

# ONE MODEL TO ENHANCE THEM ALL: ARRAY GEOMETRY AGNOSTIC MULTI-CHANNEL PERSONALIZED SPEECH ENHANCEMENT

Hassan Taherian<sup>1,2\*</sup>, Sefik Emre Eskimez<sup>1</sup>, Takuya Yoshioka<sup>1</sup>, Huaming Wang<sup>1</sup>,  
Zhuo Chen<sup>1</sup> and Xuedong Huang<sup>1</sup>

<sup>1</sup>Microsoft, Redmond, WA, USA

<sup>2</sup>The Ohio State University, Columbus, OH, USA

taherian.1@osu.edu, {seeskim, tayoshio, huawang, zhuc, xdh}@microsoft.com

## ABSTRACT

With the recent surge of video conferencing tools usage, providing high-quality speech signals and accurate captions have become essential to conduct day-to-day business or connect with friends and families. Single-channel personalized speech enhancement (PSE) methods show promising results compared with the unconditional speech enhancement (SE) methods in these scenarios due to their ability to remove interfering speech in addition to the environmental noise. In this work, we leverage spatial information afforded by microphone arrays to improve such systems' performance further. We investigate the relative importance of speaker embeddings and spatial features. Moreover, we propose a new causal array-geometry-agnostic multi-channel PSE model, which can generate a high-quality enhanced signal from arbitrary microphone geometry. Experimental results show that the proposed geometry agnostic model outperforms the model trained on a specific microphone array geometry in both speech quality and automatic speech recognition accuracy. We also demonstrate the effectiveness of the proposed approach for unseen array geometries.

**Index Terms**— multi-channel speech enhancement, target speech extraction, spatial features, microphone array.

## 1. INTRODUCTION

Video conferencing has become essential in the emerging hybrid work era catalyzed by the COVID-19 pandemic. Most modern telecommunication services are equipped with a causal/real-time speech enhancement (SE) front-end to deliver high-quality speech audio in noisy environments. Recently, "personalized" SE methods are emerging in the research field by utilizing an enrollment utterance of a target speaker as additional information to not only suppress the ambient noise and reverberation but also remove interfering speech [1, 2]. With the ability to handle overlapped speech, personalized speech enhancement (PSE) models significantly improve the perceptual speech quality and the performance of downstream tasks such as automatic speech recognition (ASR) [1].

The SE performance can be improved further by using microphone arrays. With multiple microphones, spatial information can be extracted and combined with spectral information for obtaining better SE models [3, 4, 5]. This paper extends the PSE [1] to utilize the microphone arrays for environments where strong noise, reverberation, and an interfering speaker are present. We show that the combination of the enrolled speaker's embedding and the spatial features can significantly improve the SE performance. We also

examine the impact of the speaker embedding and spatial features in challenging conditions where the target and interfering speakers have similar angles or distances to the array.

In addition, we introduce an array-geometry-agnostic PSE model that works regardless of the number of microphones and the array shape. This enables us to train the model for once and use it across multiple microphone array devices without additional operations such as adaptation or retraining, allowing the developed solution to scale significantly. Meanwhile, the geometry agnostic model also yields consistent improvements over fixed-geometry enhancement networks that are trained for matched array geometries. We also show the effectiveness of our proposed model for unseen array geometries and discuss its limitations.

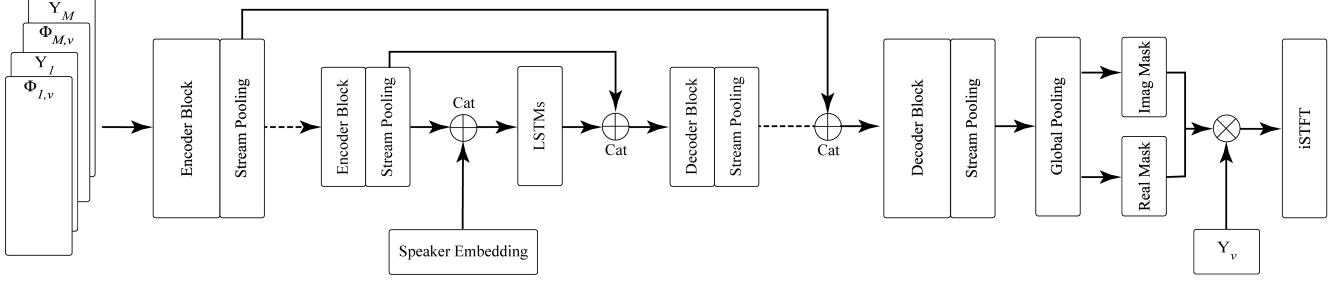
## 2. RELATED WORK

Several studies utilized speaker embeddings to extract the target speaker voice in speech separation tasks [6, 7, 8, 9]. In [8], spatial features from a binaural setup were included as an auxiliary cue to improve the separation performance for same-gender mixtures. [9] incorporated a set of fixed beamformers and an attention network to select the dominant beam using the target speaker embedding. The speaker embeddings have also been employed for SE. A perceptually motivated PSE model with low complexity was proposed in [2]. [1] introduced two real-time PSE models and tested with reverberant target speech corrupted by both noise and interfering speech. We employ their personalized DCCRN (pDCCRN) model as our single-channel baseline system.

Along a related direction, multi-channel geometry invariant modeling was explored recently [10, 11, 12, 13]. In [10], the authors proposed a transform-average-concatenate (TAC) layer for multi-channel speech separation that was invariant to the order of the microphones. An inter-channel processing layer based on self-attention was proposed in [11]. The work by [13] used deep symmetric set layers based on a Siamese network for speech dereverberation. [12] proposed a model to process a variable number of microphone pairs for speech enhancement.

Compared to the aforementioned models, our proposed geometry agnostic model has a more straightforward design. To process inputs with variable dimensions (i.e., microphones), we introduce stream pooling layers for convolutional neural networks that are only based on the averaging and concatenation of the feature maps. Furthermore, our experiments focus on real-time processing scenarios and are conducted on an extensive list of array geometries, which have not been previously explored.

\*Work done during internship at Microsoft Research.



**Fig. 1:** Geometry agnostic multi-channel PSE model is shown. Skip connections link encoders to the corresponding decoders.  $Y_M$  and  $\Phi_{M,v}$  represent STFT and IPD features of microphone  $M$ , respectively. ‘Cat’ refers to concatenation and symbol  $\otimes$  is point-wise complex multiplication.

### 3. SYSTEM DESCRIPTION

#### 3.1. Baseline: Single-channel PSE

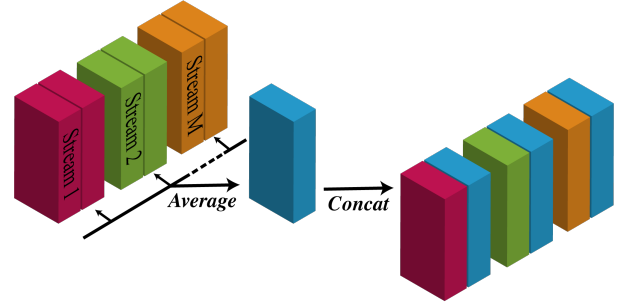
The single-channel PSE model, based on which we build our multi-channel models, performs complex spectral mapping using pDCCRN [1, 14]. pDCCRN has a U-Net architecture with encoder and decoder blocks and two complex LSTM layers in between. Each block contains complex 2-D convolutional layers followed by complex batch normalization. Complex layers are formed by two separate real-valued layers that operate on the real and imaginary parts of the layer input [15]. The pDCCRN uses a d-vector and the mixture signal  $Y$  in the short-time Fourier transform (STFT) domain. The real and imaginary parts of the mixture signal are concatenated to form the input feature  $x \in \mathbb{R}^{C \times 2F \times T}$ ,  $C = 1$ , where  $C$ ,  $F$  and  $T$  are the number of channels for convolution, frequency bins, and time frames. We replicate the d-vector through time dimension, concatenate it with the encoder’s output, and feed the concatenated vector into the first LSTM layer. The model outputs a complex ratio mask [16] which is multiplied with the input mixture to estimate the clean speech. We use 6 encoder layers with  $C = [16, 32, 64, 128, 128, 128]$ , kernel size of  $5 \times 2$  and stride of  $2 \times 1$ . The model is trained with a power-law compressed phase-aware mean-squared error (MSE) loss function [17]. We emphasize that all operations are causal and operate in real-time. For details about our d-vector extractor, we refer the reader to [18].

#### 3.2. Multi-channel PSE for Fixed Microphone Arrays

We employ two approaches to extend the single-channel PSE model to  $M$  microphones with a fixed array geometry. In the first approach, real and imaginary parts of all microphones STFT are stacked in the channel dimension ( $C$ ) to create the input  $x \in \mathbb{R}^{M \times 2F \times T}$ , which is fed to the PSE model. With this simple extension, the model can implicitly learn the spectral and spatial information [4, 19].

In the second approach, we explicitly extract the spatial information. As with [8], we use inter-channel phase difference (IPD) as the spatial feature. The IPD for the microphone pair  $(i, j)$  is defined as  $\Phi_{i,j} = \angle(Y_i/Y_j)$ . The IPD features are calculated between the first microphone and each of the other  $M - 1$  microphones. For each pair, we concatenate the cosine and sine values of the IPD features. Finally, we stack all IPD features, as well as the real and imaginary parts of the first microphone’s STFT to form the input feature  $x \in \mathbb{R}^{M \times 2F \times T}$ .

For both approaches, as with the single-channel case, the model estimates the complex-valued masks which would recover the clean signal when applied to the first microphone.



**Fig. 2:** Schematic diagram of stream pooling layer is shown. Each cube represents a stream (microphone) tensor.

#### 3.3. Geometry Agnostic Modeling

The proposed multi-channel geometry-agnostic PSE model can be described as follows. We first create a virtual microphone signal  $Y_v$  by simply taking the average of all the input microphones:

$$Y_v = \frac{1}{M} \sum_{n=1}^M Y_i. \quad (1)$$

Then, we extract the IPD features for each microphone with respect to the virtual microphone:  $\Phi_{i,v} = \angle(Y_i/Y_v)$ . We also normalize the IPD features using the unbiased exponentially weighted moving average [20] to increase the robustness of the model to arbitrary array geometries. Compared with the fixed geometry models, we observed that the IPD normalization was critical for the geometry agnostic models since the microphone arrangement could be different during training and testing.

In the geometry agnostic model, we introduce an additional dimension to all layers that contain specific stream (microphone) information and append it as the first dimension of the input tensor of each layer. Thus, the model input  $x \in \mathbb{R}^{M \times C \times 2F \times T}$  contains STFT and IPD features of each stream in the channel dimension, i.e.,  $C = 2$ . We refer to the first dimension as a stream.

Figure 1 illustrates the geometry agnostic model architecture. The model is applicable to any number of microphones, and the output is invariant to the permutation of the microphones. Given the input tensor described above, we process each stream information in parallel by using pDCCRN. To utilize the spatio-temporal patterns exhibited in the input multi-channel audio, we include a stream pooling layer after every encoder and decoder blocks of pDCCRN. In these layers, we split the channel dimension into two parts: one

**Table 1:** Evaluation results for fixed geometry multi-channel PSE models with a 7-channel circular array with a radius of 4.25 cm are shown. ‘MC PSE’ refers to multi-channel PSE. We show WER(%), SDR (dB), STOI(%). The best values are highlighted with bold font.

	A			B			B (similar angle)			B (similar distance)		
	WER	SDR	STOI	WER	SDR	STOI	WER	SDR	STOI	WER	SDR	STOI
Noisy mixture	22.38	5.65	73.54	41.07	2.74	69.07	39.75	2.63	68.41	40.97	2.64	68.58
Single-channel PSE	27.65	10.79	84.07	38.09	8.55	79.30	37.32	8.60	79.09	39.20	8.07	78.32
MC PSE (STFT)	20.73	12.22	88.52	24.91	10.14	85.04	28.35	<b>9.56</b>	83.59	25.85	9.78	84.48
MC PSE (IPD)	<b>19.88</b>	<b>12.23</b>	<b>88.78</b>	<b>22.42</b>	<b>10.44</b>	<b>86.38</b>	<b>27.26</b>	9.33	<b>83.79</b>	<b>23.24</b>	<b>9.93</b>	<b>85.44</b>
–w/o d-vector	20.27	11.82	88.11	25.80	9.40	84.37	30.00	8.53	82.23	29.17	8.19	81.65

is unique to each stream; the other is used to aggregate information across the streams. Each cross-stream convolution channel is averaged across the streams and then appended to the stream-specific channels of each stream. A diagram of the stream pooling layer is shown in Fig. 2. At the output layer of pDCCRN, we use a global-pooling layer to average across all the streams and channels to estimate complex masks. The estimated masks are applied to the virtual microphone STFT,  $Y_v$ .

#### 4. EXPERIMENTAL SETUP

We evaluated our multi-channel PSE models by using simulated data that covered various scenarios. We generated room impulse responses (RIRs) for a 7-channel circular microphone array with radius  $r = 4.25$  cm (see Table 2.IV) based on the image method with T60 in the range of 0.15 to 0.6 seconds. The microphone array was located in the room center. Target and interfering speakers were positioned randomly around the array within [0.5, 2.5] meters with the assumption of the target speaker being always closer to the array.

For the training and validation datasets, we simulated 2000 and 50 hours of audio, respectively, based on the clean speech data from the DNS challenge dataset [21]. In both datasets, 60% of utterances contained the target and interfering speakers with a signal-to-interference ratio (SIR) between 0 to 10 dB. The mixed audio was further corrupted by simulated isotropic noises and directional noises from the Audioset and Freesound datasets [22, 23] with a signal-to-noise ratio (SNR) in the range of [0, 15] dB. The sampling rate for all utterances was 16 kHz. We trained our geometry agnostic model with the 7-channel circular array and 3 other geometries derived from it: 4-channel triangular, 4-channel rectangular, and 6-channel circular arrays.

We created two 10-hour test datasets called A and B. Dataset A contained utterances mixed only with ambient noise and reverberation. In contrast, dataset B contained utterances mixed with the ambient noise, reverberation, and interfering speech. Clean utterances were selected from internal conversational speech recordings with high neural network-based mean opinion score (MOS) values [24]. The SIR and SNR ranges were the same as in the training dataset. We convolved the test utterances with RIRs from 8 different geometries. Four geometries were the same as the ones used for the training dataset. The other four geometries were unseen during the training and included a 3-channel triangular array with  $r = 4.25$  cm, a 5-channel circular array with  $r = 3$  cm, a 3-channel linear array with 6 cm length, and an 8-channel circular array with  $r = 10$  cm (as used in the AMI corpus [25]).

We evaluated the enhanced speech signal based on the word error rate (WER) and two signal quality metrics, i.e., signal-to-

distortion ratio (SDR) and short-time objective intelligibility (STOI). We followed the setup described in [17] for the ASR evaluation. We used two baselines for comparison with our geometry agnostic model. For each array geometry that was used for training, we trained a fixed geometry model based on the IPD features without normalization. The other baseline was based on processing each microphone independently with a single-channel PSE model followed by averaging the enhanced signals. Although this approach is computationally expensive, it is an acceptable alternative for unknown array geometries. Note that we also tried to use MVDR beamforming followed by the single-channel PSE; however, the results were worse, probably due to the presence of a competing speaker and inaccurate beam steering. Also, since we assumed that we did not have any knowledge about the microphone array geometries, it would be difficult, if not impossible, to do beamforming in real-time.

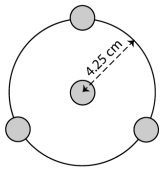
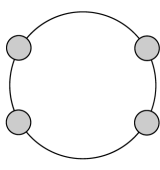
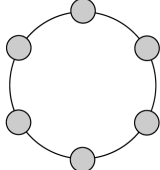
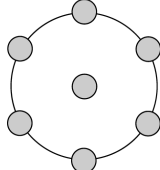
#### 5. RESULTS AND DISCUSSIONS

Table 1 shows the experimental results of different PSE models trained with the 7-channel circular array. Both STFT- and IPD-based multi-channel PSE models substantially outperformed the single-channel PSE model in all scenarios. As with previous studies [17, 26], we observed that the single-channel PSE introduced processing artifacts that yielded worse WER scores compared to the unprocessed noisy mixture for dataset A (i.e., no interfering speech). By contrast, the multi-channel PSE models improved the speech recognition performance.

With regard to the comparison between the two spatial features for the multi-channel PSE, the model trained with the IPD features performed consistently better than the model based on the stacked STFTs. We also trained a multi-channel PSE model based on the IPD features without using d-vectors. The results show that spatial information was helpful regardless of the presence of d-vectors.

To further investigate the effect of the spatial features, we created two versions of test dataset B. In the first version, speakers were randomly positioned such that the difference of their angles with respect to the array was less than 5 degrees while their distance difference was more than 1 meter. In the second version, the angle difference of the speakers was more than 45 degrees, while the distance difference was less than 10 cm. We observed that the performance gap between the single-channel and multi-channel models was smaller with the similar speaker angle setting than with the other settings. This shows the usefulness of the speakers’ directions for discriminating overlapping voices. Also, when the two speakers were at similar distances, using d-vectors substantially improved the performance of the multi-channel PSE.

**Table 2:** Performance comparison of fixed geometry and geometry agnostic multi-channel PSE models for seen microphone arrays is shown.

													
		I. Triangular (4ch)			II. Rectangular (4ch)			III. Circular (6ch)			IV. Circular (7ch)		
		WER	SDR	STOI	WER	SDR	STOI	WER	SDR	STOI	WER	SDR	STOI
A	Signal Averaging	25.30	10.75	84.89	25.80	10.50	84.34	24.93	10.60	84.69	24.86	10.73	84.92
	Fixed Geometry	21.19	11.50	86.99	22.84	11.15	86.15	20.85	11.48	87.39	19.88	12.23	88.78
	Geometry Agnostic	<b>20.75</b>	<b>12.02</b>	<b>88.01</b>	<b>21.41</b>	<b>11.79</b>	<b>87.34</b>	<b>19.99</b>	<b>12.08</b>	<b>88.41</b>	<b>19.65</b>	<b>12.24</b>	<b>88.86</b>
B	Signal Averaging	34.72	8.62	80.54	35.25	8.44	80.00	34.26	8.53	80.45	34.06	8.62	80.68
	Fixed Geometry	24.95	9.72	84.55	25.42	9.64	84.20	23.20	9.96	85.36	22.42	10.44	86.38
	Geometry Agnostic	<b>22.85</b>	<b>10.62</b>	<b>86.67</b>	<b>24.06</b>	<b>10.36</b>	<b>85.81</b>	<b>21.78</b>	<b>10.76</b>	<b>87.34</b>	<b>21.35</b>	<b>10.90</b>	<b>87.76</b>

**Table 3:** Geometry agnostic multi-channel PSE results for unseen array geometries are shown.

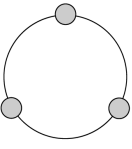
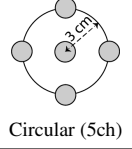
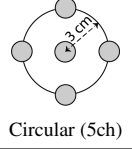
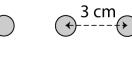
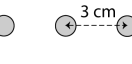
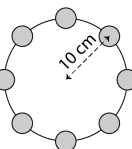
		WER	SDR	STOI
 Triangular (3ch)	A Signal Averaging	25.71	10.52	84.50
	A Fixed Geometry	23.29	11.22	85.90
	A Geometry Agnostic	<b>22.10</b>	<b>11.65</b>	<b>86.96</b>
 Circular (5ch)	B Signal Averaging	34.91	8.45	80.14
	B Fixed Geometry	25.62	9.84	84.23
	B Geometry Agnostic	<b>24.53</b>	<b>10.27</b>	<b>85.51</b>
 Circular (5ch)	A Signal Averaging	24.00	11.41	85.77
	A Geometry Agnostic	<b>20.22</b>	<b>12.51</b>	<b>88.45</b>
	B Signal Averaging	34.51	8.77	80.73
 Linear (3ch)	B Geometry Agnostic	<b>22.34</b>	<b>10.90</b>	<b>87.16</b>
	A Signal Averaging	24.61	11.37	85.49
	A Geometry Agnostic	<b>23.99</b>	<b>11.79</b>	<b>86.07</b>
 Linear (3ch)	B Signal Averaging	35.11	8.73	80.36
	B Geometry Agnostic	<b>28.82</b>	<b>9.75</b>	<b>83.51</b>
 Circular (8ch)	A Signal Averaging	<b>22.28</b>	8.93	<b>84.19</b>
	A Geometry Agnostic	25.99	<b>9.33</b>	82.27
	A Geometry Agnostic -w/o IPD Norm	30.99	8.48	78.87
	B Signal Averaging	32.82	7.03	79.62
	B Geometry Agnostic	<b>29.30</b>	<b>8.11</b>	<b>80.65</b>
	B Geometry Agnostic -w/o IPD Norm	35.41	7.36	77.01

Table 2 shows the geometry agnostic multi-channel PSE results for the arrays seen during the training. Our proposed geometry agnostic model outperformed all the fixed geometry models trained for the corresponding array geometries in both test datasets. This result suggests that not only can our approach make the model independent of the microphone array geometries, but also it can model the relationship between multiple-channel feature streams more efficiently than simply concatenating them. Without requiring to change the model architecture for individual arrays, a single model can be shared between multiple arrays with different shapes and the number of microphones. Also, training with different geometries has an

augmentation effect which improves the robustness compared with fixed geometry training.

Table 3 shows the geometry agnostic model results for the microphone arrays that were not used during the training. We observed that the geometry agnostic model still outperformed the fixed geometry model for the 3-channel triangular array, which has fewer microphones than the arrays included in the training. For the 5-channel circular array, which has a different microphone arrangement, the geometry agnostic model performed very well, achieving the performance comparable to the seen array geometries in Table 2. Regarding the 3-channel linear array, the geometry agnostic model showed consistent improvements over the average of enhanced single-channel signals despite not being exposed to the front-back ambiguity of the linear arrays during the training.

For the 8-channel circular array with  $r = 10$  cm, the geometry agnostic model improved the performance compared with the average of the enhanced signals to a smaller extent, and the results for dataset A were worse in terms of WER and STOI. We speculate the spatial aliasing was the reason for the relatively poor performance of the 8-channel circular array. A large inter-microphone distance leads to spatial aliasing, and this can introduce unseen patterns for the IPD features. For example, spatial aliasing occurs in the array geometries where its inter-microphone distance is longer than 4.28 cm with 16 kHz sampling rate [27]. We observed that if the model was trained without IPD normalization, the performance degraded significantly, suggesting the spatial aliasing problem can be mitigated by IPD normalization.

## 6. CONCLUSIONS

In this work, we utilized spatial features along with speaker embeddings for personalized speech enhancement and showed their combination significantly improved the performance for both ASR and signal quality. Furthermore, we proposed a new architecture and introduced the stream pooling layer to perform multi-channel PSE with any number and arrangement of microphones in a way where the output is invariant to the microphone order. Our proposed model consistently outperformed the geometry-dependent models. Future challenges include mitigating the spatial aliasing problem.

## 7. REFERENCES

- [1] S. E. Eskimez, T. Yoshioka, H. Wang, X. Wang, Z. Chen, and X. Huang, “Personalized speech enhancement: New models and comprehensive evaluation,” *arXiv:2110.09625*, 2022.
- [2] R. Giri, S. Venkataramani, J.-M. Valin, U. Isik, and A. Krishnaswamy, “Personalized percepnet: Real-time, low-complexity target voice separation and enhancement,” *arXiv:2106.04129*, 2021.
- [3] Z.-Q. Wang and D. Wang, “All-neural multi-channel speech enhancement,” in *Proc. Interspeech*, 2018, pp. 3234–3238.
- [4] S. Chakrabarty and E. A. Habets, “Time–frequency masking based online multi-channel speech enhancement with convolutional recurrent neural networks,” *IEEE Journal of Selected Topics in Signal Processing*, pp. 787–799, 2019.
- [5] B. Tolooshams, R. Giri, A. H. Song, U. Isik, and A. Krishnaswamy, “Channel-attention dense U-Net for multichannel speech enhancement,” in *Proc. ICASSP*, 2020, pp. 836–840.
- [6] Q. Wang, H. Muckenhirn, K. Wilson, P. Sridhar, Z. Wu, J. Hershey, R. A. Saurous, R. J. Weiss, Y. Jia, and I. L. Moreno, “Voicefilter: Targeted voice separation by speaker-conditioned spectrogram masking,” *arXiv:1810.04826*, 2018.
- [7] P. Wang, Z. Chen, X. Xiao, Z. Meng, T. Yoshioka, T. Zhou, L. Lu, and J. Li, “Speech separation using speaker inventory,” in *Proc. ASRU*, 2019, pp. 230–236.
- [8] M. Delcroix, T. Ochiai, K. Zmolikova, K. Kinoshita, N. Tawara, T. Nakatani, and S. Araki, “Improving speaker discrimination of target speech extraction with time-domain speakerbeam,” in *Proc. ICASSP*, 2020, pp. 691–695.
- [9] G. Li, S. Liang, S. Nie, W. Liu, M. Yu, L. Chen, S. Peng, and C. Li, “Direction-aware speaker beam for multi-channel speaker extraction,” in *Proc. Interspeech*, 2019, pp. 2713–2717.
- [10] Y. Luo, Z. Chen, N. Mesgarani, and T. Yoshioka, “End-to-end microphone permutation and number invariant multi-channel speech separation,” in *Proc. ICASSP*, 2020, pp. 6394–6398.
- [11] D. Wang, Z. Chen, and T. Yoshioka, “Neural speech separation using spatially distributed microphones,” in *Proc. Interspeech*, 2020, pp. 339–343.
- [12] S. Zhang and X. Li, “Microphone array generalization for multichannel narrowband deep speech enhancement,” *arXiv:2107.12601*, 2021.
- [13] Y. Yemini, E. Fetaya, H. Maron, and S. Gannot, “Scene-agnostic multi-microphone speech dereverberation,” in *Proc. Interspeech*, 2021, pp. 1129–1133.
- [14] Y. Hu, Y. Liu, S. Lv, M. Xing, S. min Zhang, Y. Fu, J. Wu, B. Zhang, and L. Xie, “DCCRN: Deep complex convolution recurrent network for phase-aware speech enhancement,” in *Proc. Interspeech*, 2020, pp. 2472–2476.
- [15] H.-S. Choi, J.-H. Kim, J. Huh, A. Kim, J.-W. Ha, and K. Lee, “Phase-aware speech enhancement with deep complex U-Net,” in *Proc. ICLR*, 2018.
- [16] D. S. Williamson, Y. Wang, and D. Wang, “Complex ratio masking for monaural speech separation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 483–492, 2016.
- [17] S. E. Eskimez, X. Wang, M. Tang, H. Yang, Z. Zhu, Z. Chen, H. Wang, and T. Yoshioka, “Human listening and live captioning: Multi-task training for speech enhancement,” in *Proc. Interspeech*, 2021, pp. 2686–2690.
- [18] T. Zhou, Y. Zhao, and J. Wu, “Resnext and res2net structures for speaker verification,” in *Proc. SLT*, 2021, pp. 301–307.
- [19] Z.-Q. Wang, P. Wang, and D. Wang, “Multi-microphone complex spectral mapping for utterance-wise and continuous speech separation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 2001–2014, 2021.
- [20] V. Chiley, I. Sharapov, A. Kosson, U. Koster, R. Reece, S. Samaniego de la Fuente, V. Subbiah, and M. James, “On-line normalization for training neural networks,” in *Proc. NIPS*, 2019, pp. 8433–8443.
- [21] C. K. Reddy, H. Dubey, K. Koishida, A. Nair, V. Gopal, R. Cutler, S. Braun, H. Gamper, R. Aichner, and S. Srinivasan, “Interspeech 2021 deep noise suppression challenge,” *arXiv:2101.01902*, 2021.
- [22] 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 *Proc. ICASSP*, 2017, pp. 776–780.
- [23] E. Fonseca, J. Pons, X. Favory, F. Font, D. Bogdanov, A. Ferraro, S. Oramas, A. Porter, and X. Serra, “Freesound datasets: A platform for the creation of open audio datasets,” in *Proc. ISMIR*, 2017, pp. 486–493.
- [24] H. Gamper, C. K. A. Reddy, R. Cutler, I. J. Tashev, and J. Gehrke, “Intrusive and non-intrusive perceptual speech quality assessment using a convolutional neural network,” in *Proc. IEEE WASPAA*, 2019, pp. 85–89.
- [25] I. Mccowan, G. Lathoud, M. Lincoln, A. Lisowska, W. Post, D. Reidsma, and P. Wellner, “The ami meeting corpus,” in *Proc. Measuring Behavior*, 2005, pp. 100–104.
- [26] H. Sato, T. Ochiai, M. Delcroix, K. Kinoshita, T. Moriya, and N. Kamo, “Should we always separate?: Switching between enhanced and observed signals for overlapping speech recognition,” *Proc. Interspeech*, 2021.
- [27] M. Wölfel and J. McDonough, *Distant speech recognition*. John Wiley & Sons, 2009.