

# Speech Sources Separation Based on Models of Interaural Parameters and Spatial Properties of Room

Muhammad Israr<sup>1</sup>, Muhammad Salman Khan<sup>2</sup>, Khushal Khan<sup>3</sup>

<sup>1,2,3</sup> Department of Electrical Engineering, University of Engineering and Technology, Peshawar, Pakistan  
muhammad.israr@uetpeshawar.edu.pk<sup>1</sup>

Received: 07 January, Revised: 16 January, Accepted: 18 January

**Abstract**— This paper presents the extended evaluation in different reverberant scenarios for different mixtures of speech having interfere at one of the six angles  $\{15^{\circ}, 30^{\circ}, 45^{\circ}, 60^{\circ}, 75^{\circ}, 90^{\circ}\}$  of the model which implementing spatial covariance with interaural parameter for improving the performance of MESSL. The binaural spatial parameters, such as interaural phase difference IPD and interaural level difference ILD and spatial covariance are modeled in the short-time Fourier transform. The parameters of the model are updated with the expectation-maximization algorithm. The performance of the model is checked in term of Signal-to-distortion ratio (SDR) and the perceptual evaluation of speech quality (PESQ), and the results confirmed that the performance of this proposed model is improved in highly reverberant rooms.

**Keywords**— MESSL, Spatial covariance, SDR, PESQ, MESSL+SC

## I. INTRODUCTION

Humans are mostly interested to focus their attention only on the speech of a single speaker, while other speakers and background noise is exist [1]. This ability is greatly reduced when someone is listening with one ear, especially in reverberant environment. For example in automatic speech recognition the machine focus on a speaker of interest and separate speaker from background noise[2][3]. The process of automated separation of sources from measured mixtures without any prior information is known as blind source separation (BSS)[4]. In blind sources separation the sound sources are separated from the mixtures in which both the numbers of mixtures and sources are unknown and only mixtures signals are available. In some situations we want to recover all individual sources from the mixture signal, or at least a specific source[5]. In laboratory conditions most of the algorithms perform fine where number of sources present in mixtures, mixing methodology and acoustic condition are already known before the separation of sources from mixtures. But in real life scenario this problem is much more complex.

The solution for source separation problem is proposed by researchers belongs to different communities[6][7][8][9]. Examples included convolution in frequency domain, blind source separation (BSS), beamforming and computational auditory scene analysis (CASA). Underdetermined problem in sound sources separation, where the number of sources is more than the number of sensors can

be handle through TF approach[10]. In timefrequency (TF) masking for source separation, by exploiting these cues, the mask for each source can be obtained and hoped that only a single source is active at each TF unit. CASA is a “machine listening” system that separate mixtures of sound sources in the same way that human listeners do. It uses cues like pitch, onset/offset time, interaural level difference (ILD) and interaural time difference (ITD) to separate the sources[11].

In this thesis speech sources are separated by using a CASA approach which perform sources segregation on the base of interaural parameters such as interaural level difference (ILD), interaural phase difference (IPD) and spatial properties of the room. The model parameters are determined by using the EM algorithm.

The observation and initial values of the parameters are used for estimation of probabilities in E-step. These parameters for the model are updated based on the measurements and the posterior probabilities from the E-step. The proposed model produces probabilistic TF masks for single sources that are used for their reconstruction.

## II. MOTIVATION AND OBJECTIVE OF THE WORK

The main motivation in Blind Source Separation (BSS) is the cocktail party problem. Imagine someone is in a loud party, he is trying to convey his voice message to his friend. During this conversation the sounds received at the listener’s ear is mixture of the voice of different musical instruments, people and so on. The signal arrived to the ear is in single wave form containing all sounds, and still listener is able to understand talk from this mixture[2][1].

For music students the sound source separation is an important tool by which he can separate a single instrument if many instruments are active simultaneously. Similarly, the ability to remove an instrument from recording is also useful, as it would allow a student to play a different instrument in the original track in place of the removed instrument[12][13].

## III. OVERVIEW OF PROPOSED ALGORITHM

Some statistical properties of the sources provide a base for separation of sound sources[14]. The common statistical assumption made in our model are that the sources are

statistically independent, statistically orthogonal and nonstationary [7][5][15].

The reverberant mixtures at left and right microphones are represented as  $l(t) = \sum_{i=1}^I s_i(t) * h_{li}(t)$ , and  $r(t) = \sum_{i=1}^I s_i(t) * h_{ri}(t)$ . In this expression  $s_i$  show sources,  $h_{li}$  represent the Room impulse response RIR from the source  $s_i$  to the left microphone,  $h_{ri}$  represent the Room impulse response RIR from the source  $s_i$  to the right sensor, and "\*" show the convolution.

### A. The ILD, IPD and spatial covariance models

The interaural spectrogram, the ratio of left and right mixture at every TF point is given by equation  $\frac{L(t,\omega)}{R(t,\omega)} = 10^{\frac{\alpha(t,\omega)}{20}} e^{j\varphi(t,\omega)}$  where  $\alpha(t,\omega)$  is the ILD measured in dB while  $\varphi(t,\omega)$  is IPD at frequency  $\omega$  and "t" is time. The IPD must exist in the range of  $[-\pi, \pi]$  [16]. For ILD and IPD each source are frequency-dependent and is given as  $\tau(\omega)$ , and  $\alpha(\omega)$  [15]. The model requires that Fourier transform window used, 1024 points (64ms) must be larger than maximum ITD of  $\approx 0.75$ ms. We use a top down approach where ITD is measured in term of IPD. The phase residual error, which is defined as difference between IPD observed and IPD predicted is given by equation  $\varphi(\omega, t, \tau) = \frac{L(t,\omega)}{R(t,\omega)} e^{-j\omega\tau}$ . The phase residual error is modeled with a normal distribution with frequency-dependent mean  $\xi(\omega)$  and variance  $\sigma^2(\omega)$  [17],

$$p(\varphi(\omega, t) | \tau(\omega), \sigma(\omega)) = N(\varphi(\omega, t; \tau) | \xi(\omega), \sigma^2(\omega)).$$

The ILD is also modeled with a normal distribution with mean  $\mu(\omega)$  and variance

$$\eta^2(\omega), p(\alpha(\omega, t) | \mu(\omega), \eta^2(\omega)) = N(\alpha(\omega, t) | \mu(\omega), \eta^2(\omega)).$$

Signal  $x(\omega, t)$ , is mixture signal consisting spatial images of  $I$  sources mixed with each other in present each channel. This can be easily modeled as a zero-mean Gaussian distribution with the covariance matrix [16].

$$R_x(\omega, t) = \sum_{i=1}^I v_i(\omega, t) R_i(\omega)$$

In above equation  $v_i(\omega, t)$  show scalar variance while  $R_i(\omega)$  is the covariance matrix containing the spatial properties of the source  $i$ . [18][19] The probability distribution of the proposed method is given by [16].

$$P(x(\omega, t) | \{v_i(\omega, t), R_i(\omega), \forall i\}) = \frac{1}{\det(\pi R_x(\omega, t))} \exp(-x^H(\omega, t) R_x^{-1}(\omega, t) x(\omega, t))$$

Where  $(\cdot)^H$  is the Hermitian transpose.

### B. Expectation maximization and source separation

The interaural time difference (ILD), interaural phase difference (IPD) discussed in MESSL model [1] and the spatial covariance models, are combined [6] in this paper and new model has been formed which is combination of MESSL and Spatial covariance model. For determining the parameters of

the models the expectation maximization algorithm is used for its solution. Their corresponding parameters as

$$p(\alpha(\omega, t), \varphi(\omega, t), x(\omega, t) | \theta) = N(\alpha(\omega, t) | \mu(\omega), \eta^2(\omega)) \cdot N(\varphi(\omega, t) | \xi(\omega), \sigma^2(\omega)) \cdot N(x(\omega, t) | 0, R_x(\omega, t))$$

All the model's parameter is given in a vector represented by  $\theta$  and its equation is given as

$$\theta = \{\mu_i(\omega), \eta^2_i(\omega), \xi_{i\tau}(\omega), \sigma^2_{i\tau}(\omega), v_i(\omega, t), \psi_{i\tau}\}$$

Where  $\mu_i$ ,  $\xi_{i\tau}$ , and  $\eta^2_i$ ,  $\sigma^2_{i\tau}$  are respectively the means and variances of the ILD, IPD models, and  $v_i$  is the scalar variance related to the spatial covariance model. The covariance matrix  $R_i(\omega)$  required to calculate  $R_x(\omega)$  is found by utilizing the posterior knowledge about the room properties. The notation  $i$  in subscript in all parameters represent that they belong to source  $i$ , and  $\tau$  and  $\omega$  show that these quantities depends on delay and frequency. The parameter  $\psi_{i\tau}$  is the mixing weight, is used for determining the probability of the point in TF spectrum belong to source  $i$  and the delay  $\tau$ . [6]

The log likelihood function ( $L$ ) given the observations can be written as

$$L(\theta) = \sum_{t,\omega} \log p(\alpha(\omega, t), \varphi(\omega, t), x(\omega, t) | \theta) = \sum_{t,\omega} \log \sum_{i,\tau} [ N(\alpha(\omega, t) | \mu_i(\omega), \eta^2_i(\omega)) \cdot N(\varphi(\omega, t; \tau) | \xi_{i\tau}(\omega), \sigma^2_{i\tau}(\omega)) \cdot N(x(\omega, t) | 0, R_x(\omega, t)) \cdot \psi_{i\tau} ]$$

And the maximum likelihood solution a vector of parameters moves the function toward maximum value. The EM algorithm is initialized with the estimated locations of the speakers. In the expectation step the probabilities are calculated given the observations and the estimates of the parameters as

$$\varepsilon_{i\tau}(\omega, t) = \psi_{i\tau} \cdot N(\alpha(\omega, t) | \mu_i(\omega), \eta^2_i(\omega)) \cdot N(\varphi(\omega, t; \tau) | \xi_{i\tau}(\omega), \sigma^2_{i\tau}(\omega)) \cdot N(x(\omega, t) | 0, R_x(\omega, t)),$$

Where  $\varepsilon_{i\tau}(\omega, t)$  is a variable which provide the information of expectation belonging to a given source of hidden variable  $m_{i\tau}(\omega, t)$ , whose value is one in case time frequency point belong to source  $i$  and delay  $\tau$  and zero otherwise. In  $M$ -step, the updating of parameters take place using previous observation of  $E$ -step. The IPD and ILD parameters and  $\psi_{i\tau}$  are re-estimated as in.

The contribution of spatial covariance model are included from 2nd iteration, because in 1st iteration the  $\varepsilon_{i\tau}$  is calculated only through ILD and IPD models, as it is dominate at the corresponding TF point for source  $i$  with delay  $\tau$ . The time frequency mask for each source is given through equation.

$$M_i(\omega, t) = \sum_{\tau} \varepsilon_{it}(\omega, \tau)$$

The masks are multiplied to mixtures and single source of interest are reconstructed. We next experimentally verify the efficiency of the proposed model through the MATLAB simulation.

#### IV. EXPERIMENTS AND RESULTS

Experiments were performed on mixtures of two sources with varying levels of reverberation and different angles. The target source was present at 00 azimuth and interferer were present at one of the six angles { 150, 300, 450, 600, 750, 900 }. Speech utterances were randomly chosen from the TIMIT acoustic phonetic continuous speech corpus to form mixtures with two speech sources. The results obtained by performing simulation in different scenarios such as different angles, RIR and RT<sub>60</sub>s for ten mixtures each containing two sources. The signal to distortion ratio (SDR) and perceptual evaluation of speech quality (PESQ) were used to evaluate the performance of the algorithms. The proposed model is named as (MESSL+SC) in all graphs and figures.

The average SDR and average PESQ is calculated for ten different two-sources mixtures and for both proposed and MESSL models at different RT<sub>60</sub>s i.e. 320ms, 470ms, 680ms and 890ms.

Figure. 1 and 2 show the average SDR and average PESQ comparisons of the proposed model with MESSL model at RT<sub>60</sub>s of 320ms.

Text heads organize the topics on a relational, hierarchical

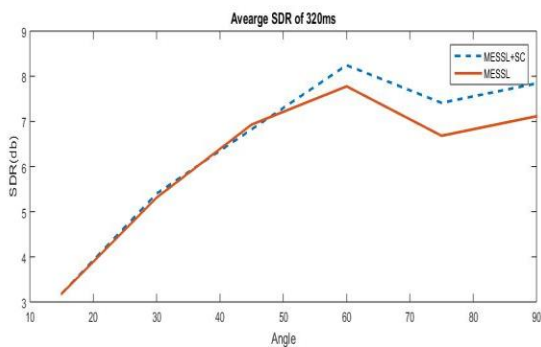


Figure.1 SDR Comparison of proposed model and MESSL Model at Reverberation of 320ms

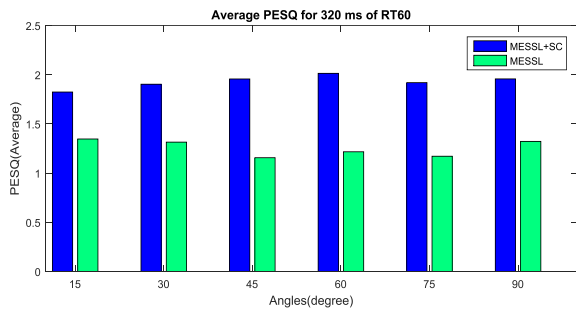


Figure. 2 PESQ Comparison of proposed model and MESSL Model at Reverberation of 320ms

The figure. 1 and 2 shows that at low reverberation the performance of proposed model and MESSL are same for average SDR and better than MESSL in the average PESQ.

Figure. 3 and 4 show the average SDR and average PESQ comparisons of the proposed model with MESSL model at RT<sub>60</sub>s of 470ms.

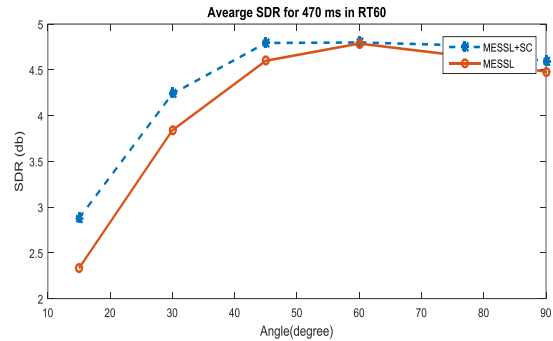


Figure. 3 SDR Comparison of proposed model and MESSL Model at Reverberation of 470ms

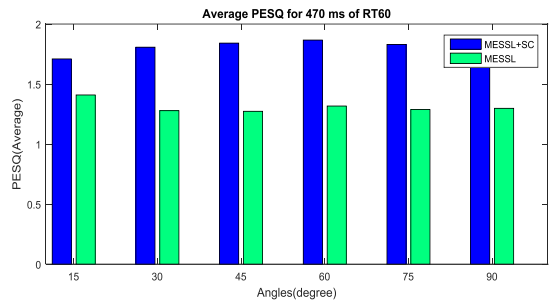


Figure. 4 PESQ Comparison of proposed model and MESSL Model at Reverberation of 470ms

The figure 3 and 4 clarifying that when reverberation increases from 320ms to 470ms the performance of the proposed model is enhanced from MESSL model in both evaluation matrix i.e. average SDR and average PESQ.

Figure.5 – 9 show the average SDR and average PESQ comparisons of the proposed model with MESSL model at RT<sub>60</sub>s of 680ms and 890ms.

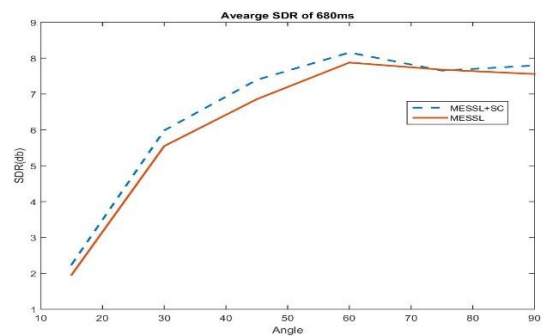


Figure.5 SDR Comparison of proposed model and MESSL Model at Reverberation of 680ms

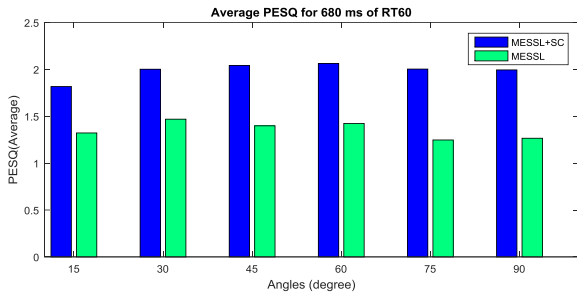


Figure.6 PESQ Comparison of proposed model and MESSL Model at Reverberation of 680ms

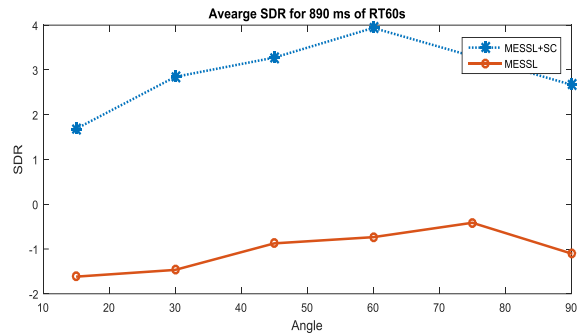


Figure.7 SDR Comparison of proposed model and MESSL Model at Reverberation of 890ms

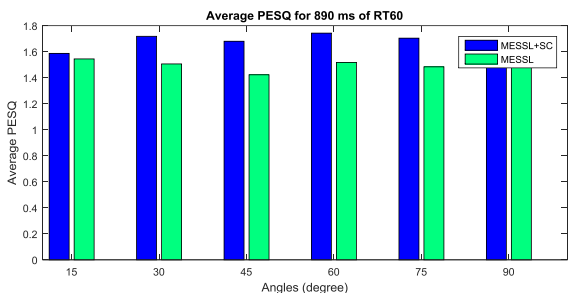


Figure.8 PESQ Comparison of proposed model and MESSL Model at Reverberation of 890ms

The above figures 5-8 clarify that whenever the reverberation time is further increased up to 680ms and 890ms the performance of the proposed model is improved more over the MESSL model. Figure. 8 shows that at the highest RT<sub>60</sub>s(890ms) under consideration, the proposed model outperforms MESSL by and Average SDR of 3.05 db.

#### CONCLUSION

A speech separation algorithm in reverberant environment is proposed, that integrate the models of the interaural parameters and spatial covariance. The observed mixtures were modeled and initialized using the EMA algorithm giving improved source estimates. Results indicate that the performance of our proposed model is much better than the existing model MESSL.

#### FUTURE WORK

The algorithm proposed in the paper can be extended in many ways. In the proposed model all the sources of sound are assumed to be stationary sources. In practical scenarios this

is not possible that all of the time the sources will be stationary they can change their positions. So case of movable sources is future challenges of this algorithm.

By applying the effect of covariance model with models of IPD and ILD assumed that the reverberation time was known. Since the estimation for finding the reverberation time was not the aim of this thesis. The future work in this thesis is that to include an algorithm to detect the reverberation time i.e. by spectral subtraction. Finally, in future work researchers could focus on the algorithm complexity reduction.

#### REFERENCES

- [1] M. I. Mandel, R. J. Weiss, and D. P. W. Ellis, "Model-Based Expectation-Maximization Source Separation and Localization," *IEEE Trans. Audio, Speech Lang. Process.*, vol. 18, no. 2, pp. 382–394, 2010.
- [2] A. Ephrat *et al.*, "Looking to Listen at the Cocktail Party: A Speaker-Independent Audio-Visual Model for Speech Separation," vol. 37, no. 4, 2018.
- [3] N. Roman, D. Wang, and G. J. Brown, "Speech segregation based on sound localization," *J. Acoust. Soc. Am.*, vol. 114, no. 4, pp. 2236–2252, 2003.
- [4] M. L. Thesis, "Blind Single Channel Sound Source Separation Mark Leddy B. Sc, M. Sc Dublin Institute of Technology Supervisors: Dan Barry, David Dorran, Eugene Coyle," 2010.
- [5] N. Hassan and D. A. Ramli, "A Comparative study of Blind source separation for Bioacoustics sounds based on FastICA, PCA and NMF," *Procedia Comput. Sci.*, vol. 126, pp. 363–372, 2018.
- [6] N. Q. K. Duong, E. Vincent, and R. Gribonval, "Under-determined reverberant audio source separation using a full-rank spatial covariance model," *IEEE Trans. Audio, Speech Lang. Process.*, vol. 18, no. 7, pp. 1830–1840, 2010.
- [7] S. Rickard, "The DUET Blind Source Separation Algorithm," pp. 217–241, 2007.
- [8] V. S. Narayanaswamy, S. Katoch, J. J. Thiagarajan, H. Song, and A. Spanias, "Audio Source Separation via Multi-Scale Learning with Dilated Dense U-Nets," 2019.
- [9] X. F. Gong, Q. H. Lin, F. Y. Cong, and L. De Lathauwer, "Double Coupled Canonical Polyadic Decomposition for Joint Blind Source Separation," *IEEE Trans. Signal Process.*, vol. 66, no. 13, pp. 3475–3490, 2018.
- [10] S. Rickard and O. Yilmaz, "Blind Separation of Speech Mixtures via Time-Frequency Masking," *IEEE Trans. Signal Process.*, vol. 52, no. 7, pp. 1830–1847, 2004.
- [11] T. Gustafsson, B. D. Rao, and M. Trivedi, "Source Localization in Reverberant Environments: Part I - Modeling," vol. 11, no. 6, pp. 1–22, 2003.
- [12] T. Esch and P. Vary, "Efficient musical noise suppression for speech enhancement systems," *ICASSP, IEEE Int. Conf. Acoust. Speech Signal Process. - Proc.*, vol. 2, pp. 4409–4412, 2009.
- [13] Z. Rafii, A. Liutkus, F. R. Stoter, S. I. Mimilakis, D. Fitzgerald, and B. Pardo, "An Overview of Lead and Accompaniment Separation in Music," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 26, no. 8, pp. 1307–1335, 2018.
- [14] M. Jia, J. Sun, C. Bao, and C. Ritz, "Separation of multiple speech sources by recovering sparse and non-sparse components from B-format microphone recordings," *Speech Commun.*, vol. 96, no. May 2017, pp. 184–196, 2018.
- [15] S. Smita, S. Biswas, and S. S. Solanki, "Audio Signal Separation and Classification: A Review Paper," *Int. J. Innov. Res. Comput. Commun. Eng.*, vol. 2, no. 11, pp. 6960–6966, 2014.
- [16] N. Q. K. Duong, E. Vincent, and R. Gribonval, "Under-determined convolutive blind source separation using spatial covariance models," *ICASSP, IEEE Int. Conf. Acoust. Speech Signal Process. - Proc.*, pp. 9–12, 2010.

- [17] M. S. Khan, S. M. Naqvi, and J. Chambers, "Two-stage audio-visual speech dereverberation and separation based on models of the interaural spatial cues and spatial covariance," *2013 18th Int. Conf. Digit. Signal Process. DSP 2013*, 2013.
- [18] B. Schuller, "【Metrics】 Performance Measurement in Blind Audio Source Separation," vol. 14, no. 4, pp. 139–147, 2013.
- [19] J. J. Thiagarajan, K. Natesan Ramamurthy, and A. Spanias, "Mixing matrix estimation using discriminative clustering for blind source separation," *Digit. Signal Process. A Rev. J.*, vol. 23, no. 1, pp. 9–18, 2013.
- [20] K. Yatabe and D. Kitamura, "Determined Blind Source Separation via Proximal Splitting Algorithm," *ICASSP, IEEE Int. Conf. Acoust. Speech Signal Process. - Proc.*, vol. 2018-April, no. 4, pp. 776–780, 2018.
- [21] A. Tsilfidis, E. Georganti, and J. Mourjopoulos, "Binaural extension and performance of single-channel spectral subtraction dereverberation algorithms," *ICASSP, IEEE Int. Conf. Acoust. Speech Signal Process. - Proc.*, no. 5, pp. 1737–1740, 2011.