

EfficientASR: Speech Recognition Network Compression via Attention Redundancy and Chunk-Level FFN Optimization

Jianzong Wang^{1†}, Ziqi Liang^{1,2‡}, Xulong Zhang^{1✉}, Ning Cheng¹, Jing Xiao¹

¹Ping An Technology (Shenzhen) Co., Ltd.

²University of Science and Technology of China

Abstract—In recent years, Transformer networks have shown remarkable performance in speech recognition tasks. However, their deployment poses challenges due to high computational and storage resource requirements. To address this issue, a lightweight model called EfficientASR is proposed in this paper, aiming to enhance the versatility of Transformer models. EfficientASR employs two primary modules: Shared Residual Multi-Head Attention (SRMHA) and Chunk-Level Feedforward Networks (CFFN). The SRMHA module effectively reduces redundant computations in the network, while the CFFN module captures spatial knowledge and reduces the number of parameters. The effectiveness of the EfficientASR model is validated on two public datasets, namely Aishell-1 and HKUST. Experimental results demonstrate a 36% reduction in parameters compared to the baseline Transformer network, along with improvements of 0.3% and 0.2% in Character Error Rate (CER) on the Aishell-1 and HKUST datasets, respectively.

Index Terms—speech recognition, attention redundancy, feed-forward network, lightweight

I. INTRODUCTION

In recent years, Transformers [1] have shown better performance than traditional sequence models [2], [3] in the ASR domain, capturing long-range dependencies through attention mechanisms. Attention mechanisms play a crucial role in Transformers by establishing connections between different positions through identity mapping. However, attention computation in Transformers is computationally expensive and contains a significant amount of redundancy. On the other hand, feedforward networks capture high-level representations through high-dimensional feature mappings, but this also leads to an increase in network parameters. These issues result in high-performing Transformer models requiring substantial storage and computational resources, making it challenging to apply them on limited computing devices.

The importance of attention mechanisms in Transformer networks has been extensively studied. For instance, He et al. [4] introduced the residual attention mechanism to improve model performance by enabling interaction and fusion of attention across different layers. Research by [5] revealed that as attention propagates from lower to higher layers, the attention distribution becomes more diagonalized, with higher-layer attention contributing less to model performance.

Additionally, the study by [6] demonstrated that randomly removing some attention heads does not significantly affect model performance, suggesting redundancy in attention computations. Shim et al. [7] found that lower-layer attention focuses more on semantic features of speech, while higher-layer attention is more concerned with language positioning. Moreover, the investigation by [8] discovered that attention distributions between adjacent layers tend to be more similar, reducing network redundancy by sharing certain attention heads. HybridFormer [9] proposed HyperMixer to replace multi-head self-attention, which can model local interaction and global interaction information respectively. Benefit from HyperMixer’s linear time and memory complexity, significantly reducing computational and memory costs.

In the domain of model lightweighting, effective methods like quantization-aware training [10] and knowledge distillation [11] [12] exist but can be complex to implement. In contrast, weight sharing [13] [14] offers a more straightforward model compression approach. However, when only the learnable parameters are shared, it doesn’t reduce the forward propagation computation of the network. By analyzing the parameter and computational costs of the attention mechanism and feed-forward networks in Transformer, we introduce an innovative approach to address attention redundancy, sharing attention scores across layers through residual connections for interactive fusion between higher and lower layers. At the same time, in order to reduce the computational redundancy introduced by cross-layer shared attention scores and enhance the ability of attention to extract local features, we apply the sliding window with deformability (SWD) method [15] to shared attention scores, reducing redundancy in single-layer attention maps. When lightweighting attention-based speech recognition models, we often focus on the point: Reduce attention redundancy and parameter count effectively.

To tackle these challenges, in this paper, we introduce an efficient ASR network, called EfficientASR, which 1) leverages sharing residual attention to reduce attention redundancy and 2) employs the Chunk-level FFN structure to reduce parameter count effectively. Specifically, the method shares attention scores among multiple layers and integrates old attention scores into new ones through residual connections, achieving interactive fusion of attention scores between higher and lower layers. Chunk-level Feed-forward Neural Network (CFFN)

[†] Equal Contributions

[✉]Corresponding author: Xulong Zhang (zhangxulong@ieee.org).

divides the feed-forward networks into multiple chunks based on embedding dimensions, using smaller FFNs in each chunk, substantially reducing the number of learnable parameters without compromising model performance. Additionally, the paper applies the previously proposed sliding window with deformability (SWD) [15] method to the shared attention scores, addressing the slow change of diagonalization degree in higher layers after shared attention and further reducing redundancy in single-layer attention maps.

II. RELATED WORK

A. Automatic speech recognition

With the application of transformer in the field of natural language processing (NLP) and its significant effect on modeling contextual information, transformer is also suitable for the field of speech [16]–[18]. Existing automatic speech recognition (ASR) works are mainly divided into three structures: CNN, RNN, and transformer. Previously, transformer and CNN were used as networks to achieve good results in ASR, but both have their limitations: transformer is not good at extracting fine-grained partial feature, while the CNN network can capture local features such as edge lines and shapes, but more convolutional layers are needed to capture global information.

In subsequent work on the ASR task, there are improvements in recognition efficiency by using depth-wise separable convolutions [19]. The model is composed of several interconnected modules, featuring residual links that facilitate the flow of information between them. Within each module, there is an implementation of one or multiple 1D temporal channel separable convolutional filters. There is also a combination of squeeze and excitation modules to capture longer contextual semantic information [20]. Compared with CNN, the significant benefit of transformer is that it can learn global dependencies based on self-attention, which is key for speech processing tasks. Gulati et al. [21] proposed the Conformer, since the Transformer model is good at capturing global interactions based on content and CNN effectively utilizes local features, Conformer integrates the capabilities of both the Transformer and CNN to effectively capture and process both local and global interdependencies present within the audio sequence.

However, in the previous Conformer structure, the feature representations derived from consecutive speech frames contained considerable repetition, resulting in unnecessary computational consumption. In this regard, Squeezeformer [22] introduces a sequential U-Net structure. The down-sampling layer halves the sampling rate of the input speech signal. At the end of the network, a lightweight up-sampling layer restores the speech signal from a low sampling rate to the original high sampling rate, which ensures that model training is stable. Branchformer [23] adopts a parallel dual-branch structure. A particular branch within the model leverages a multi-headed self-attention system to capture expansive characteristics across the input sequence, and the other branch introduces the MLP with convolutional gating structure, which

is intended to capture local features in the audio sequence. [24] proposed E-Branchformer, an advancement over the Branchformer, which applies efficient model fusion methods and stacking additional point-adding modules. As a result, E-Branchformer established a new benchmark for WER, notably without leveraging any external training data. Zipformer [25] has similar ideas to previous work on downsampling the temporal dimension. However, compared to the fixed downsampling ratio in the Squeezeformer [22], the Squeezeformer works with different downsampling ratios in different encoder stacks, and uses a larger downsampling ratio in the intermediate encoder stack.

B. Efficient attention mechanism

The self-attention mechanism is a critical component in many ASR models. It allows the model to dynamically focus on different parts of the input sequence, which is essential for understanding the context and structure of the speech data. However, the computational complexity of this mechanism is a significant concern due to its quadratic scaling with respect to the length of the input sequence. There are many ways to improve the efficiency of the transformer. Sparse attention mechanisms [26], [27], [27]–[32] offer a promising approach to address the challenges associated with the complexity of self-attention. Through the self-attention process, the mechanism transcends limitations inherent in RNNs' sequential framework, allowing individual tokens in a given sequence to engage with the entire set of tokens without dependence on their order.

However, numerous applications limit context to a very short sequence length due to current hardware and model size constraints. These constraints typically restrict input sequences to approximately 512 tokens in length, significantly diminishing direct applicability to tasks requiring broader context. To tackle this issue, [27] propose a sparse attention mechanism capable of processing sequences up to 8 times longer than previously achievable with comparable hardware. This advancement significantly enhances the performance of various Natural Language Processing (NLP) tasks, including question answering and summarization, by effectively managing longer in-context information. Longformer [28] uses sliding window attention to focus on local in-context information, and uses global attention on some preselected input positions to capture global in-context information. To mitigate potential declines in performance, [27] focuses on randomly selected pairs of tokens, but it might necessitate utilizing a substantial quantity of tokens to reduce performance degradation on longer sequences, which often limits computational speed-up.

However, even if these new methods surpass the state-of-the-art, their use in the real world often comes with significant resource costs. Hence, there is a crucial need to investigate efficient and lightweight model networks that hold increased commercial value.

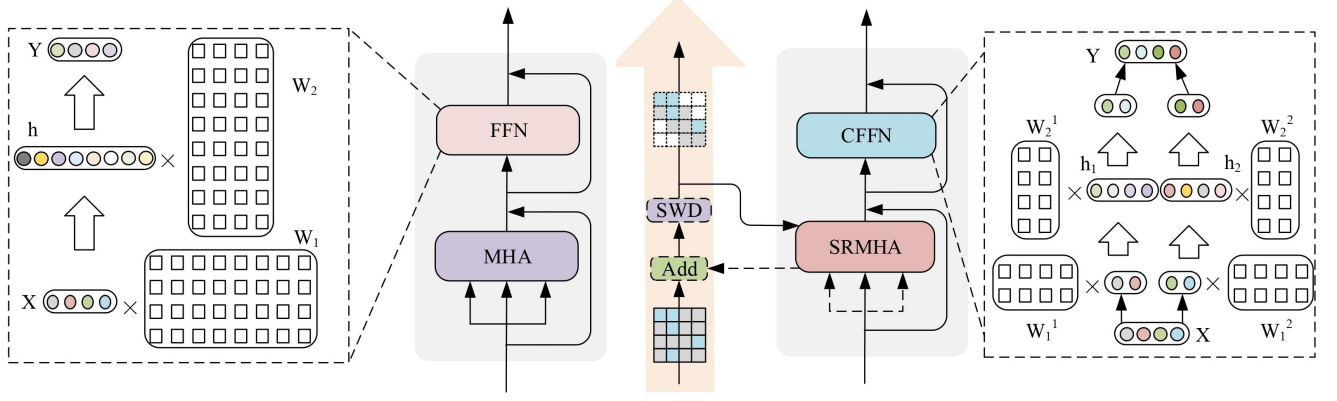


Fig. 1. (a) is a traditional Transformer structure. (b) is an EfficientASR structure, where the dashed line portion is only used in shared attention mode. The small circles represent the embedding dimension of the input features, and the small squares represent learnable parameters. SWD represents the sliding window with deformability

III. PROPOSED METHOD

A. Review Transformer

The encoder-decoder Transformer model has achieved excellent performance in ASR, a success largely attributable to its integration of an attention mechanism with a feed-forward neural network. The specific calculation process of the attention mechanism is as follows:

$$Attn = Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \quad (1)$$

The number of learnable parameters for each attention module is $4d_{model}^2 + 4d_{model}$. In addition, when the input $X \in \mathbb{R}^{B \times T \times d_{model}}$ passes through the attention module, the number of floating point calculations is $4T^2Bd_{model} + 8BTd_{model}^2$.

For the feed-forward network, it consists of two linear layers, where high-dimensional feature mapping is usually applied:

$$FFN(X) = ReLU(XW_1 + b_1)W_2 + b_2, \quad (2)$$

Where $ReLU$ serves as the activator, the matrices W_1, W_2 are belonging to the space $\mathbb{R}^{d_{model} \times d_{ff}}$ and $\mathbb{R}^{d_{ff} \times d_{model}}$, and d_{ff} is the hidden dimension of FFN. $b_1 \in \mathbb{R}^{d_{ff}}$ and $b_2 \in \mathbb{R}^{d_{model}}$ are bias vectors. It is worth noting that d_{ff} is often several times greater than d_{model} , for example, in the ESPnet [33] toolkit, the default value of d_{ff} is $8d_{model}$. Therefore, the parameter count of the FFN is $16d_{model}^2 + 9d_{model}$, and the number of floating point calculations is $32BTd_{model}^2$.

Based on the above analysis, it can be found that when $d_{ff} = 8d_{model}$, the number of learnable parameters in FFN is much larger than that of MHA, which is caused by the large weight matrices in FFN due to high-dimensional feature mapping. Even though MHA possesses a lesser count of adjustable parameters compared to FFN, our research indicates that as the size of the input feature set, represented by length T , increases the computational demand for MHA escalates significantly, surpassing that required by FFN. This is because

the computational complexity of QK^T and SV in the calculation process is proportional to the square of T . Therefore, this also inspires us to improve both MHA and FFN to reduce the overall computational complexity and number of learnable parameters of the model.

B. EfficientASR

Following the preceding discussion, we introduce an innovative lightweight network structure called EfficientASR. Illustrated on the right side of Fig 1, each encoder layer is comprised of an SRMHA component, an SWD component, and a CFFN component, while the decoder layers also include a cross multi-head attention module. Residual links and layer-wise normalization are integrated within each module, as referenced in [34], [35]. EfficientASR reduces the computational complexity and number of learnable parameters by reducing redundant information in the attention mechanism.

Previous studies have shown that there is a significant amount of redundancy in the repeated calculation of attention scores in Transformer networks [6] [8]. However, simply sharing attention heads to other layers is not the optimal choice. Inspired by the concept of [5] as the depth of the network increases, the diagonalization of attention becomes more prominent. Sharing attention heads directly to other layers slows down the formation of this diagonalization. In place of the conventional Transformer's multi-head self-attention component, we have introduced a substitute module with a shared residual multi-head attention (SRMHA) module, while retaining the cross self-attention in the decoder. The SRMHA module has two modes of operation across different sub-layers: shared attention and updated attention, as depicted on the right side of Figure 1. For any layer of SRMHA, its expression is:

$$S = W_q Z_q \cdot (W_k Z_k)^T \quad (3)$$

$$V = W_v Z_v \quad (4)$$

$$SRMHA = Softmax(S)V. \quad (5)$$

Here, S represents attention scores, and V represents the linear mapping of inputs.

C. Updated Attention Mode

In the speech recognition task, neighboring tokens have stronger correlations. Therefore, the SWD module restricts the interaction range between tokens by only retaining the data within the w -range of the attention matrix diagonal. Under the updated attention mode, the attention matrix S is updated as follows:

$$S = SWD(S', w), \quad (6)$$

$$S' = \frac{QK^T}{\sqrt{d_k}} + S'', \quad (7)$$

here, S' represents the updated attention matrix, S'' represents the previous updated attention matrix, and d_k denotes the size of either the word vector dimensions or the dimensions within the hidden layer. By using residual connections for attention, similar to [4], high-level and low-level feature distributions can be fused. Specifically, by connecting the attention matrix $\frac{QK^T}{\sqrt{d_k}}$ generated from the input features of the current layer with the old attention matrix S'' through residual connections, we obtain a new attention matrix S' . The SWD module and attention matrix residual connections are only used in sub-layers under the updated attention mode.

D. Shared Attention Mode

To reduce the redundant calculation of attention, we reduce the need to calculate attention matrices for the same input features by sharing attention. Under the shared attention mode, the attention matrix S is equal to the old attention matrix S'' , and the SRMHA module does not need to generate the Q and K matrices, thereby reducing the computational burden of two linear layers and the calculation of QK^T . In this case, $S = S''$, and the SRMHA module does not need to generate the Q and K matrices, reducing the computational cost of two linear layers and the calculation of QK^T . In this case, the SRMHA module has $2d_{model}^2 + 2d_{model}$ learnable parameters and requires $2T^2Bd_{model} + 4BTd_{model}^2$ floating-point operations in terms of computation. Compared with the standard self-attention module in Transformers, the learnable parameters and floating-point operations of the SRMHA module are reduced by half.

E. Chunk-level Feed-Forward Network

The mapping of high-dimensional features dominates most of the learnable parameters in Transformer, but it has been rarely studied in speech recognition. Based on the inspiration from [36], we treat the feed-forward network as a key-value memory and redefine its expression as:

$$FFN = f(XK^T)V. \quad (8)$$

Here, K and V correspond to the learnable matrices W_1 and W_2 , respectively, and f is utilized to signify the ReLU. The key of each layer is used to capture patterns in the input sequence, while the value generates the distribution of the next

token based on the captured patterns. In order for the feed-forward network to serve as a key-value memory and learn different feature representations, covering a broader range of semantic relationships, we propose a chunk-level feed-forward network (CFFN). As shown in Fig. 1, CFFN can project and map the input sequence differently, reducing the burden of high-dimensional computations and improving scalability. Specifically, the expression of CFFN is as follows:

$$CFFN = \text{Concat}(\text{chunk}_1, \dots, \text{chunk}_n), \quad (9)$$

$$\text{chunk}_i = FFN(X_i), \quad (10)$$

$$X_1, \dots, X_n = \text{Split}(X), \quad (11)$$

Here, n represents the number of chunks, and each chunk represents a key-value memory. When $d_{ff} = 8d_{model}$, the amount of parameters in each chunk is $16(d_{model}/n)^2 + 9d_{model}/n$. By introducing the chunk-level feed-forward network (CFFN), the model can better utilize the feed-forward network as a key-value memory to learn different feature representations, improve the performance of the speech recognition model, and reduce the burden of high-dimensional computations.

F. The global learning of EfficientASR

We selected 2 loss functions for global learning and the fitting of EfficientASR. We utilize cross-entropy loss \mathcal{L}_{CE} to calculate the probability distribution difference between the EfficientASR output and the target sequence. A CTC loss \mathcal{L}_{CTC} [37] is selected to improve the recognition accuracy of acoustic boundaries and achieve the alignment of acoustic tokens and semantic tokens.

$$\mathcal{L} = \alpha_1 \mathcal{L}_{CE} + \alpha_2 \mathcal{L}_{CTC} \quad (12)$$

Notably, α_1 and α_2 represent weight hyperparameters that adjust the relative importance of these two loss functions. Once \mathcal{L} is obtained, we compute the gradients of all parameters and update our model using backward propagation.

IV. EXPERIMENT

A. Experimental Setup

Dataset and Preprocessing: To assess the performance of the EfficientASR, we have chosen two prominent datasets: Aishell-1 [38] and HKUST [39]. The preprocessing stage of the audio input involves a series of transformations, such as segmentation into frames, application of a window function, computation via the fast Fourier transform (FFT), and subsequent processing with the discrete cosine transform (DCT). In this process, each window spans 25 ms, with a stagger of 10 ms between successive windows. The resultant feature set consists of 80 dimensions, represented by log Mel-filterbank coefficients. In the context of speech enhancement [40], the time mask window is 30 and the frequency mask window is 40. In the convolutional downsampling, the size of the convolutional kernel is 31.

Hyperparameter: The ESPNet toolkit [33] serves as the platform for executing all our experiments. Specifically, we

have implemented the label smoothing technique and a dropout rate of $p=0.1$ to mitigate the risk of overfitting in our model. For optimization, we have chosen the Adam algorithm [41], with a learning rate set to 0.002. Supplementary settings encompass a model dimension of $d_{model}=256$, a feed-forward layer size of $d_{ff}=2048$, with window dimensions denoted by $w=6$ and $h=4$, followed by the number of attention heads 12 and the number of encoder layers 6. In the attention mechanism, we have set d_q , d_k , and d_v dimensions equal to $d_{model}/h=64$. Each training batch commences with a size of 64, and the beam search strategy employs a beam width of 10. Our language model (LM), which operates on the Transformer, comprises 16 layers and is trained over 15 epochs. During training, the weight for the CTC is 0.3, and during inference, the weight for CTC is 0.6.

B. Evaluation systems

To objectively compare the automatic speech recognition (ASR) performance of the acquired representations across various noisy environments, we developed four replication-synthesis frameworks utilizing both Aishell-1 [38] and HKUST [39] corpora:

- Transformer w/o LM [1]: Transformer-based ASR model combined with pre-trained language model.
- Transformer_LM [1]: Transformer-based ASR model combined without pre-trained language model
- Conformer [42]: Combining CNN and transformer to model both local and global interdependencies within audio sequences.
- Conformer_EfficientASR: Replace multi-head attention (MHA) in Conformer with SRMHA and feedforward (FFN) network with CFFN.

C. Results

The experimental results conducted on the Aishell-1 dataset are shown in Table I. The EfficientASR model has only 19.33M parameters, which is a 36% reduction compared to the Transformer model. The EfficientASR model realized a 0.1% reduction in Character Error Rate (CER) for the development (dev) set and a 0.3% decrease for the test sets. To validate the model’s scalability, SRMHA and CFFN were applied to the Conformer model, resulting in the Conformer_EfficientASR model. Experimental results indicate that the Conformer_EfficientASR model reduces the parameter count by 38% compared to the Conformer model while maintaining a test set CER of 4.9%. Therefore, the proposed EfficientASR model significantly minimizes the quantity of parameters that require learning, as well as the computational expense of the network, without notably degrading the performance of the model.

Similarly, we further validated the proposed EfficientASR model on the HKUST dataset, and the experimental results are presented in Table II. In contrast to the Transformer model, the EfficientASR model reduces the parameter count by 36%. Employing a Language Model (LM) on the dev set led to a 0.2% enhancement in CER performance. In contrast to

TABLE I
COMPARISON OF CER RESULTS BETWEEN EFFICIENTASR AND TRANSFORMER BASELINE ON AISHELL-1.

Model	dev/test (%)	Parameter (M)
Transformer w/o LM [1]	5.5/5.9	30.35
Transformer_LM [1]	5.2/5.6	30.35
Conformer [42]	4.5/4.9	45.40
EfficientASR w/o LM	5.4/ 5.6	19.33 (36% ↓)
EfficientASR_LM	5.1/ 5.3	19.33 (36% ↓)
Conformer_EfficientASR	4.6/4.9	28.10 (38% ↓)

TABLE II
COMPARISON OF CER RESULTS BETWEEN EFFICIENTASR AND TRANSFORMER BASELINE IN HKUST DATASET.

Model	dev (%)	Parameter (M)
Transformer w/o LM [1]	21.7	29.85
Transformer_LM [1]	21.5	29.85
Conformer [42]	19.8	44.96
EfficientASR w/o LM	22.0	18.86 (36% ↓)
EfficientASR_LM	21.3	18.86 (36% ↓)
Conformer_EfficientASR	20.1	29.23 (38% ↓)

the Conformer model, the Conformer_EfficientASR achieves a reduction of 38% in the count of parameters, but lowers the CER by 0.3%. These experiments on different datasets demonstrate that the EfficientASR model maintains strong competitiveness.

TABLE III
COMPARISON OF SRMHA AND OTHER ATTENTIONAL REDUNDANCY REDUCTION METHODS ON HKUST DATASET

Model	dev (%)	Parameter (M)
Transformer [1]	21.5	29.85
Realformer [4]	21.5	29.85
SPAH [8]	21.3	29.06
SEAM [43]	21.4	28.40
RSEAM	21.3	28.40
SRMHA	21.0	28.40

D. Ablation Studies

Method comparison: We compared the sharing of partial attention heads (SPAH) [8] with sharing the entire attention map (SEAM) [43] for reducing attention redundancy. We also compared the addition of residual attention based on sharing the entire attention map (RSEAM), as well as the use of the SWD method in SRMHA.

The results of the experiments are detailed in Table III. The performance of residual attention is comparable to the baseline Transformer. We evaluated the model’s CER on the development set and the parameter size of different models. SPAH improved the CER by 0.2% while reducing only 0.79M parameters. SEAM improved the CER by only 0.1% but reduced 1.45M parameters and reduced more attention computations. Based on these experimental results, combining residual attention with sharing the entire attention map can simultaneously reduce attention redundancy and improve network performance. Furthermore, when applying the SWD method to the network, the network can reduce even more redundant computations and exhibit diagonalized features at higher layers. Therefore, compared to previous methods, SRMHA has more advantages as it not only eliminates more attention redundancy but also achieves better performance in model performance.

TABLE IV
USE SRMHA ON THE ENCODER OR DECODER ON THE HKUST DATASET.

Encoder	Decoder	dev(%)	Parameter (M)
SR_2	×	22.1	29.06
SR_3	×	21.1	28.79
SR_4	×	21.1	28.66
SR_6	×	21.3	28.53
SR_{12}	×	21.8	28.40
×	SR_2	21.2	29.45
×	SR_3	21.4	29.32
×	SR_6	21.2	29.19
SR_3	SR_2	21.0	28.40

In order to examine the effects of minimizing attention mechanism calculations on the model’s performance, we integrated the SRMHA technique into both the encoding and decoding components of the Transformer architecture. The findings from our experiments are delineated in Table IV.

Attention Redundancy Analysis: To investigate the effects of reducing attention computations on the model, we integrated the SRMHA method into both the encoding and decoding components of the Transformer. The findings from our experiments are delineated in Table IV. Here, SR_i denotes updating attention computations every i sub-layers, where a larger value of i leads to a reduction in parameters and attention computations. × represents the traditional self-attention method.

In the encoder, using SRMHA with updates every 3 or 4 layers yielded optimal results, and updating every 6 layers did not significantly impact model performance. Similarly, in the decoder, using SRMHA with values of $i = 2$, $i = 4$, and $i = 6$ all benefited network performance. This is because by employing SRMHA in the decoder instead of multi-head attention (MHA), the contextual distribution of labels can be repeatedly applied at different layers, thereby facilitating network learning.

Ultimately, in EfficientASR, we selected a configuration of $i = 3$ in the encoder and $i = 2$ in the decoder. Through this set of experiments, we found that SRMHA did not significantly reduce the majority of network parameters. However, the method analysis revealed that it could effectively reduce a substantial amount of repetitive attention computations.

TABLE V
COMPARE THE IMPACT OF DIFFERENT CHUNK-LEVEL FFNS.

Model	Dataset	CER (%)	Parameter (M)
Transformer [1]	Aishell-1	5.2/5.6	30.35
1/2FFN	Aishell-1	5.2/5.4	20.91 (31% ↓)
1/4FFN	Aishell-1	5.3/5.6	16.20 (47% ↓)
Transformer [1]	HKUST	21.5	29.85
1/2FFN	HKUST	21.6	20.44 (32% ↓)
1/4FFN	HKUST	22.4	15.74 (47% ↓)

Chunk-level Feedforward Networks: We partitioned the embedding dimension of input features into multiple chunks and analyzed the effects of varying chunk sizes on the efficiency of performance. The results are detailed in Table V, where 1/2 FFN and 1/4 FFN represent splitting the embedding dimension into two and four chunks, respectively.

On the Aishell-1 corpora, the results showed that using 1/2FFN reduced the network parameters by 31% and improved

the CER by 0.2% on the test sets. Further employing 1/4FFN reduced 47% of parameters while achieving the same CER as the baseline. However, on the HKUST dataset, using 1/2FFN reduced parameters by 32%, but performance deteriorated with a 1% increase in CER. On the other hand, using 1/4FFN reduced 47% of parameters, but performance experienced a significant decline. Based on the experimental results, we observed significant fluctuations in the effectiveness of this method across different datasets. Nevertheless, selecting 1/2 FFN allows for a significant reduction in network parameters while maintaining stable performance.

E. Memory usage of long sequences

As shown in Fig.2, we show the impact of using SRMHA and CFFN on memory usage at different input sequence lengths. Experimental results show that as the sequence length continues to increase, the architecture using EfficientASR can effectively reduce the overall memory consumption (including model loading and data reading).

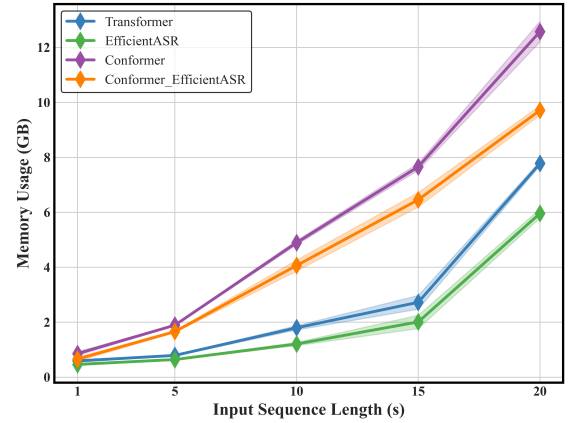


Fig. 2. Different models’ memory usage for processing long sequences on HKUST dataset.

For audio sequences longer than 10 seconds, the utilization of SRMHA and CFFN in the encoder and decoder layers of ASR model leads to a reduction in overall memory usage for the same sequence length. Regardless of whether the base model is a Transformer-based or Conformer-based model, our method demonstrates the ability to decrease memory consumption when processing long sequences.

V. CONCLUSION

The proposed EfficientASR, 1) reduces redundant computations, enhances attention diagonalization and fusion with sliding windows and residuals in the network; 2) captures spatial knowledge and reduces parameter size by dividing the input features into multiple chunks and using smaller feed-forward networks in each chunk.

VI. ACKNOWLEDGEMENT

Supported by the Key Research and Development Program of Guangdong Province (grant No. 2021B0101400003) and Corresponding author is Xulong Zhang (zhangxulong@ieee.org).

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Neural Information Processing Systems, NeurIPS*, vol. 30, 2017, pp. 5998–6008.
- [2] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *ICASSP*, 2013, pp. 6645–6649.
- [3] P. Malhotra, L. Vig, G. Shroff, and P. Agarwal, "Long short term memory networks for anomaly detection in time series," in *European Symposium on Artificial Neural Networks, ESANN*, 2015, p. 89.
- [4] R. He, A. Ravula, B. Kanagal, and J. Ainslie, "Realformer: Transformer likes residual attention," in *Findings of the Association for Computational Linguistics: ACL/IJCNLP*, 2021, pp. 929–943.
- [5] S. Zhang, E. Loweimi, P. Bell, and S. Renals, "On the usefulness of self-attention for automatic speech recognition with transformers," *IEEE Spoken Language Technology Workshop (SLT)*, pp. 89–96, 2020.
- [6] Shucong Zhang and Erfan Loweimi and Peter Bell and Steve Renals, "Stochastic attention head removal: A simple and effective method for improving transformer based ASR models," in *Interspeech*, 2021, pp. 2541–2545.
- [7] K. Shim, J. Choi, and W. Sung, "Understanding the role of self attention for efficient speech recognition," in *International Conference on Learning Representations, ICLR*, 2021.
- [8] S. Bhojanapalli, A. Chakrabarti, A. Veit, M. Lukasik, H. Jain, F. Liu, Y. Chang, and S. Kumar, "Leveraging redundancy in attention with reuse transformers," *CoRR*, vol. abs/2110.06821, 2021.
- [9] F. Mai, J. Zuluaga-Gomez, T. Parcollet, and P. Motlicek, "HyperConformer: Multi-head HyperMixer for Efficient Speech Recognition," in *Interspeech*, 2023, pp. 2213–2217.
- [10] O. Zafrir, G. Boudoukh, P. Izsak, and M. Wasserblat, "Q8bert: Quantized 8bit bert," in *Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing, EMC2@NeurIPS*. IEEE, 2019, pp. 36–39.
- [11] X. Jiao, Y. Yin, L. Shang, X. Jiang, X. Chen, L. Li, F. Wang, and Q. Liu, "Tinybert: Distilling BERT for natural language understanding," in *Empirical Methods in Natural Language Processing, EMNLP*, 2020, pp. 4163–4174.
- [12] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter," *CoRR*, vol. abs/1910.01108, 2019.
- [13] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "ALBERT: A lite BERT for self-supervised learning of language representations," in *International Conference on Learning Representations, ICLR*, 2020.
- [14] M. Dehghani, S. Gouws, O. Vinyals, J. Uszkoreit, and L. Kaiser, "Universal transformers," in *International Conference on Learning Representations, ICLR*, 2019.
- [15] G. Wei, Z. Duan, S. Li, X. Yu, and G. Yang, "Lfeformer: Local feature enhancement using sliding window with deformability for automatic speech recognition," *IEEE Signal Processing Letters*, vol. 30, pp. 180–184, 2023.
- [16] Y. Wang, A. rahman Mohamed, D. Le, C. Liu, A. Xiao, J. Mahadeokar, H. Huang, A. Tjandra, X. Zhang, F. Zhang, C. Fuegen, G. Zweig, and M. L. Seltzer, "Transformer-based acoustic modeling for hybrid speech recognition," *ICASSP*, pp. 6874–6878, 2020.
- [17] S. Watanabe, T. Hori, S. Kim, J. R. Hershey, and T. Hayashi, "Hybrid ctc/attention architecture for end-to-end speech recognition," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, pp. 1240–1253, 2017.
- [18] Q. Zhang, H. Lu, H. Sak, A. Tripathi, E. McDermott, S. Koo, and S. Kumar, "Transformer transducer: A streamable speech recognition model with transformer encoders and RNN-T loss," in *ICASSP*, 2020, pp. 7829–7833.
- [19] S. Kriman, S. Beliaev, B. Ginsburg, J. Huang, O. Kuchaiev, V. Lavrukhin, R. Leary, J. Li, and Y. Zhang, "Quartznet: Deep automatic speech recognition with 1d time-channel separable convolutions," in *ICASSP*, 2020, pp. 6124–6128.
- [20] W. Han, Z. Zhang, Y. Zhang, J. Yu, C.-C. Chiu, J. Qin, A. Gulati, R. Pang, and Y. Wu, "ContextNet: Improving Convolutional Neural Networks for Automatic Speech Recognition with Global Context," in *Interspeech*, 2020, pp. 3610–3614.
- [21] A. Gulati, J. Qin, 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 *Interspeech*, 2020, pp. 5036–5040.
- [22] S. Kim, A. Gholami, A. Shaw, N. Lee, K. Mangalam, J. Malik, M. W. Mahoney, and K. Keutzer, "Squeezeformer: An efficient transformer for automatic speech recognition," in *Neural Information Processing Systems, NeurIPS*, vol. 35, 2022, pp. 9361–9373.
- [23] Y. Peng, S. Dalmia, I. Lane, and S. Watanabe, "Branchformer: Parallel MLP-attention architectures to capture local and global context for speech recognition and understanding," in *International Conference on Machine Learning, ICML*, vol. 162, 2022, pp. 17 627–17 643.
- [24] K. Kim, F. Wu, Y. Peng, J. Pan, P. Sridhar, K. J. Han, and S. Watanabe, "E-branchformer: Branchformer with enhanced merging for speech recognition," *IEEE Spoken Language Technology Workshop (SLT)*, pp. 84–91, 2022.
- [25] Z. Yao, L. Guo, X. Yang, W. Kang, F. Kuang, Y. Yang, Z. Jin, L. Lin, and D. Povey, "Zipformer: A faster and better encoder for automatic speech recognition," *CoRR*, vol. abs/2310.11230, 2023.
- [26] Y. Tay, D. Bahri, D. Metzler, D. Juan, Z. Zhao, and C. Zheng, "Synthesizer: Rethinking self-attention for transformer models," in *International Conference on Machine Learning, ICML*, vol. 139, 2021, pp. 10 183–10 192.
- [27] M. Zaheer, G. Guruganesh, K. A. Dubey, J. Ainslie, C. Alberti, S. Ontañón, P. Pham, A. Ravula, Q. Wang, L. Yang, and A. Ahmed, "Big bird: Transformers for longer sequences," in *Neural Information Processing Systems, NeurIPS*, 2020.
- [28] I. Beltagy, M. E. Peters, and A. Cohan, "Longformer: The long-document transformer," *CoRR*, vol. abs/2004.05150, 2020.
- [29] X. Zhang, H. Tang, J. Wang, N. Cheng, J. Luo, and J. Xiao, "Dynamic alignment mask CTC: improved mask CTC with aligned cross entropy," in *ICASSP*, 2023, pp. 1–5.
- [30] M. Burchi and V. Vielzeuf, "Efficient conformer: Progressive downsampling and grouped attention for automatic speech recognition," in *Automatic Speech Recognition and Understanding Workshop, ASRU*. IEEE, 2021, pp. 8–15.
- [31] Y. Zhang, K. C. Puvvada, V. Lavrukhin, and B. Ginsburg, "Conformer-based target-speaker automatic speech recognition for single-channel audio," in *ICASSP*, 2023, pp. 1–5.
- [32] K. Dhawan, K. Rekesh, and B. Ginsburg, "Unified model for code-switching speech recognition and language identification based on concatenated tokenizer," in *Association for Computational Linguistics, ACL*, 2023, pp. 74–82.
- [33] S. Watanabe, T. Hori, S. Karita, T. Hayashi, J. Nishitoba, Y. Unno, N. E. Y. Soplin, J. Heymann, M. Wiesner, N. Chen, A. Renduchintala, and T. Ochiai, "Espnet: End-to-end speech processing toolkit," in *Interspeech*, 2018, pp. 2207–2211.
- [34] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, pp. 770–778, 2015.
- [35] L. J. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *CoRR*, vol. abs/1607.06450, 2016.
- [36] M. Geva, R. Schuster, J. Berant, and O. Levy, "Transformer feed-forward layers are key-value memories," in *Empirical Methods in Natural Language Processing, EMNLP*, 2021, pp. 5484–5495.
- [37] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *International Conference on Machine Learning, ICML*, 2006, pp. 369–376.
- [38] H. Bu, X. Du, X. Na, B. Wu, and H. Zheng, "Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline," *O-COCOSDA*, pp. 1–5, 2017.
- [39] Y. Y. Liu, P. Fung, Y. Yang, C. Cieri, S. Huang, and D. Graff, "Hkust/mts: A very large scale mandarin telephone speech corpus," in *ISCSLP*, 2006, pp. 724–735.
- [40] D. S. Park, W. Chan, Y. Zhang, C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "SpecAugment: A simple data augmentation method for automatic speech recognition," in *Interspeech*, 2019, pp. 2613–2617.
- [41] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations, ICLR*, 2015.
- [42] A. Gulati, J. Qin, 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 *Interspeech*, 2020, pp. 5036–5040.
- [43] K. Shim, J. Choi, and W. Sung, "Exploring attention map reuse for efficient transformer neural networks," *CoRR*, vol. abs/2301.12444, 2023.