) yields the following Maximum A Posteriori (MAP) problem

$$\mathcal{W} = \underset{\mathcal{W}}{\arg\max} \, \frac{\log p(\mathcal{W})}{N} + 2\sum_{f=1}^{F} \log|\det \mathbf{W}_f| \ldots$$

$$\ldots - \sum_{k=1}^{K} \hat{\mathbb{E}}\left\{G\left(\underline{\mathbf{s}}_{k,n}\right)\right\} + \hat{\mathbb{E}}\left\{\log p\left(\underline{\mathbf{z}}_n\right)\right\}. \quad (20)$$

Here, we introduced the score function $G(\underline{\mathbf{s}}_{k,n}) = -\log p(\underline{\mathbf{s}}_{k,n})$ and the averaging operator $\hat{\mathbb{E}}\left\{\cdot\right\} = \frac{1}{N}\sum_{n=1}^{N}(\cdot)$ for a concise notation.

### C. Models for SOIs

In the following, we want to introduce various widely-used models $p(\underline{\mathbf{s}}_{k,n})$ for the SOIs.

*1) Super-Gaussian PDF:* A popular and flexible source model for IVA, containing many others as a special case, is the generalized Gaussian distribution [32]

$$p\left(\underline{\mathbf{s}}_{k,n}\right) \propto \exp\left(-\|\underline{\mathbf{s}}_{k,n}\|_2^\beta\right), \quad (21)$$

where $\beta \in \mathbb{R}_+$ the shape parameter and $\|\cdot\|_2$ the Euclidean norm. The corresponding score function is given as (discarding constant terms)

$$G(\underline{\mathbf{s}}_{k,n}) = \|\underline{\mathbf{s}}_{k,n}\|_2^\beta. \quad (22)$$

*2) Time-varying Gaussian PDF:* A Gaussian PDF with time-varying broadband signal variance $\sigma_{k,n}^2$ [32]

$$p\left(\underline{\mathbf{s}}_{k,n}\right) \propto \exp\left(-\frac{\|\underline{\mathbf{s}}_{k,n}\|_2^2}{\sigma_{k,n}^2}\right), \quad (23)$$

is another popular choice, where the corresponding score function is given as (discarding constant terms)

$$G(\underline{\mathbf{s}}_{k,n}) = \frac{\|\underline{\mathbf{s}}_{k,n}\|_2^2}{\sigma_{k,n}^2}. \quad (24)$$

*3) Nonnegative Matrix Factorization:* If the source signal spectrum is structured, e.g., for music signals, or if prior knowledge about the source spectrum is available, an NMF-based source model is promising. Hereby, independence over all frequency bins is assumed [15]

$$p\left(\underline{\mathbf{s}}_{k,n}\right) = \prod_{f=1}^{F} \mathcal{N}^C\left(s_{k,f,n}|0, \sigma_{k,f,n}^2\right) \quad (25)$$

where the circularly-symmetric complex Gaussian distribution

$$\mathcal{N}^C\left(s_{k,f,n}|0, \sigma_{k,f,n}^2\right) = \frac{1}{\pi\sigma_{k,f,n}^2}\exp\left(-\frac{|s_{k,f,n}|^2}{\sigma_{k,f,n}^2}\right) \quad (26)$$

for each time-frequency bin has been chosen [16]. The frequency bin-wise signal variance $\sigma_{k,f,n}^2 = \mathbb{E}\{|s_{k,f,n}|^2\}$ is modeled as

$$\hat{\sigma}_{k,f,n}^2 = \left(\sum_{\nu=1}^{N_{\text{bases}}} t_{k,f,\nu}v_{k,\nu,n}\right)^\beta, \quad (27)$$

where $\beta \in \mathbb{R}_+$ is a user-defined parameter. Hereby, $\nu \in [N_{\text{bases}}]$ indexes the basis vectors, $t_{k,f,\nu}$ denotes the element of the $\nu$th basis vector corresponding to frequency bin $f$ and source $k$ and the associated activation at time instant $n$ is denoted by $v_{k,\nu,n}$. The resulting score function reads (discarding constant terms)

$$G\left(\underline{\mathbf{s}}_n\right) = \sum_{f=1}^{F}\sum_{k=1}^{K}\left(\log\sigma_{k,f,n}^2 + \frac{|s_{k,f,n}|^2}{\sigma_{k,f,n}^2}\right). \quad (28)$$

An in-depth discussion of different source models for ILRMA, where NMF source models are commonly used, can be found in [33].

### D. Background Model

We model the BG signals, collected in set $\mathcal{Z} = \left\{\mathbf{z}_{f,n} \in \mathbb{C}^M | f \in [F], n \in [N]\right\}$, to be independent over all frequency bins and time steps for simplicity

$$p(\mathcal{Z}) = \prod_{n=1}^{N} p\left(\underline{\mathbf{z}}_n\right) = \prod_{n=1}^{N}\prod_{f=1}^{F} p\left(\mathbf{z}_{f,n}\right). \quad (29)$$

Furthermore, we model the BG signals at each time-frequency bin to be multivariate complex Gaussian distributed

$$p\left(\mathbf{z}_{f,n}\right) = \frac{1}{\pi^{M-K}|\det\mathbf{R}_f|}\exp\left(-\mathbf{z}_{f,n}^H\mathbf{R}_f^{-1}\mathbf{z}_{f,n}\right), \quad (30)$$

where $\mathbf{R}_f$ denotes its covariance matrix. Note that we do not aim at separating the BG signals and neither aim at estimating their covariance matrix. Note that (30) puts no restrictions on the BG model except for Gaussianity, so that, e.g., spatially white noise as well as spatially correlated sound fields, notably diffuse sound fields, are captured.

To simplify the derivation of the update algorithms for the BG filters, we use an eigenvalue decomposition of the BG signal covariance matrix

$$\mathbf{T}_f^H\mathbf{R}_f^{-1}\mathbf{T}_f = \mathbf{\Lambda}_f. \quad (31)$$

Hereby, $\mathbf{T}_f \in \mathbb{C}^{(M-K)\times(M-K)}$ denotes an orthonormal matrix (i.e., $\mathbf{T}_f\mathbf{T}_f^\mathrm{H} = \mathbf{I}_{M-K}$) containing the eigenvectors of $\mathbf{R}_f$ and $\mathbf{\Lambda}_f$ denotes a diagonal matrix containing its eigenvalues. Note that such a decomposition always exists for covariance matrices. As all eigenvalues are real-valued and positive, $\mathbf{\Lambda}_f$ can be decomposed as

$$\mathbf{\Lambda}_f = \mathbf{D}_f\mathbf{D}_f, \tag{32}$$

where $\mathbf{D}_f \in \mathbb{R}^{(M-K)\times(M-K)}$ denotes the matrix square root of $\mathbf{\Lambda}_f$. Note that the entries of $\mathbf{D}_f$ are again all real-valued and positive, hence, $\mathbf{D}_f$ is invertible.

Using the relations (31) and (32), the covariance matrix $\mathbf{R}_f$ can be transformed into an identity matrix

$$\mathbf{D}_f^{-1}\mathbf{T}_f^\mathrm{H}\mathbf{R}_f^{-1}\mathbf{T}_f\mathbf{D}_f^{-1} = \mathbf{I}_{M-K}. \tag{33}$$

By using (33), we obtain

$$p(\mathbf{z}_{f,n}) = \frac{1}{\pi^{M-K}|\det \mathbf{R}_f|} \exp\left(-\tilde{\mathbf{z}}_{f,n}^\mathrm{H}\tilde{\mathbf{z}}_{f,n}\right), \tag{34}$$

with

$$\tilde{\mathbf{z}}_{f,n} = \mathbf{D}_f\mathbf{T}_f^\mathrm{H}\mathbf{z}_{f,n} = \mathbf{D}_f\mathbf{T}_f^\mathrm{H}\mathbf{B}_f\mathbf{x}_{f,n} = \tilde{\mathbf{B}}_f\mathbf{x}_{f,n}. \tag{35}$$

Here, we defined $\tilde{\mathbf{B}}_f = \mathbf{D}_f\mathbf{T}_f^\mathrm{H}\mathbf{B}_f$. Taking the i.i.d. assumption (29) w.r.t. time and frequency of the BG signals into account, the PDF of all BG signals $\mathcal{Z}$ is obtained as

$$p(\mathcal{Z}) \propto \exp\left(-\sum_{f=1}^F \sum_{n=1}^N \tilde{\mathbf{z}}_{f,n}^\mathrm{H}\tilde{\mathbf{z}}_{f,n}\right) \tag{36}$$

$$= \exp\left(-\sum_{f=1}^F \sum_{n=1}^N \sum_{k=1}^{M-K} (\tilde{\mathbf{b}}_f^k)^\mathrm{H}\mathbf{x}_{f,n}\mathbf{x}_{f,n}^\mathrm{H}\tilde{\mathbf{b}}_f^k\right) \tag{37}$$

$$= \exp\left(-N\sum_{f=1}^F \sum_{k=1}^{M-K} (\tilde{\mathbf{b}}_f^k)^\mathrm{H}\mathbf{C}_f\tilde{\mathbf{b}}_f^k\right). \tag{38}$$

Hereby, $\tilde{\mathbf{b}}_f^k$ denote the modified BG filter vectors, defined analogously to (8) and $\mathbf{C}_f = \hat{\mathbb{E}}\left\{\mathbf{x}_{f,n}\mathbf{x}_{f,n}^\mathrm{H}\right\}$ the microphone signal covariance matrix. Hence, we obtain the following term contributing to the cost function (neglecting constant terms)

$$\log p(\mathcal{Z}) = -N\sum_{f=1}^F \sum_{k=1}^{M-K} (\tilde{\mathbf{b}}_f^k)^\mathrm{H}\mathbf{C}_f\tilde{\mathbf{b}}_f^k = -N J_{\mathrm{BG}}(\mathcal{W}). \tag{39}$$

### E. Priors

The prior of the demixing matrices is chosen to be the product of marginal PDFs for each SOI filter $\mathbf{w}_f^k$, the BG filter matrix $\mathbf{B}_f$ and frequency bin $f$

$$p(\mathcal{W}) = \prod_{f=1}^F p(\mathbf{W}_f) = \prod_{f=1}^F p(\mathbf{B}_f)\prod_{k\in\mathcal{I}} p(\mathbf{w}_f^k). \tag{40}$$

In the following, we will discuss separately the priors for the SOI and the BG filters and will give the overall term contributing to the cost function.

*1) SOIs:* In many cases no prior knowledge is available for some of the channels or the optimization of the corresponding demixing filters should not be constrained. Hence, we only incorporate prior knowledge for a subset $\mathcal{I} \subseteq [K]$ of the demixing filters of the SOIs and choose uninformative priors for $k \notin \mathcal{I}$. In the following, we will present two different priors for the SOI filters based on Gaussian PDFs.

The first option for a prior for the $k$-th channel is chosen to be a zero-mean complex multivariate Gaussian PDF with precision matrix $\mathbf{P}_f^k$ and weighting factor $\tilde{\gamma}_{k,f}$

$$p(\mathbf{w}_f^k) = \frac{\sqrt{(\tilde{\gamma}_{k,f})^M \det \mathbf{P}_f^k}}{\sqrt{\pi^M}} \exp\left(-\tilde{\gamma}_{k,f}(\mathbf{w}_f^k)^\mathrm{H}\mathbf{P}_f^k\mathbf{w}_f^k\right). \tag{41}$$

The weighting factor $\tilde{\gamma}_{k,f}$ controls here and similarly for the other priors the impact of the prior on the overall model, i.e., it is a user-defined parameter. In the following, we want to discuss different choices for $\mathbf{P}_f^k$ yielding different priors for the demixing filters. To construct these priors, we use a free-field model and define the steering vector as

$$[\mathbf{h}_f(\vartheta_i)]_m = \left[\exp\left(j\frac{2\pi\mu_f}{c_s}\|\mathbf{r}_m - \mathbf{r}_1\|_2\cos\vartheta_i\right)\right]_m, \tag{42}$$

where $\mathbf{r}_m$ denotes the position of the $m$th microphone, $\mu_f$ the frequency in Hz corresponding to frequency bin $f$, $\vartheta_i$ the direction of the source and $c_s$ the speed of sound. Using this definition, we define the precision matrix yielding a spatial null

$$\mathbf{P}_{f,\mathrm{Null}}^k = \lambda_{\mathrm{Tik}}^{\mathrm{Null}}\mathbf{I}_M + \sum_{i:\vartheta_i\in\Theta_k} \lambda_i^{\mathrm{Null}}\mathbf{h}_f(\vartheta_i)\mathbf{h}_f(\vartheta_i)^\mathrm{H}, \tag{43}$$

where $\Theta_k$ denotes the set of constrained Direction of Arrivals (DOAs) and $\lambda_i^{\mathrm{Null}}$ is a weight defining the influence of the constraint in direction $\vartheta_i$, while $\lambda_{\mathrm{Tik}}^{\mathrm{Null}}$ controls the penalty on the filters energy. The intuition behind this choice can be understood if the argument of (41) is rearranged

$$(\mathbf{w}_f^k)^\mathrm{H}\mathbf{P}_{f,\mathrm{Null}}^k\mathbf{w}_f^k = \cdots \tag{44}$$
$$\cdots = \lambda_{\mathrm{Tik}}^{\mathrm{Null}}\|\mathbf{w}_f^k\|_2^2 + \sum_{i:\vartheta_i\in\Theta_k} \lambda_i^{\mathrm{Null}}\|\mathbf{h}_f(\vartheta_i)^\mathrm{H}\mathbf{w}_f^k\|_2^2.$$

The first term represents the filters power and can be seen as a Tikhonov regularizer. The second term gives the length of the projection of the filters $\mathbf{w}_f^k$ onto the steering vectors $\mathbf{h}_f(\vartheta_i)$. Hence, this prior favors solutions with small filter energy and good angular alignment to the steering vectors $\mathbf{h}_f(\vartheta_i)$. Similarly, the precision matrix yielding a spatial one is given as

$$\mathbf{P}_{f,\mathrm{One}}^k = \lambda_{\mathrm{Tik}}^{\mathrm{One}}\mathbf{I}_M - \sum_{i:\vartheta_i\in\Theta_k} \lambda_i^{\mathrm{One}}\mathbf{h}_f(\vartheta_i)\mathbf{h}_f(\vartheta_i)^\mathrm{H}, \tag{45}$$

where $\lambda_i^{\mathrm{One}}$ and $\lambda_{\mathrm{Tik}}^{\mathrm{One}}$ are weighting parameters.

As an alternative to (41), we present another prior for the channels $k \in \mathcal{I}^{\mathrm{Euc}}$ based on the Euclidean distance between the current filter estimate and the target filter vector

$$p(\mathbf{w}_f^k) = \frac{\sqrt{(\tilde{\gamma}_{k,f}^{\mathrm{Euc}})^M}}{\sqrt{\pi^M}} \cdots \tag{46}$$
$$\cdots \exp\left(-\tilde{\gamma}_{k,f}^{\mathrm{Euc}}(\mathbf{w}_f^k - \mathbf{h}_f(\vartheta_k))^\mathrm{H}(\mathbf{w}_f^k - \mathbf{h}_f(\vartheta_k))\right).$$
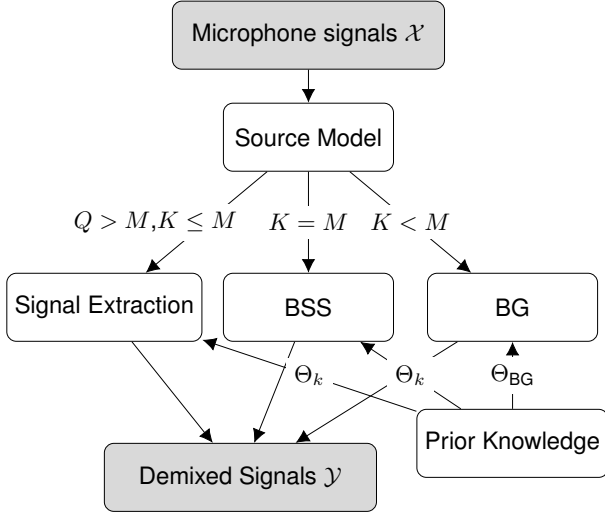
Fig. 2. Relation of proposed algorithmic variants. Depending on $Q$, $K$ and $M$, different algorithmic variants can be chosen: determined source separation, signal extraction or overdetermined BSS using a BG model.

Hereby, we used the the steering vector $\mathbf{h}_f(\vartheta_k)$ defined in (42).

In this contribution, we discuss practical realizations of the priors on the demixing vectors in the form of spatial priors which will also be the main focus in this paper. However, it should be noted that the proposed framework can be used for any prior which can be represented in the form of (41) or (46). Note that (43) and (46) have been first introduced in [27] and [30], respectively.

*2) Background:* Analogously to the priors for the SOIs (41), we choose the prior for the transformed BG filters to be

$$
p\left(\tilde{\mathbf{B}}_f\right) = \left(\frac{\sqrt{(\tilde{\gamma}_f^{\mathrm{BG}})^M \det \mathbf{P}_f^{\mathrm{BG}}}}{\sqrt{\pi^M}}\right)^{M-K} \cdots
$$
$$
\cdots \exp\left(-\tilde{\gamma}_f^{\mathrm{BG}} \sum_{k=1}^{M-K} (\tilde{\mathbf{b}}_f^k)^{\mathrm{H}} \mathbf{P}_f^{\mathrm{BG}} \tilde{\mathbf{b}}_f^k\right), \qquad (47)
$$

where we assumed independence between all channels and impose the same constraint by choosing $\mathbf{P}_f^{\mathrm{BG}}$ according to (43) for all BG channels. Note that the independence assumption applies here to the filters, not to the BG signals. This can be justified by considering filters associated with independent source positions to be independent as well. The constrained directions for the BG are collected in the set $\Theta_{\mathrm{BG}}$. Thereby, one or multiple spatial nulls can be controlled, e.g., to avoid the occurrence of the SOIs in the BG.

*3) Overall Prior:* Joining the priors for SOIs and BG yields the overall log prior term (neglecting constant terms) (cf. (40))

$$
\log p(\mathcal{W}) = -N \sum_{f=1}^{F} \left( \gamma_f^{\mathrm{BG}} \sum_{k'=1}^{M-K} (\tilde{\mathbf{b}}_f^{k'})^{\mathrm{H}} \mathbf{P}_f^{\mathrm{BG}} \tilde{\mathbf{b}}_f^{k'} \cdots \qquad (48)
$$
$$
\cdots + \sum_{k \in \mathcal{I}} \gamma_{k,f} (\mathbf{w}_f^k)^{\mathrm{H}} \mathbf{P}_f^k \mathbf{w}_f^k + \sum_{k \in \mathcal{I}^{\mathrm{Euc}}} \tilde{\gamma}_{k,f}^{\mathrm{Euc}} \|\mathbf{w}_f^k - \mathbf{h}_f(\vartheta_k)\|_2^2 \right),
$$

where we introduced the notation $\gamma_f^{\mathrm{BG}} = \frac{\tilde{\gamma}_f^{\mathrm{BG}}}{N}$, $\gamma_{k,f} = \frac{\tilde{\gamma}_{k,f}}{N}$ and $\gamma_{k,f}^{\mathrm{Euc}} = \frac{\tilde{\gamma}_{k,f}^{\mathrm{Euc}}}{N}$ for convenience in the following. The term contributing to the cost function is given by

$$
N J_{\mathrm{prior}}(\mathcal{W}) = -\log p(\mathcal{W}). \qquad (49)
$$

*F. Generic Cost Function*

Taking the negative of the MAP problem (20) and using (39) and (48) yields the generic cost function

$$
J_{\mathrm{IBSS}}(\mathcal{W}) = \underbrace{\sum_{k=1}^{K} \hat{\mathbb{E}}\left\{G\left(\underline{\mathbf{s}}_{k,n}\right)\right\} - 2\sum_{f=1}^{F} \log |\det \mathbf{W}_f| \cdots}_{J_{\mathrm{BSS}}(\mathcal{W})}
$$
$$
\cdots + J_{\mathrm{BG}}(\mathcal{W}) + J_{\mathrm{prior}}(\mathcal{W}). \qquad (50)
$$

The cost function $J_{\mathrm{IBSS}}$ consists of three parts: The BSS cost function $J_{\mathrm{BSS}}$, a component corresponding to the BG $J_{\mathrm{BG}}$ and a term representing the priors $J_{\mathrm{prior}}$ of SOIs and BG. Fig. 2 gives an overview of different tasks addressed by the generic cost function (50).

*G. Relation to BSS*

By choosing an uninformative prior over the demixing matrices $p(\mathcal{W}) = \mathrm{const.}$ and the number of SOIs equal to the number of microphones $K = M$, the cost function for non-informed determined IVA is obtained [2]

$$
J_{\mathrm{BSS}}(\mathcal{W}) = \sum_{k=1}^{K} \hat{\mathbb{E}}\left\{G\left(\underline{\mathbf{s}}_{k,n}\right)\right\} - 2\sum_{f=1}^{F} \log |\det \mathbf{W}_f|. \qquad (51)
$$

Hence, the proposed framework includes the prior work based on IVA (and ICA as a special case of IVA) [7], [9], [10], [32] and its many extensions [16], [19], [27], [30].

## III. DERIVATION OF UPDATE RULES

In the following, we develop an optimization algorithm based on the MM principle for the general informed BSS cost function $J_{\mathrm{IBSS}}(\mathcal{W})$ (50). We will start with the fundamental MM principle and then construct an upper bound of the informed BSS cost function $J_{\mathrm{IBSS}}$. Finally, we will provide update rules and summarize the proposed algorithmic framework.

*A. Majorize-Minimize Principle*

The main idea of Majorize-Minimize (MM) algorithms is to define an upper bound for the cost function which is easier to optimize than the cost function itself and which fulfills two conditions: majorization and tangency (see [34] for an accessible in-depth introduction).

Let $\mathcal{W}^{(l)}$ denote the set of estimated demixing matrices at iteration $l \in [L]$ with $L$ as the total number of iterations. Then the majorization property of the upper bound $U\left(\mathcal{W}|\mathcal{W}^{(l)}\right)$ can be expressed as

$$
J(\mathcal{W}) \le U\left(\mathcal{W}|\mathcal{W}^{(l)}\right). \qquad (52)
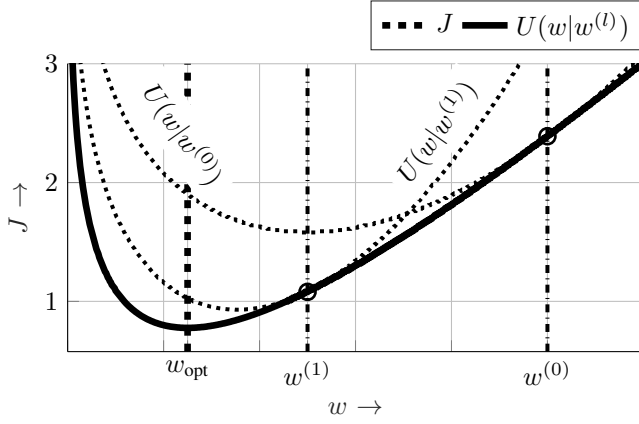$$

Fig. 3. Illustration of optimization based on the MM principle. Here, a one-dimensional cost function is used for illustration. The cost function $J$ is shown as a solid line and the upper bounds $U(w|w^{(l)})$ for $l = 0, 1$ as dotted lines. Furthermore, the global minimizer $w_{\text{opt}}$ and the minimizer of $U(w|w^{(0)})$ are shown as vertical lines.

Equality holds iff $\mathcal{W} = \mathcal{W}^{(l)}$, i.e.,

$$J\left(\mathcal{W}^{(l)}\right) = U\left(\mathcal{W}^{(l)}|\mathcal{W}^{(l)}\right), \qquad (53)$$

which represents the tangency condition. The upper bound is chosen such that its optimization is easily possible

$$\mathcal{W}^{(l+1)} = \underset{\mathcal{W}}{\arg\min}\, U\left(\mathcal{W}|\mathcal{W}^{(l)}\right), \qquad (54)$$

where $\mathcal{W}^{(l+1)}$ denotes the minimizer. As minimization does not increase the function value of the upper bound, the following downhill property [34] is obtained by using the tangency and majorization property of the upper bound

$$J\left(\mathcal{W}^{(l+1)}\right) \leq U\left(\mathcal{W}^{(l+1)}|\mathcal{W}^{(l)}\right) \qquad (55)$$
$$\leq U\left(\mathcal{W}^{(l)}|\mathcal{W}^{(l)}\right) = J\left(\mathcal{W}^{(l)}\right).$$

Hence, by iteratively optimizing the upper bound and ensuring tangency to the cost function, the cost function values are ensured to be non-increasing.

This optimization principle is illustrated in Fig. 3.

### B. Construction of Upper Bound

The problem of optimizing the informed BSS cost function $J_{\text{IBSS}}$ will now be shifted to optimizing a surrogate, an upper bound $U_{\text{IBSS}}$.

Let $\mathcal{W}_k^{(l)} = \left\{\mathbf{w}_f^{k,(l)} \in \mathbb{C}^K | f \in [F]\right\}$ be the set of all demixing vectors for channel $k$ at iteration $l$. For supergaussian PDFs (for the discussion of the time-varying Gaussian PDF see below), characterized by the score function $G(\underline{\mathbf{s}}_{k,n})$, the following inequality has been proven in [10]

$$\hat{\mathbb{E}}\left\{G(\underline{\mathbf{s}}_{k,n})\right\} \leq R_k(\mathcal{W}_k^{(l)}) + \frac{1}{2}\sum_{f=1}^{F}\left(\mathbf{w}_f^k\right)^{\text{H}}\mathbf{V}_f^k\left(\mathcal{W}_k^{(l)}\right)\mathbf{w}_f^k. \qquad (56)$$

All discussed SOI models can be written solely in dependence of the norm of the broadband SOI signal $r_{k,f,n}(\mathcal{W}_k^{(l)})$, i.e., $\tilde{G}(r_{k,f,n}(\mathcal{W}_k^{(l)})) = G(\underline{\mathbf{s}}_{k,n})$. For the supergaussian and the

time-varying Gaussian SOI model, the weighting factor depends on the estimated broadband signal energy of source $k$ at time instant $n$

$$r_{k,n}\left(\mathcal{W}_k^{(l)}\right) = \left\|\underline{\mathbf{s}}_{k,n}^{(l)}\right\|_2 = \sqrt{\sum_{f=1}^{F}\left|\left(\mathbf{w}_f^{k,(l)}\right)^{\text{H}}\mathbf{x}_{f,n}\right|^2}, \quad (57)$$

i.e., $r_{k,f,n} = r_{k,n} \quad \forall f$. The term $R_k(\mathcal{W}_k^{(l)})$ in (56) given as

$$R_k\left(\mathcal{W}_k^{(l)}\right) = \hat{\mathbb{E}}\bigg\{\tilde{G}\left(r_{k,n,f}\left(\mathcal{W}_k^{(l)}\right)\right)\dots \qquad (58)$$
$$\dots - \frac{r_{k,n,f}\left(\mathcal{W}_k^{(l)}\right)\tilde{G}'\left(r_{k,n,f}\left(\mathcal{W}_k^{(l)}\right)\right)}{2}\bigg\}$$

is independent of $\mathcal{W}$ and $\mathbf{V}_f^k\left(\mathcal{W}_k^{(l)}\right)$ denotes the weighted sensor signals' covariance matrix

$$\mathbf{V}_f^k\left(\mathcal{W}_k^{(l)}\right) = \hat{\mathbb{E}}\left\{\phi(r_{k,f,n})\mathbf{x}_{f,n}\mathbf{x}_{f,n}^{\text{H}}\right\}, \qquad (59)$$

where

$$\phi(r_{k,f,n}) = \frac{\tilde{G}'\left(r_{k,f,n}\left(\mathcal{W}_k^{(l)}\right)\right)}{r_{k,f,n}\left(\mathcal{W}_k^{(l)}\right)} \qquad (60)$$

denotes the corresponding weighting factor.

The weighting factor $\phi(r_{k,n})$ for the generalized Gaussian distribution (21) and the time-varying Gaussian PDF (23) can be expressed as (see [32])

$$\phi(r_{k,n}) = (r_{k,n})^{\beta-2}. \qquad (61)$$

For the NMF source model, we obtain for the weighting factor

$$\phi(r_{k,f,n}) = \frac{1}{\left(\sum_{\nu=1}^{N_{\text{bases}}} t_{k,f,\nu} v_{k,\nu,n}\right)^\beta}. \qquad (62)$$

Note that the weighting factor $\phi(r_{k,n,f})$ is frequency-dependent in the case of the NMF source model.

The inequality (56) transforms the optimization of a general nonlinear function dependent on all frequency bins into the optimization of the sum of quadratic functions, each of which dependent only on one frequency bin. The dependency between the frequency bins is solely expressed by the weighting $\phi(r_{k,n})$ of the microphone correlation matrix in (59).

By inserting the inequality (56) into the BSS cost function (51), we obtain the following upper bound for the BSS cost function $J_{\text{BSS}}$

$$U_{\text{BSS}}\left(\mathcal{W}|\mathcal{W}^{(l)}\right) = \sum_{f=1}^{F}\bigg[\sum_{k=1}^{K}\bigg(\frac{1}{2}\left(\mathbf{w}_f^k\right)^{\text{H}}\mathbf{V}_f^k\left(\mathcal{W}_k^{(l)}\right)\mathbf{w}_f^k\dots$$
$$\dots + \frac{1}{F}R_k\left(\mathcal{W}_k^{(l)}\right)\bigg) - 2\log|\det\mathbf{W}_f|\bigg], \qquad (63)$$

with $J_{\text{BSS}}(\mathcal{W}) = U_{\text{BSS}}\left(\mathcal{W}|\mathcal{W}^{(l)}\right)$ iff $\mathcal{W} = \mathcal{W}^{(l)}$.

For the case of a Gaussian source distribution, the upper bound is identical to the cost function (a similar relation holds for the NMF source model described in Sec. II-C3)

$$J_{\text{BSS}}\left(\mathcal{W}|\mathcal{W}^{(l)}\right) = U_{\text{BSS}}\left(\mathcal{W}|\mathcal{W}^{(l)}\right), \qquad (64)$$

where $R_k\left(\mathcal{W}_k^{(l)}\right) = 0$.

An upper bound of the cost function for informed BSS $J_{\text{IBSS}}(\mathcal{W})$ can be obtained by adding the cost function of the prior $J_{\text{prior}}$ (49) and the cost function of the BG $J_{\text{BG}}$ (39) on both sides of the inequality

$$J_{\text{IBSS}}(\mathcal{W}) \leq U_{\text{IBSS}}\left(\mathcal{W}|\mathcal{W}^{(l)}\right) \tag{65}$$
$$= U_{\text{BSS}}\left(\mathcal{W}|\mathcal{W}^{(l)}\right) + J_{\text{BG}}(\mathcal{W}) + J_{\text{prior}}(\mathcal{W}),$$

with $J_{\text{IBSS}}(\mathcal{W}) = U_{\text{IBSS}}\left(\mathcal{W}|\mathcal{W}^{(l)}\right)$ iff $\mathcal{W} = \mathcal{W}^{(l)}$, i.e., the upper bound fulfills the requirements of majorization and tangency.

### C. Optimization of Upper Bound

In the following we will derive analytic expressions for the minimum of the upper bound w.r.t. the demixing matrices

$$\mathcal{W}^{(l+1)} = \operatorname*{argmin}_{\mathcal{W}} U_{\text{IBSS}}\left(\mathcal{W}|\mathcal{W}^{(l)}\right) \tag{66}$$

and derive iterative update rules which allow the computation of the minimizer $\mathcal{W}^{(l+1)}$. To simplify the following derivation, we transform the log-det term of the upper bound (63) to have all BG filters in the transformed representation (35)

$$\log|\det \mathbf{W}_f| = \log\left|\det\begin{bmatrix} \mathbf{I}_K & \mathbf{0}_{M-K\times K} \\ \mathbf{0}_{K\times M-K} & \mathbf{T}_f\mathbf{D}_f^{-1} \end{bmatrix}\begin{bmatrix} \mathbf{W}_f^{\text{SOI}} \\ \tilde{\mathbf{B}}_f \end{bmatrix}\right|$$
$$= \log\left|\det\begin{bmatrix} \mathbf{W}_f^{\text{SOI}} \\ \tilde{\mathbf{B}}_f \end{bmatrix}\right| + \text{const.} \tag{67}$$

Hence, the transformed filters yield the same optimum as the orignal filters.

*1) Without Constraints:* For the unconstrained channels, i.e., for $k \notin \mathcal{I}$ and $k \notin \mathcal{I}^{\text{Euc}}$, we obtain the following conditions by setting the derivative of the upper bound (65) w.r.t. each of the SOI filters to zero [10]

$$\left(\mathbf{w}_f^q\right)^{\text{H}}\mathbf{V}_f^k\left(\mathcal{W}_k^{(l)}\right)\mathbf{w}_f^k \overset{!}{=} \delta_{kq}, \quad k,q \in [K] \tag{68}$$

where $\delta$ denotes the Kronecker Delta. Similarly, for the BG filters we obtain by differentiating (65) the following conditions for the relation between the SOI filters $k \in [K]$ and the BG filters $k' \in [M-K]$

$$\left(\mathbf{w}_f^k\right)^{\text{H}}\mathbf{C}_f\tilde{\mathbf{b}}_f^{k'} \overset{!}{=} 0 \tag{69}$$

and for the relation between the BG filters

$$\left(\tilde{\mathbf{b}}_f^q\right)^{\text{H}}\mathbf{C}_f\tilde{\mathbf{b}}_f^{k'} \overset{!}{=} \delta_{k'q}, \quad q \in [M-K]. \tag{70}$$

However, this condition is not investigated further in the following, as the estimation of the BG signals is not our aim. By collecting all the vector-wise constraints in (69), we can write

$$\mathbf{W}_f^{\text{SOI}}\mathbf{C}_f\tilde{\mathbf{B}}_f^{\text{H}} \overset{!}{=} \mathbf{0}_{K\times(M-K)}. \tag{71}$$

Now, we insert $\tilde{\mathbf{B}}_f = \mathbf{D}_f\mathbf{T}_f^{\text{H}}\mathbf{B}_f$

$$\mathbf{W}_f^{\text{SOI}}\mathbf{C}_f\mathbf{B}_f^{\text{H}}\mathbf{T}_f\mathbf{D}_f \overset{!}{=} \mathbf{0}_{K\times(M-K)} \tag{72}$$

and multiply with $\mathbf{D}_f^{-1}\mathbf{T}_f^{\text{H}}$ from the right, which yields the following condition between SOI and BG filters

$$\mathbf{W}_f^{\text{SOI}}\mathbf{C}_f\mathbf{B}_f^{\text{H}} \overset{!}{=} \mathbf{0}_{K\times(M-K)}. \tag{73}$$

*2) With Constraints:* For the channels constrained by the quadratic constraint (41), i.e., $k \in \mathcal{I}$, we obtain as conditions for the SOI channels by optimizing (65)

$$\left(\mathbf{w}_f^q\right)^{\text{H}}\left[\mathbf{V}_f^k\left(\mathcal{W}_k^{(l)}\right) + \gamma_{k,f}\mathbf{P}_f^k\right]\mathbf{w}_f^k \overset{!}{=} \delta_{kq}. \tag{74}$$

For the relation between the SOI and the BG channels we obtain

$$\mathbf{W}_f^{\text{SOI}}\left[\mathbf{C}_f + \gamma_f^{\text{BG}}\mathbf{P}_f^{\text{BG}}\right]\mathbf{B}_f^{\text{H}} \overset{!}{=} \mathbf{0}_{K\times(M-K)}. \tag{75}$$

Note that the conditions (74) and (75) generalize the previously known conditions (68) and (70) in the sense that the weighted correlation matrix $\mathbf{V}_f^k$ and the microphone signal correlation matrix $\mathbf{C}_f$ are regularized by the precision matrices $\mathbf{P}_f^k$ and $\mathbf{P}_f^{\text{BG}}$, which allow incorporation of many types of prior knowledge on SOIs and/or BG as discussed in Sec. II-E.

### D. Update Rules

In the following, we will present update rules which identify solutions to the conditions (68), (73), (74) and (75) presented in the previous paragraph.

*1) Demixing Filters:* In the unconstrained case the SOI filters can be optimized by ensuring orthogonality between the output signals [10]

$$\tilde{\mathbf{w}}_f^{k,(l+1)} = \left(\mathbf{W}_f^{k,(l)}\mathbf{V}_f^{k,(l)}\left(\mathcal{W}_k^{(l)}\right)\right)^{-1}\mathbf{e}_k, \tag{76}$$

where $\mathbf{e}_k$ denotes a canonical basis vector with a one at the $k$th position, and normalization

$$\mathbf{w}_f^{k,(l+1)} = \frac{\tilde{\mathbf{w}}_f^{k,(l+1)}}{\sqrt{\left(\tilde{\mathbf{w}}_f^{k,(l+1)}\right)^{\text{H}}\mathbf{V}_f^{k,(l)}\left(\mathcal{W}_k^{(l)}\right)\tilde{\mathbf{w}}_f^{k,(l+1)}}}. \tag{77}$$

This procedure is called IP and will be used to derive generalized update rules for the other algorithmic variants in the following. The channels constrained by (41), i.e., $k \in \mathcal{I}$ are updated by

$$\tilde{\mathbf{w}}_f^{k,(l+1)} = \left(\mathbf{W}_f^{(l)}\left[\mathbf{V}_f^{k,(l)}\left(\mathcal{W}_k^{(l)}\right) + \gamma_{k,f}\mathbf{P}_f^k\right]\right)^{-1}\mathbf{e}_k, \tag{78}$$

$$\mathbf{w}_f^{k,(l+1)} = \frac{\tilde{\mathbf{w}}_f^{k,(l+1)}}{\sqrt{\left(\tilde{\mathbf{w}}_f^{k,(l+1)}\right)^{\text{H}}\left[\mathbf{V}_f^{k,(l)}\left(\mathcal{W}_k^{(l)}\right) + \gamma_{k,f}\mathbf{P}_f^k\right]\tilde{\mathbf{w}}_f^{k,(l+1)}}}. \tag{79}$$

For the channels constrained by (46), i.e., $k \in \mathcal{I}^{\text{Euc}}$, we use the update rules proposed by [27]

$$\mathbf{u}_f^k = \left(\mathbf{W}_f^{(l)}\tilde{\mathbf{V}}_f^{k,(l)}\right)^{-1}\mathbf{e}_k \tag{80}$$

$$\tilde{\mathbf{u}}_f^k = \gamma_{k,f}^{\text{Euc}}\left(\tilde{\mathbf{V}}_f^{k,(l)}\right)^{-1}\mathbf{h}_f(\vartheta_k) \tag{81}$$

$$p_{k,f} = (\mathbf{u}_f^k)^{\text{H}}\tilde{\mathbf{V}}_f^{k,(l)}\mathbf{u}_f^k \tag{82}$$

$$\tilde{p}_{k,f} = (\mathbf{u}_f^k)^{\text{H}}\tilde{\mathbf{V}}_f^{k,(l)}\tilde{\mathbf{u}}_f^k \tag{83}$$

$$\tilde{\mathbf{w}}_f^{k,(l+1)} \leftarrow \begin{cases} \frac{\mathbf{u}_f^k}{\sqrt{p_{k,f}}} + \tilde{\mathbf{u}}_f^k, & \text{if } \tilde{p}_{k,f} = 0 \\ \frac{\tilde{p}_{k,f}}{2p_{k,f}}\left(-1 + \sqrt{1 + \frac{4p_{k,f}}{|\tilde{p}_{k,f}|^2}}\right)\mathbf{u}_f^k + \tilde{\mathbf{u}}_f^k, & \text{else.} \end{cases} \tag{84}$$

| | Algorithm Index → | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| $K$ | $M$ | $M$ | $M$ | $M$ | 1 | 1 | 1 | $M$ | $M$ | $M$ | 1 | 1 | 1 |
| Optimization type | GD | IP | IP | IP | IP | IP | IP | IP | IP | IP | IP | IP | IP |
| Spatial One/Null | One | One | One | Null | One | One | Null | One | One | Null | One | One | Null |
| Quadratic prior (41) | × | × | ✓ | ✓ | × | ✓ | ✓ | × | ✓ | ✓ | × | ✓ | ✓ |
| Euclidean prior (46) | × | ✓ | × | × | ✓ | × | × | ✓ | × | × | ✓ | × | × |
| BG model | × | × | × | × | ✓ | ✓ | ✓ | × | × | × | ✓ | ✓ | ✓ |
| BG prior | × | × | × | × | × | × | ✓ | × | × | × | × | × | ✓ |
| SOI model | SG | —————— SG/TVG —————— | | | | | | ————————— NMF ————————— | | | | | |
| Proposed | [26] | — New — | | [30] | ——— New ——— | | | [27] | ————— New ————— | | | | |

TABLE II

OVERVIEW OVER ALGORITHMIC VARIANTS EVALUATED IN THE EXPERIMENTS. WE USED THE FOLLOWING ABBREVIATIONS: GRADIENT DESCENT (GD), ITERATIVE PROJECTION (IP), SUPERGAUSSIAN (SG) AND TIME-VARYING GAUSSIAN (TVG).

To calculate the update of the BG filters $\mathbf{B}_f^{M,K}$ in the unconstrained case, (73) can be solved for $\mathbf{B}_f^{M,K}$ by inserting the parametrization of the BG filters, which yields

$$\mathbf{B}_f^{M,K} = \left(\mathbf{E}_2 \mathbf{C}_f (\mathbf{W}_f^{\mathrm{SOI}})^{\mathrm{H}}\right) \left(\mathbf{E}_1 \mathbf{C}_f (\mathbf{W}_f^{\mathrm{SOI}})^{\mathrm{H}}\right)^{-1}. \quad (85)$$

Hereby, we defined

$$\mathbf{E}_1 = [\mathbf{I}_K, \mathbf{0}_{K \times M-K}] \quad \text{and} \quad \mathbf{E}_2 = [\mathbf{0}_{M-K \times K}, \mathbf{I}_{M-K}]. \quad (86)$$

Note that these update rules coincide with those proposed by [19], but are rigorously derived here from the iterative projection perspective, which also makes the incorporation of priors possible. Similarly, the updates for the constrained case are obtained by generalization of (85) as

$$\mathbf{B}_f^{M,K} = \left(\mathbf{E}_2 \left[\mathbf{C}_f + \gamma_f^{\mathrm{BG}} \mathbf{P}_f^{\mathrm{BG}}\right] (\mathbf{W}_f^{\mathrm{SOI}})^{\mathrm{H}}\right) \cdots$$
$$\cdots \left(\mathbf{E}_1 \left[\mathbf{C}_f + \gamma_f^{\mathrm{BG}} \mathbf{P}_f^{\mathrm{BG}}\right] (\mathbf{W}_f^{\mathrm{SOI}})^{\mathrm{H}}\right)^{-1}. \quad (87)$$

*2) Update of Demixed Signal Variance:* The update of the variance parameter $r_{k,n,f}$ can be done directly based on the demixed signals for each iteration in case of the generalized Gaussian or time-varying Gaussian source model by (57). For the NMF source model, the elements $t_{k,f,\nu}$ of the basis vectors and the elements $v_{k,\nu,n}$ of the activation vector have to be updated in addition to the demixing filters. The update rules are given by [16]

$$t_{k,f,\nu} \leftarrow t_{k,f,\nu} \sqrt{\frac{\sum_{n \in [N]} |y_{k,f,n}|^2 v_{k,\nu,n} \left(r_{n,f}^k\right)^{-2}}{\sum_{n \in [N]} v_{k,\nu,n} \left(r_{n,f}^k\right)^{-1}}} \quad (88)$$

and

$$v_{k,\nu,n} \leftarrow v_{k,\nu,n} \sqrt{\frac{\sum_{f \in [F]} |y_{k,f,n}|^2 t_{k,f,\nu} \left(r_{n,f}^k\right)^{-2}}{\sum_{f \in [F]} t_{k,f,\nu} \left(r_{n,f}^k\right)^{-1}}}. \quad (89)$$

*E. Practical Aspects*

In this paragraph, we discuss some aspects which are relevant for a practical realization of the above algorithmic variants. To avoid distortion of the signals by the scaling ambiguity in each frequency bin, the minimal distortion principle can be applied [35]. To avoid numerical instability of the algorithmic variants relying on an NMF SOI model, [15]

proposed to normalize all estimated quantities in each iteration (see [15] for details). The proposed algorithmic framework is summarized in Alg. 1.

## IV. EXPERIMENTS

In this section, we evaluate different algorithmic variants resulting from the proposed framework and compare them with several baseline algorithms from the literature. In this experimental study, we will focus on signal extraction, i.e., the separation from one source out of the observed mixture. In addition, the challenging case of an underdetermined scenario, i.e., $Q > M$ is addressed in the experiments in the following. However, also the extraction of multiple sources from the mixture and source separation for the determined case, i.e., $K = M$, and the overdetermined case, i.e., $K > M$, are covered by the framework. We do not evaluate the determined case here as this has been subject to many experimental studies in the literature [9], [32]. We also do not investigate the overdetermined case, as this can be considered as an easier problem than the underdetermined scenario. A discussion for the overdetermined case without the incorporation of prior knowledge can be found in [19].

The discussed methods vary w.r.t. the used SOI model, the exploitation of a BG model, the optimization method and the applied priors. Method 1 is based on gradient descent and a supergaussian source model and has been proposed in [26]. The rest of the discussed algorithmic variants all use IP for optimization and are evaluated for different SOI models: the supergaussian, the time-varying Gaussian and the NMF SOI model. For each of these SOI models, we discuss the priors (41) with (45) and (46) constraining one channel by a spatial one and the prior (41) with (43) constraining all channels but one with a spatial null. Furthermore, we discuss for all source models the incorporation of the BG model in two different variants: 1) unconstrained BG with a spatial one constraint for the SOI ((41) with (45) or (46)) and 2) unconstrained SOI, but BG with a spatial null constraint (47). Tab. II summarizes the 13 algorithmic variants discussed in the following. The variants 4 and 8 are published in [30] and [27], respectively, and represent further baselines in our experimental study. Note that [19], which is a special case of the proposed framework, has been shown to be superior to [18] by comprehensive experiments. Hence, we do not repeat these experiments here.

| | 1 | | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Algorithm Index $\rightarrow$ | | | | | | |
| Step size | 0.05 | $\gamma,\gamma^{\text{Euc}},\gamma^{\text{BG}}$ | 0.5 | 1.5 | 0.5 | 2 | 2 | 50 | 5 | 3 | 5 | 2.5 | 2.5 | 100 |
| Prior Weight | 0.01 | $\lambda_{\text{Tik}}$ | 1 | 1 | $10^{-3}$ | 1 | 1 | $10^{-3}$ | 1 | 1 | $10^{-3}$ | 1 | 1 | $10^{-3}$ |
| | | $\lambda_1^{\text{One}},\lambda_1^{\text{Zero}}$ | $\times$ | 2 | 1 | $\times$ | 1.5 | 1 | $\times$ | 1.5 | 1 | $\times$ | 1 | 1 |
| | | $N_{\text{bases}}$ | $\times$ | $\times$ | $\times$ | $\times$ | $\times$ | $\times$ | 2 | 2 | 2 | 2 | 2 | 2 |
| $L$ | 2500 | $L$ | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |

TABLE III
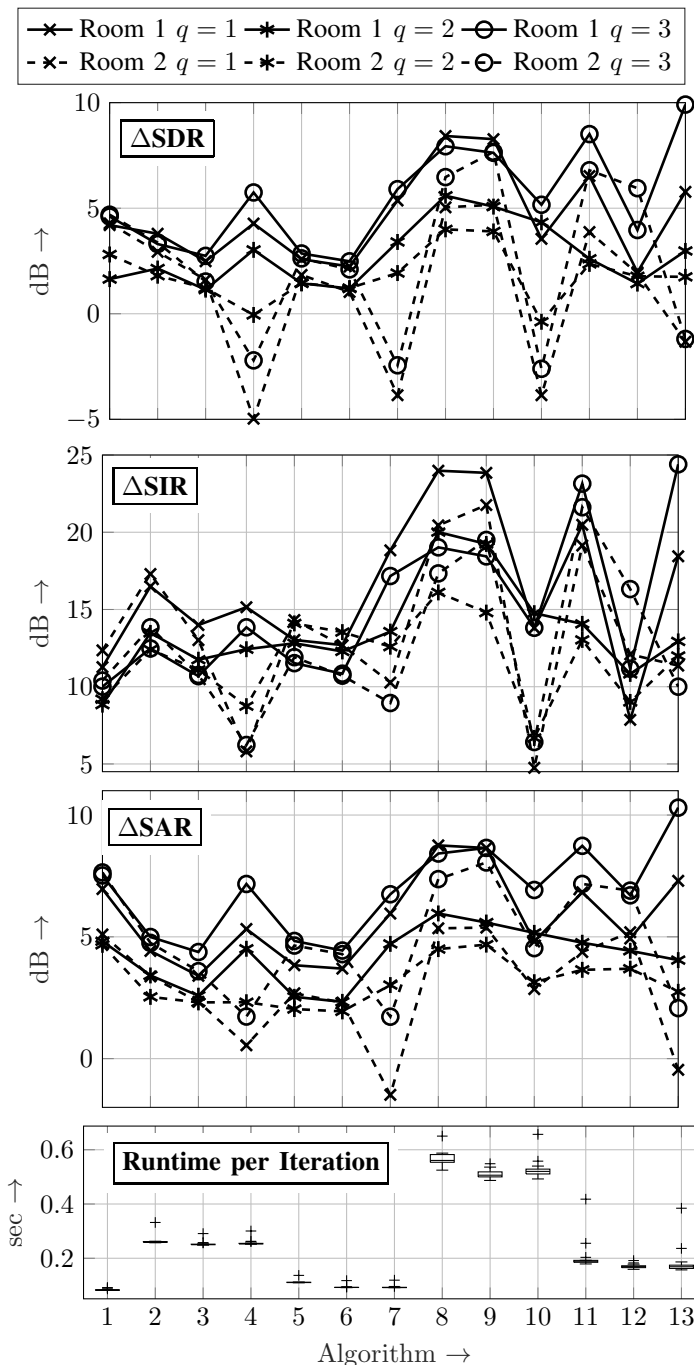PARAMETERS USED IN THE EXPERIMENTS.



Fig. 4. Improvement in performance measures [36] and average runtime per iteration for different extracted sources ($q = 1, 2, 3$, see Fig. 5 for the geometric setup) and two different rooms: Room 1 with $T_{60} = 0.2\,\text{s}$ and Room 2 with $T_{60} = 0.4\,\text{s}$.
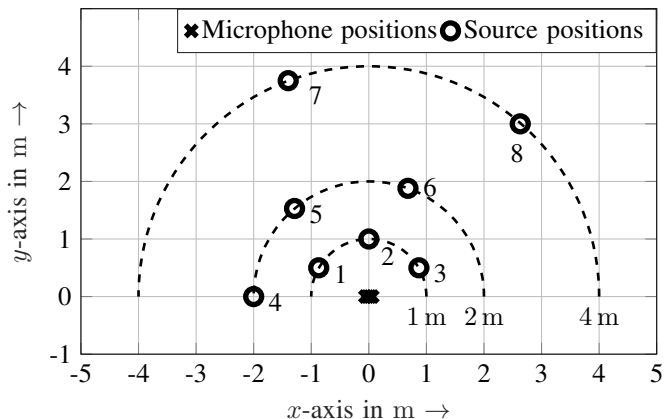


Fig. 5. Geometric setup of the scenario used in the experiments. The $M = 4$ microphone positions are marked by crosses and the $Q = 8$ source positions at $1\,\text{m}$, $2\,\text{m}$ and $4\,\text{m}$ distance from the array are marked by circles.
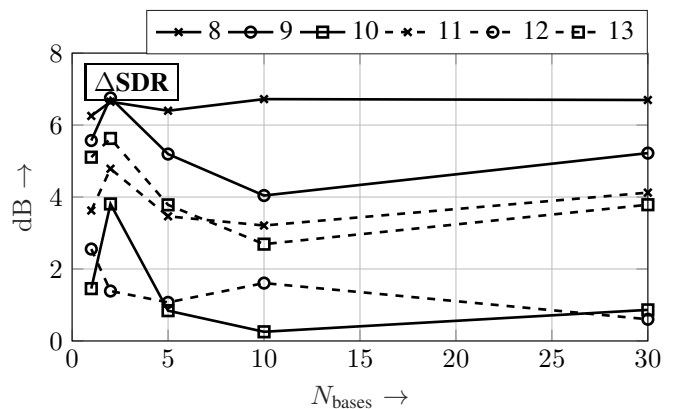


Fig. 6. Results of the number of bases $N_{\text{bases}}$ for the algorithmic variants using an NMF source model. The results for the approaches using a BG model are depicted as solid lines, the others as dashed lines.

### A. Experimental Setup

For the experiments we used a uniform linear array with $M = 4$ microphones with a spacing of $4.2\,\text{cm}$. The microphone signals are computed by convolving RIRs measured in a living room environment with male and female speech signals and adding white Gaussian noise such that an Signal-to-Noise Ratio (SNR) of $30\,\text{dB}$ at the microphones is obtained. Two enclosures are considered in the following: Room 1 with a reverberation time of $T_{60} = 0.2\,\text{s}$ and Room 2 with $T_{60} = 0.4\,\text{s}$. We placed $Q = 8$ acoustic sources at $1\,\text{m}$, $2\,\text{m}$ and $4\,\text{m}$ distance and at different angles relative to the array for measuring the RIRs (see Fig. 5 for an illustration

**Algorithm 1** Informed BSS (generic pseudo code)

**INPUT:** $\mathcal{X}$, $L$, $\{\Theta_k\}_{k\in\mathcal{I}}$, $\{\Theta_k\}_{k\in\mathcal{I}^{\mathrm{Euc}}}$, $\Theta_{\mathrm{BG}}$

---

**INITIALIZATION:**

$\mathbf{y}_{f,n} = \mathbf{x}_{f,n} \ \forall f, n$

**if** NMF Source Model **then**

   $t_{k,f,\nu}, v_{k,\nu,n} \sim \mathcal{U}(0,1) \ \forall k, f, n, \nu$

**end if**

**if** $M \leq K$ **then**

   $\mathbf{W}_f^{(0)} = \mathbf{I}_M \ \forall f$

**else**

   $\mathbf{W}_f^{(0)} = \begin{bmatrix} \mathbf{I}_K & \mathbf{0}_{K\times(M-K)} \\ \mathbf{0}_{(M-K)\times K} & -\mathbf{I}_{M-K} \end{bmatrix} \ \forall f$

**end if**

---

**for** $l = 1$ **to** $L$ **do**

  **for** $k = 1$ **to** $K$ **do**

    Calculate $\phi(r_{k,f,n}) \ \forall n$ by (61) or (62)

    **for** $f = 1$ **to** $F$ **do**

      Calculate $\mathbf{V}_f^k(\mathcal{W}_k^{(l)}) = \hat{\mathbb{E}}\left\{\phi(r_{k,f,n})\mathbf{x}_{f,n}\mathbf{x}_{f,n}^{\mathrm{H}}\right\}$

      **if** $k \in \mathcal{I}$ or $k \in \mathcal{I}^{\mathrm{Euc}}$ **then**

        Update $\mathbf{w}_f^k$ by (78), (79) or by (80)-(84)

      **else if** $k \notin \mathcal{I}$ **then**

        Update $\mathbf{w}_f^k$ by (76) and (77)

      **end if**

      **if** $M > K$ **then**

        **if** $\Theta_{\mathrm{BG}} \neq \emptyset$ **then**

          Update $\mathbf{B}_f^{M,K}$ by (87)

        **else**

          Update $\mathbf{B}_f^{M,K}$ by (85)

        **end if**

      **end if**

      Assemble $\mathbf{W}_f = \begin{bmatrix} [\mathbf{w}_f^1, \ldots, \mathbf{w}_f^K]^{\mathrm{H}} \\ [\mathbf{B}_f^{M,K} \quad -\mathbf{I}_{M-K}] \end{bmatrix}$

    **end for**

  **end for**

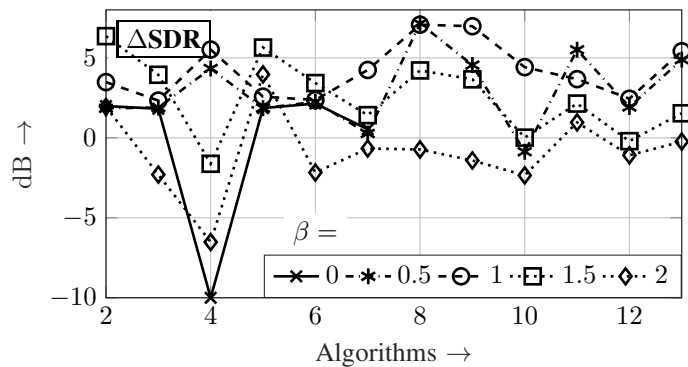  **if** NMF Source Model **then**

    Normalize [15]

  **end if**

**end for**

Scale demixing filters $\mathbf{W}_f \leftarrow \mathrm{diag}\left\{(\mathbf{W}_f)^{-1}\right\}\mathbf{W}_f$

**for** $n = 1$ **to** $N$ **do**

  **for** $f = 1$ **to** $F$ **do**

    Extract SOIs $\mathbf{s}_{f,n} = \mathbf{W}_f^{\mathrm{SOI}}\mathbf{x}_{f,n}$

  **end for**

**end for**

---

**OUTPUT:** SOIs $\mathbf{s}_{f,n} \forall f, n$



Fig. 7. Influence of the shape parameter $\beta$ of the SOI model on the performance of Methods 2-13 in terms of SDR improvement.

of the geometric setup of the measurements). All sources and microphones have been placed at the same height of $1.4\,\mathrm{m}$. The microphone signals are computed from a set of 4 female and 4 male speech signals of $20\,\mathrm{s}$ duration at a sampling frequency of $16\,\mathrm{kHz}$. The microphone signals are transformed into the STFT domain using a von Hann window of length $2048$ and $50\%$ overlap. For the SOI source models, we set $\beta = 1$ in (61) and (62). The performance of the investigated methods is measured in terms of the improvement (denoted by $\Delta$) of the Signal-to-Distortion Ratio (SDR), Signal-to-Interference Ratio (SIR) and Signal-to-Artefact Ratio (SAR) [36] w.r.t. the unprocessed microphone signals, respectively, and in terms of averaged runtime per iteration for all 20 permutations of the source signals.

In the following, we aim at extracting a source $q$ (see Fig. 5) out of the reverberant mixture of all sources. To obtain representative results, we repeat the experiment 20 times and permute the positions of the speech sources in each trial. The performance of the algorithms is assessed by using the improvement for the measures proposed by [36], where the separation of the SOI from the mixture of all other signals is evaluated. The user-defined parameters are chosen for each algorithmic variant separately by a parameter sweep such that the best results are obtained on average for the extraction of source $q = 2$ for all 20 permutations (the choice of $q = 2$ is arbitrary here). Furthermore, the parameters have been chosen such that the outer permutation has been resolved, i.e., the desired source signal indeed appeared at the selected output channel. The weighting parameters $\lambda$ and $\gamma$ have chosen to be equal for all frequency bins and channels. The obtained parameters are summarized in Tab. III.

### B. Target Direction and Acoustic Environment

The influence of different target DOAs (corresponding to sources $q = 1, 2, 3$) and of different acoustic environments is investigated in the following. To this end, the geometric setup, corresponding to Fig. 5, is used in the two different rooms described above for measuring the RIRs and for each of these acoustic conditions source $q = 1, 2, 3$ is extracted. This experiment is again repeated for 20 permutations of the association between source positions and source signals and the median of the results is taken as a statistic, which
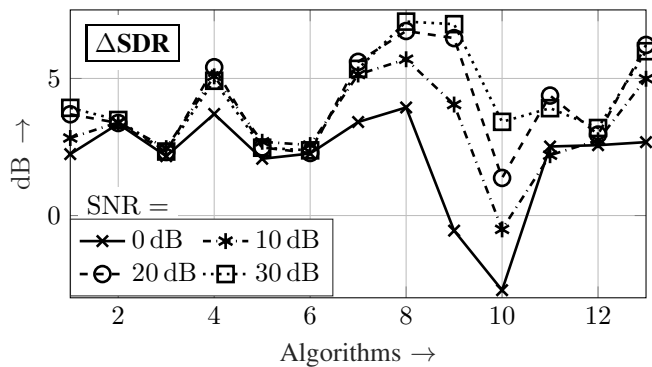
Fig. 8. Influence of different noise levels on the discussed algorithmic variants in terms of SDR improvement.

is presented in Fig. 4. The results of Room 1 are depicted as solid lines, the results of Room 2 as dashed lines. First of all, it can be seen that the extraction of source $q = 3$ yielded the best results in terms of SDR improvement for most algorithms, which may be explained by the geometric setup in which not many sources are contained in the angular region of source $q = 3$. Furthermore, the performance of all algorithms degrades for Room 2, which has a higher reverberation time. This effect is typical for algorithms which perform spatial filtering. Also the assumption of free-field propagation used for the construction of the priors is violated for an increasing reverberation time. While the performance of most of the algorithms dropped only slightly, for the Methods 4, 7, 10, 13 a large drop can be observed. These methods have in common that they rely on the prior (41) or (47) steering a spatial null. This spatial null constraint is imposed on all channels but one, instead of the priors steering a spatial one, which just impose a constraint on a single channel. As the free-field assumption is violated for increasing reverberation time, this has a larger effect on the methods using a prior steering a spatial null as this violated assumption is used multiple times. However, even for the methods with the large drop in the performance measures, SIR improvement is achieved.

### C. Runtime, Source Models and SNR

In terms of average runtime per iteration, Method 1 and 5-7 cause the lowest computational costs, followed by Methods 11-13. Hereby, the computational efficiency of the Methods 5-7 and 11-13 results from the usage of a BG model. The computational cost of the Methods 2-4 and 8-10 is much higher than their counterparts using a BG model. In terms of computational efforts to be spent until convergence, the gradient-based Method 1 is computationally much more costly as the number of iterations until convergence is much larger (about the factor $20 - 25$) than for the IP-based methods.

The influence of the number of bases $N_{\text{bases}}$ for the Methods 8-13 relying on an NMF source model is shown in Fig. 6. It can be seen that for all methods $N_{\text{bases}} = 2$ basis vectors provide satisfying results (see also, e.g., [16]).

The influence of the shape factor $\beta$ of the SOI models is discussed in terms of achieved SDR improvement in Fig. 7. The values $\beta = 0, 0.5, 1, 1.5, 2$ have been evaluated here (for

the NMF-based methods $\beta = 0$ is not evaluated as this would correspond to $\phi(r_{k,n,f}) = 1 \ \forall n, f, k$), where the value $\beta = 1$ corresponds to a Laplacian distribution and $\beta = 2$ to the time-varying Gaussian distribution (23) w.r.t. the IVA SOI models. In case of the NMF SOI model, a time-varying Gaussian SOI model is obtained for $\beta = 1$. Inspection of Fig. 7 shows that a choice of $\beta = 1$ yields good results for all algorithms. For some algorithmic variants the values of $\beta = 0.5$ or $\beta = 1.5$ are slightly better. In all cases, we obtain for the choice of $\beta = 0$ or $\beta = 2$ worse results. This is especially severe for Method 4, which relies on a prior steering a spatial one based on (41).

The performance of the discussed algorithmic variants w.r.t. varying noise levels is shown in Fig. 8. Here, we varied the additive noise, such that an SNR of $0\,\text{dB}$, $10\,\text{dB}$, $20\,\text{dB}$ and $30\,\text{dB}$ is achieved at the microphones. Unsurprisingly, for an SNR of $0\,\text{dB}$ all algorithms produce the worst results. For the other noise levels, a detrimental effect due to the additive noise can be observed for the algorithms relying on an NMF SOI model, whereas the other methods are only slightly affected by the noise level. The detrimental effect of the increasing noise level is especially severe for Methods 8, 9, 10, which are using an NMF source model and no BG model.

### D. Summary

In this experimental study, we discussed different algorithms based on IVA for source extraction, where the desired source is selected by a spatial constraint. In general, Methods 8-13 based on an NMF source model yielded better results than Methods 1-7 (see Fig. 4). As another general outcome, it can be observed that methods using a spatial null constraint degraded severely for increasing reverberation time. The influence of varying noise levels was not severe for most SNRs (see Fig. 8). The methods based on IP showed much lower computational complexity than the baseline using gradient descent [26] (see Fig. 4). The computational complexity can be further reduced significantly by the use of an BG model without sacrificing performance. By comparing the results shown in Fig. 4, it can be seen there is no single best-performing algorithm: For the TVG/SG source model, the proposed Algorithms 4 and 7 relying on a prior steering a spatial null perform especially well for $T_{60} = 0.2\,\text{s}$ and degrades for larger $T_{60}$. For the algorithmic variants relying on an NMF source model, the baseline Method 8 and the proposed Method 9, both steering a spatial one, yield similar results in all cases. However, the average runtime per iteration is slightly lower for the proposed Method 9. The proposed BG-based Methods 11-13 obtained for some acoustic setup very good results but degraded for $T_{60} = 0.4\,\text{s}$.

### V. CONCLUSION

In this contribution, we presented a unifying and flexible generic framework for systematic incorporation of prior knowledge on the demixing filters for IVA-based source separation algorithms. The potential of the framework was demonstrated for several exemplary priors representing geometric prior knowledge. As another generalization, a BG

model is incorporated into the framework, which allows for fast convergence of the corresponding algorithms at a low computational cost if the number of SOIs is smaller than the number of microphones, i.e., $M > K$. The derivation of update rules for the BG filters from this perspective had not been considered so far in the literature. For all proposed algorithmic variants, we derived stable and fast update rules with a low computational complexity based on the MM principle and the IP approach, even including most recently proposed update rules into the systematic framework.

The efficacy of the proposed algorithmic variants for real-world applications is demonstrated by experiments using measured RIRs and by comparison with established state-of-the-art baseline algorithms.

## References

[1] S. Gannot, E. Vincent, S. Markovich-Golan, and A. Ozerov, "A Consolidated Perspective on Multi-Microphone Speech Enhancement and Source Separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2017.

[2] E. Vincent, T. Virtanen, and S. Gannot, Eds., *Audio source separation and speech enhancement*. Hoboken, NJ: John Wiley & Sons, 2018.

[3] B. Van Veen and K. Buckley, "Beamforming: a versatile approach to spatial filtering," *IEEE ASSP Magazine*, vol. 5, no. 2, pp. 4–24, Apr. 1988.

[4] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent component analysis*. New York: J. Wiley, 2001.

[5] S. Makino, T.-W. Lee, and H. Sawada, Eds., *Blind speech separation*, ser. Signals and communication technology. Dordrecht: Springer, 2007.

[6] A. J. Bell and T. J. Sejnowski, "An Information-Maximization Approach to Blind Separation and Blind Deconvolution," *Neural Computation*, vol. 7, no. 6, pp. 1129–1159, Nov. 1995.

[7] P. Smaragdis, "Blind Separation of Convolved Mixtures in the Frequency Domain," *Neurocomputing Journal*, vol. 22, pp. 21–34, 1998.

[8] H. Sawada, R. Mukai, S. Araki, and S. Makino, "A Robust and Precise Method for Solving the Permutation Problem of Frequency-Domain Blind Source Separation," *IEEE Transactions on Speech and Audio Processing*, vol. 12, no. 5, pp. 530–538, Sep. 2004.

[9] T. Kim, H. T. Attias, S.-Y. Lee, and T.-W. Lee, "Blind Source Separation Exploiting Higher-Order Frequency Dependencies," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 15, no. 1, pp. 70–79, Jan. 2007.

[10] N. Ono, "Stable and fast update rules for independent vector analysis based on auxiliary function technique," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, New Paltz, NY, USA, Oct. 2011, pp. 189–192.

[11] H. Sawada, H. Kameoka, S. Araki, and N. Ueda, "Multichannel Extensions of Non-Negative Matrix Factorization With Complex-Valued Data," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 5, pp. 971–982, May 2013.

[12] D. D. Lee and H. S. Seung, "Algorithms for Non-negative Matrix Factorization," in *NIPS'00 Proceedings of the 13th International Conference on Neural Information Processing Systems*, 2000, pp. 535–541.

[13] 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, Mar. 2009.

[14] T. Haubner, A. Schmidt, and W. Kellermann, "Multichannel Nonnegative Matrix Factorization for Ego-Noise Suppression," in *13th ITG-Symposium Speech Communication*. Oldenburg, Germany: VDE, Oct. 2018.

[15] D. Kitamura, "Effective Optimization Algorithms for Blind and Supervised Music Source Separation with Nonnegative Matrix Factorization," PHD Thesis, SOKENDAI (The Graduate University for Advanced Studies), Mar. 2017.

[16] D. Kitamura, N. Ono, H. Sawada, H. Kameoka, and H. Saruwatari, "Determined Blind Source Separation Unifying Independent Vector Analysis and Nonnegative Matrix Factorization," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 9, pp. 1626–1641, Sep. 2016.

[17] T. Ono, N. Ono, and S. Sagayama, "User-guided independent vector analysis with source activity tuning," in *IEEE International Conference on Acoustic, Speech and Signal Processing (ICASSP)*, Kyoto, Japan, Mar. 2012, pp. 2417–2420.

[18] Z. Koldovský and P. Tichavský, "Gradient Algorithms for Complex Non-Gaussian Independent Component/Vector Extraction, Question of Convergence," *IEEE Transactions on Signal Processing*, vol. 67, no. 4, pp. 1050–1064, Feb. 2019.

[19] R. Scheibler and N. Ono, "Independent Vector Analysis with More Microphones than Sources," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, Oct. 2019.

[20] F. Nesta and Z. Koldovský, "Supervised independent vector analysis through pilot dependent components," in *IEEE International Conference on Acoustic, Speech and Signal Processing (ICASSP)*, Mar. 2017, pp. 536–540.

[21] T. Kounovský, Z. Koldovský, and J. Čmejla, "Recursive and Partially Supervised Algorithms for Speech Enhancement on the Basis of Independent Vector Extraction," in *2018 16th International Workshop on Acoustic Signal Enhancement (IWAENC)*, Tokyo, Sep. 2018, pp. 401–405.

[22] L. Parra and C. Alvino, "Geometric source separation: merging convolutive source separation with geometric beamforming," *IEEE Transactions on Speech and Audio Processing*, vol. 10, no. 6, pp. 352–362, Sep. 2002.

[23] Yuanhang Zheng, K. Reindl, and W. Kellermann, "BSS for improved interference estimation for Blind speech signal Extraction with two microphones," in *IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, Aruba, Dutch Antilles, Netherlands, Dec. 2009, pp. 253–256.

[24] K. Reindl, S. Meier, H. Barfuss, and W. Kellermann, "Minimum Mutual Information-Based Linearly Constrained Broadband Signal Extraction," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 6, pp. 1096–1108, Jun. 2014.

[25] Y. Zheng, K. Reindl, and W. Kellermann, "Analysis of dual-channel ICA-based blocking matrix for improved noise estimation," *EURASIP Journal on Advances in Signal Processing*, vol. 2014, no. 1, Dec. 2014.

[26] A. H. Khan, M. Taseska, and E. A. P. Habets, "A Geometrically Constrained Independent Vector Analysis Algorithm for Online Source Extraction," in *Latent Variable Analysis and Signal Separation*, E. Vincent, A. Yeredor, Z. Koldovský, and P. Tichavský, Eds. Cham: Springer International Publishing, 2015, vol. 9237, pp. 396–403.

[27] Y. Mitsui, N. Takamune, D. Kitamura, H. Saruwatari, Y. Takahashi, and K. Kondo, "Vectorwise Coordinate Descent Algorithm for Spatially Regularized Independent Low-Rank Matrix Analysis," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Calgary, AB, Apr. 2018, pp. 746–750.

[28] N. Q. K. Duong, E. Vincent, and R. Gribonval, "Under-Determined Reverberant Audio Source Separation Using a Full-Rank Spatial Covariance Model," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 7, pp. 1830–1840, Sep. 2010.

[29] Z. Koldovský, J. Málek, P. Tichavský, and F. Nesta, "Semi-Blind Noise Extraction Using Partially Known Position of the Target Source," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 10, pp. 2029–2041, Oct. 2013.

[30] A. Brendel, T. Haubner, and W. Kellermann, "Spatially guided independent vector analysis," in *submitted to: ICASSP 2020*.

[31] E. Moreau and T. Adali, *Blind identification and separation of complex-valued signals*, ser. Focus series in digital signal and image processing. London : Hoboken, NJ: ISTE ; Wiley, 2013.

[32] N. Ono, "Auxiliary-function-based independent vector analysis with power of vector-norm type weighting functions," in *Asia-Pacific Signal & Information Processing Association Annual Summit and Conference (APSIPA ASC)*, Dec. 2012.

[33] D. Kitamura, S. Mogami, Y. Mitsui, N. Takamune, H. Saruwatari, N. Ono, Y. Takahashi, and K. Kondo, "Generalized independent low-rank matrix analysis using heavy-tailed distributions for blind source separation," *EURASIP Journal on Advances in Signal Processing*, no. 1, Dec. 2018.

[34] D. R. Hunter and K. Lange, "A Tutorial on MM Algorithms," *The American Statistician*, vol. 58, no. 1, pp. 30–37, Feb. 2004.

[35] K. Matsuoka, "Minimal distortion principle for blind source separation," vol. 4. Soc. Instrument & Control Eng. (SICE), 2002, pp. 2138–2143.

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