

DUAL-PATH TRANSFORMER BASED NEURAL BEAMFORMER FOR TARGET SPEECH EXTRACTION

Aoqi Guo^{*†} Sichong Qian[†] Baoxiang Li[†] Dazhi Gao^{*}

^{*} Ocean University of China, Qingdao, China

[†] SenseTime Research, Beijing, China

ABSTRACT

Neural beamformers, which integrate both pre-separation and beamforming modules, have demonstrated impressive effectiveness in target speech extraction. Nevertheless, the performance of these beamformers is inherently limited by the predictive accuracy of the pre-separation module. In this paper, we introduce a neural beamformer supported by a dual-path transformer. Initially, we employ the cross-attention mechanism in the time domain to extract crucial spatial information related to beamforming from the noisy covariance matrix. Subsequently, in the frequency domain, the self-attention mechanism is employed to enhance the model's ability to process frequency-specific details. By design, our model circumvents the influence of pre-separation modules, delivering performance in a more comprehensive end-to-end manner. Experimental results reveal that our model not only outperforms contemporary leading neural beamforming algorithms in separation performance but also achieves this with a significant reduction in parameter count.

Index Terms— Speech separation, microphone array, beamforming, deep learning, attention

1. INTRODUCTION

Neural beamformers have demonstrated exceptional capabilities in the field of multi-channel target speech extraction [1]. One of the pioneering approaches that combines deep learning with traditional beamforming algorithms is the mask-based neural beamformer [2]. This architecture typically consists of a pre-separation module, followed by a beamforming module. Initially, the pre-separation module generates a series of time-frequency masks [3–6]. These masks are then used to calculate the covariance matrices for either the speech or noise signals, as required by the beamforming algorithm. Subsequently, these covariance matrices are fed into classical beamforming algorithms to calculate the beamforming weights. Given that traditional beamforming algorithms are typically deterministic and operate independently of the pre-separation module, any fluctuation in the output accuracy of the pre-separation module can have a substantial impact on the system's performance. To address this, newer models such as ADL-MVDR [7], GRNNBF [8] and SARNN [9] started integrating neural networks into the beamforming module. This enhances the synergy between the pre-separation and beamforming modules, thereby improving the overall performance of speech separation algorithms. However, the predictive accuracy of the pre-separation module remains a limiting factor in the overall effectiveness of the system.

Recently, there have been further developments in the integration of neural networks and beamforming algorithms. Researchers have developed models that utilize neural networks to directly characterize signals and predict beamforming weights in either the time or frequency domain. For instance, EABNet [10] directly models the time-frequency characteristics of signals captured by microphone arrays, aiming to capture more comprehensive information than mere covariance matrices in order to predict beamforming weights directly. On a parallel note, Yi Luo et al. introduced FasNet-TAC [11], which performs filtering and summation operations in the time domain. By utilizing a neural network to model input features, the model's coherence is improved. Nevertheless, it is important to note that the fundamental concept of beamforming revolves around spatial domain signal filtration. The way in which these models handle spatial signal information remains less transparent. As a result, they exhibit reduced interpretability and robustness when compared to neural beamformers that are based on spatial information obtained from covariance matrices.

In this paper, we introduce a neural beamformer that utilizes a dual-path transformer. Initially, we model both the spatial features of the input and the noisy covariance matrix. We then implement a cross-attention mechanism [12] at the narrowband level to extract spatial information that is relevant to beamforming from the noisy covariance matrix, utilizing spatial features as cues. Following this, we employ a self-attention mechanism at the broadband level, enhancing the model's capability to capture inter-frequency relationships. Ultimately, we derive beamforming weights directly from the spatial information that we have modeled. Our approach avoids the estimation of intermediate variables, and facilitates a more end-to-end prediction of beamforming weights. Most notably, our model maintains a level of interpretability while significantly reducing the total number of parameters.

The remainder of this paper is organized as follows. Section 2 presents the signal model and the theoretical basis for our proposed model. Section 3 details our proposed neural beamforming algorithm. Section 4 provides an overview of the experimental setup and presents an analysis of the experimental results. Finally, Section 5 concludes the paper.

2. SIGNAL MODELS AND METHODS

Consider the far-field frequency-domain signal model in the real scene, described as

$$Y(t, f) = X(t, f) + S(t, f) + N(t, f) \quad (1)$$

where $Y(t, f) = [Y^{(0)}(t, f), Y^{(1)}(t, f), \dots, Y^{(M-1)}(t, f)]^T$ indicates the frequency-domain signal received by the M -channel microphone array. $X(t, f)$, $S(t, f)$ are the reverberated speech signals

Work conducted when the first author was intern at SenseTime Research. Dazhi Gao is the corresponding author.

of the target speaker and the interference speaker respectively, and $N(t, f)$ represents the background noise. When not focusing on the dereverberation task, the task goal of multi-channel target speech extraction is to extract the monaural speech $X'(t, f)$ of the target speaker from the noisy signal $Y(t, f)$.

The purpose of the beamforming algorithm is to obtain the filter weight w for the array observation signal, and extract the desired signal by employing spatial filtering, that is:

$$X'(t, f) = w^H \cdot Y(t, f) \quad (2)$$

where $(\cdot)^H$ is the hermitian transpose, and " \cdot " denotes the matrix multiplication. Typically, beamforming algorithms calculate the covariance matrix for speech or noise signals to determine beamforming weights. Hence, the prediction accuracy of the covariance matrix has a great influence on the overall performance.

Generally, considering the covariance matrix of the noisy signal:

$$\Phi_{YY}(t, f) = Y(t, f) \cdot Y^H(t, f) \in \mathbb{C}^{M \times M} \quad (3)$$

when the speech and interference noise signals are independent of each other and the influence of background noise is ignored, we can get:

$$\Phi_{YY}(t, f) \approx \Phi_{XX}(t, f) + \Phi_{NN}(t, f) \quad (4)$$

It is evident that the covariance matrix of the noisy signal encompasses all the information from the covariance matrices of both speech and noise signals. With this understanding, we employ a neural network to model the noisy covariance matrix and input features, thereby predicting the beamforming weights, described as:

$$w(t, f) = \text{Model}(\text{Features}(t, f), \Phi_{YY}(t, f)) \quad (5)$$

3. OUR PROPOSED NEURAL BEAMFORMER

Fig. 1 illustrates the overall structure of our proposed model. Initially, the model computes the input spatial features as well as the noisy covariance matrix. These are then transformed to a uniform dimensionality via a DNN-RNN architecture. Finally, the beamforming weights are predicted using a dual-path transformer. The specific steps are as follows:

3.1. Feature Extraction

We select the first channel of the microphone array as the reference channel and compute the magnitude spectrum to serve as the input feature.

$$Y_{mag}^{(0)}(t, f) = |Y^{(0)}(t, f)| \quad (6)$$

We then choose P pairs of microphones to calculate the phase difference between each pair, which serves as a spatial feature.

$$\cos IPD_p(t, f) = \cos(\angle Y_p^1(t, f) - \angle Y_p^0(t, f)) \quad (7)$$

while $\angle(\cdot)$ outputs the angle of the input argument. $Y_p^0(t, f)$ and $Y_p^1(t, f)$ represent the spectrum of the signal received by each microphone in the p -th microphone pair, respectively. Ultimately, we compute the angle feature [13], given that the Direction of Arrival(DOA) of the target speaker is known.

$$AF_\theta(t, f) = \sum_{p=1}^P \cos(IPD_p(t, f) - \Delta_{\theta,p}(t, f)) \quad (8)$$

where $\Delta_{\theta,p}(t, f)$ represents the ground truth phase difference given the direction of arrival θ and the p -th microphone pair, and for

a speaker at a fixed position, its characteristics remain consistent across all time frames. All these features are stacked along the channel dimension to serve as the model's input features, that is $\text{Features}(t, f) \in \mathbb{R}^{P+1+1}$. After computing the complex-valued covariance matrix of the noisy signal through Eq. (3), we concatenate the real and imaginary components along the channel dimension, that is

$$\Phi(t, f) = [\Phi_{YY}^r(t, f), \Phi_{YY}^i(t, f)] \in \mathbb{R}^{M \times M \times 2} \quad (9)$$

where $\Phi_{YY}^r(t, f)$ and $\Phi_{YY}^i(t, f)$ represent the real and imaginary components, respectively.

3.2. Dimensional modeling

We first process each frequency point independently at the narrow-band level and then transform the input features and noisy covariance matrix to the same dimension through two different Conv1D layers:

$$E_{Feat}(t, f) = \text{Conv1D}_1(\text{Features}(t, f)) \in \mathbb{R}^D \quad (10)$$

$$E_\Phi(t, f) = \text{Conv1D}_2(\Phi(t, f)) \in \mathbb{R}^D \quad (11)$$

where D represents the embedding dimensions. Then they are concatenated along the channel dimension, and then input to the GRU to enhance the modeling and mapping capabilities of the network.

$$E_{mix}(t, f) = [E_{Feat}(t, f), E_\Phi(t, f)] \in \mathbb{R}^{D \times 2} \quad (12)$$

$$E'_{mix}(t, f) = \text{GRU}(E_{mix}(t, f)) \in \mathbb{R}^{D \times 2} \quad (13)$$

After splitting the modeled $E'_{mix}(t, f)$, we get $E'_{Feat}(t, f) \in \mathbb{R}^D$ and $E'_\Phi(t, f) \in \mathbb{R}^D$, which represents the spatial information of the input features and the noisy covariance matrix, respectively.

$$[E'_{Feat}(t, f), E'_\Phi(t, f)] = \text{Chunk}(E'_{mix}(t, f)) \quad (14)$$

3.3. Dual-path Transformer Based Module

Initially, we leverage a cross-attention module in the time domain to extract beamforming-related information. We use the modeled input features, $E'_{Feat}(t, f)$, as the query, and the noisy covariance matrix, $E'_\Phi(t, f)$, as both key and value for performing the operation of the cross-attention mechanism. This approach enables us to focus on beamforming-relevant spatial information.

$$E_w(t, f) = \text{MHCA}(E'_{Feat}(t, f), E'_\Phi(t, f), E'_\Phi(t, f)) \quad (15)$$

The prior operations of the model are executed independently at each frequency. To leverage inter-frequency information more effectively, we transform the dimension of the embedding and employ a self-attention module in the frequency domain, and apply it frame by frame at the broadband level.

$$E'_w(t, f) = \text{MHSA}(E_w(t, f), E_w(t, f), E_w(t, f)) \quad (16)$$

Following these steps, the model has extracted beamforming-related information from the noisy covariance matrix. Ultimately, we employ a Conv1D layer to predict the real and imaginary parts of the beamforming weights, thereby enabling the application of beamforming to the array observation signals.

$$w(t, f) = \text{Conv1D}_3(E'_w(t, f)) \quad (17)$$

The time-domain speech signal is restored after the processed frequency spectrum is subjected to Inverse Short-Time Fourier Transform(ISTFT). Note that only the target speech extraction is considered and dereverberation is not addressed in this paper.

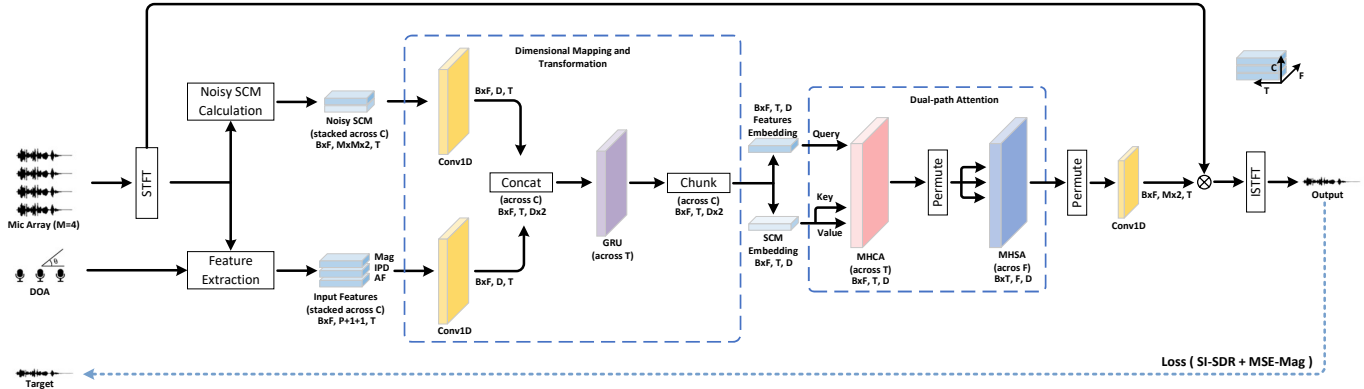


Fig. 1. The overall structure of our proposed model. MHCA is a time-domain multi-head cross-attention module, and MHSA is a frequency-domain multi-head self-attention module. The MHCA and Conv1D layers process each frequency independently at the narrowband level. The MHSA layer models frequency-domain information frame by frame at wideband level. B, C, F, T, D and M represent batch size, channel dimensions, frequency dimensions, time dimensions, embedding dimensions and channels, respectively.

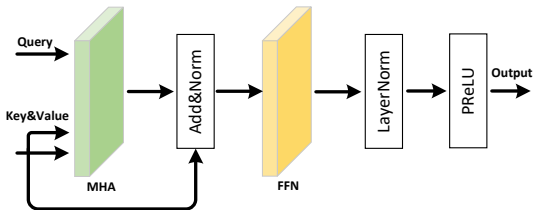


Fig. 2. Structural diagram of the attention module. At the narrowband level, the spatial information is used as Query, and the covariance matrix is used as Key and Value. At the broadband level, Query, Key, and Value are all embedding itself.

4. DATASET AND EXPERIMENTAL SETUP

4.1. Dataset

The source speech signals for target and interfering speakers are obtained from AISHELL-1 [14], with background noise data sourced from DNS2023 [15]. For indoor reverberation scenes, we randomly set room dimensions between [3, 3, 1.5] and [8, 8, 2.5] meters, with a reverberation time (rt60) of 0.1 to 0.6 seconds. The simulation uses the image source method, producing room impulse responses for target and interference signals [16]. The microphone array consists of a four-element linear array with 3 cm spacing. To maintain spatial independence, we set a minimum 5° angle between the two sources relative to the microphone array. We adjust the Signal-Interference Ratio (SIR) between target and interference signals from [-6, 6] dB, and add background noise from [-5, 20] dB to enhance model robustness. During training, all data is divided into 4 seconds. This process yields approximately 133.3 hours of training data (120,000 pieces), 15.6 hours of validation data (14,000 pieces), and 7.8 hours of test data (7,000 pieces), all down-sampled to 16kHz.

4.2. Implementation Details

During the training process, a 512-point STFT is utilized to extract audio features using a 32ms Hann window with a 50% over-

lap. We select three microphone pairs ($P=3$), specifically (0,1), (0,2) and (0,3), to calculate spatial features. The unidirectional GRU layer includes 256 hidden layer units, while the cross-attention and self-attention module dimensions are set at 128 ($D=128$). The final linear layer predicts complex-valued beamforming weights, hence the output dimension is configured to $8(M \times 2)$. More details of the model are provided here.¹

The network undergoes 60 epochs of training with a batch size of 20. We use the Adam optimizer with an initial learning rate of $2e-3$, decaying exponentially at 0.98 per epoch. A maximum gradient clipping of 10 accelerates network convergence.

We use the MIMO Conv-TasNet [17], IRM MVDR [2] and GRNNBF [8] models as our experimental baselines for comparison. In line with the settings from the original paper, we employ a 3x8 Temporal Convolutional Network (TCN) [18] Block for the MIMO Conv-TasNet to model input features, which predicts the Complex-Ratio-Mask (CRM) [5] to restore the spectrum of the reference channel. Concurrently, using the same configuration as MIMO Conv-TasNet, we predict the Ideal-Ratio-Mask (IRM) [4] for both speech and noise signals and integrate it into the mask-based MVDR for computation. For the GRNNBF, a 4x8 TCN block is used to estimate the Complex-Ratio-Filter (CRF) [19]. This prediction aids in calculating the covariance matrices for both speech and noise. This matrix is subsequently fed into a GRU layer comprising 500 hidden units to predict the beamforming weights. For our loss function, we combine the Si-SDR [20] of the time-domain signal with the MSE of the magnitude spectrum, assigning them equal weights to derive the composite loss function.

4.3. Experiment Results

Table 1 presents a comparison between our proposed DPTBF model, the baseline model, and the results from ablation experiments. We evaluate the performance of the model using PESQ [21], STOI [22], Si-SDR, and measure the Word Error Rate (WER) using the ASR model from a specific source [23]. The IRM MVDR model experiences a significant performance drop under conditions of background noise and strong reverberation due to the inherent residual

¹<https://github.com/Aworselife/DPTBF>

Table 1. PESQ, STOI, Si-SDR and WER of several baselines and the proposed DPTBF model.

Systems	GMACs (per sec.)	Para.(M)	Simulated Data			
			PESQ↑	STOI↑	Si-SDR↑	WER(%)↓
Reverberant Clean	—	—	4.5	1.0	∞	2.25
No processing	—	—	1.148	0.563	-1.76	68.02
IRM MVDR [2]	0.35	5.33	1.586	0.757	5.25	25.13
MISO Conv-TasNet [17]	0.33	5.4	1.566	0.759	5.96	28.15
GRNNBF [8]	50.17	15.73	2.176	0.845	8.42	13.15
DPTBF (proposed)	13.52	0.96	2.313	0.861	9.34	9.45
DPTBF (less)	3.44	0.24	2.244	0.855	8.96	12.52
-FA	13.24	0.88	2.096	0.83	8.25	14.33
+SC	13.52	0.96	2.242	0.854	8.96	12.44

"FA" means Frequency-domain Attention module and "SC" means Skip Connection for GRU.

noise issue in MVDR. The MISO Conv-TasNet, which predicts the target speech mask using TCN, eliminates beamforming operations and reduces computational load. However, it does come with a large number of parameters due to the use of multi-layer stacked TCN blocks. Building on the stacked TCN block, GRNNBF integrates an RNN for covariance matrix modelling, leading to a considerable increase in the parameter count. In terms of WER, MISO Conv-TasNet delivers the worst performance, primarily due to the significant spectral distortion it induces. Although the no-distortion constraint of IRM MVDR ensures some level of performance, it does not fully achieve the desired effect due to notable residual noise. In contrast, our proposed DPTBF significantly reduces parameters and computations compared to the baseline models, while also enhancing performance.

We further refined our model into a more streamlined version of DPTBF by reducing GRU hidden layer units from 256 to 128. Even though there is a performance decrease with fewer parameters, the experimental results remain encouraging. Ablation experiments on DPTBF show that using a frequency-domain self-attention module enhances the model's capability to capture frequency-domain information, thereby improving the performance of separation. However, adding skip connections to the GRU layer, which models spatial information and noisy covariance matrix information, did not improve the performance. This could potentially be due to the network's shallow depth and minimal loss of information.

4.4. Spectrogram Analysis

We selected a challenging audio sample from the test set to evaluate the performance of different models in challenging conditions. This sample had a low signal-to-noise ratio and a small angle between the two speakers. The noise signal received by the reference channel is depicted in Fig. 3(a), while the reverberated speech of the target speaker is illustrated in Fig. 3(f). Both the baseline and DPTBF models processed this data, with their results presented in Fig. 3(b-e). The IRM MVDR model, limited by its inherent algorithm, produces speech output with substantial residual noise. The MISO Conv-TasNet model, which utilizes TCN for direct mask prediction, shows robust noise suppression but suffers from significant spectral distortion, which will seriously affect the recognition results of the ASR model. On the other hand, the GRNNBF model has less spectral distortion, but it still struggles with residual noise. In comparison to these baseline models, our proposed DPTBF model excels in reducing interfering noise, mitigating background noise, and im-

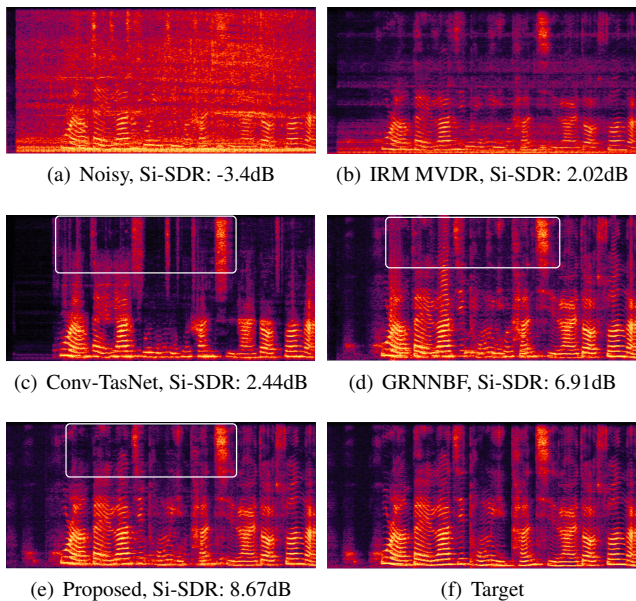


Fig. 3. Spectrum of a piece of data in the test set after processing. SIR=-2.2dB, SNR=-4.8dB. The angle between the speaker is 9°.

proving spectral distortion.

5. CONCLUSIONS

In summary, we propose a neural beamformer based on a dual-path transformer architecture. By incorporating both cross-attention and self-attention mechanisms, our model efficiently eliminates the need to estimate intermediate variables and overcomes the limitations of pre-separation modules. Empirical results underscore the superior performance of the model in target speech extraction tasks, further confirming the effectiveness of the attention mechanism in extracting spatial information relevant to beamforming from the covariance matrix. Looking ahead, our goal is to reduce the computational complexity of the model without compromising separation effectiveness, while also expanding the model to support the multiple-input multiple-output paradigm for additional speakers.

6. REFERENCES

- [1] Katerina Zmolikova, Marc Delcroix, Tsubasa Ochiai, Keisuke Kinoshita, Jan Černocký, and Dong Yu, “Neural target speech extraction: An overview,” *IEEE Signal Processing Magazine*, vol. 40, no. 3, pp. 8–29, 2023.
- [2] Jahn Heymann, Lukas Drude, and Reinhold Haeb-Umbach, “Neural network based spectral mask estimation for acoustic beamforming,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 196–200.
- [3] Hakan Erdogan, John R Hershey, Shinji Watanabe, Michael I Mandel, and Jonathan Le Roux, “Improved mvdr beamforming using single-channel mask prediction networks.,” in *Interspeech*, 2016, pp. 1981–1985.
- [4] Takuya Higuchi, Nobutaka Ito, Takuya Yoshioka, and Tomohiro Nakatani, “Robust mvdr beamforming using time-frequency masks for online/offline asr in noise,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 5210–5214.
- [5] Donald S Williamson, Yuxuan Wang, and DeLiang Wang, “Complex ratio masking for monaural speech separation,” *IEEE/ACM transactions on audio, speech, and language processing*, vol. 24, no. 3, pp. 483–492, 2015.
- [6] Wolfgang Mack and Emanuël AP Habets, “Deep filtering: Signal extraction and reconstruction using complex time-frequency filters,” *IEEE Signal Processing Letters*, vol. 27, pp. 61–65, 2019.
- [7] Zhuohuang Zhang, Yong Xu, Meng Yu, Shi-Xiong Zhang, Lianwu Chen, and Dong Yu, “Adl-mvdr: All deep learning mvdr beamformer for target speech separation,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6089–6093.
- [8] Yong Xu, Zhuohuang Zhang, Meng Yu, Shi-Xiong Zhang, and Dong Yu, “Generalized spatial-temporal rnn beamformer for target speech separation.,” in *Interspeech*, 2021, pp. 3076–3080.
- [9] Xiyun Li, Yong Xu, Meng Yu, Shi-Xiong Zhang, Jiaming Xu, Bo Xu, and Dong Yu, “Mimo self-attentive rnn beamformer for multi-speaker speech separation,” in *Interspeech*, 2021, pp. 1119–1123.
- [10] Andong Li, Wenzhe Liu, Chengshi Zheng, and Xiaodong Li, “Embedding and beamforming: All-neural causal beamformer for multichannel speech enhancement,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 6487–6491.
- [11] Yi Luo, Zhuo Chen, Nima Mesgarani, and Takuya Yoshioka, “End-to-end microphone permutation and number invariant multi-channel speech separation,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6394–6398.
- [12] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [13] Zhuo Chen, Xiong Xiao, Takuya Yoshioka, Hakan Erdogan, Jinyu Li, and Yifan Gong, “Multi-channel overlapped speech recognition with location guided speech extraction network,” in *2018 IEEE Spoken Language Technology Workshop (SLT)*, 2018, pp. 558–565.
- [14] Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng, “Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline,” in *2017 20th conference of the oriental chapter of the international coordinating committee on speech databases and speech I/O systems and assessment (O-COCOSDA)*. IEEE, 2017, pp. 1–5.
- [15] Harishchandra Dubey, Ashkan Aazami, Vishak Gopal, Babak Naderi, Sebastian Braun, Ross Cutler, Alex Ju, Mehdi Zohourian, Min Tang, Hannes Gamper, et al., “Icassp 2023 deep speech enhancement challenge,” *arXiv preprint arXiv:2303.11510*, 2023.
- [16] Robin Scheibler, Eric Bezzam, and Ivan Dokmanić, “Py-roomacoustics: A python package for audio room simulation and array processing algorithms,” in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2018, pp. 351–355.
- [17] Rongzhi Gu, Shi-Xiong Zhang, Yong Xu, Lianwu Chen, Yuxian Zou, and Dong Yu, “Multi-modal multi-channel target speech separation,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 3, pp. 530–541, 2020.
- [18] Shaojie Bai, J Zico Kolter, and Vladlen Koltun, “An empirical evaluation of generic convolutional and recurrent networks for sequence modeling,” *arXiv preprint arXiv:1803.01271*, 2018.
- [19] Hendrik Schroter, Alberto N Escalante-B, Tobias Rosenkranz, and Andreas Maier, “Deepfilternet: A low complexity speech enhancement framework for full-band audio based on deep filtering,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 7407–7411.
- [20] Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R. Hershey, “Sdr – half-baked or well done?,” in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 626–630.
- [21] A.W. Rix, J.G. Beerends, M.P. Hollier, and A.P. Hekstra, “Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs,” in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221)*, 2001, vol. 2, pp. 749–752 vol.2.
- [22] Cees H. Taal, Richard C. Hendriks, Richard Heusdens, and Jesper Jensen, “A short-time objective intelligibility measure for time-frequency weighted noisy speech,” in *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2010, pp. 4214–4217.
- [23] Zhifu Gao, Shiliang Zhang, Ian McLoughlin, and Zhijie Yan, “Paraformer: Fast and accurate parallel transformer for non-autoregressive end-to-end speech recognition,” in *INTER-SPEECH*, 2022.