

# BLIND SPEECH SEPARATION USING PARAFAC ANALYSIS AND INTEGER LEAST SQUARES

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

Dept. of Electronic and Computer Engineering Technical University of Crete, 73100 Chania, Greece

{ nikos,kleanthis,potam}@telecom.tuc.gr

## ABSTRACT

We propose a new two-step frequency domain algorithm for blind speech separation (BSS) for unknown channel order. This new approach employs parallel factor analysis (PARAFAC) to separate the speech signals and a novel integer-least-squares-based method for matching the arbitrary permutations in the frequency domain. The proposed algorithm offers guaranteed convergence and good separation performance, measured both quantitatively and in subjective tests. Performance gains in signal-to-interference ratio of up to 10 db are achieved for certain source-sensor geometries.

## 1. INTRODUCTION

One of the key problems in speech processing for teleconferencing and mobile telephony applications is that of speaker separation from multichannel measurements. 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 source signals, using only the measured microphone signals, i.e., without exploring any additional information about the properties of the sources or the mixing channels.

In recent years, several methods have been developed to tackle the BSS problem, but only few of them exploit the inherent non-stationarity of the speech signals. The majority of the proposed methods that exploit non-stationarity, consider the simple case of instantaneous linear mixtures [3]. For the more general case of convolutive linear mixtures, the approaches can be divided into frequency domain [1], [2] and time domain methods [4]. In this paper we propose a new frequency domain technique that exploits non-stationarity of speech signals and deals with the more general case of convolutive linear mixtures. Similar approaches have been proposed in [1], [2]. The method in [1], however, employs a least-squares (LS) gradient descent procedure in order to separate signals, which is not guaranteed to converge. On the contrary, our technique is based on the formulation of the BSS problem as a conjugate symmetric PARAFAC model that is fitted optimally using a fast and monotonically convergent alternating LS algorithm. Furthermore, 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 [1], [2]. The key problem in the BSS context is the resolution of the frequency-dependent permutation ambiguity. In [1] the frequency-dependent permutation ambiguity 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, novel approach based on integer-least-squares (ILS) to adjust the arbitrary permutations.

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

To support our arguments we present numerical simulations using real speech data and compare the performance of our algorithm with the one proposed in [1].

## 2. PROBLEM STATEMENT

Assume  $I$  mutually uncorrelated (not necessarily statistically independent) sources  $\mathbf{s}(t) = [s_1(t), \dots, s_I(t)]^T$ . These sources are convolved and mixed in a linear medium leading to  $J$  sensor signals  $\mathbf{x}(t) = [x_1(t), \dots, x_J(t)]^T$

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

for  $t = 1, \dots, N$ , where  $\mathbf{A}(\tau) = [\alpha_1(\tau), \dots, \alpha_I(\tau)]^T \in \mathbb{R}^{J \times I}$  for  $\tau = 0, \dots, L$  is the mixing impulse response matrix,  $\alpha_i(\tau) = [\alpha_{1,i}(\tau), \dots, \alpha_{J,i}(\tau)]^T \in \mathbb{R}^{J \times 1}$  is the spatial signature of the  $i$ th source for lag  $\tau$ ,  $\mathbf{n}(t) = [n_1(t), \dots, n_J(t)]^T$  is the additive noise vector,  $L$  is the maximum (unknown) channel length,  $N$  is the number of snapshots,  $*$  denotes convolution and  $(\cdot)^T$  denotes the transpose. The objective of the blind separation problem is to estimate the inverse-channel impulse response matrix  $\mathbf{W}(\tau)$  from the observed signals  $\mathbf{x}(t)$ , such that their convolution provides us with an estimate of the source signals  $\hat{\mathbf{s}}(t)$

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

Before proceeding further we make the following main assumptions: (i) The sources  $\mathbf{s}(t)$  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, and (ii) The noise  $\mathbf{n}(t)$  is zero mean, uncorrelated at each sensor and independent of the sources.

Following, e.g. [1], 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 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 convolutional sum [1]. In such a case we can write approximately

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

where  $\mathbf{x}(f, t) = \sum_{\tau=0}^{T-1} \mathbf{x}(t+\tau) e^{-\frac{j2\pi f \tau}{T}}$  is the DFT of the frame of size  $T$  starting at  $t$ ,  $[\mathbf{x}(t), \dots, \mathbf{x}(t+T-1)]$  and corresponding expressions apply for  $\mathbf{s}(f, t)$  and  $\mathbf{A}(f)$ . The  $i$ th column of  $\mathbf{A}(f)$  represents now the spatial signature of the  $i$ th source in the frequency domain, at frequency  $f$ . Let us focus on a time interval over which the measured signals can be assumed

stationary (see assumption (i)). The autocorrelation matrix of the vector of microphone outputs at frequency  $f$  is then

$$\begin{aligned} R_x(f, t) &= E[\mathbf{x}(f, t)\mathbf{x}^H(f, t)] \\ &= \mathbf{A}(f)E[\mathbf{s}(f, t)\mathbf{s}^H(f, t)]\mathbf{A}^H(f) + E[\mathbf{n}(f, t)\mathbf{n}^H(f, t)] \\ &= \mathbf{A}(f)D_s(f, t)\mathbf{A}^H(f) + D_n(f, t) \end{aligned} \quad (4)$$

where  $E[\cdot]$  denotes the expectation operator and  $(\cdot)^H$  denotes the Hermitian transpose. Since we assume mutually uncorrelated sources we postulate diagonal autocorrelation matrix of the sources  $D_s(f, t)$ , while by assumption the autocorrelation matrix of the additive noise  $D_n(f, t)$  is also diagonal.

Now assume that the matrices  $\mathbf{A}(f)$ ,  $f = 0, \dots, T-1$  are available. Then, by taking the 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  $\mathbf{A}^\dagger(f)$  for  $f = 0, \dots, T-1$ , where  $(\cdot)^\dagger$  denotes the Moore-Penrose pseudo-inverse, we could determine estimates of the inverse-channel matrices  $\widehat{\mathbf{W}}(\tau)$ , where  $\tau = 0, \dots, T-1$ . Therefore the BSS problem boils down to the problem of estimating the matrices  $\mathbf{A}(f)$ ,  $f = 0, \dots, T-1$ .

It is well known that the matrices  $\mathbf{A}(f)$ ,  $f = 0, \dots, T-1$  can only be specified up to an individual permutation and scaling of their columns (spatial signatures), that vary arbitrarily from one frequency bin to another. This constitutes a serious problem since only consistent permutations and scaling across frequencies will correctly reconstruct the sources. Hence a second major problem arises, that of resolving the frequency-dependent permutation and scaling issue.

### 3. PARALLEL FACTOR ANALYSIS

By dividing the whole data block of  $N$  snapshots into  $P$  sub-blocks, with each sub-block corresponding to a time interval over which the source signals are assumed stationary, the measured snapshots within any  $p$ th sub-block correspond to the following autocorrelation matrix

$$R_x(f, t_p) = \mathbf{A}(f)D_s(f, t_p)\mathbf{A}^H(f) + D_n(f, t_p) \quad (5)$$

for  $f = 0, \dots, T-1$ , in accordance with (4). Using all  $P$  sub-blocks, we will have  $P$  different autocorrelation matrices  $\{R_x(f, t_1), \dots, R_x(f, t_P)\}$  for each frequency. By neglecting the noise component for the moment, we observe that for each frequency, these matrices differ from each other only because the source signal autocorrelation matrices  $D_s(f, t_p)$  differ from one sub-block to another.

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

$$r_{j,l,p} \triangleq [\underline{R}_x(f)] = \sum_{i=1}^I \alpha_{j,i}(f) v_i(f, p) \alpha_{l,i}^* \quad (6)$$

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

$$\mathbf{P}(f) \triangleq \begin{bmatrix} v_1(f, 1) & \dots & v_I(f, 1) \\ \vdots & \ddots & \vdots \\ v_1(f, P) & \dots & v_I(f, P) \end{bmatrix} \quad (7)$$

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

$$D_s(f, t_p) = \mathcal{D}_p\{\mathbf{P}(f)\} \quad (8)$$

for all  $p = 1, \dots, P$  and all  $f = 0, \dots, T-1$ . In (8),  $\mathcal{D}_p\{\cdot\}$  is the operator which makes a diagonal matrix by selecting the  $p$ th row and putting it on the main diagonal while putting zeros elsewhere.

Equation (6) implies that  $r_{j,l,p}$  is a sum of rank-1triple products; this equation is known as (conjugate-symmetric) *parallel factor* (PARAFAC) analysis of  $r_{j,l,p}$  [5]. If  $I$  is sufficiently small (6) represents a low-rank decomposition of  $\underline{R}_x(f)$ . 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  $\underline{R}_x(f)$ . By solving a similar problem separately for every frequency we obtain the entire collection of the frequency-domain mixing matrices  $\mathbf{A}(f)$ ,  $f = 0, \dots, T-1$ .

The uniqueness (up to inherently unresolvable source permutation and scale ambiguities) of all source spatial signatures, given the exact frequency-domain autocorrelation data is an issue known as the identifiability of the model. Identifiability conditions i.e., conditions under which the trilinear decomposition of  $\underline{R}_x(f)$  is unique, are discussed in [5].

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

$$\widehat{R}_x(f, t_p) = \frac{1}{K} \sum_{k=0}^{K-1} \mathbf{x}(f, t_p + kT) \mathbf{x}^H(f, t_p + kT) \quad (9)$$

for  $p = 1, \dots, P$ .

In our problem, PARAFAC fitting at each frequency was based on the implementation of a fast and monotonically convergent least squares separation algorithm, known as trilinear alternating least squares (TALS) technique [5], which is used to estimate the matrices  $\mathbf{A}(f)$ , 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 4. For more information on TALS we refer the interested reader to the above references.

### 4. RESOLVING THE SCALING AND PERMUTATION AMBIGUITIES

In this section, we propose a novel approach for resolving the scaling and permutation ambiguities that arise in the speech separation context. In the first subsection, our method of dealing with the scaling problem is presented, while the second subsection describes how the frequency-dependent permutation ambiguity is resolved.

#### 1. Resolving the frequency-dependent scaling problem

Assume, for the moment, that the frequency-dependent permutation ambiguity is resolved, that is, the only remaining problem is the frequency-dependent scaling ambiguity of columns of the estimated mixing matrices. The resolution of the latter is based upon the following result, whose proof is omitted for brevity

*Result 1:* Consider a full rank matrix  $\mathbf{A} \in \mathbb{C}^{J \times I}$ , with  $J > I$ . Let  $\mathbf{A}_n \in \mathbb{C}^{J \times I}$  denote the matrix whose  $i$ th column  $\mathbf{A}_{n(:,i)}$  is the  $i$ th column of  $\mathbf{A}$ ,  $\mathbf{A}_{(:,i)}$ , divided by the element  $\alpha_{j_i,i}$  of  $\mathbf{A}_{(:,i)}$ , with  $j_i \in \{1, \dots, J\}$ , i.e.,  $\mathbf{A}_{n(:,i)} = \frac{1}{\alpha_{j_i,i}} \mathbf{A}_{(:,i)}$ ,  $i = 1, \dots, I$ . Then  $\mathbf{A}_n^\dagger \mathbf{A} = \text{diag}\{\alpha_{j_1,1}, \dots, \alpha_{j_I,I}\}^T$ , where  $\text{diag}\{\mathbf{x}\}$  forms a diagonal matrix from vector  $\mathbf{x}$ .  
Let

$$\begin{aligned} \mathbf{A}_s(f) &= [c_1(f)\alpha_1(f), \dots, c_I(f)\alpha_I(f)] = \\ &= \mathbf{A}(f) \text{diag}\{[c_1(f), \dots, c_I(f)]^T\} = \mathbf{A}(f)\mathbf{C}(f) \end{aligned} \quad (10)$$

where  $c_i(f) \in \mathbb{C}$  for all  $i = 1, \dots, I$ , represent a 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  $\{c_1(f)\alpha_{j_1,1}(f), \dots, c_I(f)\alpha_{j_I,I}(f)\}$ , with  $j_i \in \{1, \dots, J\}$ , to yield matrix  $\mathbf{A}_{s,n}(f)$

$$\begin{aligned} \mathbf{A}_{s,n}(f) &= \mathbf{A}_s(f) \text{diag} \left\{ \frac{1}{c_1(f)\alpha_{j_1,1}(f)}, \dots, \frac{1}{c_I(f)\alpha_{j_I,I}(f)} \right\}^T = \\ &= \mathbf{A}(f) \mathbf{C}(f) \mathbf{C}^{-1}(f) \text{diag} \left\{ \frac{1}{\alpha_{j_1,1}(f)}, \dots, \frac{1}{\alpha_{j_I,I}(f)} \right\}^T = \\ &= \mathbf{A}(f) \text{diag} \left\{ \frac{1}{\alpha_{j_1,1}(f)}, \dots, \frac{1}{\alpha_{j_I,I}(f)} \right\}^T = \mathbf{A}_n(f) \quad (11) \end{aligned}$$

where  $\mathbf{A}_n(f)$  denotes the matrix whose  $i$ th column  $\mathbf{A}_n(f)_{(:,i)}$  is the  $i$ th column of the true mixing matrix  $\mathbf{A}(f)$ ,  $\mathbf{A}(f)_{(:,i)}$ , divided by element  $\alpha_{j_i,i}(f)$  of  $\mathbf{A}(f)_{(:,i)}$ .

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

$$\begin{aligned} \mathbf{A}_{s,n}^\dagger(f) \mathbf{x}(f, t) &\approx \mathbf{A}_n^\dagger(f) \mathbf{A}(f) \mathbf{s}(f, t) = \\ &= \text{diag} \{ [\alpha_{j_1,1}(f), \dots, \alpha_{j_I,I}(f)]^T \} \mathbf{s}(f, t) = \\ &= [\alpha_{j_1,1}(f) s_1(f, t), \dots, \alpha_{j_I,I}(f) s_I(f, t)]^T \quad (12) \end{aligned}$$

We see from (12), that each source signal  $s_i(f, t)$  is multiplied by the element  $\alpha_{j_i,i}(f)$ , for all  $i = 1, \dots, I$ . But  $\alpha_{j_i,i}(f)$  is just the total frequency response that corresponds to the link between the  $i$ th speaker and the  $j_i$ th microphone, at frequency  $f$ , which is approximately constant across frequency in (quasi-) non reverberant environments. Thus, the frequency-dependent scaling of the individual sources that would result if we had used the pseudo-inverses of the scaled matrices (10) in (12), now is reduced to a single, frequency-independent scaling, i.e., the frequency-dependent scaling ambiguity problem is resolved.

Intuitively, the filtering process of the  $i$ th speaker's speech signal  $s_i(t)$  through  $\alpha_{j_i,i}(f)$ , that yields the separated signal  $\hat{s}_i(t) = \alpha_{j_i,i}(t) * s_i(t)$ , means that after separation, we essentially listen to the signal  $s_i(t)$  the way it is captured by the  $j_i$ th microphone, as if there were no other interfering speakers.

## 2. Resolving the frequency-dependent permutation problem

In the previous subsection we assumed that the frequency-dependent permutation ambiguity was already resolved. Thus, prior to applying the method presented above for dealing with the scaling problem, the frequency-dependent permutation ambiguity of the columns of the matrices  $\mathbf{A}(f)$ ,  $f = 0, \dots, T-1$ , must be lifted. 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 also formulate a joint optimization criterion that incorporates these two constituent criteria.

•**Criterion 1:** Let the collection of vectors  $\{\alpha_{n_i}(f), f = 0, \dots, T-1, i = 1, \dots, I\}$  denote the entire collection of the (normalized with respect to an arbitrary reference microphone) spatial signatures that are acquired after the application of the TALS algorithm. In non-reverberant (or, mildly reverberant) environments, where the recordings have been conducted using microphones with (almost) flat frequency responses, the magnitudes of the elements of the normalized spatial signatures, are approximately independent of frequency. Thus the aforementioned collection of normalized spatial signatures can be divided into

$I$  separate clusters, each of the latter being associated with a different source.

We use vector quantization [7] to determine the  $I$  centroids  $\{\mathbf{c}_1, \dots, \mathbf{c}_I\}$ , with  $\mathbf{c}_i \in \mathbb{C}^J$ , but other clustering techniques could also be used. With the  $I$  centroids at our disposal, the frequency-dependent permutation problem boils down to the following Integer Least Squares (ILS) minimization problem

$$\begin{aligned} \min_{\{s_f\}_{f=1}^T \in \{1, \dots, I\}^T} & \sum_{f=1}^T \|\mathbf{C} - \mathbf{A}_n(f) \mathbf{\Pi}_1 1_{\{s_f=1\}} \\ & - \mathbf{A}_n(f) \mathbf{\Pi}_2 1_{\{s_f=2\}} - \dots - \mathbf{A}_n(f) \mathbf{\Pi}_I 1_{\{s_f=I\}}\|_F^2 \quad (13) \end{aligned}$$

where  $\mathbf{C} \in \mathbb{C}^{J \times I}$  is the matrix whose columns are the centroids  $\{\mathbf{c}_i, i = 1, \dots, I\}$ ,  $\mathbf{A}_n(f) \in \mathbb{C}^{J \times I}$  is the matrix whose columns are the arbitrarily permuted vectors  $\{\alpha_{n_i}(f), i = 1, \dots, I\}$ ,  $\mathbf{\Pi}_l \in \mathbb{C}^{I \times I}$  for  $l = 1, \dots, I!$  is one of the total number of  $I!$  permutation matrices of dimensionality  $I \times I$ . Note that  $1_{\{x=\alpha\}} = 0$  if  $x \neq \alpha$  and  $1_{\{x=\alpha\}} = 1$  if  $x = \alpha$ .

•**Criterion 2:** For adjacent frequency bins,  $f_1, f_2 = f_1 + 1$ ,  $\alpha_i(f_1) \approx \alpha_i(f_2)$  and consequently  $\alpha_{n_i}(f_1) \approx \alpha_{n_i}(f_2)$ . This is valid for dense FFT grids even in reverberant environments, by continuity of the Fourier transform. Based on this fact, we formulate our second ILS minimization criterion

$$\begin{aligned} \min_{\{s_f\}_{f=1}^T \in \{1, \dots, I\}^T} & \sum_{f=1}^T \|\mathbf{A}_n(f) \mathbf{\Pi}_1 1_{\{s_f=1\}} + \dots + \mathbf{A}_n(f) \mathbf{\Pi}_{I!} 1_{\{s_f=I!\}} \\ & - \mathbf{A}_n(f-1) \mathbf{\Pi}_1 1_{\{s_{f-1}=1\}} - \dots - \mathbf{A}_n(f-1) \mathbf{\Pi}_{I!} 1_{\{s_{f-1}=I!\}}\|_F^2 \quad (14) \end{aligned}$$

where  $\mathbf{A}_n(0) = \mathbf{0}^{J \times I}$ , by convention.

The criteria of (13) and (14) can be combined into one overall ILS minimization criterion

$$\begin{aligned} \min_{\{s_f\}_{f=1}^T \in \{1, \dots, I\}^T} & \sum_{f=1}^T \{ \|\mathbf{C} - \mathbf{A}_n(f) \mathbf{\Pi}_1 1_{\{s_f=1\}} - \dots \\ & - \mathbf{A}_n(f) \mathbf{\Pi}_{I!} 1_{\{s_f=I!\}}\|_F^2 + \lambda \|\mathbf{A}_n(f) \mathbf{\Pi}_1 1_{\{s_f=1\}} + \dots \\ & + \mathbf{A}_n(f) \mathbf{\Pi}_{I!} 1_{\{s_f=I!\}} - \mathbf{A}_n(f-1) \mathbf{\Pi}_1 1_{\{s_{f-1}=1\}} - \dots \\ & - \mathbf{A}_n(f-1) \mathbf{\Pi}_{I!} 1_{\{s_{f-1}=I!\}}\|_F^2 \quad (15) \end{aligned}$$

where  $\lambda$  is a weighting factor. When  $\lambda \rightarrow 0$  (practically, for  $\lambda \ll 1$ ), the combined criterion (15) coincides with criterion 1, while for  $\lambda \rightarrow \infty$  (practically, for  $\lambda \gg 1$ ), (15) coincides with criterion 2.

The ILS minimization criterion of (15) proves to be satisfactory for resolving the permutation problem that arises in the majority of real world BSS problems, as will be shown in section 5. The ILS problem in (15) is NP-hard, however, a plethora of polynomial complexity approximate solutions have been developed for it in the multiuser detection literature. We chose to employ the Successive Interference Cancellation - Iterative Least Squares (SIC-ILS) method [6], since it exhibits the best performance-complexity trade-off in our context, where the decision vectors are of very high dimensionality (=FFT length).

## 5. EXPERIMENTAL RESULTS

In our experiments, we used a subset of the PEACH multi-microphone database [8] using  $I = 2$  loudspeakers as sources and  $J = 4$  omni-directional microphone as sensors. The measurements were separated into three sessions of three experiments each, with each session corresponding to a different rectangular geometric configuration. The loudspeakers and the microphones were uniformly-spaced on two parallel lines at distance  $d_1$ . The

**Table I.** BSS performance of proposed PARAFAC-ILS method compared to Parra's method [1] in terms of output SIR and SIR improvement in dB for each session, speaker pair and utterance (source S1 or S2).

Experiment:	Session I: $d_1=150\text{cm}$ $d_2=10\text{cm}$						Session II: $d_1=150\text{cm}$ $d_2=30\text{cm}$						Session III: $d_1=100\text{cm}$ $d_2=30\text{cm}$					
	1st pair		2nd pair		3rd pair		1st pair		2nd pair		3rd pair		1st pair		2nd pair		3rd pair	
	S1	S2	S1	S2	S1	S2	S1	S2	S1	S2	S1	S2	S1	S2	S1	S2	S1	S2
Raw Data Mean SIR	-3.8	3.8	-4.1	4.1	0.5	-0.5	-4.6	4.6	-3.0	3.0	-0.2	0.2	-6.7	6.7	-1.0	1.0	-1.2	1.2
Parra's Method SIR	-3.0	5.5	3.3	6.0	5.3	4.8	-4.2	6.4	1.9	-0.1	3.8	3.8	-3.3	4.6	1.0	1.8	1.8	0.4
SIR Improvement	<b>0.8</b>	1.7	<b>7.4</b>	<b>1.9</b>	<b>4.8</b>	<b>5.3</b>	0.4	1.8	4.9	-3.1	<b>4.0</b>	<b>3.6</b>	3.4	-2.1	2.0	0.8	3.0	-0.8
PARAFAC-ILS SIR	-3.9	6.9	0.7	5.9	1.5	1.8	-1.2	9.7	2.5	5.3	<b>2.1</b>	<b>3.2</b>	-3.2	12.3	3.5	3.4	3.5	9.5
SIR Improvement	-0.1	<b>3.1</b>	4.8	1.8	1.0	2.3	<b>3.4</b>	<b>5.1</b>	<b>5.5</b>	<b>2.3</b>	2.3	3.0	<b>3.5</b>	<b>5.6</b>	<b>4.5</b>	<b>2.4</b>	<b>4.7</b>	<b>8.3</b>

**Table II.** Overall BSS performance of PARAFAC-ILS and Parra's methods in terms of average SIR improvement in dB for each experimental session.

Session	PARAFAC-ILS SIR impr. mean	Parra's Method SIR impr. mean
I: 150cm-10cm	2.15	3.65
II: 150cm-30cm	3.60	1.93
III: 100cm-30cm	4.83	1.05
Overall	3.53	2.21

distance between microphones was  $d_2$ . For session I:  $d_1 = 150\text{cm}$ ,  $d_2 = 10\text{cm}$ , for session II:  $d_1 = 150\text{cm}$ ,  $d_2 = 30\text{cm}$  and for session III:  $d_1 = 100\text{cm}$ ,  $d_2 = 30\text{cm}$ . For all experiments, the transcription of the utterance played back from each loudspeaker was the same. The same three sets of speakers were used for each of the three session giving a total of 9 experiments. For more information on the measurement setup see [8].

In order to quantify the separation performance, we measured the average Signal to Interference Ratio (SIR) across microphones (mean SIR of raw data in Table I) 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 [2]. A separate SIR value was measured for each source ("S1" and "S2" in Table I).

Throughout all BSS experiments using our proposed PARAFAC-ILS algorithm, the number of time segments the speech signals were divided into, was set to  $P = 5$ . For the solution of the permutation problem, the two ILS minimization criteria in Eq. (15) were weighted equally, i.e.,  $\lambda = 1$ . The frequency-domain sample autocorrelation matrix estimates were calculated according to Eq. (9) using  $T = 512$  FFT points in each segment. The performance of our method was compared with the method presented in [1] referred to as "Parra's method" in Table I. The parameters for Parra's method are identical to the PARAFAC-ILS parameters; in addition, the length of the inverse channel impulse response in Parra's method was set to  $Q = 128$  as suggested in [1]. In Table I, the raw data SIR is presented for each experiment and signal source along with the SIR of the two BSS algorithms. Improvement over raw data SIR is also shown for the two methods (marked as "SIR Improvement"). The overall performance of the BSS algorithms in terms of average SIR improvement in dB for each session is reported in Table II<sup>1</sup>.

We see from Table II, that our method outperforms the method proposed in [1] in sessions II and III. Furthermore, the overall score of the PARAFAC-ILS method is higher by more than 1.3 dB compared to Parra's method. In the first experimental set, however, Parra's method exhibits better separation performance.

<sup>1</sup>A 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](http://www.speech.tuc.gr/res/research_bss.html)

Bearing in mind that the experimental sets correspond to different geometric configurations, we conclude that the separation performance is a function of the geometry of the measurement setup.

## 6. CONCLUSIONS

We have presented a new two-step approach for the blind speech separation problem. The new approach employs PARAFAC to separate signals in the frequency domain and a context-specific ILS method to resolve the frequency-dependent permutation ambiguity, which is the key problem in this setting. Our approach offers guaranteed convergence, unlike [1], and often considerably better separation performance, measured both quantitatively and in subjective tests. 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 [1]. 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  $\lambda$ .

## 7. REFERENCES

- [1] L. Parra and C. Spence, "Convolutional blind separation of non-stationary sources," *IEEE Transactions on Speech and Audio Processing*, vol. 8, pp. 320-327, May 2000.
- [2] K. Rahbar and J.P. Reilly, "A Frequency Domain Method for Blind Source Separation of Convolutional Audio Mixtures," *IEEE Transactions on Speech and Audio Processing*, vol. 13, pp. 832-844, Sep. 2005.
- [3] 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.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] A. Gersho and R.M. Gray, *Vector Quantization and Signal Compression*, Kluwer Academic Publishers, 1992.
- [8] PEACH multi-microphone database: <http://shine.itc.it>