

# FLatten Transformer: Vision Transformer using Focused Linear Attention

Dongchen Han\* Xuran Pan\* Yizeng Han Shiji Song Gao Huang†

Department of Automation, BNRist, Tsinghua University

## Abstract

The quadratic computation complexity of self-attention has been a persistent challenge when applying Transformer models to vision tasks. Linear attention, on the other hand, offers a much more efficient alternative with its linear complexity by approximating the Softmax operation through carefully designed mapping functions. However, current linear attention approaches either suffer from significant performance degradation or introduce additional computation overhead from the mapping functions. In this paper, we propose a novel **Focused Linear Attention** module to achieve both high efficiency and expressiveness. Specifically, we first analyze the factors contributing to the performance degradation of linear attention from two perspectives: the focus ability and feature diversity. To overcome these limitations, we introduce a simple yet effective mapping function and an efficient rank restoration module to enhance the expressiveness of self-attention while maintaining low computation complexity. Extensive experiments show that our linear attention module is applicable to a variety of advanced vision Transformers, and achieves consistently improved performances on multiple benchmarks. Code is available at <https://github.com/LeapLabTHU/FLatten-Transformer>.

## 1. Introduction

Recent years have witnessed the vast development of Transformer and self-attention in the field of computer vision. With the advent of Vision Transformer [11, 39], self-attention techniques have shown great potential in a variety of vision tasks including image classification [41, 43, 30, 46], semantic segmentation [6, 49], object detection [4, 61, 22], and multi-modal tasks [35, 31].

However, applying Transformer to vision models is a non-trivial task. Unlike lightweight convolution neural networks [37, 16, 44, 33], the quadratic computation complexity  $\mathcal{O}(n^2)$  with respect to sequence length  $n$  leads to high computation costs when employing self-attention with

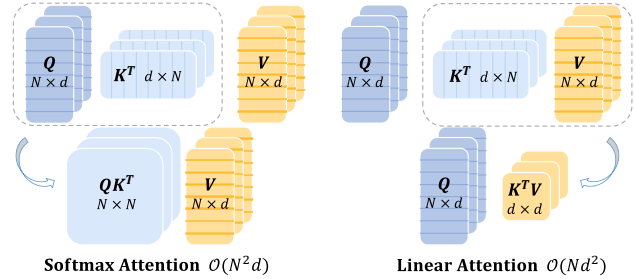


Figure 1. **Difference between Softmax attention and Linear attention.**  $Q, K, V \in \mathbb{R}^{N \times d}$  denote query, key and value matrix respectively. Softmax attention compels to compute the pairwise similarity between queries and keys, and results in the complexity of  $\mathcal{O}(N^2d)$ . Linear attention manages to decouple the Softmax operation with proper approximation and change the computation order by computing  $K^T V$  first, which leads to the complexity of  $\mathcal{O}(Nd^2)$ . Considering that channel dimension  $d$  is usually smaller than token number  $N$  in modern vision Transformer designs, e.g.,  $d = 64, N = 196$  in DeiT [39] and  $d = 32, N = 49$  in Swin Transformer [24], linear attention modules practically save the overall computation cost while can also enjoy the benefits of a larger receptive field and higher throughput.

a global receptive field. Previous works have sought to mitigate this challenge by confining the global receptive field to a smaller region, such as designing sparse global attention patterns [41, 46] or applying smaller attention windows [24, 17]. Albeit effective, these methods are either prone to disregarding informative features in other regions due to their attention patterns or inevitably sacrifice the ability to model long-range dependencies.

Linear attention, on the other hand, has been considered a simple yet effective alternative to address the computation dilemma by reducing the general complexity. Early research leverages a locally-sensitive hashing scheme [21] that compresses the computation complexity from  $\mathcal{O}(n^2)$  to  $\mathcal{O}(n \log(n))$ . Nevertheless, it introduces a large constant before the complexity term, which makes it still unaffordable under common cases. More recent studies have noticed that the utilization of Softmax function in the self-attention operation practically compels a pairwise computation between all queries and keys, resulting in the predominant  $\mathcal{O}(n^2)$  complexity. To tackle this, several approaches adopt sim-

\*Equal contribution.

†Corresponding Author.

ple activation functions [19, 38] or tailored mapping functions [7, 26] to approximate the original Softmax function. As illustrated in Fig. 1, by changing the computation order from (query·key)-value to query·(key·value), the overall computation complexity can be reduced to  $\mathcal{O}(n)$ . However, compared to Softmax attention, current linear attention approaches still suffer from severe performance drop and may involve additional computation overhead from the mapping function, thereby constraining their practical application.

In this paper, we target on the limitations of current linear attention approaches and propose a novel **Focused Linear Attention** module, which achieves both high efficiency and expressiveness. Specifically, we undertake a dual-pronged analysis of the factors contributing to the performance decline in linear attention and subsequently propose corresponding solutions. First, the distribution of attention weight in the former linear attention modules is relatively smooth, lacking the focus ability to address the most informative features. As a remedy, we propose a simple mapping function to adjust the feature direction of queries and keys, making the attention weights more distinguishable. Second, we notice that the diminished rank of the attention matrix curtails the diversity of features in linear attention. To address this, we propose a rank restoration module by applying an additional depthwise convolution (DWC) to the original attention matrix, which helps to restore the matrix rank and keeps the output feature of different positions diversified. Leveraging these improved techniques, our module demonstrates comparable or superior performance to its Softmax counterparts, while enjoying the benefits of low computation complexity.

We empirically validate the effectiveness of our module on image classification, semantic segmentation, and object detection tasks using five advanced vision Transformer models. The results demonstrate consistent improvements over all baselines and other linear attention approaches.

## 2. Related Works

### 2.1. Vision Transformer

Transformer and self-attention mechanism are first introduced in the field of natural language processing and have earned wide research interest in computer vision. Nevertheless, the high computation complexity of self-attention set constraints on the direct application to vision tasks. Previous works have attempted to address this concern from several perspectives. The pioneer Vision Transformer [11] considers reducing the input resolution by merging neighbouring pixels into a single token. Similar insights have been adopted in the following researches [55, 54] and also extend to downstream tasks [22]. Another line of research reduces the feature resolution gradually and adopts carefully designed attention patterns to constrain the number

of attentive tokens. For instance, PVT [41, 42] uses a sparse attention pattern and selects attentive tokens from a global perspective. DAT [46] follows the path and designs a deformable attention module to achieve data-dependent attention pattern. Swin Transformer [24] selects attentive tokens locally by dividing input into isolated windows. NAT [17] follows the query-centric pattern in convolution and designs independent attentive tokens for all queries. Some researches also notice that convolution operations are valuable to Transformer models and may help to improve the overall efficiency [48]. CMT [12] combines Transformer blocks with efficient convolution operators like depthwise convolution [37], and achieves better efficiency-performance trade-off. ACmix [30] shares the computation overhead of convolution and self-attention, and integrates both modules with limited cost. Methods have also been proposed for the efficient training of Transformers [45, 29]. In application scenarios demanding high efficiency, MobileFormer [5] maintains two paths for convolution and Transformer respectively and enjoys the benefit from both modules. Dyn-Perceiver [13] achieves efficient visual recognition through dynamic early exiting [15, 14, 51]. MobileViT [28] takes advantage of the success of MobileNets [37] and uses the combination of mobilenet blocks and Transformer blocks to achieve light-weight and low latency.

However, these approaches still relied on the Softmax operator, whose inherit high computation complexity inevitably results in the inconvenience in model architecture design and practical application.

### 2.2. Linear Attention

Apart from the above methods, another line of research addresses high computation complexity with linear attention [19]. Specifically, linear attention replaces the Softmax function in self-attention with separate kernel functions. In this case, linear attention does not have to compute the pairwise similarity  $QK^T$  first. As illustrated in Fig. 1, based on the associative property of matrix multiplication, linear attention can change the computation order by computing  $K^T V$  first, thus reducing the computation complexity from  $\mathcal{O}(N^2 d)$  to  $\mathcal{O}(Nd^2)$ . Though efficient, how to design linear attention module as effective as softmax attention is a nontrivial problem. Performer [7] approximates the Softmax operation with orthogonal random features. Efficient attention [38] applies Softmax function to  $Q$  and  $K$  respectively, which naturally ensures each row of  $QK^T$  sums up to 1. Nyströmformer [50] and SOFT [26] approximate the full self-attention matrix via matrix decomposition. Hydra attention [1] replaces Softmax with cosine similarity and proposes hydra trick which reduces the computation complexity to  $\mathcal{O}(Nd)$ . EfficientVit [2] uses depth-wise convolution to improve linear attention’s local feature extraction capacity. Castling-ViT [52] proposes linear angular kernel

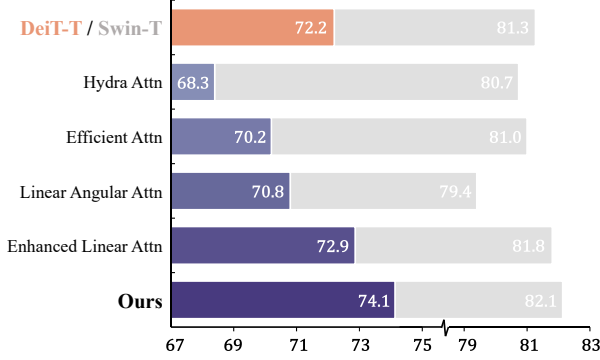


Figure 2. Comparison of different linear attention designs on DeiT-Tiny and Swin-Tiny structures.

to measure spectral similarity between each  $Q_i$  and  $K_j$ .

Nevertheless, current linear attention designs either do not have enough expressive capability to catch up with Softmax attention or involve additional computation overhead from the complex kernel function. In this work, we analyze the reasons for the performance drop of linear attention from the focus ability and feature diversity perspectives. Based on these analyses, we propose a novel linear attention module called focused linear attention which achieves better performance than Softmax attention with lower computation complexity (Fig. 2).

### 3. Preliminaries

#### 3.1. Vision Transformer and Self-Attention

We first revisit the general form of self-attention in Vision Transformers. Given the input  $N$  tokens  $x \in \mathbb{R}^{N \times C}$ , within each head, self-attention can be written as:

$$Q = xW_Q, K = xW_K, V = xW_V, \quad (1)$$

$$O_i = \sum_{j=1}^N \frac{\text{Sim}(Q_i, K_j)}{\sum_{j=1}^N \text{Sim}(Q_i, K_j)} V_j,$$

where  $W_Q, W_K, W_V \in \mathbb{R}^{C \times C}$  are projection matrices and  $\text{Sim}(\cdot, \cdot)$  denotes the similarity function. Modern vision Transformers mainly adopt Softmax attention [40] where similarity is measured as  $\text{Sim}(Q, K) = \exp(QK^T / \sqrt{d})$ . In this case, the attention map is obtained by computing the similarity between *all* query-key pairs, which leads to the computation complexity of  $\mathcal{O}(N^2)$ .

Due to the quadratic computation complexity, simply using self-attention with global receptive field becomes intractable, which usually leads to excessive computation costs. Previous works either addressed this concern by designing sparse global attention pattern [41, 46] or applying smaller attention windows [24, 10]. Though effective, these approaches become susceptible to the carefully-designed attention patterns, or inevitably sacrifice the ability to model long-range dependencies.

#### 3.2. Linear Attention

Comparably, linear attention [19] is considered as an effective alternative which restricts the computation complexity from  $\mathcal{O}(N^2)$  to  $\mathcal{O}(N)$ . Specifically, carefully designed kernels are introduced as the approximation of the original similarity function, *i.e.*,

$$\text{Sim}(Q, K) = \phi(Q)\phi(K)^T, \quad (2)$$

where the self-attention module can be rewritten as:

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

In this way, we can change the computation order from  $(QK^T)V$  to  $Q(K^TV)$  based on the associative property of matrix multiplication (as illustrated in Fig. 1):

$$O_i = \frac{\phi(Q_i) \left( \sum_{j=1}^N \phi(K_j)^T V_j \right)}{\phi(Q_i) \left( \sum_{j=1}^N \phi(K_j)^T \right)}, \quad (4)$$

where the computation complexity with respect to token number is reduced to  $\mathcal{O}(N)$ .

However, current linear attention approaches also face the dilemma between model complexity and expressiveness. On one hand, simple approximations, *e.g.*, using ReLU activation [2], are too loose and lead to significant performance drop. On the other hand, carefully designed kernel functions [7] or matrix decomposition approaches [26, 50] may incur additional computation overhead. In general, there is still a gap between the practical performance of linear attention and Softmax attention.

### 4. Focused Linear Attention

Although enjoying linear computational complexity, various previous works have also proved that simply replacing Softmax attention with linear attention usually results in severe performance drop [34, 2, 7, 27]. In this section, we first perform a detailed analysis of the inferior performances of linear attention from two perspectives: focus ability and feature diversity. Then, we introduce our **Focused Linear Attention** which adequately addresses these concerns and achieves high efficiency and expressive capability.

#### 4.1. Focus ability

Softmax attention practically provides a nonlinear re-weighting mechanism, which makes it easy to concentrate on important features [34, 2, 58]. As shown in Fig. 3, the distribution of attention map from Softmax attention is especially sharp on certain regions, *e.g.*, foreground objects. Comparably, the distribution in linear attention is relatively

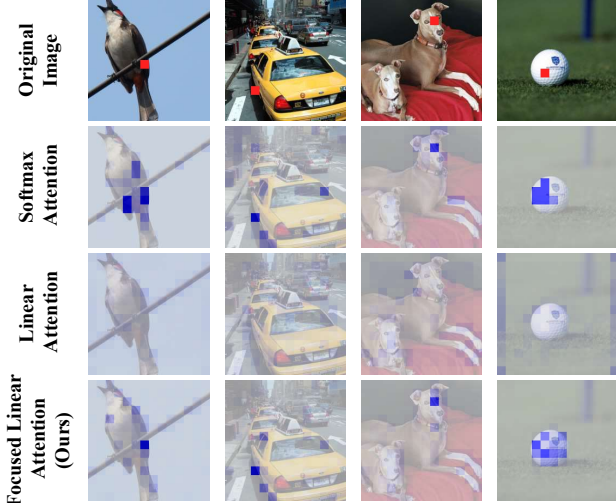


Figure 3. The distribution of Softmax attention, linear attention and our focused linear attention from DeiT-tiny. Softmax attention can produce sharp distribution, while linear attention’s distribution is relatively smooth. Our module restores the sharp distribution as the original Softmax attention. Feature corresponding to the red block is used as query. See more visualizations in Appendix.

smooth, making its output closer to the average of all features and failing to focus on more informative regions.

As a remedy, we propose a simple yet effective solution by adjusting the direction of each query and key features, driving similar query-key pairs closer while pushing dissimilar query-key pairs away. Specifically, we present a simple mapping function  $f_p$  called **Focused Function**:

$$\text{Sim}(Q_i, K_j) = \phi_p(Q_i) \phi_p(K_j)^T, \quad (5)$$

$$\text{where } \phi_p(x) = f_p(\text{ReLU}(x)), \quad f_p(x) = \frac{\|x\|}{\|x^{**p}\|} x^{**p}, \quad (6)$$

and  $x^{**p}$  represents element-wise power  $p$  of  $x$ . We follow previous linear attention modules to use the ReLU function first to ensure the non-negativity of input and validity of denominator in Eq.(4). A direct observation is that the norm of the feature is preserved after the mapping, *i.e.*,  $\|x\| = \|f_p(x)\|$ , indicating that only feature direction is adjusted.

On this basis, we show that under mild assumptions, the proposed mapping function  $f_p$  practically affects the distribution of attention.

**Proposition 1** (Feature direction adjustment with  $f_p$ ) *Let  $x = (x_1, \dots, x_n)$ ,  $y = (y_1, \dots, y_n) \in \mathbb{R}^n$ ,  $x_i, y_j \geq 0$ . Assume  $x$  and  $y$  have the **single** largest value  $x_m$  and  $y_n$  respectively. For a pair of feature  $\{x, y\}$  with  $m = n$ :*

$$\exists p > 1, \text{ s.t. } \langle \phi_p(x), \phi_p(y) \rangle > \langle x, y \rangle. \quad (7)$$

*For a pair of feature  $\{x, y\}$  with  $m \neq n$ :*

$$\exists p > 1, \text{ s.t. } \langle \phi_p(x), \phi_p(y) \rangle < \langle x, y \rangle. \quad (8)$$

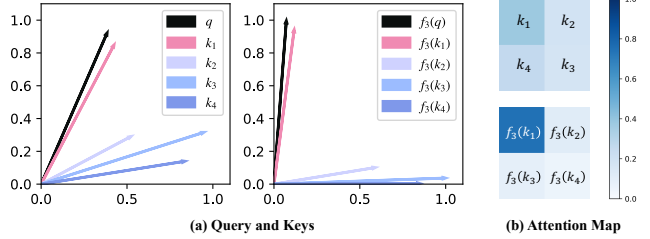


Figure 4. (a)  $f_p$  “pulls” each vector to its nearest axis, thus helping linear attention focus on similar features. (b) The vanilla linear attention scores are  $[0.37, 0.19, 0.26, 0.18]$ , while the attention scores after  $f_3$  are  $[0.75, 0.11, 0.09, 0.05]$ .

*Proof.* Please refer to Appendix for complete proof.  $\square$

Therefore, with a proper  $p$ , our focused function  $f_p(\cdot)$  practically achieves a more distinguished difference between similar query-key pairs (Eq. (7)) and dissimilar query-key pairs (Eq. (8)), restoring the sharp attention distribution as the original Softmax function.

For better understanding, we give an example to show the effects of  $f_p$  in Fig. 4. It can be seen that  $f_p$  actually “pulls” each vector to its nearest axis, and  $p$  determines the degree of this “pulling”. By doing so,  $f_p$  helps divide the features into several groups according to their nearest axes, improving the similarity within each group while reducing the similarity between the groups. The visualizations are in accordance with our analysis above.

## 4.2. Feature diversity

Apart from focus ability, feature diversity is also one of the factors that set restriction on the expressive power of linear attention. One of the possible reasons may give credit to the rank of the attention matrix [36, 53], where a significant difference can be seen. Take one of the Transformer layers from DeiT-Tiny [39] with  $N = 14 \times 14$  for example, we can see from Fig. 5 (a) that the attention matrix has the full rank (196 out of 196), showing the diversity when aggregating features from values.

Nevertheless, this can be hardly achieved in the case of linear attention. As a matter of fact, the rank of the attention matrix in linear attention is bounded by the number of tokens  $N$  and the channel dimension  $d$  for each head:

$$\begin{aligned} \text{rank}(\phi(Q)\phi(K)^T) &\leq \min\{\text{rank}(\phi(Q)), \text{rank}(\phi(K))\} \\ &\leq \min\{N, d\}, \end{aligned} \quad (9)$$

where  $d$  is usually smaller than  $N$  in common vision Transformer designs, *e.g.*,  $d = 64, N = 196$  in DeiT [39] and  $d = 32, N = 49$  in Swin Transformer [24]. In this case, the upper bound of attention matrix rank is restricted at a lower ratio, which indicates that many rows of the attention map are seriously homogenized. As the output of self-attention

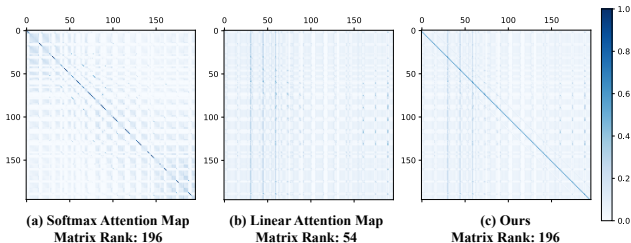


Figure 5. Attention map ( $196 \times 196$ ) from the 3rd block of DeiT-Tiny. (a) Softmax attention can learn a full-rank attention map. (b) Linear attention can not learn an attention map with a rank greater than head dim 64. Many rows of the attention map are seriously homogenized, resulting in the resemblance among output features. (c) The lightweight DWC helps linear attention learn an equivalent attention map with a high rank and maintain feature diversity. Both (b) and (c) involve focused function  $f_p$ .

is the weighted sum of the same set of  $V$ , the homogenization of attention weights inevitably leads to the resemblance among the aggregated features.

To better illustrate, we substitute the original Softmax attention in DeiT-Tiny with linear attention, and show the rank of the attention map in Fig. 5 (b). It can be observed that the rank is greatly decreased (54 out of 196) and many rows of the attention matrix are similar.

As a remedy, we present a simple yet effective solution to address this limitation of linear attention. Specifically, a depthwise convolution (DWC) module is added to the attention matrix and the output can be formulated as:

$$O = \phi(Q)\phi(K)^T V + \text{DWC}(V). \quad (10)$$

To better understand the effect of this DWC module, we can consider it as a kind of attention, in which each query will only focus on several adjacent features in space instead of all features  $V$ . This locality ensures that even if the linear attention values corresponding to two queries are the same, we can still get different outputs from different local features, thus maintaining feature diversity. The effect of DWC can also be explained from the perspective of matrix rank. Based on Eq.(10), we have:

$$O = (\phi(Q)\phi(K)^T + M_{\text{DWC}}) V = M_{\text{eq}} V, \quad (11)$$

where we denote  $M_{\text{DWC}}$  as the sparse matrix corresponding to the depthwise convolution function, and denote  $M_{\text{eq}}$  as the equivalent full attention map. As  $M_{\text{DWC}}$  has the potential to be a full rank matrix, we practically increase the upper bound of the rank of the equivalent attention matrix, which incurs little computation overhead while greatly improving the linear attention’s performance.

To better illustrate, we conduct similar modifications on DeiT-Tiny. With the additional DWC module, the rank of the attention map in the linear attention can be restored to full rank (196 out of 196 as shown in Fig. 5 (c)), which keeps the feature diversity as the original Softmax attention.

### 4.3. Focused linear attention module

Based on the aforementioned analysis, we propose a novel linear attention module, dubbed *focused linear attention*, which reduces the computation complexity while maintaining the expressive power. Specifically, we first design a novel mapping function to imitate the sharp distribution of the original Softmax attention. On this basis, we focus on the low-rank dilemma in previous linear attention modules, and adopt a simple depthwise convolution to restore feature diversity. In this way, our new module can enjoy benefits from both linear complexity and high expressiveness. Specifically, our module can be formulated as:

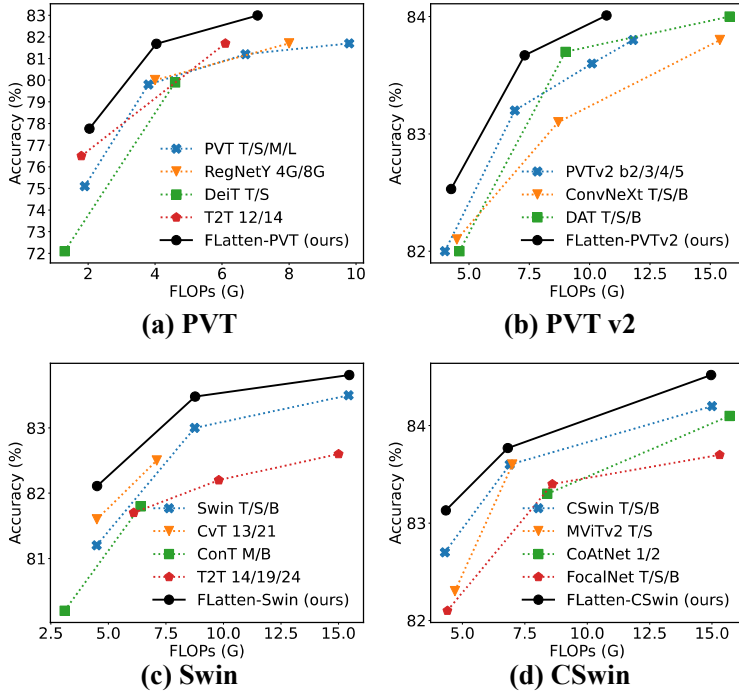
$$O = \text{Sim}(Q, K)V = \phi_p(Q)\phi_p(K)^T V + \text{DWC}(V). \quad (12)$$

In general, our module has the following advantages:

(1) **Low computation complexity as linear attention.** By changing the computation order of self-attention, the complexity is transformed from  $\mathcal{O}(N^2 d)$  to  $\mathcal{O}(Nd^2)$ , where  $N$  and  $d$  denote the token number and channel dimension of each head respectively.  $d$  is usually smaller than  $N$  in common vision Transformer designs, e.g.,  $d = 64, N = 196$  in DeiT [39] and  $d = 32, N = 49$  in Swin Transformer [24], the overall computation is practically decreased. Also, compared to previous linear attention modules [7] that design complex kernel function, our proposed focused function  $f_p$  only adopts simple operators which achieves approximation with minimum computation overhead.

(2) **High expressive capability as Softmax attention.** As we have analyzed above, previous kernel-based linear attention designs are generally inferior to the Softmax counterpart from the focus ability and feature diversity perspective. With the proposed focused function  $f_p$  and depthwise convolution, our *focused linear attention* can achieve even better performance than Softmax attention.

In addition, our module also has the potential of adapting to larger receptive field and different model architectures. Modern Transformer models based on Softmax attention mainly use a limited number of key/value pairs because of the quadratic complexity towards token numbers. Nevertheless, the linear complexity of our module endows us to expand the receptive field to a larger region while maintaining the same amount of computation, and enjoying the advantage of modeling long-range dependencies. Also, our module can serve as a plug-in module and be easily adopted on a variety of modern vision Transformer architectures. We empirically implement our module on five advanced models including DeiT [39], PVT [41], PVT-v2 [42], Swin Transformer [24] and CSwin Transformer [10]. Considering the advantage of enlarged receptive field, we adopt the focused linear attention block at early stages of the vision Transformers, and keep the rest of blocks unchanged. Detailed model architectures are shown in Appendix.



Method	Reso	#Params	Flops	Top-1
DeiT-T [39]	224 <sup>2</sup>	5.7M	1.2G	72.2
<b>FLatten-DeiT-T</b>	224 <sup>2</sup>	6.1M	1.1G	<b>74.1 (+1.9)</b>
PVT-T [41]	224 <sup>2</sup>	13.2M	1.9G	75.1
<b>FLatten-PVT-T</b>	224 <sup>2</sup>	12.2M	2.0G	<b>77.8 (+2.7)</b>
PVT-S	224 <sup>2</sup>	24.5M	3.8G	79.8
<b>FLatten-PVT-S</b>	224 <sup>2</sup>	21.7M	4.0G	<b>81.7 (+1.9)</b>
PVTv2-B1 [42]	224 <sup>2</sup>	13.1M	2.1G	78.7
<b>FLatten-PVTv2-B1</b>	224 <sup>2</sup>	12.9M	2.2G	<b>79.5 (+0.7)</b>
PVTv2-B2	224 <sup>2</sup>	25.4M	4.0G	82.0
<b>FLatten-PVTv2-B2</b>	224 <sup>2</sup>	22.6M	4.3G	<b>82.5 (+0.5)</b>
Swin-T [24]	224 <sup>2</sup>	29M	4.5G	81.3
<b>FLatten-Swin-T</b>	224 <sup>2</sup>	29M	4.5G	<b>82.1 (+0.8)</b>
Swin-S	224 <sup>2</sup>	50M	8.7G	83.0
<b>FLatten-Swin-S</b>	224 <sup>2</sup>	51M	8.7G	<b>83.5 (+0.5)</b>
Swin-B	224 <sup>2</sup>	88M	15.4G	83.5
<b>FLatten-Swin-B</b>	224 <sup>2</sup>	89M	15.4G	<b>83.8 (+0.3)</b>
Swin-B	384 <sup>2</sup>	88M	47.0G	84.5
<b>FLatten-Swin-B</b>	384 <sup>2</sup>	91M	46.5G	<b>85.0 (+0.5)</b>
CSwin-T [10]	224 <sup>2</sup>	23M	4.3G	82.7
<b>FLatten-CSwin-T</b>	224 <sup>2</sup>	21M	4.3G	<b>83.1 (+0.4)</b>
CSwin-S	224 <sup>2</sup>	35M	6.9G	83.6
<b>FLatten-CSwin-S</b>	224 <sup>2</sup>	35M	6.9G	<b>83.8 (+0.2)</b>
CSwin-B	224 <sup>2</sup>	78M	15.0G	84.2
<b>FLatten-CSwin-B</b>	224 <sup>2</sup>	75M	15.0G	<b>84.5 (+0.3)</b>
CSwin-B	384 <sup>2</sup>	78M	47.0G	85.4
<b>FLatten-CSwin-B</b>	384 <sup>2</sup>	78M	46.4G	<b>85.5 (+0.1)</b>

Figure 6. Comparison of different models on ImageNet-1K. See the full comparison table in Appendix.

## 5. Experiments

To verify the effectiveness of our method, we conduct experiments on ImageNet-1K classification [9], ADE20K semantic segmentation [60], and COCO object detection [23]. We also provide a detailed comparison with other linear attention modules based on two representative model structures. In addition, we perform comprehensive ablation studies to analyze each important design element.

### 5.1. ImageNet-1K Classification

ImageNet-1K [9] contains 1.28M images for training and 50K images for validation. We practically implement our module on five advanced Vision Transformer models, and report the Top-1 accuracy on the validation split to compare with various state-of-the-art models.

For fair comparison, we use the exact same settings as the corresponding baseline model to train our FLatten model. Specifically, we use AdamW [25] optimizer to train all our models for 300 epochs with a cosine learning rate decay and 20 epochs of linear warm-up. The basic learning rate for a batch size of 1024 is set to  $1 \times 10^{-3}$ , and then linearly scaled *w.r.t.* the batch size. We follow DeiT [39] and apply RandAugment [8], Mixup [57], CutMix [56] and random erasing [59] to avoid overfitting. In addition, a weight decay of 0.05 is used. To be consistent with [10], we also adopt EMA [32] in the training of our FLatten-CSwin models. In terms of larger resolution finetuning, we follow the

### Semantic Segmentation on ADE20K

Backbone	Method	FLOPs	#Params	mIoU	mAcc
PVT-T	S-FPN	158G	17M	36.57	46.72
<b>FLatten-PVT-T</b>	S-FPN	169G	16M	<b>37.21</b>	48.95
Swin-T	UperNet	945G	60M	44.51	55.61
<b>FLatten-Swin-T</b>	UperNet	946G	60M	<b>44.82</b>	57.01
Swin-S	UperNet	1038G	81M	47.64	58.78
<b>FLatten-Swin-S</b>	UperNet	1038G	82M	<b>48.14</b>	59.31

Table 1. Results of semantic segmentation. The FLOPs are computed over encoders and decoders with an input image at the resolution of  $512 \times 2048$ . S-FPN is short for SemanticFPN [20] model.

setting in [24, 10] that finetunes the models for 30 epochs.

The classification results are provided in Fig. 6. It is shown that our method achieves consistent improvements against baseline models under comparable FLOPs or parameters. For example, our FLatten-PVT-T/S surpass PVT-T/S by 2.7% and 1.9% respectively with similar FLOPs. Based on Swin, our model achieves comparable performance with 60% FLOPs. Our model based on PVT-v2 and CSwin also achieves a better trade-off between computation cost and model performance. These results demonstrate that our module has high expressive capability and is applicable to various model structures.

(a) Mask R-CNN Object Detection & Instance Segmentation on COCO															
Method	FLOPs	#Param	Schedule	AP <sup>b</sup>	AP <sub>50</sub> <sup>b</sup>	AP <sub>75</sub> <sup>b</sup>	AP <sub>s</sub> <sup>b</sup>	AP <sub>m</sub> <sup>b</sup>	AP <sub>l</sub> <sup>b</sup>	AP <sup>m</sup>	AP <sub>50</sub> <sup>m</sup>	AP <sub>75</sub> <sup>m</sup>	AP <sub>s</sub> <sup>m</sup>	AP <sub>m</sub> <sup>m</sup>	AP <sub>l</sub> <sup>m</sup>
PVT-T	240G	33M	1x	36.7	59.2	39.3	21.6	39.2	49.0	35.1	56.7	37.3	19.5	37.4	48.5
<b>FLatten-PVT-T</b>	244G	32M	1x	38.2	61.6	41.9	24.1	40.7	51.0	37.0	57.6	39.0	19.4	39.0	52.1
Swin-T	267G	48M	1x	43.7	66.6	47.7	28.5	47.0	57.3	39.8	63.3	42.7	24.2	43.1	54.6
<b>FLatten-Swin-T</b>	268G	49M	1x	44.2	67.3	48.5	29.4	47.5	57.0	40.2	63.8	43.0	24.5	43.8	54.7
Swin-T	267G	48M	3x	46.0	68.1	50.3	31.2	49.2	60.1	41.6	65.1	44.9	25.9	45.1	56.9
<b>FLatten-Swin-T</b>	268G	49M	3x	46.5	68.5	50.8	31.2	49.6	60.4	42.1	65.4	45.1	25.4	45.4	56.8

(b) Cascade Mask R-CNN Object Detection & Instance Segmentation on COCO															
Method	FLOPs	#Param	Schedule	AP <sup>b</sup>	AP <sub>50</sub> <sup>b</sup>	AP <sub>75</sub> <sup>b</sup>	AP <sub>s</sub> <sup>b</sup>	AP <sub>m</sub> <sup>b</sup>	AP <sub>l</sub> <sup>b</sup>	AP <sup>m</sup>	AP <sub>50</sub> <sup>m</sup>	AP <sub>75</sub> <sup>m</sup>	AP <sub>s</sub> <sup>m</sup>	AP <sub>m</sub> <sup>m</sup>	AP <sub>l</sub> <sup>m</sup>
Swin-T	745G	86M	3x	50.4	69.2	54.7	33.8	54.1	65.2	43.7	66.6	47.3	27.3	47.5	59.0
<b>FLatten-Swin-T</b>	747G	87M	3x	50.8	69.6	55.1	34.2	54.6	65.5	44.1	67.0	48.1	27.6	48.1	59.0
Swin-S	838G	107M	3x	51.9	70.7	56.3	35.2	55.7	67.7	45.0	68.2	48.8	28.8	48.7	60.6
<b>FLatten-Swin-S</b>	841G	108M	3x	52.2	71.2	56.8	35.6	56.4	67.6	45.4	68.3	49.4	29.3	49.0	60.8

Table 2. Results on COCO dataset. The FLOPs are computed over backbone, FPN and detection head with input resolution of 1280×800.

(a) Comparison on DeiT-T Setting			
Linear Attention	FLOPs	#Param	Acc.
Hydra Attn [1]	1.1G	5.7M	68.3
Efficient Attn [38]	1.1G	5.7M	70.2
Linear Angular Attn [52]	1.1G	5.7M	70.8
Enhanced Linear Attn [2]	1.1G	5.8M	72.9
<b>Ours</b>	1.1G	6.1M	<b>74.1</b>

(b) Comparison on Swin-T Setting			
Linear Attention	FLOPs	#Param	Acc.
Hydra Attn [1]	4.5G	29M	80.7
Efficient Attn [38]	4.5G	29M	81.0
Linear Angular Attn [52]	4.5G	29M	79.4
Enhanced Linear Attn [2]	4.5G	29M	81.8
<b>Ours</b>	4.5G	29M	<b>82.1</b>

Table 3. Comparison of different linear attention designs on DeiT-Tiny and Swin-Tiny structures.

## 5.2. Semantic Segmentation

ADE20K [60] is a widely adopted benchmark for semantic segmentation with 20K/2K training/validation images. We employ our model on two representative segmentation models, SemanticFPN [20] and UperNet [47]. As shown in Tab. 1, our model achieves consistently better results under all settings. Specifically, we can see a 0.5 ~ 1% mIoU improvement with comparable computation cost and parameters. The improvements in mAcc are more significant.

## 5.3. Object Detection

COCO [23] object detection and instance segmentation dataset has 118K training and 5K validation images. We

use ImageNet pretrained model as the backbone in Mask R-CNN [18] and Cascade Mask R-CNN [3] frameworks to evaluate the effectiveness. We conduct experiments on 1x and 3x schedules with different detection heads and show results in Tab. 2. Taking advantage of larger receptive field, our model shows better results under all settings.

## 5.4. Comparison with Other Linear Attention

To show a fair comparison with other linear attention modules, we conduct experiments based on two representative Vision Transformer structures, DeiT and Swin Transformer respectively. Based on these two models, we compare our focused linear attention module with four previous linear attention designs, including hydra attention [1], efficient attention [38], linear angular attention [52] and enhanced linear attention [2].

As shown in Tab. 3, previous linear attention modules are generally inferior to the Softmax counterpart, while our model significantly outperforms all other designs and the Softmax baseline. This indicates that our module has high expressive capability and can achieve better performance than Softmax attention with lower computation complexity.

## 5.5. Inference Time

We further evaluate the practical efficiency of our model and compare it with two competitive baselines. The results are presented in Fig. 7. We test the inference latency on multiple hardware platforms, including a desktop CPU (Intel i5-8265U) and two server GPUs (RTX2080Ti and RTX3090). It can be observed that our model achieves a significantly better trade-off between runtime and accuracy on both CPU and GPU, enjoying up to 2.1x faster inference speed with on par or even better performances.

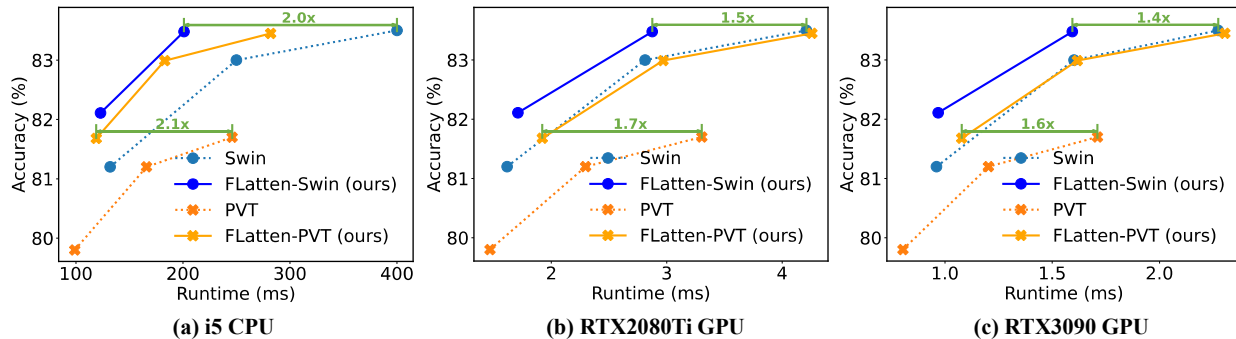


Figure 7. Accuracy-Runtime curve on ImageNet. Runtime is tested with image resolution 224×224.

	FLOPs	#Param	Acc.	Diff.
Vanilla Linear Attention	1.1G	5.7M	70.5	-3.6
+ Focused Function	1.1G	5.7M	71.8	-2.3
+ DWC	1.1G	6.1M	<b>74.1</b>	<b>Ours</b>
DeiT-T	1.2G	5.7M	72.2	-1.9

Table 4. Ablation on each module based on DeiT-T.

focused factor $p$	2	3	4	8	32
Acc.	82.03	82.11	81.94	81.99	81.88

Table 5. Ablation on focused factor  $p$  based on FLatten-Swin-T.

## 5.6. Ablation Study

In this section, we ablate the key components in our focused linear attention to verify the effectiveness of these designs. We report the results on ImageNet-1K classification based on FLatten-DeiT-T and FLatten-Swin-T.

**Focused function  $f_p$  and DWC.** We first evaluate the effectiveness of our proposed focused function  $f_p$  and depth-wise convolution. We start from the vanilla linear attention and introduce  $f_p$  and DWC in turn. As shown in Tab. 4, adopting the proposed focused function  $f_p$  provides +1.3 improvement. Using DWC to maintain feature diversity further leads to an accuracy gain of +2.3, achieving an overall accuracy of 74.1. These results prove that our proposed  $f_p$  and DWC can greatly improve the expressive capability of linear attention, thus helping our focused linear attention module achieve better performance than Softmax attention.

**Ablation on different  $p$ .** In Tab. 5, we study the effect of focused factor  $p$  on the model performance. We find that when  $p$  changes between 2 and 32, the Top-1 classification accuracy does not change much, implying the robustness of our module to this hyper-parameter. Practically, we choose  $p = 3$  for all models in the paper without additional tuning.

**Receptive field.** We also study the impact of receptive field based on FLatten-Swin-tiny. As illustrated in Tab. 6, with the increase of window size, the performance of our model improves progressively. This further proves that though window attention is effective, it inevitably sacrifices the long-range dependency of self-attention from the global perspective and is still inferior to global attention. With lin-

	Window	FLOPs	#Param	Acc.	Diff.
FLatten-Swin-T	$7^2$	4.5G	29M	81.6	-0.5
	$14^2$	4.5G	29M	81.8	-0.3
	$28^2$	4.5G	29M	81.9	-0.2
	$56^2$	4.5G	29M	<b>82.1</b>	<b>Ours</b>
Swin-T	$7^2$	4.5G	29M	81.3	-0.8

Table 6. Ablation on window size based on FLatten-Swin-T.

Stages w/ FLatten				FLOPs	#Param	Acc.	Diff.
Stage1	Stage2	Stage3	Stage4				
✓				4.5G	29M	81.7	-0.4
✓	✓			4.5G	29M	<b>82.1</b>	<b>Ours</b>
✓	✓	✓		4.5G	30M	82.0	-0.1
✓	✓	✓	✓	4.5G	30M	81.9	-0.2
Swin-T				4.5G	29M	81.3	-0.8

Table 7. Ablation on applying focused linear attention on different stages of the Swin-T structure.

ear complexity, it is possible for our module to realize a large even global receptive field while maintaining the same amount of computation.

**Focused linear attention at different stages.** We replace the shift-window attention of Swin-T with our module at different stages. As shown in Tab. 7, we can see that replacing the first two stages leads to a performance gain of 0.8, while replacing the last two stages slightly decreases the overall accuracy. We attribute this result to the fact that the first two stages of Swin have larger resolutions and are more suitable for our module with large receptive field.

## 6. Conclusion

In this paper, we propose a novel *focused linear attention* module. By addressing the limitations of previous linear attention methods from focus ability and feature diversity perspectives, our module achieves an impressive combination of high efficiency and expressive capability. Extensive experiments on image classification, object detection and semantic segmentation demonstrated that our module can be widely applied to a variety of vision Transformers and achieve a better trade-off between computation efficiency and model performance.



## Acknowledgement

This work is supported in part by National Key R&D Program of China (2021ZD0140407), the National Natural Science Foundation of China (62022048, 62276150) and THU-Bosch JCML. We appreciate the generous donation of computing resources by High-Flyer AI.

## References

- [1] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, and Judy Hoffman. Hydra attention: Efficient attention with many heads. In *Computer Vision–ECCV 2022 Workshops: Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VII*, pages 35–49. Springer, 2023. 2, 7
- [2] Han Cai, Chuang Gan, and Song Han. Efficientvit: Enhanced linear attention for high-resolution low-computation visual recognition. *arXiv preprint arXiv:2205.14756*, 2022. 2, 3, 7
- [3] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Conference on Computer Vision and Pattern Recognition*, 2018. 7
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020. 1
- [5] Yinpeng Chen, Xiyang Dai, Dongdong Chen, Mengchen Liu, Xiaoyi Dong, Lu Yuan, and Zicheng Liu. Mobileformer: Bridging mobilenet and transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5270–5279, 2022. 2
- [6] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems*, 34:17864–17875, 2021. 1
- [7] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. In *International Conference on Learning Representations*, 2021. 2, 3, 5
- [8] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition workshops*, pages 702–703, 2020. 6
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on Computer Vision and Pattern Recognition*, pages 248–255. Ieee, 2009. 6
- [10] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12124–12134, 2022. 3, 5, 6, 13, 14
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 1, 2
- [12] Jianyuan Guo, Kai Han, Han Wu, Yehui Tang, Xinghao Chen, Yunhe Wang, and Chang Xu. Cmt: Convolutional neural networks meet vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12175–12185, 2022. 2
- [13] Yizeng Han, Dongchen Han, Zeyu Liu, Yulin Wang, Xuran Pan, Yifan Pu, Chao Deng, Junlan Feng, Shiji Song, and Gao Huang. Dynamic perceiver for efficient visual recognition. In *International Conference on Computer Vision*, 2023. 2
- [14] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 2
- [15] Yizeng Han, Yifan Pu, Zihang Lai, Chaofei Wang, Shiji Song, Junfeng Cao, Wenhui Huang, Chao Deng, and Gao Huang. Learning to weight samples for dynamic early-exiting networks. In *European Conference on Computer Vision*, 2022. 2
- [16] Yizeng Han, Zhihang Yuan, Yifan Pu, Chenhao Xue, Shiji Song, Guangyu Sun, and Gao Huang. Latency-aware spatial-wise dynamic networks. *Advances in Neural Information Processing Systems*, 35:36845–36857, 2022. 1
- [17] Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. *arXiv preprint arXiv:2204.07143*, 2022. 1, 2
- [18] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *International Conference on Computer Vision*, 2017. 7
- [19] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pages 5156–5165. PMLR, 2020. 2, 3
- [20] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *Conference on Computer Vision and Pattern Recognition*, 2019. 6, 7
- [21] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020. 1
- [22] Yanghao Li, Hanzi Mao, Ross Girshick, and Kaiming He. Exploring plain vision transformer backbones for object detection. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, pages 280–296. Springer, 2022. 1, 2
- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, 2014. 6, 7
- [24] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 1, 2, 3, 4, 5, 6, 13

- [25] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. 6
- [26] Jiachen Lu, Jinghan Yao, Junge Zhang, Xiatian Zhu, Hang Xu, Weiguo Gao, Chunjing Xu, Tao Xiang, and Li Zhang. Soft: Softmax-free transformer with linear complexity. *Advances in Neural Information Processing Systems*, 34:21297–21309, 2021. 2, 3
- [27] Xuezhe Ma, Xiang Kong, Sinong Wang, Chunting Zhou, Jonathan May, Hao Ma, and Luke Zettlemoyer. Luna: Linear unified nested attention. *Advances in Neural Information Processing Systems*, 34:2441–2453, 2021. 3
- [28] Sachin Mehta and Mohammad Rastegari. Mobilevit: lightweight, general-purpose, and mobile-friendly vision transformer. In *International Conference on Learning Representations*, 2022. 2
- [29] Zanlin Ni, Yulin Wang, Jiangwei Yu, Haojun Jiang, Yue Cao, and Gao Huang. Deep incubation: Training large models by divide-and-conquering. In *International Conference on Computer Vision*, 2023. 2
- [30] Xuran Pan, Chunjiang Ge, Rui Lu, Shiji Song, Guanfu Chen, Zeyi Huang, and Gao Huang. On the integration of self-attention and convolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 815–825, 2022. 1, 2
- [31] Xuran Pan, Tianzhu Ye, Dongchen Han, Shiji Song, and Gao Huang. Contrastive language-image pre-training with knowledge graphs. *Advances in Neural Information Processing Systems*, 35:22895–22910, 2022. 1
- [32] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM Journal on Control and Optimization*, 30(4):838–855, 1992. 6
- [33] Yifan Pu, Yiru Wang, Zhuofan Xia, Yizeng Han, Yulin Wang, Weihao Gan, Zidong Wang, Shiji Song, and Gao Huang. Adaptive rotated convolution for rotated object detection. In *International Conference on Computer Vision*, 2023. 1
- [34] Zhen Qin, Weixuan Sun, Hui Deng, Dongxu Li, Yunshen Wei, Baohong Lv, Junjie Yan, Lingpeng Kong, and Yiran Zhong. cosformer: Rethinking softmax in attention. In *International Conference on Learning Representations*, 2022. 3
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 1
- [36] Hongyu Ren, Hanjun Dai, Zihang Dai, Mengjiao Yang, Jure Leskovec, Dale Schuurmans, and Bo Dai. Combiner: Full attention transformer with sparse computation cost. *Advances in Neural Information Processing Systems*, 2021. 4
- [37] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018. 1, 2
- [38] Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Efficient attention: Attention with linear complexities. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3531–3539, 2021. 2, 7
- [39] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. 1, 4, 5, 6, 13
- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017. 3
- [41] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021. 1, 2, 3, 5, 6, 13
- [42] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022. 2, 5, 6, 13
- [43] Yulin Wang, Rui Huang, Shiji Song, Zeyi Huang, and Gao Huang. Not all images are worth 16x16 words: Dynamic transformers for efficient image recognition. *Advances in Neural Information Processing Systems*, 34:11960–11973, 2021. 1
- [44] Yulin Wang, Kangchen Lv, Rui Huang, Shiji Song, Le Yang, and Gao Huang. Glance and focus: a dynamic approach to reducing spatial redundancy in image classification. *Advances in Neural Information Processing Systems*, 33:2432–2444, 2020. 1
- [45] Yulin Wang, Yang Yue, Rui Lu, Tianjiao Liu, Zhao Zhong, Shiji Song, and Gao Huang. Efficienttrain: Exploring generalized curriculum learning for training visual backbones. In *International Conference on Computer Vision*, 2023. 2
- [46] Zhuofan Xia, Xuran Pan, Shiji Song, Li Erran Li, and Gao Huang. Vision transformer with deformable attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4794–4803, 2022. 1, 2, 3
- [47] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *European Conference on Computer Vision*, 2018. 7
- [48] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions help transformers see better. *Advances in Neural Information Processing Systems*, 34:30392–30400, 2021. 2
- [49] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. 1
- [50] Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and Vikas Singh.

Nyströmformer: A nyström-based algorithm for approximating self-attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14138–14148, 2021. 2, 3

- [51] Le Yang, Yizeng Han, Xi Chen, Shiji Song, Jifeng Dai, and Gao Huang. Resolution adaptive networks for efficient inference. In *Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [52] Haoran You, Yunyang Xiong, Xiaoliang Dai, Bichen Wu, Peizhao Zhang, Haoqi Fan, Peter Vajda, and Yingyan Lin. Castling-vit: Compressing self-attention via switching towards linear-angular attention during vision transformer inference. *arXiv preprint arXiv:2211.10526*, 2022. 2, 7
- [53] Tong Yu, Ruslan Khalitov, Lei Cheng, and Zhirong Yang. Paramixer: Parameterizing mixing links in sparse factors works better than dot-product self-attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 4
- [54] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 558–567, 2021. 2
- [55] Li Yuan, Qibin Hou, Zihang Jiang, Jiashi Feng, and Shuicheng Yan. Volo: Vision outlooker for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 2
- [56] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6023–6032, 2019. 6
- [57] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 6
- [58] Guangxiang Zhao, Junyang Lin, Zhiyuan Zhang, Xuancheng Ren, Qi Su, and Xu Sun. Explicit sparse transformer: Concentrated attention through explicit selection. *arXiv preprint arXiv:1912.11637*, 2019. 3
- [59] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 13001–13008, 2020. 6
- [60] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. In *International Journal of Computer Vision*. Springer, 2019. 6, 7
- [61] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 1

## Appendix

### A. Proof of Proposition 1

As mentioned in the main paper, with the aim to restore the sharp distribution in linear attention, we present our **Focused Function**  $f_p$ :

$$\text{Sim}(Q_i, K_j) = \phi_p(Q_i) \phi_p(K_j)^T, \quad (13)$$

$$\text{where } \phi_p(x) = f_p(\text{ReLU}(x)), \quad f_p(x) = \frac{\|x\|}{\|x^{**p}\|} x^{**p}, \quad (14)$$

and  $x^{**p}$  represents the power  $p$  of  $x$  bit by bit. We follow previous linear attention modules to use the ReLU function first to ensure the non-negativity of input. Therefore, when proving the effects of  $f_p$ , we only consider  $x, y \geq 0$ .

**Proposition 1** (Feature direction adjustment with  $f_p$ ) *Let  $x = (x_1, \dots, x_n), y = (y_1, \dots, y_n) \in \mathbb{R}^n, x_i, y_j \geq 0$ . Assume  $0 < \langle x, y \rangle < \|x\| \|y\|$  and  $x, y$  have the **single largest value**  $x_m, y_n$  respectively.*

*For a pair of feature  $\{x, y\}$  with  $m = n$ :*

$$\exists p > 1, \text{ s.t. } \langle \phi_p(x), \phi_p(y) \rangle > \langle x, y \rangle. \quad (15)$$

*For a pair of feature  $\{x, y\}$  with  $m \neq n$ :*

$$\exists p > 1, \text{ s.t. } \langle \phi_p(x), \phi_p(y) \rangle < \langle x, y \rangle. \quad (16)$$

*Proof.*

$$\phi_p(x) = f_p(\text{ReLU}(x)) = f_p(x), \quad (17)$$

$$\phi_p(y) = f_p(\text{ReLU}(y)) = f_p(y).$$

$$\|f_p(x)\| = \frac{\|x\|}{\|x^{**p}\|} \|x^{**p}\| = \|x\|, \quad \|f_p(y)\| = \|y\|. \quad (18)$$

Therefore, we have:

$$\begin{aligned} \langle \phi_p(x), \phi_p(y) \rangle &= \langle f_p(x), f_p(y) \rangle \\ &= \|f_p(x)\| \|f_p(y)\| \langle u, v \rangle \\ &= \|x\| \|y\| \langle u, v \rangle, \end{aligned} \quad (19)$$

where

$$\begin{aligned} \langle u, v \rangle &= \left\langle \frac{f_p(x)}{\|f_p(x)\|}, \frac{f_p(y)}{\|f_p(y)\|} \right\rangle \\ &= \frac{\sum_{i=1}^n x_i^p y_i^p}{\sqrt{\left(\sum_{i=1}^n x_i^{2p}\right) \left(\sum_{i=1}^n y_i^{2p}\right)}} \\ &= \frac{\sum_{i=1}^n a_i^p b_i^p}{\sqrt{\left(\sum_{i=1}^n a_i^{2p}\right) \left(\sum_{i=1}^n b_i^{2p}\right)}}, \end{aligned} \quad (20)$$

and

$$\begin{aligned} a_i &= x_i / \max_{1 \leq i \leq n} (x_i), \quad b_i = y_i / \max_{1 \leq i \leq n} (y_i), \\ a_i, b_i &\in [0, 1]. \end{aligned} \quad (21)$$

Based on our assumption, we have:

$$\exists! m, \text{ s.t. } a_m = 1, \quad \exists! n, \text{ s.t. } b_n = 1. \quad (22)$$

Therefore,

$$\lim_{p \rightarrow \infty} a_i^p = \begin{cases} 1, & i = m \\ 0, & i \neq m \end{cases}, \quad \lim_{p \rightarrow \infty} b_j^p = \begin{cases} 1, & j = n \\ 0, & j \neq n \end{cases}. \quad (23)$$

Then we consider the following two cases:

(1)  $m = n$ :

$$\begin{aligned} \lim_{p \rightarrow \infty} \langle u, v \rangle &= \lim_{p \rightarrow \infty} \frac{\sum_{i=1}^n a_i^p b_i^p}{\sqrt{\left(\sum_{i=1}^n a_i^{2p}\right) \left(\sum_{i=1}^n b_i^{2p}\right)}} \\ &= \frac{1 \times 1}{\sqrt{1 \times 1}} = 1. \end{aligned} \quad (24)$$

Eq. (19), Eq. (24)  $\Rightarrow$

$$\begin{aligned} \lim_{p \rightarrow \infty} \langle \phi_p(x), \phi_p(y) \rangle &= \lim_{p \rightarrow \infty} \|x\| \|y\| \langle u, v \rangle \\ &= \|x\| \|y\| > \langle x, y \rangle. \end{aligned} \quad (25)$$

Thus we have,

$$\exists p > 1, \text{ s.t. } \langle \phi_p(x), \phi_p(y) \rangle > \langle x, y \rangle. \quad (26)$$

(2)  $m \neq n$ :

$$\begin{aligned} \lim_{p \rightarrow \infty} \langle u, v \rangle &= \lim_{p \rightarrow \infty} \frac{\sum_{i=1}^n a_i^p b_i^p}{\sqrt{\left(\sum_{i=1}^n a_i^{2p}\right) \left(\sum_{i=1}^n b_i^{2p}\right)}} \\ &= \frac{1 \times 0 + 0 \times 1}{\sqrt{1 \times 1}} = 0. \end{aligned} \quad (27)$$

Eq. (19), Eq. (27)  $\Rightarrow$

$$\begin{aligned} \lim_{p \rightarrow \infty} \langle \phi_p(x), \phi_p(y) \rangle &= \lim_{p \rightarrow \infty} \|x\| \|y\| \langle u, v \rangle \\ &= 0 < \langle x, y \rangle. \end{aligned} \quad (28)$$

Thus we have,

$$\exists p > 1, \text{ s.t. } \langle \phi_p(x), \phi_p(y) \rangle < \langle x, y \rangle. \quad (29)$$

□

Therefore, with a proper  $p$ , our focused function  $f_p(\cdot)$  practically achieves a more distinguished difference between similar query-key pairs (Eq. (15)) and dissimilar query-key pairs (Eq. (16)). Actually,  $f_p$  divides the features into several groups according to their nearest axes, improving the similarity within each group while reducing the similarity between the groups, thus restoring the sharp attention distribution as the original Softmax function.

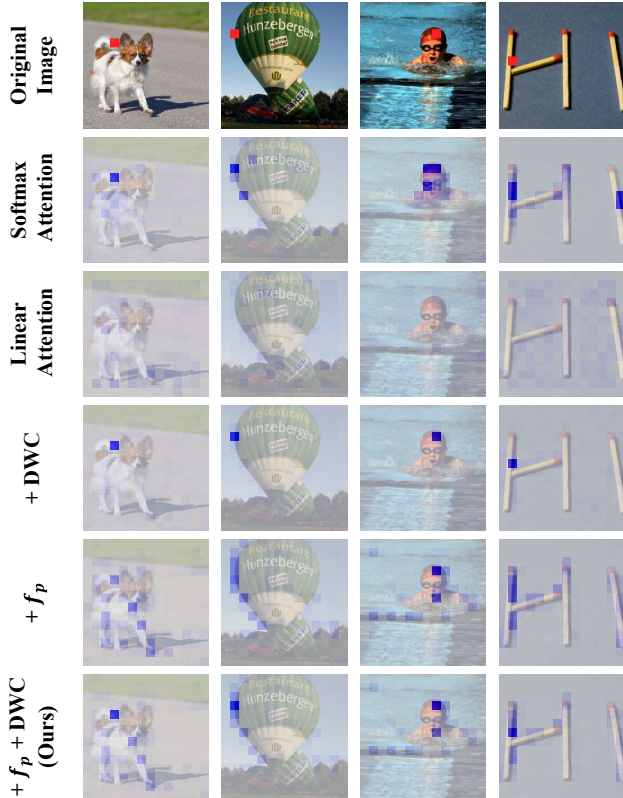


Figure 8. The distribution of attention weights from DeiT-tiny. Feature corresponding to the red block is used as query.

## B. More Visualizations

We visualize more examples of attention weights in Fig. 8. To better show the contribution of our focused function and DWC, we start from the vanilla linear attention and introduce  $f_p$  and DWC separately. As demonstrated in the last three rows, DWC improves local focus ability but cannot focus on any position, while  $f_p$  practically enhances model’s focus ability, helping model focus on more informative regions. Combining  $f_p$  and DWC, our focused linear attention module restores the sharp distribution as the original Softmax attention.

## C. Full Classification Results

Due to the page limit, we only present representative ImageNet classification results in Figure 6 of main paper. Here, we give all the classification results when applying our focused linear attention module on various sizes of the five baseline models in Tab.8.

## D. Model Architectures

We summarize the architectures of five Transformer models adopted in the main paper, including DeiT [39], PVT [41], PVTv2 [42], Swin Transformer [24], CSwin

Method	Reso	#Params	Flops	Top-1
DeiT-T [39]	224 <sup>2</sup>	5.7M	1.2G	72.2
<b>FLatten-DeiT-T</b>	224 <sup>2</sup>	6.1M	1.1G	<b>74.1 (+1.9)</b>
PVT-T [41]	224 <sup>2</sup>	13.2M	1.9G	75.1
<b>FLatten-PVT-T</b>	224 <sup>2</sup>	12.2M	2.0G	<b>77.8 (+2.7)</b>
PVT-S	224 <sup>2</sup>	24.5M	3.8G	79.8
<b>FLatten-PVT-S</b>	224 <sup>2</sup>	21.7M	4.0G	<b>81.7 (+1.9)</b>
PVT-M	224 <sup>2</sup>	44.2M	6.7G	81.2
<b>FLatten-PVT-M</b>	224 <sup>2</sup>	37.2M	7.0G	<b>83.0 (+1.8)</b>
PVT-L	224 <sup>2</sup>	61.4M	9.8G	81.7
<b>FLatten-PVT-L</b>	224 <sup>2</sup>	50.6M	10.4G	<b>83.4 (+1.7)</b>
PVTv2-B0 [42]	224 <sup>2</sup>	3.4M	0.6G	70.5
<b>FLatten-PVTv2-B0</b>	224 <sup>2</sup>	3.6M	0.6G	<b>71.1 (+0.6)</b>
PVTv2-B1	224 <sup>2</sup>	13.1M	2.1G	78.7
<b>FLatten-PVTv2-B1</b>	224 <sup>2</sup>	12.9M	2.2G	<b>79.5 (+0.7)</b>
PVTv2-B2	224 <sup>2</sup>	25.4M	4.0G	82.0
<b>FLatten-PVTv2-B2</b>	224 <sup>2</sup>	22.6M	4.3G	<b>82.5 (+0.5)</b>
PVTv2-B3	224 <sup>2</sup>	45.2M	6.9G	83.2
<b>FLatten-PVTv2-B3</b>	224 <sup>2</sup>	38.3M	7.3G	<b>83.7 (+0.5)</b>
PVTv2-B4	224 <sup>2</sup>	62.6M	10.1G	83.6
<b>FLatten-PVTv2-B4</b>	224 <sup>2</sup>	51.8M	10.7G	<b>84.0 (+0.4)</b>
Swin-T [24]	224 <sup>2</sup>	29M	4.5G	81.3
<b>FLatten-Swin-T</b>	224 <sup>2</sup>	29M	4.5G	<b>82.1 (+0.8)</b>
Swin-S	224 <sup>2</sup>	50M	8.7G	83.0
<b>FLatten-Swin-S</b>	224 <sup>2</sup>	51M	8.7G	<b>83.5 (+0.5)</b>
Swin-B	224 <sup>2</sup>	88M	15.4G	83.5
<b>FLatten-Swin-B</b>	224 <sup>2</sup>	89M	15.4G	<b>83.8 (+0.3)</b>
Swin-B	384 <sup>2</sup>	88M	47.0G	84.5
<b>FLatten-Swin-B</b>	384 <sup>2</sup>	91M	46.5G	<b>85.0 (+0.5)</b>
CSwin-T [10]	224 <sup>2</sup>	23M	4.3G	82.7
<b>FLatten-CSwin-T</b>	224 <sup>2</sup>	21M	4.3G	<b>83.1 (+0.4)</b>
CSwin-S	224 <sup>2</sup>	35M	6.9G	83.6
<b>FLatten-CSwin-S</b>	224 <sup>2</sup>	35M	6.9G	<b>83.8 (+0.2)</b>
CSwin-B	224 <sup>2</sup>	78M	15.0G	84.2
<b>FLatten-CSwin-B</b>	224 <sup>2</sup>	75M	15.0G	<b>84.5 (+0.3)</b>
CSwin-B	384 <sup>2</sup>	78M	47.0G	85.4
<b>FLatten-CSwin-B</b>	384 <sup>2</sup>	78M	46.4G	<b>85.5 (+0.1)</b>

Table 8. Comparisons of focused linear attention with other vision transformer backbones on the ImageNet-1K classification task.

stage	output	FLatten-DeiT-T	
		FLatten	DeiT Block
res1	$14 \times 14$	$\begin{bmatrix} \text{win } 14 \times 14 \\ \text{dim } 192 \\ \text{head } 3 \end{bmatrix} \times 12$	None

Table 9. Architectures of FLatten-DeiT models.

stage	output	FLatten-PVT-M		FLatten-PVT-L	
		FLatten	PVT Block	FLatten	PVT Block
res1	$56 \times 56$	Conv1×1, stride=4, 64, LN			
		$\begin{bmatrix} \text{win } 56 \times 56 \\ \text{dim } 64 \\ \text{head } 1 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 56 \times 56 \\ \text{dim } 64 \\ \text{head } 1 \end{bmatrix} \times 3$	None
res2	$28 \times 28$	Conv1×1, stride=2, 128, LN			
		$\begin{bmatrix} \text{win } 28 \times 28 \\ \text{dim } 128 \\ \text{head } 2 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 28 \times 28 \\ \text{dim } 128 \\ \text{head } 2 \end{bmatrix} \times 3$	None
res3	$14 \times 14$	Conv1×1, stride=2, 320, LN			
		$\begin{bmatrix} \text{win } 14 \times 14 \\ \text{dim } 320 \\ \text{head } 5 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 14 \times 14 \\ \text{dim } 320 \\ \text{head } 5 \end{bmatrix} \times 6$	None
res4	$7 \times 7$	Conv1×1, stride=2, 512, LN			
		$\begin{bmatrix} \text{win } 7 \times 7 \\ \text{dim } 512 \\ \text{head } 8 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 7 \times 7 \\ \text{dim } 512 \\ \text{head } 8 \end{bmatrix} \times 3$	None

Table 10. Architectures of FLatten-PVT models (Part1).

Transformer [10] in Tab.9-15. In practice, we substitute the original self-attention blocks at all stages of the DeiT, PVT and PVTv2 with the focused linear attention block, but only adopt our module at early stages of Swin and CSwin. The model structure (width and depth) are kept unchanged, except for CSwin-T and CSwin-B, where we increase the depth of the first and second stages and correspondingly reduce the depth of the third stage to better reflect our module’s advantage of enlarged receptive field.

stage	output	FLatten-PVT-M		FLatten-PVT-L	
		FLatten	PVT Block	FLatten	PVT Block
res1	$56 \times 56$	Conv1 $\times$ 1, stride=4, 64, LN			
		$\begin{bmatrix} \text{win } 56 \times 56 \\ \text{dim } 64 \\ \text{head } 1 \end{bmatrix} \times 3$	None	$\begin{bmatrix} \text{win } 56 \times 56 \\ \text{dim } 64 \\ \text{head } 1 \end{bmatrix} \times 3$	None
res2	$28 \times 28$	Conv1 $\times$ 1, stride=2, 128, LN			
		$\begin{bmatrix} \text{win } 28 \times 28 \\ \text{dim } 128 \\ \text{head } 2 \end{bmatrix} \times 3$	None	$\begin{bmatrix} \text{win } 28 \times 28 \\ \text{dim } 128 \\ \text{head } 2 \end{bmatrix} \times 8$	None
res3	$14 \times 14$	Conv1 $\times$ 1, stride=2, 320, LN			
		$\begin{bmatrix} \text{win } 14 \times 14 \\ \text{dim } 320 \\ \text{head } 5 \end{bmatrix} \times 18$	None	$\begin{bmatrix} \text{win } 14 \times 14 \\ \text{dim } 320 \\ \text{head } 5 \end{bmatrix} \times 27$	None
res4	$7 \times 7$	Conv1 $\times$ 1, stride=2, 512, LN			
		$\begin{bmatrix} \text{win } 7 \times 7 \\ \text{dim } 512 \\ \text{head } 8 \end{bmatrix} \times 3$	None	$\begin{bmatrix} \text{win } 7 \times 7 \\ \text{dim } 512 \\ \text{head } 8 \end{bmatrix} \times 3$	None

Table 11. Architectures of FLatten-PVT models (Part2).

stage	output	FLatten-PVTv2-B0		FLatten-PVTv2-B1		FLatten-PVTv2-B2	
		FLatten	PVTv2 Block	FLatten	PVTv2 Block	FLatten	PVTv2 Block
res1	$56 \times 56$	Conv4 $\times$ 4, stride=4, 32, LN		Conv4 $\times$ 4, stride=4, 64, LN			
		$\begin{bmatrix} \text{win } 56 \times 56 \\ \text{dim } 32 \\ \text{head } 1 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 56 \times 56 \\ \text{dim } 64 \\ \text{head } 1 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 56 \times 56 \\ \text{dim } 64 \\ \text{head } 1 \end{bmatrix} \times 3$	None
res2	$28 \times 28$	Conv1 $\times$ 1, stride=2, 64, LN		Conv1 $\times$ 1, stride=2, 128, LN			
		$\begin{bmatrix} \text{win } 28 \times 28 \\ \text{dim } 64 \\ \text{head } 2 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 28 \times 28 \\ \text{dim } 128 \\ \text{head } 2 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 28 \times 28 \\ \text{dim } 128 \\ \text{head } 2 \end{bmatrix} \times 3$	None
res3	$14 \times 14$	Conv2 $\times$ 2, stride=2, 160, LN		Conv2 $\times$ 2, stride=2, 320, LN			
		$\begin{bmatrix} \text{win } 14 \times 14 \\ \text{dim } 160 \\ \text{head } 5 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 14 \times 14 \\ \text{dim } 320 \\ \text{head } 5 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 14 \times 14 \\ \text{dim } 320 \\ \text{head } 5 \end{bmatrix} \times 6$	None
res4	$7 \times 7$	Conv2 $\times$ 2, stride=2, 256, LN		Conv2 $\times$ 2, stride=2, 512, LN			
		$\begin{bmatrix} \text{win } 7 \times 7 \\ \text{dim } 256 \\ \text{head } 8 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 7 \times 7 \\ \text{dim } 512 \\ \text{head } 8 \end{bmatrix} \times 2$	None	$\begin{bmatrix} \text{win } 7 \times 7 \\ \text{dim } 512 \\ \text{head } 8 \end{bmatrix} \times 3$	None

Table 12. Architectures of FLatten-PVTv2 models (Part1).

stage	output	FLatten-PVTv2-B3				FLatten-PVTv2-B4			
		<b>FLatten</b>		PVTv2 Block		<b>FLatten</b>		PVTv2 Block	
res1	$56 \times 56$	Conv4×4, stride=4, 64, LN							
		win $56 \times 56$ dim 64 head 1	×3	None		win $56 \times 56$ dim 64 head 1	×3	None	
res2	$28 \times 28$	Conv2×2, stride=2, 128, LN							
		win $28 \times 28$ dim 128 head 2	×3	None		win $28 \times 28$ dim 128 head 2	×8	None	
res3	$14 \times 14$	Conv2×2, stride=2, 320, LN							
		win $14 \times 14$ dim 320 head 5	×18	None		win $14 \times 14$ dim 320 head 5	×27	None	
res4	$7 \times 7$	Conv1×1, stride=2, 512, LN							
		win $7 \times 7$ dim 512 head 8	×3	None		win $7 \times 7$ dim 512 head 8	×3	None	

Table 13. Architectures of FLatten-PVTv2 models (Part2).

stage	output	FLatten-Swin-T				FLatten-Swin-S				FLatten-Swin-B			
		<b>FLatten</b>		Swin Block		<b>FLatten</b>		Swin Block		<b>FLatten</b>		Swin Block	
res1	$56 \times 56$	concat $4 \times 4$ , 96, LN											
		win $56 \times 56$ dim 96 head 3	×2	None		win $56 \times 56$ dim 96 head 3	×2	None		win $56 \times 56$ dim 128 head 3	×2	None	
res2	$28 \times 28$	concat $4 \times 4$ , 192, LN											
		win $28 \times 28$ dim 192 head 6	×2	None		win $28 \times 28$ dim 192 head 6	×2	None		win $28 \times 28$ dim 256 head 6	×2	None	
res3	$14 \times 14$	concat $4 \times 4$ , 384, LN											
		None		win $7 \times 7$ dim 384 head 12	×6	None		win $7 \times 7$ dim 384 head 12	×18	None		win $7 \times 7$ dim 512 head 12	×18
res4	$7 \times 7$	concat $4 \times 4$ , 768, LN											
		None		win $7 \times 7$ dim 768 head 24	×2	None		win $7 \times 7$ dim 768 head 24	×2	None		win $7 \times 7$ dim 1024 head 24	×2

Table 14. Architectures of FLatten-Swin models.



stage	output	FLatten-CSwin-T			FLatten-CSwin-S			FLatten-CSwin-B		
		FLatten		CSwin Block	FLatten		CSwin Block	FLatten		CSwin Block
res1	$56 \times 56$	Conv7×7, stride=4, 64, LN						Conv7×7, stride=4, 96, LN		
		win 3×3 dim 64 head 2	×2	None	win 3×3 dim 64 head 2	×2	None	win 3×3 dim 96 head 4	×3	None
res2	$28 \times 28$	Conv7×7, stride=4, 128, LN						Conv7×7, stride=4, 192, LN		
		win 3×3 dim 128 head 4	×4	None	win 3×3 dim 128 head 4	×4	None	win 3×3 dim 192 head 8	×6	None
res3	$14 \times 14$	Conv7×7, stride=4, 256, LN						Conv7×7, stride=384, LN		
		None		win 3×3 dim 256 head 8	×18	None	win 3×3 dim 256 head 8	×32	None	win 3×3 dim 384 head 16
res4	$7 \times 7$	Conv7×7, stride=4, 512, LN						Conv7×7, stride=4, 768, LN		
		None		win 7×7 dim 512 head 16	×1	None	win 7×7 dim 512 head 16	×2	None	win 7×7 dim 768 head 32

Table 15. Architectures of FLatten-CSwin models.