

 Open access • Proceedings Article • DOI:10.1109/ICASSP39728.2021.9415109

SSLIDE: Sound Source Localization for Indoors Based on Deep Learning

— [Source link](#) 

Yifan Wu, Roshan Ayyalasomayajula, Michael J. Bianco, Dinesh Bharadia ...+1 more authors





Institutions: University of California, San Diego

Published on: 06 Jun 2021 - International Conference on Acoustics, Speech, and Signal Processing

Topics: Acoustic source localization, Reverberation, Deep learning, Microphone and Artificial neural network

Related papers:

- [Blind Sound Source Localization based on Deep Learning](#)
- [Indoor Sound Source Localization based on Sparse Bayesian Learning and Compressed Data](#)
- [Configuration-Invariant Sound Localization Technique Using Azimuth-Frequency Representation and Convolutional Neural Networks](#)
- [Detection Sound Source Direction in 3D Space Using Convolutional Neural Networks](#)
- [Deep Residual Network for Sound Source Localization in the Time Domain.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/sslidle-sound-source-localization-for-indoors-based-on-deep-45teegti0d>

SSLIDE: SOUND SOURCE LOCALIZATION FOR INDOORS BASED ON DEEP LEARNING

Yifan Wu, Roshan Ayyalasomayajula, Michael J. Bianco, Dinesh Bharadia, and Peter Gerstoft

University of California, San Diego, La Jolla, CA, USA

ABSTRACT

This paper presents *SSLIDE*, Sound Source Localization for Indoors using DEep learning, which applies deep neural networks (DNNs) with encoder-decoder structure to localize sound sources with random positions in a continuous space. The spatial features of sound signals received by each microphone are extracted and represented as likelihood surfaces for the sound source locations in each point. Our DNN consists of an encoder network followed by two decoders. The encoder obtains a compressed representation of the input likelihoods. One decoder resolves the multipath caused by reverberation, and the other decoder estimates the source location. Experiments based on both the simulated and experimental data show that our method can not only outperform multiple signal classification (MUSIC), steered response power with phase transform (SRP-PHAT), sparse Bayesian learning (SBL), and a competing convolutional neural network (CNN) approach in the reverberant environment but also achieve a good generalization performance.

Index Terms— Indoor sound source localization, multipath, encoder-decoder structure, deep neural networks

1. INTRODUCTION

Sound source localization (SSL) has widespread applications in human-robot interaction [1], ocean acoustics [2], teleconferencing [3], and automatic speech recognition [4]. For example, in a hospital, attending robots can locate and attend to patients based on their voices [5]. However, SSL in reverberant environments is challenging due to multipath artifacts in received signals. This effect degrades SSL performance. Thus, it is important to develop SSL methods that are robust to reverberation [6].

While traditional SSL algorithms [7–12] rely on estimation theory or statistics, they fail in dynamic and reverberant environments. A well-known subspace based technique, multiple signal classification (MUSIC) [8] is known to suffer from correlated sources which are prevalent in reverberant environments. Another classical SSL method, steered response power with phase transform (SRP-PHAT) [9–11] has been shown to not be robust to non-stationary signal like speech. Recently, SSL approaches based on deep neural networks (DNNs) have been proposed [13–20]. Most of the approaches are based on supervised learning. In [13], a multilayer perceptron DNN is proposed for DOA estimation. In [14], a SSL framework based on convolutional neural network (CNN) is proposed. A learning based SSL approach using discriminative training is presented in [15]. The authors in [16] propose a convolutional recurrent DNN for SSL and sound event detection. In [17], a robust SSL guided by deep learning based time-frequency masking framework was presented. There are also some works using unsupervised learning [18] and semi-supervised learning methods based on manifold learning [19], and deep generative modeling [20]. But all of these methods can only work well when the sensor-source distance is small, which limits their implementation in real-world settings.

In this work, we present *SSLIDE*, a SSL method based on DNN with encoder-decoder structure. Our method can resolve randomly located sources in the room and can achieve a good generalization performance. Inspired by [21], the major novelty of the architecture lies in the two parallel decoders that help in solving two distinct and independent problems. One decoder is designed to resolve the multipath artifacts, and the other to predict the locations of the sound sources. By training these decoders in parallel, the DNN learns to jointly predict the locations of sound sources and remove the multipath artifacts on range offsets. We compare our approach with other baseline SSL methods, including multiple signal classification (MUSIC) [8], steered response power with phase transform (SRP-PHAT) [9–11], Sparse Bayesian Learning (SBL) [22], and CNN [14]. Based on the experiment results, we find *SSLIDE* outperforms the baseline methods and generalizes well across space, perturbations of reverberation time, microphone spacing, and input speech.

2. PROPOSED METHOD

To understand how the proposed DNN solves for the reverberation problem and helps in efficient SSL, let us first look into the fundamentals of sound transmissions in a given environment. Consider the acoustics signals in the time domain

$$y_i = s * h_i + n_i \quad (1)$$

where $y_i \in \mathbb{R}^L$ is the signal received by i th microphone ($i \in \{1, \dots, M\}$, M is the number of microphones), s the source signal, and n_i the noise for the i th microphone. h_i is the room impulse responses (RIRs), which characterizes the reverberation of the room. Denote $\mathbf{y} = [y_1, \dots, y_M]^T \in \mathbb{R}^{M \times L}$ as the collections of the received signal of all sensors with audio length L .

For N arrays with K microphones in each array (i.e. $M = NK$), \mathbf{y} can be reshaped as a tensor with dimension $K \times N \times L$. So, for a given input received signal y with S snapshots and T datapoints (number of independent measurements), there are $C = TS$ frames for y .

2.1. Features extraction

One of the key components in designing a DNN model is to understand the input data and represent it appropriately for the network to be able to learn the required application using the input data. While there have been existing works like MUSIC and SRP-PHAT that enable accurate SSL for environments with low reverberation, their SSL performance degrades significantly in dynamic and reverberant environments.

We first use standard beamforming to obtain source location likelihoods in a 2D space. We obtain $\mathbf{Y} \in \mathbb{C}^{S \times F \times K \times N}$, the STFT output of \mathbf{y} with S snapshots and F frequency bins, where $L = 2FS$. While $2F$ is the total number of frequency bins in the STFT,

we only consider the positive half of the frequency bins. For T datapoints, then there are $C = TS$ frames for \mathbf{Y} .

Assume a uniform linear array (ULA) and a broadband signal. Inspired by [23], we can define a 2-D function [23] which can indicate the likelihood of the signal coming from the angle θ and distance d for array $n \in \{1, \dots, N\}$ and frame $s \in \{1, \dots, S\}$

$$P_{ns}(\theta, d) = \left| \sum_{i=1}^K \sum_{l=1}^F Y_{il} e^{j \frac{2\pi i u \sin \theta f_l}{c}} e^{j \frac{2\pi l f_l d}{c}} \right| \quad (2)$$

where $j = \sqrt{-1}$ and u, f_0, f_l, c stand for the spacing between microphones, median frequency, the frequency corresponding for the l th frequency bin, the speed of sound, respectively. $P_{ns}(\theta, d)$ is beam power. Y_{il} represents the STFT output for the i th microphone and l th frequency bin. When the sound source is from angle θ and distance d , then $P_{ns}(\theta, d)$ have a high value. If we have U and V grid points for θ and d , then we will obtain a likelihood surface with dimension $U \times V$ which can indicate the likelihood of the signal in the given θ and d . Fig. 1 (a) is one of the examples.

For reverberation and noise free data, the localization is simply identifying the θ and d that correspond to the maximum likelihood [21]. Due to the reverberation, much of the sound received by the microphones is a result of multipath, which is a complicated function of the different microphone locations relative to the source. Therefore, peaks in the likelihood surface may no longer indicate the correct result in terms of their predicted distance d as depicted in Fig. 1 (a).

2.2. Range compensation

To help overcome challenges of source localization in reverberant environments, we design a second decoder to explicitly correct for variation in multipath artifacts due to differences in microphone location. Details on the decoder and the loss function are further described in Section 2.4.

To enable this decoder to learn to alleviate range offsets cause by multipath artefacts, we will artificially generate likelihood surfaces with range compensation as labels. To do so, we first identify the direct path as the path with the least range measurement, \hat{d} in the incorrect range image as shown in Fig. 1 (a). We then use the actual range measurement expected range measurement, d , from the given ground truth location for that specific measurement. We then compensate this offset in the given RIR measurement to get the expected likelihood profile as seen in Fig. 1 (b). More formally, for the STFT output in the s th frame and k th microphone of the n th array $Y_k^{ns} \in \mathbb{C}^{F \times 1}$, the range is compensated by

$$\bar{Y}_k^{ns} = Y_k^{ns} \circ e^{j2\pi\vartheta \frac{\hat{d}_n - d_n}{c}} \quad (3)$$

where $\vartheta = [f_1, \dots, f_F]^T \in \mathbb{R}^{F \times 1}$ is a collection of all frequencies. Scalar d_n and \hat{d}_n are the estimated ranges for the direct path and true ranges of the n th array, and \circ represents the Hadamard product. Fig. 1(b) shows the likelihood surface after range compensation. Our results show that the range compensation will make DNN easier to identify the correct location of the sound source.

We have generated two categories of likelihood surfaces with dimension $U \times V$. While we can perform single point identification based object detection tasks on these images, each of these images are with respect to their own microphones and lack the context of the global coordinates. To overcome that problem, we convert these range-angle images into 2D Cartesian images which show the coordinate with respect to one of the arrays. We perform a coordinate

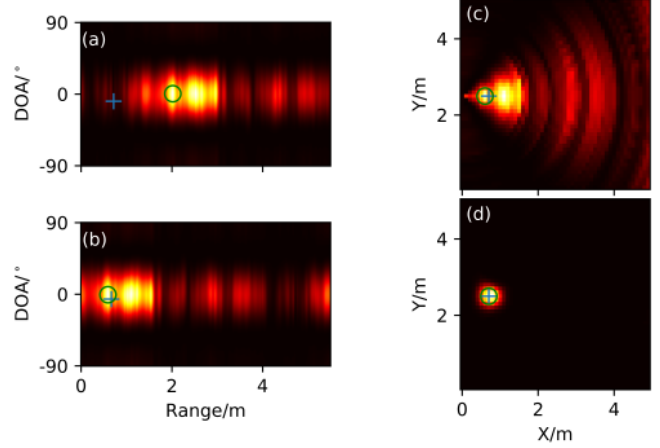


Fig. 1. Likelihood surfaces (a) before and (b) after range compensation, and (c) after range compensation and coordinate transformation, (d) Output likelihood surface from localization decoder. The correct source position (+), and maximum value (for (a) - (c)) / predicted location (for (d)) (o) are indicated. All plots are based on the same testing data.

transform on these images to convert to a 2-D Cartesian plane with dimension $Y \times X$ as shown in Fig. 1(c) to encode the locations of these multiple arrays.

2.3. Target processing

Now that we have defined the images for us to perform the single point identification, we need to define the targets for the network to learn the SSL task. One naive way to generate the target images is to only mark the target position as one and the rest positions as zeros. Unfortunately, this method will make the loss extremely small, which will bring about the gradient underflows. The network cannot learn how to predict the locations due to the almost vanishing gradients.

Thus, we use a negative exponential label to represent the target position. The target of the network will also be a likelihood surface with dimension $Y \times X$. The distance between a random position (x', y') in the likelihood surface and ground-truth position (x, y) is

$$d(x', y') = \sqrt{(x' - x)^2 + (y' - y)^2}; \quad (4)$$

Then its value in the likelihood surface $l(x', y')$ will be marked as

$$l(x', y') = e^{-d(x', y')^2 / \sigma^2} \quad (5)$$

where σ is a hyperparameter controlling the rate of decay. For $d(x', y') = 0$, then $l(x', y') = 1$ its maximum value. Far from the target position, the value will decay significantly. For most of the points in the heatmap, the values is close to 0. These output representations is helpful for a smoother gradient flow.

2.4. SSLIDE architecture

Now that we have the inputs and targets for performing single point identification, we utilize the network architecture as shown in Fig. 2 and based on encoder-decoder architecture with one encoder and two parallel decoders inspired from [21]. The input to the encoder is the likelihood surface without range compensation indicated as (2).

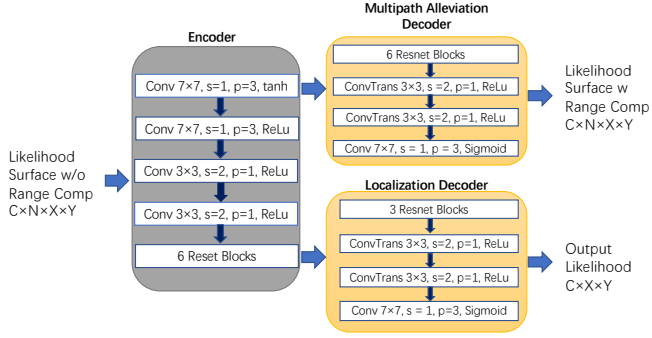


Fig. 2. SSLIDE architecture. C , N , X and Y are followed by the definitions in Sec 2.1–2.4. s and p stand for stride and padding, respectively

The encoder compresses this representation and then feeds to both of the decoders. The two decoders will focus on two different tasks simultaneously. The multipath alleviation decoder will use the likelihood surface with range compensation mentioned in Sec. 2.2 as targets to train the network to generate the likelihood surface without range offsets. With the help of this decoder, the neural network will learn the multipath profile and how to alleviate such an artifact with respect to the range estimation, which will facilitate the localization decoder to identify the source locations. The localization decoder aims to predict the location of the sound source by using the target likelihood surface mentioned in 2.3 as labels. The output for the localization decoder is also a likelihood surface with dimension $Y \times X$. The location with the highest value in this output image will be marked as the predicted location. Note that since we have used the ground-truth position to generate the target images with range compensation, the multipath alleviation decoder will only appear during the training phase, and it will be turned off during the testing phase.

The loss function for the multipath alleviation decoder is l_2 -loss

$$L_{multipath} = \frac{1}{N} \sum_{i=1}^N \|I_{out}^i - I_{target}^i\|_2 \quad (6)$$

where I_{out}^i and I_{target}^i are the decoder outputs and the targets (likelihood surfaces with range compensation) of the i th array, separately. All of the outputs and targets are likelihood surfaces with the same dimension. N is the number of arrays. The advantage of averaging across multiple receiver arrays is that we can enforce consistency of peaks across all the target images, and the network will learn the consistency across these multiple receiver arrays.

For the localization decoder, we use l_2 -norm loss with l_1 regularization to enforce the sparsity as there only exists one global maxima in the output likelihood surface. The loss function of that decoder can be expressed as

$$L_{localization} = \|T_{out} - T_{target}\|_2 + \lambda \|T_{out}\|_1 \quad (7)$$

where T_{out} and T_{target} are the decoder outputs and targets (target images with negative exponential labels), respectively. λ is the regularization term. The loss functions from these two decoders are summed and back-propagated to the input. Fig. 1 (d) shows the output likelihood surfaces from the localization decoder based on the same data as Fig. 1 (a).

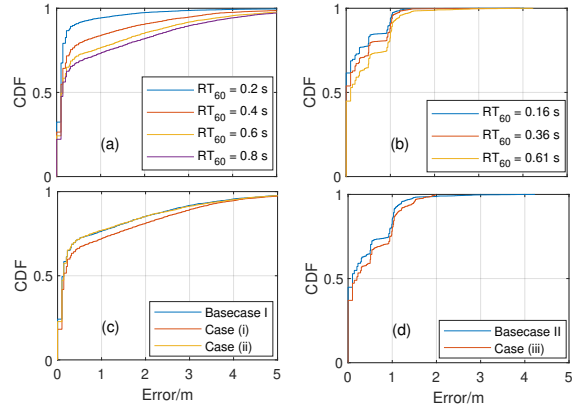


Fig. 3. CDF for (a) simulated data, (b) MIRD data, and generalization experiments for (c) Base Case I and (d) Base Case II.

3. EXPERIMENTS

The localization performance of SSLIDE and other baseline methods is evaluated with different levels of reverberation by using both simulated and real RIRs. The real RIRs are from Multi-Channel Impulse Responses Database (MIRD) [24]. For each case, the networks for the learning-based methods are trained and evaluated on both datasets.

3.1. Datasets

3.1.1. Simulated Data

Simulated RIRs are synthesized by the RIR generator [25], which models the reverberation using the image method [26]. The room is $8 \times 5 \times 4$ m with reverberation times (RT_{60}) from 0.2–0.8 s and speed of sound $c = 340$ m/s. There are $N = 3$ identical ULA with array centers (0, 2.5, 2), (4, 0, 2), and (8, 2.5, 2) m and $K = 4$ sensors in each array with identical space 2.6 cm. To train the generalization across space, the sources have random (x, y) on the boundary of the room in the array plane ($z = 2$ m). We generate $T = 600$ RIRs with random source positions. The sampling frequency is 16 kHz. The input speech signal is a 1 s clean segment randomly chosen from LibriSpeech corpus [27]. The microphone signals are obtained by convolving the RIRs with the speech signal. White noise from the Audio Set [28] is added to give a 35 dB signal-to-noise ratio (SNR).

3.1.2. MIRD Data

These methods are also evaluated on MIRD [24] which provides recorded RIRs for 8-microphone ULA with spacing 8 cm for 3 RT_{60} s. We downsample to the audio frequency from 48 kHz to 16 kHz. All reverberation times (0.16, 0.36, and 0.61 s) are applied to assess the localization performance. There are 2 ranges (1 and 2 m) and 13 candidate DOAs, $[-90, 90^\circ]$ in 15° steps. The sound source is located in one of the 26 candidate positions. We use 20 recordings with 2 s duration and half female/male voices, resulting in $T = 520$ RIRs (datapoints). The noise is generated in the same way as simulation.

RT_{60}/s	Simulated				MIRD		
	.20	.40	.60	.80	.16	.36	.61
Testing	.24	.56	.77	.90	.23	.29	.39
Ablation	.54	1.1	1.8	1.5	.35	.43	.59

Table 1. MAE (m) of localization for testing data (first row) and ablation study (second row).

3.2. Parameters and implementation details

SSLIDE is compared with MUSIC [8], SRP-PHAT [9–11], SBL [22] and CNN [14]. The MUSIC and SRP-PHAT are implemented by Pyroomacoustics [29]. The spectrogram is used as input to train the CNN for classification. We use 1° resolution for MUSIC, SRP-PHAT and SBL in simulations and 15° for MIRD.

For CNN, based on the architecture suggested by [14], we have $M - 1$ convolutional layers with kernel size 2 ($M = 12$ for simulation and 8 for MIRD), and 64 filters per layer. Then, two fully connected layers (512 units for both MIRD and simulation) are added following the convolutional layers. To reduce overfitting, we apply dropout (0.50 dropout rate) in output layer [14].

For both simulated and MIRD data, we use $N_{FFT} = 256$ with no overlap for the STFT implementation, the number of snapshots $S_{sim} = 63$ and $S_{MIRD} = 125$. We only consider positive frequency bins, thus $F = 128$. The size for the simulation and MIRD dataset are $C_{sim} = T \times S_{sim} = 600 \times 63 = 37,800$ and $C_{MIRD} = T \times S_{MIRD} = 520 \times 125 = 65,000$, separately. $\sigma = 0.25$ (See (5)) is chosen for generating the target likelihoods. Fig. 1 (d) is one of the targets likelihood examples with $\sigma = 0.25$, and we can see that it provides a sparse likelihood surface and only a small region of points that are near the target have significant values. For the simulations, the likelihood surface dimension is $101(Y) \times 161(X)$, and for MIRD 121×121 .

The model of SSLIDE is implemented by Pytorch [30] with learning rate 10^{-5} , batch size 32, and weight decay regularization $\lambda = 5 \times 10^{-4}$, and Adam is the optimizer with weight decay 10^{-5} . The data is split based on 70% for training, 15% for validation and 15% for testing. The model is trained for 50 epochs.

3.3. Results and discussions

Testing Performance The cumulative distribution function (CDF) of the localization is shown for 4 RT_{60} s for the simulated data and 3 RT_{60} s for MIRD data. Fig. 3 (a) and (b) show the localization error distribution for simulation and MIRD in the *testing phase*. The localization error for ground-truth (x, y) and predicted (\hat{x}, \hat{y}) location can be expressed as $e = \sqrt{(\hat{x} - x)^2 + (\hat{y} - y)^2}$. The mean absolute error (MAE) of our approach for testing data is listed in the first row of Table 1.

Comparison with Other Baseline Methods Since we can obtain the estimated coordinates of the sound sources, the DOA can also be accessed by simple computations. The comparison of MAE for DOA estimation with other baseline methods is in Table 2. From Table 2, we can see that our approach outperforms all baseline methods for both the simulated and MIRD data in all levels of reverberation. For the simulated data, due to the randomness of the source positions, the localization performance of the baseline methods degrades significantly, especially when the sources are far from the sensors. In contrast, our approach can still have satisfactory performance and thus generalize across space well. For MIRD evaluation, all of the other baseline methods leverage the prior information of the candidate DOAs. In specific, for SRP-PHAT, MUSIC, and SBL,

Method	RT_{60}/s (simulation)				RT_{60}/s (MIRD)		
	.20	.40	.60	.80	.16	.36	.61
SRP-PHAT	14	21	25	25	13	16	19
MUSIC	12	22	27	29	12	17	18
SBL	7.6	13	17	16	11	16	18
CNN	10	14	16	20	4.6	8.0	9.8
SSLIDE	3.0	5.0	6.7	7.9	4.3	5.7	8.1

Table 2. MAE ($^\circ$) of DOA estimation for SSLIDE and other baseline methods

Training	Test Setup	MAE/m	MAE/ $^\circ$	
			SSLIDE	CNN
Base Case I	Base Case I	0.77	6.66	15.6
Base Case I	(i)	0.93	7.85	17.6
Base Case I	(ii)	0.78	6.67	16.9
Base Case II	Base Case II	0.39	8.06	9.81
Base Case II	(iii)	0.45	12.0	23.9

Table 3. Generalization performance of SSLIDE and CNN

when generating the steering vectors, the distribution of DOA (-90° to 90° in 15° steps) is used. For CNN, it also “knows” that there are 13 potential classes. Only our method does not rely on the prior information about candidate DOAs and achieves a competitive localization performance.

Ablation Study To validate the function of multipath alleviation decoder, we conduct ablation study which removes that decoder during the training phase. The MAEs for the testing data are listed in the second row of Table 1. Compared with the first row, we can see that The localization error increases when that decoder is removed for all of the cases, which verifies the role of that decoder in resolving the multipath artifacts.

Generalization Performance First, we evaluate the generalization under perturbations of RT_{60} and microphone locations (see Fig. 3(c)). The model is trained with 0.6 s RT_{60} (Base Case I) but tested with the following two cases: (i) RT_{60} increases to 0.7 s (ii) the same RT_{60} , but microphone spacing increases from 2.6 to 2.7 cm.

Besides, we evaluate generalization performance across different speech for the MIRD data (See Fig. 3(d)). The model is trained with $RT_{60} = 0.61$ s and 20 speech signals (Base Case II) but tested with (iii) 3 new recordings that are not used in training. We also compare the generalization performance with CNN and list the MAEs in Table 3. From that table, we can see that our method has a more robust performance under the perturbations of RT_{60} , microphone positions, and speech signal than CNN.

4. CONCLUSIONS

We developed SSLIDE, a SSL method based on a DNN with an encoder and two decoders which can localize the sources in a continuous space. This enables the DNN to simultaneously predict the locations of sound sources and mitigate multipath artifacts. Experiments indicate our method outperforms MUSIC, SRP-PHAT, SBL, and CNN in environments with different reverberation levels in a continuous space. The ablation study shows the importance of multipath alleviation decoder to reduce multipath and the generalization experiments show strong generalization abilities across space, perturbations of reverberation time and microphone locations, and unseen input recordings.

5. REFERENCES

- [1] K. Nakadai, T. Lourens, H. G. Okuno, and H. Kitano, "Active audition for humanoid," in *Proc. Natl. Conf. Artif. Intell.*, 2000, pp. 832–839.
- [2] K. L. Gemba, S. Nannuru, and P. Gerstoft, "Robust ocean acoustic localization with sparse bayesian learning," *IEEE J. Sel. Topics Signal Process.*, vol. 13, no. 1, pp. 49–60, 2019.
- [3] S. Zhao, S. Ahmed, Y. Liang, K. Rupnow, D. Chen, and DL. Jones, "A real-time 3D sound localization system with miniature microphone array for virtual reality," in *Proc. IEEE ICIEA*, 2012, pp. 1853–1857.
- [4] X. Xiao, S. Zhao, DHH. Nguyen, X. Zhong, DL. Jones, E. S. Chng, and H. Li, "The NTU-ADSC systems for reverberation challenge 2014," in *REVERB Challenge Workshop*, 2014.
- [5] S. Argentieri, P. Danès, and P. Souères, "A survey on sound source localization in robotics: From binaural to array processing methods," *Computer Speech & Language*, vol. 34, no. 1, pp. 87–112, 2015.
- [6] M. J. Bianco, P. Gerstoft, J. Traer, E. Ozanich, M. A. Roch, S. Gannot, and C.-A. Deledalle, "Machine learning in acoustics: Theory and applications," *J. Acoust. Soc. Am.*, vol. 146, no. 5, pp. 3590–3628, 2019.
- [7] J. H. DiBiase, H. F. Silverman, and M. S. Brandstein, "Robust localization in reverberant rooms," in *Microphone arrays*, pp. 157–180. Springer, 2001.
- [8] R. O. Schmidt, "Multiple emitter location and signal parameter estimation," *IEEE Trans. Antennas Propagat.*, vol. 34, no. 3, pp. 276–280, 1986.
- [9] J. Dmochowski, J. Benesty, and S. Affes, "A generalized steered response power method for computationally viable source localization," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 15, no. 8, pp. 2510–2526, 2007.
- [10] J. Velasco, C. J. Martn-Arguedas, J. Macias-Guarasa, D. Pizarro, and M. Mazo, "Proposal and validation of an analytical generative model of SRP-PHAT power maps in reverberant scenarios," *Signal Process.*, vol. 119, pp. 209–228, 2016.
- [11] D. Yook, T. Lee, and Y. Cho, "Fast sound source localization using two-level search space clustering," *IEEE Trans. Cybern.*, vol. 46, no. 1, pp. 20–26, 2015.
- [12] T. Gustafsson, B. D. Rao, and M. Trivedi, "Source localization in reverberant environments: Modeling and statistical analysis," *IEEE Trans. Speech Audio Process.*, vol. 11, no. 6, pp. 791–803, 2003.
- [13] X. Xiao, S. Zhao, X. Zhong, DL. Jones, E. S. Chng, and H. Li, "A learning-based approach to direction of arrival estimation in noisy and reverberant environments," in *IEEE ICASSP*, 2015, pp. 2814–2818.
- [14] S. Chakrabarty and E. A. P. Habets, "Multi-speaker DOA estimation using deep convolutional networks trained with noise signals," *IEEE J. Sel. Topics Signal Process.*, vol. 13, no. 1, pp. 8–21, 2019.
- [15] R. Takeda and K. Komatani, "Sound source localization based on deep neural networks with directional activate function exploiting phase information," in *IEEE ICASSP*, 2016, pp. 405–409.
- [16] 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. Topics Signal Process.*, vol. 13, no. 1, pp. 34–48, 2018.
- [17] Z. Q. Wang, X. Zhang, and D.-L. Wang, "Robust speaker localization guided by deep learning-based time-frequency masking," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 1, pp. 178–188, 2018.
- [18] R. Takeda and K. Komatani, "Unsupervised adaptation of deep neural networks for sound source localization using entropy minimization," in *IEEE ICASSP*, 2017, pp. 2217–2221.
- [19] B. Laufer-Goldshtein, R. Talmon, and S. Gannot, "Semi-supervised sound source localization based on manifold regularization," *IEEE/ACM Trans. Audio, Speech, and Lang. Process.*, vol. 24, no. 8, pp. 1393–1407, 2016.
- [20] M. J. Bianco, S. Gannot, and P. Gerstoft, "Semi-supervised source localization with deep generative modeling," *arXiv preprint arXiv:2005.13163*, 2020.
- [21] R. Ayyalasomayajula, A. Arun, C. Wu, S. Sharma, A. Sethi, D. Vasisht, and D. Bharadia, "Deep learning based wireless localization for indoor navigation," in *MobiCom*, 2020, pp. 1–14.
- [22] A. Xenaki, J. B. Boldt, and M. G. Christensen, "Sound source localization and speech enhancement with sparse Bayesian learning beamforming," *J. Acoust. Soc. Am.*, vol. 143, no. 6, pp. 3912–3921, 2018.
- [23] R. Ayyalasomayajula, D. Vasisht, and D. Bharadia, "BLoc: CSI-based accurate localization for BLE tags," in *Int. Conf. Emerg. Netw. Exp. Technol. (CoNext)*, 2018, pp. 126–138.
- [24] E. Hadad, F. Heese, P. Vary, and S. Gannot, "Multichannel audio database in various acoustic environments," in *Proc. Int. Workshop Acoust. Signal Enh.*, 2014, pp. 313–317.
- [25] E. A. P. Habets, "Room Impulse Response (RIR) generator," (*Online*) Available: <https://github.com/ehabets/RIR-Generator>, 2016.
- [26] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," *J. Acoust. Soc. Am.*, vol. 65, no. 4, pp. 943–950, 1979.
- [27] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *IEEE ICASSP*, 2015, pp. 5206–5210.
- [28] J. F. Gemmeke, D. P.W. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in *IEEE ICASSP*, 2017, pp. 776–780.
- [29] R. Scheibler, E. Bezzam, and I. Dokmanić, "Pyroomacoustics: A python package for audio room simulation and array processing algorithms," in *IEEE ICASSP*, 2018, pp. 351–355.
- [30] A. Paszke, S. Gross, F. Massa, A. Lerer, and J. Bradbury et al., "Pytorch: An imperative style, high-performance deep learning library," in *Adv. Neural Info. Process. Sys.*, 2019, pp. 8026–8037.