

# Direction of Arrival Estimation for Multiple Sound Sources Using Convolutional Recurrent Neural Network

Sharath Adavanne<sup>\*1</sup>, Archontis Politis<sup>\*2</sup>, Tuomas Virtanen<sup>1</sup>

<sup>1</sup>Laboratory of Signal Processing, Tampere University of Technology, Finland

<sup>2</sup>Department of Signal Processing and Acoustics, Aalto University, Finland

**Abstract**—This paper proposes a deep neural network for estimating the directions of arrival (DOA) of multiple sound sources. The proposed stacked convolutional and recurrent neural network (DOAnet) generates a spatial pseudo-spectrum (SPS) along with the DOA estimates in both azimuth and elevation. We avoid any explicit feature extraction step by using the magnitudes and phases of the spectrograms of all the channels as input to the network. The proposed DOAnet is evaluated by estimating the DOAs of multiple concurrently present sources in anechoic, matched and unmatched reverberant conditions. The results show that the proposed DOAnet is capable of estimating the number of sources and their respective DOAs with good precision and generate SPS with high signal-to-noise ratio.

## I. INTRODUCTION

Direction of arrival (DOA) estimation is the task of identifying the relative position of the sound sources with respect to the microphone. DOA estimation is a fundamental operation in microphone array processing and forms an integral part of speech enhancement [1], multichannel sound source separation [2] and spatial audio coding [3]. Popular approaches to DOA estimation are based on time-delay-of-arrival (TDOA) [4], the steered-response-power (SRP) [5], or on subspace methods such as multiple signal classification (MUSIC) [6] and the estimation of signal parameters via rotational invariance technique (ESPRIT) [7].

The aforementioned methods differ from each other in terms of algorithmic complexity, and their suitability to various arrays and sound scenarios. MUSIC specifically is very generic with regards to array geometry, directional properties and can handle multiple simultaneously active narrowband sources. On the other hand, MUSIC and subspace methods in general, require a good estimate of the number of active sources, which are often unavailable or difficult to obtain. Furthermore, MUSIC can suffer at low signal to noise ratio (SNR) and in reverberant scenarios [8]. In this paper, we propose to overcome the above shortcomings with a deep neural network (DNN) method, referred to as DOAnet, that learns the number of sources from the input data, generates high precision DOA estimates and is robust to reverberation. The proposed DOAnet

also generates a spatial acoustic activity map similar to the MUSIC pseudo-spectrum (SPS) as an intermediate output. The SPS has numerous applications that rely on a directional map of acoustic activity such as soundfield visualizations [9], and room acoustics analysis [10]. In comparison, the proposed DOAnet outputs the SPS and DOA's of multiple overlapping sources similar to any popular DOA estimators like MUSIC, ESPRIT or SRP without requiring the critical information of the number of active sound sources. A successful implementation of this will enable the integration of such DNN methods to higher-level learning based end-to-end sound analysis and detection systems.

Recently, several DNN-based approaches have been proposed for DOA estimation [11], [12], [13], [14], [15], [16]. There are six significant differences between them and the proposed method: a) All the aforementioned works focused on azimuth estimation, with the exception of [15] where the 2-D Cartesian coordinates of sound sources in a room were predicted, and [11] trained separate networks for azimuth and elevation estimation. In contrast, we demonstrate the estimation of both azimuth and elevation for the DOA by sampling the unit sphere uniformly and predicting the probability of sound source at each direction. b) The past works focused on the estimation of a single DOA at every time frame, with the exception of [13] where localization of azimuth for up to two sources simultaneously was proposed. On the other hand, the proposed DOAnet does not algorithmically limit the number of directions to be estimated, i.e., with a higher number of audio channels input, the DOAnet can potentially estimate a larger number of sound events.

c) Past works were evaluated with different array geometries making comparison difficult. Although the DOAnet can be applied to any array geometry, we evaluate the method using real spherical harmonic input signals, which is an emerging popular spatial audio format under the name Ambisonics. Microphone signals from various arrays, such as spherical, circular, planar or volumetric, can be transformed to Ambisonic signals by an appropriate transform [17], resulting in a common representation of the 3-D sound recording. Although the DOAnet is scalable to higher-order Ambisonics, in this paper we evaluate it using the compact four-channel first-order Ambisonics (FOA).

d) Regarding classifiers, earlier methods have used fully

<sup>\*</sup>Equally contributing authors in this paper. The research leading to these results has received funding from the European Research Council under the European Unions H2020 Framework Programme through ERC Grant Agreement 637422 EVERYSOUND. The authors also wish to acknowledge CSC-IT Center for Science, Finland, for computational resources

connected (FC) neural networks [11], [12], [13], [14], [15] and convolutional neural networks (CNN) [16]. In this work, along with the CNNs we use recurrent neural network (RNN) layers. The usage of RNN allows the network to learn long-term temporal information. Such an architecture is referred to as a convolutional recurrent neural network (CRNN) in literature and is the state-of-the-art method in many single- [18], [19] and multichannel [20], [21] audio tasks. e) Previous methods used inter-channel features such as generalized cross-correlation with phase transform (GCC-PHAT) [15], [12], eigen-decomposition of the spatial covariance matrix [13], inter-channel time delay (ITD) and inter-channel level differences (ILD) [11], [14]. More recently, Chakrabarty et al. [16] proposed to use only the phase component of the spectrogram, avoiding explicit feature extraction. In the proposed method, we use both the magnitude and the phase component. Contrary to [16], which employed omnidirectional sensors only, general arrays with directional microphones additionally encode the DOA information in magnitude differences, while Ambisonics format especially encode directional information mainly in the magnitude component. f) All previous methods were evaluated on speech recordings that were synthetically spatialized and spatially static. We continue to use the static sound sources in the present work and extend them to a larger variety of sound events, such as impulsive and transient sounds.

## II. METHOD

The block diagram of the proposed DOAnet is presented in Figure 1. The DOAnet takes multichannel audio as the input and first extracts the spectrograms of all the channels. The phases and the magnitudes of the spectrograms are mapped using a CRNN to two outputs sequentially. The first output, spatial pseudo-spectrum (SPS) is generated as a regression task, followed by the DOA estimates as a classification task. The DOA is defined by the azimuth  $\phi$  and elevation  $\lambda$  with respect to the microphone and the SPS is the intensity of sound along the DOA given by  $S(\phi, \lambda)$ .

In this paper, we use discrete  $\phi$  and  $\lambda$  by uniformly sampling the 2-D polar coordinate space, with a resolution of 10 degrees in both azimuth and elevation, resulting in 614 sampled directions. The SPS is computed at each sampled direction, whereas, a subset of 432 directions is used for DOA, where the elevations are limited between -60 and 60 degrees.

### A. Feature extraction

The spectrogram is calculated for each of the audio channels whose sampling frequencies are 44100 Hz. A 2048-point discrete Fourier transform (DFT) is calculated on Hamming windows of 40 ms with 50 % overlap. We keep 1024 values of the DFT corresponding to the positive frequencies, without the zeroth bin.  $L$  frames of features, each containing 1024 magnitude and phase values of the DFT extracted in all the  $C$  channels, are stacked in a  $L \times 1024 \times 2C$  3-D tensor and used as the input to the proposed neural network. The  $2C$  dimension results from ordering the magnitude component of

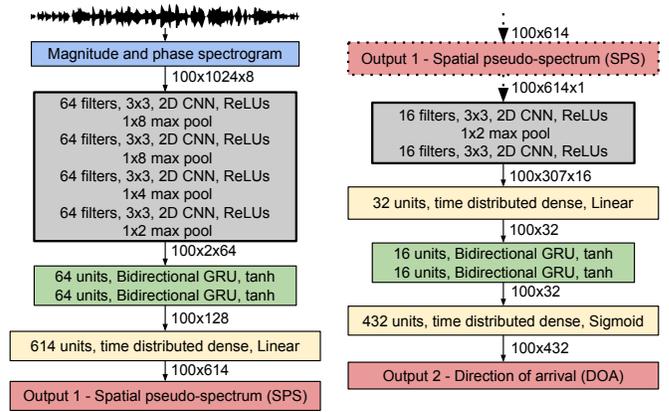


Fig. 1. DOAnet - neural network architecture for direction of arrival estimation of multiple sound sources.

all channels first, followed by the phase. We use a sequence length  $L$  of 100 ( $= 2$  s) in this work.

### B. Direction of arrival estimation network (DOAnet)

Local shift-invariant features are extracted from the input spectrogram tensor ( $L \times 1024 \times 2C$  dimension) using CNN layers. In every CNN layer, the intra-channel time-frequency features are processed using a receptive field of  $3 \times 3$ , rectified linear unit (ReLU) activation and pad zeros to the resulting activation map to keep the output dimension equal to input. Batch normalization and max-pooling operation along frequency axis are performed after every CNN layer to reduce the final dimension to  $L \times 2 \times N_C$ , where  $N_C$  is the number of CNN filters in the last CNN layer. The CNN activations are reshaped to  $L \times 2N_C$  keeping the time axis length unchanged and fed to RNN layers in order to learn temporal structure. Specifically, the bi-directional gated recurrent units (GRU) with tanh activation are used. Further, the RNN output is mapped to the first output, the SPS, in regression manner using FC layers with linear activation.

The SPS is further mapped to DOA estimates—the final output of the proposed method—using a similar CRNN network as above with two minor architectural changes. An FC layer is introduced between the CNN and RNN layers to reduce the dimension of the RNN output. Additionally, the output layer which predicts the DOA uses sigmoid activation in order to estimate more than one DOA for a given time frame. Each node in this output layer represents a direction in 2-D polar space. During testing, the probabilities at these nodes are thresholded with a value of 0.5, so that anything greater suggests the presence of a source in the direction or otherwise absence of source.

We refer to the combined architecture of SPS and DOA estimation in this work as DOAnet. The DOAnet is trained using the target SPS computed at each sampled direction, and for every time frame applying MUSIC (see Section III-B), and is represented using nonnegative real numbers. For the DOA output, the DOAnet aims to make a discrete decision about the presence of a source in a certain direction; and during training,

the DOAnet uses the ground truth DOAs utilized to synthesize the audio (see Section III-A).

The DOAnet was trained for 1000 epochs using Adam optimizer, mean squared error loss for SPS output and binary cross entropy loss for DOA output. The sum of the two losses was used for back propagation. Dropout was used after every layer and early stopping was used if the DOA metric (Section III-C) did not improve for 100 epochs. The DOAnet was implemented using Keras framework with Theano backend.

### III. EVALUATION

#### A. Dataset

In order to evaluate the proposed DOAnet, there are no publicly available real or synthetic datasets which consist of general sound events each associated with a 2D spatial coordinate. Since DNN-based methods need sufficiently large datasets to train on, most DNN-based methods proposed [11], [12], [14], [15], [16] have studied the performance on synthetic datasets. In similar fashion, we evaluate the proposed DOAnet on synthetic datasets about the same size as in the previous works.

We synthesize datasets consisting of static point sources associated with a spatial coordinate in the space in two contexts - anechoic and reverberant. For each context, three datasets are generated with no temporally overlapping sources (*O1*), maximum two overlapping sources (*O2*), and maximum three overlapping sound sources (*O3*). We refer to the anechoic context dataset as *OxA* and reverberant as *OxR*, where  $x$  denotes the number of overlapping sources. Each of these datasets has three cross-validation (CV) splits with 240 recordings for training and 60 for testing. Recordings are sampled at 44.1 kHz and 30 s long.

In order to generate these datasets, we use the isolated real-life sound event recordings from the DCASE 2016 task 2 [22]. This dataset consists of 11 sound event classes, each with 20 examples. The classes in this dataset included speech, coughing, door slam, page-turning, phone ringing and keyboard sounds. During CV, for each of the splits, we randomly chose disjoint sets of 16 and 4 examples for training and testing, amounting to 176 examples for training and 44 for testing. In order to synthesize a recording, a random subset of the 176 or 44 sound examples was chosen from the respective split. The subset size varied for each recording based on the chosen sound examples. We start synthesizing a recording by randomly choosing the beginning time of the first randomly chosen sound example within the first second of the recording. The next randomly chosen sound example is placed 250-500 ms after the end of the first sound example. On reaching the maximum recording length of 30 s, the process is repeated as many times as the number of required overlapping sound events.

Each of the sound examples were assigned a DOA randomly using the following conditions. All sound events were placed in a spatial grid of ten degrees resolution along both azimuth and elevation. Two temporally overlapping sound events have at least ten degrees of spatial separation to avoid spatial

overlapping. The elevation was constrained within the range of [-60, 60] degrees, as most natural sound events occur in this range. Finally, for the anechoic dataset, the sound sources were randomly placed at a distance  $d$  in the range 1-10 m. For the reverberant dataset, the sound events were randomly placed inside a room of dimensions  $10 \times 8 \times 4$  m with the microphone in the center of the room.

Spatialization for the anechoic case was done as following. Each point source signal  $s_i$  with DOA  $(\phi_i, \lambda_i)$ , was converted to Ambisonics format by multiplying the signal with the vector  $\mathbf{y}(\phi_i, \lambda_i) = [Y_{00}(\phi_i, \lambda_i), Y_{1(-1)}(\phi_i, \lambda_i), Y_{10}(\phi_i, \lambda_i), Y_{11}(\phi_i, \lambda_i)]^T$  of real orthonormalized spherical harmonics  $Y_{nm}(\phi, \lambda)$ . The complete anechoic sound scene multichannel recording  $\mathbf{x}_A$  was generated as  $\mathbf{x}_A = \sum_i g_i s_i \mathbf{y}(\phi_i, \lambda_i)$ , with the gains  $g_i < 1$  modeling the distance attenuation. Each entry of  $\mathbf{x}_A$  corresponds to one channel and  $g_i = \sqrt{1/10^{d/d_{max}}}$ , where  $d_{max} = 10$  m is the maximum distance.

In the reverberant case, a fast geometrical acoustics simulator was used to model natural reverberation based on the rectangular room image-source model [23]. For each point source  $s_i$  with DOA in the dataset,  $K$  image sources were generated modeling reflections up to a predefined time-limit. Based on the room and its propagation properties, each image source was associated with a propagation filter  $h_{ik}$  and DOA  $(\phi_k, \lambda_k)$  resulting in the spatial impulse response  $\mathbf{h}_i = \sum_{k=1}^K h_{ik} \mathbf{y}(\phi_k, \lambda_k)$ . The reverberant scene signal was finally generated by  $\mathbf{x}_R = \sum_i s_i * \mathbf{h}_i$ , where  $(*)$  denotes convolution of the source signal with the spatial impulse responses. The room absorption properties were adjusted to match reverberation times of typical office spaces. Three sets of testing data were generated with similar room size as training data (Room 1), 80% of room size ( $8 \times 8 \times 4$  m) and reverberation time (Room 2), and 60% of room size ( $8 \times 6 \times 4$  m) and reverberation time (Room 3).

#### B. Baseline

The proposed method to our knowledge is the first DNN-based implementation for 2D DOA estimation of multiple overlapping sound events. Thus in order to evaluate the complete features of the proposed DOAnet, we compare the performance with the conventional, high-resolution DOA estimator based on MUSIC. Similar to the SPS and DOA outputs estimated by the DOAnet, the MUSIC method also estimates SPS and DOA, thus allowing a direct one-to-one comparison.

The MUSIC SPS is based on a measure of orthogonality between the signal subspace (dominated by the source signals) of the spatial covariance matrix  $\mathbf{C}_s$  and the noise subspace (dominated by diffuse and ambient sounds, late reverberation, and microphone noise). The spatial covariance matrix is calculated as  $\mathbf{C}_s = \mathbb{E}_{f,t} [\mathbf{X}(f,t)\mathbf{X}(f,t)^H]$ , where spectrogram  $\mathbf{X}(f,t)$  is a frequency  $f$  and time  $t$  dependent  $C$ -dimensional vector, where  $C$  is the number of channels,  $^H$  is the conjugate transpose and  $\mathbb{E}_{f,t}$  denotes the expectation over  $f$  and  $t$ . For a sound scene with  $O$  number of sources, the MUSIC SPS

$S_{GT}$  is obtained from  $\mathbf{C}_s$  by first performing an eigenvalue decomposition on  $\mathbf{C}_s = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^H$ . The sorted eigenvectors  $\mathbf{E}$  (according to eigenvalues with decreasing magnitude) are further partitioned into the two aforementioned subspaces  $\mathbf{E} = [\mathbf{U}_s \ \mathbf{U}_n]$ , where  $\mathbf{U}_s$  denotes the signal subspace and will be composed of  $O$  eigenvectors corresponding to the higher eigenvalues and the rest will form the noise subspace  $\mathbf{U}_n$ . The  $S_{GT}$  along the direction  $(\phi_i, \lambda_i)$  is now given by  $S_{GT}(\phi_i, \lambda_i) = 1/(\mathbf{y}^T(\phi_i, \lambda_i)\mathbf{U}_n\mathbf{U}_n^H\mathbf{y}(\phi_i, \lambda_i))$ . Finally, the source DOAs are found by selecting the directions  $(\phi_i, \lambda_i)$  corresponding to the  $O$  largest peaks from  $S_{GT}$ .

### C. Metric

The DOAnet estimated SPS ( $S_E(\phi, \lambda)$ ) is evaluated with respect to the baseline MUSIC estimated ground truth ( $S_{GT}(\phi, \lambda)$ ) using the SNR metric calculated as  $SNR = 10 \log_{10}(\sum_{\phi} \sum_{\lambda} S_{GT}(\phi, \lambda)^2 / \sum_{\phi} \sum_{\lambda} (S_E(\phi, \lambda) - S_{GT}(\phi, \lambda))^2)$ .

As the DOA metric we use the angle between the estimate DOA (defined by azimuth  $\phi_E$  and elevation  $\lambda_E$ ) and the ground truth DOA ( $\phi_{GT}, \lambda_{GT}$ ) used to synthesize the dataset in degrees. This is calculated as  $\sigma = \arccos(\sin \phi_E \sin \phi_{GT} + \cos \phi_E \cos \phi_{GT} \cos(\lambda_{GT} - \lambda_E)) \cdot 180.0/\pi$ . Further, to accommodate the scenario of unequal number of estimated and ground truth DOAs we calculate and report the minimum distance between them using the Hungarian algorithm [24] along with the percentage of frames in which the number of DOAs estimated were correct. The final metric for the entire dataset, referred as DOA error, is calculated by normalizing the minimum distance with the total number of estimated DOA's.

### D. Evaluation procedure

The parameter tuning for DOAnet was performed on the *O1A* test data, and the best configuration is as shown in Figure 1. This configuration has 677 K weights, and the same configuration is used in all of the following studies.

At test time, the SNR metric for SPS output of the DOAnet ( $S_E$ ) is calculated with respect to SPS of baseline MUSIC ( $S_{GT}$ ). The DOA metric for the DOAs predicted by DOAnet and baseline MUSIC are calculated with respect to the ground truth DOA used to synthesize the dataset.

In the above experiment, the baseline MUSIC algorithm uses the knowledge of the number of active sources. In order to have a fair evaluation, we test the DOAnet in a similar scenario where the number of sources is known. We use this knowledge to choose the top probabilities in prediction layer of the DOAnet instead of thresholding it with a value of 0.5.

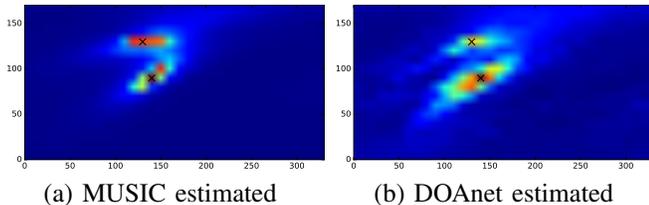


Fig. 2. SPS for two closely located sound sources. The black-cross markers represent the ground truth DOA. The horizontal axis is azimuth and vertical axis is elevation angle (in degrees)

TABLE I  
EVALUATION METRIC SCORES FOR THE SPATIAL POWER MAP AND DOAS ESTIMATED BY THE DOANET FOR DIFFERENT DATASETS.

	Anechoic			Reverberant (Room 1)		
	1	2	3	1	2	3
Max. no. of overlapping sources						
SPS SNR (in dB)	9.90	3.35	-0.26	3.11	1.24	0.13
DOA error with unknown number of active sources (threshold of 0.5)						
DOAnet	0.57	8.03	18.34	6.31	11.46	38.41
Correctly predicted frames (in %)	95.4	42.7	1.8	59.3	15.8	1.2
DOA error with known number of active sources						
DOAnet	1.14	27.52	49.30	12.61	38.98	67.07
MUSIC	2.29	8.60	28.66	25.80	57.33	91.72

## IV. RESULTS AND DISCUSSION

The results of the evaluations are presented in Table I. The high SNRs for SPS in both the contexts, with up to one and two overlapping sound events show that the SPS generated by DOAnet ( $S_E$ ) is comparable with the baseline MUSIC SPS ( $S_{GT}$ ). Figure 2 shows the  $S_E$  and the respective  $S_{GT}$  when two active sources are closely located. In the case of up to three overlapping sound events, the baseline MUSIC is already at its theoretical limit of estimating  $N - 1$  sources from  $N$ -dimensional signal space [25]. In practice, for  $N - 1$  sources only one noise subspace vector  $\mathbf{U}_n$  is used to generate SPS, which for real signals is too weak for stable estimation. In the present evaluation of DOAnet which is trained with four-channel audio features and MUSIC SPS, for the case of three overlapping sound sources the SPS used is an unstable estimate resulting in poor training and consequently the results. With more than four-channels input, which the proposed DOAnet can easily extend to, it can potentially localize more than two sound sources simultaneously.

The DOA error for the proposed DOAnet when the number of active sources are unknown is presented in Table I. The DOAnet error is considerably better in comparison to the baseline MUSIC that uses the active sources knowledge for all datasets. However, the number of frames in which DOAnet produced the correct number of active sources were few. For example, in the case of anechoic recordings with up to two overlapping sound events, only 42.7% of the estimated frames had the correct number of DOA predictions. This prediction drops even drastically when the number of sources is three, due to the theoretical limit of MUSIC as explained previously, and consequently for the DOAnet as MUSIC SPS is used for training. Finally, the confusion matrix for the number of DOA estimates per frame for *O1* and *O2* datasets are visualized

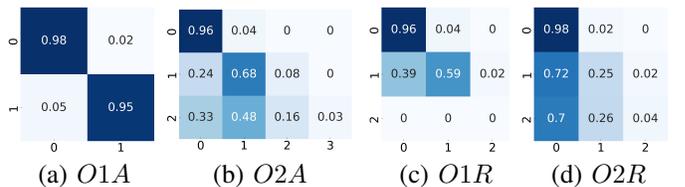


Fig. 3. Confusion matrix for the number of DOA estimated per frame by the DOAnet. The horizontal axis is the DOAnet estimate, and the vertical axis is the ground truth.

TABLE II  
EVALUATION SCORES FOR UNMATCHED REVERBERANT ROOM.

	Room 2		Room 3	
	1	2	1	2
Max. no. of overlapping sources	1	2	1	2
SPS SNR (in dB)	3.53	1.49	3.49	1.46
DOAnet error (Unknown number of sources)				
DOAnet	3.44	6.88	4.59	10.89
Correctly predicted frames (in %)	46.2	14.3	49.7	14.1
DOA error (Known number of sources)				
DOAnet	8.60	32.10	9.17	33.82
MUSIC	31.52	58.47	33.25	60.76

in Figure 3. We skipped the confusion matrices for the  $O3$  datasets as they were not meaningful for similar reasons as explained above.

With the knowledge of the number of active sources (Table I), the DOAnet performs considerably better than baseline MUSIC for all datasets other than the  $O2A$  and  $O3A$ . The MUSIC DOA's were chosen using a 2D peak finder on the MUSIC SPS, whereas the DOA's in DOAnet were chosen by simply picking the top probabilities in the final DOA prediction layer. A smarter peak picking method from the DOAnet, or using the number of sources as an additional input can potentially result in better scores across all datasets. Further, the DOAnet error on unmatched reverberant data is presented in Table II. The performance of DOAnet is seen to be consistent in comparison to the matched reverberant data in Table I, and significantly better than the performance of MUSIC.

In this paper, since the baseline was chosen to be MUSIC, for a fair comparison the DOAnet was also trained using MUSIC SPS. In an ideal scenario, considering the DOAnet is trained using datasets for which the ground truth DOAs are known, we can generate accurate high-resolution SPS from the ground truth DOA's as per the required application and use them for training. Alternatively, the DOAnet can be trained without the SPS to directly generate the DOAs, it was only used in this paper to present the complete potential of the method in the limited paper space. In general, the above results show that the proposed DOAnet has the potential to learn the 2D direction information of multiple overlapping sound sources directly from the spectrogram of the input audio without the knowledge of the number of active sound sources. An exhaustive study with more detailed experiments including both synthetic and real datasets are planned for future work.

## V. CONCLUSION

A convolutional recurrent neural network (DOAnet) was proposed for multiple source localization. The DOAnet was shown to learn the number of active sources directly from the input spectrogram, and estimate precise DOA in 2-D polar space. The method was evaluated on anechoic, matched and unmatched reverberant dataset. The proposed DOAnet performed considerably better than baseline MUSIC in most scenarios. Thereby showing the potential of DOAnet in learning highly computational algorithm without prior knowledge of the number of sources.

## REFERENCES

- [1] M. Woelfel and J. McDonough, "Distant speech recognition," in Wiley, 2009.
- [2] J. Nikunen and T. Virtanen, "Direction of arrival based spatial covariance model for blind sound source separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 3, 2014.
- [3] A. Politis *et al.*, "Sector-based parametric sound field reproduction in the spherical harmonic domain," *IEEE Journal of Selected Topics in Signal Processing*, vol. 9, no. 5, pp. 852–866, 2015.
- [4] Y. Huang *et al.*, "Real-time passive source localization: a practical linear-correction least-squares approach," in *IEEE Transactions on Speech and Audio Processing*, vol. 9, no. 8, 2001.
- [5] M. S. Brandstein and H. F. Silverman, "A high-accuracy, low-latency technique for talker localization in reverberant environments using microphone arrays," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1997.
- [6] R. O. Schmidt, "Multiple emitter location and signal parameter estimation," in *IEEE Transactions on Antennas and Propagation*, vol. 34, no. 3, 1986.
- [7] R. Roy and T. Kailath, "ESPRIT-estimation of signal parameters via rotational invariance techniques," in *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 37, no. 7, 1989.
- [8] J. H. DiBiase *et al.*, "Robust localization in reverberant rooms," in *Microphone Arrays*, 2001, pp. 157–180.
- [9] A. O'Donovan *et al.*, "Imaging concert hall acoustics using visual and audio cameras," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2008.
- [10] D. Khaykin and B. Rafaely, "Acoustic analysis by spherical microphone array processing of room impulse responses," *The Journal of the Acoustical Society of America*, vol. 132, no. 1, 2012.
- [11] R. Roden *et al.*, "On sound source localization of speech signals using deep neural networks," in *Deutsche Jahrestagung für Akustik (DAGA)*, 2015.
- [12] X. Xiao *et al.*, "A learning-based approach to direction of arrival estimation in noisy and reverberant environments," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015.
- [13] R. Takeda and K. Komatani, "Discriminative multiple sound source localization based on deep neural networks using independent location model," in *IEEE Spoken Language Technology Workshop (SLT)*, 2016.
- [14] A. Zermini *et al.*, "Deep neural network based audio source separation," in *International Conference on Mathematics in Signal Processing*, 2016.
- [15] F. Vesperini *et al.*, "A neural network based algorithm for speaker localization in a multi-room environment," in *IEEE International Workshop on Machine Learning for Signal Processing (MLSP)*, 2016.
- [16] S. Chakrabarty and E. A. P. Habets, "Broadband DOA estimation using convolutional neural networks trained with noise signals," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2017.
- [17] H. Teutsch, *Modal array signal processing: principles and applications of acoustic wavefield decomposition*. Springer, 2007, vol. 348.
- [18] T. N. Sainath *et al.*, "Convolutional, long short-term memory, fully connected deep neural networks," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015.
- [19] M. Malik *et al.*, "Stacked convolutional and recurrent neural networks for music emotion recognition," in *Sound and Music Computing Conference (SMC)*, 2017.
- [20] T. Sainath *et al.*, "Multichannel signal processing with deep neural networks for automatic speech recognition," in *IEEE Transactions on Audio, Speech, and Language Processing*, 2017.
- [21] S. Adavanne *et al.*, "Sound event detection using spatial features and convolutional recurrent neural network," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017.
- [22] E. Benetos *et al.*, "Sound event detection in synthetic audio," <http://www.cs.tut.fi/sgn/arg/dc2016/>, 2016.
- [23] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," in *The Journal of the Acoustical Society of America*, vol. 65, no. 4, 1979.
- [24] H. W. Kuhn, "The hungarian method for the assignment problem," in *Naval Research Logistics Quarterly*, no. 2, 1955, p. 8397.
- [25] B. Ottersten *et al.*, "Exact and large sample maximum likelihood techniques for parameter estimation and detection in array processing," in *Radar Array Processing. Springer Series in Information Sciences*, 1993.