

Plug n’ Play Channel Shuffle Module for Enhancing Tiny Vision Transformers

1st Xuwei Xu

*School of Electrical Engineering and Computer Science
The University of Queensland
Brisbane, Australia
xuwei.xu@uq.edu.au*

2nd Sen Wang

*School of Electrical Engineering and Computer Science
The University of Queensland
Brisbane, Australia
sen.wang@uq.edu.au*

3rd Yudong Chen

*School of Electrical Engineering and Computer Science
The University of Queensland
Brisbane, Australia
yudong.chen@uq.edu.au*

4th Jiajun Liu

*DATA61
Commonwealth Scientific and Industrial Research Organisation
Brisbane, Australia
ryan.liu@data61.csiro.au*

Abstract—Vision Transformers (ViTs) have demonstrated remarkable performance in various computer vision tasks. However, the high computational complexity hinders ViTs’ applicability on devices with limited memory and computing resources. Although certain investigations have delved into the fusion of convolutional layers with self-attention mechanisms to enhance the efficiency of ViTs, there remains a knowledge gap in constructing tiny yet effective ViTs solely based on the self-attention mechanism. Furthermore, the straightforward strategy of reducing the feature channels in a large but outperforming ViT often results in significant performance degradation despite improved efficiency. To address these challenges, we propose a novel channel shuffle module to improve tiny-size ViTs, showing the potential of pure self-attention models in environments with constrained computing resources. Inspired by the channel shuffle design in ShuffleNetV2 [1], our module expands the feature channels of a tiny ViT and partitions the channels into two groups: the *Attended* and *Idle* groups. Self-attention computations are exclusively employed on the designated *Attended* group, followed by a channel shuffle operation that facilitates information exchange between the two groups. By incorporating our module into a tiny ViT, we can achieve superior performance while maintaining a comparable computational complexity to the vanilla model. Specifically, our proposed channel shuffle module consistently improves the top-1 accuracy on the ImageNet-1K dataset for various tiny ViT models by up to 2.8%, with the changes in model complexity being less than 0.03 GMACs.

Index Terms—vision transformer, channel shuffle, efficiency

I. INTRODUCTION

Vision Transformers (ViTs) have dominated the computer vision area since the success of [2], demonstrating remarkable performance in image classification [2]–[4], object detection [5]–[7] and segmentation [8]–[10]. However, the high computational burden of the self-attention mechanism makes ViTs less efficient compared to traditional convolutional neural networks (CNNs) [1], [11]–[14] on devices with constrained memory and computing resources. As a result, there is a growing interest in the research community to develop lightweight and efficient ViT models.

TABLE I: Comparisons of pure ViT models on ImageNet [27] validation set. * indicates that the model is re-trained on our machine. The number of feature channels of Swin Transformer [5] only indicates the first stage’s feature channels, which would expand as the layer goes deep

Methods	Feature Channels	Layers	Param.	MACs	Top-1 Acc.
T2T-ViT-7 [26]	256	7	4.3M	1.1G	71.7%
T2T-ViT-14 [26]	384	14	21.5M	4.8G	81.7%
DeiT-Tiny [25]	192	12	5.7M	1.3G	72.2%
DeiT-Small [25]	384	12	21.8M	4.6G	79.9%
Swin-ExtraTiny* [5]	48	12	6.8M	1.1G	74.8%
Swin-Tiny [5]	96	12	28.8M	4.5G	80.8%

Various approaches have been proposed to address this challenge. Some methods integrate efficient convolution operations with computationally expensive self-attentions to create hybrid efficient ViTs [15]–[20]. However, these methods do not fully exploit the potential of pure self-attention models to achieve both high performance and efficiency. Alternatively, certain studies revisit the design principles of efficient CNNs and transfer them to the design of efficient ViTs, such as window-based attention [5], [21], hierarchical network architecture [5], [22], bottleneck structure [23] and spatially separable self-attention [24]. It is worth noting that the channel shuffle design introduced by [13] is less explored in this context. In addition, some powerful ViT models construct their lightweight versions by simply reducing the number of feature channels, layers, or self-attention heads [5], [22], [25], [26]. However, as demonstrated in Table I, such a naive model size reduction often leads to a significant performance drop. For example, DeiT-Tiny [25] suffers a 7.7% top-1 accuracy drop on the ImageNet [27] compared to DeiT-Small when the number of feature channels declines from 384 to 192.

We figure out that one of the main reasons for the per-

formance degradation in tiny ViTs is the limited number of feature channels. The insufficient number of feature channels makes tiny ViTs unable to represent the image effectively. To mitigate the similar issue of insufficient image representations, previous studies in efficient CNNs leverage the concept of grouped convolution [28]–[30], which reduces computational complexity and memory footprint without compromising the total number of feature channels. Besides, [13] proposes a channel shuffle operation to help the information flow across groups. And [1] extensively explores the architecture design and introduces a strategy that splits the channels into two groups, allowing one group to remain idle throughout the layer and shuffling channels between the two groups. It is worth noting that these designs in efficient CNNs are seldom introduced to efficient pure self-attention models.

Hence, in this paper, we present a channel shuffle module specifically designed for tiny ViT models to address the aforementioned challenges. Inspired by [1], [13], our module expands the feature channels of a compressed ViT model and separates them into two groups, namely the *Attended* group and the *Idle* group. In each layer, the *Attended* group performs self-attention computation like a conventional ViT while the *Idle* group remains inactive during the computation. At the end of each layer, a channel shuffle operation is employed to interleave the two feature channel groups and facilitate information exchange. This module serves as a plug-and-play enhancement to tiny ViTs, which improves the performance with merely a bit more computations. Meanwhile, our module is generic and can be applied to both plain and hierarchical ViTs. In this paper, we select DeiT [25] and T2T-ViT [26] as representatives for plain ViTs, and Swin Transformer [5] as the representative for hierarchical ViTs. Moreover, we observe that the *Idle* channels may exhibit different scales compared to the *Attended* channels, resulting in many trivial channels after Layer Normalization. To address this issue, we propose a simple channel re-scaling optimization to alleviate the problem. Extensive experiments have demonstrated the efficacy and efficiency of our module.

We summarize the key contributions of our work as follows:

- We develop an efficient channel shuffle module to enhance tiny ViTs with very few additional computations, satisfying the environment with constrained computing resources.
- Our module can work as an independent plug-and-play component to the vanilla tiny ViTs and is generic for both plain and hierarchical ViTs.
- We introduce a simple channel re-scaling method to mitigate the problem of distinct scales between the *Attended* and *Idle* groups.
- Extensive experimental results have shown the efficiency and efficacy of our proposed module.

To our best knowledge, this is the first work improving tiny-scale efficient ViTs by enriching channel-wise information while maintaining computational complexity.

II. RELATED WORK

A. Efficient convolutional neural networks

Efficient convolutional neural networks (CNNs) have attracted significant attention due to the need for deployment on devices with constrained computing resources. AlexNet [28] proposes grouped convolution to distribute the model over multiple GPUs and consequently reduce the computational complexity and memory footprint on each single GPU. ResNeXt [29] validates the efficacy of grouped convolution, showing an improved accuracy by grouped convolution. GoogLeNet [31] leverages grouped convolution to establish the inception module, which successfully expands the width and depth of a CNN model while keeping the computational budget constant. Xception [30] enforces the number of channel groups the same as the number of feature channels and brings up the concept of depth-wise separable convolution. MobileNets [11], [12] extensively utilize depth-wise separable convolution to reduce the computational complexity of CNNs for mobile vision applications. ShuffleNet [13] puts forth the concept of channel shuffle for grouped convolution, which realizes efficient information exchange between channel groups by shuffling the channels. These efficiency designs are less explored in ViTs than in CNNs.

B. Vision Transformers

Vision Transformers (ViTs) have gained significant attention in the field of computer vision as a promising alternative to CNNs. The original ViT architecture [2] demonstrates the effectiveness of the self-attention mechanism for image classification, which is capable of capturing global relationships. In general, there are two types of ViT architectures, namely plain and hierarchical, which are distinguished by whether token downsampling is adopted in the network. Plain ViTs [3], [25], [26] have the same backbone architecture as the vanilla ViT that the number of tokens and feature channels keeps static throughout the network. Hierarchical ViTs [5], [6], [21], [22] apply token downsampling between stages to enable multi-scale self-attention. As a result, the number of tokens decreases and the number of feature channels increases as the layer goes deep in hierarchical ViTs. In this work, we choose DeiT [25] and T2T-ViT [26] as representatives for plain ViTs, and Swin Transformer [5] as the representative for hierarchical ViTs.

C. Efficient Vision Transformers

One direction of ViT research focuses on improving the efficiency of ViTs, as their computational complexity can be prohibitive for resource-constrained devices. Several approaches have been explored, such as distillation methods that transfer knowledge from large ViT models to smaller ones [25], spatial-wise token pruning [32]–[34], and regional self-attention design [5], [21]. These methods mainly concentrate on compressing a powerful ViT into a smaller counterpart without compromising much performance, while our method attempts to enhance the tiny ViTs. Besides, some studies [15]–[20] integrate convolution with self-attention to achieve efficient ViTs. However, these methods fail to reveal the

potential of pure self-attention models in an environment with limited computing resources. In contrast, our module demonstrates the possibility of efficient and high-performance pure self-attention models.

III. METHODS

A. Preliminaries

The vanilla ViT [2] first splits the input image into patches and then linearly projects the image patches into image tokens. These tokens serve as the input for subsequent computations, enabling the model to capture global contextual information. We denote the input feature map of layer i as $X_i \in \mathbb{R}^{N \times C}$, where N and C are the numbers of tokens and feature channels, respectively. Each ViT layer comprises a multi-head self-attention (MHSA) module and a feed-forward network (FFN) module. For the MHSA module, it linearly transforms the input feature map into three matrices called *Key* (K_i), *Query* (Q_i) and *Value* (V_i) by

$$K_i = X_i W_{k_i}, \quad Q_i = X_i W_{q_i}, \quad V_i = X_i W_{v_i}, \quad (1)$$

where W_{k_i} , W_{q_i} and W_{v_i} are the corresponding weights and the bias terms are omitted. Next, it computes the attention map $A_i \in \mathbb{R}^{N \times N}$ by a dot production with the softmax activation between *Key* and *Query* as

$$A_i = \text{softmax}\left(\frac{Q_i K_i^\top}{\sqrt{d_k}}\right), \quad (2)$$

where d_k is the dimension of K and is usually the same as the channel dimension C . The attention map reflects the similarities between each pair of tokens. Eventually, MHSA calculates the attended output with residual connection as

$$X'_i = \text{MHSA}(X_i) + X_i = A_i V_i W_m + X_i, \quad (3)$$

where $W_m \in \mathbb{R}^{C \times C}$ is the learnable weight. After the MHSA module, a two-layer Multi-Layer Perceptron (MLP) is utilized as the FFN module to self-activate each token by

$$X_{i+1} = \text{FFN}(X'_i) + X'_i = X'_i W_{f1} W_{f2} + X'_i, \quad (4)$$

where $W_{f1} \in \mathbb{R}^{C \times \mu C}$ and $W_{f2} \in \mathbb{R}^{\mu C \times C}$ are two learnable projection weights, μ is the channel expansion ratio, $X_{i+1} \in \mathbb{R}^{N \times C}$ is the output of the i^{th} ViT layer and the bias terms are omitted. Moreover, ViT introduces pre-module Layer Normalization (LN) [35] on the feature map to improve both the training time and the generalization performance.

B. Channel shuffle module

Fig. 1 provides an overview of our proposed channel shuffle module, which aims to enhance the capabilities of tiny ViT models. To improve the feature representation ability, we start by doubling the feature channels of the input feature map $X_i \in \mathbb{R}^{N \times C}$ during the token embedding phase, resulting in $X_i^{\text{Double}} \in \mathbb{R}^{N \times 2C}$. In the i^{th} layer, the feature map X_i^{Double} is partitioned into two groups along the feature channel dimension, namely the *Attended* group $X_i^{\text{Attn}} \in \mathbb{R}^{N \times C}$ and *Idle* group $X_i^{\text{Idle}} \in \mathbb{R}^{N \times C}$. The *Attended* group only occupies half of the channels and participates in this layer's

calculations, resulting in the attended output $X_{i+1}^{\text{Attn}} \in \mathbb{R}^{N \times C}$. On the contrary, the *Idle* group holds the rest half channels and maintains the same until the end of this layer so that $X_i^{\text{Idle}} = X_{i+1}^{\text{Idle}} \in \mathbb{R}^{N \times C}$. At the end of layer i , the two groups are concatenated together as

$$X_{i+1}^{\text{Double}} = \text{concat}(X_{i+1}^{\text{Attn}}, X_{i+1}^{\text{Idle}}), \quad (5)$$

where the $X_{i+1}^{\text{Double}} \in \mathbb{R}^{N \times 2C}$ is the output of layer i . Finally, we apply channel shuffle on X_{i+1}^{Double} to enforce information exchange between the two groups.

C. Channel re-scaling

We have identified a potential issue related to the residual connections [36] employed in the MHSA and FFN modules, which can result in significant scale differences between X_{i+1}^{Attn} and X_{i+1}^{Idle} . This discrepancy can lead to many trivial values after Layer Normalization, particularly in deeper layers. To address this issue, we devise a channel re-scaling approach for X_{i+1}^{Attn} . Specifically, we modify the Transformer layer described in Equations 3 and 4 as follows:

$$\begin{aligned} X_i^{\text{Attn}'} &= \text{MHSA}(X_i^{\text{Attn}}) + \alpha_1 X_i^{\text{Attn}}, \\ X_{i+1}^{\text{Attn}'} &= \text{FFN}(X_i^{\text{Attn}'}) + \alpha_2 X_i^{\text{Attn}'}, \end{aligned} \quad (6)$$

where $\alpha_1, \alpha_2 \in \mathbb{R}^C$ are two learnable coefficients. We can use this simple modification to enable our module automatically control the scale of the feature representation in the *Attended* group.

D. Computational complexity analysis

We begin by comparing the theoretical computational complexity of a vanilla ViT and a tiny ViT equipped with our proposed module, assuming both models have an equal number of total feature channels. In the case of the vanilla ViT, the computational complexity per layer can be represented as

$$\Omega(\text{vanilla}) = (4 + 2\mu)NC^2 + 2N^2C, \quad (7)$$

where N, C and μ denote the number of tokens, feature channels and the expansion ratio of the MLP, respectively. On the other hand, the computational complexity of a tiny ViT with our channel shuffle module is given by

$$\Omega(\text{shuffle}) = \left(1 + \frac{\mu}{2}\right)NC^2 + N^2C + NC. \quad (8)$$

Notably, the computational complexity of our shuffle module is significantly smaller than that of the vanilla ViT. Therefore, when considering models with an equal number of total feature channels, indicating a similar feature representation capacity, our proposed module proves to be more computationally efficient and suitable for deployment on resource-constrained devices.

Furthermore, we analyze the overall computational complexity of a tiny ViT model with and without the channel shuffle module. Since the *Idle* group does not participate in computations, the increase in computational complexity per layer in our module solely stems from the channel re-scaling, which amounts to NC . Consequently, the total increase in computational complexity over L layers is LNC .

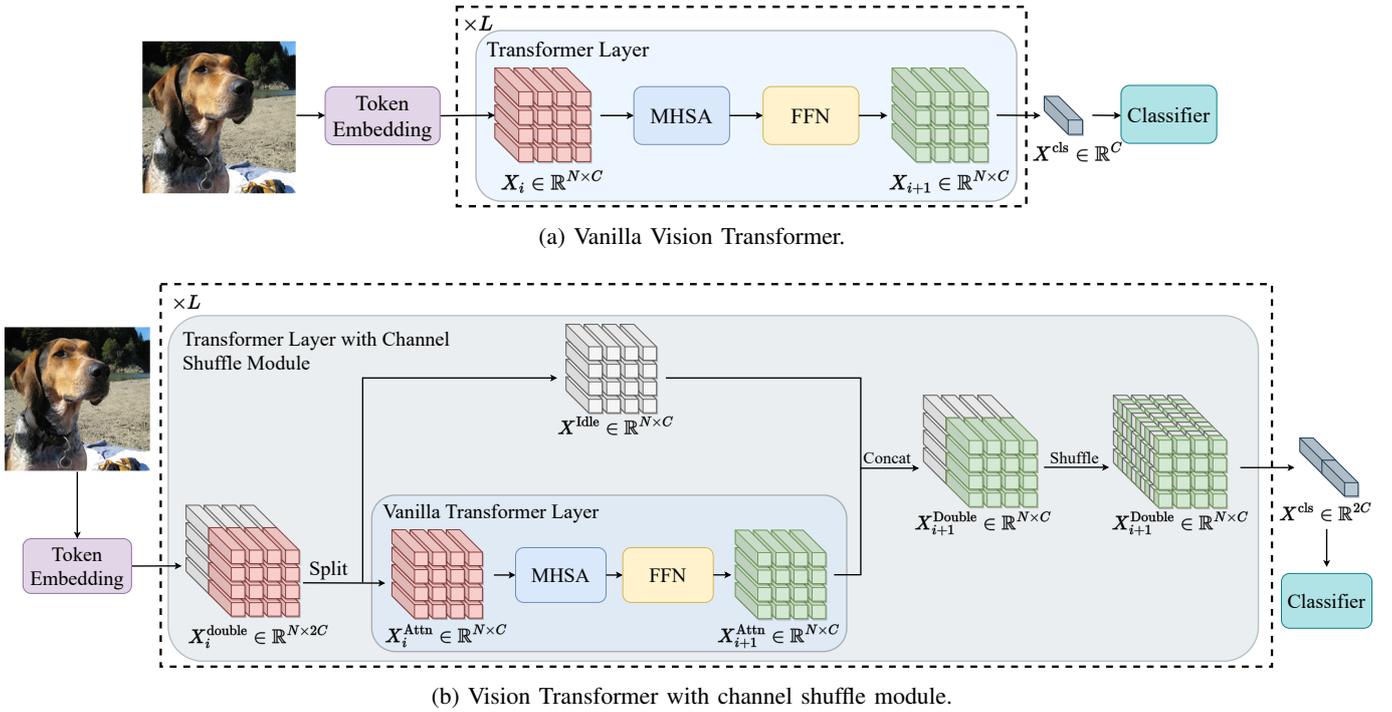


Fig. 1: Overview of the channel shuffle module. In each ViT layer, the *Attended* group participates in the computations while the *Idle* group retains the same until the end of the layer. At the end of each layer, we concatenate the two groups and shuffle the channels to facilitate information exchange.

Additionally, our module doubles the channels during the token embedding phase, resulting in additional $3P^2NC$ computations, where P represents the size of the image patches. Moreover, in the final layer, our module introduces an extra SC computations, with S denoting the number of classes. As a result, the total increase in computational complexity is given by $LNC + 3P^2NC + SC$, which is relatively small compared to the overall computational complexity of $3P^2NC + (4 + 2\mu)LNC^2 + 2LN^2C + SC$. Taking DeiT-Tiny [25] as an example, the total computational complexity is approximately 1.25GMACs while the extra computational complexity brought up by our channel shuffle module is merely 0.03GMACs, which is about 2% of the total.

E. Plug-and-play module

An important contribution of this module is its ability to address the performance degradation of tiny ViT models. The channel shuffle module can be easily incorporated into a tiny ViT without introducing significant modifications to the backbone architecture. As depicted in Fig. 1, the *Idle* group does not participate in calculations, making it straightforward to apply the module to different variants of ViT architectures.

When integrating the channel shuffle module into a plain ViT, the number of feature channels is doubled during the image token embedding phase and remains constant throughout the network. However, in hierarchical ViTs, the computational complexity can be proportional to the square of the number of feature channels in downsampling layers. Our channel

shuffle module could lead to a non-negligible increase in computations. Taking Swin Transformer [5] as an instance, the computational complexity of a fully-connected downsampling layer is $2NC^2$. However, if the feature channels are doubled in the channel shuffle module, the computational complexity of downsampling would increase to $8NC^2$ per downsampling layer. To address this issue, we perform separate downsampling operations on the *Attended* and *Idle* groups to prevent an excessive computational burden when applying the channel shuffle module to hierarchical ViTs.

IV. EXPERIMENTS

A. Dataset and settings

Dataset. We choose ImageNet-1K [27] as the target dataset, which contains around 1.28 million images for training and 50 thousand images for validation. It is an acknowledged standard dataset for model benchmarking.

Base model. DeiT-Tiny [25] and T2T-ViT-7 [26] are selected as representatives for tiny ViTs with a plain architecture. Swin Transformer [5] is chosen to be the representative for hierarchical ViTs. Since Swin Transformer does not officially provide a mobile-level version whose computational complexity is approximately 1GMACs, we scale down the Swin-Tiny by reducing the number of its feature channels in the first stage from 96 to 48 and subsequently construct a Swin-ExtraTiny model. These are well-known in ViT families for their excellent standard performance and data-efficient

TABLE II: Effectiveness of the channel shuffle module. We compare the top-1 accuracy on ImageNet [27] of tiny ViT models equipped with the channel shuffle module and its vanilla version.

Methods	Feature Channels	Layers	Param.	MACs	Top-1 Acc.
T2T-ViT-7 [26]	256	7	4.3M	1.1G	71.7%
Shuffled T2T-ViT-7	512	7	4.7M	1.2G	74.4% (+2.7%)
DeiT-Tiny [25]	192	12	5.7M	1.3G	72.2%
Shuffled DeiT-Tiny	384	12	6.1M	1.3G	74.4% (+2.2%)
Swin-ExtraTiny [5]	48	12	6.8M	1.1G	74.8%
Shuffled Swin-ExtraTiny	96	11	7.2M	1.0G	77.8% (+3.0%)

training. Furthermore, we also reproduce the experiments on Swin-Ti to explore the influence of base model size.

Training configurations. We follow the image augmentations and training recipes in [25] and its official GitHub repository for all the models, except setting the learning rate to 5e-3, the total batch size to 4096.

B. Main results

Table II provides comparisons between pure self-attention models with and without our channel shuffle module in terms of accuracy, the number of parameters and computational complexity. The results clearly show that the channel shuffle module significantly enhances the performance of tiny ViTs in the classification task, improving the accuracy by 2.2~3.0% without a significant increase in computational budgets. In the case of DeiT-Tiny, the shuffled version achieves a 2.2% higher top-1 accuracy while maintaining an equivalent computational cost to the vanilla version. As for Swin-ExtraTiny, since we eliminate one layer, the shuffled version (1.0GMACs) runs faster than the unshuffled version (1.1GMACs) with a remarkable 3.0% higher top-1 accuracy. These results highlight the effectiveness and generalizability of our simple yet powerful design module.

Table III provides comprehensive comparisons between our channel shuffle module, other efficient CNNs, ViTs, and hybrid models. Firstly, our channel shuffle module consistently outperforms all other pure ViT models in terms of the trade-off between model complexity and accuracy, highlighting the effectiveness of our module. Secondly, when compared to hybrid models that combine convolution and self-attention, our module achieves comparable or even better performance. For instance, at the same model complexity of 1.1GMACs, Swin-ExtraTiny with the channel shuffle module outperforms EdgeNeXt-XS by 2.8%. This result demonstrates the potential of efficient pure self-attention networks. However, it is worth noting that when compared to mobile-friendly CNNs, self-attention-based networks still exhibit lower efficiency when achieving the same performance.

C. Feature channel analysis

The channel shuffle module plays a critical role in enhancing the ability of a tiny ViT model to represent image features. Even though half of the channels do not actively participate in

TABLE III: Comparisons against tiny-size CNNs and ViTs.

Methods	Type	Param.	MACs	Top-1 Acc.
MobileNet V2 [12]	CNN	6.9M	0.6G	74.7%
ShuffleNet V2 [1]	CNN	7.4M	0.6G	74.9%
EfficientNet-B0 [14]	CNN	5.3M	0.4G	76.3%
EfficientNet-B1 [14]	CNN	7.8M	0.7G	79.1%
RegNetY-800MF [37]	CNN	6.3M	0.8G	76.3%
RegNetY-1.6GF [37]	CNN	11.2M	1.6G	78.0%
T2T-ViT-7 [26]	ViT	4.3M	1.1G	71.7%
HVT-Ti-1 [38]	ViT	5.6M	0.7G	69.6%
DeiT-Tiny [25]	ViT	5.7M	1.3G	72.2%
PVTv2-B0 [39]	ViT	3.4M	0.6G	70.5%
PVTv2-B1 [39]	ViT	13.1M	2.1G	78.7%
PVTv1-Tiny [22]	ViT	13.2M	1.9G	75.1%
AutoFormer-Tiny [40]	ViT	5.7M	1.3G	74.7%
PiT-Ti [41]	Hybrid	4.9M	0.7G	74.6%
ConViT-Ti [20]	Hybrid	6.0M	1.0G	73.1%
Visformer-Ti [15]	Hybrid	10.3M	1.3G	78.6%
MobileViTv1-XS [18]	Hybrid	2.3M	0.9G	74.8%
MobileViTv1-S [18]	Hybrid	5.6M	2.0G	78.4%
EdgeNeXt-XS [42]	Hybrid	2.3M	1.1G	75.0%
EdgeNeXt-S [42]	Hybrid	5.6M	2.6G	78.4%
Shuffled T2T-ViT-7 (ours)	ViT	4.7M	1.2G	74.4%
Shuffled DeiT-Tiny (ours)	ViT	6.1M	1.3G	74.4%
Shuffled Swin-ExtraTiny (ours)	ViT	7.2M	1.0G	77.8%

the Transformer layer, they still serve two important functions: propagating gradients across layers and enriching the available image features.

To validate this argument, we visualize and compare the feature channel distributions with respect to the four stages in the shuffled and vanilla Swin-ExtraTiny in Fig. 2. We employ t-distributed stochastic neighbour embedding (t-SNE) [43] for dimensionality reduction. In Fig. 2(a), it is evident that in the early stage with a small number of channels (e.g., 48 for Swin-ExtraTiny and 96 for shuffled Swin-ExtraTiny), the shuffled model (orange) exhibits a more diverse distribution compared to the unshuffled model (blue). This indicates that the shuffled model benefits from our module by incorporating richer information in the early stage. However, as the model progresses to deeper stages and the number of feature channels increases, the distribution variances between the shuffled and unshuffled models become less distinct, as shown in Fig. 2(c) and 2(d). We think this is the main reason why our channel shuffle module can improve the performance of tiny ViT models.

However, we also point out that the improvement of this module diminishes as the model size increases. Larger ViT models already possess sufficient feature channels to effectively represent the image, and excessively oversized channels may even have a negative influence. For example, the shuffled Swin-Ti (28.8M, 4.3G MACs) only reaches a slightly higher accuracy at 81.4% than its vanilla version (28.8M, 4.5G MACs) at 80.8%. We argue that this module is most beneficial for tiny models that lack feature representations. As the feature representations get more complicated when the model scales up, the influence of this module vanishes.

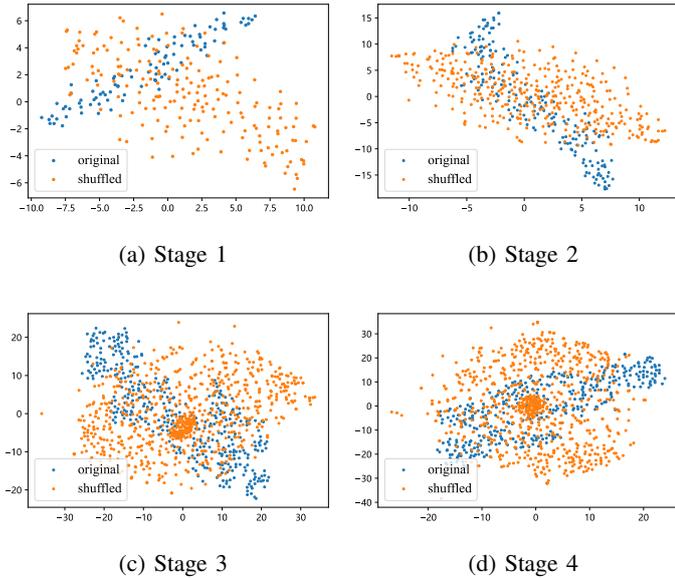


Fig. 2: Channel distributions of feature maps at the end of the four stages in Swin-ExtraTiny. The blue dots represent the distribution of the vanilla Swin-ExtraTiny while the orange dots stand for the shuffled version. Since Swin-ExtraTiny is a hierarchical vision Transformer, the number of channels (i.e., the number of dots in this figure) increases in deeper stages.

TABLE IV: Ablation study on the two components

Method	channel shuffle	channel re-scale	#Params	#MACs	Top-1 Acc.
T2T-ViT-7	✓		4.3M	1.1G	71.7%
	✓		4.7M	1.3G	72.5%
		✓	4.7M	1.3G	74.4%
DeiT-Tiny	✓		5.7M	1.3G	72.2%
	✓		6.1M	1.3G	72.9%
		✓	6.1M	1.3G	74.4%
Swin-ExtraTiny	✓		6.9M	1.1G	74.8%
	✓		7.2M	1.0G	75.8%
		✓	7.2M	1.0G	77.8%

D. Ablation study

We propose two crucial components in the channel shuffle module: the independent channel shuffle process and the channel re-scaling. In this section, we conduct ablation studies to evaluate the impact of these two design choices.

Channel shuffle Table IV demonstrates that without re-scaling, simply integrating the channel shuffle module into the vision Transformer leads to minor performance improvements. This is attributed to the fact that the channel shuffle process can introduce more imbalanced features when the number of feature channels is small. Despite marginal performance enhancement, merely adopting the channel shuffle module still helps to reach higher accuracy.

Channel re-scaling Table IV highlights the importance of channel re-scaling in this module, as it improves the outcome of the channel shuffle. For example, the shuffled Swin-

ExtraTiny model with channel re-scaling surpasses both the original version and the non-scaling version, achieving a 3.0% and 2.0% increase in top-1 accuracy, respectively.

V. CONCLUSION

This paper presents an efficient module designed to enhance tiny ViT models. The proposed channel shuffle module expands the number of feature channels of a tiny ViT to improve its feature representation ability. In each layer, the feature channels are partitioned into two groups called the *Attended* group and the *Idle* group. Only the *Attended* group participates in each layer’s calculation while the *Idle* group maintains the same until the end of the layer. The channel shuffle module effectively leverages channel shuffle operation to exchange information between the two groups, which contributes to enriched channel-wise information without introducing significant additional computational complexity. Experimental results demonstrate the effectiveness and generalizability of our module in improving both plain and hierarchical tiny ViTs.

REFERENCES

- [1] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, “Shufflenet v2: Practical guidelines for efficient cnn architecture design,” in *ECCV*, 2018.
- [2] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *ICLR*, 2021.
- [3] X. Zhai, A. Kolesnikov, N. Houlsby, and L. Beyer, “Scaling vision transformers,” in *CVPR*, 2022.
- [4] M. Ding, B. Xiao, N. Codella, P. Luo, J. Wang, and L. Yuan, “Davvit: Dual attention vision transformers,” in *ECCV*, 2022.
- [5] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *ICCV*, 2021.
- [6] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong *et al.*, “Swin transformer v2: Scaling up capacity and resolution,” in *CVPR*, 2022.
- [7] Y. Li, H. Mao, R. Girshick, and K. He, “Exploring plain vision transformer backbones for object detection,” in *ECCV*, 2022.
- [8] J. Yang, C. Li, P. Zhang, X. Dai, B. Xiao, L. Yuan, and J. Gao, “Focal attention for long-range interactions in vision transformers,” in *NeurIPS*, 2021.
- [9] Z. Chen, Y. Duan, W. Wang, J. He, T. Lu, J. Dai, and Y. Qiao, “Vision transformer adapter for dense predictions,” in *ICLR*, 2023.
- [10] Y. Fang, W. Wang, B. Xie, Q. Sun, L. Wu, X. Wang, T. Huang, X. Wang, and Y. Cao, “Eva: Exploring the limits of masked visual representation learning at scale,” in *CVPR*, 2023.
- [11] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, “Mobilenets: Efficient convolutional neural networks for mobile vision applications,” 2017, unpublished.
- [12] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *CVPR*, 2018.
- [13] X. Zhang, X. Zhou, M. Lin, and J. Sun, “Shufflenet: An extremely efficient convolutional neural network for mobile devices,” in *CVPR*, 2018.
- [14] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *ICML*, 2019.
- [15] Z. Chen, L. Xie, J. Niu, X. Liu, L. Wei, and Q. Tian, “Visformer: The vision-friendly transformer,” in *ICCV*, 2021.
- [16] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, “Cvt: Introducing convolutions to vision transformers,” in *ICCV*, 2021.
- [17] Z. Dai, H. Liu, Q. V. Le, and M. Tan, “Coatnet: Marrying convolution and attention for all data sizes,” in *NeurIPS*, 2021.
- [18] S. Mehta and M. Rastegari, “Mobilevit: light-weight, general-purpose, and mobile-friendly vision transformer,” in *ICLR*, 2022.
- [19] B. Graham, A. El-Nouby, H. Touvron, P. Stock, A. Joulin, H. Jégou, and M. Douze, “Levit: a vision transformer in convnet’s clothing for faster inference,” in *ICCV*, 2021.

- [20] S. d’Ascoli, H. Touvron, M. L. Leavitt, A. S. Morcos, G. Biroli, and L. Sagun, “Convit: Improving vision transformers with soft convolutional inductive biases,” in *ICML*, 2021.
- [21] X. Dong, J. Bao, D. Chen, W. Zhang, N. Yu, L. Yuan, D. Chen, and B. Guo, “Cswin transformer: A general vision transformer backbone with cross-shaped windows,” in *CVPR*, 2022.
- [22] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, “Pyramid vision transformer: A versatile backbone for dense prediction without convolutions,” in *ICCV*, 2021.
- [23] A. Srinivas, T.-Y. Lin, N. Parmar, J. Shlens, P. Abbeel, and A. Vaswani, “Bottleneck transformers for visual recognition,” in *CVPR*, 2021.
- [24] X. Chu, Z. Tian, Y. Wang, B. Zhang, H. Ren, X. Wei, H. Xia, and C. Shen, “Twins: Revisiting the design of spatial attention in vision transformers,” in *NeurIPS*, 2021.
- [25] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, “Training data-efficient image transformers & distillation through attention,” in *ICML*, 2021.
- [26] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z.-H. Jiang, F. E. Tay, J. Feng, and S. Yan, “Tokens-to-token vit: Training vision transformers from scratch on imagenet,” in *ICCV*, 2021.
- [27] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *CVPR*, 2009.
- [28] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *NeurIPS*, 2012.
- [29] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *CVPR*, 2017.
- [30] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *CVPR*, 2017.
- [31] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *CVPR*, 2015.
- [32] Y. Rao, W. Zhao, B. Liu, J. Lu, J. Zhou, and C.-J. Hsieh, “Dynamicvit: Efficient vision transformers with dynamic token sparsification,” in *NeurIPS*, 2021.
- [33] Y. Liang, G. Chongjian, Z. Tong, Y. Song, J. Wang, and P. Xie, “Evit: Expediting vision transformers via token reorganizations,” in *ICLR*, 2021.
- [34] Y. Xu, Z. Zhang, M. Zhang, K. Sheng, K. Li, W. Dong, L. Zhang, C. Xu, and X. Sun, “Evo-vit: Slow-fast token evolution for dynamic vision transformer,” in *AAAI*, 2022.
- [35] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” 2016, unpublished.
- [36] K. He *et al.*, “Deep residual learning for image recognition,” in *CVPR*, 2016.
- [37] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. He, and P. Dollár, “Designing network design spaces,” in *CVPR*, 2020.
- [38] Z. Pan, B. Zhuang, J. Liu, H. He, and J. Cai, “Scalable vision transformers with hierarchical pooling,” in *ICCV*, 2021.
- [39] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, “Pvt v2: Improved baselines with pyramid vision transformer,” *Computational Visual Media*, 2022.
- [40] M. Chen, H. Peng, J. Fu, and H. Ling, “Autoformer: Searching transformers for visual recognition,” in *ICCV*, 2021.
- [41] B. Heo, S. Yun, D. Han, S. Chun, J. Choe, and S. J. Oh, “Rethinking spatial dimensions of vision transformers,” in *ICCV*, 2021.
- [42] M. Maaz, A. Shaker, H. Cholakkal, S. Khan, S. W. Zamir, R. M. Anwer, and F. Shahbaz Khan, “Edgenext: efficiently amalgamated cnn-transformer architecture for mobile vision applications,” in *ECCV*, 2022.
- [43] G. Hinton and S. T. Roweis, “Stochastic neighbor embedding,” in *NeurIPS*, 2002.