

BLIND SEPARATION OF MULTICHANNEL SPEECH MIXTURES USING PARAFAC ANALYSIS AND INTEGER LEAST SQUARES

By
Kleanthis Mokios

SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE OF
MASTER OF SCIENCE
AT
TECHNICAL UNIVERSITY OF CRETE
CHANIA, GREECE
FEBRUARY 2006

Blind Separation of Multichannel Speech Mixtures Using PARAFAC Analysis and Integer Least Squares

Kleanthis N. Mokios, Nicholas D. Sidiropoulos, *Senior Member, IEEE*,
and Alexandros Potamianos, *Member, IEEE*

Abstract

We propose two new low-complexity frequency-domain algorithms for blind speech separation (BSS) for unknown channel order. The proposed methods consist of two separate stages: in the first stage parallel factor analysis (PARAFAC) is employed in order to separate the speech signals, while in the second stage the task of matching the arbitrary permutations in the frequency domain is performed, via a novel integer-least-squares-based method. The proposed algorithms offer guaranteed convergence and good separation performance at considerably reduced complexity relative to previous methods, when applied to real-world speech data. Our approach reveals the much broader identifiability potential of joint-diagonalization-based BSS methods, which went unrecognized in the past. The new algorithms are compared with previous BSS methods.

Index Terms

Blind speech separation, multipath channel, nonstationary signals, parallel factor analysis, integer least squares.

I. INTRODUCTION

ONE of the key problems in speech processing for teleconferencing and mobile telephony applications is that of speaker separation from multichannel measurements. In a typical teleconferencing application, multiple microphones can be deployed in a room, each providing a different linear combination of the speaker signals. In a mobile telephony application, two microphones can be deployed, e.g., one on the headset, another on the dashboard or elsewhere in a car. Such multichannel measurements provide the opportunity for speaker separation or speaker-background noise separation, which can be crucial for intelligibility. The objective of the “blind” speech separation problem is to separate the multiple speaker signals, using only the measured microphone signals, i.e., without assuming knowledge of the mixing channels.

A significant attribute of speech signals that one can exploit in order to separate multichannel speech mixtures, is their inherent non-stationarity. This attribute, that apart from speech signals is encountered in most real world signals (e.g., biological signals), has led several researchers towards developing methods for tackling the BSS problem which rely on proper utilization of the signals’ non-stationary nature. The majority of the proposed methods that exploit non-stationarity, consider the simple case of instantaneous linear mixtures [8], [11], [21]. In recent years, however, the focus of the blind source separation community has shifted towards the more general and realistic case of convolutive linear mixtures. Various approaches have been proposed, which can be divided into time-domain [12], [26], [17] and frequency-domain methods [30], [19], [18], [23], [24].

Time-domain BSS methods are applied directly to the convolutive mixture model. This approach may achieve good separation once the algorithm converges, but has the disadvantage of being computationally

The authors are with the Department of Electronic and Computer Engineering, Technical University of Crete, 73100 Chania, Greece (e-mails: kleanthis@telecom.tuc.gr, nikos@telecom.tuc.gr, potam@telecom.tuc.gr)

K. Mokios and A. Potamianos were partially supported by the IST-FP6 HIWIRE Project.

Earlier conference version of part of this work appears in Proc. ICASSP 2006, May 14-19, 2006, Toulouse, France.

expensive due to calculating many convolutions. In addition the mixing channel order need be known, which is rather impractical in real-world BSS problems.

On the other hand, frequency-domain BSS methods suggest moving to the frequency domain in order to decompose the time-domain convolutive BSS problem into multiple independent instantaneous BSS problems, one at each frequency bin. In general the frequency-domain algorithms have a simpler implementation as well as the mixing channel order need not be known, in contrast with their time domain counterparts. However, there is an inherent frequency-dependent permutation and scaling ambiguity problem in all frequency-domain BSS methods which does not exist in time-domain BSS methods.

In this paper we propose two new frequency domain BSS techniques that exploit non-stationarity of speech signals and deal with the more general case of convolutive linear mixtures. Both methods presented in this paper consist of two separate stages: the first stage performs the separation of the speech signals while the second one resolves the frequency-dependent permutation and scaling ambiguity, which is a key problem when the separation of the signals is conducted in the frequency-domain. As mentioned earlier, various frequency domain methods [30], [19], [18], [23], [24] have been proposed over the years. The methods in [30], [19], [18], however, employ a least-squares (LS) gradient descent procedure in order to separate signals while our techniques are based on the formulation of the BSS problem as a conjugate symmetric PARAFAC model that is fitted optimally, using a more efficient alternating LS algorithm that converges monotonically. The method presented in [24] is also based on an alternating LS algorithm in order to separate the speech signals, but doesn't establish the link to the conjugate-symmetric PARAFAC model. The conjugate-symmetric PARAFAC model exhibits strong identifiability properties that allows the separation of a much higher number of signals than what is suggested in any of the above methods. Regarding the frequency-dependent permutation ambiguity, e.g., in [19] the latter is dealt with by imposing a constraint on the length of the inverse-channel impulse response $\mathbf{W}(\tau)$, an approach that is loosely motivated. We propose an interpretable approach that follows naturally, by considering the physical properties of the system under study. Our approach is based on an integer-least-squares (ILS) formulation to adjust the arbitrary permutations, and its implementation assures that the occurrence of an error at any point of the permutation correction procedure, does not have a catastrophic impact on the quality of the final separated signals, unlike e.g. [30], [18], [23]. Finally, it should be stressed that our methods are characterized by low computational cost, unlike [19] or [24], which is a feature that is highly desired, especially in applications where real-time processing is essential. In [24] for example, the number of frequency bins used, for the method to attain good separation performance, is at least 4096 as apposed to 16 or 32 frequency bins in our methods. Thus, the alternating LS algorithm in [24] is iterated for at least 4064 times more, relative to our methods, which can be very costly computationally. To support our arguments we present numerical simulations using real speech data and compare the performance of our algorithm with the one proposed in [19]¹.

Our contributions in this paper can be summarized as follows:

- Improved performance relative to the method in [19].
- Complexity reduced by 1-2 orders of magnitude relative to the method in [19], bringing execution times within range for on-line application, at least for teleconference, possibly also cellular telephony applications.
- Guaranteed convergence of overall algorithm.
- Interpretable criteria for resolving the permutation-scaling ambiguities, leading to well-known ILS problem, for which many good approximate solutions exist.
- Our approach reveals the much broader identifiability potential of joint-diagonalization-based BSS methods, which went unrecognized in the past. Prior references assumed $J \geq I$, where J is the number of microphones and I the number of speakers.
- Our permutation correction scheme prevents catastrophic errors that can happen as a result of a wrongly adjusted permutation.

¹We downloaded the MATLAB implementation of the algorithm proposed in [19] from <http://newton.bme.columbia.edu/~lparra/publish/>.

- We used a much larger corpus, validated with real-world data in our simulations.

The paper is organized as follows: Section II states the BSS problem along with our assumptions and some preliminaries. Necessary background in PARAFAC theory is presented in section III. Section IV reviews the Trilinear Alternating Least Squares technique that is subsequently used to fit the PARAFAC model. The proposed method for resolving the permutation and scaling ambiguities is discussed in section V. Our algorithms are summarized in section VI. The results we obtained after application of our algorithms to real-world speech data are presented in section VII, whereas our conclusions and final remarks are presented in section VIII.

Some notation conventions that will be used in this paper follow.

$*$	Convolution operator.
$(\cdot)^*$	Complex conjugate.
$E[\cdot]$	Expectation operator.
$\lfloor x \rfloor$	Integer floor of x .
\mathbf{a}_i	i th column of \mathbf{A} .
\mathbf{A}^T	Transpose of \mathbf{A} .
\mathbf{A}^H	Conjugate transpose of \mathbf{A} .
\mathbf{A}^\dagger	Moore-Penrose pseudo-inverse of \mathbf{A} .
$\mathcal{D}_p\{\mathbf{A}\}$	Diagonal matrix constructed from the p th row of \mathbf{A} .
$\mathbf{A} \odot \mathbf{B}$	Khatri-Rao (column-wise Kronecker) product of \mathbf{A}, \mathbf{B} .
$\text{diag}\{\mathbf{x}\}$	Forms a diagonal matrix from vector \mathbf{x} .
$\ \cdot\ _F$	Frobenius norm.
$\ \cdot\ _2$	Euclidean vector norm.

II. BLIND SPEECH SEPARATION

A. Data Model and Problem Statement

Assume I mutually uncorrelated (not necessarily statistically independent) speaker signals $\mathbf{s}(t) = [s_1(t), \dots, s_I(t)]^T$. These signals are convolved and mixed in a linear medium leading to J microphone signals $\mathbf{x}(t) = [x_1(t), \dots, x_J(t)]^T$. We consider the following I -speaker J -microphone multi-input multi-output (MIMO) linear model for the received signal for the convolutive mixing problem

$$\mathbf{x}(t) = \mathbf{A} * \mathbf{s}(t) + \mathbf{n}(t) = \sum_{\tau=0}^L \mathbf{A}(\tau) \mathbf{s}(t - \tau) + \mathbf{n}(t), \quad t = 1, \dots, N \quad (1)$$

where $\mathbf{A}(\tau) = [\boldsymbol{\alpha}_1(\tau), \dots, \boldsymbol{\alpha}_I(\tau)] \in \mathbb{R}^{J \times I}$ for $\tau = 0, \dots, L$ is the mixing impulse response matrix, $\boldsymbol{\alpha}_i(\tau) = [\alpha_{1,i}(\tau), \dots, \alpha_{J,i}(\tau)]^T \in \mathbb{R}^{J \times 1}$ is the spatial signature of the i th speaker for lag τ , $\mathbf{n}(t) = [n_1(t), \dots, n_J(t)]^T$ is the additive noise vector, L is the maximum (unknown) channel length and N is the number of snapshots. The objective of the blind separation problem is to estimate the inverse-channel impulse response matrix $\mathbf{W}(\tau)$ from the observed signals $\mathbf{x}(t)$, since

$$\mathbf{s}(t) = \mathbf{W} * \mathbf{x}(t) \quad (2)$$

Following, e.g. [19], one approach towards solving the problem is to transform it into the frequency domain and to solve the joint separation problem across frequency bins.

Recall the well known property of the Discrete Fourier Transform (DFT) that allows us to express circular convolutions as products. In (1), however, we assumed a linear convolution. A linear convolution can be approximated by a circular convolution if the size T of the DFT's frame is much larger than the length of the convolution sum. In such a case we can write approximately²

$$\mathbf{x}(f, t) \approx \mathbf{A}(f) \mathbf{s}(f, t) + \mathbf{n}(f, t), \text{ for } L \ll T \quad (3)$$

²For simplicity, we use the same symbol to denote time-domain and frequency-domain representation of a signal or a filter, depending on its argument.

where $\mathbf{x}(f, t) = \sum_{\tau=0}^{T-1} \mathbf{x}(t + \tau) e^{-\frac{i2\pi f \tau}{T}}$ is the DFT of the frame of size T starting at t , $[\mathbf{x}(t), \dots, \mathbf{x}(t + T - 1)]$. Accordingly $\mathbf{s}(f, t) = \sum_{\tau=0}^{T-1} \mathbf{s}(t + \tau) e^{-\frac{i2\pi f \tau}{T}}$ and $\mathbf{A}(f) = \sum_{\tau=0}^{T-1} \mathbf{A}(\tau) e^{-\frac{i2\pi f \tau}{T}} = \sum_{\tau=0}^L \mathbf{A}(\tau) e^{-\frac{i2\pi f \tau}{T}}$, because $\mathbf{A}(\tau) = \mathbf{0}^{J \times I}$ for $\tau > L$. The i th column of $\mathbf{A}(f)$ represents now the spatial signature of the i th speaker in the frequency domain, at frequency f .

Before proceeding further, we list our main assumptions:

Assumption 2.1: The speaker signals $\mathbf{s}(t)$ are zero mean, second-order quasi-stationary signals; i.e. the variances of the signals are slowly varying with time such that over short time intervals they can be assumed approximately stationary.

Assumption 2.2: The number of speakers is known.

Assumption 2.3: The contribution of the noise term $\mathbf{n}(t)$ is negligible compared to the contribution of the speaker signals. In our context, this is an accurate approximation in many real-world cases. When noise is not negligible, its power can be estimated from silence periods and subtracted from correlation matrix estimates, which is the only place where noise comes into play in our algorithms.

B. Direct Channel Estimation versus Inverse Channel Estimation

• **Direct channel estimation (DCE) approach:** Equation (3) dictates that the BSS problem could have been solved, had we had in our disposal the frequency domain counterparts of the multipath channel matrices $\mathbf{A}(\tau)$, that is $\mathbf{A}(f)$, for all frequency bins. Provided that $\text{rank}[\mathbf{A}(f)] = I$, by taking the left Moore-Penrose pseudo-inverse of each frequency's corresponding matrix $\mathbf{A}(f)$, and applying the Inverse DFT to the collection of the acquired pseudo-inverses $\{\mathbf{A}^\dagger(f), f = 0, \dots, T - 1\}$, we could determine estimates of the inverse-channel matrices $\{\mathbf{W}(\tau), \tau = 0, \dots, T - 1\}$, and subsequently estimate the speaker signals via (2). The same could be achieved by computing $\hat{\mathbf{s}}(f, t) \approx \mathbf{A}^\dagger(f) \mathbf{x}(f, t)$ for all frequency bins, and applying Inverse DFT to $\{\hat{\mathbf{s}}(f, t), f = 0, \dots, T - 1\}$. When $\text{rank}[\mathbf{A}(f)] < I$ (in particular, when $J < I$) perfect separation is not possible; substantial reduction of crosstalk is still possible, however, using matched filtering to the columns of $\mathbf{A}(f)$, or more sophisticated array processing methods, like Capon beamforming. Therefore, the blind speech separation problem boils down to the problem of estimating the matrices $\{\mathbf{A}(f), f = 0, \dots, T - 1\}$.

Let us focus on a time interval over which the measured signals can be assumed stationary (see assumption 2.1). Furthermore we consider one specific frequency bin f , ignoring for the time being the rest $T - 1$ frequencies. The autocorrelation matrix of the vector of microphone outputs at frequency f is then

$$\mathbf{R}_x(f, t) = E[\mathbf{x}(f, t) \mathbf{x}^H(f, t)] \approx \mathbf{A}(f) E[\mathbf{s}(f, t) \mathbf{s}^H(f, t)] \mathbf{A}^H(f) = \mathbf{A}(f) \mathbf{D}_s(f, t) \mathbf{A}^H(f) \quad (4)$$

Since we assume mutually uncorrelated speaker signals we postulate diagonal autocorrelation matrix $\mathbf{D}_s(f, t)$, while by assumption 2.3 the autocorrelation matrix of the noise vector $\mathbf{n}(f, t)$ has been neglected. As we will see in Section III, proper processing of the autocorrelation data (4) via PARAFAC analysis, for all frequency bins and all intervals over which the measured signals are assumed stationary, enables us to specify the matrices $\{\mathbf{A}(f), f = 0, \dots, T - 1\}$ up to an individual permutation and scaling of their columns (spatial signatures), that vary arbitrarily from one frequency bin to another.

• **Inverse channel estimation (ICE) approach:** Instead of first estimating and then inverting the forward-channel matrices $\{\mathbf{A}(f), f = 0, \dots, T - 1\}$, we can directly estimate the inverse-channel matrices in the frequency domain, $\{\mathbf{W}(f), f = 0, \dots, T - 1\}$, perform the Inverse DFT to obtain $\{\mathbf{W}(\tau), \tau = 0, \dots, T - 1\}$ and subsequently estimate the speaker signals via (2), or, again we can compute $\hat{\mathbf{s}}(f, t) \approx \mathbf{W}(f) \mathbf{x}(f, t)$ for all frequency bins, and then apply Inverse DFT to $\{\hat{\mathbf{s}}(f, t), f = 0, \dots, T - 1\}$.

Assume, for the moment, that the number of microphones equals the number of the speakers ($J = I$) and $\text{rank}\{\mathbf{A}(f)\} = I$. In this case, the inverse of the autocorrelation matrix $\mathbf{R}_x(f, t)$ in (4) can be written as

$$\mathbf{R}_x^{-1}(f, t) = [\mathbf{A}(f) \mathbf{D}_s(f, t) \mathbf{A}^H(f)]^{-1} = [\mathbf{A}^{-1}(f)]^H \mathbf{D}_s^{-1}(f, t) \mathbf{A}^{-1}(f) = \mathbf{W}^H(f) \mathbf{D}_s^{-1}(f, t) \mathbf{W}(f) \quad (5)$$

where the last equality follows from the fact that the inverse-channel matrix $\mathbf{W}(f)$ equals the inverse of the forward-channel matrix $\mathbf{A}^{-1}(f)$, in this case. Employing parallel factor analysis of the “inverse” autocorrelation data (5), for all frequency bins and all intervals over which the speaker signals are assumed stationary, we may obtain the beamforming matrices $\{\mathbf{W}(f), f = 0, \dots, T-1\}$, up to a frequency-dependent scaling and permutation ambiguity of their rows.

In the case of fewer speakers than microphones ($J > I$), however, the inversion of the autocorrelation matrix $R_x(f, t)$ cannot be performed, and therefore PARAFAC analysis cannot be applied. In order to deal with this problem, we reduce our model for the received signal back to the $J = I$ case, by resorting to an appropriate dimension-reducing prefiltering of the data in (3). The following paragraphs describe how this is done.

Consider one specific frequency bin f , ignoring for the time being the rest $T-1$ frequencies. Let $N_T = \lfloor \frac{N}{T} \rfloor$ and $\mathbf{X}(f) = [\mathbf{x}(f, 0), \mathbf{x}(f, T), \dots, \mathbf{x}(f, T(N_T - 1))]$, then from (3), the $J \times N_T$ matrix $\mathbf{X}(f)$ has the following factorization

$$\mathbf{X}(f) = \mathbf{A}(f)\mathbf{S}(f) \quad (6)$$

where $\mathbf{S}(f) = [\mathbf{s}(f, 0), \mathbf{s}(f, T), \dots, \mathbf{s}(f, T(N_T - 1))] \in \mathbb{C}^{I \times N_T}$, and the noise term has been neglected in accordance with assumption 2.3. We choose $\mathbf{F}(f)$ to be the matrix whose columns are the I dominant eigenvectors which correspond to the eigenvalue decomposition (ED) of the autocorrelation matrix $\bar{\mathbf{R}}_x(f) = \frac{1}{N_T} \mathbf{X}(f)\mathbf{X}^H(f)$. It is known that $\mathbf{F}(f)$ is an asymptotically unbiased estimate of $\mathbf{U}_A(f)$ where $\mathbf{A}(f) = \mathbf{U}_A(f)\mathbf{\Sigma}_A(f)\mathbf{V}_A(f)$ is the economy-size singular value decomposition (SVD) of $\mathbf{A}(f)$.

Let us now consider a prefiltered version of the data in (6), using the matrix $\mathbf{F}^H(f)$ as the filtering matrix, i.e.,

$$\tilde{\mathbf{X}}(f) \triangleq \mathbf{F}^H(f)\mathbf{X}(f) = \tilde{\mathbf{A}}(f)\mathbf{S}(f) \quad (7)$$

with $\tilde{\mathbf{A}}(f) \triangleq \mathbf{F}^H(f)\mathbf{A}(f) \in \mathbb{C}^{I \times I}$. By properties of the SVD, $\tilde{\mathbf{A}}(f)$ is a nonsingular $I \times I$ matrix. Hence, the autocorrelation matrix $\tilde{\mathbf{R}}_x(f)$ that corresponds to the prefiltered data model in (7), can be inverted to yield the following factorization

$$\tilde{\mathbf{R}}_x^{-1}(f, t) = [\tilde{\mathbf{A}}(f)\mathbf{D}_s(f, t)\tilde{\mathbf{A}}^H(f)]^{-1} = [\tilde{\mathbf{A}}^{-1}(f)]^H \mathbf{D}_s^{-1}(f, t) \tilde{\mathbf{A}}^{-1}(f) = \tilde{\mathbf{W}}^H(f)\mathbf{D}_s^{-1}(f, t)\tilde{\mathbf{W}}(f) \quad (8)$$

with $\tilde{\mathbf{W}}(f) = \tilde{\mathbf{A}}^{-1}(f)$. Employing parallel factor analysis of the “inverse” autocorrelation data (8), for all frequency bins and all intervals over which the speaker signals are assumed stationary, we may obtain the matrices $\{\tilde{\mathbf{W}}(f), f = 0, \dots, T-1\}$, up to a frequency-dependent scaling and permutation ambiguity of their rows. In this case, the beamforming matrices $\{\mathbf{W}(f), f = 0, \dots, T-1\}$ can be written as

$$\mathbf{W}(f) = \tilde{\mathbf{W}}(f)\mathbf{F}^H(f) \quad (9)$$

Note that the ICE approach cannot be employed when we have fewer microphones than sensors ($J < I$). In this case we have to resort to the DCE approach.

Of course, the beamforming matrices $\{\mathbf{W}(f), f = 0, \dots, T-1\}$ in the ICE approach, just as the acquired mixing matrices $\{\mathbf{A}(f), f = 0, \dots, T-1\}$ in the DCE approach, are specified up to a frequency-dependent individual permutation and scaling of their columns. This constitutes a serious problem since only consistent permutations and scaling across frequencies will correctly reconstruct the source signals. Hence a second major problem arises, that of resolving the frequency-dependent permutation and scaling issue. The following section reviews the basics of PARAFAC analysis, and shows how it is employed on the BSS problem, via our two approaches.

III. PARALLEL FACTOR ANALYSIS

The present section is devoted to the presentation of the necessary background in PARAFAC theory, and also shows how blind speech separation comes under this category of problems. We only consider the model in (4), since the analysis for the model of (8) follows by substituting $\tilde{\mathbf{R}}_x^{-1}(f, t)$, $\tilde{\mathbf{W}}^H(f)$ and $\mathbf{D}_s^{-1}(f, t)$ for $\mathbf{R}_x(f, t)$, $\mathbf{A}(f)$ and $\mathbf{D}_s(f, t)$, respectively.

A. PARAFAC Model

Let us divide the whole data block of N snapshots into P sub-blocks, so that each sub-block contains $N_P = \lfloor \frac{N}{P} \rfloor$ snapshots. We assume that each sub-block corresponds to a time interval over which the speaker signals are stationary, in accordance with assumption 2.1. Under this framework, the measured snapshots within any p th sub-block correspond to the following autocorrelation matrix

$$\mathbf{R}_x(f, t_p) = \mathbf{A}(f) \mathbf{D}_s(f, t_p) \mathbf{A}^H(f), f = 0, \dots, T - 1 \quad (10)$$

where $\mathbf{D}_s(f, t_p)$ is the diagonal autocorrelation matrix of the speaker signals in the p th sub-block. Using all P sub-blocks, we will have P different autocorrelation matrices $\{\mathbf{R}_x(f, t_1), \dots, \mathbf{R}_x(f, t_P)\}$ for each frequency. We observe that for each frequency, these matrices differ from each other only because the source signal autocorrelation matrices $\mathbf{D}_s(f, t_p)$ differ from one sub-block to another.

Let us stack the P matrices $\{\mathbf{R}_x(f, t_p), p = 1, \dots, P\}$ together to form a three-way array $\underline{\mathbf{R}}_x(f)$. The (j, l, p) th element of such an array can be written as

$$r_{j,l,p}(f) \triangleq [\underline{\mathbf{R}}_x(f)]_{j,l,p} = \sum_{i=1}^I \alpha_{j,i}(f) v_i(f, p) \alpha_{l,i}^*, f = 0, \dots, T - 1 \quad (11)$$

where $v_i(f, p) \triangleq [\mathbf{D}_s(f, t_p)]_{i,i}$ is the power-spectral-density of the i th source in the p th sub-block. Defining the matrices $\mathbf{P}(f) \in \mathbb{C}^{P \times I}$, $f = 0, \dots, T - 1$ as

$$\mathbf{P}(f) \triangleq \begin{bmatrix} v_1(f, 1) & \dots & v_I(f, 1) \\ \vdots & \ddots & \vdots \\ v_1(f, P) & \dots & v_I(f, P) \end{bmatrix} \quad (12)$$

we can write the following relationship between $\mathbf{D}_s(f, t_p)$ and $\mathbf{P}(f)$

$$\mathbf{D}_s(f, t_p) = \mathcal{D}_p\{\mathbf{P}(f)\} \quad (13)$$

for all $p = 1, \dots, P$ and all $f = 0, \dots, T$.

Equation (11) implies that $r_{j,l,p}(f)$ is a sum of rank-1triple products; this equation is known as (conjugate-symmetric) *parallel factor* (PARAFAC) analysis of $r_{j,l,p}$ [9], [13], [28], [29]. If I is sufficiently small (11) represents a low-rank decomposition of $\underline{\mathbf{R}}_x(f)$. Therefore, the problem of estimating the matrix $\mathbf{A}(f)$ for a specific frequency f can be reformulated as the problem of low-rank decomposition of the three-way autocorrelation array $\underline{\mathbf{R}}_x(f)$. By solving a similar problem separately for every frequency we obtain the entire collection of the frequency-domain mixing matrices $\{\mathbf{A}(f), f = 0, \dots, T - 1\}$, up to certain inherent indeterminacies, discussed next.

B. Identifiability

By identifiability we mean the uniqueness (up to inherently unresolvable source permutation and scale ambiguities) of all speaker spatial signatures at a given frequency, given the exact frequency-domain autocorrelation data at that frequency.

In this subsection, we briefly review identifiability properties of the PARAFAC-model based spatial signature estimation. Towards this end, we discuss conditions under which the trilinear decomposition of $\underline{\mathbf{R}}_x(f)$ is unique. Identifiability conditions on the number of sub-blocks and the number of microphones

are presented. In the following we consider the frequency index fixed, and omit it for convenience since the results presented do not depend on the frequency f .

We start with the definition of the Kruskal rank of a matrix [13].

Definition: The Kruskal rank (or k -rank) of a matrix \mathbf{C} is k_C if and only if every k_C columns of \mathbf{C} are linearly independent, and either \mathbf{C} has k_C columns or \mathbf{C} contains a set of $k_C + 1$ linearly dependent columns. Note that k -rank is always less than or equal to the conventional matrix rank. It can be easily checked that if \mathbf{C} is full column rank, then it is also full k -rank.

Using (13), we can rewrite (10) as

$$\mathbf{R}_x(t_p) = \mathbf{A} \mathcal{D}_p \{\mathbf{P}\} \mathbf{A}^H, p = 1, \dots, P \quad (14)$$

To establish identifiability, we have to obtain under which conditions matrices \mathbf{P} and \mathbf{A} are the unique (up to the scaling and permutation ambiguities) matrices that give rise to the data $\{\mathbf{R}_x(t_p), p = 1, \dots, P\}$ of (14). In the context of our present application, which involves a conjugate-symmetric PARAFAC model the following theorem follows from [13] (also cf. [25]).

Theorem 1: Consider the set of matrices (14). For $I > 1$, if

$$k_A + k_P + k_{A^*} = 2k_A + k_P \geq 2I + 2 \quad (15)$$

then \mathbf{A} and \mathbf{P} are unique up to inherently unresolvable permutation and scaling of columns, i.e. if there exists any other pair $\{\bar{\mathbf{A}}, \bar{\mathbf{P}}\}$ which satisfies (15), then this pair is related to the pair $\{\mathbf{A}, \mathbf{P}\}$ via

$$\bar{\mathbf{A}} = \mathbf{A} \mathbf{\Pi} \mathbf{\Delta}_1, \bar{\mathbf{P}} = \mathbf{P} \mathbf{\Pi} \mathbf{\Delta}_2 \quad (16)$$

where $\mathbf{\Pi}$ is a permutation matrix, and $\mathbf{\Delta}_1, \mathbf{\Delta}_2$ are diagonal scaling matrices satisfying

$$\mathbf{\Delta}_1 \mathbf{\Delta}_1^* \mathbf{\Delta}_2 = \mathbf{I} \quad (17)$$

For $I = 1$, \mathbf{A} and \mathbf{P} are always unique, irrespective of (15).

It is worth noting that condition (15) is sufficient but not necessary for identifiability [2]. For $I > 1$ the condition $k_P \geq 2$ becomes necessary [10]. In terms of the number of sub-blocks the latter condition requires that $P \geq 2$.

Alternative uniqueness conditions can be derived if random component matrices are considered, giving rise to the concept of almost sure-identifiability. The best known result on almost sure-identifiability conditions for the conjugate-symmetric PARAFAC model has been established in [31] and is presented in the next theorem

Theorem 2: Suppose that the elements of \mathbf{A} and \mathbf{P} are drawn from a jointly continuous distribution. If

$$\frac{I(I-1)}{2} \leq \frac{J(J-1)}{4} \left[\frac{J(J-1)}{2} + 1 \right] - \binom{J}{4} 1_{\{J \geq 4\}} \quad (18)$$

where $1_{\{J \geq 4\}} = \begin{cases} 0 & \text{if } J < 4 \\ 1 & \text{if } J \geq 4 \end{cases}$, then for $I \leq 30$, \mathbf{A} and \mathbf{P} are almost surely unique up to inherently unresolvable permutation and scaling of columns.

In our context J corresponds to the number of microphones we use while I corresponds to the number of speakers that can be separated. In Table I below, we give for values $J = 2, \dots, 8$ the upper bound for I , according to the formula given in (18). Observe that 6 microphones it is possible to estimate the steering vectors of up to 15 speakers, whereas with 8 microphones it is possible to estimate the steering vectors of up to 26 speakers - more than 3 times the number of microphones, in theory.

Note that the formula in (18) provides theoretical bounds on the number of speaker signals that can be separated, relying on the link to the conjugate symmetric PARAFAC model established in the previous subsection. The ICE approach, however, requires the number of speakers to be at most equal to the number of microphones used. Thus, practically in the ICE approach, the maximum number of speakers that can be separated, cannot surpass the number of microphones used to achieve this separation, regardless of the bounds resulting from (18) and shown in Table I.

TABLE I
VALUES OF J AND UPPER BOUNDS FOR I , WHICH FOLLOW FROM (18)

J	2	3	4	5	6	7	8
Upper bound for I	2	4	6	10	15	20	26

IV. TRILINEAR ALTERNATING LEAST SQUARES ESTIMATOR

We will now review the Trilinear Alternating Least Squares (TALS) technique [10], [28], [29] for fitting the PARAFAC model of Section III. Again, we solely consider the PARAFAC model of the DCE approach (5), since by letting $\tilde{\mathbf{x}}(f, t)$, $\tilde{\mathbf{R}}_x^{-1}(f, t)$, $\tilde{\mathbf{W}}^H(f)$ and $\tilde{\mathbf{D}}_s^{-1}(f, t)$ stand in for $\mathbf{x}(f, t)$, $\mathbf{R}_x(f, t)$, $\mathbf{A}(f)$ and $\mathbf{D}_s(f, t)$ respectively, the following analysis is identical for the model of the ICE approach. TALS will be subsequently used to estimate the matrices $\mathbf{A}(f)$ and $\tilde{\mathbf{W}}^H(f)$ up to a frequency-dependent permutation and scaling ambiguity. Exactly how these ambiguities are resolved (to yield our final separation solution), will be discussed in section V.

A. Preliminaries

In practice, the exact autocorrelation matrices $\mathbf{R}_x(f, t_p)$ are unavailable but can be estimated from the array snapshots $\mathbf{x}(t)$, $t = 1, \dots, N$. If we define $K = \lfloor \frac{N_P}{T} \rfloor$, the sample autocorrelation matrix estimates are given by

$$\hat{\mathbf{R}}_x(f, t_p) = \frac{1}{K} \sum_{k=0}^{K-1} \mathbf{x}(f, t_p + kT) \mathbf{x}^H(f, t_p + kT), p = 1, \dots, P \quad (19)$$

We can choose t_p such that we have non-overlapping times for $\hat{\mathbf{R}}_x(f, t_p)$, i.e. $t_p = pTK$, but overlapping averaging times can be chosen, as well. The matrices $\hat{\mathbf{R}}_x(f, t_p)$, $p = 1, \dots, P$ can be used to form a sample three-way autocorrelation array denoted as $\underline{\hat{\mathbf{R}}}_x(f)$.

In order to explain TALS, we will need to define “slices” of the three-way arrays $\underline{\mathbf{R}}_x(f)$ and $\underline{\hat{\mathbf{R}}}_x(f)$ along different dimensions [29]. Towards this end, let us define the “slice” matrices³

$$\mathbf{R}_\alpha^{(j)}(f) \triangleq [r_{j,:,:}(f)], \mathbf{R}_b^{(l)}(f) \triangleq [r_{:,l,:}(f)], \mathbf{R}_c^{(p)}(f) \triangleq [r_{:,:,p}(f)] \quad (20)$$

$$\hat{\mathbf{R}}_\alpha^{(j)}(f) \triangleq [\hat{r}_{j,:,:}(f)], \hat{\mathbf{R}}_b^{(l)}(f) \triangleq [\hat{r}_{:,l,:}(f)], \hat{\mathbf{R}}_c^{(p)}(f) \triangleq [\hat{r}_{:,:,p}(f)] \quad (21)$$

where $j, l = 1, \dots, J$; $p = 1, \dots, P$; $r_{j,l,p}(f) \triangleq [\underline{\mathbf{R}}_x(f)]_{j,l,p}$; and $\hat{r}_{j,l,p}(f) \triangleq [\underline{\hat{\mathbf{R}}}_x(f)]_{j,l,p}$.

For the sake of convenience, let us introduce $\mathbf{B}(f) \triangleq \mathbf{A}^H(f)$. Define the matrices

$$\mathbf{R}_\alpha(f) \triangleq \begin{bmatrix} \mathbf{R}_\alpha^{(1)}(f) \\ \mathbf{R}_\alpha^{(2)}(f) \\ \vdots \\ \mathbf{R}_\alpha^{(J)}(f) \end{bmatrix} = [\mathbf{P}(f) \odot \mathbf{A}(f)] \mathbf{B}(f) \quad (22)$$

$$\mathbf{R}_b(f) \triangleq \begin{bmatrix} \mathbf{R}_b^{(1)}(f) \\ \mathbf{R}_b^{(2)}(f) \\ \vdots \\ \mathbf{R}_b^{(J)}(f) \end{bmatrix} = [\mathbf{B}^T(f) \odot \mathbf{P}(f)] \mathbf{A}^T(f) \quad (23)$$

³We use MATLAB notation.

$$\mathbf{R}_c(f) \triangleq \begin{bmatrix} \mathbf{R}_c^{(1)}(f) \\ \mathbf{R}_c^{(2)}(f) \\ \vdots \\ \mathbf{R}_c^{(P)}(f) \end{bmatrix} = [\mathbf{A}(f) \odot \mathbf{B}^T(f)] \mathbf{P}^T(f) \quad (24)$$

and similarly for the updates of their respective sample estimates $\hat{\mathbf{R}}_\alpha(f)$, $\hat{\mathbf{R}}_b(f)$, $\hat{\mathbf{R}}_c(f)$.

Note that for the sake of algorithm simplicity, we will not exploit the fact that our specific PARAFAC model is symmetric. That is, the algorithm that follows treats $\mathbf{A}(f)$ and $\mathbf{B}(f)$ as independent variables; symmetry will only be exploited in the calculation of the final estimate of $\mathbf{A}(f)$.

B. Algorithm

The basic idea behind the TALS procedure for PARAFAC fitting is to update each time a subset of parameters using LS regression while keeping the previously obtained estimates of the rest of parameters fixed. This alternating projections-type procedure is iterated for all subsets of parameters until convergence is achieved [9], [3], [28], [29].

In application to our problem, the PARAFAC TALS procedure can be formulated as follows

- **Step 1:** Initialize $\mathbf{P}(f)$ and $\mathbf{A}(f)$.
- **Step 2:** Find the estimate $\mathbf{B}(f)$ by solving the following LS problem

$$\hat{\mathbf{B}}(f) = \arg \min_{\mathbf{B}(f)} \|\hat{\mathbf{R}}_\alpha(f) - [\mathbf{P}(f) \odot \mathbf{A}(f)] \mathbf{B}(f)\|_F^2 \quad (25)$$

whose solution is given by

$$\hat{\mathbf{B}}(f) = [\mathbf{P}(f) \odot \mathbf{A}(f)]^\dagger \hat{\mathbf{R}}_\alpha(f) \quad (26)$$

Set $\mathbf{B}(f) = \hat{\mathbf{B}}(f)$.

- **Step 3:** Find the estimate of $\mathbf{A}(f)$ by solving the following LS problem

$$\hat{\mathbf{A}}(f) = \arg \min_{\mathbf{A}(f)} \|\hat{\mathbf{R}}_b(f) - [\mathbf{B}^T(f) \odot \mathbf{P}(f)] \mathbf{A}^T(f)\|_F^2 \quad (27)$$

whose solution is given by

$$\hat{\mathbf{A}}(f) = \hat{\mathbf{R}}_b^T(f) \left[[\mathbf{B}^T(f) \odot \mathbf{P}(f)]^\dagger \right]^T \quad (28)$$

Set $\mathbf{A}(f) = \hat{\mathbf{A}}(f)$.

- **Step 4:** Find the estimate of $\mathbf{P}(f)$ by solving the following LS problem

$$\hat{\mathbf{P}}(f) = \arg \min_{\mathbf{P}(f)} \|\hat{\mathbf{R}}_c(f) - [\mathbf{A}(f) \odot \mathbf{B}^T(f)] \mathbf{P}^T(f)\|_F^2 \quad (29)$$

whose solution is given by

$$\hat{\mathbf{P}}(f) = \hat{\mathbf{R}}_c^T(f) \left[[\mathbf{A}(f) \odot \mathbf{B}^T(f)]^\dagger \right]^T \quad (30)$$

Set $\mathbf{P}(f) = \hat{\mathbf{P}}(f)$.

- **Step 5:** Repeat steps 2-4 until convergence is achieved and then compute the final estimate of $\mathbf{A}(f)$ as $\hat{\mathbf{A}}(f) = [\mathbf{A}(f) + \mathbf{B}^H(f)]/2$.

The TALS algorithm is repeated for every frequency. The complexity of the TALS algorithm is $\mathcal{O}(I^3 + J^2IP)$ ($\mathcal{O}(I^3 + I^3P)$ for the model of the ICE approach) per iteration. It is worth noting that when I is small relative to J and P , only a few iterations of the algorithm are usually required to achieve convergence [28]. The TALS algorithm can be initialized randomly, or via algebraic eigenvalue-based solutions, similar to ESPRIT, in certain cases; see [3].

V. RESOLVING THE SCALING AND PERMUTATION AMBIGUITIES

As already mentioned, in order to achieve adequate separation performance, the resolution of the frequency-dependent permutation and scaling ambiguity problems must be addressed. In this section, we suggest a novel method for solving the permutation and scaling problems that arise in the blind speech separation context. In the first two subsections, our method of dealing with the frequency-dependent permutation ambiguities is presented, while the third subsection describes how the scaling problem is resolved.

A column-permuted and scaled version of the true mixing matrix $\mathbf{A}(f)$, can be written as

$$\mathbf{A}_{s,p}(f) = \mathbf{A}_s(f)\mathbf{\Pi}(f) = \mathbf{A}(f)\mathbf{D}(f)\mathbf{\Pi}(f) \quad (31)$$

where $\mathbf{\Pi}(f)$ is a permutation matrix and $\mathbf{D}(f)$ is a diagonal scaling matrix. This is the output matrix of the TALS algorithm at frequency f , when following the DCE approach.

Similarly, a column-permuted and scaled version of the matrix $\mathbf{U}(f) = \mathbf{W}^H(f)$, where $\mathbf{W}(f)$ is the true beamforming matrix, can be written as

$$\mathbf{U}_{s,p}(f) = \mathbf{U}_s(f)\mathbf{\Pi}(f) = \mathbf{U}(f)\mathbf{D}(f)\mathbf{\Pi}(f) \quad (32)$$

This is the output matrix of the TALS algorithm at frequency f , when following the ICE approach.

In order to separate the individual speech signals, it is imperative to reduce this frequency-dependent permutation and scaling to a single, frequency-independent permutation and scaling. How this is done, is discussed next.

A. Resolving the frequency-dependent permutation problem

First, let us consider the data model of the first approach. We resolve the permutation problem by exploiting two inherent attributes of the system under study. Proper interpretation and utilization of each of these attributes allows us to transform the frequency dependent permutation problem into an optimization problem with each attribute providing us with a different optimization criterion. We then formulate a joint optimization criterion that incorporates these two constituent criteria.

•**Criterion 1:** Based on the underlying physical model of the electro-acoustic $J \times I$ (microphones \times speakers) system, the elements of the mixing matrices $\{\mathbf{A}(f), f = 0, \dots, T-1\}$ may be decomposed into products of the form

$$\alpha_{j,i}(f) = G_{j,i}(f)H_j(f)e^{-\frac{2\pi\sqrt{-1}f}{T}\tau_{j,i}} \quad (33)$$

where

- $G_{j,i}(f)$ represents the frequency response of the acoustic channel between the j th sensor and the i th source.
- $H_j(f)$ represents the j th microphone's frequency response.
- $\tau_{j,i}$ represents the delay, due to the time needed for the sound to cover the distance between the i th speaker and the j th microphone.

For this subsection only (derivation of criterion 1), we make the following assumptions:

Assumption 5.1: The speech signals travel in an almost non-reverberant (unechoic) environment. This means that, aside from a real multiplicative factor that is associated with the power loss over distance of the signals, the frequency responses of the acoustic channels are almost identical. Hence, $G_{j,i}(f) \approx n_{j,i}G(f)$, where $n_{j,i}$ is the aforementioned multiplicative factor, that depends solely on the distance between the i th speaker and the j th microphone.

Assumption 5.2: The microphones' frequency responses are almost flat over the frequency band of interest (i.e., 20Hz-3.4KHz approximately, which is the frequency band of human speech). Hence, $H_j(f) \approx H_j$ for all $j = 1, \dots, J$, with $H_j \in \mathbb{C}$.

Under the assumptions made, (33) can be rewritten as

$$\alpha_{j,i}(f) \approx n_{j,i}G(f)H_je^{-\frac{2\pi\sqrt{-1}f}{T}\tau_{j,i}} \quad (34)$$

Now, normalize the spatial signatures of the output matrices (31) with respect to an arbitrary reference microphone (e.g., the first microphone for $j = 1$). Consider the vector whose elements are the magnitudes of the respective elements of the i th normalized spatial signature. Using (34), the latter can be written as

$$\boldsymbol{\alpha}_i^{(n)}(f) \approx \begin{bmatrix} 1 \\ \frac{n_{2,i}|H_2|}{n_{1,i}|H_1|} \\ \vdots \\ \frac{n_{J,i}|H_J|}{n_{1,i}|H_1|} \end{bmatrix} \quad (35)$$

for all $i = 1, \dots, I$.

Hence, in non-reverberant (or, mildly reverberant) environments, where the recordings have been conducted using microphones with (almost) flat frequency responses, the magnitude of the elements of the normalized spatial signature, is approximately independent of frequency, as is readily seen from (35). Thus, the collection of vectors $\{\boldsymbol{\alpha}_i^{(n)}(f), f = 0, \dots, T-1, i = 1, \dots, I\}$, can be divided into I separate clusters, each of the latter being associated with a different speaker. Each cluster consists of T ressemblant (in the Euclidean distance sense) vectors $\{\boldsymbol{\alpha}_i^{(n)}(f), f = 0, \dots, T-1\}$, i.e., the vectors (35) which correspond to the same specific speaker (fixed i).

We apply a vector quantization (VQ) clustering procedure [7], over all the available vectors $\{\boldsymbol{\alpha}_i^{(n)}(f), f = 0, \dots, T-1, i = 1, \dots, I\}$, in order to determine the I center vectors (centroids) $\{\mathbf{c}_1, \dots, \mathbf{c}_I\}$, with $\mathbf{c}_i \in \mathbb{R}^J$, of the I clusters. With the I centroids at our disposal, the frequency-dependent permutation problem boils down to the following Integer Least Squares (ILS) minimization problem

$$\min_{\{s_f\}_{f=0}^{T-1} \in \{1, \dots, I\}^T} \sum_{f=0}^{T-1} \|\mathbf{C} - \mathbf{A}_n(f)\mathbf{\Pi}_1 1_{\{s_f=1\}} - \dots - \mathbf{A}_n(f)\mathbf{\Pi}_{I!} 1_{\{s_f=I!\}}\|_F^2 \quad (36)$$

where $\mathbf{C} \in \mathbb{R}^{J \times I}$ is the matrix whose columns are the centroids $\{\mathbf{c}_i, i = 1, \dots, I\}$, $\mathbf{A}_n(f) \in \mathbb{R}^{J \times I}$ is the matrix whose columns are the arbitrarily permuted⁴ vectors $\{\boldsymbol{\alpha}_i^{(n)}(f), i = 1, \dots, I\}$, $\mathbf{\Pi}_l \in \mathbb{R}^{I \times I}$ for $l = 1, \dots, I!$ is one of the total number of⁵ $I!$ permutation matrices of dimensionality $I \times I$, $1_{\{x=\alpha\}} = \begin{cases} 0 & \text{if } x \neq \alpha \\ 1 & \text{if } x = \alpha \end{cases}$, and s_f takes on the integer values $\{1, \dots, I!\}$.

The solution $\{s_1, \dots, s_T\}$ to the minimization problem (36) provides us with all the information we need for dealing with the frequency-dependent permutation problem. Suppose that $s_f = l$ with $l \in \{1, \dots, I!\}$. Then, by permuting the columns of $\mathbf{A}_{s,p}(f)$ in accordance with the permutation matrix $\mathbf{\Pi}_l$, the permutation ambiguity at frequency f is lifted. Thus, based on the values $\{s_f, f = 1, \dots, T\}$ that criterion (36) provides, and performing the appropriate permutations across all frequency bins, the frequency-dependent permutation problem can be adequately resolved.

•**Criterion 2:** For adjacent frequency bins, $f_1, f_2 = f_1 + 1$, $\boldsymbol{\alpha}_i(f_1) \approx \boldsymbol{\alpha}_i(f_2)$ and consequently $\boldsymbol{\alpha}_i^{(n)}(f_1) \approx \boldsymbol{\alpha}_i^{(n)}(f_2)$. This is valid for dense FFT grids irrespective of whether or not assumptions 5.1 and 5.2 hold. Based on this fact, we formulate our second ILS minimization criterion

$$\min_{\{s_f\}_{f=0}^{T-1} \in \{1, \dots, I!\}^T} \sum_{f=0}^{T-1} \|\mathbf{A}_n(f)\mathbf{\Pi}_1 1_{\{s_f=1\}} + \dots + \mathbf{A}_n(f)\mathbf{\Pi}_{I!} 1_{\{s_f=I!\}} - \mathbf{A}_n(f-1)\mathbf{\Pi}_1 1_{\{s_{f-1}=1\}} - \dots - \mathbf{A}_n(f-1)\mathbf{\Pi}_{I!} 1_{\{s_{f-1}=I!\}}\|_F^2 \quad (37)$$

where $\mathbf{A}_n(-1) = \mathbf{0}^{J \times I}$, by convention. Criterion (37) should also provide adequate resolution of the permutation problem, by reordering the columns of the matrices $\mathbf{A}_{s,p}(f)$ of (31), as outlined in the previous paragraph for our first criterion.

⁴in accordance with the arbitrary permutations of the columns of the output matrices $\{\mathbf{A}_{s,p}(f), f = 0, \dots, T-1\}$

⁵ $I!$ equals the total number of possible permutations of the columns or the rows of an $I \times I$ matrix or equivalently the total number of permutation matrices of dimensionality $I \times I$

The criteria of (36) and (37) can be combined into one overall ILS minimization criterion

$$\min_{\{s_f\}_{f=0}^{T-1} \in \{1, \dots, I\}^T} \sum_{f=0}^{T-1} \left\{ \left\| \mathbf{C} - \mathbf{A}_n(f) \mathbf{\Pi}_1 1_{\{s_f=1\}} - \dots - \mathbf{A}_n(f) \mathbf{\Pi}_{I!} 1_{\{s_f=I!\}} \right\|_F^2 + \lambda \left\| \mathbf{A}_n(f) \mathbf{\Pi}_1 1_{\{s_f=1\}} + \dots + \mathbf{A}_n(f) \mathbf{\Pi}_{I!} 1_{\{s_f=I!\}} - \mathbf{A}_n(f-1) \mathbf{\Pi}_1 1_{\{s_{f-1}=1\}} - \dots - \mathbf{A}_n(f-1) \mathbf{\Pi}_{I!} 1_{\{s_{f-1}=I!\}} \right\|_F^2 \right\} \quad (38)$$

where λ is a weighting factor. When $\lambda \rightarrow 0$ (practically, for $\lambda \ll 1$), the combined criterion (38) coincides with criterion 1, while for $\lambda \rightarrow \infty$ (practically, for $\lambda \gg 1$), (38) coincides with criterion 2.

Note that the ILS minimization criterion (38), fixes the permutation ambiguity globally across all frequency bins. This is an important advantage versus a sequential approach to the permutation problem, where starting from the first frequency bin f_1 , we adjust the permutation of each bin relative to its previous bin. This is because the sequential approach, although simpler, has a major drawback as explained next. Consider the situation where an error is made in estimating the correct permutation matrix for frequency bin f_k . Using the sequential approach, it is quite possible that many of the estimated permutation matrices for the frequency bins after f_k will be wrong. In the worst case scenario we will have half of the frequency bins with one permutation and the other half with a different permutation, which will result in very poor or no separation. Such a catastrophic situation is prevented by the global ILS minimization criterion (38).

The criterion of (38) can be properly modified in order to resolve the frequency-dependent permutation problem that arises when the second approach is followed. In this case we have

$$\min_{\{s_f\}_{f=0}^{T-1} \in \{1, \dots, I\}^T} \sum_{f=0}^{T-1} \left\{ \left\| \mathbf{C} - \mathbf{U}_n(f) \mathbf{\Pi}_1 1_{\{s_f=1\}} - \dots - \mathbf{U}_n(f) \mathbf{\Pi}_{I!} 1_{\{s_f=I!\}} \right\|_F^2 + \lambda \left\| \mathbf{U}_n(f) \mathbf{\Pi}_1 1_{\{s_f=1\}} + \dots + \mathbf{U}_n(f) \mathbf{\Pi}_{I!} 1_{\{s_f=I!\}} - \mathbf{U}_n(f-1) \mathbf{\Pi}_1 1_{\{s_{f-1}=1\}} - \dots - \mathbf{U}_n(f-1) \mathbf{\Pi}_{I!} 1_{\{s_{f-1}=I!\}} \right\|_F^2 \right\} \quad (39)$$

where $\mathbf{U}_n(-1) = \mathbf{0}^{J \times I}$, by convention. In (39), $\mathbf{U}_n(f) = [\mathbf{u}_1^{(n)}(f), \dots, \mathbf{u}_I^{(n)}(f)] \in \mathbb{R}^{J \times I}$ with $\mathbf{u}_i^{(n)}(f)$ representing the vector whose elements are the magnitudes of the respective elements of the (normalized with respect to an arbitrary microphone) i th column of the corresponding frequency's output matrix (32). Again, $\mathbf{C} \in \mathbb{R}^{J \times I}$ is the matrix whose columns are the centroids $\{\mathbf{c}_i, i = 1, \dots, I\}$, obtained after application of a clustering procedure via VQ over all the available vectors $\{\mathbf{u}_i^{(n)}(f), f = 0, \dots, T-1, i = 1, \dots, I\}$.

The ILS minimization criteria of (38), (39) proved to be satisfactory for resolving the permutation problem of the real-world BSS problems we examine in section VI. The following subsection considers the 2-sources case and discusses the selection of the algorithm we used for the implementation of the criterion (38). The pertinent analysis for the criterion (39) is identical, thus is omitted for brevity.

B. Special case: $J \times 2$ mixture

We will henceforth focus on the 2-sources (i.e., $(J, I) = (J, 2)$) case, but all our results readily extend to the case of general (J, I) . By letting $s_f \in \{-1, 1\}$, $\mathbf{\Pi}_{-1} = \mathbf{I} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ and $\mathbf{\Pi}_1 = \mathbf{\Pi} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$, the optimization criterion of (38) simplifies to

$$\min_{\{s_f\}_{f=1}^T \in \{\pm 1\}^T} \sum_{f=1}^T \left\{ \left\| \mathbf{C} - \frac{1-s_f}{2} \mathbf{A}_n(f) - \frac{1+s_f}{2} \mathbf{A}_n(f) \mathbf{\Pi} \right\|_F^2 + \lambda \left\| \frac{1-s_f}{2} \mathbf{A}_n(f) + \frac{1+s_f}{2} \mathbf{A}_n(f) \mathbf{\Pi} - \frac{1-s_{f-1}}{2} \mathbf{A}_n(f-1) - \frac{1+s_{f-1}}{2} \mathbf{A}_n(f-1) \mathbf{\Pi} \right\|_F^2 \right\} \quad (40)$$

It is shown in Appendix A that (40) can be put in the following more familiar form:

$$\min_{\mathbf{s} \in \{\pm 1\}^T} \|\mathbf{x} - \mathbf{B}\mathbf{s}\|_2^2 \quad (41)$$

where $\mathbf{s} = [s_1, s_2, \dots, s_T]^T$.

The ILS optimization problem of (41) is known to be NP-hard. Over the years, several methods have been developed to obtain high-quality solutions at manageable complexity cost. The trade-off between performance and complexity makes the selection of the most suitable algorithm an application-dependent issue.

Sphere Decoding (SD) [32] along with numerous variants of it [1], [5], [33], [34], and Semi-Definite Relaxation (SDR) [16], are methods that can provide near-optimal performance at often acceptable complexity. Unfortunately, they are too complex when the dimensionality of \mathbf{s} exceeds 100 (e.g. 128 in our context). Block Minimum Mean Squared Error - Decision Feedback Equalization (BMMSE-DFE) [6] and Probabilistic Data Association (PDA) [15], [22] are two simple methods that sometimes perform well, but they need variance estimates for “signal” and “noise” at their input, which are not obvious in the context of our application. A very suboptimal solution of the ILS optimization problem (41) is to use the sign of the unconstrained LS solution, known as Quantized Zero-Forcing (QZF) in the digital communications community, but performance leaves much to be desired. The best performance-complexity trade-off in our context was exhibited by the Successive Interference Cancellation - Iterative Least Squares (SIC-ILS) method [14]. In our specific application, SIC-ILS performs equally well as the near optimal Sphere Decoding method for low dimensionalities (e.g., 32 or 64). We therefore used SIC-ILS throughout.

C. Resolving the frequency-dependent scaling problem

After the resolution of the frequency-dependent permutation ambiguity, as described in the previous subsections, the only remaining problem is the frequency-dependent scaling ambiguity: i) of columns of the estimated mixing matrices $\{\mathbf{A}(f), f = 0, \dots, T-1\}$ for the DCE approach, ii) of rows of the estimated unmixing matrices $\{\mathbf{W}(f), f = 0, \dots, T-1\}$ for the ICE approach. First, let us consider the resolution of the scaling problem for the DCE approach, that is based upon the following simple fact whose proof can be found in Appendix B.

Fact 1: Consider a full rank matrix $\mathbf{A} \in \mathbb{C}^{J \times I}$, with $J \geq I$. Let $\mathbf{A}_n \in \mathbb{C}^{J \times I}$ denote the matrix whose i th column $\mathbf{a}_i^{(n)}$ is the i th column \mathbf{a}_i of \mathbf{A} , divided by the element $a_{j_i, i}$ of \mathbf{a}_i , i.e., $\mathbf{a}_i^{(n)} = \frac{1}{a_{j_i, i}} \mathbf{a}_i$, with $i = 1, \dots, I, j_i \in \{1, \dots, J\}$. Then $\mathbf{A}_n^\dagger \mathbf{A} = \text{diag}\{[a_{j_1, 1}, \dots, a_{j_I, I}]^T\}$.

Let

$$\mathbf{A}_s(f) = [d_1(f)\boldsymbol{\alpha}_1(f), \dots, d_I(f)\boldsymbol{\alpha}_I(f)] = \mathbf{A}(f) \begin{bmatrix} d_1(f) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & d_I(f) \end{bmatrix} = \mathbf{A}(f)\mathbf{D}(f) \quad (42)$$

where $d_i(f) \in \mathbb{C}$ for all $i = 1, \dots, I$, represent a column scaled version of the true mixing matrix $\mathbf{A}(f)$ at frequency f . Divide the columns of matrix $\mathbf{A}_s(f)$ by their respective elements $\{d_1(f)\alpha_{j_1, 1}(f), \dots, d_I(f)\alpha_{j_I, I}(f)\}$, with $j_i \in \{1, \dots, J\}$, to yield matrix $\mathbf{A}_{s,n}(f)$

$$\begin{aligned} \mathbf{A}_{s,n}(f) = \mathbf{A}_s(f) \begin{bmatrix} \frac{1}{d_1(f)\alpha_{j_1, 1}(f)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{d_I(f)\alpha_{j_I, I}(f)} \end{bmatrix} &= \mathbf{A}(f)\mathbf{D}(f)\mathbf{D}^{-1}(f) \begin{bmatrix} \frac{1}{\alpha_{j_1, 1}(f)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\alpha_{j_I, I}(f)} \end{bmatrix} = \\ &= \mathbf{A}(f) \begin{bmatrix} \frac{1}{\alpha_{j_1, 1}(f)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\alpha_{j_I, I}(f)} \end{bmatrix} = \mathbf{A}_n(f) \quad (43) \end{aligned}$$

where $\mathbf{A}_n(f)$ denotes the matrix whose i th column $\mathbf{a}_i^{(n)}(f)$ is the i th column $\mathbf{a}_i(f)$ of the true mixing matrix $\mathbf{A}(f)$, divided by element $a_{j_i, i}(f)$ of $\mathbf{a}_i(f)$.

Let us multiply (3) (we neglect the noise term) from the left with the pseudoinverse $\mathbf{A}_{s,n}^\dagger(f) = \mathbf{A}_n^\dagger(f)$. Invoking Fact 1, we have

$$\mathbf{A}_{s,n}^\dagger(f)\mathbf{x}(f,t) \approx \mathbf{A}_n^\dagger(f)\mathbf{A}(f)\mathbf{s}(f,t) = \text{diag}\{\alpha_{j_1,1}(f), \dots, \alpha_{j_I,I}(f)\}^T \mathbf{s}(f,t) = [\alpha_{j_1,1}(f)s_1(f,t), \dots, \alpha_{j_I,I}(f)s_I(f,t)]^T \quad (44)$$

Regarding the ICE approach, it is shown in Appendix C that by performing the appropriate manipulations on the row scaled version $\mathbf{W}_s(f)$ of the true unmixing matrix $\mathbf{W}(f)$, we can construct matrix $\mathbf{W}_{s,n}(f)$ given in (53). Multiplying (3) from the left with $\mathbf{W}_{s,n}(f)$ yields

$$\mathbf{W}_{s,n}(f)\mathbf{x}(f,t) \approx [\alpha_{j_1,1}(f)s_1(f,t), \dots, \alpha_{j_I,I}(f)s_I(f,t)]^T \quad (45)$$

We see from (44) and (45) that each speaker signal $s_i(f,t)$ is multiplied by the element $a_{j_i,i}(f)$, for all $i = 1, \dots, I$. But $a_{j_i,i}(f)$ is just the total frequency response⁶ that corresponds to the link between the i th speaker and the j_i th microphone, at frequency f , which is almost constant across all frequencies of interest (see assumption 5.2). Thus, the frequency-dependent scaling of the individual speaker signals that would result if we had used the pseudoinverses of the matrices $\{\mathbf{A}_s(f), f = 0, \dots, T-1\}$ in (44) or the matrices $\{\mathbf{W}_s(f), f = 0, \dots, T-1\}$ in (45), now is reduced to a single, frequency-independent scaling, i.e., the frequency-dependent scaling ambiguity problem is resolved. Intuitively, equations (44) and (45) imply that after separation we essentially listen to the i th speaker's speech signal the way it is captured by the j_i th microphone, as if there were no other interfering speakers.

VI. ALGORITHMS

In this section, we summarize our two BSS algorithms, namely DCE-PARAFAC-ILS and ICE-PARAFAC-ILS for $J \geq I$. The two algorithms are formulated as follows

- **Step 1**

DCE-PARAFAC-ILS: Compute the autocorrelation data (4) for P time intervals over which the measured signals are assumed stationary, using the sample average (19) on the frequency-domain data (3).

ICE-PARAFAC-ILS: If $J = I$ compute the “inverse” autocorrelation data (5) for P time intervals over which the measured signals are assumed stationary, using the sample average (19) on the frequency-domain data (3). If $J > I$ perform a dimension reducing prefiltering of the data in (3), as outlined in section II. Subsequently, compute the “inverse” autocorrelation data (8) for P time intervals over which the measured signals are assumed stationary, using the sample average (19) on the prefiltered frequency-domain data.

- **Step 2**

DCE-PARAFAC-ILS (ICE-PARAFAC-ILS): Apply the TALS algorithm to the autocorrelation data calculated in Step 1, iteratively for all frequency bins $f = 0, \dots, T-1$, to obtain the matrices $\{\mathbf{A}_{s,p}(f), f = 0, \dots, T-1\}$ ($\{\mathbf{W}_{s,p}(f), f = 0, \dots, T-1\}$) i.e., the estimates of the true mixing (unmixing) matrices up to a frequency-dependent permutation and scaling ambiguity of their columns (rows).

- **Step 3**

DCE-PARAFAC-ILS (ICE-PARAFAC-ILS): Apply a VQ clustering procedure over all the available vectors $\{\mathbf{a}_i^{(n)}(f), f = 0, \dots, T-1, i = 1, \dots, I\}$ ($\{\mathbf{u}_i^{(n)}(f), f = 0, \dots, T-1, i = 1, \dots, I\}$) as described in subsection A of section V, to obtain the associated with the I speakers centroids $\{\mathbf{c}_1, \dots, \mathbf{c}_I\}$. Form matrix \mathbf{C} whose columns are the centroids $\{\mathbf{c}_1, \dots, \mathbf{c}_I\}$.

- **Step 4**

⁶which equals the triple product of the frequency response of the acoustic channel between i th speaker- j_i th microphone, the frequency response of the j_i th microphone and an exponential term that is dependent on the delay, as verified from (33).

DCE-PARAFAC-ILS (ICE-PARAFAC-ILS): Apply the SIC-ILS algorithm in order to solve the ILS optimization problem in (38) (in (39)), and fix the arbitrary permutations as described in subsection A of section V, thus determining the matrices $\{\mathbf{A}_s(f), f = 0, \dots, T-1\}$ ($\{\mathbf{W}_s(f), f = 0, \dots, T-1\}$), i.e., the estimates of the true mixing (unmixing) matrices up to a frequency-dependent scaling ambiguity of their columns (rows).

- **Step 5**

DCE-PARAFAC-ILS (ICE-PARAFAC-ILS): Compute the matrices $\{\mathbf{A}_{s,n}(f), f = 0, \dots, T-1\}$ ($\{\mathbf{W}_{s,n}(f), f = 0, \dots, T-1\}$), as outlined in subsection C of section V.

- **Step 6**

DCE-PARAFAC-ILS (ICE-PARAFAC-ILS): Apply the Inverse DFT to the collection of matrices $\{\mathbf{A}_{s,n}^\dagger(f), f = 0, \dots, T-1\}$ ($\{\mathbf{W}_{s,n}, f = 0, \dots, T-1\}$) to yield the estimate of the inverse-channel impulse response matrix $\{\hat{\mathbf{W}}(\tau), \tau = 0, \dots, T-1\}$.

- **Step 7**

Estimate the speaker signals via $\hat{\mathbf{s}}(t) = \hat{\mathbf{W}} * \mathbf{x}(t)$

VII. EXPERIMENTAL RESULTS

In the present section, we report the results obtained after application of our two algorithms to real-world data of the PEACH multi-microphone database [20].

A. Experimental setup and performance measurement considerations

The mixed speech signals of the PEACH database were recorded using $I = 2$ loudspeakers and $J = 4$ omni-directional microphones. Analytically, the speech mixtures' acquisition procedure was conducted as follows: First, three utterances were acquired with a close-talk microphone for each one of eight different speakers. The speakers were separated into two groups of four speakers each, with speakers of the same group pronouncing the same three sentences every time. After the completion of the close-talk signals' acquisition procedure, the speakers were divided into four pairs, with each pair consisting of speakers that belong to different groups. The close-talk signals of each pair were reproduced together in a quiet room, by two loudspeakers in front of a microphone array consisting of four microphones, for four different geometric configurations. Specifically, utterance i of one speaker was played back simultaneously with utterance i of the pair's other speaker by the two loudspeakers, for $i = 1, 2, 3$. Thus, each pair of speakers provided three different speech mixtures for each geometric configuration. The measurements can be separated into four sessions of twelve experiments each, giving a total of 48 experiments, and with each session corresponding to a different geometric configuration. As shown in Fig. 1, the loudspeakers and the microphones were uniformly-spaced on two parallel lines at distance d_1 . The distance between microphones was d_2 , which was also the distance between the loudspeakers. For session I: $d_1 = 150\text{cm}$, $d_2 = 10\text{cm}$, for session II: $d_1 = 150\text{cm}$, $d_2 = 30\text{cm}$, for session III: $d_1 = 100\text{cm}$, $d_2 = 10\text{cm}$ and for session IV: $d_1 = 100\text{cm}$, $d_2 = 30\text{cm}$.

In order to quantify the separation performance, we measured the average input Signal-to-Interference Ratio (SIR) across microphones and the SIR for each of the unmixed signals that are the output of the BSS algorithms. SIR is measured following the method proposed in [24]: The close-talk signals of each speaker were played back through the corresponding loudspeaker one at a time, i.e. at each time only one source (loudspeaker) was active. Let $P_{j,i}^x = \sum_{t=0}^{T-1} x_j^2(t)$, where T is the number of samples, represent the power of the recorded signal at the j th microphone when only source i is active and all the other sources are inactive. The signal-to-interference ratio in dB of the recorded signal $x_j(t)$ at the j th microphone can be estimated using

$$SIR_x(j, i) = 10 \log \frac{\max_i P_{j,i}^x}{\sum_{i=1}^I P_{j,i}^x - \max_i P_{j,i}^x} \quad (46)$$

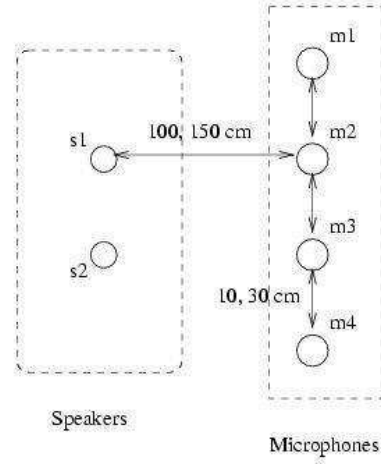


Fig. 1. Geometry of the measurement setup

for $j = 1, \dots, J$ and $i = 1, \dots, I$. In our case the number of microphones $J = 4$ and the number of sources $I = 2$. The formula (46) can also be used to measure $SIR_y(j, i)$, that is the SIR of the j th output of the separating algorithm, simply by substituting $P_{j,i}^x$ with $P_{j,i}^y$, the latter representing the power of the signal at the j th output of the algorithm when only source i is active, for $j, i = 1, \dots, I$. For our experiments $I = 2$. Finally, the improvement on the SIR for the i th source is defined as $SIR^+(i) = \max_j SIR_y(j, i) - \underline{SIR}_x(i)$, where $\underline{SIR}_x(i)$ is the average input SIR across microphones for the i th source, i.e. $\underline{SIR}_x(i) = \frac{\sum_{j=1}^J SIR_x(j, i)}{J}$.

The experiments have been conducted using speech data of approximately equal duration in time (about 19 seconds), sampled at 16KHz. Throughout all BSS experiments using our proposed algorithms, the number of time segments the speech mixtures were divided into, was set to $P = 10$. Due to the approximately equal length of all speech signals (about 300,000 samples), the size of each segment was about 30,000 data samples long, for all experiments. For the solution of the permutation problem, the two ILS minimization criteria in (38) and (39) were weighted equally, i.e., $\lambda = 1$.

B. Discussion on the results

Figures 2 and 3 show how the performance of the proposed algorithms changes for various numbers of frequency bins used to estimate the autocorrelation matrices, for all geometric configurations. The plots of each figure were obtained by averaging the improvements on the SIR over all twelve experiments and speakers of the corresponding session, at each point. It should be mentioned that depending on the experiment the SIR improvement reached up to more than 12 dBs in certain cases, but we chose to report the aggregate behavior of the proposed algorithms instead of focusing on specific experiments⁷. Note that on the average, the performance seems to decrease as the size of the Fast Fourier Transform (FFT) grid gets bigger; this is the reason we don't report the performance for more than 256 frequency bins. In most cases, the highest separation performance for both algorithms, is achieved for small FFT grid sizes of 16 and 32 bins. A second remark is that the separation performance is a function of the geometry of the measurement setup, since the performance improves in the following cases: i) when the distance between microphones increases, ii) when the distance between loudspeakers and microphones decreases. The rationale for

⁷A subset of the recorded speech mixtures along with the corresponding separated signals, can be found in the following website: http://www.speech.tuc.gr/res/research_bss.html.

TABLE II

COMPARISON TABLE FOR DCE-PARAFAC-ILS, ICE-PARAFAC-ILS, AND PARRA'S METHOD: AVERAGE SIR IMPROVEMENT IN DB FOR EACH EXPERIMENTAL SESSION, ALONG WITH THE AVERAGE EXECUTION TIME FOR EACH OF THE ABOVE METHODS IN SECONDS. THE NUMBER OF TIME SEGMENTS $P = 10$ WHILE THE NUMBER OF FREQUENCY BINS $T = 32$ FOR DCE-, ICE-PARAFAC-ILS AND $T = 4096$ FOR PARRA'S METHOD.

<i>Session</i>	<i>DCE-PARAFAC-ILS SIR impr. mean</i>	<i>ICE-PARAFAC-ILS SIR impr. mean</i>	<i>Parra's Method SIR impr. mean</i>
<i>I: 150cm-10cm</i>	2.9	3.2	3.5
<i>II: 150cm-30cm</i>	4.2	5.0	3.1
<i>III: 100cm-10cm</i>	5.6	5.2	6.8
<i>IV: 100cm-30cm</i>	6.3	5.3	1.0
<i>Overall Average</i>	4.750	4.675	3.600
<i>Execution Time</i>	4.5	6.4	351.3

this, partially lies in the method that resolves the frequency-dependent permutation problem; the method virtually exploits the diversification in the distances between the speakers-microphones positions in each geometric configuration, via the first criterion of (38) (see Assumption 5.1). In addition, PARAFAC is sensitive to diversity. In the limit as $d_1 \rightarrow \infty$ or as $d_2 \rightarrow 0$ there is no diversity, which results in poor or no estimation of the matrices $\{\mathbf{A}(f), f = 0, \dots, T-1\}$ and $\{\tilde{\mathbf{W}}(f), f = 0, \dots, T-1\}$ (up to a frequency-dependent permutation and scaling ambiguity) in the PARAFAC step of DCE-PARAFAC-ILS and ICE-PARAFAC-ILS, respectively. Figure 4 depicts the overall behavior of the proposed separating algorithms in terms of SIR improvement averaged over all four geometric configurations versus the number of frequency bins used. Finally, bearing in mind that the total running time of BSS algorithms is a crucial parameter (especially in the case where the algorithms are used in real-time applications), we report in figure 5 the average time needed for each of the proposed algorithms, in order to perform the source signals' separation for several FFT grid sizes. The times reported correspond to the MATLAB implementations of the algorithms, which were tested on a Pentium M-based personal computer running at 2 GHz with 1.5 MB of RAM. This figure demonstrates convincingly the strength of the proposed algorithms: High separation performance is achieved at a low computational cost, using 16 or 32 FFT grid points.

To support our arguments, the performance of our methods is compared with the method presented in [19], henceforth referred to as "Parra's method", for the corpus of the PEACH database. Table II reports the results for DCE-PARAFAC-ILS, ICE-PARAFAC-ILS and Parra's method for each experimental session along with the average execution time for each method, in order to perform the separation of the speaker signals. The parameters for "Parra's method" were chosen so that the method exhibits the highest possible performance. Specifically, the number of frequency bins $T = 4096$ while the length of the inverse channel impulse response was set to $Q = 512$, as suggested in [19]. For DCE-, ICE-PARAFAC-ILS $T = 32$. The number of time segments the speech signals were divided into, was set to $P = 10$, for all three methods. We see from Table II, that our methods outperform the method proposed in [19] in sessions II and IV, while Parra's method exhibits slightly better separation performance in the first and third experimental session. However, the overall separation performance in Table II, which can be thought of as a index of a method's effectiveness regardless of the geometric configuration, favors our proposed algorithms over Parra's method by more than 1 dB. In addition, the robustness of the new methods in changes of the experimental setup's geometry, seems to be higher when compared to Parra's method: The separation improvement of the PARAFAC-ILS methods varies from 2.9 to 6.3 dB's while Parra's method varies from 1 to 6.8 dB's. Finally, our algorithms are faster than Parra's by two orders of magnitude on average.

VIII. CONCLUSIONS

We have presented two new approaches for the blind speech separation problem. The new approaches employ PARAFAC to separate signals in the frequency domain and a context-specific ILS method to

resolve the frequency-dependent permutation ambiguity, which is a key problem in this setting. Our approaches offer guaranteed convergence, unlike [19], and better separation performance in most cases, measured both quantitatively and in subjective tests, at considerably reduced complexity. The link to the conjugate symmetric PARAFAC model established here, also shows that a much higher number of signals can be separated than what was suggested in [19]. In the future, we will investigate the separation performance of our algorithm as a function of the source-sensor geometry and the permutation minimization criterion weight λ .

IX. ACKNOWLEDGMENT

The authors would like to thank Mr. Christian Zieger and Mr. Luca Cristoforetti at ITC-IRST for collecting the corpus of the PEACH multi-microphone database. Thanks are also due to Prof. Vassilis Digalakis for useful suggestions.

APPENDIX A

It suffices to bring (40) in the form

$$\min_{\mathbf{s} \in \{\pm 1\}^T} \{ \|\mathbf{x}_1 - \mathbf{B}_1 \mathbf{s}\|_2^2 + \lambda \|\mathbf{x}_2 - \mathbf{B}_2 \mathbf{s}\|_2^2 \} \quad (47)$$

since $\|\mathbf{x}_1 - \mathbf{B}_1 \mathbf{s}\|_2^2 + \lambda \|\mathbf{x}_2 - \mathbf{B}_2 \mathbf{s}\|_2^2 = \|\mathbf{x} - \mathbf{B} \mathbf{s}\|_2^2$ for $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \sqrt{\lambda} \mathbf{x}_2 \end{bmatrix}$ and $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \sqrt{\lambda} \mathbf{B}_2 \end{bmatrix}$.

Considering the first term of (40) we have

$$\begin{aligned} \sum_{f=0}^{T-1} \left\| \mathbf{C} - \frac{1-s_f}{2} \mathbf{A}_n(f) - \frac{1+s_f}{2} \mathbf{A}_n(f) \mathbf{\Pi} \right\|_F^2 &= \\ \sum_{f=0}^{T-1} \left\| \mathbf{C} - \frac{\mathbf{A}_n(f)}{2} - \frac{\mathbf{A}_n(f) \mathbf{\Pi}}{2} + \left(\frac{\mathbf{A}_n(f)}{2} - \frac{\mathbf{A}_n(f) \mathbf{\Pi}}{2} \right) s_f \right\|_F^2 &= \\ \sum_{f=0}^{T-1} \left\| \mathbf{c}_1 - \frac{\mathbf{a}_1^{(n)}(f)}{2} - \frac{\mathbf{a}_2^{(n)}(f)}{2} + \left(\frac{\mathbf{a}_1^{(n)}(f)}{2} - \frac{\mathbf{a}_2^{(n)}(f)}{2} \right) s_f \right\|_2^2 + \\ \sum_{f=0}^{T-1} \left\| \mathbf{c}_2 - \frac{\mathbf{a}_2^{(n)}(f)}{2} - \frac{\mathbf{a}_1^{(n)}(f)}{2} + \left(\frac{\mathbf{a}_2^{(n)}(f)}{2} - \frac{\mathbf{a}_1^{(n)}(f)}{2} \right) s_f \right\|_2^2 &= \|\mathbf{x}_1 - \mathbf{B}_1 \mathbf{s}\|_2^2 \quad (48) \end{aligned}$$

where $\mathbf{x}_1 = \begin{bmatrix} \mathbf{c}_1 - \frac{\mathbf{a}_1^{(n)}(1)}{2} - \frac{\mathbf{a}_2^{(n)}(1)}{2} \\ \vdots \\ \mathbf{c}_1 - \frac{\mathbf{a}_1^{(n)}(T)}{2} - \frac{\mathbf{a}_2^{(n)}(T)}{2} \\ \mathbf{c}_2 - \frac{\mathbf{a}_1^{(n)}(1)}{2} - \frac{\mathbf{a}_2^{(n)}(1)}{2} \\ \vdots \\ \mathbf{c}_2 - \frac{\mathbf{a}_1^{(n)}(T)}{2} - \frac{\mathbf{a}_2^{(n)}(T)}{2} \end{bmatrix}$, $\mathbf{B}_1 = \begin{bmatrix} \frac{\mathbf{a}_1^{(n)}(1)}{2} - \frac{\mathbf{a}_2^{(n)}(1)}{2} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \frac{\mathbf{a}_1^{(n)}(T)}{2} - \frac{\mathbf{a}_2^{(n)}(T)}{2} \\ \frac{\mathbf{a}_2^{(n)}(1)}{2} - \frac{\mathbf{a}_1^{(n)}(1)}{2} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \frac{\mathbf{a}_2^{(n)}(T)}{2} - \frac{\mathbf{a}_1^{(n)}(T)}{2} \end{bmatrix}$ and $\mathbf{a}_i^{(n)}(f)$ denotes the i th column of $\mathbf{A}_n(f)$.

Similarly, the second term of (40) can be written as

$$\sum_{f=1}^T \left\| \frac{1-s_f}{2} \mathbf{A}_n(f) + \frac{1+s_f}{2} \mathbf{A}_n(f) \mathbf{\Pi} - \frac{1-s_{f-1}}{2} \mathbf{A}_n(f-1) - \frac{1+s_{f-1}}{2} \mathbf{A}_n(f-1) \mathbf{\Pi} \right\|_F^2 = \|\mathbf{x}_2 - \mathbf{B}_2 \mathbf{s}\|_2^2 \quad (49)$$

with

$$\mathbf{x}_2 = \begin{bmatrix} \frac{\mathbf{a}_1^{(n)}(1)}{2} + \frac{\mathbf{a}_2^{(n)}(1)}{2} \\ \frac{\mathbf{a}_1^{(n)}(2)}{2} + \frac{\mathbf{a}_2^{(n)}(2)}{2} - \frac{\mathbf{a}_1^{(1)}(1)}{2} - \frac{\mathbf{a}_2^{(n)}(1)}{2} \\ \vdots \\ \frac{\mathbf{a}_1^{(n)}(T)}{2} + \frac{\mathbf{a}_2^{(n)}(T)}{2} - \frac{\mathbf{a}_1^{(n)}(T-1)}{2} - \frac{\mathbf{a}_2^{(n)}(T-1)}{2} \\ \frac{\mathbf{a}_1^{(n)}(1)}{2} + \frac{\mathbf{a}_2^{(n)}(1)}{2} \\ \frac{\mathbf{a}_1^{(n)}(2)}{2} + \frac{\mathbf{a}_2^{(n)}(2)}{2} - \frac{\mathbf{a}_1^{(1)}(1)}{2} - \frac{\mathbf{a}_2^{(n)}(1)}{2} \\ \vdots \\ \frac{\mathbf{a}_1^{(n)}(T)}{2} + \frac{\mathbf{a}_2^{(n)}(T)}{2} - \frac{\mathbf{a}_1^{(n)}(T-1)}{2} - \frac{\mathbf{a}_2^{(n)}(T-1)}{2} \end{bmatrix} \text{ and}$$

$$\mathbf{B}_2 = \begin{bmatrix} \frac{\mathbf{a}_2^{(n)}(1)}{2} - \frac{\mathbf{a}_1^{(n)}(1)}{2} & \mathbf{0} & \dots & \mathbf{0} \\ \frac{\mathbf{a}_1^{(n)}(1)}{2} - \frac{\mathbf{a}_2^{(n)}(1)}{2} & \frac{\mathbf{a}_2^{(n)}(2)}{2} - \frac{\mathbf{a}_1^{(n)}(2)}{2} & \dots & \mathbf{0} \\ \vdots & \ddots & \ddots & \vdots \\ \mathbf{0} & \dots & \frac{\mathbf{a}_1^{(n)}(T-1)}{2} - \frac{\mathbf{a}_2^{(n)}(T-1)}{2} & \frac{\mathbf{a}_2^{(n)}(T)}{2} - \frac{\mathbf{a}_1^{(n)}(T)}{2} \\ \frac{\mathbf{a}_1^{(n)}(1)}{2} - \frac{\mathbf{a}_2^{(n)}(1)}{2} & \mathbf{0} & \dots & \mathbf{0} \\ \frac{\mathbf{a}_2^{(n)}(1)}{2} - \frac{\mathbf{a}_1^{(n)}(1)}{2} & \frac{\mathbf{a}_1^{(n)}(2)}{2} - \frac{\mathbf{a}_2^{(n)}(2)}{2} & \dots & \mathbf{0} \\ \vdots & \ddots & \ddots & \vdots \\ \mathbf{0} & \dots & \frac{\mathbf{a}_2^{(n)}(T-1)}{2} - \frac{\mathbf{a}_1^{(n)}(T-1)}{2} & \frac{\mathbf{a}_1^{(n)}(T)}{2} - \frac{\mathbf{a}_2^{(n)}(T)}{2} \end{bmatrix}$$

By setting $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \sqrt{\lambda}\mathbf{x}_2 \end{bmatrix}$ and $\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \sqrt{\lambda}\mathbf{B}_2 \end{bmatrix}$, as already mentioned, the optimization problem of (47) is reformulated as in (41).

APPENDIX B

Proof of Fact 1: Matrix \mathbf{A}_n can be written as

$$\mathbf{A}_n = \mathbf{A}\mathbf{V} \quad (50)$$

where the diagonal matrix $\mathbf{V} = \begin{bmatrix} \frac{1}{\alpha_{j_1,1}} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\alpha_{j_I,I}} \end{bmatrix}$. Hence, the left Moore-Penrose pseudoinverse of \mathbf{A}_n is written

$$\mathbf{A}_n^\dagger = (\mathbf{A}_n^H \mathbf{A}_n)^{-1} \mathbf{A}_n^H = ((\mathbf{A}\mathbf{V})^H \mathbf{A}\mathbf{V})^{-1} (\mathbf{A}\mathbf{V})^H = (\mathbf{V}^H \mathbf{A}^H \mathbf{A} \mathbf{V})^{-1} \mathbf{V}^H \mathbf{A}^H = \mathbf{V}^{-1} (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{V}^{-H} \mathbf{V}^H \mathbf{A}^H = \mathbf{V}^{-1} (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H \quad (51)$$

and

$$\mathbf{A}_n^\dagger \mathbf{A} = \mathbf{V}^{-1} (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H \mathbf{A} = \mathbf{V}^{-1} = \begin{bmatrix} \alpha_{j_1,1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \alpha_{j_I,I} \end{bmatrix} = \text{diag}\{[\alpha_{j_1,1}, \dots, \alpha_{j_I,I}]^T\} \quad (52)$$

This completes the proof.

APPENDIX C

Consider the resolution of the scaling problem for the ICE approach. In this case let

$$\mathbf{W}_s(f) = \begin{bmatrix} d_1(f) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & d_I(f) \end{bmatrix} \mathbf{A}^\dagger(f) = \mathbf{D}(f) \mathbf{A}^\dagger(f) \quad (53)$$

where $d_i(f) \in \mathbb{C}$ for all $i = 1, \dots, I$, represent a row scaled version of the true unmixing matrix $\mathbf{W}(f) = \mathbf{A}^\dagger(f)$ at frequency f . The right Moore-Penrose pseudoinverse of $\mathbf{W}_s(f)$ can be written as

$$\mathbf{W}_s^\dagger(f) = (\mathbf{D}(f) \mathbf{A}^\dagger(f))^\dagger = \mathbf{A}(f) \mathbf{D}^{-1}(f) \quad (54)$$

with entries $\mathbf{W}_s^\dagger(f)_{(j,i)} = \frac{\alpha_{j,i}(f)}{d_i(f)}$, where $\mathbf{X}_{(j,i)}$ denotes the (j, i) th entry of \mathbf{X} .

With $\mathbf{W}_s^\dagger(f)$ in our disposal, and using (53), (54) we can form the following matrix

$$\begin{aligned} \mathbf{W}_{s,n}(f) &= \begin{bmatrix} \frac{\alpha_{j_1,1}(f)}{d_1(f)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{\alpha_{j_I,I}(f)}{d_I(f)} \end{bmatrix} \mathbf{W}_s(f) = \begin{bmatrix} \alpha_{j_1,1}(f) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \alpha_{j_I,I}(f) \end{bmatrix} \mathbf{D}^{-1}(f) \mathbf{D}(f) \mathbf{A}^\dagger(f) = \\ & \begin{bmatrix} \alpha_{j_1,1}(f) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \alpha_{j_I,I}(f) \end{bmatrix} \mathbf{A}^\dagger(f) \quad (55) \end{aligned}$$

Now, by multiplying (3) (we neglect the noise term) from the left with the matrix $\mathbf{W}_{s,n}(f)$ given in (55) we have

$$\begin{aligned} \mathbf{W}_{s,n}(f) \mathbf{x}(f, t) &\approx \text{diag}\{\alpha_{j_1,1}(f), \dots, \alpha_{j_I,I}(f)\}^T \mathbf{A}^\dagger(f) \mathbf{A}(f) \mathbf{s}(f, t) = \\ &\text{diag}\{\alpha_{j_1,1}(f), \dots, \alpha_{j_I,I}(f)\}^T \mathbf{s}(f, t) = [\alpha_{j_1,1}(f) s_1(f, t), \dots, \alpha_{j_I,I}(f) s_I(f, t)]^T \quad (56) \end{aligned}$$

Again, we have reached to the result of (44).

REFERENCES

- [1] E. Agrell, T. Eriksson, A. Vardy, and K. Zeger, "Closest Point search in lattices," *IEEE Transactions on Information Theory*, vol. 48, no. 8, pp. 2201-2214, Aug. 2002.
- [2] J.M.F. ten Berge and N.D. Sidiropoulos, "On uniqueness in CANDECOMP/PARAFAC," *Psychometrica*, vol. 67, no. 3, Sept. 2002
- [3] R. Bro, N.D. Sidiropoulos, and G.B. Giannakis, "A fast least squares algorithm for separating trilinear mixtures," in *Proc. Int. Workshop Independent Component Analysis and Blind Signal Separation*, Aussois, France, Jan. 1999.
- [4] J.D. Carroll and J. Chang, "Analysis of individual differences in multidimensional scaling via an N-way generalization of "Eckart-Young" decomposition," *Psychometrica*, vol. 35, no. 3, pp. 283-319, 1970.
- [5] A. Chan and I. Lee, "A new reduced-complexity sphere decoder for multiple antenna systems," in *Proc. ICC*, vol. 1, New York, pp. 460-464, Apr. 28-May 2 2002.
- [6] A. Duel-Hallen, "A family of multiuser decision-feedback detectors for asynchronous code-division multiple-access channels," *IEEE Trans. on Communications*, vol. 43, issue 234, pp. 421-434, Feb. 1995.
- [7] A. Gersho and R.M. Gray, *Vector Quantization and Signal Compression*, Kluwer Academic Publishers, 1992.
- [8] S. Van Gerven and D. Van Compernelle, "Signal separation by symmetric adaptive decorrelation: Stability, convergence, and uniqueness," *IEEE Transactions on Signal Processing*, vol. 43, pp. 1602-1612, July 1995.
- [9] R.A. Harshman, "Foundation of the PARAFAC procedure: Model and conditions for an 'explanatory' multi-mode factor analysis," *UCLA Working Papers in Phonetics*, vol. 16, pp. 1-84, Dec. 1970.
- [10] R.A. Harshman, "Determination and proof of minimum uniqueness conditions for PARAFAC1," *UCLA Working Papers in Phonetics*, vol. 22, pp. 111-117, 1972.
- [11] S. Ikeda and N. Murata, "A method of ICA in time-frequency domain," in *Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, New Paltz, NY, Oct. 1997.
- [12] M. Kawamoto, K. Matsuoka, and N. Ohnishi, "A method of blind separation for convolved nonstationary signals," *Neurocomputing*, vol. 22, pp. 157-171, 1998.
- [13] J.B. Kruskal, "Three-way arrays: Rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics," *Linear Algebra Applications*, vol. 18, pp. 95-138, 1977.

- [14] T. Li and N.D. Sidiropoulos, "Blind digital signal separation using successive interference cancellation iterative least squares," *IEEE Transactions on Signal Processing*, vol. 48, no. 3, pp. 810-823, Apr. 2000.
- [15] J. Luo, K. Pattipati, P. Willet, and F. Hasegawa, "Near-optimal multiuser detection in synchronous CDMA using probabilistic data association," *IEEE Commun. Lett.*, vol. 5, no. 8, pp. 361-363, Sep. 2001.
- [16] W.-K. Ma, T.N. Davidson, K.M. Wong, Z.-Q. Luo, and P.-C. Ching, "Quasi-ML multiuser detection using semi-definite relaxation with application to synchronous CDMA," *IEEE Transactions on Signal Processing*, vol. 50, no. 4, pp. 912-922, Apr. 2002.
- [17] K. Matsuoka and S. Nakashima, "Minimal distortion principle for blind source separation," in *Proceeding of ICA 2001 Conference*, pp. 722-727, Dec. 2001.
- [18] N. Mitianoudis and M. Davies, "Audio source separation of convolutive mixtures," *IEEE Transactions on Speech and Audio Processing*, vol. 11, no. 5, pp. 489-497, Sep. 2003.
- [19] L. Parra and C. Spence, "Convolutive blind separation of non-stationary sources," *IEEE Transactions on Speech and Audio Processing*, vol. 8, pp. 320-327, May 2000.
- [20] PEACH multi-microphone database: <http://shine.itc.it>
- [21] D.T. Pham and J. Cardoso, "Blind source separation of instantaneous mixtures of nonstationary sources," *IEEE Transactions on Signal Processing*, vol. 49, no. 9, pp. 1837-1848, Sep. 2001.
- [22] D. Pham, K.R. Pattipati, P.K. Willet, and J. Luo, "A generalized probabilistic data association detector for multiple antenna systems," *IEEE Commun. Lett.*, vol. 8, no. 4, pp. 205-207, Apr. 2004.
- [23] D.T. Pham, C. Serviere, and H. Boumaraf, "Blind separation of convolutive audio mixtures using nonstationarity," in *Proceeding of ICA 2003 Conference*, Nara, Japan, Apr. 2003.
- [24] K. Rahbar and J.P. Reilly, "A Frequency Domain Method for Blind Source Separation of Convolutive Audio Mixtures," *IEEE Transactions on Speech and Audio Processing*, vol. 13, pp. 832-844, Sep. 2005.
- [25] Y. Rong, S.A. Vorobyov, A.B. Gershman, and N.D. Sidiropoulos, "Blind Spatial Signature Estimation via Time-Varying User Power Loading and Parallel Factor Analysis," *IEEE Transactions on Signal Processing*, vol. 53, pp. 1697-1710, May 2005.
- [26] H. Sahlin and H. Broman, "MIMO signal separation for FIR channels: a criterion and performance analysis," *IEEE Transactions on Signal Processing*, vol. 48, no. 3, pp. 642-649, Mar. 2000.
- [27] N.D. Sidiropoulos and R. Bro, "On the uniqueness of multilinear decomposition of N-way arrays," *J. Chemometrics*, vol. 14, pp. 229-239, 2000.
- [28] N.D. Sidiropoulos, R. Bro, and G.B. Giannakis, "Parallel factor analysis in sensor array processing," *IEEE Transactions on Signal Processing*, vol. 48, no. 8, pp. 2377-2388, Aug. 2000.
- [29] N.D. Sidiropoulos, G.B. Giannakis, and R. Bro, "Blind PARAFAC receivers for DS-CDMA systems," *IEEE Transactions on Signal Processing*, vol. 48, no. 3, pp. 810-823, Mar. 2000.
- [30] P. Smaragdis, "Blind separation of convolved mixtures in the frequency domain," *Neurocomputing*, vol. 22, pp. 21-34, 1998.
- [31] A. Stegeman, J.M.F. ten Berge, and L. De Lathauwer, "Sufficient conditions for uniqueness in CANDECOMP/PARAFAC and INDSCAL with random component matrices," *Psychometrika*, to appear.
- [32] E. Viterbo and J. Boutros, "A universal lattice code decoder for fading channels," *IEEE Transactions on Information Theory*, vol. 45, no. 7, pp. 1639-1642, July 1999.
- [33] R. Wang and G.B. Giannakis, "Approaching MIMO capacity with reduced-complexity soft sphere-decoding," in *Proc. WCNC*, Atlanta, GA, Mar. 21-25, 2004.
- [34] W. Zhao and G.B. Giannakis, "Sphere decoding algorithms with improved radius search," in *Proc. WCNC*, Atlanta, GA, Mar. 21-25, 2004.

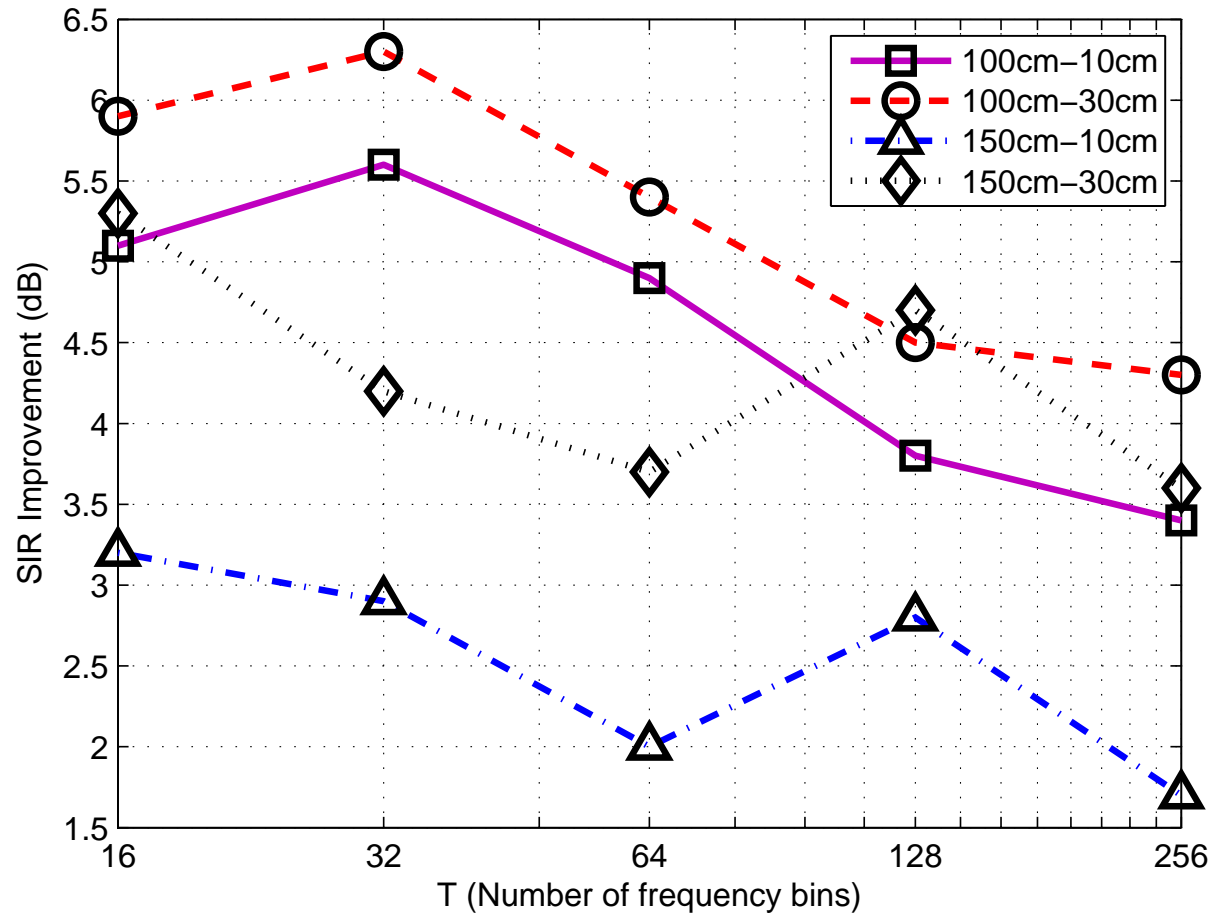


Fig. 2. Average SIR improvement versus number of frequency bins using DCE-PARAFAC-ILS.

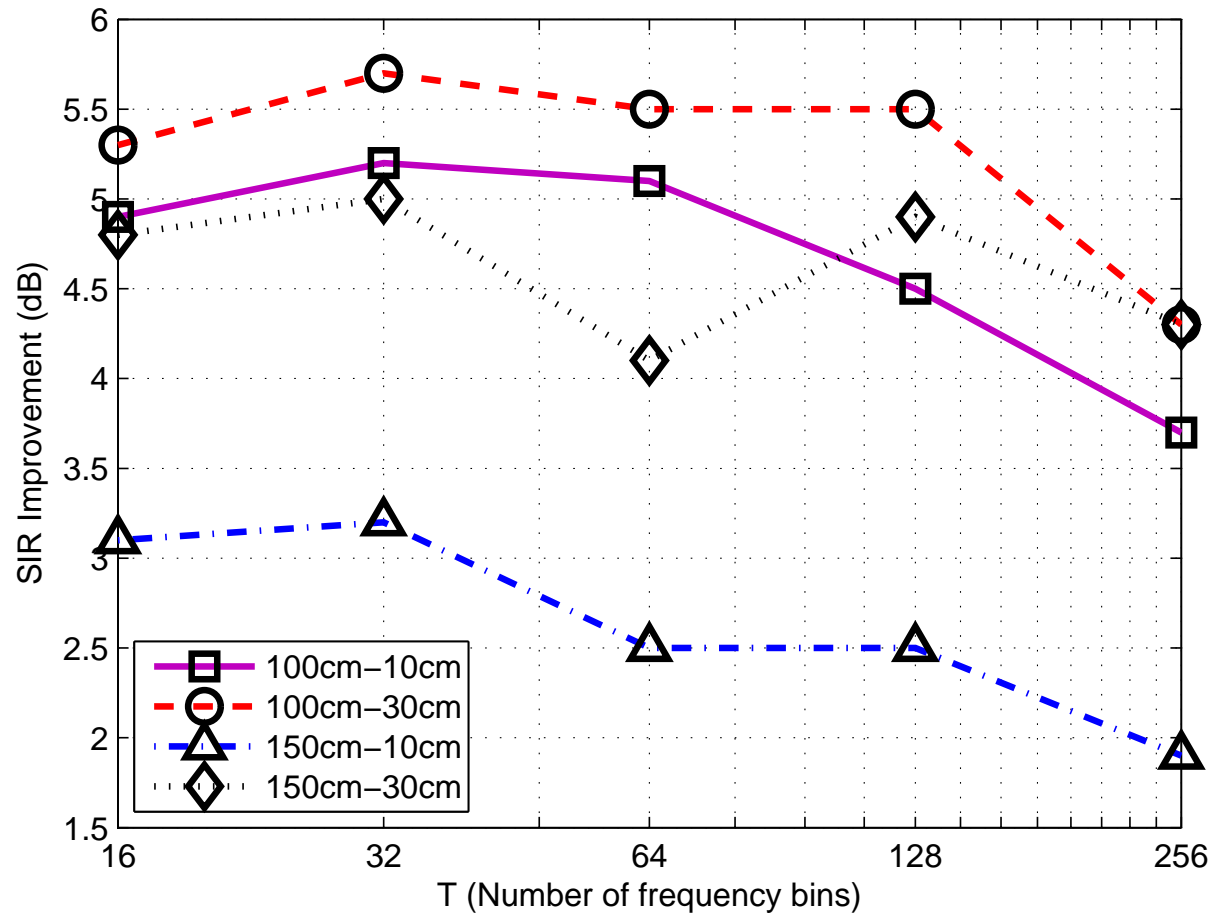


Fig. 3. Average SIR improvement versus number of frequency bins for ICE-PARAFAC-ILS.

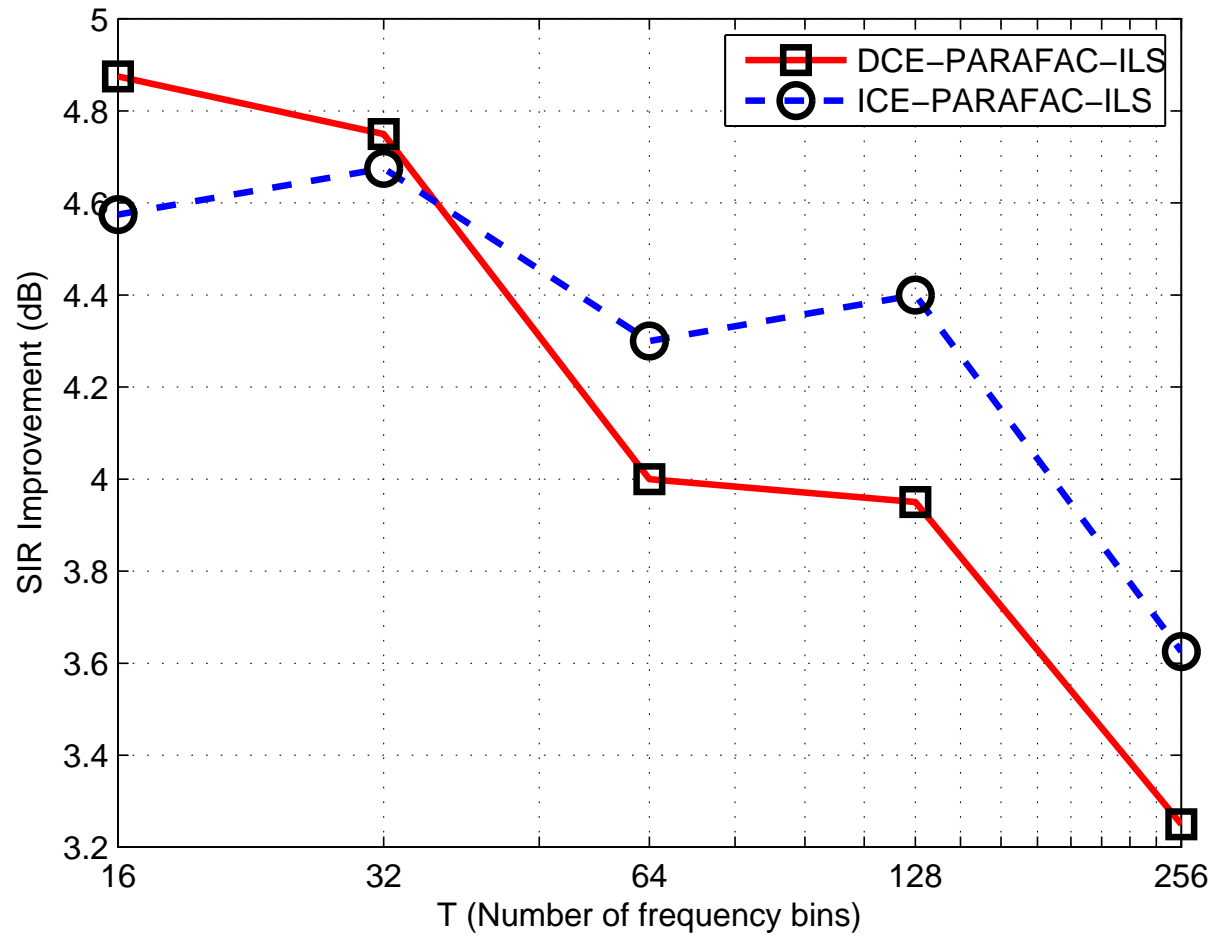


Fig. 4. Separation performance of our algorithms versus number of frequency bins, averaged over all geometric configurations

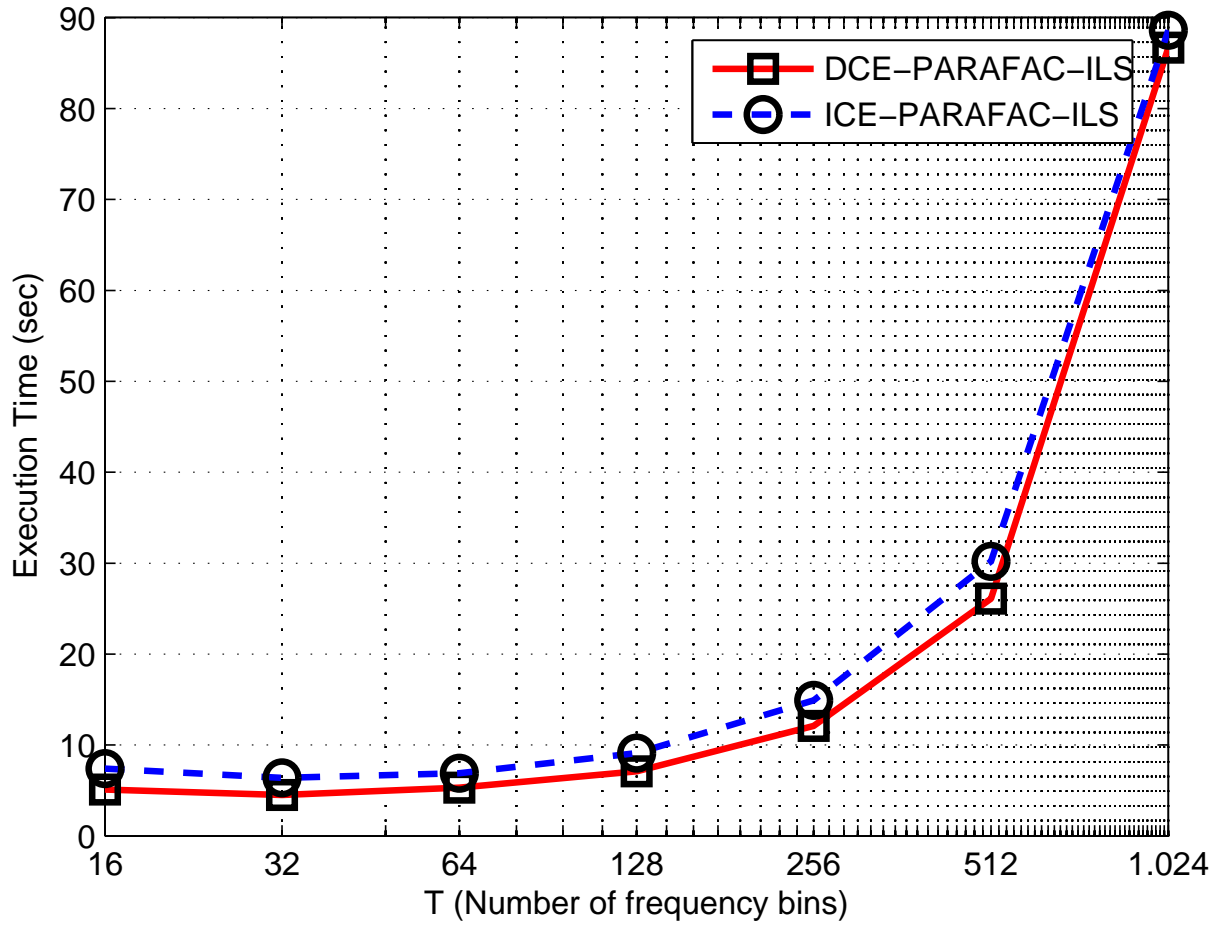


Fig. 5. Average execution time of our algorithms. The average length of the speech mixtures processed is about 300,000 samples.