# FLASepformer: Efficient Speech Separation with Gated Focused Linear Attention Transformer

*Haoxu Wang, Yiheng Jiang, Gang Qiao, Pengteng Shi, Biao Tian*

Tongyi Lab, Alibaba Group, China

{wanghaoxu.whx,jiangyiheng.jyh,songjiang.qg,pengteng.spt,tianbiao.tb}@alibaba-inc.com

## Abstract

Speech separation always faces the challenge of handling prolonged time sequences. Past methods try to reduce sequence lengths and use the Transformer to capture global information. However, due to the quadratic time complexity of the attention module, memory usage and inference time still increase significantly with longer segments. To tackle this, we introduce Focused Linear Attention and build FLASepformer with linear complexity for efficient speech separation. Inspired by SepReformer and TF-Locoformer, we have two variants: FLA-SepReformer and FLA-TFLocoformer. We also add a new Gated module to improve performance further. Experimental results on various datasets show that FLASepformer matches state-of-the-art performance with less memory consumption and faster inference. FLA-SepReformer-T/B/L increases speed by 2.29x, 1.91x, and 1.49x, with 15.8%, 20.9%, and 31.9% GPU memory usage, proving our model's effectiveness.

**Index Terms**: speech separation, focused linear attention, efficient

## 1. Introduction

Monaural speech separation (SS) extracts individual speech sources from a single-channel mixture, which is crucial for addressing the cocktail party problem [1, 2] and improving speech applications. While previous SS methods achieve good results using neural networks [3–5], SS still faces the challenge of modeling prolonged time sequences, resulting in high computational complexity and slow inference. Past methods use various methods to reduce sequence lengths. TF domain models use Short-Time Fourier transform (STFT) by downsampling FFT bins [6, 7]. Some time domain models like TasNet [8] use fixed receptive fields. Dual-Path Modeling, as introduced by DPRNN [9], chunks the sequences for eariser training with Long-Short Term Memory (LSTM) [10]. DPTNet [11] introduces Transformer [12] to improve long-range modeling. In addition, some efficient speech processing models often use U-Net for downsampling and lightweight designs [13–17].

However, these methods often lose global information through STFT, fixed receptive fields, chunking, downsampling, etc. Even with downsampling, using attention mechanisms still results in quadratic complexity, causing processing time to grow as $O(N^2)$ with longer audio. On the other hand, $O(N)$ models like LSTM struggle to capture global information effectively. Mossformer [18] and Mossformer2 [19] use GAU [20] with approximate linear time complexity to capture global information. Linear attention, on the other hand, is considered a simple and effective alternative by reducing the general complexity. Vanilla Linear Attention [21] (VLA) achieves $O(N)$ complexity by
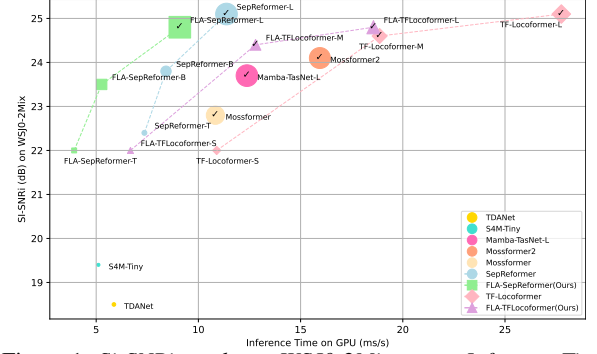


Figure 1: *Si-SNRi results on WSJ0-2Mix versus Inference Time on RTX A800 GPU (ms/s) for 30s audio mixture. The check mark indicates using DM method for training. The radius of circle is proportional to the parameter size of the model.*

computing $K^T V$ first compared to $QK^T$ in Softmax Attention. In SS tasks $d << N$ (channel dimension is much smaller than sequence length), using an efficient linear attention scheme to model long sequences is important.

Therefore, we introduce a new Focused Linear Attention (FLA) [22] to the SS task, replacing commonly used downsampling and quadratic complexity attention mechanisms used in previous SS work [23]. FLA improves VLA by using a simple pull-push mapping method to address overly smooth attention weights. It also uses a depthwise convolution (DWC) module to alleviate the loss of feature diversity from the lower rank of VLA attention weights. We apply FLA to the SS task, modifying the original 2D DWC module to a 1D DWC suitable for speech tasks. We also introduce a Gated module to control the output of FLA and improve model performance. Using the Gated FLA module, we build FLASepformer, primarily based on the latest models, SepReformer [5] and TF-locoformer [7]. We replace the MHSA w d/u in SepReformer and the attention mechanism in TF-locoformer's Temporal Modeling. Ultimately, we propose FLASepformer with two variants, FLA-SepReformer and FLA-TFLocoformer, achieving comparable state-of-the-art (SOTA) performance to SepReformer and TF-Locoformer on the WSJ0-2Mix dataset while greatly reducing inference complexity. As shown in Fig. 1 and Fig. 3, FLA-SepReformer-T/B/L increases speed by 2.29x, 1.91x, and 1.49x for 30s audio mixture, with 15.8%, 20.9%, and 31.9% GPU memory usage, proving our model's effectiveness.

## 2. Methods

### 2.1. Vanilla Linear Attention

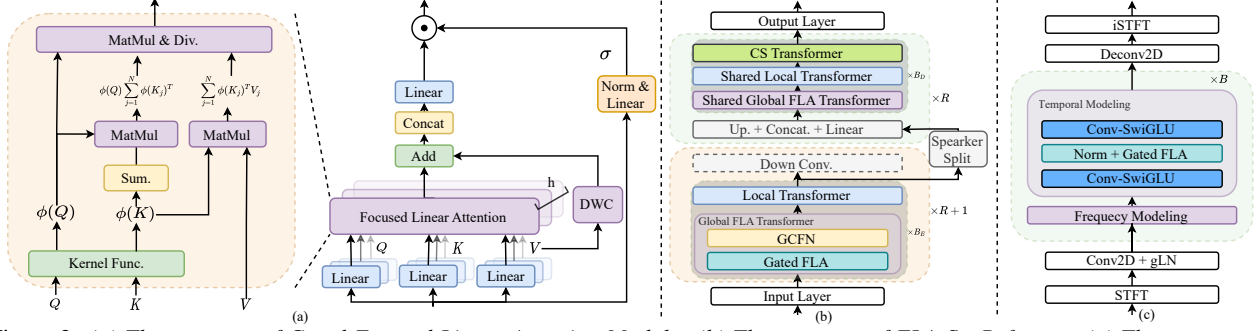Let's first review the original attention module, which can be defined as:

Figure 2: *(a) The structure of Gated Focused Linear Attention Module. (b) The structure of FLA-SepReformer. (c) The structure of FLA-TFLocoformer.*

$$Q = XW^q, K = XW^k, V = XW^v,$$
$$O_i = \sum_{j=1}^{N} \frac{Sim(Q_i, K_j)}{\sum_{j=1}^{N} Sim(Q_i, K_j)} V_j \quad (1)$$

which $X \in R^{N \times d}$ is the input sequence, $W^* \in R^{d \times d}$ is the linear layer, and $i$ is the sequence index. Traditional attention uses Softmax attention as the similarity function $Sim(Q, K) = exp(QK^T/sqrt(d))$ with $O(N^2)$ complexity, which is inefficient for long sequences in speech separation. Linear attention adjusts the similarity calculation to $Sim(Q, K) = \phi(Q)\phi(K)^T$, where $\phi(\cdot)$ is a kernel function, providing a more efficient alternative. The linear attention module is defined as:

$$O_i = \sum_{j=1}^{N} \frac{\phi(Q)\phi(K)^T}{\sum_{j=1}^{N} \phi(Q)\phi(K)^T} V_j \quad (2)$$

by adjusting the similarity function, we can first compute $\phi(K)^T V$ using matrix multiplication law, defined as:

$$O_i = \frac{\phi(Q)(\sum_{j=1}^{N} \phi(K_j)^T V_j)}{\phi(Q) \sum_{j=1}^{N} \phi(K)^T} \quad (3)$$

which reduces the computational complexity from $O(N^2)$ to $O(N)$, allowing for parallel training and linear inference time compared to Softmax Attention.

However, current linear attention mechanisms face a trade-off between accuracy and complexity. Simple kernels like RELU fail to model long-range relationships between $Q$ and $K$, leading to smooth attention weights and degrading performance [24]. Complex kernel functions may increase computational complexity [25]. Additionally, in the SS task, VLA may not effectively capture speech feature diversity in very long sequences.

**2.2. Gated Focused Linear Attention**

The whole architecture of Gated Focused Linear Attention is shown in Fig.2(a). We introduce Focused Linear Attention (FLA) into the SS task to better model long speech sequence features. FLA improves smooth attention weights from VLA. FLA introduces a novel kernel function, Focused Function $\phi_p(\cdot)$, to enhance the model's focus ability on different features. This helps effectively combine important features from long sequences for various queries. The new kernel function, Focused Function $\phi_p(\cdot)$, is defined as:

$$\phi_p(x) = f_p(RELU(x)), f_p = \frac{||x||}{||x^{**p}||} x^{**p} \quad (4)$$

where $x^{**p}$ denotes the element-wise power $p$ of $x$. By designing an appropriate focus factor $p$, FLA pulls similar query-key pairs closer and pushes dissimilar pairs apart. This de-

sign enhances the aggregation of similar speech token features while reducing the gathering of unrelated ones. Consequently, FLA improves VLA to mimic the sharp distribution of Softmax Attention weights, effectively modeling important long-range speech features.

Compared to Softmax Attention, VLA needs to build both sharper attention weights and feature weights with high diversity. According to [22], the maximum rank of VLA's attention weights depends on the number of speech features $N$ and the channel dimension $d$, given by $rank(\phi(Q)\phi(K)^T) \leq min\{rank(\phi(Q)), rank(\phi(K))\} \leq min\{N, d\}$. In SS tasks, $d$ is often much smaller than $N$; for instance, in SepReformer-B [23], $d = 16$ and $N = 2000$ for 1 second of audio. This leads to significantly less feature diversity in VLA than in Softmax Attention. The original FLA uses a 2D depthwise convolution (DWC2d) module to enhance attention weight rank, defined as $O = \phi(Q)\phi(K)^T V + DWC2d(V)$. While DWC2d targets local information in images, for speech separation, we switch it to DWC1d with kernel size $k$ for one-dimensional speech features. The DWC1d module improves attention weight diversity and, with a full-rank attention matrix similar to Softmax Attention, compensates for the rank lost in $\phi(Q)\phi(K)^T$, matching the original upper bound. Additionally, it also acts as a local attention mechanism, enhancing information about adjacent features in different heads within the time domain.

Compared to the VLA, FLA modifies only the kernel function to the Focused Function and adds a fixed-cost DWC1d module, maintaining $O(N)$ complexity. This enables the modeling of long-range speech features and builds an efficient speech separation model while reducing inference complexity to linear compared to the $O(N^2)$ complexity of Softmax Attention.

To enhance FLA's modeling capability further, we introduce a gating mechanism inspired by [20]. We use a Gated MultiLayer Perceptron (MLP), consisting of a Norm layer, a linear layer, and an activation function to obtain the gating result, which is then multiplied by FLA's output. This Gated FLA module uses current token features to control global feature interaction better.

**2.3. The architecture of the FLASepformer**

We build FLASepformer, with two variants: FLA-SepReformer and FLA-TFLocoformer, based on SepReformer [5] and TF-Locoformer [7], respectively. In Fig.2(b), for SepReformer, we replace the original efficient global attention (GLA) in all global transformers with the Gated FLA from Sec. 2.2. Unlike GLA, which uses downsampled MHSA to reduce the sequence length to $\frac{N}{2^R}$ while still maintaining $O(N^2)$ complexity, Gated FLA can model long-range features without downsampling and re-

Table 1: *Comparisons with other methods on WSJ0-2Mix. Results in [dB]. ∗ represents the results of our reproduced model.*

| Methods | Param [M] | MACs [G] | w/o DM SI-SNRi | w/o DM SDRi | w/ DM SI-SNRi | w/ DM SDRi |
|---|---|---|---|---|---|---|
| Conv-TasNet [8] | 5.1 | 5.1 | 15.3 | 15.6 | - | - |
| SuDoRM-RF [13] | 6.4 | - | 18.9 | - | - | - |
| TDANet [15] | 2.3 | - | 18.5 | 18.7 | - | - |
| Sandglasset [14] | 2.3 | - | 20.8 | 21.0 | - | - |
| S4M-Tiny [26] | 1.8 | - | 19.4 | 19.7 | - | - |
| S4M [26] | 3.6 | - | 20.5 | 20.7 | - | - |
| DPRNN [9] | 2.6 | 42.2 | 18.8 | 19.0 | - | - |
| DPTNet [11] | 2.7 | - | 20.2 | 20.6 | - | - |
| SepFormer [5] | 25.7 | 59.5 | 20.4 | 20.5 | 22.3 | 22.4 |
| TF-GridNet [6] | 14.4 | 231.1 | 23.5 | 23.6 | - | - |
| Mamba-TasNet(L) [27] | 59.6 | - | - | - | 23.7 | 23.8 |
| MossFormer [18] | 42.1 | 42.7* | - | - | 22.8 | - |
| MossFormer2 [19] | 55.7 | 56.4* | - | - | 24.1 | - |
| SepReformer-T [23] | 3.7 | 5.2* | 22.4 | 22.6 | - | - |
| SepReformer-B [23] | 14.2 | 19.9* | 23.8 | 23.9 | - | - |
| SepReformer-L [23] | 59.4 | 77.7* | - | - | 25.1 | 25.2 |
| SepReformer-B (Rep.) | 14.2 | 19.9* | 23.6 | 23.7 | - | - |
| FLA-SepReformer-T | 3.7 | 5.6 | 22.0 | 22.1 | - | - |
| FLA-SepReformer-B | 14.2 | 21.6 | 23.5 | 23.7 | - | - |
| FLA-SepReformer-L | 59.4 | 84.6 | - | - | 24.7 | 24.8 |
| TF-Locoformer-S [7] | 5.0 | 43.7 | 22.0 | 22.1 | 22.8 | 23.0 |
| TF-Locoformer-M [7] | 15.0 | 127.8 | 23.6 | 23.8 | 24.6 | 24.7 |
| TF-Locoformer-L [7] | 22.5 | 191.7 | 24.2 | 24.3 | 25.1 | 25.2 |
| FLA-TFLocoformer-S | 5.2 | 44.0 | 22.1 | 22.3 | 22.8 | 22.9 |
| FLA-TFLocoformer-M | 15.1 | 128.7 | 23.4 | 23.5 | 24.4 | 24.5 |
| FLA-TFLocoformer-L | 22.6 | 193.0 | 24.2 | 24.3 | 24.8 | 24.9 |

Table 2: *Comparisons with others methods on various dataset. DM is not used on all Libri2Mix-100 results. Results in [dB].*

| Methods | WHAM! SI-SNRi /SDRi | WHAMR! SI-SNRi /SDRi | Libri2Mix-100 SI-SNRi /SDRi |
|---|---|---|---|
| Conv-TasNet [8] | 12.7/- | 8.3/- | 12.2/12.7 |
| SuDoRM-RF [13] | 13.7/14.1 | -/- | 14.0/14.4 |
| TDANet [15] | 15.2/15.4 | -/- | 17.4/17.9 |
| S4M-Tiny [26] | -/- | -/- | 16.2/16.6 |
| S4M [26] | -/- | -/- | 16.9/17.4 |
| TF-GridNet [6] | -/- | 17.1/15.6 | -/- |
| Sepformer + DM [5] | 16.4/16.7 | 14.0/13.0 | -/- |
| MossFormer(L) + DM [18] | 17.3/- | 16.3/- | -/- |
| MossFormer2 + DM [19] | 18.1/- | 17.0/- | -/- |
| SepReformer-T | 17.2/17.5 | -/- | 19.7/20.2 |
| SepReformer-B | 17.6/18.0 | -/- | 21.7/22.1 |
| SepReformer-L + DM | 18.5/18.7 | 17.1/16.0 | -/- |
| FLA-SepReformer-T | 16.9/17.2 | 14.8/13.6 | 19.1/19.5 |
| FLA-SepReformer-B | 17.4/17.8 | 15.6/14.4 | 20.3/20.7 |
| FLA-SepReformer-L + DM | 18.1/18.5 | 16.4/15.2 | -/- |
| TF-Locoformer-S | -/- | 17.4/15.9 | -/- |
| TF-Locoformer-M | -/- | 18.5/16.9 | -/- |
| FLA-TFLocoformer-S | -/- | 17.7/16.0 | -/- |
| FLA-TFLocoformer-M | 17.5/17.7 | 18.7/17.0 | -/- |

and FLA-TFLocoformer models on train-360.

duce the complexity to $O(N)$. Other components like the Separation Encoder, Reconstruction Decoder, Local Transformer, CS Transformer, and Speaker Split remain unchanged.

In Fig.2(c), for TF-Locoformer, we replace the original MHSA in Temporal modeling with Gated FLA. Although STFT helps reduce the length of speech features, the original MHSA maintains $O(N^2)$ complexity. Using Gated FLA, we convert quadratic complexity to linear, improving inference efficiency. Other components like Frequency Modeling and Conv-SwiGLU remain unchanged.

## 3. Experimental Setup

### 3.1. Dataset

We validate our model's performance using four popular speech separation datasets: WSJ0-2Mix [3], WHAM! [28], WHAMR! [29], and Libri2Mix [30]. We train and test all datasets using the full overlap min version with a sampling rate of 8kHz.

**WSJ0-2Mix** is a commonly used benchmark for SS, created from the WSJ0 corpus to generate 2-speaker clean mixtures. It consists of 30h, 10h, and 5h training, validation and test data.

**WHAM!/WHAMR!** are the noisy and noisy-reverberant versions of WSJ0-2Mix. WHAM! introduces noise from urban environments, mixed with speech at SNRs between -6 and +3 dB. WHAMR! further adds reverberation to the clean sources in WHAM!, which is used to train models for dereverberation, denoising, and speech separation.

**Libri2Mix** includes two training sets: train-360 with 106 hours and train-100 with 29 hours. Source speakers are drawn from LibriSpeech [31] sets (train-100, train-360). Both validation and test sets are 5.5 hours each. For a fair comparison with previous work, we train FLA-SepReformer models on train-100

### 3.2. Training and Model Configuration

We use the ESPnet-SE [32] toolkits for all experiments. For FLA-SepReformer, we mainly study the T/B/L model scale with settings similar to SepReformer: up to 200 training epochs, batch size of 2, and an initial learning rate (LR) of 1e-3 with a 1k-step warm-up. The LR is fixed for the first 50 epochs and then decays 0.8 if validation loss doesn't improve for 2 epochs. We use the same multi-loss as SepReformer. LR adjustments for dynamic mixing (DM) follow [5].

For FLA-TFLocoformer, we mainly study the S/M/L model scale, with settings similar to TF-Locoformer: up to 150 training epochs, batch size of 4, and an initial LR of 1e-3 with a 4k-step warm-up. The LR is halved if validation loss doesn't improve for 3 epochs. LR for S model remains unchanged for 50 epochs. With DM, training extends to 200 epochs, and the LR remains unchanged for 75/75/65 epochs for S/M/L, respectively.

We set FLA's focused factor $p$ to 3 and the DWC1d kernel size $k$ to 7. Gated MLP uses LayerNorm for SepReformer and RMSGroupNorm for TF-Locoformer. Both models utilize the AdamW optimizer with a 0.01 weight decay. Audio length is set to 4 seconds for training, and gradients are clipped at an L2 norm of 5. We use speed perturbation when doing DM. All experiments are conducted on a GeForce RTX A800. We plan to release the source code at a later time.

## 4. Results and Discussion

### 4.1. Comparison with previous models

We use SI-SNR improvement (SI-SNRi) and SDR improvement (SDRi) [33] to evaluate model performance. Also, we report the number of multiply-accumulate operations (MACs) for 8k samples using pytorch-OpCounter[1]. Table 1 shows results
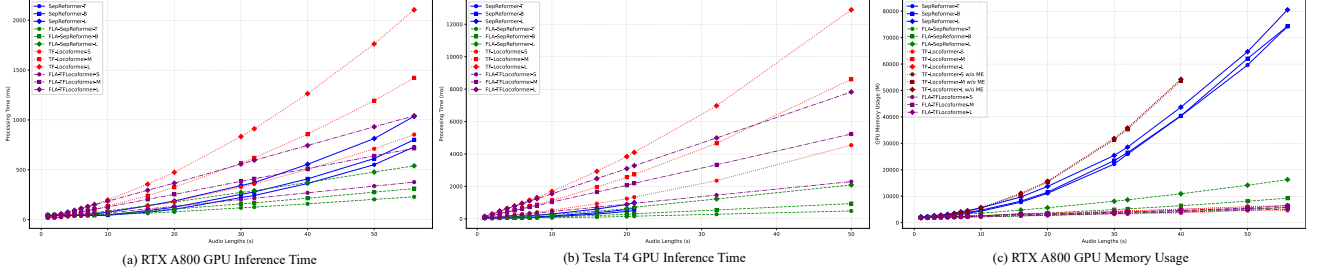
---

[1] https://github.com/Lyken17/pytorch-OpCounter

Figure 3: *Inference Time or Memory Usage using various GPU.*

(a) RTX A800 GPU Inference Time     (b) Tesla T4 GPU Inference Time     (c) RTX A800 GPU Memory Usage

Table 3: *Comparisons with other methods on Libri2Mix-360. DM is not used. Results in [dB].*

| Methods | SI-SNRi | SDRi |
|---|---|---|
| Conv-TasNet [8] | 14.7 | - |
| Sepformer [5] | 19.2 | 19.4 |
| MossFormer2 [19] | 21.7 | - |
| TF-Locoformer-M [7] | 22.1 | 22.2 |
| FLA-TFLocoformer-M | 22.2 | 22.4 |

Table 4: *Ablation study on WSJ0-2Mix using FLA-SepReformer-B. Results in [dB].*

| Gated | $p$ | $k$ | SI-SNRi | SDRi |
|---|---|---|---|---|
| ✓ | 3 | 7 | 23.5 | 23.7 |
|   | 3 | 7 | 23.4 | 23.5 |
| ✓ | 2/4 | 7 | 23.4/23.3 | 23.6/23.4 |
| ✓ | 8/16 | 7 | 23.1/22.9 | 23.3/23.0 |
| ✓ | 3 | 15/65 | 23.4/23.4 | 23.6/23.6 |

on the clean WSJ0-2Mix dataset, divided into four sections: efficient SS Model, normal SS Model, SepReformer/FLA-SepReformer, and TF-Locoformer/FLA-TFLocoformer. Our FLA-SepReformer-T surpasses previous efficient SS models like TDANet and S4M. Despite slight reductions on SI-SNRi and SDRi, it matches SepReformer in parameter count while reducing computational complexity from $O(N^2)$ to $O(N)$, leading to much faster inference and lower memory usage for long sequences, as shown in Figure 3.

Additionally, we reproduce SepReformer-B (Rep.), and our FLA-SepReformer-B achieves similar results. FLA-TF-Locoformer also matches TF-Locoformer's performance across all sizes (S/M/L) while having faster inference and lower memory usage for long sequences, proving our model's efficiency. Our largest models, FLA-SepReformer-L and FLA-TFLocoformer-L, outperform prior SOTA models like TF-GridNet, MossFormer2, and those using Mamba mechanisms. This shows our models' capability to deliver faster inference with enhanced performance, highlighting the effectiveness of the Gated FLA module.

Tables 2 and 3 show results on the WHAM!, WHAMR!, and Libri2Mix-100/360 datasets, demonstrating the robust generalization of FLA-SepReformer and FLA-TFLocoformer in noisy, reverberant conditions. Our models achieve performance similar to SepReformer and TF-Locoformer, even better in some metrics, such as Libri2Mix-360 FLA-TFLocoformer-M with an SI-SNRi of 22.2.

### 4.2. Inference time

Our key contribution is matching the results of SepReformer and TF-Locoformer while reducing the attention mechanism's complexity from quadratic to linear for long sequence SS. This ensures linear growth in inference time and memory usage,

preventing quadratic increases for lengthy speech. As Figure 1 shows, for 30s speech, FLA-SepReformer-T/B/L is faster by 2.29x/1.91x/1.49x than SepReformer-T/B/L, and FLA-TFLocoformer-S/M/L is faster compared to TF-Locoformer-S/M/L. Additionally, FLA-SepReformer-T surpasses other efficient SS Models.

We discover inaccuracies MACs for recent models using Pytorch-OpCounter and don't align to the inference time. We analyze GPU inference time and memory usage for FLA-SepReformer/TFLocoformer, observing linear growth in inference time and memory for FLA-SepReformer versus SepReformer. The test environment for GPU are RTX A800 , Tesla T4, single-threaded. SepReformer-L meets Out-of-Memory at 57s on an 80GB setup, while models with linear complexity achieve better efficiency for long speech. Though TF-Locoformer cuts memory usage to $O(\sqrt{N})$ with Memory-Efficient (ME) Attention [34], its inference time grows quadratically $O(N^2)$. And without ME, memory usage grows quadratically. In contrast, FLA-TFLocoformer maintains stable linear growth. Inference time calculations are averages from 100 iterations for each audio length.

### 4.3. Ablation Study

Table 4 shows the FLA-SepReformer-B results with various module and parameter modifications. Removing the Gated MLP drops SI-SNRi from 23.5 to 23.4 and SDRi from 23.7 to 23.5. Testing different focused factor $p$ values reveals that a lower value ($p = 2$) reduces SI-SNRi to 23.4, and higher $p$ values also decrease performance. Adjusting the DWC1d kernel size $k$ shows that larger sizes (like $k = 15/65$, similar to Local-Transformer) does not improve results. These results highlight the Gated MLP's role in gating global features and suggest that optimal $p$ and $k$ values enhance performance.

## 5. Conclusion

In this paper, we build FLASepformer, an efficient speech separation model with linear complexity. Although previous methods use STFT and downsampling to reduce speech sequence length, the attention module in those still has $O(N^2)$ time complexity. We improve SepReformer and TF-Locoformer by integrating Focused Linear Attention, creating two variants: FLA-SepReformer and FLA-TFLocoformer. We also add a new Gated module to improve performance. Experimental results on various datasets show that our models achieve similar SOTA results with reduced memory consumption and improved inference speed. FLA-SepReformer achieves speedups of 2.29x/1.91x/1.49x, and FLA-TFLocoformer-L also demonstrates significant speed gains. In the future, we will explore scenarios with more speakers.

# 6. References

[1] E. C. Cherry, "Some experiments on the recognition of speech, with one and with two ears," *Journal of the acoustical society of America*, vol. 25, pp. 975–979, 1953.

[2] A. S. Bregman, *Auditory Scene Analysis: The Perceptual Organization of Sound.* MIT Press, 1994.

[3] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, "Deep clustering: Discriminative embeddings for segmentation and separation," in *Proc. ICASSP.* IEEE, 2016, pp. 31–35.

[4] M. Kolbæk, D. Yu, Z.-H. Tan, and J. Jensen, "Multitalker speech separation with utterance-level permutation invariant training of deep recurrent neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 10, pp. 1901–1913, 2017.

[5] C. Subakan, M. Ravanelli, S. Cornell, M. Bronzi, and J. Zhong, "Attention Is All You Need In Speech Separation," in *Proc. ICASSP*, 2021, pp. 21–25.

[6] Z.-Q. Wang, S. Cornell, S. Choi, Y. Lee, B.-Y. Kim, and S. Watanabe, "TF-GRIDNET: Making Time-Frequency Domain Models Great Again for Monaural Speaker Separation," in *Proc. ICASSP*, 2023, pp. 1–5.

[7] K. Saijo, G. Wichern, F. G. Germain, Z. Pan, and J. Le Roux, "TF-Locoformer: Transformer with Local Modeling by Convolution for Speech Separation and Enhancement," in *Proc. International Workshop on Acoustic Signal Enhancement (IWAENC)*, Sep. 2024.

[8] Y. Luo and N. Mesgarani, "Conv-tasnet: Surpassing ideal time–frequency magnitude masking for speech separation," *IEEE/ACM transactions on audio, speech, and language processing*, vol. 27, no. 8, pp. 1256–1266, 2019.

[9] Y. Luo, Z. Chen, and T. Yoshioka, "Dual-Path RNN: Efficient Long Sequence Modeling for Time-Domain Single-Channel Speech Separation," in *Proc. ICASSP*, 2020, pp. 46–50.

[10] S. Hochreiter, "Long Short-term Memory," *Neural Computation MIT-Press*, 1997.

[11] J. Chen, Q. Mao, and D. Liu, "Dual-Path Transformer Network: Direct Context-Aware Modeling for End-to-End Monaural Speech Separation," in *Interspeech 2020*, 2020, pp. 2642–2646.

[12] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention Is All You Need," in *Proc. NIPS 2017*, vol. 30, 2017.

[13] E. Tzinis, Z. Wang, and P. Smaragdis, "Sudo RM -RF: Efficient Networks for Universal Audio Source Separation," in *Proc. MLSP.* IEEE, 2020, pp. 1–6.

[14] M. W. Y. Lam, J. Wang, D. Su, and D. Yu, "Sandglasset: A Light Multi-Granularity Self-Attentive Network for Time-Domain Speech Separation," in *Proc. ICASSP.* IEEE, 2021, pp. 5759–5763.

[15] K. Li, R. Yang, and X. Hu, "An efficient encoder-decoder architecture with top-down attention for speech separation," in *Proc. ICLR.* OpenReview.net, 2023.

[16] M. Xu, K. Li, G. Chen, and X. Hu, "TIGER: Time-frequency Interleaved Gain Extraction and Reconstruction for Efficient Speech Separation," in *Proc. ICLR.* OpenReview.net, 2025.

[17] H. Wang and B. Tian, "ZipEnhancer: Dual-Path Down-Up Sampling-based Zipformer for Monaural Speech Enhancement," in *Proc. ICASSP*, 2025, pp. 1–5.

[18] S. Zhao and B. Ma, "MossFormer: Pushing the Performance Limit of Monaural Speech Separation Using Gated Single-Head Transformer with Convolution-Augmented Joint Self-Attentions," in *Proc. ICASSP*, 2023, pp. 1–5.

[19] S. Zhao, Y. Ma, C. Ni, C. Zhang, H. Wang, T. H. Nguyen, K. Zhou, J. Q. Yip, D. Ng, and B. Ma, "MossFormer2: Combining Transformer and RNN-Free Recurrent Network for Enhanced Time-Domain Monaural Speech Separation," in *Proc. ICASSP*, 2024, pp. 10 356–10 360.

[20] W. Hua, Z. Dai, H. Liu, and Q. V. Le, "Transformer Quality in Linear Time," in *Proc. ICML*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvári, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 2022, pp. 9099–9117.

[21] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret, "Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention," in *Proc. ICML*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 5156–5165.

[22] D. Han, X. Pan, Y. Han, S. Song, and G. Huang, "FLatten Transformer: Vision Transformer using Focused Linear Attention," in *Proc. ICCV.* IEEE, 2023, pp. 5938–5948.

[23] U. Shin, S. Lee, T. Kim, and H. Park, "Separate and Reconstruct: Asymmetric Encoder-Decoder for Speech Separation," in *Proc. NeurIPS*, A. Globersons, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. M. Tomczak, and C. Zhang, Eds., 2024.

[24] H. Cai, J. Li, M. Hu, C. Gan, and S. Han, "EfficientViT: Lightweight Multi-Scale Attention for High-Resolution Dense Prediction," in *Proc. ICCV.* IEEE, 2023, pp. 17 256–17 267.

[25] K. M. Choromanski, V. Likhosherstov, D. Dohan, X. Song, A. Gane, T. Sarlós, P. Hawkins, J. Q. Davis, A. Mohiuddin, L. Kaiser, D. B. Belanger, L. J. Colwell, and A. Weller, "Rethinking Attention with Performers," in *Proc. ICLR.* OpenReview.net, 2021.

[26] C. Chen, C.-H. H. Yang, K. Li, Y. Hu, P.-J. Ku, and E. S. Chng, "A Neural State-Space Modeling Approach to Efficient Speech Separation," in *Interspeech 2023*, 2023, pp. 3784–3788.

[27] X. Jiang, Y. A. Li, A. N. Florea, C. Han, and N. Mesgarani, "Speech Slytherin: Examining the Performance and Efficiency of Mamba for Speech Separation, Recognition, and Synthesis," *arXiv preprint arXiv:2407.09732*, 2024. [Online]. Available: https://arxiv.org/abs/2407.09732

[28] 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 *Interspeech 2019*, 2019, pp. 1368–1372.

[29] M. Maciejewski, G. Wichern, E. McQuinn, and J. L. Roux, "WHAMR!: Noisy and Reverberant Single-Channel Speech Separation," in *Proc. ICASSP.* IEEE, 2020, pp. 696–700.

[30] J. Cosentino, M. Pariente, S. Cornell, A. Deleforge, and E. Vincent, "LibriMix: An Open-Source Dataset for Generalizable Speech Separation," 2020.

[31] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. ICASSP*, 2015, pp. 5206–5210.

[32] Y.-J. Lu, X. Chang, C. Li, W. Zhang, S. Cornell, Z. Ni, Y. Masuyama, B. Yan, R. Scheibler, Z.-Q. Wang, Y. Tsao, Y. Qian, and S. Watanabe, "ESPnet-SE++: Speech Enhancement for Robust Speech Recognition, Translation, and Understanding," in *Interspeech 2022*, 2022, pp. 5458–5462.

[33] J. L. Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR – Half-baked or Well Done?" in *Proc. ICASSP*, 2019, pp. 626–630.

[34] M. N. Rabe and C. Staats, "Self-attention Does Not Need $O(n^2)$ Memory," *arXiv preprint arXiv:2112.05682*, 2021. [Online]. Available: https://arxiv.org/abs/2112.05682