

# MODEL-BASED SPARSE COMPONENT ANALYSIS FOR REVERBERANT SPEECH LOCALIZATION

Afsaneh Asaei<sup>1</sup>, Hervé Bourlard<sup>1,2</sup>, Mohammad J. Taghizadeh<sup>1,2</sup> and Volkan Cevher<sup>2</sup>

<sup>1</sup>Idiap Research Institute, Martigny, Switzerland

<sup>2</sup>Ecole Polytechnique Fédérale de Lausanne (EPFL), Switzerland

{afsaneh.asaei, herve.bourlard, mohammad.taghizadeh}@idiap.ch, volkan.cevher@epfl.ch

## ABSTRACT

This paper studies the problem of multiple speaker localization via speech separation based on model-based sparse recovery. We compare and contrast computational sparse optimization methods incorporating harmonicity and block structures as well as autoregressive dependencies underlying spectrographic representation of speech signals. The results demonstrate the effectiveness of block sparse Bayesian learning framework incorporating autoregressive correlations to achieve a highly accurate localization performance. Furthermore, significant improvement is achieved using ad-hoc microphones for data acquisition set-up compared to the compact microphone array.

**Index Terms**— Structured sparsity, Reverberant speech localization, Autoregressive modeling, Ad-hoc microphone array

## 1. INTRODUCTION

Speech localization in the clutter of voice and acoustic multipath is an active area of research on microphone arrays for hands-free speech communication. The accurate knowledge of the speaker location is essential for an effective beamforming steering and interference suppression [1, 2]. We briefly review the main approaches to address this problem.

*High Resolution Spectral Estimation:* These approaches are based on analysis of the received signals' covariance matrix and impose a stationarity assumption for accurate estimation [3]. Important techniques applied for speech localization include minimum variance spectral estimation as well as eigen-analysis methods such as multiple signal classification (MUSIC). The underlying hypotheses are not quite realistic in reverberant speech localization and alternative strategies have been usually considered [4, 5].

*Time Difference Of Arrival (TDOA) Estimation:* Another approach is based on TDOA estimation of the sources with respect to a pair of sensors. The generalized cross correlation (GCC) is the most common technique for TDOA estimation where the idea is basically to map the peak location of the cross-correlation function of the signal of two microphones to an angular spectrum. A weighting scheme is usually employed to increase the robustness of this approach to noise and multi-path effects. Maximum likelihood estimation of the weights has been considered as an optimal approach in the presence of uncorrelated noise, while the phase transform (PHAT) has been shown to be effective to overcome reverberation ambiguities [6, 7]. In addition to the GCC-PHAT, iden-

tification of the speaker-microphone acoustic channel has been incorporated for TDOA estimation and reverberant speech localization [8, 9]. However, despite of being practical and robust, TDOA-based techniques do not offer a high update rate. Alternative strategies have thus been sought for multiple-target tracking and adaptive beam-steering [10, 11].

*Beamformer Steered Response Power (SRP):* In this approach, the space is scanned by steering a microphone array beam-pattern and finding the direction associated to the maximum power. Delay-and-sum, minimum variance beamformers, and generalized sidelobe canceler have been the most effective methods for speaker localization [12]. The SRP-based approaches have a higher effective update rate compared to TDOA-based methods, and are applicable in multi-party scenarios using phase-transform weighting scheme [13].

In this paper, we adopt our speech separation framework using sparse component analysis [14, 15] and conduct the evaluations in terms of speech localization [16, 17]. We analyze the reverberant mixtures of speech signals in spectro-temporal domain. The planar area of the room is discretized into a dense grid such that the speakers are located at particular cells exclusively. A spatio-spectral sparse representation is obtained by concatenating the spectral components attributed to the sources located on the grid. The compressive acoustic measurements associated to the microphone array recordings are characterized using Image model of multipath propagation. The spatio-spectral sparse representation is estimated from the compressive array measurements using sparse optimization methods where the supports of high energy components indicate the source locations. The computational approaches to model-based sparse recovery of spectrographic speech are compared and contrasted considering block, harmonic as well as autoregressive dependencies.

The rest of the paper is organized as follows: Section 2 explains the premises underlying model-based sparse component analysis of reverberant recordings, and sets up the formulation of reverberant speech source localization. The structured sparsity models underlying speech components are elaborated in Section 3 followed by the computational approaches to model-based sparse recovery in Sections 4. Section 5 presents the details of the experiments. Conclusions are drawn in Section 6. The notations used in this paper are as follows

- ◇  $g \in \{1, \dots, G\}$ : number of a cell on a grids.
- ◇  $n \in \{1, \dots, N\}$ : number of source;  $N \ll G$ .
- ◇  $m \in \{1, \dots, M\}$ : number of microphones;  $M < N$ .
- ◇  $f \in \{1, \dots, F\}$ : number of spectral coefficients.
- ◇  $\{S, S\}$ : spectral representation of single/all source signals.
- ◇  $\{X, X\}$ : spectral representation of single/all micro. signals.
- ◇  $\Phi$ : microphone array manifold matrix.

AA and MT are supported by the Swiss National Science Foundation under the National Center of Competence in Research (NCCR) on Interactive Multi-modal Information Management (IM2). VC is supported in part by the European Commission under Grant MIRG-268398, ERC Future Proof and SNF 200021-132548.

## 2. SPARSE COMPONENT ANALYSIS OF REVERBERANT SPEECH MIXTURES

### 2.1. Spatio-Spectral Sparse Representation

We consider a scenario in which  $N$  speakers are distributed in a planar area spatially discretized into a grid of  $G$  cells. We assume to have a sufficiently dense grid so that each speaker is supposed to be located at the center of a cell, and  $N \ll G$ . The signals corresponding to each cell are concatenated to form a spatial representation of sources. Hence, the energy of the signals on the grid define a spatial spectrum with a sparse support denoting the location of the sources. We consider the spectro-temporal representation of speech signals and entangle the spatial representation of the sources with the spectral representation of the speech signal to form the complex vector  $\mathcal{S} = [S_1^T \dots S_G^T]^T \in \mathbb{C}^{G \times F \times 1}$  where  $\cdot^T$  corresponds to the transpose operator. Each  $S_g \in \mathbb{C}^{F \times 1}$  denotes the spectral representation of the  $g^{\text{th}}$  source (located at cell number  $g$ ) in Fourier domain. We express the signal ensemble at microphone array as a single vector  $\mathcal{X} = [X_1^T \dots X_M^T]^T$  where each  $X_m \in \mathbb{C}^{F \times 1}$  denotes the spectral representation of recorded signal at microphone number  $m$ . The sparse vector  $\mathcal{S}$  generates the underdetermined ( $M < G$ ) microphone mixture observations as  $\mathcal{X} = \Phi \mathcal{S}$  where  $\Phi$  is the microphone array measurement matrix consisted of the acoustic projections associated to the acquisition of the spatio-spectral sources.

### 2.2. Acoustic Measurement Characterization

We assume the room to be a rectangular enclosure consisting of finite impedance walls. The point source-to-microphone impulse responses are calculated using *Image Model* technique [18] where a reverberant signal is modeled as superposition of the signals attributed to the source images with respect to the reflective surfaces. Taking into account the physics of the signal propagation and multi-path effects, the projections associated with the source located on the cell  $g$  where  $\mathbf{v}_g$  represents the position of the center of the cell and captured by microphone  $m$  located at position  $\mu_m$  are characterized by the media Green's function through

$$\xi_{\mathbf{v}_g \rightarrow \mu_m}^f : \mathcal{X}(f) = \sum_{r=1}^R \frac{\tau_r}{\|\mu_m - \mathbf{v}_g^r\|^\alpha} \exp(-\sqrt{-1} \frac{\|\mu_m - \mathbf{v}_g^r\| f}{c}) S(f), \quad (1)$$

where  $\tau_r$  is the reflection ratio associated to the  $r^{\text{th}}$  image source located at  $\mathbf{v}_g^r$ . The attenuation constant  $\alpha$  depends on the nature of the propagation and is considered in our model to equal 1 which corresponds to the spherical propagation. This formulation assumes that if  $s_1(l) = s(l)$  and  $s_2(l) = s(l - \rho)$ , then  $S_2(f) \approx \exp(-j f \rho) S_1(f)$ .

Given the source-sensor projection defined in (1), we construct matrix  $\Xi_{\mathbf{v}_g \rightarrow \mu_m}$  for the measurement of the  $F$  consecutive frequencies as  $\Xi_{\mathbf{v}_g \rightarrow \mu_m} = \text{diag}(\xi_{\mathbf{v}_g \rightarrow \mu_m}^1 \dots \xi_{\mathbf{v}_g \rightarrow \mu_m}^F)$ . Hence, the projections associated to the acquisition of the source signals located on the grid by microphone  $m$  is  $\phi_m = [\Xi_{\mathbf{v}_1 \rightarrow \mu_m} \dots \Xi_{\mathbf{v}_g \rightarrow \mu_m} \dots \Xi_{\mathbf{v}_G \rightarrow \mu_m}]$  and the measurement matrix of  $M$ -channel microphone array is characterized as  $\Phi = [\phi_1 \dots \phi_M]^T$ . To fully identify this model, the location of the source images as well as the associated reflected ratios have been estimated and incorporated for sparse recovery of the reverberant speech signals  $\mathcal{S}$  [19]. We cast the underdetermined reverberant speech localization problem as sparse approximation where we exploit the underlying structure of the sparse coefficients for efficient recovery using fewer number of measurements [20, 14, 16]. The source locations are determined from the support of the high energy components of  $\mathcal{S}$  corresponding to the cells on the grid.

### 2.3. Computational Approaches to Sparse Recovery

Defining a set  $\mathbb{M}$  as the union of all vectors with a particular support structure, estimation of the sparse coefficient vector  $\hat{\mathcal{S}}$  from the microphone recordings  $\mathcal{X}$  can be expressed as

$$\hat{\mathcal{S}} = \underset{\mathcal{S} \in \mathbb{M}}{\text{argmin}} \|\mathcal{S}\|_0 \quad \text{s.t.} \quad \mathcal{X} = \Phi \mathcal{S} \quad (2)$$

where the counting function  $\|\cdot\|_0 : \mathbb{R}^G \rightarrow \mathbb{N}$  returns the number of non-zero components in its argument.

The major classes of computational techniques for solving sparse approximation problem are *Greedy pursuit*, *Convex optimization* and *Sparse Bayesian learning* [21].

*Greedy pursuit*: The nonzero components of  $\mathcal{S}$  are estimated in an iterative procedure by modifying one or several coefficients chosen to yield a substantial improvement in quality of the estimated signal. The present work considers an extension of the iterative hard thresholding [22, 23] to incorporate sparsity structures underlying spectrographic speech.

*Convex optimization*: The counting function in (2) is replaced with a sparsity inducing convex norm that exploits the structure underlying  $\mathcal{S}$ . Therefore, a convex objective is obtained which can be solved using convex optimization. The present work considers extension of basis pursuit algorithm which relies on  $L_1$  recovery [24].

*Sparse Bayesian learning*: A prior distribution is associated to  $\mathcal{S}$  with sparsity inducing hyperparameters and a maximum a posteriori estimation is derived. The present work considers the Bayesian framework proposed in [25, 26].

## 3. STRUCTURED SPARSITY MODELS

We consider three types of structures underlying the spectral coefficients: *harmonicity*, *block structure* as well as *AR dependency*. These structures are supported by the evidences from the studies on computational auditory scene analysis [27, 16].

*Harmonic structure* is exhibited if there are some interconnections between frequencies which are the harmonics of a fundamental frequency. In voiced speech, most of the energy in the speech signal occurs at harmonics of a fundamental frequency. The harmonicity model captures this structure as indicates that at any cell of the grid, energy is present in all frequencies that can be expressed as harmonics of a fundamental frequency. To state it more precisely, the support of vector  $\mathcal{S}$  is recovered imposing the structure of  $K$  harmonics of a fundamental frequency  $f_0$  defined as

$$\mathcal{F}_H \triangleq \{k f_0 | 1 < k < K\}, \quad (3)$$

*Block structure* is exhibited if some interconnections between the adjacent frequencies exist. In the case of vector  $\mathcal{S}$ , the block dependency model indicates that the spatial sparsity structure is the same at all neighboring discrete frequencies. In other words, a block of  $B$  consecutive frequencies corresponds to the same cell so the signal of the individual sources is recovered with a structure of independent blocks of size  $B$  defined as

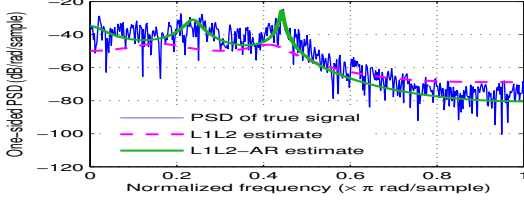
$$\mathcal{F}_B \triangleq \{[f_1, \dots, f_B], \dots, [f_{F-B+1}, \dots, f_F]\} \quad (4)$$

*AR dependency*: An additional inter-dependency is exhibited due to the correlation among the block entries corresponding to each source, which we model using an auto regressive (AR) process of order  $\mathcal{R}$  characterized by the following model

$$\mathcal{F}_{AR} \triangleq [1, \beta_g(1), \beta_g(2), \dots, \beta_g(\mathcal{R})] \quad (5)$$

where  $\beta_g \in (-1, 1)$  denotes the AR coefficients. The sources  $\mathcal{S}_g$  are mutually independent, but each source satisfies an AR model as

$$\mathcal{S}_g(\mathbf{b}) = \mathcal{F}_{AR} [\mathbf{u}(\mathbf{b}), \mathcal{S}_g(\mathbf{b} - 1), \dots, \mathcal{S}_g(\mathbf{b} - \mathcal{R})]^T, \quad \mathbf{b} \in \{1, \dots, B\} \quad (6)$$



**Fig. 1:** Incorporating AR dependencies in basis pursuit sparse recovery.

where  $\mathbf{u}(\mathbf{b})$  denotes an input sequence. From (6) we can see that the covariance matrix  $\mathcal{B}_g$  of each source is a Toeplitz matrix identified by the AR coefficients (5).

#### 4. STRUCTURED SPARSE RECOVERY

We consider different model-based sparse recovery algorithms to recover the sparse vector incorporating the structures defined above. In particular, we employ Iterative hard thresholding *IHT* [28],  $L_1L_2$  convex optimization [24] as well as Block Sparse Bayesian Learning framework, *BSBL* [26].

*IHT*: Iterative hard thresholding (IHT) offers a simple yet effective approach to estimate the sparse vectors. It seeks an  $N$ -sparse approximation  $\hat{\mathcal{S}}$  matching the observation  $\mathcal{X}$  by minimizing the residual error. We use the algorithm proposed in [23] which is an accelerated scheme for hard thresholding methods with the following recursion

$$\begin{aligned} \hat{\mathcal{S}}^0 &= 0, \quad \mathcal{R}^i = \mathcal{X} - \Phi \hat{\mathcal{S}}^i \\ \hat{\mathcal{S}}^{i+1} &= \mathcal{M}^{\mathcal{F}}(\hat{\mathcal{S}}^i + \kappa \Phi^T \mathcal{R}^i) \end{aligned} \quad (7)$$

where the step-size  $\kappa$  is the Lipschitz gradient constant to guarantee the fastest convergence speed. To incorporate for the underlying structure of the sparse coefficients, the model approximation operator  $\mathcal{M}^{\mathcal{F}}$  is defined as reweighting and thresholding the energy of the components of  $\hat{\mathcal{S}}$  with either  $\mathcal{F}_B$  or  $\mathcal{F}_H$  structures.

$L_1L_2$ : Another fundamental approach to sparse approximation replaces the combinatorial counting function in the mathematical formulation stated in (2) with the  $L_1$  norm, resulting in a convex optimization problem that admits a tractable algorithm referred to as basis pursuit [24]. We use a group version of basis pursuit algorithm with the number of group components  $n^{\mathcal{F}}$  determined by each structure. The optimization problem to recover the block sparse coefficients  $\hat{\mathcal{S}}$  is formulated as follows:

$$\hat{\mathcal{S}} = \underset{\mathcal{S}}{\operatorname{argmin}} \{ \|\mathcal{S}\|_{L_1, L_2} \quad \text{s.t.} \quad \mathcal{X} = \Phi \mathcal{S}, \|\mathcal{S}\|_{L_1, L_2} = \sum_{g=1}^G \sqrt{\sum_{b=1}^{n^{\mathcal{F}}} (S_g(\mathbf{b}))^2} \} \quad (8)$$

To incorporate the AR dependencies of the block coefficients of  $\mathcal{S}$ ,  $\mathcal{X} = \tilde{\Phi} \mathcal{S}$  is solved by (8) where  $\tilde{\Phi}$  constitutes of  $\tilde{\mathcal{S}}_{v_g \rightarrow \mu_m}$  where the diagonal elements are multiplied by  $\mathcal{F}_{A, R}$ . Fig. 1 demonstrates an example of an AR signal of order 4 recovered using the proposed procedure. More details are discussed in Section 5.2.

*BSBL*: The correlation among the coefficients modeled as AR dependencies is incorporated by [26] in the framework of SBL [25]. The sources  $S_g$  are assumed to be Gaussian and mutually independent. The AR dependency model indicates that the linear combination of the univariate Gaussian holds a Gaussian distribution. More precisely, the joint distribution of  $S_g = [S_g^1, \dots, S_g^B]$  is a multivariate Gaussian, expressed by  $p(S_g; \gamma_g, \beta_g) \sim \mathcal{N}(0, \gamma_g \mathcal{B}_g)$ , where  $\gamma_g$  is a non-negative hyper-parameter controlling the block-sparsity of  $\mathcal{S}$  and  $\mathcal{B}_g \in \mathbb{R}^{B \times B}$  is a positive definite matrix that captures the correlation structure of  $S_g$  as defined in (6). Under the assumption that blocks are mutually uncorrelated, the prior

for  $\mathcal{S}$  is given by  $p(\mathcal{S}; \gamma_g, \mathcal{B}_g, \forall g) \sim \mathcal{N}(0, \Sigma_0)$ , where  $\Sigma_0 = \operatorname{diag}([\gamma_1 \mathcal{B}_1 \dots \gamma_G \mathcal{B}_G])$ . Assuming the Gaussian likelihood for the block sparse model as  $p(\mathcal{X}|\mathcal{S}; \sigma^2) \sim \mathcal{N}(\Phi \mathcal{S}, \sigma^2 \mathbf{I})$  and applying the Bayes rule, we obtain the posterior density of  $\mathcal{S}$ , which is also Gaussian,  $p(\mathcal{S}|\mathcal{X}; \sigma^2, \{\gamma_g, \mathcal{B}_g\}_{g=1}^G) = \mathcal{N}(\mu_s, \Sigma_s)$  with the covariance matrix  $\Sigma_s = (\Sigma_0^{-1} + \frac{1}{\sigma^2} \Phi^T \mathbf{I} \Phi)^{-1} \mathcal{X}$ . Having all the hyper-parameters  $\sigma^2$ ,  $\gamma_g$ ,  $\mathcal{B}_g$ , the MAP estimate of  $\mathcal{S}$  is given by the mean defined as [26]

$$\hat{\mathcal{S}} \triangleq \mu_s = \Sigma_0 \Phi^T (\sigma^2 \mathbf{I} + \Phi \Sigma_0 \Phi^T)^{-1} \mathcal{X}, \quad (9)$$

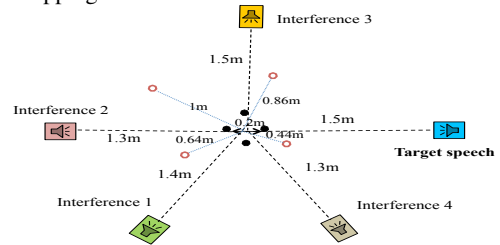
Clearly, the block sparsity of  $\hat{\mathcal{S}}$  is controlled by  $\gamma_g$  in  $\Sigma_0$ . During the estimation procedure,  $\gamma_g = 0$  indicates that the associated block in  $\hat{\mathcal{S}}$  is zeros and no source is located on the corresponding cell. The framework proposed in [29], derives the EM-based learning rule to learn the hyperparameters. We will see in Section 5.2 that the AR-dependency matrix can be estimated offline for the specific task of speech localization.

#### 5. EXPERIMENTAL STUDY

The experiments are conducted to quantify the performance of different structured sparse recovery algorithms on different microphone array geometric settings in terms of speech localization accuracy.

##### 5.1. Acoustic and Analysis Setup

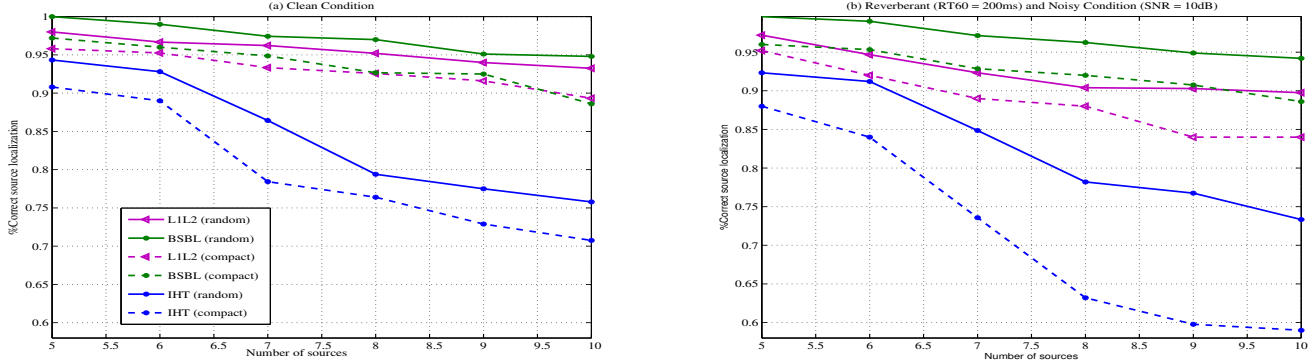
The overlapping speech was synthesized by mixing speech utterances taken from the Wall Street Journal (WSJ) corpus [30]. The WSJ corpus is a 20000-word corpus consisting of read Wall Street Journal sentences. The sentences are read by a range of speakers (34 in total) with varying accents. All the files are normalized prior to mixing. The microphone array recording set-up is consisted of four channels microphones. The planar area of the room with dimension  $3\text{m} \times 3\text{m} \times 3\text{m}$  is divided into cells with 50 cm spacing. The data collection setup is depicted in Figure 2. The scenarios include *random* and *compact* topologies of microphone array in clean as well as reverberant and noisy conditions. Room impulse responses are generated with the Image model technique [18] using intra-sample interpolation, up to 15<sup>th</sup> order reflections and omni-directional microphones for a room reverberation time equal to 180 ms. The number of source is known in our experiments. The speech signals of length one second are recorded at 16 kHz sampling frequency and the spectro-temporal representation for source separation is obtained by windowing the signal in 250 ms frames using Hann function with 50% overlapping.



**Fig. 2:** Overhead view of the room set-up for uniform (black dots) and random microphone array (red dots)

##### 5.2. Speech Localization Performance

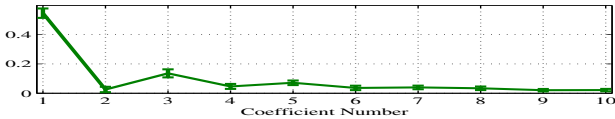
The probabilistic performance bounds of multi-speaker localization are obtained by averaging the results over an exhaustive and exclusive set of configurations. The results are evaluated over all configurations consisted of  $N \in \{5 - 10\}$  sources. The probabilistic evaluations are necessary to form a realistic expectation of our sparse recovery framework as the deterministic performance bounds



**Fig. 3:** Speaker localization performance evaluated for 5-10 sources exploiting spectral structured sparsity models

are derived for the worst case scenario which is not likely to occur [31, 32]. The localization accuracy is measured as the number of times sources are localized correctly (the support of the recovered signal corresponds to the cell on the grid where the source is located) divided by the number of all sources.

The block sparse Bayesian learning (BSBL) algorithm can learn the AR parameters during the optimization, although, the procedure is very expensive in terms of computational cost. Hence, we carry out some studies on an average AR model for speech signal which can be exploited for source localization. To estimate the AR coefficients, the frequency band (number-of-FFT-points =  $2048 \times 4$ ) is split into blocks of size 16 and processed independently. Fig. 4 illustrates the frequency domain average AR model for 10 min speech signal. The first-order coefficient is estimated as 0.45. We can see that the higher order coefficients are small so the blocks are modeled as a first-order AR process to incorporate the intra-block correlation.



**Fig. 4:** 10-order AR coefficients estimated for 10min speech signal. The cross lines illustrate the variance of estimates

The results of multi-speaker localization exploiting spectral structure are illustrated in Figure 3 for  $B = 16$ . All the algorithms are run for the stopping threshold fixed to  $1e-2$  and the maximum iteration of 150. The BSBL algorithm assumes that all sources have similar correlation structure. The experimental analysis on speech-specific average AR model as depicted in Fig. 4 shows very small variance around the AR coefficients and supports this approximation. We can see that exploiting the frequency structures yield very strong results. The number of microphones is only 4 whereas we can localize up to 9 sources with 95% accuracy. The orthogonality or disjointness of spectrographic speech signals is a key property to achieve this bound of performance [14, 33].

The BSBL algorithm employs the AR dependency model to replace their row norms with Mahalanobis distance measure and it plays a role of temporally whitening the sources during the learning of hyperparameters [29]. On the other hand, incorporation of AR model in the framework of  $L_1L_2$  enables preserving the structural dependencies. As a basic example, an AR signal is generated by filtering a white Gaussian noise. The formulation of the  $L_1L_2$ -AR enables recovery of the input signal  $u$  (6) along with the signal coefficients. The speech signal can be constructed by filtering the recovered  $u$  [34] while the AR model parameters can be estimated

from the initial estimates of block sparse recovery and refined in an iterative manner. We can see that AR dependency is better preserved using the proposed procedure as illustrated in Fig. 1. However, this approach did not outperform the standard basis pursuit in terms of speech localization. Furthermore, the results of the harmonic sparse recovery were comparable to the block-sparse recovery, hence they are not further elaborated here [16, 35].

The other important observation is that the ad-hoc layout of microphone array improves the results for all sparse recovery algorithms. It can be justified as the theoretical analysis of the performance bounds of sparse recovery algorithms is entangled with the spectral properties of  $\Phi$ . A key property to guarantee the theoretical performance bounds is the coherence  $\vartheta$  of the measurement matrix defined as the smallest angle between any pairs of the columns of  $\Phi$ . The number  $N$  of recoverable non-zero coefficients using either convexified or greedy sparse recovery is inversely proportional to the coherence as  $N < \frac{1}{2}(\vartheta^{-1} + 1)$  [21]. Therefore, to guarantee the performance of sparse recovery algorithms, it is desired to minimize the coherence. As the measurement matrix is constructed of the location-dependent projections, this property implies that the performance of our localization framework is entangled with the microphone array layout. A large-aperture random design of microphone array yields the projections to be mutually incoherent, so the projections are spread across all the acoustic scene and each microphone captures the information about all components of  $S$  [36]. Furthermore, the coherence of the acoustic measurements is smaller at the high frequencies of the broadband speech spectrum, hence, the bands below 100Hz are discarded from our localization scheme [16].

## 6. CONCLUSIONS

In this paper, we incorporated the speech-specific models for structured sparse recovery of reverberant speech sources. We outlined the fundamental computational approaches to model-based sparse recovery and evaluated their performance in terms of source localization accuracy. The numerical assessments show the block sparse Bayesian learning framework yields the best performance and an average AR model can be learned for speaker localization and specified to the algorithm to reduce the computational cost. Furthermore, we considered the impact of construction layout of the microphone array in the performance of sparse recovery algorithms. The theoretical insights suggest that an ad-hoc design of microphone array can better preserve the acoustic information by reducing the coherence of the acoustic measurements. The empirical evaluations confirm that considering the design specifications acknowledged by the generic theory of sparse signal recovery leads to significant improvement in speech localization performance.

## References

- [1] A. Waibel, M. Bett, F. Metzger, K. Ries, T. Schaaf, T. Schultz, H. Soltan, H. Yu, and K. Zechner, "Advances in automatic meeting record creation and access," in *Proc. of ICASSP*, 2001.
- [2] A. Asaei, M. J. Taghizadeh, M. Bahrololum, and M. Ghanbari, "Verified speaker localization utilizing voicing level in split-bands," *Signal Processing*, vol. 89(6), 2009.
- [3] J. Dmochowski, S. Benesty, and S. Affes, "Broadband music: opportunities and challenges for multiple source localization," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2007.
- [4] S. Araki, H. Sawada, R. Mukai, and S. Makino, "DOA estimation for multiple sparse sources with arbitrarily arranged multiple sensors," *Journal of Signal Processing Systems*, vol. 63(2), 2011.
- [5] K. Kumatani, J. McDonough, and B. Raj, "Microphone array processing for distant speech recognition: From close-talking microphones to far-field sensors," *IEEE Signal Processing Magazine, Special Issue on Fundamental Technologies in Modern Speech Recognition*, 2012.
- [6] M. Omologo and P. Svaizer, "Acoustic source localization in noisy and reverberant environments using CSP analysis," in *Proc. of ICASSP*, 1996.
- [7] C. Blandin, A. Ozerov, and E. Vincent, "Multi-source TDOA estimation in reverberant audio using angular spectra and clustering," *Signal Processing*, vol. 92, 2012.
- [8] F. Ribeir, C. Zhang, D. Florencio, and D. Ba, "Using reverberation to improve range and elevation discrimination in sound source localization," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18(7), 2010.
- [9] S. Nam and R. Gribonval, "Physics-driven structured cosparsity modeling for source localization," in *Proc. of ICASSP*, 2012.
- [10] M. S. Brandstein and H. F. Silverman, "A practical methodology for speech source localization with microphone arrays," *Computer Speech and Language*, vol. 11(2), 1997.
- [11] M. Omologo, F. Nesta, P. Svaizer, "Cumulative state coherence transform for a robust two-channel multiple source localization," in *Int. Conf. on Independent Component Analysis and Signal Separation (ICA)*, 2009.
- [12] J. H. DiBiase, *A High-Accuracy, Low-Latency Technique for Talker Localization in Reverberant Environments Using Microphone Arrays*, Ph.D. thesis, 1993.
- [13] M. J. Taghizadeh, P. N. Garner, H. Bourlard, H. R. Abutalebi, and A. Asaei, "An integrated framework for multi-channel multi-source speaker localization and source activity detection," in *Proceedings of Hands-free Speech Communication and Microphone Arrays (HSCMA)*, 2011.
- [14] A. Asaei, H. Bourlard, and V. Cevher, "Model-based compressive sensing for distant multi-party speech recognition," in *Proc. of ICASSP*, 2011.
- [15] A. Asaei, M. J. Taghizadeh, H. Bourlard, and V. Cevher, "Multi-party speech recovery exploiting structured sparsity models," in *Proceeding of INTERSPEECH*, 2011.
- [16] Afsaneh Asaei, *Model-based Sparse Component Analysis for Multi-party Distant Speech Recognition*, Ph.D. thesis, Ecole Polytechnique Federal de Lausanne (EPFL), 2013.
- [17] A. Asaei, H. Bourlard, and V. Cevher, "A method, apparatus and computer program for determining the location of a plurality of speech sources," *2012US-13/654055, US Patent, October*, 2012.
- [18] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," *Journal of Acoustical Society of America*, vol. 60(s1), 1979.
- [19] A. Asaei, M. Golbabaee, H. Bourlard, and V. Cevher, "Structured sparsity models for multiparty speech recovery from reverberant recordings," *IEEE Trans. on Speech and Audio Processing (accepted)*, 2013.
- [20] R. G. Baraniuk, V. Cevher, M. F. Duarte, and C. Hegde, "Model-based compressive sensing," *IEEE Transactions in Information Theory*, 2010.
- [21] J. A. Tropp and S. J. Wright, "Computational methods for sparse solution of linear inverse problems," *Proceedings of the IEEE*, vol. 98, 2010.
- [22] T. Blumensath and M. E. Davies, "Iterative hard thresholding for compressed sensing," *Applied and Computational Harmonic Analysis*, vol. 27(3), pp. 265–274, 2009.
- [23] A. Kyrillidis and V. Cevher, "Recipes on hard thresholding methods," in *Proceedings of CAMSAP*, 2011.
- [24] E. V. D. Berg and M. P. Friedlander, "Probing the pareto frontier for basis pursuit solutions," *SIAM Journal on Scientific Computing*, 2008, <http://www.cs.ubc.ca/labs/scl/spg11>.
- [25] D. P. Wipf and B. D. Rao, "An empirical bayesian strategy for solving the simultaneous sparse approximation problem," *IEEE Transactions on Signal Processing*, vol. 55(7), 2007.
- [26] Z. Zhang and B. D. Rao, "Extension of SBL algorithms for the recovery of block sparse signals with intra-block correlation," *IEEE Transactions on Signal Processing*, 2012.
- [27] D. Wang and G. J. Brown, *Computational Auditory Scene Analysis: Principles, Algorithms, and Applications*, Wiley-IEEE Press, 2006.
- [28] T. Blumensath and M. E. Davies, "Gradient pursuits," *IEEE Transactions on Signal Processing*, vol. 56, pp. 2370–2382, 2008.
- [29] Z. Zhang and B. D. Rao, "Sparse signal recovery with temporally correlated source vectors using sparse bayesian learning," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5(5), 2011.
- [30] M. Lincoln, I. McCowan, I. Vepa, and H. K. Maganti, "The multi-channel wall street journal audio visual corpus (mc-wsj-av): Specification and initial experiments," in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2005.
- [31] P. Boufounos, P. Smaragdis, and B. Raj, "Joint sparsity models for wideband array processing," in *Wavelets and Sparsity XIV, SPIE Optics and Photonics*, 2011.
- [32] J. Le Roux, P. T. Boufounos, K. Kang, and J. R. Hershey, "Source localization in reverberant environments using sparse optimization," in *Proc. of ICASSP*, 2013.
- [33] A. Asaei, H. Bourlard, and P. N. Garner, "Sparse component analysis for speech recognition in multi-speaker environment," in *Proceeding of INTERSPEECH*, 2010.
- [34] B. Yegnanarayana, C. Avendano, H. Hermansky, and P. S. Murthy, "Speech enhancement using linear prediction residual," *Speech Communication*, vol. 28(1), 1999.
- [35] A. Asaei, M. Davies, H. Bourlard, and V. Cevher, "Computational methods for structured sparse recovery of convolutive speech mixtures," in *Proc. of ICASSP*, 2012.
- [36] L. Carin, "On the relationship between compressive sensing and random sensor arrays," *IEEE Antennas and Propagation Magazine*, vol. 51, pp. 72–81, 2009.