

MOSSFORMER: PUSHING THE PERFORMANCE LIMIT OF MONAURAL SPEECH SEPARATION USING GATED SINGLE-HEAD TRANSFORMER WITH CONVOLUTION-AUGMENTED JOINT SELF-ATTENTIONS

Shengkui Zhao, Bin Ma

Alibaba Group
{shengkui.zhao, b.ma}@alibaba-inc.com

ABSTRACT

Transformer based models have provided significant performance improvements in monaural speech separation. However, there is still a performance gap compared to a recent proposed upper bound. The major limitation of the current dual-path Transformer models is the inefficient modelling of long-range elemental interactions and local feature patterns. In this work, we achieve the upper bound by proposing a gated single-head transformer architecture with convolution-augmented joint self-attentions, named *MossFormer* (Monaural speech separation TransFormer). To effectively solve the indirect elemental interactions across chunks in the dual-path architecture, MossFormer employs a joint local and global self-attention architecture that simultaneously performs a full-computation self-attention on local chunks and a linearised low-cost self-attention over the full sequence. The joint attention enables MossFormer model full-sequence elemental interaction directly. In addition, we employ a powerful attentive gating mechanism with simplified single-head self-attentions. Besides the attentive long-range modelling, we also augment MossFormer with convolutions for the position-wise local pattern modelling. As a consequence, MossFormer significantly outperforms the previous models and achieves the state-of-the-art results on WSJ0-2/3mix and WHAM!/WHAMR! benchmarks. Our model achieves the SI-SDRi upper bound of 21.2 dB on WSJ0-3mix and only 0.3 dB below the upper bound of 23.1 dB on WSJ0-2mix.

Index Terms— speech separation, transformer, attention, convolution, deep learning

1. INTRODUCTION

Monaural speech separation that aims to separate individual source speeches from a single overlapped mixture is a fundamental and important task. Recent end-to-end deep learning speech separation models have seen large performance improvements [1–8]. The time-domain Conv-TasNet [1] modelled on an encoded representation eventually surpasses the time-frequency domain counterparts. DPRNN [2] provides an effective dual-path framework for handling extreme long encoded input sequences by splitting into smaller chunks and processing the intra- and inter-chunk separately. With capability of learning long-term temporal dependency, DPRNN outperforms Conv-TasNet with a big margin. Building on the dual-path architecture, VSUNOS [3] proposes gated RNN modules to further improve the separation performance. However, RNN based models inherently pass history information recurrently through many intermediate states, leading to suboptimal performance. Recently, the Transformer architecture based on self-attention [9] has been successfully integrated into the dual-path speech separation pipeline. Unlike the recurrent learning of RNN, Transformer provide ability

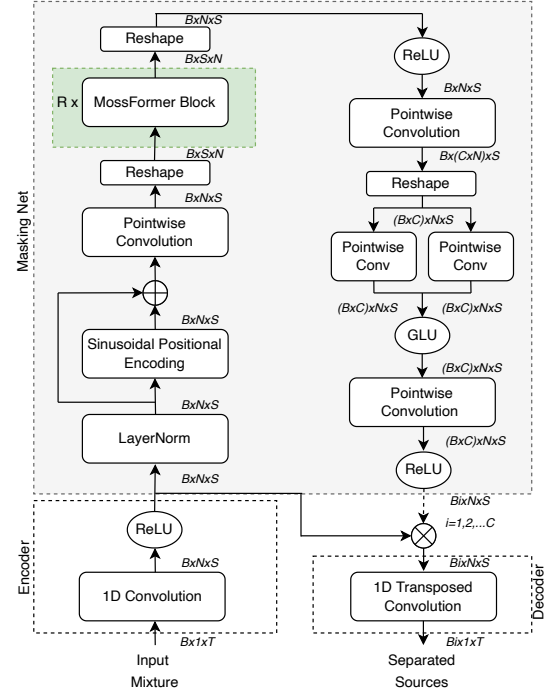


Fig. 1. MossFormer model architecture. MossFormer comprises of a convolutional encoder-decoder structure and a masking net. The masking net is a stack of MossFormer blocks and convolutions.

to capture long-range elemental interactions directly. DPTNet [4] uses an amended Transformer architecture with embedded RNN for preserving sequence positional information and shows superior performance than DPRNN. SepFormer [8] completely eliminates RNN recurrence by building on standard Transformer with multi-head self-attention (MHSA) and achieves the state-of-the-art (SOTA) performance. Due to the quadratic complexity over the input sequence in attention computations, the self-attentions in DPTNet and SepFormer are still limited to short context size. The long-range elemental dependencies across chunks are still implicitly modelled through intermediate states. This fact may impose negative impacts on long-range modelling capability. Compared to the recent Cramer Rao bound for non-linear methods [10], there is still a large performance gap. In addition, convolutions are not well exploited in the existing dual-path Transformer models to learn local feature patterns.

In this work, we propose a novel Monaural speech separation TransFormer (*MossFormer*) model as illustrated in Figure 1. For the dominant MossFormer block, we propose a gated single-head

transformer (GSHT) architecture with convolution-augmented joint self-attentions as illustrated in Figure 2. The GSHT architecture employs a powerful attentive gating mechanism such that only a weakened single-head self-attention (SHSA) is required. This facilitates a joint local and global self-attention for effective long-range direct interaction modelling. To further model position-wise local feature patterns, we propose a convolution module as illustrated in Figure 3 and integrate the convolution module into the attentive gating architecture. Our work is mainly motivated by [11, 12]. Our proposed model outperforms SepFormer and the other previous models, and redefines the state-of-the-art on the WSJ0-2/3mix and WHAM!/WHAMR! benchmarks. Moreover, we achieve the SI-SDRi upper bound [10] on WSJ0-3mix.

2. THE MOSSFORMER MODEL

Given a speech mixture $x = \sum_{i=1}^C s_i$, we aim to estimate C individual sources $s_i \in \mathbb{R}^{1 \times T}$, $i = 1, 2, \dots, C$ based on a deep learning model. Our overall model architecture is built on the time-domain masking-net framework [1] as illustrated in Figure 1. It comprises of a convolutional encoder-decoder structure and a masking net. The encoder-decoder structure responds for feature extraction and waveform reconstruction. The masking net maps the encoded output to a group of masks.

2.1. Encoder and Decoder

The encoder responds for feature extraction and consists of a one-dimensional (1D) convolutional layer (Conv1D) and a rectified linear unit (ReLU), which constrains the encoded output to be non-negative values. Let the kernel size of the encoder be K_1 with stride of $K_1/2$ and the number of filters be N . The input sequence $\mathbf{X} \in \mathbb{R}^{B \times 1 \times T}$ is encoded to the output \mathbf{X}' as follows:

$$\mathbf{X}' = \text{ReLU}(\text{Conv1D}(\mathbf{X})) \quad (1)$$

where $\mathbf{X}' \in \mathbb{R}^{B \times N \times S}$ and $S = 2(T - K_1)/K_1 + 1$. The batch size B is omitted in the follows for ease of presentation. The sequence \mathbf{X}' is multiplied element-wisely by each individual speaker's mask to obtain the separated feature sequence: $\mathbf{X}''_i = \mathbf{M}_i \otimes \mathbf{X}'$. The separated feature sequence is finally decoded into waveform by the decoder:

$$\hat{s}_i = \text{Transposed_Conv1D}(\mathbf{X}''_i) \quad (2)$$

The decoder is a 1D transposed convolutional layer and it uses the same kernel size and stride as the encoder.

2.2. Masking Net

The masking net performs a non-linear mapping from the encoder output to C groups of masks as shown in Figure 1. To achieve this, the encoded sequence \mathbf{X}' is first normalized and added with positional encodings for global order information. And the sequence is then passed through a pointwise convolution and after reshaping passed to the MossFormer block for sequential processing.

In the MossFormer block, the sequence is processed by the convolution modules and the attentive gating mechanism. The convolution modules process the sequence with linear projections and depthwise convolutions. The attentive gating mechanism performs a joint local and global self-attention and gating operations. The MossFormer block learns only the residual and applies the skip connection

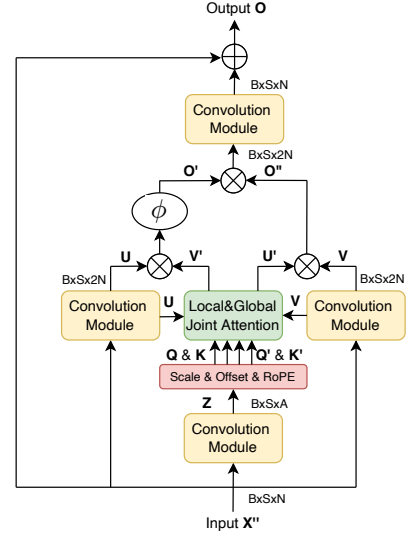


Fig. 2. MossFormer block. It consists of four convolution modules, scale and offset operations, a joint local and global single-head self-attention, and three element-wise gating operations (ϕ is an element-wise activation function, \oplus is element-wise summation, \otimes is element-wise multiplication).

from the input for ease of training. The output of the current MossFormer block is fed as input to the next MossFormer block. The process of the MossFormer block is repeated R times.

The output of the final MossFormer block is processed by a ReLU followed by another pointwise convolution, which expands the sequence dimension from $\mathbb{R}^{N \times S}$ to $\mathbb{R}^{(C \times N) \times S}$. It is then passed through a parallel pointwise convolutions and a GLU. Finally, the sequence is passed through pointwise convolution one more time followed by a ReLU to obtain the mask sequence $\mathbf{M} \in \mathbb{R}^{C \times N \times S}$. The mask sequence \mathbf{M} is reformed for each individual speaker $\mathbf{M}_i \in \mathbb{R}^{N \times S}$ and is then fed to the decoder separately.

2.3. MossFormer Block

The architecture of our proposed MossFormer block is shown in Figure 2, which is developed based on the recent proposed gated attention unit (GAU) [12] for long sequence modelling. A MossFormer block comprises of four convolution modules, scale and offset operations, a joint local and global SHSA, and three gating operations. We aim to boost the modelling capability of the MossFormer block by incorporating convolution modules and a triple-gating structure. The use of gates allows a much simpler SHSA that facilitates a joint local and global attention for effective long-range modelling.

2.3.1. Convolution Module

We propose a convolution module to replace the dense layers in GAU for extracting fine-grained local feature patterns in the MossFormer block. Figure 3 illustrates the architecture of the convolution module. In the convolution module, the sequence is first normalized and projected by a linear layer followed by a SiLU. And then it is feature-wise convoluted by 1D depthwise convolution. Skip-connection and dropout are used to help training and regularizing the network.

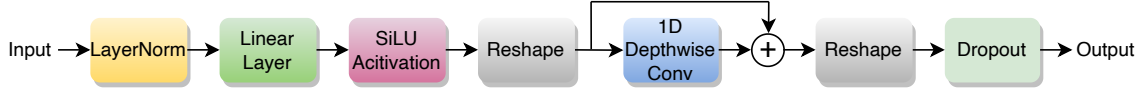


Fig. 3. Convolution module. It contains a linear layer with an expansion factor followed by a SiLU activation layer, and then followed by a 1-D depthwise convolution layer with a skip connection.

2.3.2. Attentive Gating Mechanism

Our attentive gating mechanism combines attention into the triple-gating process to enhance model capability. The mechanism is formulated as follows. Let the input sequence of the current MossFormer block be $\mathbf{X}'' \in \mathbb{R}^{S \times N}$. It is processed by the convolution module to obtain the values $\mathbf{U} \in \mathbb{R}^{S \times 2N}$ and $\mathbf{V} \in \mathbb{R}^{S \times 2N}$ as follows:

$$\mathbf{U} = \text{ConvM}(\mathbf{X}''), \quad \mathbf{V} = \text{ConvM}(\mathbf{X}'') \quad (3)$$

where ConvM refers to the convolution module. Here we increase the feature dimensions from N to $2N$ in the linear layer of ConvM using an expansion factor of 2. By denoting the attention matrix as $\mathbf{A} \in \mathbb{R}^{S \times S}$, the output sequence $\mathbf{O} \in \mathbb{R}^{S \times N}$ of the MossFormer block can be expressed as follows:

$$\mathbf{O}' = \phi(\mathbf{U} \otimes \mathbf{V}') \quad \text{where} \quad \mathbf{V}' = \mathbf{A}\mathbf{V} \quad (4)$$

$$\mathbf{O}'' = \mathbf{U}' \otimes \mathbf{V} \quad \text{where} \quad \mathbf{U}' = \mathbf{A}\mathbf{U} \quad (5)$$

$$\mathbf{O} = \mathbf{X}'' + \text{ConvM}(\mathbf{O}' \otimes \mathbf{O}'') \quad (6)$$

where the linear layer of ConvM decreases the feature dimension from $2N$ to N , and ϕ is an element-wise activation function.

2.3.3. Joint Local and Global Single-Head Self-Attention

For long sequence where S is large, computing the attentions in Eqs. (4) and (5) directly is very expensive. Fortunately, the presence of gating allows us compute \mathbf{V}' and \mathbf{U}' based on the joint local and global attention in an efficient and effective way. We first project the input sequence \mathbf{X}'' through the convolution module into a shared representation: $\mathbf{Z} = \text{ConvM}(\mathbf{X}'') \in \mathbb{R}^{S \times D}$, where $D \ll N$. We then apply low-cost per-dim scalars and offsets, and RoPE [13] to the shared \mathbf{Z} to obtain the queries $\mathbf{Q}, \mathbf{Q}' \in \mathbb{R}^{S \times D}$ and the keys $\mathbf{K}, \mathbf{K}' \in \mathbb{R}^{S \times D}$ for both the local and global attentions. For the global attention, we employ the following low-cost linearised form to capture long-range global interactions for both sequences \mathbf{V} and \mathbf{U} :

$$\mathbf{V}'_{\text{global}} = \mathbf{Q}' \left(\beta \mathbf{K}'^T \mathbf{V} \right), \quad \mathbf{U}'_{\text{global}} = \mathbf{Q}' \left(\beta \mathbf{K}'^T \mathbf{U} \right) \quad (7)$$

where $\beta = 1/S$ is a scaling factor. To compute the local quadratic attention, we chunk $\mathbf{V}, \mathbf{U}, \mathbf{Q}$, and \mathbf{K} into H non-overlapping chunks of size P , where zero-padding is used when $S < H \times P$. The quadratic attention is therefore independently applied to each chunk as follows:

$$\mathbf{V}'_{\text{local},h} = \text{ReLU}^2(\gamma \mathbf{Q}_h \mathbf{K}_h^T) \mathbf{V}_h, \quad \mathbf{U}'_{\text{local},h} = \text{ReLU}^2(\gamma \mathbf{Q}_h \mathbf{K}_h^T) \mathbf{U}_h \quad (8)$$

where $\gamma = 1/P$ is a scaling factor. Here we adopt the squared ReLU instead of the softmax in MHSA for optimizing performance [12]. Note that $\mathbf{Q}_h \mathbf{K}_h^T$ only needs to be computed once as it is shared by $\mathbf{V}'_{\text{local},h}$ and $\mathbf{U}'_{\text{local},h}$. We concatenate all outputs of $\mathbf{V}'_{\text{local},h}$ and $\mathbf{U}'_{\text{local},h}$ along the time dimension to form back the full sequences: $\mathbf{V}'_{\text{local}} = [\mathbf{V}'_{\text{local},1}, \dots, \mathbf{V}'_{\text{local},H}]$ and $\mathbf{U}'_{\text{local}} = [\mathbf{U}'_{\text{local},1}, \dots, \mathbf{U}'_{\text{local},H}]$.

We add the local attention and the global attention together to form the final joint attentions of \mathbf{V}' and \mathbf{U}' in Eqs. (4) and (5):

$$\mathbf{V}' = \mathbf{V}'_{\text{local}} + \mathbf{V}'_{\text{global}}, \quad \mathbf{U}' = \mathbf{U}'_{\text{local}} + \mathbf{U}'_{\text{global}} \quad (9)$$

Table 1. Hyper-parameters for MossFormer S, M, and L models for optimized performance within parameter and GPU resource limits.

Model	MossFormer (S)	MossFormer (M)	MossFormer (L)
No. Parameters	10.8M	25.3M	42.1M
No. MossFormer Blocks (R)	22	25	24
Encoder Output Dimension (N)	256	384	512
Encoder Kernel Size (K_1) / Stride	8/4	16/8	16/8
Depthwise Conv Kernel Size (K_2)	31	17	17
Chunk Size (P)	256	256	256
Attention Dimension (D)	128	128	128
Gating Activation Function (ϕ)	Sigmoid	Sigmoid	Sigmoid

3. EXPERIMENTS

3.1. Dataset

We evaluate the proposed model on both clean and noisy/reverberated settings using the speech separation benchmarks of WSJ0-2/3mix [14] and WHAM!/WHAMR! [15, 16] datasets. We rely on the 8kHz version of the data. The utterances are randomly segmented into 4s long during training and validation. Beside the standard versions of the data, we also consider dynamic mixing (DM) with speed perturbation for data augmentation as described in [8].

3.2. Training Setup

Our evaluations are made for three models with parameters of 11M (S), 25M (M), and 42M (L), respectively. The models are chosen based on our best settings of network depth, model dimensions, convolution kernel sizes, chunk size, attention dimensions, as well as gating activation function ϕ within parameter size and training resource constraints. We use a single NVIDIA V100 GPU with 16 GB of memory for training. Table 1 provides the hyper-parameters.

Our models are implemented based on the SpeechBrain toolkit¹ and optimized using the SI-SDR training loss [17]. We train our models for a maximum of 200 epochs with the Adam optimizer [18] using learning rate of $15e^{-5}$ and batch size of 1. During training, the learning rate holds for 85 epochs and then is reduced by a factor of 0.5 with patience of 2. We limit the l_2 norm of the training gradients to 5 with gradient clipping. The dropout rate is set to 0.1 for all models. We release audio samples online² and source code later.

3.3. Results

We use SI-SDR improvement (SI-SDRi) as evaluation metric. Table 2 reports the results on the clean settings of WSJ0-2mix/3mix datasets and the noisy/reverberated WHAM!/WHAMR! datasets. Our models are compared with the best reported results in the literature. We report results of the small MossFormer(S) model on

¹speechbrain.github.io/

²https://github.com/alibabasglab/MossFormer

Table 2. Performance comparisons on the WSJ0-2mix/3mix and WHAM!/WHAMR! benchmark datasets.

Model	Para.(M)	SI-SDRi	
		WSJ0-2mix/3mix	WHAM!/WHAMR!
TasNet [19]	-	10.8 / -	- / -
Chimera++ [20]	-	11.5 / -	9.9 / -
SignPredictionNet [21]	55.2	15.3 / -	- / -
Conv-TasNet [1]	5.1	15.3 / 12.7	12.7 / 8.3
DeepCASA [5]	12.8	17.7 / -	- / -
Learnable fbank [22]	-	- / -	12.9 / -
Two-Step CTN [23]	8.6	16.1 / -	- / -
MGST [24]	-	17.0 / -	13.1 / -
DPRNN [2]	2.6	18.8 / 14.7	13.9 / 10.3
SuDoRMRF [6]	2.6	18.9 / -	- / -
VSUNOS [3]	7.5	20.1 / 16.9	15.2 / 12.2
DPTNet [4]	2.6	20.2 / -	- / -
Wavesplit [7]	29	21.0 / 17.3	- / -
Wavesplit + DM	29	22.2 / 17.8	16.0 / 13.2
SepFormer [8]	25.7	20.4 / 17.6	- / -
SepFormer + DM	25.7	22.3 / 19.5	16.4 / 14.0
MossFormer(S)	10.8	20.9 / 17.8	- / -
MossFormer(M) + DM	25.3	22.5 / 20.8	17.1 / 15.9
MossFormer(L) + DM	42.1	22.8 / 21.2	17.3 / 16.3
Upper bound [10]	-	23.1 / 21.2	- / -

standard data version and results of the MossFormer(M) and MossFormer(L) models on the augmented data version. For the clean settings, our MossFormer(S) outperforms all the previous models except a very competitive result against Wavesplit on WSJ0-2mix. Note that Wavesplit has 29M parameters, more than 2 times of MossFormer(S), and uses additional speaker identity labels for training. With DM, our MossFormer(M) outperforms all the previous models. On WSJ0-2mix/3mix, MossFormer(M) achieves 22.5 dB and 20.8 dB SI-SNRi compared to 22.3 dB and 19.5 dB for the 26M SepFormer, and 22.2 dB and 17.8 dB for Wavesplit. Our MossFormer(L) makes further performance improvements on top of MossFormer(M) with an increase of the model dimension. Compared to the upper bounds of 23.1 dB and 21.2 dB on WSJ0-2mix/3mix reported in [10], MossFormer(L) achieves 22.8 dB and 21.2 dB on WSJ0-2mix/3mix. Therefore, MossFormer(L) achieves not only an upper bound but also new state-of-the-art results on WSJ0-2mix/3mix.

For the noisy and reverberated settings, Table 2 shows that MossFormer(M) and MossFormer(L) outperform the previous models with big margins and MossFormer(L) achieves new state-of-the-art results on WHAM! and WHAMR!, respectively. For instance, MossFormer(L) achieves 0.9 dB and 2.3 dB more compared to SepFormer. Note that the WHAM!/WHAMR! datasets are built on top of WSJ0-2mix by introducing additional noise and reverberation. Therefore the WHAM!/WHAMR! tasks become harder as the models need to address not only speech separation but also denoising and dereverberation. We observe that the reverberation affects Wavesplit and SepFormer more than MossFormer as their performance drop more from WHAM! to WHAMR!.

3.4. Ablation Studies

We base MossFormer(S) model and WSJ0-2mix to make ablation studies. Table 3 shows the effects of the convolution module, the gating mechanism, and the joint-attention. We observe that the convolution modules has an important impact on the performance. When replacing the convolution modules with the dense layers used in GAU [12] in the values as well as in the queries and keys, it is seen that

Table 3. Ablation studies on the convolution module, the triple-gating mechanism, and the joint attention.

Model	SI-SDRi
MossFormer(S)	20.9
Quadratic Local Attention Only	17.8
Linear Global Attention Only	19.6
Remove output \mathbf{O}' in Eq.(6)	20.4
Replace ConvM with Dense for \mathbf{U} & \mathbf{V}	20.5
Replace ConvM with Dense for \mathbf{Q}_s & \mathbf{K}_s	20.3
Replace ConvM with Dense for both \mathbf{U} & \mathbf{V} and \mathbf{Q}_s & \mathbf{K}_s	19.9

Table 4. Ablation study on different choice of the gating activation function ϕ .

Model	Gating Activation Function ϕ				
	ReLU	GELU	Swish	Bilinear	Sigmoid
MossFormer(S)	20.0	19.9	19.9	20.1	20.9

the performance becomes worse. It also shows that the convolution module affects the queries and keys more than the values. When removing the output \mathbf{O}' in Eq.(6), the triple-gating structure becomes a single-gating structure as in GAU [12] and the result drops from 20.9 dB to 20.4 dB. It shows the effectiveness of the triple-gating design. We also test the model using the quadratic local attention only or the linear global attention only. The results show that none of them performs well individually. It demonstrates the impact of the joint-attention scheme. Another observation is that the global attention alone performs better than the local attention alone, implying the importance of global modelling.

Table 5. Ablation study on settings of attention dimension and chunk size.

Kernel Size K_2			Attention Dim. D			Chunk Size P		
21	31	65	64	128	256	128	256	384
20.6	20.9	20.7	20.5	20.9	20.8	20.5	20.9	20.9

In Table 4, we validate the choice of the activation function ϕ according to the study [25] and it shows that the Sigmoid function works the best. Table 5 studies the effects of kernel size K_2 in the depthwise convolution in ConvM, the attention dimension D , and the chunk size P . We find that the performance improves with larger kernel sizes till 31 but worsens for 65. In addition, $D = 128$ performs the best. A larger chunk size tends to perform better, but further increasing P from 256 to 384 has no more improvement.

4. CONCLUSIONS

In this work, we introduced MossFormer, a Transformer model for monaural speech separation. Contrary to prior models, we learn the local feature patterns and the global long-range dependencies in a unified attentive gating model. Unlike the dual-path framework that models the long-rang interactions implicitly, we employed a quadratic local attention and a low-cost global attention in a joint form to model the long-rang interactions directly. We also combine convolution modules to model the local feature patterns. Our studies demonstrated the importance of each component and achieved much better results than previous models with new state-of-the-art on the benchmarks of WSJ0-2/3mix and WHAM!/WHAMR!.

5. REFERENCES

- [1] Y. Luo and N. Mesgarani, "Conv-TasNet: Surpassing ideal time–frequency magnitude masking for speech separation," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [2] Y. Luo, Z. Chen, and T. Yoshioka, "Dual-Path RNN: Efficient long sequence modeling for time-domain single-channel speech separation," in *Proc. of ICASSP*, 2020.
- [3] E. Nachmani, Y. Adi, and L. Wolf, "Voice separation with an unknown number of multiple speakers," in *Proc. of ICML*, 2020.
- [4] J. Chen, Q. Mao, and D. Liu, "Dual-Path Transformer Network: Direct context-aware modeling for end-to-end monaural speech separation," in *Proc. of Interspeech*, 2020.
- [5] Y. Liu and D. Wang, "Divide and conquer: A deep casa approach to talker-independent monaural speaker separation," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 12, 2019.
- [6] E. Tzinis, Z. Wang, and P. Smaragdis, "Sudo rm -rf: Efficient networks for universal audio source separation," in *Proc. of MLSP*, 2020.
- [7] N. Zeghidour and D. Grangier, "Wavesplit: End-to-end speech separation by speaker clustering," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, 2021.
- [8] C. Subakan, M. Ravanelli, S. Cornell, M. Bronzi, and J. Zhong, "Attention is all you need in speech separation," in *Proc. of ICASSP*, 2021.
- [9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. of NIPS*, 2017.
- [10] S. Lutati, E. Nachmani, and L. Wolf, "SepIt: Approaching a single channel speech separation bound," *arXiv preprint arXiv:2205.11801*, 2022.
- [11] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, and R. Pang, "Conformer: Convolution-augmented transformer for speech recognition," in *Proc. Interspeech*, 2020, pp. 5036–5040.
- [12] W. Hua, Z. Dai, H. Liu, and Q. Le, "Transformer quality in linear time," in *Proc. of ICML*, 2022.
- [13] J. Su, Y. Lu, S. Pan, A. Murtadha, B. Wen, and Y. Liu, "RoFormer: Enhanced transformer with rotary position embedding," *arXiv:2104.09864*, 2021.
- [14] J. R. Hershey, Z. Chen, J. L. Roux, and S. Watanabe, "Deep clustering: Discriminative embeddings for segmentation and separation," in *Proc. of ICASSP*, 2016.
- [15] G. Wichern, J. Antognini, M. Flynn, L. R. Zhu, E. McQuinn, D. Crow, E. Manilow, and J. L. Roux, "Wham!: Extending speech separation to noisy environments," in *Proc. of Interspeech*, 2019.
- [16] M. Maciejewski, G. Wichern, E. McQuinn, and J. L. Roux, "Whamr!: Noisy and reverberant single-channel speech separation," *arXiv:1910.10279*, 2019.
- [17] J. L. Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR – half-baked or well done?," in *Proc. of ICASSP*, 2019.
- [18] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [19] Y. Luo and N. Mesgarani, "TasNet: time-domain audio separation network for real-time, single-channel speech separation," in *Proc. of ICASSP*, 2018.
- [20] Z.-Q. Wang, J. L. Roux, and J. R. Hershey, "Alternative objective functions for deep clustering," in *Proc. of ICASSP*, 2018.
- [21] Z.-Q. Wang, K. Tan, and D. Wang, "Deep learning based phase reconstruction for speaker separation: A trigonometric perspective," in *Proc. of ICASSP*, 2019.
- [22] M. Pariente, S. Cornell, A. Deleforge, and E. Vincent, "Filterbank design for end-to-end speech separation," in *Proc. of ICASSP*, 2020.
- [23] E. Tzinis, S. Venkataramani, Z. Wang, C. Subakan, and P. Smaragdis, "Two-step sound source separation: Training on learned latent targets," in *Proc. of ICASSP*, 2020.
- [24] Y. Zhao, C. Luo, Z.-J. Zha, and W. Zeng, "Multi-scale group transformer for long sequence modeling in speech separation," in *Proc. of IJCAI*, 2020.
- [25] N. Shazeer, "GLU variants improve transformer," *arXiv:2002.05202*, 2020.