

Multi-Source DOA Estimation through Pattern Recognition of the Modal Coherence of a Reverberant Soundfield

Abdullah Fahim, Prasanga N. Samarasinghe, Thushara D. Abhayapala
Audio & Acoustic Signal Processing Group, The Australian National University, Canberra, Australia

Abstract—We propose a novel multi-source direction of arrival (DOA) estimation technique using a convolutional neural network algorithm which learns the modal coherence patterns of an incident soundfield through measured spherical harmonic coefficients. We train our model for individual time-frequency bins in the short-time Fourier transform spectrum by analyzing the unique snapshot of modal coherence for each desired direction. The proposed method is capable of estimating simultaneously active multiple sound sources on a 3D space using a single-source training scheme. This single-source training scheme reduces the training time and resource requirements as well as allows the reuse of the same trained model for different multi-source combinations. The method is evaluated against various simulated and practical noisy and reverberant environments with varying acoustic criteria and found to outperform the baseline methods in terms of DOA estimation accuracy. Furthermore, the proposed algorithm allows independent training of azimuth and elevation during a full DOA estimation over 3D space which significantly improves its training efficiency without affecting the overall estimation accuracy.

Index Terms—Convolutional neural network, DOA estimation, spatial audio processing, spherical harmonics

I. INTRODUCTION

THE task of a direction of arrival (DOA) estimator is to identify sound source directions from the output of a microphone array. DOA estimation is typically a prerequisite for several signal processing algorithms such as beamforming, power spectral density (PSD) estimation, and spatial audio coding, which in turn are integral parts of many practical applications such as multi-source separation [1]–[3], speech recognition [4], robot audition [5], [6], bio-diversity monitoring [7], [8], audio surveillance and smart home applications [9], [10]. In many practical environments, we face a scenario where multiple sound sources from different directions contribute to the acoustic scene. In such mixed recording scenarios, different sound sources can overlap on each other partially (e.g., teleconference) or fully (e.g., cocktail party problem, bioacoustic) over time. In this paper, we utilize the modal coherence of the spherical harmonic coefficients of a reverberant soundfield to train a convolutional neural network (CNN) for DOA estimation of multiple sound sources irrespective of their overlapping nature.

This work is supported by the Australian National University Strategic Research Funds.

A. Literature review

DOA estimation is a decade-old problem with a number of algorithms developed over the years to accurately estimate sound source locations. However, while different algorithms have shown their usefulness under certain environments, they all have their own constraints and limitations and hence, DOA estimation remains an active problem in acoustic signal processing. A large number of DOA estimation techniques have been developed in the parametric domain [11]. There are subspace-based methods like multiple signal classification (MUSIC) [12] or the estimation of signal parameters via rotational invariance technique (ESPRIT) [13] which utilizes the orthogonality between the signal and noise subspaces to estimate the source DOAs. MUSIC algorithm was originally developed for narrowband signals, however, it has been extensively used with wideband processing using a frequency smoothing technique [14] or by decomposing the signal into multiple narrowband subspaces [15]. It is common knowledge that the performance of the subspace-based methods are susceptible to strong reverberation and background noise [16]. Recently a variation of MUSIC was proposed in [17] to improve its robustness in a reverberant room assuming the prior knowledge of room coupling coefficients.

There also exist beamforming-based methods for DOA estimation where the output power of a beamformer is scanned in all possible directions to find out when it reaches the maximum. A popular formulation of the beamformer-based technique is the steered response power (SRP) method which formulates the output power as a sum of cross-correlations between the received signals. Dibiase proposed an improvement to SRP in [18] using the phase transform (PHAT) variant of the generalized cross-correlation (GCC) model [19]. The beamforming-based methods experience degradation in their performance for closely-spaced sources due to the limitation of the spatial resolution. Furthermore, both subspace and beamforming based techniques require to scan for all possible DOA angles during the run time which can be both time and resource intensive. Several modifications have been proposed to reduce the computational cost of SRP-PHAT by replacing the traditional grid search with region-based search [20]–[23], however, this increases the probability of missing a desired source in reverberant conditions.

Another group of parametric approaches to DOA estimation uses the maximum likelihood (ML) optimization with the

statistics of the observed data which usually requires accurate statistical modeling of the noise field [24]–[26]. In more recent works, DOA estimation, posed as a ML problem, was separately solved for reverberant environments [27], [28] and with unknown noise power [29] using expectation-maximization technique. A large number of localization techniques are based on the assumption of non-overlapping source mixture in the short-time Fourier transform (STFT) domain, known as W-disjoint orthogonality (WDO) [30]. Li *et al.* adopted a Gaussian mixture model to estimate source locations using ML method on the basis of WDO [31]. The sparsity of speech signals in the STFT domain was exploited in [32], [33] to localize broadside sources by mapping phase difference histogram of the STFT coefficients. The works in [34], [35] imply sparsity on both signals and reflections to isolate time-frequency (TF) bins that contain only direct path contributions from a single source and subsequently estimate source DOAs based on the selected TF bins. Recently, there has been an increase in efforts for intensity-based approaches where both sound pressure and particle velocity are measured and used together for DOA estimation [36]–[39].

Lately, the application of spatial basis functions, especially the spherical harmonics, is gaining researchers’ attention in solving a wide variety of acoustic problems including DOA estimation. Among the works we have referred so far in this paper, [14], [17], [38], [39] were implemented in the spherical harmonic domain. Tervo *et al.* proposed a technique for estimating DOA of the room reflections based on a ML optimization using a spherical microphone array [40]. Kumar *et al.* implemented MUSIC and beamforming-based techniques with the spherical harmonic decomposition of a soundfield for nearfield DOA estimation in a non-reverberant environment [41]. A free-field model of spherical harmonic decomposition was used in [42] to perform an optimized grid search for acoustic source localization. The spherical harmonics are the natural basis functions for spatial signal processing and consequently offers convenient ways for recognizing the spatial pattern of a soundfield. Furthermore, the spherical harmonic coefficients are independent of the array structure, hence, the same DOA algorithm can be used with different shapes and designs of sensor arrays as long as they meet a few basic criteria of harmonic decomposition [43]–[48].

Over the past decade, the rapid technology advances in storage and processing capabilities led researchers to lean towards machine learning in solving many practical problems including DOA estimation. Being a data-driven approach, neural networks can be trained for different acoustic environments and source distributions. In the area of single source localization, significant progresses have been made in solving the limitations in the parametric approaches by incorporating machine learning-based algorithms. The authors of [49]–[51] derived features from different variations of the GCC model to train a neural network for single source localization. Ferguson *et al.* used both cepstrogram and GCC to propose a single source DOA estimation technique for under-water acoustics [52]. Inspired by the MUSIC algorithm, the authors of [53] utilized the eigenvectors of a spatial correlation matrix to train a deep neural network. Conversely, multi-source localization

poses a more challenging problem to solve, especially with overlapping sources. In the recent past, a few algorithms have been proposed for multi-source localization based on CNN. A CNN-based multi-source DOA estimation technique was proposed in [54] where the authors used the phase spectrum of the microphone array output as the learning feature. The method in [54] was implemented in the STFT domain and all the STFT bins for each time frame were stacked together to form the feature snapshot. On the contrary, Adavanne *et al.* considered both magnitude and phase information of the STFT coefficients and used consecutive time frames to form the feature snapshot to train a convolutional recurrent neural network (CRNN) and performed a joint sound event detection and localization [55]. Both [54] and [55] require the model to be trained for unique combinations of sound sources from different angles in order to accurately estimate the DOA of simultaneously active multiple sound sources.

B. Contribution of this paper

In this paper, we propose a CNN-based framework to estimate DOAs of simultaneously active multiple sound sources. We use a novel feature to train a CNN which utilizes the modal coherence of a reverberant soundfield in the spherical harmonic domain. We show that modal coherence represents a unique pattern for each direction which can be learned and used for estimating source DOAs in a composite acoustic scene. Note that, previously we developed a parametric multi-source separation algorithm in [56] using the modal coherence model, hence, the proposed method offers an efficient way of joint multi-source localization and separation in a reverberant environment.

We design the algorithm to perform multi-source DOA estimation while being trained for the single-source case only. To the best of our knowledge, this approach has not been taken in the learning-based DOA estimation so far as most of the existing techniques require the training data to contain different angular combinations from the desired DOA grid. The single-source training model of the proposed method offers multiple advantages during training and testing stages. As the number of unique combinations increases rapidly with the number of audio sources and size of the DOA grid¹, a requirement for multi-source training causes a significant increment of time and resource requirements compared to the single-source training model. Furthermore, the adoption of single-source training allows us to train the model for once and use the same model in the testing environment irrespective of the number of sound sources in the mixed acoustic scenario.

Unlike the existing methods in the CNN domain, we treat individual STFT bins separately for training. The proposed algorithm uses the W-disjoint orthogonality assumption [30] for multi-source separation and proposes a probability-based post-filtering of the CNN output to address any violation of W-disjoint orthogonality. Furthermore, as the training stage of the proposed algorithm involves single-source scenarios

¹For a perspective, 250 distinct DOA points contribute to approximately 2.5 million and 160 million unique combinations for 3 and 4 source-mixtures, respectively.

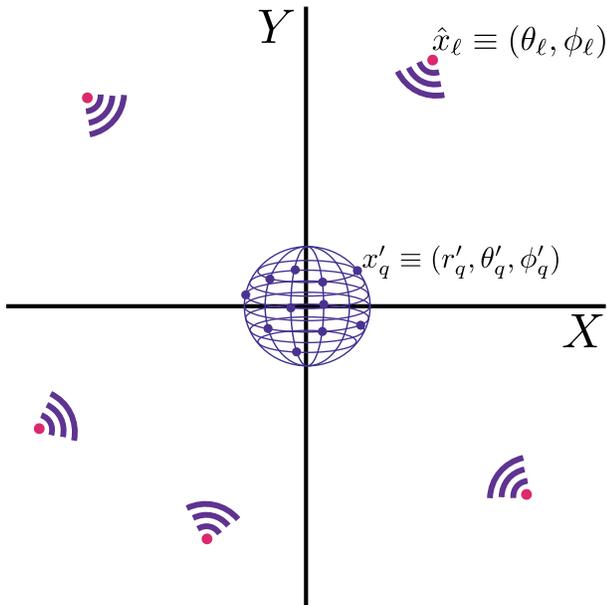


Fig. 1. A graphical impression of a spherical microphone array setup in the presence of multiple sound sources. Array shape may differ depending on the spherical harmonic decomposition technique.

only, we can independently train our model for azimuths and elevations based on the same input dataset provided that we measure the soundfield for various source positions in each intended azimuth and elevation planes. Hence, we can share the same convolutional layers and perform joint azimuth and elevation estimation with two separate fully connected heads. In the results section, we include a demonstration for such a joint estimation along with a complete performance evaluation of the proposed method with a wide range of criteria such as different practical and simulated room conditions with a variable number of sources. We also compare and analyze the results with another state of the art CNN-based algorithm for DOA estimation.

The remainder of the paper is structured as follows. Section II contains the problem statement and defines the objective of the work. In Section III, we briefly describe the existing frameworks for spherical harmonic decomposition and modal coherence of a reverberant soundfield. Section IV contains a detailed description of the proposed model including different aspects of feature selection. Finally, in Section V, we evaluate and analyze the performance of the proposed algorithm and compare it with a contemporary method based on objective metrics and graphical aid.

II. PROBLEM FORMULATION

Consider L sound sources concurrently emitting sound in a reverberant room. The sound pressure observed by an omnidirectional microphone placed at a coordinate $\mathbf{x}'_q \equiv (r'_q, \theta'_q, \phi'_q)$ inside the room, where r'_q , θ'_q , and ϕ'_q are the radius, elevation, and azimuth of point \mathbf{x}'_q in the spherical coordinate system, respectively, is expressed by

$$p(\mathbf{x}'_q, t) = \sum_{\ell=1}^L h_{\ell}(\mathbf{x}'_q, t) * s_{\ell}(t) \quad (1)$$

where t is the discrete time index, $h_{\ell}(\mathbf{x}'_q, t)$ is the room impulse response (RIR) between the ℓ^{th} source position and \mathbf{x}'_q , $s_{\ell}(t)$ is the ℓ^{th} source signal, and $*$ denotes the convolution operation. The corresponding frequency domain representation of (1) in STFT domain can be obtained using the multiplicative model of convolution and is formulated as

$$P(\mathbf{x}'_q, k, \tau) = \sum_{\ell=1}^L S_{\ell}(k, \tau) H_{\ell}(\mathbf{x}'_q, k) \quad (2)$$

where $\{P, S, H\}$ represent the corresponding signals of $\{p, s, h\}$ in the STFT domain, τ is the timeframe index, $k = 2\pi f/c$, f denotes the frequency, and c is the speed of sound propagation. Henceforth, τ is omitted for brevity as we shall treat each of the time frames independently.

In this work, we intend to estimate the individual DOAs in the presence of multiple concurrent sound sources, i.e., we want to estimate $\hat{\mathbf{x}}_{\ell} \equiv (\hat{\theta}_{\ell}, \hat{\phi}_{\ell}) \forall \ell \in [1, L]$, given a set of measured sound pressure $p(\mathbf{x}'_q, t) \forall q \in [1, Q]$ or the corresponding spherical harmonic coefficients² of a mixed soundfield. We pose the DOA estimation as a CNN classification problem where we sample the intended DOA range into discrete sets $\Theta = \{\theta_a\}_{a \in [1, I]}$ for elevations and $\Phi = \{\phi_b\}_{b \in [1, J]}$ for azimuths. Thereafter, we propose a feature unique to each angle and train a CNN framework individually for each of the members of Θ and Φ . Finally, during the evaluation, the CNN model finds the closest match of the true DOA $\hat{\mathbf{x}}_{\ell} \equiv (\theta_{\ell}, \phi_{\ell}) \forall \ell$ in the DOA sets Θ and Φ based on its learning and accurately combines the independent estimations $\hat{\theta}_{\ell}$ and $\hat{\phi}_{\ell}$ for each individual source to achieve full DOA estimation.

III. MODAL FRAMEWORK

In this section, we describe a few established concepts in the spherical harmonic domain that is used as the framework for the proposed DOA estimation technique.

A. Spherical harmonic decomposition of a soundfield

A continuous soundfield on a sphere can be decomposed using the spherical harmonic basis functions as [57]

$$P(\mathbf{x}'_q, k) = \sum_{nm}^{\infty} \alpha_{nm}(k) b_n(kr'_q) Y_{nm}(\hat{\mathbf{x}}'_q) \quad (3)$$

where $\sum_{nm}^{(\cdot)} \equiv \sum_{n=0}^{(\cdot)} \sum_{m=-n}^n$, $\alpha_{nm}(k)$ is the sensor-independent soundfield coefficient of order n and degree m , the position vector $\mathbf{x}'_q \equiv (r'_q, \hat{\mathbf{x}}'_q)$, and the unit vector $\hat{\mathbf{x}}'_q \equiv (\theta'_q, \phi'_q)$. The infinite summation of (3) is often truncated at the soundfield order $N = \lceil kr'_q \rceil$ [58], [59], where $\lceil \cdot \rceil$ denotes the ceiling operation, due to the high-pass nature of the higher-order Bessel functions. The complex spherical harmonic basis function $Y_{nm}(\cdot)$ is defined as

$$Y_{nm}(\hat{\mathbf{x}}'_q) = \sqrt{\frac{(2n+1)(n-|m|)!}{4\pi(n+|m|)!}} \mathcal{P}_{n|m|}(\cos \theta'_q) e^{im\phi'_q} \quad (4)$$

²The spherical harmonic decomposition technique is described in Section III-A.

where $|\cdot|$ denotes absolute value, $(\cdot)!$ represents factorial, $\mathcal{P}_{n|m}(\cdot)$ is an associated Legendre polynomial, and $i = \sqrt{-1}$. Furthermore, the dependency on array radius comes through the function $b_n(\cdot)$ which is defined as

$$b_n(\xi) = \begin{cases} j_n(\xi) & \text{for an open array} \\ j_n(\xi) - \frac{j'_n(\xi)}{h'_n(\xi)} h_n(\xi) & \text{for a rigid spherical array.} \end{cases} \quad (5)$$

The spherical harmonic coefficients α_{nm} can be estimated using different kinds of arrays where the process and the formulation depend on the array geometry. E.g., utilizing the orthonormal property of the spherical harmonics, a spherical microphone array allows us to calculate α_{nm} from (3) as [43], [46]

$$\alpha_{nm}(k) = \frac{1}{b_n(kr)} \int_{\mathbb{S}^2} P(\mathbf{x}, k) Y_{nm}^*(\hat{\mathbf{x}}) d\hat{\mathbf{x}} \quad (6)$$

$$\approx \frac{1}{b_n(kr)} \sum_{q=1}^Q w_q P(\mathbf{x}'_q, k) Y_{nm}^*(\hat{\mathbf{x}}'_q) \quad (7)$$

where r is the array radius, $\mathbf{x} = (r, \hat{\mathbf{x}})$, and $\hat{\mathbf{x}}$ denotes an arbitrary direction on the spherical shell \mathbb{S}^2 . Obviously, it is impractical to realize a spatially continuous microphone array as required by (6), hence, (7) is used as an approximation of (6) with Q microphones. $w_q \forall q$ are suitable microphone weights that ensure the validity of the orthonormal property of the spherical harmonics with a limited number of sampling points, i.e.,

$$\sum_{q=1}^Q w_q Y_{nm}(\hat{\mathbf{x}}'_q) Y_{n'm'}^*(\hat{\mathbf{x}}'_q) \approx \delta_{nn'} \delta_{mm'} \quad (8)$$

where $\delta_{nn'}$ and $\delta_{mm'}$ are the Kronecker delta functions.

The same spherical harmonic decomposition can be achieved using alternate array geometries and formulation, [43]–[48] are a few examples of such procedures.

B. RTF in the spatial domain

The room transfer function (RTF) can be decomposed into two parts

$$H_\ell(\mathbf{x}'_q, k) = H_\ell^{(d)}(\mathbf{x}'_q, k) + H_\ell^{(r)}(\mathbf{x}'_q, k) \quad (9)$$

where $H_\ell^{(d)}(\cdot)$ and $H_\ell^{(r)}(\cdot)$ are the corresponding direct and reverberant path components of the RTF. The RTF components are modeled in the spatial domain as

$$H_\ell^{(d)}(\mathbf{x}'_q, k) = G_\ell^{(d)}(k) e^{ik \hat{\mathbf{x}}_\ell \cdot \mathbf{x}'_q} \quad (10)$$

$$H_\ell^{(r)}(\mathbf{x}'_q, k) = \int_{\mathbb{S}^2} G_\ell^{(r)}(k, \hat{\mathbf{x}}) e^{ik \hat{\mathbf{x}} \cdot \mathbf{x}'_q} d\hat{\mathbf{x}} \quad (11)$$

where $G_\ell^{(d)}(k)$ represents the direct path gain between the origin and the ℓ^{th} source and $G_\ell^{(r)}(k, \hat{\mathbf{x}})$ is the reflection gain at the origin along the direction of $\hat{\mathbf{x}}$ for the ℓ^{th} source. Hence,

we obtain the spatial domain equivalent of (2) by substituting the spatial domain RTF from (9) - (11) as

$$P(\mathbf{x}'_q, k) = \sum_{\ell=1}^L S_\ell(k) \left(G_\ell^{(d)}(k) e^{ik \hat{\mathbf{x}}_\ell \cdot \mathbf{x}'_q} + \int_{\mathbb{S}^2} G_\ell^{(r)}(k, \hat{\mathbf{x}}) e^{ik \hat{\mathbf{x}} \cdot \mathbf{x}'_q} d\hat{\mathbf{x}} \right). \quad (12)$$

C. Modal coherence model

The spherical harmonic expansion of the far-field approximation of the Green's function is given by [60, pp. 27–33]

$$e^{ik \hat{\mathbf{x}}_\ell \cdot \mathbf{x}'_q} = \sum_{nm}^{\infty} 4\pi i^n Y_{nm}^*(\hat{\mathbf{x}}_\ell) b_n(kr) Y_{nm}(\hat{\mathbf{x}}'_q). \quad (13)$$

Using (13) in (12) and then by comparing it with (3), we obtain an analytical expression for α_{nm} in a reverberant room as [56]

$$\alpha_{nm}(k) = 4\pi i^n \sum_{\ell=1}^L S_\ell(k) \left(G_\ell^{(d)}(k) Y_{nm}^*(\hat{\mathbf{x}}_\ell) + \int_{\mathbb{S}^2} G_\ell^{(r)}(k, \hat{\mathbf{x}}) Y_{nm}^*(\hat{\mathbf{x}}) d\hat{\mathbf{x}} \right). \quad (14)$$

Taking into account the autonomous behavior of the reflective surfaces in a room (i.e., the reflection gains from the reflective surfaces are independent) and imposing the assumptions of uncorrelated sources, we previously established [56] a closed form expression for the modal coherence of a reverberant soundfield as

$$\mathbb{E}\left\{ \alpha_{nm}(k) \alpha_{n'm'}^*(k) \right\} = \sum_{\ell=1}^L \mathbb{E}\left\{ |S_\ell(k)|^2 \right\} \left(\Upsilon_{nm}^{n'm'}(\hat{\mathbf{x}}_\ell) \mathbb{E}\left\{ |G_\ell^{(d)}(k)|^2 \right\} + \sum_{vu}^V \mathbb{E}\left\{ \gamma_{vu}^{(\ell)}(k) \right\} \Psi_{n,n',v}^{m,m',u} \right) \quad (15)$$

where $\mathbb{E}\{\cdot\}$ denotes expected value and

$$\sum_{vu}^V \mathbb{E}\left\{ \gamma_{vu}^{(\ell)}(k) \right\} Y_{vu}(\hat{\mathbf{x}}) = \mathbb{E}\left\{ |G_\ell^{(r)}(k, \hat{\mathbf{x}})|^2 \right\} \quad (16)$$

$$\Upsilon_{nm}^{n'm'}(\hat{\mathbf{x}}_\ell) = C_{nn'} Y_{nm}^*(\hat{\mathbf{x}}_\ell) Y_{n'm'}(\hat{\mathbf{x}}_\ell) \quad (17)$$

$$\Psi_{n,n',v}^{m,m',u} = C_{nn'} W_{v,n,n'}^{u,m,m'} \quad (18)$$

$$C_{nn'} = 16\pi^2 i^{n-n'} \quad (19)$$

$$W_{v,n,n'}^{u,m,m'} = (-1)^m \sqrt{\frac{(2v+1)(2n+1)(2n'+1)}{4\pi}} \times \begin{pmatrix} v & n & n' \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} v & n & n' \\ u & -m & m' \end{pmatrix} \quad (20)$$

where (\cdot) in (20) represents Wigner-3j symbol [61]. Note that, though (15) was developed using a far-field sound propagation model, it can easily be written for near-field sound sources by replacing (13) with its near-field counterpart [56]. Furthermore, for the temporal processing, it is common to estimate the

expected value by applying the exponential moving average technique on the instantaneous measurements, i.e.,

$$\mathbb{E}\left\{\alpha_{nm}(k, \tau)\alpha_{n'm'}^*(k, \tau)\right\} = (1 - \beta) \alpha_{nm}(k, \tau) \times \alpha_{n'm'}^*(k, \tau) + \beta \mathbb{E}\left\{\alpha_{nm}(k, \tau - 1) \alpha_{n'm'}^*(k, \tau - 1)\right\} \quad (21)$$

where $\beta \in [0, 1]$ is a smoothing factor.

IV. CNN-BASED DOA ESTIMATION

CNN is a popular technique in the deep learning domain, and is predominantly used in computer vision applications. The input, often a 2D or 3D tensor, goes through multiple convolution filters followed by a traditional fully-connected neural network. In this work, we pose the DOA estimation problem as an image-classification problem where the input image represents the modal coherence of the soundfield.

A. Feature selection

Intuitively, the soundfield coefficients α_{nm} work as natural beamformers in the modal domain due to the inherent properties of the spherical harmonic functions. Hence, the energy distribution of α_{nm} among different modes can be used as a clue for understanding the source directionality. However, there only exists a limited number of active modes in the low frequencies which might prove insufficient to train a neural network for high spatial resolution scenario, especially in a reverberant environment. Therefore, we use the modal coherence model of (15) to construct our input feature. For a multi-source scenario, it is common to assume W-disjoint orthogonality [30] in the STFT domain, i.e., only a single sound source remains active in each TF bin of the STFT spectrum. Under the W-disjoint orthogonality assumption, (15) takes the following form

$$\mathbb{E}\left\{\alpha_{nm}(k)\alpha_{n'm'}^*(k)\right\} = \mathbb{E}\left\{\left|S_{\ell'}(k)\right|^2\right\} \left(\mathbb{E}\left\{\left|G_{\ell'}^{(d)}(k)\right|^2\right\} \Upsilon_{nm}^{n'm'}(\hat{\mathbf{x}}_{\ell'}) + \sum_{vu} \mathbb{E}\left\{\gamma_{vu}^{(\ell')}(k)\right\} \Psi_{n,n',v}^{m,m',u}\right) \quad (22)$$

where $\ell' \in [1, L]$. Note that, (22) remains true for a single-source scenario as well.

When dealing with audio signals having variable spectral densities, e.g., speech signals, it is intuitive to select a feature based on the relative transfer function to make the feature independent to the variations in the input signal [62]. Following

a similar reasoning, we define the relative modal coherence (RMC) as

$$\begin{aligned} & \mathbb{E}\left\{\alpha_{nm}(k)\alpha_{n'm'}^*(k)\right\} / \mathbb{E}\left\{\alpha_{00}(k)\alpha_{00}^*(k)\right\} \\ &= \frac{\mathbb{E}\left\{\left|G_{\ell'}^{(d)}(k)\right|^2\right\} \Upsilon_{nm}^{n'm'}(\hat{\mathbf{x}}_{\ell'}) + \sum_{vu} \mathbb{E}\left\{\gamma_{vu}^{(\ell')}(k)\right\} \Psi_{n,n',v}^{m,m',u}}{\mathbb{E}\left\{\left|G_{\ell'}^{(d)}(k)\right|^2\right\} \Upsilon_{00}^{00}(\hat{\mathbf{x}}_{\ell'}) + \sum_{vu} \mathbb{E}\left\{\gamma_{vu}^{(\ell')}(k)\right\} \Psi_{0,0,v}^{0,0,u}} \\ &= \frac{\mathbb{E}\left\{\left|G_{\ell'}^{(d)}(k)\right|^2\right\} \Upsilon_{nm}^{n'm'}(\hat{\mathbf{x}}_{\ell'}) + \sum_{vu} \mathbb{E}\left\{\gamma_{vu}^{(\ell')}(k)\right\} \Psi_{n,n',v}^{m,m',u}}{4\pi \mathbb{E}\left\{\left|G_{\ell'}^{(d)}(k)\right|^2\right\} + \frac{16\pi^2}{\sqrt{4\pi}} \mathbb{E}\left\{\gamma_{00}^{(\ell')}(k)\right\}} \end{aligned} \quad (23)$$

where (23) is derived using (17) - (20). From (23) it is evident that the relative modal coherence has a direct relation with the source position in a particular room. However, with multiple active sources in a strong reverberant environment, (23) introduces additional complexity due to the additive terms in the denominator. On the other hand, the modal coherence of (22) offers a simpler alternative to train a CNN due to the fact that the mode-independent term $\left\{\left|S_{\ell'}(k)\right|^2\right\}$ of (22) acts merely as a constant scaling factor across different TF bins without altering the relative strength between different modes inside a TF bin. The use of (22) as a feature reduces the complexity by eliminating location-dependency from the denominator compared to RMC. Hence, we pose the DOA estimation problem as an image-identification problem for CNN where the feature snapshot is defined as the modal coherence of the soundfield in the individual TF bins of the STFT spectrum

$$\hat{\mathcal{F}}_{\text{mc}}(k) = \left\{ \mathbb{E}\left\{\alpha_{nm}(k)\alpha_{n'm'}^*(k)\right\} : n \in [0, N], m \in [-n, n], n' \in [0, N], m' \in [-n', n'] \right\} \quad (24)$$

where $\hat{\mathcal{F}}_{\text{mc}}$ is considered as an image consisting of $[\mathcal{N} \times \mathcal{N}]$ complex-valued pixels with $\mathcal{N} = (N+1)^2$ is the total number of modes. Note that, $\hat{\mathcal{F}}_{\text{mc}}$ is a frequency-dependent function due to the frequency dependency of α_{nm} , hence, we need to collect $\hat{\mathcal{F}}_{\text{mc}}$ from different frequency bands for training so that the model can learn the frequency variations of the feature for the same source position. This deviation is analogous to the transformed image conundrum in an image-identification problem.

We also need to train the CNN model for various amplification levels due to the presence of the source PSD term $\mathbb{E}\left\{\left|S_{\ell'}(k)\right|^2\right\}$ in $\hat{\mathcal{F}}_{\text{mc}}$ (analogous to train a neural network to accommodate the differences in brightness of the same image). This can be achieved through training the CNN model with any non-white random audio signal such that the signal has a variable PSD in both time and frequency directions of the STFT spectrum.

Finally, since a CNN model is best suited to work with real data, we convert our 2D complex-valued feature $\hat{\mathcal{F}}_{\text{mc}}$ into

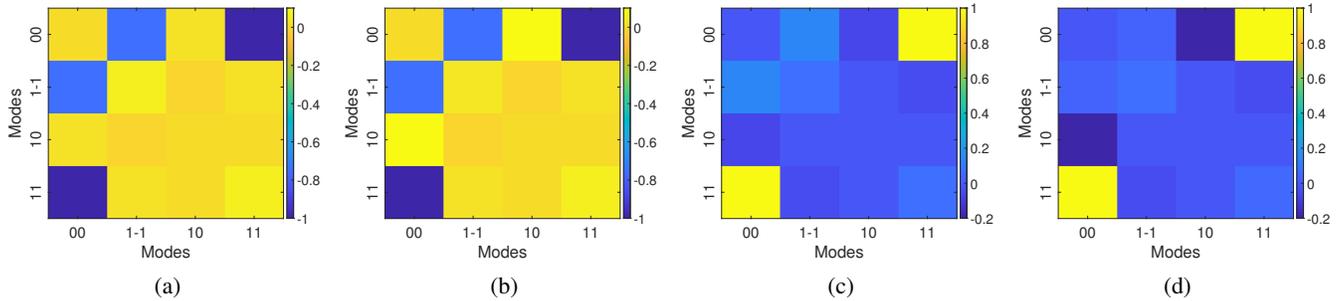


Fig. 2. Normalized $|\hat{\mathcal{F}}_{\text{mc}}|$ at different time instants. The snapshot is taken at 1500 Hz with a speech source present at (a)-(b) $(\theta, \phi) = (60^\circ, 60^\circ)$ and (c)-(d) $(\theta, \phi) = (60^\circ, 120^\circ)$.

corresponding 3D tensor \mathcal{F}_{mc} of $[\mathcal{N} \times \mathcal{N} \times 2]$ dimension such that

$$\mathcal{F}_{\text{mc}} = \left\langle \left\langle \mathcal{R}\{\hat{\mathcal{F}}_{\text{mc}}\}, \mathcal{I}\{\hat{\mathcal{F}}_{\text{mc}}\} \right\rangle \right\rangle_3 \quad (25)$$

where $\langle \langle \cdot, \cdot \rangle \rangle_3$ stacks two matrices in the 3rd dimension and $\mathcal{R}\{\cdot\}$ and $\mathcal{I}\{\cdot\}$ denote the real and imaginary part, respectively.

Fig. 2 shows the normalized snapshots of \mathcal{F}_{mc} captured at random time instants at 1500 Hz. For a CNN model to work with our input features, we want them to be time-independent for the same source position in a room irrespective of the nature of the audio signal. Indeed, as we observe from Fig. 2, \mathcal{F}_{mc} changes as a function of source angle and remains fairly constant across time.

B. TF bin processing

During both training and evaluation phases, the proposed CNN framework processes each TF bin independently, i.e., it learns the directional patterns based on the spatial distribution of the TF bin energy. Hence, it is important to consider only the TF bins with a significant energy level to avoid misleading the neural network. However, due to the sparse nature of speech signals in both time and frequency, a large proportion of the TF bins usually ends up having low energy. The sparsity in time can be addressed with a suitably designed voice activity detector, but the sparsity along the frequency can still mislead the CNN. Hence, to exclude the low-energy TF bins from the training and evaluation datasets, a energy-based pre-selection of TF bins is required where we drop all the TF bins with an average energy below a certain threshold. If \mathcal{T}_{all} is denoted as the collection of all the TF bins, we can define a new set $\mathcal{T}_{\text{act}} \subseteq \mathcal{T}_{\text{all}}$ such that

$$\mathcal{T}_{\text{act}} = \{\kappa \in \mathcal{T}_{\text{all}} : E_{\kappa} \geq E_{\text{min}}\} \quad (26)$$

where E_{κ} is the average energy of the spatial coherence matrix for the κ^{th} TF bin and E_{min} is the minimum energy threshold. This lowest energy threshold can be a preset based on empirical measurements, or it can be set dynamically based on the average energy of all the TF bins in the processing block. However, the average energy of the processing block can be low when the number of low-energy TF bins is high. We can also set the minimum energy level at the \mathcal{K}^{th} percentile of the average energy distribution of all TF bins, where \mathcal{K} is

usually large for speech signal, as long as we are able to make at least a high-level prediction about the energy distribution among the TF bins. Note that, (26) should be applied at both training and evaluation stage, however, E_{min} does not need to be the same.

A second issue may arise when a TF bin violates the W-disjoint orthogonality principle. In the proposed algorithm, CNN predicts the most dominant source in each TF bin which are later combined in a clustered histogram to reach a global outcome. However, as shown in [30], the number of TF bins violating the W-disjoint orthogonality increases as the number of simultaneous sources increases. When a TF bin contains significant energy from multiple sources, the prediction of the CNN model can be arbitrary. To circumvent the uncertainty in prediction due to the violation of W-disjoint orthogonality principle, we only consider the predictions in the TF bins where the CNN model predicts a single DOA with a high confidence level. Hence, if we define the probability score for each TF bin as

$$\mathcal{P}_{\kappa, \Theta} = \{\mathcal{P}_{\kappa}(\theta)\}_{\theta \in \Theta} \quad (27)$$

$$\mathcal{P}_{\kappa, \Phi} = \{\mathcal{P}_{\kappa}(\phi)\}_{\phi \in \Phi} \quad (28)$$

where $\mathcal{P}_{\kappa}(\theta)$ and $\mathcal{P}_{\kappa}(\phi)$ are the probability scores of the corresponding elevation and azimuth classes at the κ^{th} TF bin, the final DOA estimation at the evaluation stage should be based on the set $\mathcal{T}_{\text{test}} \subseteq \mathcal{T}_{\text{act}}$ such that

$$\mathcal{T}_{\text{test}} = \left\{ \kappa \in \mathcal{T}_{\text{act}} : \max \{\mathcal{P}_{\kappa, \Theta}\} \geq \mathcal{P}_{\text{min}} \text{ and } \max \{\mathcal{P}_{\kappa, \Phi}\} \geq \mathcal{P}_{\text{min}} \right\} \quad (29)$$

where \mathcal{P}_{min} denotes the minimum confidence level.

C. CNN Architecture

We utilize a CNN to estimate source DOA based on the local connectivity of the modal coherence coefficients. A CNN topology typically consists of multiple convolution layers followed by fully-connected networks. For DOA estimation, we perform multi-output multi-class classification where we share the same convolution layer structure to predict both the azimuth and elevation using separate fully connected heads at the last stage - each responsible for predicting either azimuth

or elevation. We opt for a classification-based approach due to the limited resolution of the practical dataset. However, the proposed technique can be studied with a regression-based model subject to the availability of a denser training grid to learn the evolution of dynamic reverberation characteristics.

In each convolutional layer, we use 64 spatial filters of 2×2 size to learn the spatial coherence pattern for each desired point in a predefined DOA grid. As our feature is defined as the modal coherence for each TF bin, it is important to consider 2D filters in the convolution layers. A rectified linear unit (ReLU) activation follows the convolution layer at each stage. The evaluation is done with 8 convolution layers with zero padding to keep the output size same for each layer. The final convolution layer is connected to 2 fully-connected layers that use ReLU activation. Finally, two separate fully connected heads responsible for azimuth and elevation estimation, respectively, are used with Sigmoid activation. A Sigmoid activation is chosen over Softmax in the last stage as it allows us to perform prediction-based TF bin selection to remove the bins with low confidence, as described in Section IV-B.

Due to the W-disjoint orthogonality assumption, ideally each TF bin is designated with a single DOA and can be classified using a multi-class classification network using a Softmax activation-based categorical cross-entropy loss. However, in a practical environment with multiple simultaneously active sources, it is unrealistic to expect each of the TF bins to honor the W-disjoint orthogonality. Therefore, it is possible to find occasional TF bins whose energy are contributed by multiple sound sources. Such a TF bin produces a feature snapshot which does not match with any of the patterns learned by the model during the single-source training stage. In such cases it is expected that the output will not have a large prediction score for any of the classes. Hence, we use binary cross-entropy loss function with Sigmoid activation instead of categorical cross-entropy in order to independently predict the probability of each individual class in every TF bin. This approach allows us to enforce the criterion mentioned in (29).

Detailed parameter settings are included in the experimental results section.

D. Training the model

In the training stage, we train the model based on the feature snapshots in a single source scenario. Each training data is labeled independently for azimuths and elevations. The model is trained for the elevation set Θ and the azimuth set Φ in different azimuth and elevation planes, respectively. Once the model learns the patterns for each of the intended directions, we independently predict the elevations and azimuths for any number of concurrent sources as long as the W-disjoint orthogonality principle majorly holds. This is a more realistic approach than training the model for each possible angular combination [54], [63] which becomes a resource-intensive operation as the number of classes or the number of simultaneous sources increases. Furthermore, the proposed method does not require retraining the model every time an additional source appears in the mixture.

Algorithm 1: Algorithm for DOA estimation - training stage

- Data:** $\Theta, \Phi, \alpha_{nm}(\theta, \phi) \forall nm, \forall \theta \in \Theta, \forall \phi \in \Phi$
- 1 Calculate spatial coherence $\mathbb{E}\{\alpha_{nm}(k)\alpha_{n'm'}^*(k)\}$ in each TF bin using (21);
 - 2 Get $\hat{\mathcal{F}}_{mc} \forall \theta \in \Theta, \forall \phi \in \Phi$ using (24);
 - 3 Apply (26) to filter out low energy TF bins and get \mathcal{T}_{act} ;
 - 4 Use \mathcal{T}_{act} to train the model using the parameters in Table I independently for Θ and Φ ;
 - 5 Save the model
-

Algorithm 2: Algorithm for DOA estimation - evaluation stage

- Data:** $\alpha_{nm} \forall nm$
- 1 Calculate spatial coherence $\mathbb{E}\{\alpha_{nm}(k)\alpha_{n'm'}^*(k)\}$ in each TF bin using (21);
 - 2 Get $\hat{\mathcal{F}}_{mc}$ using (24);
 - 3 Apply (26) to filter out low energy TF bins and get \mathcal{T}_{act} ;
 - 4 Calculate the probability of each classes in Θ and Φ for the TF bins in \mathcal{T}_{act} using the model saved during training;
 - 5 Apply (29) to get \mathcal{T}_{test} ;
 - 6 Apply (30) to form the prediction multiset \mathcal{X} ;
 - 7 **if** $L == 1$ **then**
 - 8 | Use (31) to estimate DOA;
 - 9 **else**
 - 10 | Using a suitable clustering algorithm, divide \mathcal{X} into L clusters;
 - 11 | Use (32) to estimate L source directions.
 - 12 **end**
-

E. DOA estimation

First, we jointly pick the highest probable elevation and azimuth classes for each TF bin in \mathcal{T}_{test} to form prediction multiset \mathcal{X}

$$\mathcal{X} = \left\{ \left(\arg \max_{\theta \in \Theta} \{f : \theta \mapsto \mathcal{P}_\kappa(\theta)\}, \arg \max_{\phi \in \Phi} \{f : \phi \mapsto \mathcal{P}_\kappa(\phi)\} \right) : \kappa \in \mathcal{T}_{test} \right\}. \quad (30)$$

Subsequently, the simplest way of multi-source DOA estimation is to pick L largest peaks from the 2D histogram of \mathcal{X} , i.e.,

$$\{\hat{\mathbf{x}}\} = \{(\theta, \phi) \text{ of } L \text{ largest peaks in } \mathcal{X}\}. \quad (31)$$

However, in case of a noisy prediction for a multi-source environment, the aforementioned technique can cause erroneous results. For example, in a 2-source environment, if the true DOA of the prominent source lies between two adjacent classes, both the adjacent classes for the prominent source might occur more frequently than the true class corresponding to the weaker source. To avoid such a scenario, a more robust technique is to apply a suitable clustering algorithm, such as

TABLE I
EXPERIMENTAL PARAMETER SETTINGS

Name	Value
Model parameters	
\mathcal{P}_{\min}	0.5
E_{\min}	Percentile-based
\mathcal{K}	90
N	1
CNN parameters	
Input size	$[4 \times 4 \times 2]$
# Convolution layers	8
# Conv. filters	64 ($[2 \times 2]$)
# Dense layers	2 (512)

k-means [64] or density-based [65] clustering, to divide \mathcal{X} into L clusters and pick the peak in each cluster, i.e.,

$$\{\hat{\mathbf{x}}\} = \{(\theta, \phi) \text{ of the peak of } L \text{ clusters in } \mathcal{X}\}. \quad (32)$$

The training and evaluation steps are outlined in Algorithm 1 and 2, respectively.

V. EXPERIMENTAL RESULTS AND DISCUSSION

In this section, we present the experimental results, comparison, and discussion of the proposed algorithm with the contemporary counterparts.

A. Experimental methodology

We evaluated the proposed method in simulated and practical environments under different room conditions. The parameter settings used in the evaluation are listed in Table I. We assessed the performance of the model in 3 simulated room environments (room S1, S2, and S3 in Table II) generated using a RIR Generator [66] as well as with the recordings from a practical room (room P1 in Table II) in the presence of babble noise. The training and the tests were performed under the same room environment. The reverberation time (T_{60}) and direct to reverberation ratio (DRR) shown in Table II for room P1 were calculated using the techniques outlined in [67]. We only considered the first-order harmonic coefficients which need at least 4 microphones to calculate, however, in the result section, we used recordings from 9 microphones oriented on a spherical grid suggested in [68] for evaluating the proposed as well as the competing methods.

We used random mixed-gender speech signals from the TIMIT corpus [69] to synthesize the reverberant signals from the measured/simulated RIRs. The spherical harmonic decomposition, as described in Section III-A, was performed on the reverberant signals to calculate the spherical harmonic coefficients up to first order. Note that, the proposed method is independent of the type and shape of the sensor array as long as the array is capable of performing spherical harmonic decomposition³.

The CNN architecture has been discussed and presented in Section IV-C and Table I. The implementation was done in

³A number of alternate array structures are available in the literature for capturing spherical harmonic coefficients, a few can be found in [43]–[48].

TABLE II
TEST ENVIRONMENTS. d_{sm} DENOTES SOURCE TO MICROPHONE DISTANCE.

Room	Dimension	T_{60}	DRR	d_{sm}
P1	$[11 \times 7.5 \times 2.75]$ m	640 ms	-0.6 dB	2.8 m
S1	$[6 \times 4 \times 3]$ m	200 ms	-	1 m
S2	$[7 \times 6 \times 3]$ m	300 ms	-	1 m
S3	$[8 \times 6 \times 3]$ m	500 ms	-	1 m

Python using Keras [70] running on top of TensorFlow [71]. For the proposed method, the TF bin-level predictions were accumulated and clustered using k-means algorithm (step 10 in Algorithm 2) assuming that the number of active sources was known as *a priori*, however, certain algorithms offer to cluster the data without the advance knowledge of the number of sources [65]. Furthermore, as k-means algorithm works with the Euclidean geometry, we converted the predicted DOAs to corresponding Cartesian coordinates on a unit sphere before clustering the data.

The processing was done at 16 kHz sampling frequency. The STFT used a 16 ms Hanning window, 50% overlap, and 256-point discrete Fourier transform (DFT). We only utilized the frequencies ranged 500 – 2000 Hz for DOA estimation. A 30s long speech was used to synthesize the data for training, however, the actual number of features reduced significantly after applying (26) to filter out low-energy TF bins. The majority of the results and discussions in this section are presented for azimuth estimation only considering the fact that the estimation of elevation is independent of the azimuth estimation and follows the same mechanism. However, in Section V-C3, we have demonstrated how a joint azimuth and elevation estimation can be performed using the proposed method.

B. Baseline methods and evaluation metrics

The performance of the proposed algorithm is compared with a recent CNN-based DOA estimation method proposed in [54] (subsequently denoted as “CNN-PH”) where it was already shown that “CNN-PH” outperforms conventional parametric methods like MUSIC and SRP-PHAT. For a fair comparison, we kept the CNN architecture and other evaluation criteria same in all possible ways. We used the same 9-microphone setup as described in Section V-A for the competing methods unless mentioned otherwise. The convolution filter size for “CNN-PH” was set to $[2 \times 1]$ as per the recommendation of the authors [54] whereas we applied $[2 \times 2]$ filters with the proposed method. The difference in filter size between the competing methods comes from the fact that the feature used in “CNN-PH” spans across the frequency band where multiple active sources can be present in the horizontal dimension. On the other hand, the proposed method uses the modal coherence snapshot of a single TF bin as a feature where only one active source is expected due to the assumption of W-disjoint orthogonality.

To evaluate the performance, we first defined the prediction error for the ℓ^{th} source in a single test by the angular difference

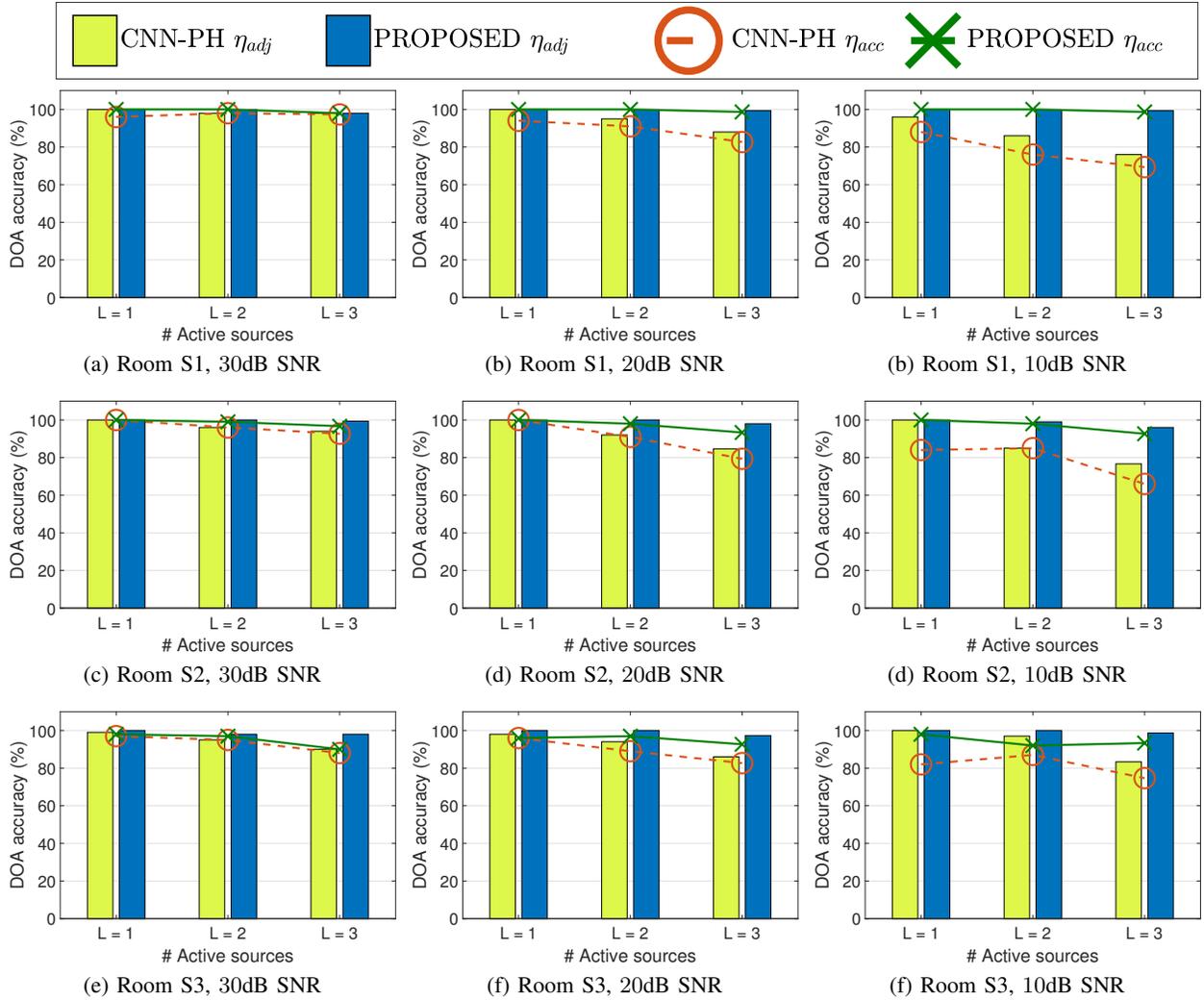


Fig. 3. Azimuth estimation under different simulated reverberant and noisy environments on a 45° elevation plane.

between the true and the estimated points at the origin of a unit sphere, i.e.,

$$\Delta_\ell = \cos^{-1} \left[\cos(\hat{\theta}_\ell) \cos(\theta_\ell) + \sin(\hat{\theta}_\ell) \sin(\theta_\ell) \cos(\hat{\phi}_\ell - \phi_\ell) \right]. \quad (33)$$

As we are posing the DOA estimation as a classification problem, the mean error can be misleading unless the angular difference between adjacent classes are very small. Hence, instead we propose to use performance metrics based on estimation accuracy. At first, we define the multi-source DOA classification accuracy as the percentage of the correct predictions, i.e.,

$$\eta_{\text{acc}} = \frac{\mathcal{M}_{\{\{\Delta_\ell\}\}}(0)}{\#\{\{\Delta_\ell\}\}} \times 100\% \quad (34)$$

where $\mathcal{M}_{\mathcal{X}}(\vartheta)$ denotes the multiplicity of ϑ in the multiset \mathcal{X} , $\{\{\Delta_\ell\}\}$ is a multiset containing $\Delta_\ell \forall \ell$ for all the tests, and $\#\{\{\Delta_\ell\}\}$ denotes the cardinality of the underlying multiset. Note that, for a single test, the sequence of the true and estimated DOAs need not be in the same order, hence, we map them in such a way that $\mathcal{M}_{\{\{\Delta_\ell\}\}}(0)$ is maximized.

Occasionally, the definition of (34) may fail to offer the full picture as it does not take it into consideration how far a wrong prediction deviates from the true value although the adjacent classes are highly correlated in a DOA classification task. Hence, we define another accuracy metric, termed as adjacent accuracy, where we consider the predictions for adjacent classes as true positives as well, i.e.,

$$\eta_{\text{adj}} = \frac{\mathcal{M}_{\{\{\max[0, \Delta_\ell - \Delta_\Omega]\}\}}(0)}{\#\{\{\Delta_\ell\}\}} \times 100\% \quad (35)$$

where Δ_Ω is the angular separation between two adjacent classes. A high η_{adj} with low η_{acc} indicates that the transition of the feature pattern is not very sharp between the adjacent classes, a phenomenon expected in a noisy environment.

All the results presented in the subsequent sections are based on the accumulation of the results of 50 random experiments in each test case. Each experiment was evaluated with random source positions and subsequently added with random Gaussian noise unless specified otherwise.

C. Results and discussions

In this section, we discuss and compare DOA estimation performances under different criteria and room environments. During the experiments, the microphone array was placed at the center of the room at 1 m height. For the proposed method, we completed the training once per room considering a single-source scenario and used the same trained model at the testing stage irrespective of the number of simultaneously active sources. A 30s long speech data were synthesized during the training stage, however, we only used the top 10% STFT bins based on TF bin energy to train the network⁴.

1) **Azimuth estimation for a fixed elevation:** The first set of experiments considered the same elevation plane at 45° for both training and testing. We considered uniformly spaced azimuth points at 10° interval (i.e., $J = 36$) which makes the angular separation between adjacent classes $\Delta_\Omega = 7.07^\circ$ on the 45° elevation plane. For each room, we emulated 2 different signal to noise ratio (SNR) by adding white Gaussian noise and evaluated the performance for up to 3 active sources, i.e., $L = [1, 3]$. As “CNN-PH” was originally designed to be trained for all possible angular combinations in an L -source DOA estimation, we trained “CNN-PH” for 36 and 1260 unique angular combinations based on 36 azimuth classes for $L = 1$ and 2, respectively⁵. However, for testing with $L = 3$, we trained “CNN-PH” for 2-source mixture (1260 angular combinations) to understand the performance in a dynamic acoustic scenario. In contrast, the proposed method was always trained for the single-source scenario, e.g., 36 unique cases for this experiment, irrespective of the number of sources in the testing environment.

Fig. 3 shows DOA accuracy of the competing methods under different scenarios. At SNR = 30dB in Fig. 3, we observe that both the methods perform well for $L = 1$ and 2 although the proposed method consistently exhibits slightly better performance. For 3-source combination under stronger reverberation, the proposed method holds the level of adjacent accuracy η_{adj} while “CNN-PH” shows performance degradation for both the metrics. For the noisy environments of SNR = 20dB and 10dB in Fig. 3, the performance distinctions are more prominent where the proposed algorithm outperforms “CNN-PH” in each scenario. The use of modal coherence as learning feature ensures steady performance of the proposed algorithm at low SNR. We can also observe “CNN-PH” suffers significant performance issues for $L = 3$ due to the fact that we did not train “CNN-PH” for all possible 3-source combinations.

For reference, Fig. 4 plots the TF bin prediction histogram in room S3 for $L = 2$ and 3 along with the true azimuths. The histogram shows a clear peak at each true azimuth location which can be separated using a suitable clustering algorithm.

2) **Performance in a practical room with babble noise:** We conducted the next set of experiments in a big hall with strong reverberation, we named it room P1 in Table II. The recording was performed with an *Eigenmike* [72], however,

⁴As an instance, the average size of the training dataset in Section V-C1 was 10,505 samples per class, each being a $[4 \times 4 \times 2]$ tensor.

⁵The training process for “CNN-PH” is outlined in [54, pp. 13]

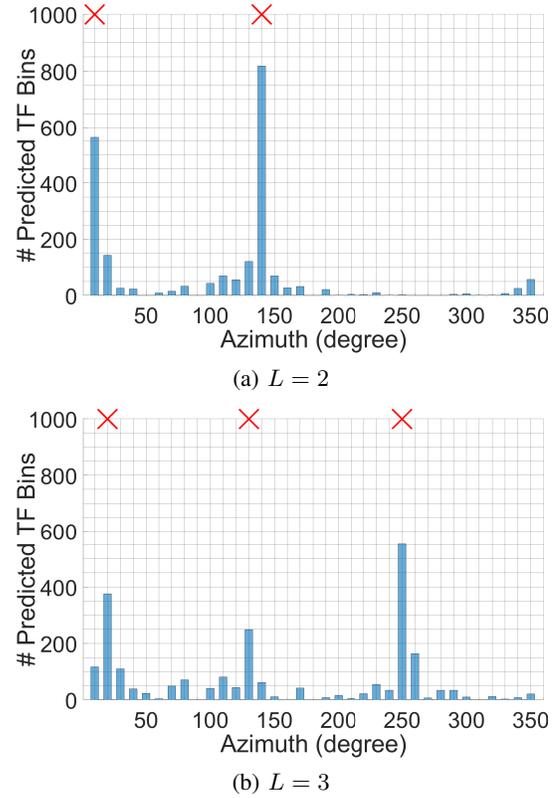


Fig. 4. TF bin prediction histogram in room S3 ($T_{60} = 500$ ms). Red crosses denote the ground truths.

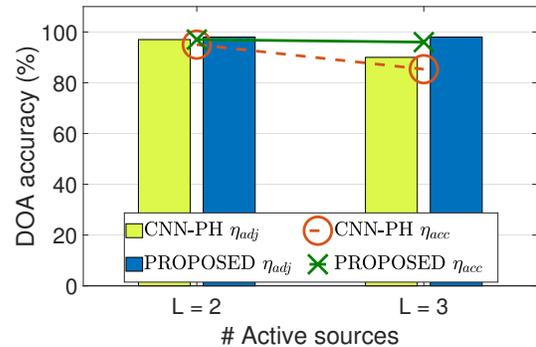


Fig. 5. Azimuth estimation accuracy with practical recordings with babble noise at 10dB SNR. The tests were performed on 95° elevation plane.

only first-order harmonics were captured for this task. The source was placed at a 2.8 m distance from the array in a uniform azimuth grid of 30° interval ($J = 12$) on a 95° elevation plane. Directional babble noise was added to the recordings at 10dB SNR from multiple random locations. This time we trained “CNN-PH” with all possible angular combinations for both $L = 2$ (132 angular combinations) and $L = 3$ (1320 angular combinations) while the proposed method used the same strategy of single-source training for 12 classes. The comparative performance is shown in Fig. 5 where the proposed algorithm shows a significantly better accuracy than “CNN-PH”, especially for $L = 3$, despite “CNN-PH” being trained for all possible angular combinations

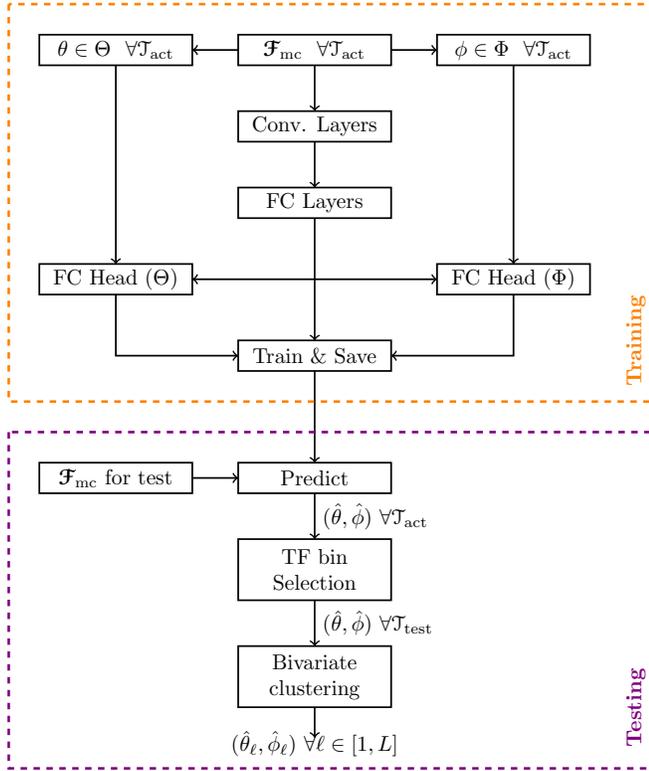


Fig. 6. A block diagram for joint estimation of azimuth and elevation.

in each case. Note that, we found no significant performance improvement for higher harmonic orders. This can be due to low spatial resolution of the training data. The higher order modes can be useful with a denser source distribution, at high frequencies, or when a regression-based model is used.

3) **Joint estimation of azimuth and elevation:** So far, we have shown results only for azimuth estimation based on the proposition that the proposed algorithm can estimate azimuth and elevation simultaneously without interfering with each other. In this section, we are going to validate this proposition by performing a full DOA estimation in room S2. We designed a 3D uniform spatial grid with 30° resolution for azimuths ($J = 12$) and 20° resolution for elevations. Furthermore, we considered the elevation range $30^\circ - 150^\circ$. That makes a total of 7 unique elevation classes ($I = 7$) and a total 84 points on the 3D DOA grid. The rest of the simulation criteria remain the same as Section V-C1.

We slightly modified the CNN architecture for this section to accommodate the joint estimation in an efficient manner. As in the previous experiments, we calculated the feature snapshot for each TF bin, but this time we labeled them separately for azimuth and elevation. The CNN architecture remains the same for the most part except at the last layer when we branched out 2 identical but separated fully connected heads and supplied them with azimuth and elevation labels, respectively. Hence, at the testing stage, the system outputs two separate prediction sets for azimuth and elevation - one from each separated head. Note that, due to the independent estimation strategy, it is important to jointly pack the predicted

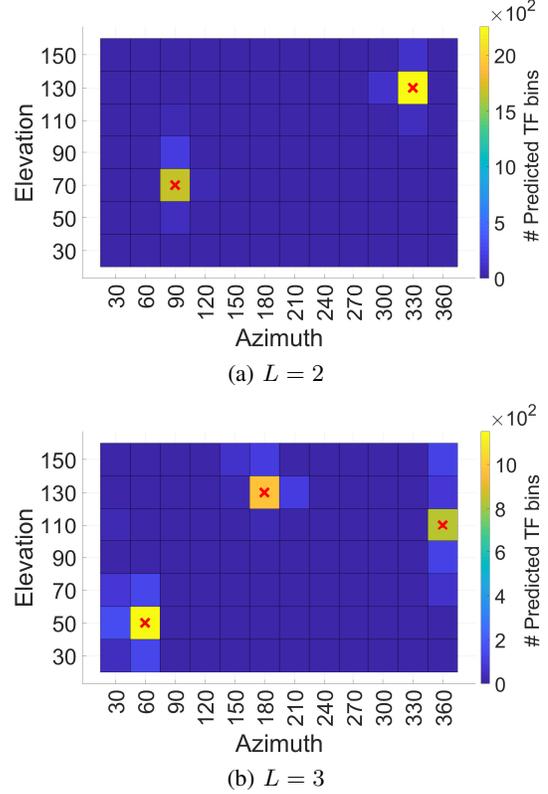


Fig. 7. Color map for joint estimation of azimuth and elevation with the proposed method in room S2 at 30dB SNR. Red crosses denote the ground truths. The accuracy of a joint estimation over 50 tests was found to be similar compared to the standalone azimuth estimation.

azimuths and elevations for each TF bin so that the estimated angles from the same source remain together. Subsequently, we clustered the prediction outcomes using the bivariate k-means clustering algorithm. A block diagram for the joint estimation of azimuth and elevation is shown in Fig. 6⁶. It is worth mentioning that the proposed algorithm can readily be expanded for full source localization through additional training for radius-dependency or corresponding Cartesian coordinates due to its ability of independent estimation of different location parameters.

The outcome of the experiments is shown in Fig. 7 in terms of colored heat map based on the number of predicted TF bins for each DOA. It is clear from the figure that the proposed method had no difficulties in predicting full DOA in the same manner as with the azimuth predictions. More importantly, the accuracy of the joint DOA estimation was found to be similar to that of a standalone azimuth estimation of Fig. 3 (and hence, the accuracy plots are not shown separately for this section). This is expected behavior as the azimuth and elevation estimation processes are independent and should not be affected due to the joint processing.

4) **Azimuth estimation on a different elevation plane:** In this section, we analyze the performance of the proposed algorithm when training and testing were performed on different

⁶Fig. 6 does not necessarily depict the actual processing flow, rather a visual aid for understanding the task.

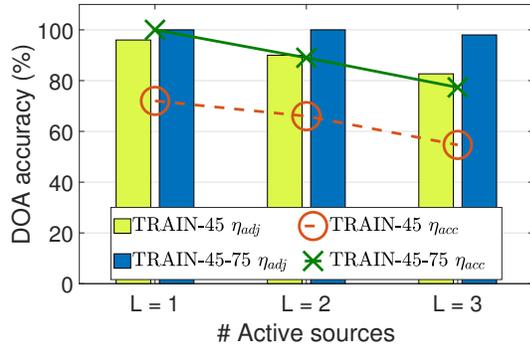


Fig. 8. Azimuth estimation on 60° elevation plane when training was performed on different elevation plane. TRAIN-45 denotes the case when training was performed on 45° elevation only whereas for TRAIN-45-75, training was performed with data from 45° and 75° elevations.

elevation planes. For the purpose of this section, we used room S2 at 30dB SNR. The tests were performed for sources on 60° elevation plane while the training data were obtained from a different elevation. In Fig. 8, we show the estimation accuracy for 2 distinct cases - when training data were obtained from (case 1) 45° elevation plane only and (case 2) 45° and 75° elevation planes. We observe a clear improvement in case 2 over case 1 due to the fact that when we trained the network on 2 different elevation planes, the model learned the evolution of feature for change in elevation and predicted azimuths in an arbitrary elevation plane more accurately. As the machine learning algorithms take a data-driven approach, it is possible to further improve the performance by training on additional planes.

It is worth noting that the proposed feature snapshot \mathcal{F}_{mc} is mostly comprised of the spherical harmonics where the dependency on θ and ϕ come through independent Legendre and exponential functions, respectively, as shown in (4). Therefore, the impact of elevation change on \mathcal{F}_{mc} comes mainly as a constant scaling factor. For this reason, even for case 1 when training and testing were done in separate individual elevation planes, the model didn't entirely fail, rather gets confused by the reverberation, noise and other non-linear distortions. This is apparent from Fig. 8 where we observe a better accuracy in terms of η_{adj} but a significant difference with η_{acc} .

5) **Impact of source to microphone distance:** We investigate the impact of the varying source to microphone distance on the proposed algorithm. We used the same simulated room S2 with the exception that we increased the dimension of the room to $[8 \times 8 \times 4]$ m for this particular section in order to have a larger range for distance variation. The microphone array position remained at the center of the room, however, we varied the source position between 0.5 – 3 m from the microphone array. The training was performed at a fixed distance of 1 m. The plots in Fig. 9 suggests that there is no significant change in estimation accuracy for varying source to microphone distances during the test. To understand the behavior, we examine the analytical expression of α_{nm} for the direct path for ℓ^{th} source [60, pp. 31]

$$\alpha_{nm}^{(\ell)}(k, r) = ikh_n(kr_\ell)Y_{nm}^*(\hat{\mathbf{x}}_\ell) \quad (36)$$

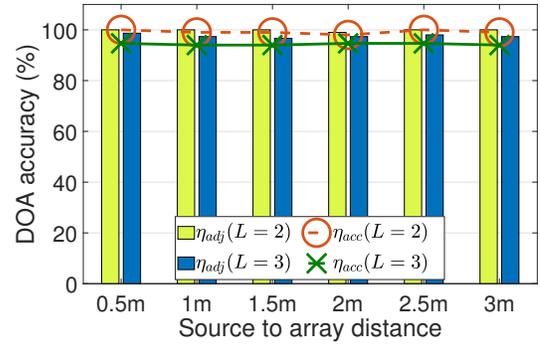


Fig. 9. Azimuth estimation performance of the proposed algorithm with different source to microphone distances. The training was performed with sources at 1 m distance from the array center on a 45° elevation plane.

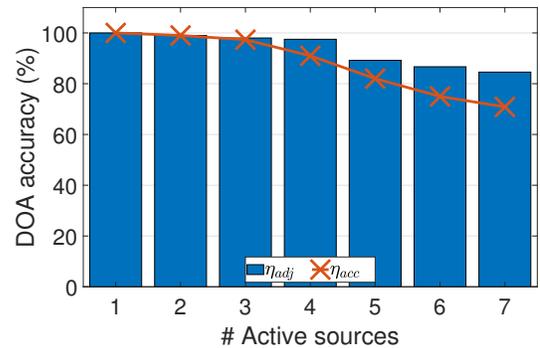


Fig. 10. Azimuth estimation performance of the proposed algorithm with different number of active sources on a 45° elevation plane.

where $h_n(\cdot)$ is the spherical Hankel function of the first kind. From (36) it is clear that the radial dependency comes through the Hankel function $h_n(kr_\ell)$ together with the frequency-dependent k . As we trained our model for different frequencies, the impact of varying $h_n(kr_\ell)$ on the feature pattern is already captured during the training even with a fixed radius, hence, any radial change does not pose a major threat to the performance of the proposed algorithm.

6) **Number of active sources:** In the last set of experiments, we tried increasing the number of sources on the same 45° elevation plane in the acoustic scene of room S2 at 30dB SNR. As we observe in Fig. 10, the accuracy gradually decreases with the increasing number of sources. The performance issue can be contributed by an increased violation of W-disjoint orthogonality with an increased number of sources. However, examining the histograms of random individual tests, we also found many instances when the performance degradation was caused by the failure of the k-means clustering algorithm and the ambiguity in the histogram for nearby sources. It is possible to improve the performance with a careful selection of a more robust clustering algorithm, however, the investigation for a better clustering algorithm is out of the scope of this work. To avoid ambiguity due to nearby sources, we can impose a restriction for maintaining a minimum distance between two sources before applying this algorithm.

VI. CONCLUSIONS

In this paper, we proposed a modal coherence based feature to train a convolutional neural network for DOA estimation realized in the spherical harmonic domain. We offered a new perspective by introducing a single-source training scheme for multi-source localization in a reverberant environment. The proposed strategy saves significant time and resources during the training stage as well as allows us to reuse the same trained model during the testing stage irrespective of the number sources in the acoustic mixture. Furthermore, the proposed method is capable of performing parallel azimuth and elevation estimation, hence, allows us to perform full DOA estimation without affecting the estimation accuracy compared to standalone azimuth estimation. Several existing works have already shown that the application of deep learning algorithms in DOA estimation can alleviate the limitations of the parametric approaches such as MUSIC. We further contribute to the cause by proposing a method that performs better than the contemporary CNN-based methods in dynamic acoustic environments, requires significantly less resource and time for training, and predict the DOA based on a single training model for a room irrespective of the number of sources.

This work presents multiple future research directions. The proposed technique can be studied with a regression-based prediction model to achieve DOA estimation over a continuous grid. It can also be investigated to develop a generalized model for different reverberant scenarios which is capable of predicting DOA in versatile room environments.

REFERENCES

- [1] Y. Hioka, K. Furuya, K. Kobayashi, K. Niwa, and Y. Haneda, "Underdetermined sound source separation using power spectrum density estimated by combination of directivity gain," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 21, no. 6, pp. 1240–1250, 2013.
- [2] A. A. Nugraha, A. Liutkus, and E. Vincent, "Multichannel audio source separation with deep neural networks," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 24, no. 9, pp. 1652–1664, 2016.
- [3] S. Gannot, E. Vincent, S. Markovich-Golan, and A. Ozerov, "A consolidated perspective on multimicrophone speech enhancement and source separation," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 4, pp. 692–730, 2017.
- [4] C. Busso *et al.*, "Smart room: Participant and speaker localization and identification," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, vol. II, 2005, pp. ii–1117.
- [5] K. Nakadai, T. Lourens, G. O. Hiroshi, and H. Kitano, "Active audition for humanoid," in *Natl. Conf. Artif. Intell.*, 2000, pp. 832–839.
- [6] S. Yamamoto *et al.*, "Design and implementation of a robot audition system for automatic speech recognition of simultaneous speech," in *IEEE Work. Autom. Speech Recognit. Underst.*, 2007, pp. 111–116.
- [7] D. T. Blumstein *et al.*, "Acoustic monitoring in terrestrial environments using microphone arrays: applications, technological considerations and prospectus," *J. Appl. Ecol.*, vol. 48, no. 3, pp. 758–767, 2011.
- [8] S. Chu, S. Narayanan, and C.-C. J. Kuo, "Environmental sound recognition with time-frequency audio features," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 17, no. 6, pp. 1142–1158, 2009.
- [9] B. W. Chen, C. Y. Chen, and J. F. Wang, "Smart homecare surveillance system: Behavior identification based on state-transition support vector machines and sound directivity pattern analysis," *IEEE Trans. Syst. Man, Cybern. Syst.*, vol. 43, no. 6, pp. 1279–1289, 2013.
- [10] M. Crocco, M. Cristani, A. Trucco, and V. Murino, "Audio surveillance: A systematic review," *ACM Comput. Surv.*, vol. 48, no. 4, 2016.
- [11] H. Krim and M. Viberg, "Two decades of array signal processing research: The parametric approach," *IEEE Signal Process. Mag.*, vol. 13, pp. 67–94, 1996.
- [12] R. O. Schmidt, "Multiple emitter location and signal parameter estimation," *IEEE Trans. Antennas Propag.*, vol. AP-34, no. 3, pp. 276–280, 1986.
- [13] R. Roy and T. Kailath, "ESPRIT-estimation of signal parameters via rotational invariance techniques," *Adapt. Antennas Wirel. Commun.*, vol. 37, no. 7, pp. 984–995, 1989.
- [14] D. Khaykin and B. Rafaely, "Coherent signals direction-of-arrival estimation using a spherical microphone array: Frequency smoothing approach," in *IEEE Work. Appl. Signal Process. to Audio Acoust.*, 2009, pp. 221–224.
- [15] H. Wang and M. Kaveh, "Coherent signal-subspace processing for the detection and estimation of angles of arrival of multiple wide-band sources," *IEEE Trans. Acoust.*, vol. 33, no. 4, pp. 823–831, 1985.
- [16] J. H. Dibiase, H. F. Silverman, and M. S. Brandstein, "Robust localization in reverberant rooms," in *Microphone Arrays*, 2001, pp. 157–180.
- [17] L. Birnie, T. D. Abhayapala, H. Chen, and P. N. Samarasinghe, "Sound source localization in a reverberant room using harmonic based MUSIC," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, 2019, pp. 651–655.
- [18] J. H. Dibiase, "A high-accuracy, low-latency technique for talker localization in reverberant environments using microphone arrays," Ph.D. dissertation, Brown University, 2000.
- [19] C. Knapp and G. Carter, "The generalized correlation method for estimation of time delay," *IEEE Trans. Acoust.*, vol. 24, no. 4, pp. 320–327, 1976.
- [20] A. Marti, M. Cobos, J. J. Lopez, and J. Escolano, "A steered response power iterative method for high-accuracy acoustic source localization," *J. Acoust. Soc. Am.*, vol. 134, no. 4, pp. 2627–2630, 2013.
- [21] L. O. Nunes *et al.*, "A steered-response power algorithm employing hierarchical search for acoustic source localization using microphone arrays," *IEEE Trans. Signal Process.*, vol. 62, no. 19, pp. 5171–5183, 2014.
- [22] M. V. Lima *et al.*, "A volumetric SRP with refinement step for sound source localization," *IEEE Signal Process. Lett.*, vol. 22, no. 8, pp. 1098–1102, 2015.
- [23] H. Do, H. F. Silverman, and Y. Yu, "A real-time SRP-PHAT source location implementation using stochastic region contraction(SRC) on a large-aperture microphone array," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, vol. 1, 2007, pp. 121–124.
- [24] H. Ye and R. D. DeGroat, "Maximum likelihood DOA estimation and asymptotic Cramer-Rao bounds for additive unknown colored noise," *IEEE Trans. Signal Process.*, vol. 43, no. 4, pp. 938–949, 1995.
- [25] J. C. Chen, R. E. Hudson, and K. Yao, "Maximum-likelihood source localization and unknown sensor location estimation for wideband signals in the near-field," *IEEE Trans. Signal Process.*, vol. 50, no. 8, pp. 1843–1854, 2002.
- [26] P. Stoica and K. C. Sharman, "Maximum likelihood methods for direction-of-arrival estimation," *IEEE Trans. Acoust.*, vol. 38, no. 7, pp. 1132–1143, 1990.
- [27] M. I. Mandel, D. P. Ellis, and T. Jebara, "An EM algorithm for localizing multiple sound sources in reverberant environments," in *Adv. Neural Inf. Process. Syst.*, 2007, pp. 953–960.
- [28] O. Schwartz, Y. Dorfan, E. A. P. Habets, and S. Gannot, "Multi-speaker DOA estimation in reverberation conditions using expectation-maximization," in *Int. Work. Acoust. Signal Enhanc.*, 2016, pp. 1–5.
- [29] O. Schwartz, Y. Dorfan, M. Taseska, E. A. P. Habets, and S. Gannot, "DOA estimation in noisy environment with unknown noise power using the EM algorithm," in *Jt. Work. Hands-free Speech Commun. Microphone Arrays*, 2017, pp. 86–90.
- [30] O. Yilmaz and S. Rickard, "Blind separation of speech mixtures via time-frequency masking," *IEEE Trans. Signal Process.*, vol. 52, no. 7, pp. 1830–1846, 2004.
- [31] X. Li, L. Girin, R. Horaud, and S. Gannot, "Multiple-speaker localization based on direct-path features and likelihood maximization with spatial sparsity regularization," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 10, pp. 1997–2012, 2017.
- [32] W. Zhang and B. D. Rao, "A two microphone-based approach for source localization of multiple speech sources," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 18, no. 8, pp. 1913–1928, 2010.
- [33] C. Liu, B. C. Wheeler, W. D. O'Brien, R. C. Bilger, C. R. Lansing, and A. S. Feng, "Localization of multiple sound sources with two microphones," *J. Acoust. Soc. Am.*, vol. 108, no. 4, pp. 1888–1905, 2000.
- [34] B. Rafaely and D. Kolossa, "Speaker localization in reverberant rooms based on direct path dominance test statistics," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, 2017, pp. 6120–6124.

- [35] O. Nadiri and B. Rafaely, "Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 22, no. 10, pp. 1494–1505, 2014.
- [36] D. Levin, E. A. P. Habets, and S. Gannot, "On the angular error of intensity vector based direction of arrival estimation in reverberant sound fields," *J. Acoust. Soc. Am.*, vol. 128, no. 4, pp. 1800–1811, 2010.
- [37] A. H. Moore, C. Evers, P. A. Naylor, D. L. Alon, and B. Rafaely, "Direction of arrival estimation using pseudo-intensity vectors with direct-path dominance test," in *Eur. Signal Process. Conf.*, 2015, pp. 2296–2300.
- [38] A. H. Moore, C. Evers, and P. A. Naylor, "Direction of arrival estimation in the spherical harmonic domain using subspace pseudointensity vectors," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 1, pp. 178–192, 2017.
- [39] S. Hafezi, A. H. Moore, and P. A. Naylor, "Augmented intensity vectors for direction of arrival estimation in the spherical harmonic domain," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 10, pp. 1956–1968, 2017.
- [40] S. Tervo and A. Politis, "Direction of arrival estimation of reflections from room impulse responses using a spherical microphone array," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 23, no. 10, pp. 1539–1551, 2015.
- [41] L. Kumar and R. M. Hegde, "Near-field acoustic source localization and beamforming in spherical harmonics domain," *IEEE Trans. Signal Process.*, vol. 64, no. 13, pp. 3351–3361, 2016.
- [42] S. Hafezi, A. H. Moore, and P. A. Naylor, "3D acoustic source localization in the spherical harmonic domain based on optimized grid search," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, vol. 2016-May, 2016, pp. 415–419.
- [43] T. D. Abhayapala and D. B. Ward, "Theory and design of high order sound field microphones using spherical microphone array," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, vol. 2, 2002, pp. 1949–1952.
- [44] Z. Li and R. Duraiswami, "Flexible and optimal design of spherical microphone arrays for beamforming," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 15, no. 2, pp. 702–714, 2007.
- [45] H. Chen, T. D. Abhayapala, and W. Zhang, "Theory and design of compact hybrid microphone arrays on two-dimensional planes for three-dimensional soundfield analysis," *J. Acoust. Soc. Am.*, vol. 138, no. 5, pp. 3081–3092, 2015.
- [46] J. Meyer and G. Elko, "A highly scalable spherical microphone array based on an orthonormal decomposition of the soundfield," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, vol. 2, 2002, pp. 1781–1784.
- [47] P. N. Samarasinghe, H. Chen, A. Fahim, and T. D. Abhayapala, "Performance analysis of a planar microphone array for three dimensional soundfield analysis," in *IEEE Work. Appl. Signal Process. to Audio Acoust.*, 2017, pp. 249–253.
- [48] T. D. Abhayapala and A. Gupta, "Spherical harmonic analysis of wavefields using multiple circular sensor arrays," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 18, no. 6, pp. 1655–1666, 2010.
- [49] X. Xiao, S. Zhao, X. Zhong, D. L. Jones, E. S. Chng, and H. Li, "A learning-based approach to direction of arrival estimation in noisy and reverberant environments," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, 2015, pp. 2814–2818.
- [50] F. Vesperini, P. Vecchiotti, E. Principi, S. Squartini, and F. Piazza, "A neural network based algorithm for speaker localization in a multi-room environment," in *IEEE Int. Work. Mach. Learn. Signal Process.*, 2016, pp. 1–6.
- [51] Y. Sun, J. Chen, C. Yuen, and S. Rahardja, "Indoor sound source localization with probabilistic neural network," *IEEE Trans. Ind. Electron.*, vol. 65, no. 8, pp. 6403–6413, 2018.
- [52] E. L. Ferguson, S. B. Williams, and C. T. Jin, "Sound source localization in a multipath environment using convolutional neural networks," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, 2018, pp. 2386–2390.
- [53] R. Takeda and K. Komatani, "Sound source localization based on deep neural networks with directional activate function exploiting phase information," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, 2016, pp. 405–409.
- [54] S. Chakrabarty and E. A. P. Habets, "Multi-speaker DOA estimation using deep convolutional networks trained with noise signals," *IEEE J. Sel. Top. Signal Process.*, vol. 13, no. 1, pp. 8–21, 2019.
- [55] S. Adavanne, A. Politis, J. Nikunen, and T. Virtanen, "Sound event localization and detection of overlapping sources using convolutional recurrent neural networks," *IEEE J. Sel. Top. Signal Process.*, vol. 13, no. 1, pp. 34–48, 2019.
- [56] A. Fahim, P. N. Samarasinghe, and T. D. Abhayapala, "PSD estimation and source separation in a noisy reverberant environment using a spherical microphone array," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 26, no. 9, pp. 1594–1607, 2018.
- [57] E. G. Williams, *Fourier acoustics: sound radiation and nearfield acoustical holography*. London, UK: Academic press, 1999, vol. 108.
- [58] H. M. Jones, R. A. Kennedy, and T. D. Abhayapala, "On dimensionality of multipath fields: Spatial extent and richness," in *IEEE Int. Conf. Acoust. Speech, Signal Process.*, vol. 3, 2002, pp. 2837–2840.
- [59] D. B. Ward and T. D. Abhayapala, "Reproduction of a plane-wave sound field using an array of loudspeakers," *IEEE Trans. speech audio Process.*, vol. 9, no. 6, pp. 697–707, 2001.
- [60] D. Colton and R. Kress, *Inverse acoustic and electromagnetic scattering theory*, 3rd ed. Berlin, Germany: Springer-Verlag, 2012, vol. 93.
- [61] F. W. J. Olver, "[3], [6], [9] symbols," in *NIST Handb. Math. Funct.* Cambridge, UK: Cambridge University Press, 2010, ch. 34, pp. 755–766.
- [62] B. Laufer-Goldshtein, R. Talmon, and S. Gannot, "Semi-supervised source localization on multiple manifolds with distributed microphones," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 7, pp. 1477–1491, 2017.
- [63] S. Adavanne, A. Politis, and T. Virtanen, "Direction of arrival estimation for multiple sound sources using convolutional recurrent neural network," in *Eur. Signal Process. Conf.*, vol. 2018-Sept, 2018, pp. 1462–1466.
- [64] B. Loesch and B. Yang, "Source number estimation and clustering for underdetermined blind source separation," in *Int. Work. Acoust. Echo Noise Control*, no. 5, 2008, pp. 751–758.
- [65] S. Hafezi, A. H. Moore, and P. A. Naylor, "Robust source counting and acoustic DOA estimation using density-based clustering," in *IEEE Sens. Array Multichannel Signal Process. Work.*, 2018, pp. 395–399.
- [66] E. A. P. Habets, "Room impulse response generator," 2006. [Online]. Available: <https://github.com/ehabets/RIR-Generator>
- [67] J. Eaton, N. D. Gaubitch, A. H. Moore, and P. A. Naylor, "The ACE challenge - corpus description and performance evaluation," in *IEEE Work. Appl. Signal Process. to Audio Acoust.*, 2015, pp. 1–5.
- [68] J. Fliege and U. Maier, "A two-stage approach for computing cubature formulae for the sphere," in *Math. 139T, Univ. Dortmund, Fachbereich Math. Univ. Dortmund, 44221*, 1996, pp. 1–31.
- [69] J. S. Garofolo et al., "TIMIT acoustic-phonetic continuous speech corpus LDC93S1," *Linguist. data Consort.*, 1993.
- [70] F. Chollet et al., "Keras," 2015. [Online]. Available: <https://keras.io>
- [71] M. Abadi et al., "TensorFlow: Large-scale machine learning on heterogeneous distributed systems," 2015. [Online]. Available: <https://www.tensorflow.org/>
- [72] MH Acoustics, "EM32 Eigenmike microphone array release notes (v17. 0)," 25 Summit Ave, Summit, NJ 07901, USA, 2013. [Online]. Available: <https://mhacoustics.com/>



Abdullah Fahim received the B.Sc.(Hons.) degree in electrical and electronic engineering from Bangladesh University of Engineering and Technology, Dhaka, Bangladesh, in 2007. From 2007-2015, he was involved in different projects with Ericsson and Nokia SN. He is currently pursuing his Ph.D. degree in spatial audio signal processing from the Australian National University, Canberra, Australia. In 2018, he worked with the spatial audio team of Apple Inc. in California for a 6-months' internship program. His research interests include spatial audio processing techniques, virtual audio, and soundfield separation and enhancement.



Prasanga N. Samarasinghe received her Ph.D. degree from the Australian National University (ANU), Australia in 2014 and her B.E. degree (with Hons.) in electronic and electrical engineering from the University of Peradeniya, Sri Lanka in 2009. In 2019 she received a prestigious Fulbright Future Scholarship to visit the University of Maryland, USA. She is currently a Senior Lecturer at ANU, and her research interests include spatial sound recording, reproduction and analysis, room acoustics, spatial noise cancellation and virtual audio.



Thushara D. Abhayapala is a Professor of Signal Processing at the Australian National University (ANU), Canberra. He received his B.E. degree in engineering in 1994 and his Ph.D. degree in telecommunications engineering in 1999, both from the ANU. He held a number of leadership positions including Deputy Dean of the College of Engineering and Computer Science (2015-19), Head of the Research School of Engineering at ANU (2010-14) and the leader of the Wireless Signal Processing Program at the National ICT Australia (NICTA)

from 2005-07. His research interests are in the areas of spatial audio and acoustic signal processing, and multichannel signal processing. Among many contributions, he is one of the first researchers to use spherical harmonic based eigen-decomposition in microphone arrays and to propose the concept of spherical microphone arrays; novel contributions on the multi-zone soundfield reproduction problem; was one of the first to show the fundamental limits of spatial soundfield reproduction using arrays of loudspeakers and spherical harmonics. He worked in industry for two years, before his doctoral study and has active collaboration with a number of companies. He has supervised 36 PhD students and co-authored more than 280 peer-reviewed papers. He was an associate editor of IEEE/ACM Transactions on Audio, Speech, and Language Processing and was a member of the Audio and Acoustic Signal Processing Technical Committee (2011-2016) of the IEEE Signal Processing Society. He is a fellow of Engineers Australia (IEAust).