# IPDnet: A Universal Direct-Path IPD Estimation Network for Sound Source Localization

Yabo Wang, Bing Yang, and Xiaofei Li

*Abstract*—Extracting direct-path spatial feature is crucial for sound source localization in adverse acoustic environments. This paper proposes the IPDnet, a neural network that estimates direct-path inter-channel phase difference (DP-IPD) of sound sources from microphone array signals. The estimated DP-IPD can be easily translated to source location based on the known microphone array geometry. First, a full-band and narrow-band fusion network is proposed for DP-IPD estimation, in which alternating narrow-band and full-band layers are responsible for estimating the rough DP-IPD information in one frequency band and capturing the frequency correlations of DP-IPD, respectively. Second, a new multi-track DP-IPD learning target is proposed for the localization of flexible number of sound sources. Third, the IPDnet is extend to handling variable microphone arrays, once trained which is able to process arbitrary microphone arrays with different number of channels and array topology. Experiments of multiple-moving-speaker localization are conducted on both simulated and real-world data, which show that the proposed full-band and narrow-band fusion network and the proposed multi-track DP-IPD learning target together achieves excellent sound source localization performance. Moreover, the proposed variable-array model generalizes well to unseen microphone arrays. Code is available on our github page [1].

*Index Terms*—Sound source localization, direct-path IPD, full-band and narrow-band fusion network, microphone array generalization, multi-source.

## I. INTRODUCTION

SOUND source localization (SSL) aims to estimate the position of one or multiple sound sources from microphone array signals. SSL is widely used in video conferencing and human-computer interaction. The spatial cues of SSL can also be used to boost the performance of speech enhancement and source separation tasks [1], [2].

Traditional SSL methods typically rely on estimating spatial features that are associated with the direct-path signal propagation in order to establish a mapping between features and source locations. Commonly used spatial features include time delay, inter-channel phase/level difference (IPD/ILD) [3], [4], and relative transfer function (RTF) [5], [6]. Actually, the aforementioned spatial features can be straightforwardly estimated under ideal acoustic conditions (noise-free and anechoic condition). While estimating reliable features from microphone signals becomes challenging in real-world

scenarios where noise, reverberation, and the presence of multiple moving sources introduce complexity. Noise and speech overlapping introduce uncertain acoustic distortions to the microphone signals and reverberation causes an overlap-masking effect or coloration of the originally anechoic signal [7]. Meanwhile, source movement introduces time-varying characteristics of spatial cues. These all result a significant decline in localization performance.

In recent years, deep learning-based SSL approaches have been extensively developed, and show performance superiority over conventional methods [8]–[20], particularly in challenging environments. The superiority stems from the capacity of neural networks to learn complex patterns and subtle differences in acoustic signals. This work proposes a new SSL network for multiple-moving-source localization in the presence of noise and reverberation, which is an extended version of our previous conference paper, i.e. FN-SSL [20]. FN-SSL is a full-band and narrow-band fusion network proposed for single-source localization with two microphones, which achieved excellent SSL performance due to its efficient network architecture. In this work, we extend FN-SSL for multi-microphone and multi-source localization, and propose a new learning target. Moreover, based on the extended network, a new variable-array model is proposed which can be applied to variable microphone arrays with different array topology.

Specifically, the contributions of this work are as follows:

*1) Full-band and narrow-band fusion SSL network.:* The proposed network is motivated by the recently proposed speech enhancement networks [21]–[23], in which full-band layers and narrow-band layers are cascaded for predicting the clean speech signal, which shows a large performance superiority for speech enhancement, and is now becoming a new research trend. The proposed network takes as input the (short-time Fourier transform, STFT of) multichannel microphone signals, and predicts the direct-path IPD (DP-IPD) as localization feature. Compared to processing the noisy IPD [24] or the noisy spatial spectrum [12], [13], the microphone signals preserve the natural properties of noise and reverberation, such as the spatial-diffuseness of noise and late reverberation, and it is more effective to leverage these properties to remove them. In the proposed network, full-band/narrow-band layers process the time frames/frequencies independently, and all the time frames/frequencies share the same network weights. In narrow-band, there are rich information for extracting localization features, which are largely leveraged in conventional methods. For example, localization features are extracted by narrow-band channel identification in [25], by coherence test in [26], and by direct-path dominance

Yabo Wang is with Zhejiang University and also with Westlake University, Hangzhou, China, e-mail: wangyabo@westlake.edu.cn. Bing Yang and Xiaofei Li are with the School of Engineering, Westlake University and also with the Institute of Advanced Technology, Westlake Institute for Advanced Study, Hangzhou, China (e-mail: yangbing@westlake.edu.cn; lixiaofei@westlake.edu.cn). Yabo Wang and Bing Yang: equal contribution; Xiaofei Li: corresponding author.

[1] https://github.com/Audio-WestlakeU/FN-SSL

test in [27]. The narrow-band layers processes the along-time sequences to focus on learning these narrow-band information. The full-band layers processes the along-frequency sequence to focus on learning the full-band correlation of spatial cues, such as the linear relation of DP-IPD to frequency.

*2) Mutli-track DP-IPD learning target.:* Popular SSL learning targets include location classification [10], location regression [28] and spatial spectrum regression [9]. In this work, we propose to use multi-track DP-IPD as the learning target for multi-source localization. DP-IPD stands for the IPD of the direct-path propagation, which is theoretically related to the microphone array geometry and source location, and thus it can be easily translated to source location with known array geometry. In this work, we estimate the source location by simply matching the estimated DP-IPD with theoretical DP-IPD of candidate locations. Multi-track DP-IPD means we let the network output/estimate the DP-IPD of multiple sources simultaneously. DP-IPD is a signal-level localization feature, which can be estimated from microphone signals based on only signal-level information. By contrast, other targets are array-dependent, and require one further step of conversion from localization feature to target, which may arise more difficulty. This is verified by our experiments conducted in Section V-E.

DP-IPD is well defined for active sound source, but not for non-source frames. One straightforward representation for non-source is an all-zero vector. However, learning the output space combined by DP-IPDs and the all-zero vector may be not easy. This work proposes taking the DP-IPD mean point of the whole localization space as the representation for non-source, with which it is more easier to switch between source and non-source frames, as shown in Fig. 7.

*3) Variable Array SSL:* Most of existing SSL networks are array-dependent, namely training and test using the same array. For a new array, the network needs to be retrained, which is time-consuming. Especially, when using real data for training, collecting a large amount of annotated data for a new array is even much more time consuming. In this paper, a variable-array SSL model is proposed, once trained which can be directly used for any unseen microphone array. Specifically, microphones are processed pair-wisely, and the mean pooling of pair-wise hidden units is used for the communication between microphone pairs. This pair-wise processing plus mean pooling scheme can handle variable number of microphones, and it is motivated by the variable-array speech enhancement networks [29]–[32]. The network outputs the DP-IPD estimation for varying number of microphone pairs. As DP-IPD estimation is a signal-level task, one network can easily handle the task for different arrays. At the test stage, the estimated DP-IPDs can be used for source localization with the known array geometry.

Experiments have been conducted on both simulated and real dataset, which demonstrate that the proposed learning target outperforms all comparison targets, and the proposed method as a whole outperforms recently proposed baseline methods by a large margin. Moreover, the proposed variable-array model can generalize well to unseen simulated and real microphone arrays.

The rest of this paper is organized as follows. Section II presents an overview of related works in the literature. Section III defines the problem of multiple moving source localization. Section IV details the proposed method. Section V gives the experimental results and discussions. Finally, Section VI concludes the paper and suggesting directions for future research.

## II. RELATED WORKS

### A. Deep Learning Based Sound Source Localization

In recent years, significant research progress has been made in the field of sound source localization using neural networks [8]–[20]. Table I provides a chronological overview of some representative sound source localization methods. Although these methods normally achieve promising SSL performance, they can only handle certain limited tasks, in terms of the number of sources, frame-wise or block-wise SSL, fixed array or variable array. Where multi-source localization sometimes requires especially designed output format, frame-wise methods output SSL result for each time frame and are suitable for online and moving source localization, variable-array models can be applied to unseen microphone arrays.

Various network architectures have been adopted for SSL, among which convolutional neural networks (CNN) [9], [13], [14], [19] and convolutional recurrent neural Networks [11], [16], [18] (CRNN) are the most commonly used networks. These networks are all designed to process all the frequencies together. The network input can be in the signal level, such as the time-domain signal [33], the STFT coefficients [14], [20] or the magnitude and phase of STFT coefficients [9], [11], [16], or in the feature level, such as IPD, IID, the generalized cross-correlation (GCC) function [34]–[36] and noisy spatial spectrum [12], [13], [19].

According to the learning target, SSL methods are classified as feature/location regression or location classification methods. Feature/location regression methods estimate the localization feature (such as DP-RTF, DP-IPD and inter-channel time difference (ITD)) [11], [16], [20], [36]–[39] or directly estimate source location [9], [12]–[14], [19] from the noisy signal or noisy localization features. Most works output the feature/location for one source, and few works study how to extend feature/location regression to multiple sources. Location classification methods [10], [15], [17], [18], [40]–[42] take candidate locations as classes, and multi-source localization can be easily conducted as a multi-class classification task [10], [40]–[42]. Due to the non-orthogonal relationship between adjacent locations, the classification output often exhibits a blurred response around the main peak which degrades the localization performance.

Based on Table I and the above overview, it is clear that the proposed method is totally different from existing works in both network architecture and learning target.

### B. Variable Array SSL

Mapping from microphone signals and/or localization features to source location is intrinsically an array-dependent problem, which requires to know the geometry of microphone

TABLE I: Brief overview of deep-learning-based sound source localization methods.

| Method | Year | Input Feature | Target | Multi-source (vs. Single-source) | Frame-wise (vs. Chunk-wise) | Variable-array (vs. Fixed-array ) |
|---|---|---|---|---|---|---|
| [9] | 2021 | Mag + Phase | Spatial spectrum regression | ✓ | ✗ | ✗ |
| [10] | 2021 | Intensity vector | Multi-class location classification | ✓ | ✓ | ✗ |
| [11] | 2021 | Mag + Phase | DP-RTF regression | ✗ | ✗ | ✓ (2-channel) |
| [12] | 2021 | SRP-PHAT Spectrogram | Location regression | ✗ | ✓ | ✗ |
| [13] | 2022 | SRP-PHAT Spectrogram | Location regression | ✗ | ✓ | ✗ |
| [14] | 2022 | STFT Coefficients | Spatial spectrum regression | ✓ | ✗ | ✗ |
| [15] | 2022 | Mag + IPD | Multi-track spatial spectrum regression | ✓ | ✗ | ✗ |
| [16] | 2022 | Mag + Phase | Mixed DP-IPD regression | ✓ | ✓ | ✗ |
| [17] | 2022 | GCC-PHAT + Array Geometry | Location classification | ✗ | ✗ | ✓ (constant-channel) |
| [18] | 2023 | MFCC and Mel features | Multi-class location classification | ✓ | ✓ | ✗ |
| [19] | 2023 | SRP-PHAT Spectrogram | Location regression | ✗ | ✓ | ✗ |
| [20] | 2023 | STFT Coefficients | DP-IPD regression | ✗ | ✓ | ✗ |
| Proposed | - | STFT Coefficients | Multi-track DP-IPD regression | ✓ | ✓ | ✓ |

array or uses a fixed microphone array. In [17], by also taking as input the microphone array geometry along the localization feature to the network, the network can perform SSL for variable arrays. However, limited by the fixed input size, one network can only process variable arrays with the same number of microphones. In our previous works [11], [16], the clean localization feature, i.e. DP-RTF or DP-IPD, is taken as the network output, and 2-channel array is processed, for which case one network can be directly trained with variable arrays and test on unseen arrays. In this work, the proposed variable-array model can handle any microphone array without the limit of number of microphones.

## III. PROBLEM FORMULATION

Assuming there are multiple sound sources in a closed environment with noise and reverberation. The multichannel signals recorded by a microphone array are denoted as

$$x_m(t) = \sum_{k=1}^{K} a_m(t, \theta_k) * s_k(t) + v_m(t), \quad (1)$$

where $m \in [1, M]$, $k \in [1, K]$ and $t \in [1, T]$ represent the indices of microphones, sound sources and time samples, respectively. As for the $k$-th source, $s_k(t)$, $\theta_k$, and $a_m(t, \theta_k)$ represent the source signal, the direction of arrival (DOA), and the direct-path response (within the room impulse response, RIR) to the $m$-th microphone, respectively, and * denotes convolution. The noise signal $v_m(t)$ includes both ambient noise and the reflections/reverberation of sources.

Applying the short-time Fourier transform (STFT), the multichannel signals are expressed as

$$X_m(n, f) = \sum_{k=1}^{K} A_m(f, \theta_k) S_k(n, f) + V_m(n, f), \quad (2)$$

where $n \in [1, N]$, $f \in [1, F]$ represent the time frame index and frequency index, respectively. Here $X_m(n, f)$, $S_k(n, f)$ and $V_m(n, f)$ are the STFT coefficients of microphone, source

and noise signals, respectively. $A_m(f, \theta_k)$ is the transfer function (Fourier transform) of the direct-path response. The direct-path relative transfer function (DP-RTF) of two microphones encodes the direct-path IPD and ILD within its phase and amplitude, respectively, and it is thus a reliable localization feature

$$B_m(f, \theta_k) = A_m(f, \theta_k)/A_r(f, \theta_k), \quad (3)$$

where $r$ is the index of one selected reference channel. For simplicity, only the DP-IPD, i.e. the phase part $\angle B_m(f, \theta_k)$, is employed and learned for SSL in this work.

In the free and far field, for one given microphone pair and source DOA $\theta$ (relative to the microphone pair), the complex-valued DP-IPD can be theoretically computed as

$$\tilde{B}(f, \theta) = e^{-j2\pi v_f d\cos(\theta)/c} \quad (4)$$

where $v_f$ is the frequency in Hz, $d$ is the microphone distance, $c$ is sound speed in air, and $d\cos(\theta)/c$ is the time difference of arrival (TDOA) from the direction of $\theta$ to the two microphones.

In this work, sound source localization amounts to using a neural network to estimate the DP-IPDs from the multichannel microphone signals. As for $M$ microphones, one of microphone is selected as the reference channel, and the network predicts the DP-IPDs of the $M - 1$ microphone pairs (other channels relative to the reference channel). Then, the DOA estimation can be obtained by simply matching the predicted DP-IPDs with the DP-IPD templates (namely the theoretical DP-IPDs of a set of pre-defined candidate directions), as shown in Fig. 1. Specifically, the inner product between the predicted DP-IPD vector with the theoretical DP-IPD vector are computed, and the candidate direction with the maximum product value is taken as the DOA estimation. The training targets of DP-IPD and the DP-IPD templates are all computed using Eq. (4).

Moreover, the proposed DP-IPD estimation network is designed to handle more realistic and complex applications in the following aspects:
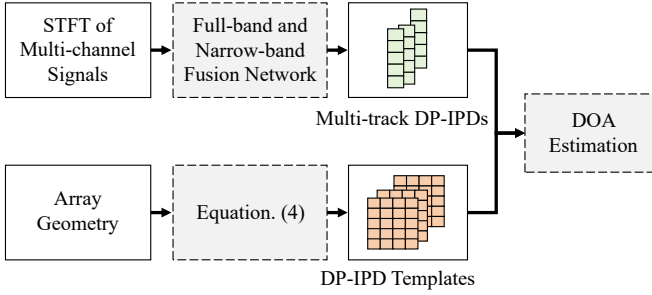
Fig. 1: Block diagram of the proposed method.

- Flexible number of sound sources. We consider the number of source sources is flexible and time varying. However, we rely on a strong and reasonable assumption that at most $K$ (e.g. 2 in this work) sources present at one time, based on which a fixed number of $K$ tracks of DP-IPD are predicted. The DP-IPD estimations could belong to different sources at different times in one track. A trivial DP-IPD value is set for non-source. Overall, the network can detect and localize 0 to $K$ sources at one time.
- Moving sound source. For moving source, the DOA, i.e. $\theta_k$, and also its transfer function $A_m(f, \theta_k)$ are time-dependent/varying. To address this, we frame-wisely predict the DP-IPD and localize the sources, in either an online (causal) or an offline (non-causal) way.
- Variable microphone arrays. The network can be designed to work for variable arrays with different topology and number of channels. One network is trained using many different arrays, and it predicts the DP-IPDs for the varying $M-1$ microphone pairs. Then, the network can be used for DP-IPD estimation with an arbitrary test array, for which the DP-IPD templates can be theoretically computed during test. This way disentangles the DP-IPD estimation step and the localization step (namely template matching), and thus forms an universal sound source localization network. This is reasonable considering the facts that the DP-IPD estimation step could be an array-independent signal-level task, and the localization step is array-dependent but theoretically simple. The accuracy of DP-IPD estimation is critical for localization, as perfect DP-IPD estimation leads to nearly perfect localization.

The number of microphones and sources, i.e. $M$ and $K$, and the DOA $\theta_k$ could all be varied either for different settings or along time, but we will not specify the variation in the following for notational simplicity, unless otherwise stated.

## IV. METHOD

This section presents the proposed IPDNet. Two versions are proposed for a fixed microphone array and variable arrays, whose network architectures are shown in Fig. 3 (a) and (b), respectively.

### A. Learning Target

The proposed network takes as input the STFT of multi-channel recordings, and outputs/predicts the DP-IPD features.

To enable the optimization of real-value networks, following the setting in our prevision works [11], [16], [20], the learning target is set as the real ($\mathcal{R}$) and imaginary ($\mathcal{I}$) parts of complex-valued DP-IPD (concatenated along frequencies, for one microphone pair) as

$$\mathbf{q}(\theta) = [\mathcal{R}\{\tilde{B}(1, \theta)\}, \ldots, \mathcal{R}\{\tilde{B}(F, \theta)\}, \\ \mathcal{I}\{\tilde{B}(1, \theta)\}, \ldots, \mathcal{I}\{\tilde{B}(F, \theta)\}]^\top \in \mathbb{R}^{2F}, \quad (5)$$

where $^\top$ denotes vector transpose. The output activation layer is set as *tanh* to predict DP-IPD.

$\mathbf{q}(\theta)$ defines the DP-IPD target vector for one sound source presents at DOA $\theta$, and across all DOAs it forms an DP-IPD manifold. However, it is not straightforward to define the target for non-source frames. In our previous works, an all-zero vector is used as the target for non-source frames, which however seems not meaningfully correlated to the DP-IPD vector. The network needs to learn the output space expanded by the DP-IPD manifold and the all-zero vector, and rapidly switch between the manifold (for source frames) and the all-zero vector (for non-source frames), which is possibly not easy. In this work, to facilitate the network learning, we propose to define the target for non-source frames as the mean point of the (complex-valued) DP-IPD manifold, which can be derived as

$$\bar{q}(f) = \frac{1}{2\pi} \int_0^{2\pi} \tilde{B}(f, \theta) d\theta \\ = \frac{1}{2\pi} \int_0^{2\pi} e^{-j2\pi v_f d \cos(\theta)/c} d\theta \quad (6) \\ = J_0(2\pi v_f d/c),$$

where $J_0(\cdot)$ denotes the zero-order Bessel function of the first kind. DP-IPDs are integrated/averaged over all possible $\theta \in [0, 2\pi)$ for one microphone pair, which is analogous to computing the spatial coherence of cylindrically isotropic diffuse sound field [43]. The non-source target values are real numbers as a function of frequency and microphone distance. For one given microphone distance, the non-source target values are computed for the $F$ discrete frequencies using Eq. (6) and then concatenated to form the target vector. Fig. 2 shows the non-source target value as a function of frequency for two different microphone distances.

To determine whether one source is active or not in one frame, namely conducting the frame-wise per-source activity detection, we compute the following direct-path to noisy magnitude ratio at the reference channel:

$$v_k(n) = \frac{1}{F} \sum_{f=1}^F \frac{|A_r(f, \theta_k) S_k(n, f)|}{|X_r(n, f)|}, \quad (7)$$

where $|\cdot|$ denotes absolute value. If $v_k(n)$ is larger than a pre-set threshold, we consider the $k$-th source is active in frame $n$, otherwise inactive. For active source, we use the DP-IPD target vector, otherwise we use the non-source target vector.

### B. Full-Narrow Network Block

The (two versions of) proposed network consists of a cascade of full-narrow network blocks. One full-band layer plus one narrow-band layer make up the full-narrow block.
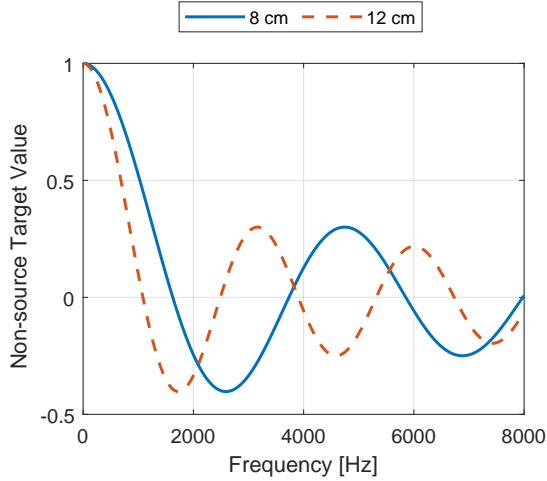
Fig. 2: Examples of non-source target for two different microphone distances.

*1) Full-band BLSTM layer:* The full-band BLSTM layers process the time frames independently, and all the time frames share the same network weights. The input is a sequence along the frequency axis of one single time frame:

$$H^{\text{full}}(n) = (\mathbf{h}(n,1),\dots,\mathbf{h}(n,f),\dots,\mathbf{h}(n,F)), \quad (8)$$

where the superscript $^{\text{full}}$ indicates full-band layer, $\mathbf{h}(n,k) \in \mathbb{R}^D$ represents the hidden vector for one time-frequency (TF) bin, and $D$ is the hidden dimension. Note that $\mathbf{h}(n,k)$ is the microphone signals for the first full-band layer as shown in Eq. (10), while it is the output of the previous layer for other layers. The full-band layers focus on learning the inter-frequency dependencies of spatial/localization cues. DP-IPD of different frequencies has a strong correlation, as they are all derived from the same TDOA. In addition, the spatial/localization cues of those frequencies with low direct-path energy can be well predicted with the help of other frequencies. The full-band layers do not learn any temporal information, which is left for the following narrow-band layers.

*2) Narrow-band (B)LSTM layer:* The narrow-band (B)LSTM layers process the frequencies independently, and all the frequencies share the same network weights. The input is a sequence along the time axis of one single frequency:

$$H^{\text{narrow}}(f) = (\mathbf{h}(1,f),\dots,\mathbf{h}(n,f),\dots,\mathbf{h}(N,f)), \quad (9)$$

where the superscript $^{\text{narrow}}$ indicates narrow-band layer. The input vector $\mathbf{h}(n,f)$ is the output vector of the previous full-band layer. Estimating the direct-path localization features in narrow-band has been studied in many conventional methods, such as by channel identification in [25], coherence test in [26], direct-path dominance test in [27]. The proposed narrow-band layers focus on exploiting these narrow-band inter-channel information. In addition, DP-IPD is time-varying for moving sound source, and the narrow-band layers learn the temporal evolution of DP-IPD as well.

*3) Skip connections:* The full-band and narrow-band layers are tailored to emphasize their specific information domains.

However, there's a risk of losing narrow-band information after processing through a full-band layer, and vice versa. To circumvent this, skip connections are incorporated to prevent such information loss. As illustrated in Fig. 3, this entails concatenating the input sequence of each full-band layer and narrow-band layer with the original input signal's STFT coefficients (after proper dimension transformation).

### C. Fixed-Array Model

For the scenario with a fixed microphone array for training and test, the model architecture is shown in the Fig. 3 (a). Without loss of generality, the reference channel is set as $r = 1$. The model takes as input the $M$-channel microphone signals, and outputs $M - 1$ DP-IPD vectors. Specifically, the input is formed by concatenating the real and imaginary parts of multichannel microphone signals as:

$$\mathbf{x}[n,f,:] = [\mathcal{R}(X_1(n,f)), \mathcal{I}(X_1(n,f)),\dots, \\ \mathcal{R}(X_M(n,f)), \mathcal{I}(X_M(n,f))] \in \mathbb{R}^{2M}, \quad (10)$$

where $[:]$ is an operator to take all values of one dimension of a tensor. The input is first processed by full-narrow blocks to extract the spatial feature embedding of multiple sources which can be formulated as:

$$\mathbf{h}^{\text{fn}} = FNBlocks(\mathbf{x}) \in \mathbb{R}^{N \times F \times D}. \quad (11)$$

Then, a convolutional block is used to separate and extract the DP-IPD vector of multiple sources. The structure of the convolutional block is also shown in Fig. 3 (a) which consists of three 2-dimensional causal convolutional layers and two temporal average pooling layers. Convolutional layers are used to capture the local features to separate the microphone pairs and sources. The average pooling layers are used to compress the frame rate. The activation function of the first two convolutional layers is $relu$, and the activation function of the last convolutional layer is set as $tanh$ for DP-IPD estimation. The final output is obtained as:

$$\mathbf{Q} = ConvBlock(\mathbf{h}^{\text{fn}}) \in \mathbb{R}^{N \times F \times O}, \quad (12)$$

where $O = 2$ (real and imaginary parts) $\times (M - 1)$ (microphone pairs) $\times K$ (sources).

The network output is defined as a fixed number of $K$ tracks, each track represents one source (can be non-source) and contains $M - 1$ estimated DP-IPD vectors. For one time frame, say $n$, the output can be written as

$$\mathbf{Q}[n,:] = [\hat{\mathbf{q}}_{1,2}(\theta_1),\dots,\hat{\mathbf{q}}_{1,M}(\theta_1)), \\ \dots,\hat{\mathbf{q}}_{1,2}(\theta_K),\dots,\hat{\mathbf{q}}_{1,M}(\theta_K)], \quad (13)$$

where $\hat{\mathbf{q}}_{1,m}(\theta_k)$ denotes the estimated DP-IPD for the $1,m$ microphone pair and the $k$-th source, at time frame $n$. Regarding the source permutation problem, the frame-level permutation invariant training (frame-level PIT) is used. The order of microphone pair is fixed according to the microphone index.
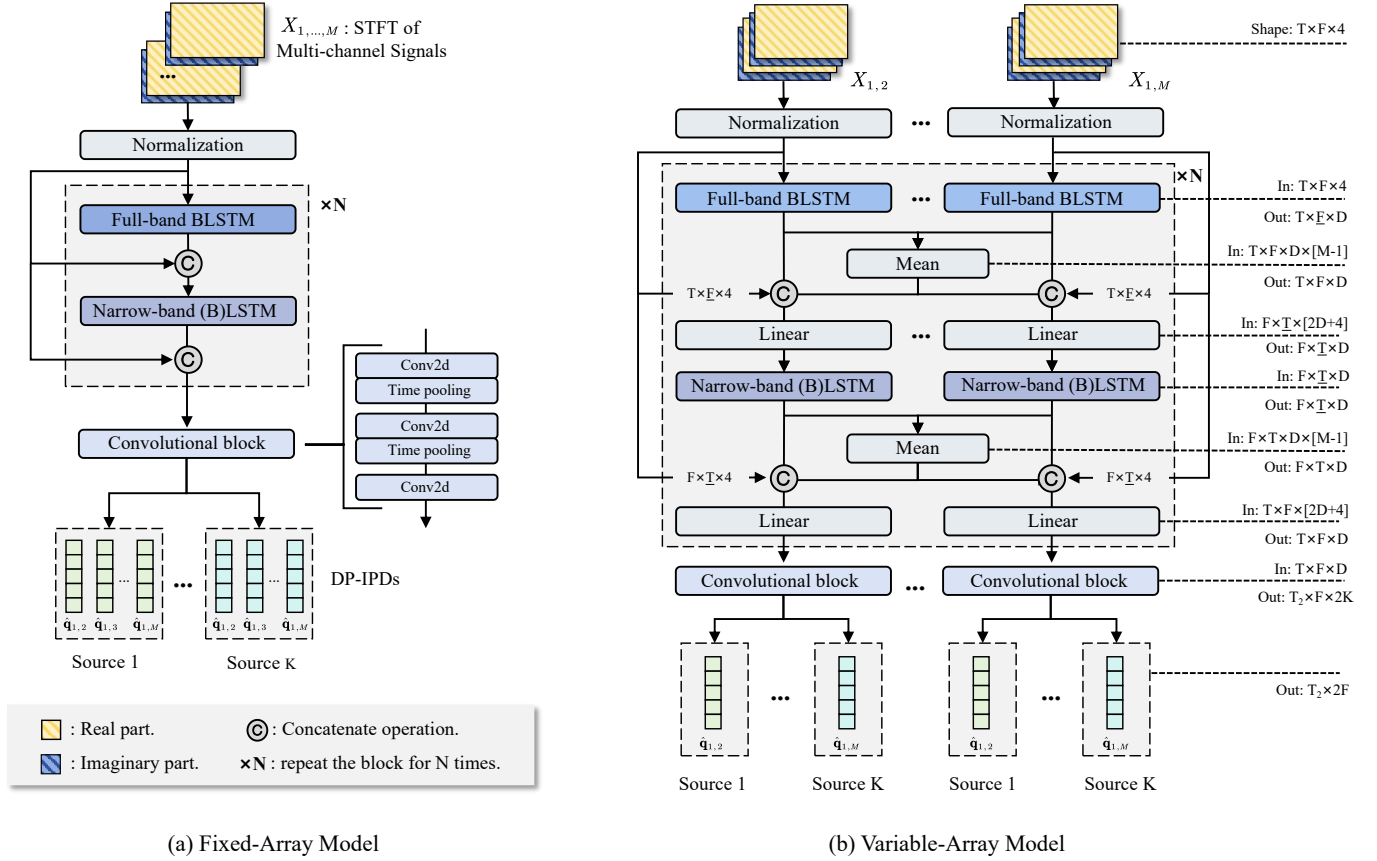
Fig. 3: Network architecture for the proposed fix-array and variable-array models. The data organization is in the format of: number of sequences × sequence length × feature dimension.

## D. Variable-Array Model

The proposed IPDNet is an end-to-end DP-IPD estimation network that leverages the network to fully learn useful information from the multi-channel microphone signals for removing the interference of such as noise, reverberation and overlapping speech. The network can be easily designed when the microphone array is fixed, namely the input data has a fixed size. However, when processing different arrays using one network with variable input sizes, implementing the end-to-end DP-IPD estimation network is not trivial. This work proposes a pair-wise processing scheme, namely microphone pairs (pair the reference channel with other channels) are individually processed by weight-shared networks, which makes it feasible for processing arbitrary number of microphones (pairs). The intermediate hidden units of all microphone pairs are mean pooled and broadcast to every pair, which makes microphone pairs communicate to each other. Although such microphone pair communication scheme is not as efficient as the one performed in the fixed-array model that directly learns the dependency of multiple microphones, it largely improves the DP-IPD estimation accuracy compared to when there is no communication between microphone pairs.

The network architecture of the variable-array model is shown in Fig. 3 (b). Specifically, the reference channel is set as $r = 1$, and other microphones are individually paired with

the reference channel to form the network input as:

$$
\begin{aligned}
\mathbf{x}_{1,m}[n, f, :] = [\mathcal{R}\left(X_1(n, f)\right), \mathcal{I}\left(X_1(n, f)\right), \\
\mathcal{R}\left(X_m(n, f)\right), \mathcal{I}\left(X_m(n, f)\right)] \in \mathbb{R}^4,
\end{aligned}
\tag{14}
$$

which is then processed by weight-shared full-narrow blocks:

$$
\mathbf{h}_{1,m}^{\text{fn}} = FNBlocks(\mathbf{x}_{1,m}) \in \mathbb{R}^{N \times F \times D}.
\tag{15}
$$

Let $\mathbf{h}_{1,m}^{\text{f/n}} \in \mathbb{R}^{N \times F \times D}$ denotes the hidden units of any one of full-band BLSTM or narrow-band LSTM for the $1, m$ microphone pair. An average operation after each full-band BLSTM and narrow-band (B)LSTM is conducted as:

$$
\mathbf{c}^{\text{f/n}}[n, f, :] = \frac{1}{M-1} \sum_{m=2}^{M} \mathbf{h}_{1,m}^{\text{f/n}}[n, f, :],
\tag{16}
$$

which contains the intermediate information extracted from all microphone pairs. As the input to the next layer, $\mathbf{c}^{\text{f/n}}$ is first concatenated onto each $\mathbf{h}_{1,m}^{\text{f/n}}$, and then transformed back to $D$-dimensional from $2D$-dimensional via a Linear layer. This interaction between microphone pairs provides the dependency across the entire array, and thus helps each microphone pair learn better.

Finally, the same convolutional block as in the fixed-array model is used to separate the DP-IPD of multiple sources from the output of full-narrow blocks:

$$
\mathbf{Q}_{1,m} = ConvBlock(\mathbf{h}_{1,m}^{\text{fn}}) \in \mathbb{R}^{F \times T \times O'},
\tag{17}
$$

and now $O' = 2$ (real and imaginary parts) $\times K$ (sources).

Similarly to the fixed-array model, for each of the $M-1$ microphone pairs, the network output is defined as a fixed number of $K$ tracks, and each track represents one source (can be non-source). For one time frame $n$, the output can be written as

$$\mathbf{Q}_{1,m}[n,:] = [\hat{\mathbf{q}}_{1,m}(\theta_1), \ldots, \hat{\mathbf{q}}_{1,m}(\theta_K)]. \qquad (18)$$

The frame-level PIT is also used for training. The order of microphone pair now can be arbitrarily set.

### E. Frame-level PIT and Sound Source Localization

The source permutation problem in the training stage is solved through frame-level permutation invariant training (PIT). Let $\alpha \in P$ denote one of the all possible source orders. For one time frame $n$, the frame-level PIT loss can be computed as:

$$\mathcal{L}(n) = \min_{\alpha \in P} \frac{1}{K(M-1)} \sum_{k=1}^{K} \sum_{m=2}^{M} \text{MSE}\left(\mathbf{q}_{1,m}(\theta_{\alpha(k)}), \hat{\mathbf{q}}_{1,m}(\theta_k)\right)$$
$$(19)$$

Note that all the $M-1$ microphone pairs share the same source order. The frame-wise PIT loss is averaged over frames as the overall training loss. The mean squared error (MSE) is used as the loss function.
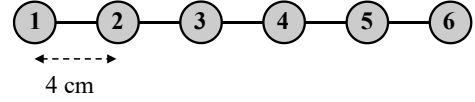
At inference, each track of DP-IPD estimation is associated to one source (or non-source). Computing the inner product of the estimated DP-IPDs with the DP-IPD templates, we can derive one source's spatial spectrum:

$$s_k(i) = \frac{1}{M-1} \sum_{m=2}^{M} \hat{\mathbf{q}}_{1,m}(\theta_k)^\top \mathbf{q}_{1,m}(\theta_i), \qquad (20)$$
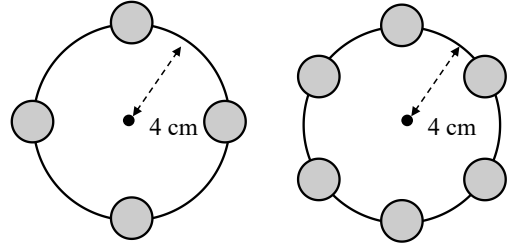
where $\theta_i$ denotes the $i$-th ($i \in [1, I]$) candidate location. The spatial spectrum $s_k(i)$ is independent for each source, as a function of candidate locations. If the maximum value of $s_k(i)$ exceeds a pre-defined threshold, a source is deemed present at the corresponding candidate location, otherwise, it is considered as non-source.

### F. Configurations

The proposed network can be easily implemented for both offline or online SSL, by setting the narrow-band LSTMs to be bidirectional or unidirectional and the convolutional layers to be non-causal or causal, respectively. To make the model easier to optimize, Laplace normalization is performed on the network input as $X_m(n,f)/\mu(n)$, where $\mu(n)$ is a normalization factor. For offline SSL, $\mu(n)$ is computed as $\frac{1}{NF} \sum_{n=1}^{N} \sum_{f=1}^{F} |X_m(n,f)|$. For online SSL, $\mu(n)$ is recursively calculated as $\mu(n) = \beta\mu(n-1) + (1-\beta)\frac{1}{F}\sum_{f=1}^{F}|X_m(n,f)|$ [22] to ensure the causality of the method. Here, $\beta = (L-1)/(L+1)$ denotes the smoothing weight of the historical time frames, and $L$ represents the length of smoothing window.



(a) Linear array (2-/4-/6-channel)



(b) Circle array (4-/6-channel)

Fig. 4: Test microphone arrays.

## V. EXPERIMENTS AND DISCUSSIONS

In this section, the performance of the proposed method on both simulated and real-world data are presented. We first describe the experimental setup and then give the experimental results with detailed discussions.

### A. Datasets

We train the model on simulated datasets, and evaluated it on both simulated and real-world data.

**Simulated dataset:** Multichannel microphone signals are simulated through convolving RIRs with speech source signals. RIRs (of moving sources) are generated using the gpuRIR toolbox [44]. Clean speech signals are randomly selected from the training, dev and test sets of the LirbriSpeech corpus [45]. Single-source microphone signals are added to obtain multi-source microphone signals. The number of sound sources is set to 2. The room reverberation time (RT60) is randomly set in the range of [0.2, 1.3] s. The room size is randomly set in the range from 6×6×2.5 m to 10×8×6 m. Diffuse noise signals generated following [43] are added to speech signals with a signal-to-noise ratio (SNR) randomly selected from -5 dB to 15 dB. For the variable-array model, the microphone arrays used for training are randomly generated. We pre-defined 6 categories of microphone array, namely uniform/non-uniform linear microphone array, circular microphone array, circular with center microphone array, 2D/3D ad-hoc microphone array. Each category of microphone array is equally presented in proportion. The number of microphones is randomly set in the range of [2, 8]. The distance of any two microphones is limited to [3, 25] cm.

Five commonly-used microphone array topology are set for test, including 2-/4-/6-channel linear arrays and 4-/6-channel circular arrays, which are shown in Fig. 4. For the fixed-array model, one network is trained for each test array. The number of utterances (of both the fixed-array and variable-array models) for training, validation and test are 300,000,

TABLE II: Results of ablation studies, FB and NB represent full-band and narrow-band, respectively.

| Method | #Params. [M] | FLOPs [G/s] | Tolerance: 5° | | | Tolerance: 10° | | |
|---|---|---|---|---|---|---|---|---|
| | | | MDR [%] | FAR [%] | MAE [°] | MDR [%] | FAR [%] | MAE [°] |
| IPDnet (prop.) | 0.7 | 19.4 | **17.1** | **16.8** | **1.6** | **8.3** | **7.7** | **2.1** |
| w/o FB BLSTM | 0.5 | 12.7 | 64.4 | 33.5 | 2.1 | 51.9 | 21.0 | 3.6 |
| w/o NB LSTM | 0.4 | 10.6 | 19.3 | 21.4 | 1.7 | 9.5 | 11.5 | 2.3 |
| w/o Conv. block | 0.4 | 12.8 | 20.3 | 17.5 | 1.7 | 10.5 | 7.8 | 2.2 |

4,000, and 4,000, respectively. Each utterance includes two sources, the length of each source's audio clip is in the range of [5, 25] s, and the two audio clips are overlapped according to an overlap rate in the range of [0, 1]. The moving speed of speakers is in the range of [0, 1] m/s.

**Real-world dataset:** The LOCATA [46] dataset provides audio signals recorded in the computing laboratory of the Department of Computer Science at the Humboldt University Berlin. The room size is 7.1×9.8×3 m and the reverberation time is 0.55 s. This dataset includes tasks for localizing single/multiple static sound sources (task 1 and 2) and single/multiple moving sound sources (task 3, 4, 5, and 6) using four different microphone arrays. We evaluate on all the 6 tasks with the commonly used benchmark2 (12-mic pseudo-spherical array) and DICIT (15-mic planar array) [47] arrays. Note that we only consider the LOCATA utterances with no more than 2 speakers.

### B. Configurations

The sampling rate of microphone signals is 16 kHz. The window length of STFT is 512 samples (32 ms) with a frame shift of 256 samples (16 ms). The length of training audio clips is 4.5 s. The modeld output a localization result every 12 frames (192 ms). The maximum number of sources is set to $K = 2$. The threshold for source activity detection is set to 0.001. The threshold for the estimated spatial spectrum to determine the presence of a source is set to 0.5.

The number of cascaded Full-Narrow blocks is set to 2. The hidden dimension of the network has been well searched. For the variable-array model, the hidden dimension of every (B)LSTM layers are all set to 128. For the 2-channel and 4-/6-channel fixed-array models, the hidden dimension of every (B)LSTM layers are all set to 128 and 256, respectively, where a higher number of channels carries more information and thus requires a larger network. The Adam [48] is used as the optimizer for training. The batch size of the fixed-array model and variable-array model are set to 16 and 4, respectively. The learning rate is initially set to 0.0005, and exponentially decays with a decaying factor of 0.975. We train the model for almost 30 epochs for the fixed-array model and 15 epochs for the variable-array model.

### C. Evaluation Metrics

The performance is only evaluated on voice-active periods. The resolution of candidate azimuths and elevations are both
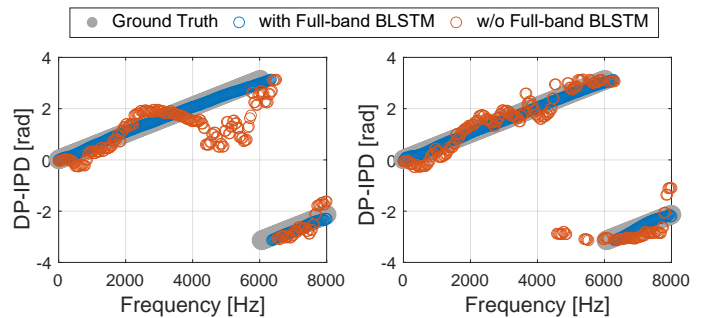


Fig. 5: DP-IPD estimations for the proposed model with or without full-band layers. The acoustic condition for left figure is: RT60 = 0.6 s, SNR = 0 dB (white noise) and for right figure is: RT60 = 0.6 s, SNR = 10 dB (white noise).

1°. The angle estimation error is computed as the difference of estimated and true angles. *tolerance: n* means that the source is considered to be successfully localized if the azimuth estimation error is not larger than n°. Evaluation metrics include miss detection rate (MDR), false alarm rate (FAR) and mean absolute error (MAE). MDR and FAR represent the frame proportions of source active but not successfully localized and source detected but not active, respectively. MAE represents the absolute angle estimation error of all successfully localized sources and time frames.

### D. Ablation study

To analyze the contribution of various components of the proposed network, ablation experiments are conducted with the fixed-array model for the simulated 2-channel array, and 180°-azimuth localization is performed. Table II presents the localization performance, model size, and the number of floating point operations (FLOPs) [2] for the proposed network and its variants achieved by removing one sub-network. It is shown that SSL performance noticeably degrades when anyone of the three sub-networks is removed, which indicates the important contribution of the sub-networks to the overall performance. Especially, the full-band BLSTM module seems playing an extremely important role. The full-band BLSTM learns the cross-frequency dependency of localization information. As shown in Eq. (4), DP-IPD is basically a linear function of frequency, and leveraging such a strong cross-frequency relationship of DP-IPD would be very helpful for improving the estimation accuracy. Fig. 5 shows two examples of DP-IPD estimation. It can be seen that, when SNR is 0 dB, the DP-IPD estimates are highly biased for those frequencies possibly have very low SNR, and the full-band BLSTM helps to correct them based on the cross-frequency relationship of DP-IPD.

### E. Comparison with Different Learning Targets

Deep learning based sound source localization methods use various different learning targets, among which location

---

[2]The FLOPs, expressed in Giga per second (G/s), are calculated based on 4.5 s long utterances, and then divided by 4.5. For the computation of FLOPs, we utilize the official tool available in PyTorch (torch.utils.flop_counter.FlopCounterMode on the meta device).
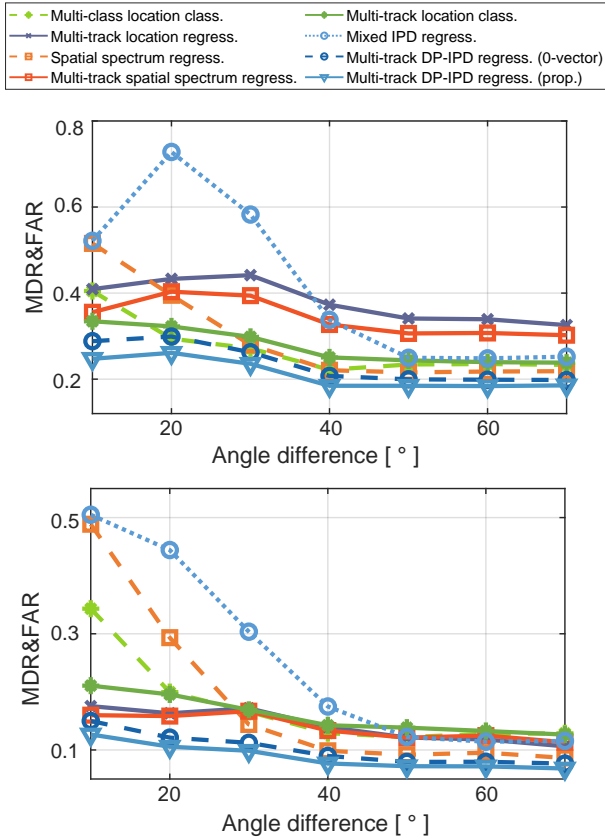
Fig. 6: SSL performance for different learning targets, as a function of the angle difference of two sources. The error tolerance is (up) 5°, and (bottom) 10°.

classification and regression are two commonly used targets. We compare the proposed DP-IPD target with other targets. Comparison experiments are conducted with the simulated 2-channel array, and 180°-azimuth localization is performed. The same backbone network as the proposed 2-channel fixed-array model is used, on top of which an extra head will be added if necessary. The following targets are compared, which are all well tuned to achieve their optimal performance on the present task. (1) *Multi-class location classification* [10]. One candidate location is considered as one class. Multi-class classification simultaneously localize multiple sources. At inference, classes with output probability larger than a threshold are detected as active sources. (2) *Multi-track location regression* [28] predicts the location of multiple sources. The 2D coordinate on a unit circle is used to represent the azimuth angle, and [0,0] represents non-source. Frame-level PIT is used for training. (3) *Spatial spectrum regression* [9] predicts a multi-source spatial spectrum, within which each spectral peak represents one source. At inference, the spectral peaks larger than a threshold are detected as active sources. (4) *Multi-track spatial spectrum regression* [15] uses an independent spatial spectrum to represent each source. Two spatial spectrum tracks are arranged in the descending order of the azimuth angles. (5) *Multi-track location classification*. Inspired by the *multi-track spatial spectrum regression* method, we also test one new learning target, i.e. multi-track location classification which

outputs two classification tracks, each track represents one source. The two tracks are also arranged in the descending order of the azimuth angles. (6) *Mixed IPD regression* [16] predicts the DP-IPD as in the proposed method, but the DP-IPD of multiple sources are mixed together with a mixing weight (based on the power proportion of each source) for each source. This way circumvents the source permutation problem. At inference, an iterative source detection and localization technique is applied to iteratively extract the DP-IPD of each source. (7) *Multi-track DP-IPD regression (0-vector)*, namely the proposed DP-IPD target, but uses the all-zero vector to represent non-source. (8) The proposed *multi-track DP-IPD regression*. For all regression-based methods, the loss function is MSE. Except *Multi-track location regression*, all these targets need to set candidate locations, and they all use every 1° as one candidate location.

The results (square root of the sum of squared MDR and squared FAR, denoted as MDR&FAR) as a function of angle difference of the two sources are illustrated in Fig. 6. It can be observed that:

- Multi-track methods outperform their single-track counterparts when the angle difference of two sources is small, for example *multi-track* versus *multi-source spatial spectrum regression*, *multi-track* versus *multi-class location classification*, *multi-track* versus *mixed DP-IPD regression*. When the two sources are close, the peak of multiple sources in the multi-source spatial spectrum tend to merge into one peak, which however can be well separated by setting two independent tracks. It's worth noting that, the single-track methods can handle a flexible number of sources, while the multi-track methods can only output a fixed maximum number of sources.

- In comparison with classification-based methods, regression-based methods gain a larger accuracy improvement when the error tolerance is increased. For example, comparing *multi-track spatial spectrum regression* with *multi-track location classification*, the latter performs better when the error tolerance is 5°, while the former performs better when the error tolerance is 10°. As source location is continuous in space, when the network fails to predict an accurate localization result, using MSE regression loss ensures that the erroneous result does not deviate significantly from the true value. In contrast, the cross-entropy loss lacks such constraint, possibly leading to larger localization error.

- The proposed *multi-track DP-IPD regression* consistently outperforms *multi-track location* and *spatial spectrum regression*, while the latter two performs similarly. DP-IPD estimation is a signal-level task, which can be readily learned from microphone signals. By contrast, location and spatial spectrum estimation require one further step of array-dependent conversion, which may arise more difficulty when directly learned from microphone signals.

- Compared with using the all-zero vector for non-source, the proposed non-source target achieves better performance, which indicates that it is indeed easier for the network to learn the proposed non-source target. To further testify this,

TABLE III: SSL (azimuth localization) performance on simulated data. Error tolerance is 10°.

| | Methods | FLOPs [G/s] | #Params. [M] | 2-CH | | | 4-CH LA | | | 4-CH CA | | | 6-CH LA | | | 6-CH CA | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | MDR [%] | FAR [%] | MAE [°] | MDR [%] | FAR [%] | MAE [°] | MDR [%] | FAR [%] | MAE [°] | MDR [%] | FAR [%] | MAE [°] | MDR [%] | FAR [%] | MAE [°] |
| Online / Fixed-Array | SRP-DNN [16] | 2.3 | 0.8 | 19.9 | 13.1 | 2.9 | 12.5 | 8.2 | 1.9 | 17.1 | 9.0 | 2.5 | 12.2 | 6.7 | 1.8 | 15.9 | 6.4 | 2.3 |
| | IPDnet (prop.) | 19.4/62.8 | 0.7/2.1 | **8.3** | **7.7** | **2.1** | **5.4** | **3.7** | **1.4** | **4.7** | **4.5** | **1.8** | **4.7** | **3.5** | **1.2** | **3.5** | **4.3** | **1.6** |
| Online / Variable-Array | IPDnet (prop.) | 23.2 | 1.1 | 10.5 | 12.1 | 2.3 | 7.6 | 6.7 | 1.5 | 5.5 | 7.8 | 1.8 | 7.7 | 5.5 | 1.3 | 5.1 | 7.7 | 1.7 |
| | w/o mean. | 21.1 | 0.7 | 9.2 | 13.1 | 2.2 | 8.3 | 9.8 | 1.5 | 6.6 | 13.0 | 1.9 | 8.8 | 8.7 | 1.4 | 6.2 | 13.4 | 1.8 |
| Offline / Fixed-Array | SALADNet [10] | 2.5 | 1.1 | 13.5 | 10.2 | 2.1 | 9.4 | 7.9 | 1.5 | 10.3 | 6.0 | 1.7 | 7.3 | 7.4 | 1.4 | 9.6 | 5.0 | 1.5 |
| | SALSA-Lite [24] | 7.5 | 14 | 10.5 | 11.1 | 3.2 | 8.4 | 8.4 | 3.0 | 10.3 | 7.6 | 3.2 | 8.1 | 7.2 | 2.9 | 8.2 | 9.1 | 3.1 |
| | SE-Resnet [49] | 3.3 | 10.2 | 13.6 | 16.5 | 3.3 | 10.1 | 14.4 | 2.9 | 9.8 | 9.2 | 2.9 | 9.0 | 13.6 | 3.0 | 10.3 | 8.9 | 3.0 |
| | IPDnet (prop.) | 34.6/54.3 | 0.6/1.8 | **4.6** | **5.7** | **1.7** | **3.3** | **3.6** | **1.1** | **3.5** | **3.4** | **1.3** | **3.2** | **3.5** | **0.9** | **2.6** | **3.9** | **1.2** |
| Offline / Variable-Array | IPDnet (prop.) | 44.5 | 0.9 | 6.9 | 9.1 | 1.9 | 4.2 | 6.9 | 1.2 | 3.7 | 6.9 | 1.4 | 4.3 | 6.0 | 1.2 | 3.4 | 6.8 | 1.2 |



(a) 0-vector Target
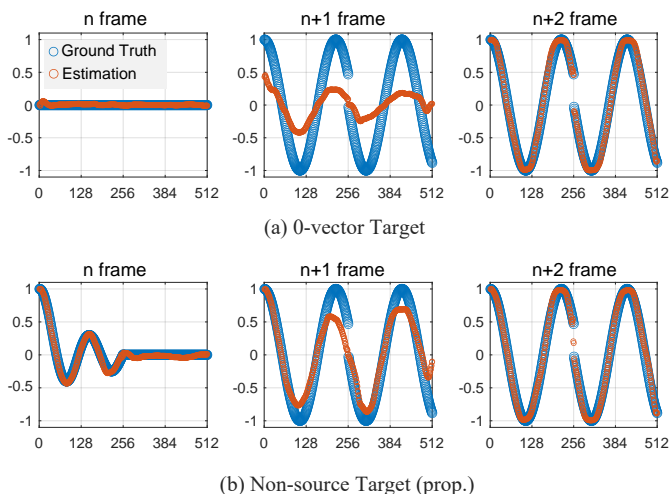


(b) Non-source Target (prop.)

Fig. 7: An example of estimated 512-dimensional DP-IPD vector for three consecutive frames where non-source becomes active source at frame $n + 1$.

in Fig. 7, we plot an example when switching from one non-source frame to one source frame. It can be seen that the proposed non-source target can achieve a quicker transition from the non-source frame to the source frame.

### F. Results on Simulated Data

We conduct both online and offline multi-source localization experiments on the simulated dataset. The following advanced SSL methods are compared with the proposed method: **(1) SRP-DNN [16]** is an online multi-source localization method. It uses a causal CRNN network to estimate the mixed DP-IPD of multiple sources, then the iterative source detection and localization method is used to get the DP-IPD estimate of each source [3]. **(2) SALADNet [10]** uses cascade convolutional and self-attention modules to perform multi-source localization. SALADNet originally takes as input feature the intensity

vector from the first-order Ambisonics, we change it as the microphone signals in this experiment. **(3) SALSA-Lite [24]** is designed for joint sound events localization and detection. It uses the frequency normalized IPD concatenated with the magnitude spectrum as the network input, and uses a ResNet-GRU network. We only use its localization branch and change the target from the sound-class-aligned DOA tracks to simply the multi-source DOAs [4]. **(4) SE-Resnet [49]** is a top-ranked method for joint sound events localization and detection in DCASE22. It uses a squeeze-and-excitation residual network (as encoder) and a Gated Recurrent Unit network (as decoder). As is done for SALSA-Lite, we change its target to the multi-source DOAs. SALADNet, SALSA-Lite and SE-Resnet perform offline localization.

Table III presents the localization performance on five simulated test arrays. The error tolerance is set to 10°. Fixed-array models are independently trained for each test array, while the proposed variable-array model is trained once and used for all test arrays. In the fixed-array experiments, it can be seen that the proposed method prominently outperforms all the comparison methods under all conditions. The advantages of the proposed method are as follows: i) the proposed method (and SRP-DNN and SALADNet) takes as input the microphone signals, while SALSA-Lite and SE-Resnet take as input the noisy IPD (concatenated with the magnitude spectrum). Directly processing the microphone signals is more effective for suppressing the interference of noise and reverberation, as the natural properties of noise and reverberation (such as their spatial-diffuseness) presented in the original signals can be better leveraged; ii) the proposed learning target, i.e. DP-IPD, is more effective than other learning targets, as discussed in the previous section; iii) the proposed full-band and narrow-band fusion network is efficient to exploit the temporal evolution of narrow-band spatial information and the cross-band correlation of localization cues.

We lack direct comparison methods for the proposed variable-array model, but comparing the proposed variable-array model with the fixed-array model still evaluate its efficiency. For the 2-CH array, both the variable-array and fixed-

---

[3] https://github.com/BingYang-20/SRP-DNN

[4] https://github.com/thomeou/SALSA-Lite

array models directly learn the DP-IPD of one microphone pair within one network. The fixed-array model performs better, since it only handle the fixed 4 cm microphone distance. By contrast, the variable-array model handles a large range of microphone distances, i.e. [3, 25] cm, which requires to cover a more large and difficult learning space. For the 4-/6-CH arrays, the pair-wise variable-array model conducts inter-channel communication using the hidden units' mean pooling, which could be sub-optimal compared to the direct inter-channel communication in the fixed-array models. Fortunately, the variable-array model has a reasonable performance degradation relative to the fixed-array models. This indicates that the mean pooling scheme is somehow effective, which can be further testified by the large performance loss when the mean pooling scheme is removed ('w/o mean') as shown in Table III. In addition, the proposed variable-array model outperforms all other fixed-array comparison methods. This showcases the broader applicability and economic training requirements. Especially, when we consider to use real-recorded data to train the SSL model in the future, there will be no need to collect new data for every new microphone array.

Comparing with other methods, the proposed models have a small model size but a large computational complexity. The proposed full-band/narrow-band network processes frames/frequencies independently, which requires a small model as there is no too much information in one frame/frequency, but the network is run many times and thus has a large complexity. The large computational complexity may limit its use in some real-time applications, and this problem will be resolved in our future work.

### G. Results on the LOCATA dataset

*1) Comparison experiments:* We first evaluate the proposed method and comparison methods for azimuth localization on all the six tasks. We only perform online (causal) SSL according to the setting of LOCATA. Two sub-arrays are used: i) microphone 5, 8, 11, and 12 in the Benchmark2 array, which forms a nearly rectangular array located on the top of robot head; i) microphone 6, 7 and 9 in the DICIT array, which forms a 3-channel linear array with a 4 cm microphone distance. The fixed-array models are re-trained for these two arrays using simulated data. The same variable-array model as used in the previous section is directly used in this experiment.

For online SSL, besides SRP-DNN [16], two extra methods are compared: **(1) Cross3D** [12] takes the SRP-PHAT spatial spectrum as input, and uses a causal 3D CNN network to perform moving-source localization [5]. **(2) IcoDOA** [13] uses an Icosahedral CNN to extract localization feature from the SRP-PHAT spatial spectrum, and uses a casual CNN to combine the temporal context for moving-source localization [6]. Different from the proposed model that automatically detects the number of active speakers, Cross3D and IcoDOA localize fixed one speaker, thus they are compared only on the single-speaker tasks, i.e. task 1, 3 and 5, and uses the localization

[5]https://github.com/DavidDiazGuerra/Cross3D
[6]https://github.com/DavidDiazGuerra/IcoDOA

TABLE IV: Azimuth localization performance on all the six tasks of the LOCATA dataset. Error tolerance is 10°.

| Method | Benchmark2 | | | DICIT | | |
|---|---|---|---|---|---|---|
| | MDR[%] | FAR[%] | MAE[°] | MDR[%] | FAR[%] | MAE[°] |
| **Task 1** | | | | | | |
| Cross3D [12] | 23.1 | | 3.7 | 8.8 | | 2.7 |
| IcoDOA [13] | 21.0 | | 3.6 | 14.7 | | 4.0 |
| SRP-DNN [16] | 0.0 | 1.5 | 1.2 | 4.2 | 4.5 | 2.4 |
| IPDnet (fixed) | 0.0 | 0.0 | 1.5 | 1.6 | 1.6 | 2.5 |
| IPDnet (variable) | 0.0 | 0.0 | 2.5 | 0.0 | 0.0 | 2.9 |
| **Task 2** | | | | | | |
| SRP-DNN | 27.8 | 3.8 | 2.4 | 24.4 | 10.0 | 3.0 |
| IPDnet (fixed) | 5.5 | 8.0 | 2.7 | 1.1 | 13.0 | 1.3 |
| IPDnet (variable) | 4.8 | 8.9 | 4.3 | 1.6 | 15.9 | 1.5 |
| **Task 3** | | | | | | |
| Cross3D | 13.9 | | 3.5 | 15.5 | | 3.2 |
| IcoDOA | 11.7 | | 3.3 | 10.6 | | 4.2 |
| SRP-DNN | 1.4 | 5.8 | 1.8 | 1.7 | 1.7 | 2.5 |
| IPDnet (fixed) | 1.8 | 3.4 | 2.0 | 2.6 | 4.8 | 2.0 |
| IPDnet (variable) | 1.5 | 2.5 | 2.8 | 1.2 | 4.2 | 2.1 |
| **Task 4** | | | | | | |
| SRP-DNN | 17.3 | 10.0 | 2.4 | 17.6 | 21.1 | 3.4 |
| IPDnet (fixed) | 9.2 | 8.5 | 2.4 | 7.5 | 14.5 | 2.4 |
| IPDnet (variable) | 11.3 | 8.9 | 3.4 | 8.7 | 16.2 | 2.7 |
| **Task 5** | | | | | | |
| Cross3D | 12.0 | | 3.6 | 4.5 | | 3.5 |
| IcoDOA | 11.8 | | 3.6 | 6.9 | | 3.8 |
| SRP-DNN | 2.3 | 15.2 | 2.1 | 0.4 | 5.3 | 2.6 |
| IPDnet (fixed) | 1.6 | 2.2 | 2.0 | 0.6 | 2.6 | 1.6 |
| IPDnet (variable) | 3.8 | 1.2 | 3.6 | 1.9 | 6.7 | 2.1 |
| **Task 6** | | | | | | |
| SRP-DNN | 8.0 | 12.6 | 2.7 | 33.3 | 40.6 | 3.6 |
| IPDnet (fixed) | 7.0 | 5.2 | 2.5 | 27.4 | 6.4 | 2.5 |
| IPDnet (variable) | 11.3 | 5.7 | 3.6 | 25.4 | 27.8 | 2.7 |
| **AVG. (Single Source)** | | | | | | |
| Cross3D | 18.1 | | 3.6 | 10.1 | | 3.0 |
| IcoDOA | 16.2 | | 3.5 | 12.4 | | 3.9 |
| SRP-DNN | **0.9** | 6.0 | 1.6 | 2.9 | 3.8 | 2.5 |
| IPDnet (fixed) | **0.9** | 1.4 | **1.8** | 1.7 | 2.7 | **2.2** |
| IPDnet (variable) | 1.3 | **0.9** | 2.9 | **0.6** | **2.2** | 2.6 |
| **AVG. (Multi-source)** | | | | | | |
| SRP-DNN | 7.1 | 7.5 | 2.0 | 12.6 | 13.4 | 2.9 |
| IPDnet (fixed) | **3.6** | 3.7 | **2.1** | 6.9 | **9.6** | **2.2** |
| IPDnet (variable) | 4.7 | **3.6** | 3.1 | **6.0** | 10 | 2.5 |

error rate as evaluation metric. The localization error rate is translated to the equal MDR and FAR.

Table IV presents the localization performance. Across all methods, a consistent bias of approximately 4° was observed in the DOA estimation on DICIT data, likely stemming from the annotation bias. To mitigate this effect, we adjusted all DOA estimations by subtracting this bias. It can be seen that the proposed fixed-array model still achieves superior performance compared to other methods, and the proposed variable-array model achieves comparable performance with the proposed fixed-array model. This verifies that the proposed models trained with simulated data can well generalize to real data. Moreover, the proposed variable-array model can well generalize to unseen real microphone arrays.

*2) Generalization across different number of channels and to elevation estimation:* The maximum number of microphones we used for training the variable-array model is 8. In this experiment, we test the variable-array model on a 8-channel and a 12-channel sub-arrays of Benchmark2. The 8-channel sub-array includes the microphones 1, 3, 4, 5, 8, 10, 11, and 12, which forms a nearly cubic array. The 12-channel array includes all the 12 microphones of Benchmark2,

TABLE V: Azimuth and elevation localization performance using different numbers of microphones of Benchmark2.

| Method | Error tolerance: 5° | | | | Error tolerance: 10° | | | |
|---|---|---|---|---|---|---|---|---|
| | MDR [%] | FAR [%] | AZI [°] | ELE [°] | MDR [%] | FAR [%] | AZI [°] | ELE [°] |
| IPDnet (4-mic) | 37.1 | 36.1 | 1.9 | - | **4.7** | 3.6 | 3.1 | - |
| IPDnet (8-mic) | 25.1 | 22.7 | **1.8** | **3.4** | 5.5 | **3.2** | **2.5** | 3.6 |
| IPDnet (12-mic) | **23.1** | **21.6** | **1.8** | **3.4** | 5.1 | 3.6 | **2.5** | **3.4** |



Fig. 8: Azimuth and elevation (trajectory) estimations for a LOCATA example with two moving sources.

which is a nearly spherical array. In addition, as the 8-channel and 12-channel arrays are both 3D and provide the discrimination ability of vertical direction, the elevation angle is also localized.

Table V presents the localization performance, where the results of 4-channel sub-array used in the previous section is also given for comparison. The average performance of all six tasks is reported. It can be observed that the performance measures can be gradually improved with the increase of the number of microphones, especially when the error tolerance is 5°, which indicates that more accurate DP-IPD estimation can be obtained with more microphones. This verifies that the proposed variable-array model can well generalize to microphone array with more channels than training arrays. The elevation angle is also well localized. This demonstrates that, by separating the feature estimation step and the localization step, the proposed method can be flexibly adapted to various SSL configurations. Fig. 8 illustrates the localization result of azimuth and elevation for an example with two moving sources, where the 8-channel array is used. It can be seen that both the azimuth and elevation angles can be well localized, but a larger localization error is obtained for elevation.

## VI. CONCLUSION

This paper proposes a multi-track DP-IPD learning network, named IPDnet, for localization of multiple moving sound sources. The proposed network architecture, i.e. full-band and narrow-band fusion network, is efficient to learn the properties of noise and reverberation and thus to extract reliable DP-IPD of sound sources. The proposed multi-track DP-IPD regression target well disentangles the feature extraction step and the source localization step, and thus outperforms other commonly used SSL targets. Moreover, the proposed variable-array model facilitates the training of SSL network. In this work, the proposed models are trained with pure simulation data in terms of simulated RIR and multi-channel noise, which may has the simulation-to-real problem. To resolve this problem, in future works, the proposed variable-array model can be trained with cross-dataset/array real-recorded data, which provides a reasonable way for alleviating the annotation difficulty and data scarcity when collecting real-world data.

## REFERENCES

[1] H.-Y. Lee, J.-W. Cho, M. Kim, and H.-M. Park, "DNN-based feature enhancement using doa-constrained ica for robust speech recognition," *IEEE Signal Process. Lett.*, vol. 23, no. 8, pp. 1091–1095, 2016.
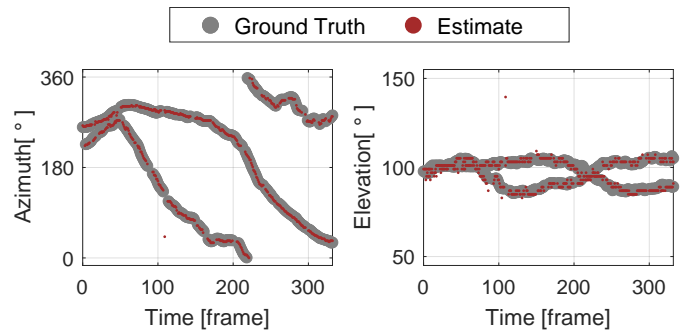
[2] S. E. Chazan, H. Hammer, G. Hazan, J. Goldberger, and S. Gannot, "Multi-microphone speaker separation based on deep DOA estimation," in *Proc. Euro. Signal Process. Conf.*, 2019, pp. 1–5.

[3] M. Raspaud, H. Viste, and G. Evangelista, "Binaural source localization by joint estimation of ILD and ITD," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 18, no. 1, pp. 68–77, 2010.

[4] W. Zhang and B. D. Rao, "A two microphone-based approach for source localization of multiple speech sources," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 18, no. 8, pp. 1913–1928, 2010.

[5] S. Braun, W. Zhou, and E. Habets, "Narrowband direction-of-arrival estimation for binaural hearing aids using relative transfer functions," in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust.*, 2015, pp. 1–5.

[6] Z. Wang, J. Li, Y. Yan, and E. Vincent, "Semi-supervised learning with deep neural networks for relative transfer function inverse regression," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2018, pp. 191–195.

[7] M. Jeub, M. Schaefer, T. Esch, and P. Vary, "Model-based dereverberation preserving binaural cues," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 18, pp. 1732–1745, 2010.

[8] P.-A. Grumiaux, S. Kitić, L. Girin, and A. Guérin, "A survey of sound source localization with deep learning methods," *The Journal of the Acoustical Society of America*, vol. 152, no. 1, pp. 107–151, 2022.

[9] W. He, P. Motlícek, and J.-M. Odobez, "Neural network adaptation and data augmentation for multi-speaker direction-of-arrival estimation," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, pp. 1303–1317, 2021.

[10] P.-A. Grumiaux, S. Kitic, P. Srivastava, L. Girin, and A. Gu'erin, "SALADNet: Self-attentive multisource localization in the ambisonics domain," *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust.*, pp. 336–340, 2021.

[11] B. Yang, H. Liu, and X. Li, "Learning deep direct-path relative transfer function for binaural sound source localization," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, pp. 3491–3503, 2021.

[12] D. Diaz-Guerra, A. Miguel, and J. R. Beltrán, "Robust sound source tracking using SRP-PHAT and 3D convolutional neural networks," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, pp. 300–311, 2020.

[13] D. Diaz-Guerra, A. Miguel, and J. R. Beltran, "Direction of arrival estimation of sound sources using Icosahedral CNNs," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 31, pp. 313–321, 2022.

[14] Y. Fu, M. Ge, H. Yin, X. Qian, L. Wang, G. Zhang, and J. Dang, "Iterative sound source localization for unknown number of sources," in *Proc. INTERSPEECH*, 2022.

[15] H. Yin, M. Ge, Y. Fu, G. Zhang, L. Wang, L. Zhang, L. Qiu, and J. Dang, "MIMO-DoAnet: Multi-channel input and multiple outputs doa network with unknown number of sound sources," in *Proc. INTERSPEECH*, 2022.

[16] B. Yang, H. Liu, and X. Li, "SRP-DNN: Learning direct-path phase difference for multiple moving sound source localization," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, Singapore, 2022, pp. 721–725.

[17] U. Kowalk, S. Doclo, and J. Bitzer, "Geometry-Aware DOA estimation using a deep neural network with mixed-data input features," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2022, pp. 1–5.

[18] L. Wang, Z. Jiao, Q. Zhao, J. Zhu, and Y. Fu, "Framewise multiple sound source localization and counting using binaural spatial audio signals," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2023, pp. 1–5.

[19] J.-H. Cho and J. Chang, "SR-SRP: Super-Resolution based SRP-PHAT for sound source localization and tracking," in *Proc. INTERSPEECH*, 2023, pp. 3769–3773.

[20] Y. Wang, B. Yang, and X. Li, "FN-SSL: Full-band and narrow-band fusion for sound source localization," in *Proc. INTERSPEECH*, 2023, pp. 3779–3783.

[21] K. Tesch and T. Gerkmann, "Insights into deep non-linear filters for improved multi-channel speech enhancement," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 31, pp. 563–575, Nov. 2022.

[22] Y. Yang, C. Quan, and X. Li, "McNet: Fuse multiple cues for multi-channel speech enhancement," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2023, pp. 1–5.

[23] C. Quan and X. Li, "SpatialNet: Extensively learning spatial information for multichannel joint speech separation, denoising and dereverberation," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 32, pp. 1310–1323, 2024.

[24] T. N. T. Nguyen, D. L. Jones, K. N. Watcharasupat, H. A. Phan, and W. Gan, "SALSA-Lite: A fast and effective feature for polyphonic sound event localization and detection with microphone arrays," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2022, pp. 716–720.

[25] X. Li, L. Girin, R. Horaud, and S. Gannot, "Estimation of the direct-path relative transfer function for supervised sound-source localization," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 24, no. 11, pp. 2171–2186, 2016.

[26] S. Mohan, M. E. Lockwood, M. L. Kramer, and D. L. Jones, "Local-ization of multiple acoustic sources with small arrays using a coherence test." *J. Acoust. Soc. Amer.*, vol. 123, no. 4, pp. 2136–2147, 2008.

[27] O. Nadiri and B. Rafaely, "Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 22, no. 10, pp. 1494–1505, Oct. 2014.

[28] K. Shimada, Y. Koyama, S. Takahashi, N. Takahashi, E. Tsunoo, and Y. Mitsufuji, "Multi-ACCDOA: Localizing and detecting overlapping sounds from the same class with auxiliary duplicating permutation in-variant training," *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, pp. 316–320, 2021.

[29] Y. Luo, Z. Chen, N. Mesgarani, and T. Yoshioka, "End-to-end micro-phone permutation and number invariant multi-channel speech separa-tion," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6394–6398.

[30] H. Taherian, S. E. Eskimez, T. Yoshioka, H. Wang, Z. Chen, and X. Huang, "One model to enhance them all: Array geometry agnostic multi-channel personalized speech enhancement," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2021, pp. 271–275.

[31] S. Zhang and X. Li, "Microphone array generalization for multichannel narrowband deep speech enhancement," in *Proc. INTERSPEECH*, 2021, pp. 666–670.

[32] T. Yoshioka, X. Wang, D. Wang, M. Tang, Z. Zhu, Z. Chen, and N. Kanda, "VarArray: Array-geometry-agnostic continuous speech sepa-ration," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2021, pp. 6027–6031.

[33] P. Vecchiotti, N. Ma, S. Squartini, and G. J. Brown, "End-to-end binaural sound localisation from the raw waveform," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2019, pp. 451–455.

[34] X. Xiao, S. Zhao, X. Zhong, D. L. Jones, C. E. Siong, and H. Li, "A learning-based approach to direction of arrival estimation in noisy and reverberant environments," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2015, pp. 2814–2818.

[35] N. Ma, T. May, and G. J. Brown, "Exploiting deep neural networks and head movements for robust binaural localization of multiple sources in reverberant environments," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 25, pp. 2444–2453, 2017.

[36] P. Pertila and M. Parviainen, "Time difference of arrival estimation of speech signals using deep neural networks with integrated time-frequency masking," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2019, pp. 436–440.

[37] B. Yang, R. Ding, Y. Ban, X. Li, and H. Liu, "Enhancing direct-path relative transfer function using deep neural network for robust sound source localization," *CAAI Trans. Intell. Technol.*, vol. 7, pp. 446–454, 2022.

[38] D. Tang, M. Taseska, and T. van Waterschoot, "Supervised contrastive embeddings for binaural source localization," in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust.*, 2019, pp. 358–362.

[39] J. Pak and J. W. Shin, "Sound localization based on phase difference enhancement using deep neural networks," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, pp. 1335–1345, 2019.

[40] S. Adavanne, A. Politis, and T. Virtanen, "Direction of arrival estima-tion for multiple sound sources using convolutional recurrent neural network," in *Proc. Euro. Signal Process. Conf.*, 2017, pp. 1462–1466.

[41] T. N. T. Nguyen, W. S. Gan, R. Ranjan, and D. L. Jones, "Robust source counting and doa estimation using spatial pseudo-spectrum and convolutional neural network," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 28, pp. 2626–2637, 2020.

[42] S. Adavanne, A. Politis, J. Nikunen, and T. Virtanen, "Sound event localization and detection of overlapping sources using convolutional recurrent neural networks," *IEEE J. Selected Topics Signal Process.*, vol. 13, pp. 34–48, 2018.

[43] E. A. P. Habets, I. Cohen, and S. Gannot, "Generating nonstationary multisensor signals under a spatial coherence constraint," *J. Acoust. Soc. Amer.*, vol. 124, no. 5, pp. 2911–2917, 2008.

[44] D. Diaz-Guerra, A. Miguel, and J. R. Beltran, "gpuRIR: A python library for room impulse response simulation with GPU acceleration," *Multimedia Tools Appl.*, vol. 80, no. 4, pp. 5653–5671, 2021.

[45] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2015, pp. 5206–5210.

[46] H. W. Löllmann, C. Evers, A. Schmidt, H. Mellmann, H. Barfuss, P. A. Naylor, and W. Kellermann, "The LOCATA challenge data corpus for acoustic source localization and tracking," in *Proc. IEEE Sensor Array and Multichannel Signal Processing Workshop*, 2018, pp. 410–414.

[47] A. Brutti, L. Cristoforetti, W. Kellermann, L. Marquardt, and M. Omol-ogo, "WOZ acoustic data collection for interactive TV," in *Proc. International Conference on Language Resources and Evaluation*, 2008.

[48] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, 2015.

[49] J. S. Kim, H. J. Park, W. Shin, and S. W. Han, "A robust framework for sound event localization and detection on real recordings," in *Proc. Detect. and Classification of Acoust. Scenes Events Workshop*, 2022.