

CrossFormer++: A Versatile Vision Transformer Hinging on Cross-scale Attention

Wenxiao Wang[†], Wei Chen[†], Qibo Qiu, Long Chen, Boxi Wu, Binbin Lin,
Xiaofei He, *Senior Member, IEEE*, and Wei Liu, *Fellow, IEEE*

Abstract—While features of different scales are perceptually important to visual inputs, existing vision transformers do not yet take advantage of them explicitly. To this end, we first propose a cross-scale vision transformer, CrossFormer. It introduces a cross-scale embedding layer (CEL) and a **long-short distance attention** (LSDA). On the one hand, CEL blends each token with multiple patches of different scales, providing the self-attention module itself with cross-scale features. On the other hand, LSDA splits the self-attention module into a short-distance one and a long-distance counterpart, which not only reduces the computational burden but also keeps both small-scale and large-scale features in the tokens. Moreover, through experiments on CrossFormer, we observe another two issues that affect **vision transformers** performance, *i.e.*, the enlarging self-attention maps and amplitude explosion. Thus, we further propose a progressive group size (PGS) paradigm and an amplitude cooling layer (ACL) to alleviate the two issues, respectively. The CrossFormer incorporating with PGS and ACL is called CrossFormer++. Extensive experiments show that CrossFormer++ outperforms the other vision transformers on image classification, object detection, instance segmentation, and semantic segmentation tasks. The code will be available at: <https://github.com/cheerss/CrossFormer>.

Index Terms—Vision Transformer, Image Classification, Object Detection, Semantic Segmentation

1 INTRODUCTION

TRANSFORMER based vision backbones have achieved great success in many computer vision tasks such as image classification, object detection, semantic segmentation, etc. Compared with convolutional neural networks (CNNs), vision transformers enable long range dependencies by introducing a **self-attention module**, endowing the models with higher capacities.

A transformer requires a sequence of tokens (*e.g.*, word embeddings) as input. To adapt this requirement to typical vision tasks, most existing vision transformers [1], [2], [3], [4] produce tokens by splitting an input image into equally-sized patches. For example, a 224×224 image can be split into 56×56 patches of size 4×4 , and these patches are embedded through a linear layer to yield a token sequence. Inside a certain transformer, self-attention is engaged to build the interactions between any two tokens. Thus, the computational or memory cost of the self-attention module is $O(N^2)$, where N is the length of a token sequence. Such a cost is too big for a visual input because its token sequence is much longer than that of NLP tasks. Therefore, the recently proposed vision transformers [3], [4], [5] develop multiple substitutes (*e.g.*, Swin [4] restricts the self-attention in a small local region instead of a global region) to approximate the vanilla self-attention module with a lower cost.

Though the aforementioned vision transformers have made progress, they still face challenges in utilizing visual

features of different scales, whereas multi-scale features are very vital for a lot of vision tasks. Particularly, an image often contains many objects of different sizes, and to fully understand the image, the model is required to extract features at different scales (*i.e.*, different ranges and sizes). Existing vision transformers fail to deal with the above case due to two reasons: (1) Those models' input tokens are generated from **equally-sized patches**. Though these patches theoretically have a chance to extract any scale features if only the patch size is large enough, it is difficult to promise that they can learn appropriate multi-scale features automatically in practice. (2) Some vision transformers such as Swin [4] restrict the attention in a small local region, giving up the long range attention.

In this paper, we propose a cross-scale embedding layer (CEL) and a long-short distance attention (LSDA) to fill the cross-scale gap from two perspectives:

Cross-scale Embedding Layer (CEL). Following PVT [3], we also employ a pyramid structure for our designed transformer, which naturally splits the vision transformer model into multiple stages. CEL appears at the start of each stage, which receives last stage's output (or an input image) as input and samples patches with multiple kernels of different scales (*e.g.*, 4×4 or 8×8). Then, each token is constructed by embedding and concatenating these patches. Through this way, we enforce some dimensions (*e.g.*, from 4×4 patches) to focus on small-scale features only, while others (*e.g.*, from 8×8 patches) have a chance to learn large-scale features, leading to a token with explicit cross-scale features.

Long-Short Distance Attention (LSDA). Instead of using a shifted window based self-attention like Swin, we argue that both short distance and long distance attentions are imperative for a visual input and thereby propose our LSDA. In particular, we split the self-attention module into

- [†] means equal contributions.
- Wenxiao Wang, Boxi Wu, and Binbin Lin are with Zhejiang University.
- Wei Chen and Xiaofei He are with State Key Lab of CAD&CG, Zhejiang University.
- Qibo Qiu is with Zhejiang Lab.
- Long Chen is with Hong Kong University of Science and Technology.
- Wei Liu is with Data Platform, Tencent.

Manuscript received January 31, 2023.

a Short Distance Attention (SDA) and a Long Distance Attention (LDA). SDA builds the dependencies among neighboring embeddings, while LDA takes charge of the dependencies among embeddings far away from each other. SDA and LDA appear alternately in consecutive layers of a CrossFormer, which also reduce the cost of the self-attention module while keep both the short distance and long distance attentions.

Dynamic Position Bias (DPB). Besides, following prior work [4], [6], we employ a relative position bias for tokens’ position representations. The Relative Position Bias (RPB) only supports fixed image/group size. However, image size for many vision tasks such as object detection is variable, so does group size for many architectures, including ours. To make the RPB more flexible, we further introduce a trainable module called **Dynamic Position Bias (DPB)**, which receives two tokens’ relative distance as input and outputs their position bias. The DPB module is optimized end-to-end in the training phase, inducing an ignorable cost but making RPB apply to variable group size.

Armed with the above three modules, we name our architecture as CrossFormer and design four variants of different depths and channels. While these variants achieve great performance, we also find that a further expansion of CrossFormer does not bring a persistent accuracy gain, but even a degradation. To this end, we analyze the output of each layer (*i.e.*, the self-attention maps and the MLPs) and observe another two issues which affect the models’ performance, *i.e.*, the enlarging self-attention maps and amplitude explosion. Thus, we further propose a progressive group size (PGS) and an amplitude cooling layer (ACL) to alleviate the problems.

Progressive Group Size (PGS). In terms of the self-attention maps, we observe that tokens at shallow layers always attend to tokens around themselves. Whereas tokens at deep layers pay nearly equal attention to all other tokens. This phenomenon shows that vision transformers perform similarly as CNN, *i.e.*, extracting local features at shallow layers and global features at deep layers, respectively. It will hinder the models’ performance if adopting a fixed group size for all stages like most existing vision transformers. Hence, we propose to enlarge the group size progressively (PGS) from shallow to deep layers and implement a manually designed group size policy under this guide. Though it is a simple policy, PGS is actually a general paradigm. We hope to present its importance and appeal other researchers to explore more automated and adaptive group size policies.

Amplitude Cooling Layer (ACL). In the vision transformers, the activation’s amplitude grows dramatically as the layer goes deeper. For example, in a CrossFormer-B, the maximal amplitude in the 22_{nd} layer is 300 times larger than that in the 1_{st} layer. The tremendous amplitude discrepancy among layers makes the training process unstable. However, we also find that our proposed CEL can effectively suppress the amplitude. Considering that it is a bit cumbersome to put more CELs into a CrossFormer, we design a CEL-like but more lightweight layer dubbed Amplitude Cooling Layer (ACL). ACL is inserted after some CrossFormer blocks to cool down the amplitude.

We improve CrossFormer by introducing PGS and ACL, yielding CrossFormer++, and propose four new variants.

Extensive experiments on four downstream tasks (*i.e.*, image classification, object detection, semantic segmentation, and instance segmentation) show that CrossFormer++ outperforms CrossFormer and other existing vision transformers on all these tasks.

The following sections are organized as follows: the background and related works will be introduced in Sec. 2. Then, we will retrospect CrossFormer¹ [7] and its CEL, LSDA, and DPB in Sec. 3. PGS, ACL, and CrossFormer++ are introduced in detail in Sec. 4. Thereafter, the experiments will be shown in Sec. 5. Finally, we will present our conclusions and future work in Sec. 6.

2 BACKGROUND AND RELATED WORKS

Vision Transformers. Motivated by the great success achieved by transformers in NLP, researchers have tried to design specific visual transformers for vision tasks to take advantage of the powerful attention mechanism. Early vision transformers like ViT [1] and DeiT [2] transfer the original transformer to vision tasks, achieving impressive results and demonstrating great potentials. Later works like PVT [3], Swin [4], and MViTv2 [8] combine the pyramid structure with the transformer and remove the class token used in the original architecture, enabling the use in further vision tasks like object detection and image segmentation. As is mentioned in MViTv2 [8], such a pyramid structure is partly inspired by CNN, which hierarchically expands the channel capacity while reducing the spatial resolution. Our work proposes that other insights from CNN may also be valuable for transformers. To be specific, the lower layers in the network tend to refine local features while the higher layers focus more on global information communications.

Efficient Self-Attention and Universal Vision Backbone. Vanilla self-attention module in the transformer suffers from a quadratic complexity with respect to the image size, which is unacceptable for dense vision tasks like object detection and image segmentation. Therefore, in order to build transformer-based universal vision backbones, researchers have proposed ways to do efficient self-attention. The first choice is to do sparse self-attention. Instead of doing self-attention among all embeddings, some transformers divide the embeddings into groups and do self-attention within each group. The ways of choosing the group while also enabling global information interactions become the core design for works like Swin [4], CAT [5], CvT [9], and CSwin [10]. Another choice is to reduce the cost of global self-attention through reducing the size of input queries and keys. PVT [3] uses average pooling to reduce keys size, while Scalable ViT [11] uses convolution with large intervals. MViTv2 [8] inserts a pooling layer after the linear transformation layer of the self-attention module. It should also be noted that the results reported in the above works are sometimes obtained under different settings such as different depths and embedding methods, which disables the direct comparison between different ways to divide the group. Our work proposes a novel way of group division and uses the same experiment setting and architecture to compare the effectiveness of different division ways.

1. An earlier version of this paper has appeared in ICLR 2022: https://openreview.net/pdf?id=_PHymLLxul

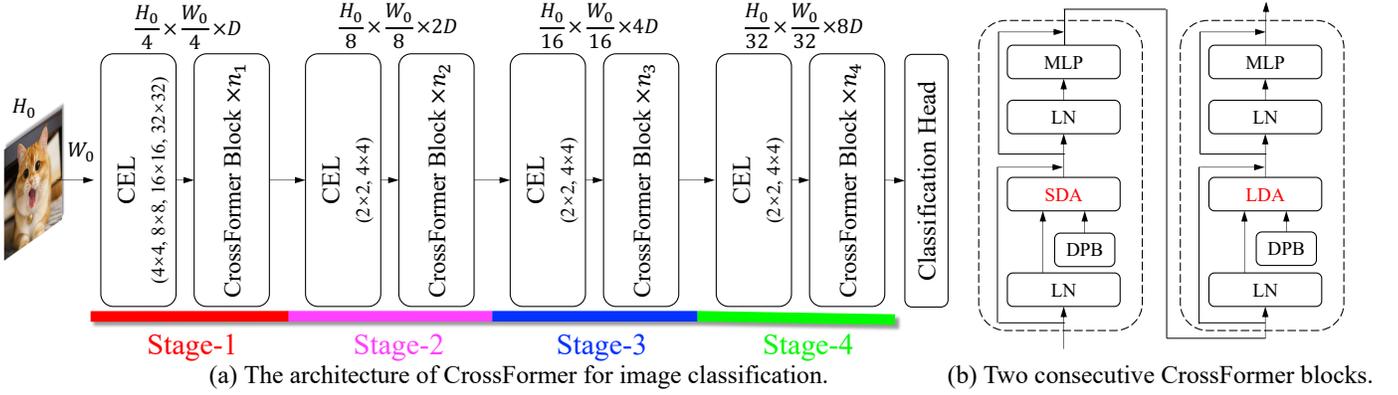


Fig. 1: (a) The architecture of CrossFormer for classification. The input size is $H_0 \times W_0$, and the size of feature maps in each stage is shown on the top. $Stage-i$ consists of a CEL and n_i CrossFormer blocks. Numbers in CELs represent kernels’ sizes used for sampling patches. (b) The inner structure of two consecutive CrossFormer blocks. SDA and LDA appear alternately in different blocks.

Position Representations. Transformers are permutation invariant models, which are unsatisfying for both NLP and vision tasks, where the permutation of the embeddings affects the semantic information. To make the model aware of position information, many different position representations are proposed. For example, Dosovitskiy et al. (2021) [1] directly added the embeddings with the vectors that contain absolute position information. In contrast, Relative Position Bias (RPB) (Shaw et al., 2018) [6] shows that relative position information is more important for vision tasks and resorts to position information indicating the relative distance of two embeddings. Besides, Xiangxiang Chu et al. [12] pointed out that a successful positional encoding for vision tasks should meet the requirements that: (1) being permutation-variant but translation-invariant; (2) being able to handle different lengths of inputs; (3) containing absolute position information. Based on such requirements, they further proposed conditional positional encoding (CPE) that is dynamically generated and conditioned on the local neighborhood of the input tokens. CPE is used in latest transformers like MaxViT. Furthermore, since it could be implemented by a depthwise convolution layer, CPE is able to be combined with other designs like Squeeze-and-Excitation layer to remove explicit positional encodings layers.

Neural Network Design. There have been a number of works in CNN to design architectures that achieve a good trade-off between efficiency and accuracy. Works such as X3D [13] have proposed greedy methods to search for good hyperparameters like spatial resolution and depth for each stage of a CNN, which attach importance to the design choice other than a novel convolution design. However, to the best of our knowledge, such an exploration in vision transformers is clearly insufficient. AutoFormer [14] uses Neural Architecture Search (NAS) methods to find good embedding dimension for each self-attention layer. GLiT [15] searches a good permutation of global and local self-attention modules in the transformer. Anyway, most works that focus on transformer architecture design fix their group size as 7 in order to keep pace with Swin [4]. Such choices are made out of convenience and clearly not the best. Our work makes a small step forward through pointing out

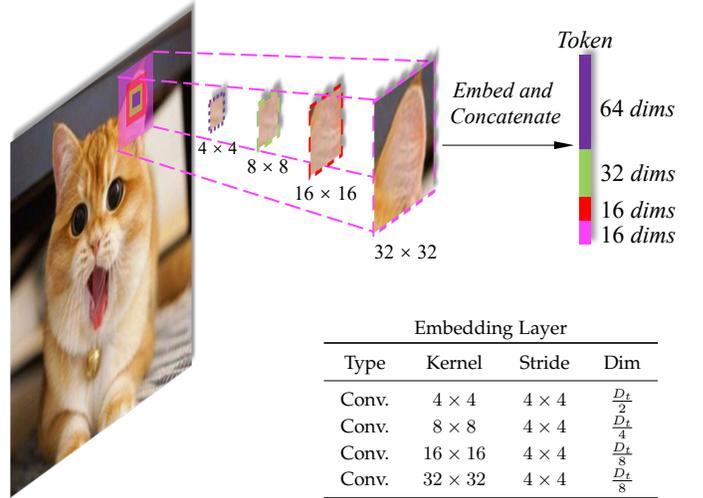


Fig. 2: Illustration of the CEL layer. The input image is sampled by four different kernels (*i.e.*, 4×4 , 8×8 , 16×16 , 32×32) with same stride 4×4 . Each embedding is constructed by embedding and concatenating the four patches. D_t means the total dimension of the embedding.

that choosing good group size for embeddings allows the transformers to achieve a better efficiency-accuracy tradeoff.

3 CROSSFORMER

The overall architecture of CrossFormer is plotted in Fig. 1. Following [3], [4], [5], CrossFormer also employs a pyramid structure, which naturally splits the transformer model into four stages. Each stage consists of a cross-scale embedding layer (CEL, Sec. 3.1) and several CrossFormer blocks (Sec. 3.2). A CEL receives last stage’s output (or an input image) as input and generates cross-scale tokens through an embedding layer. In this process, CEL (except that in $Stage-1$) reduces the number of embeddings to a quarter while doubles their dimensions for a pyramid structure. Then, several CrossFormer blocks, each of which involves long short distance attention (LSDA) and dynamic position bias (DPB), are set up after CEL. A specialized head (*e.g.*,

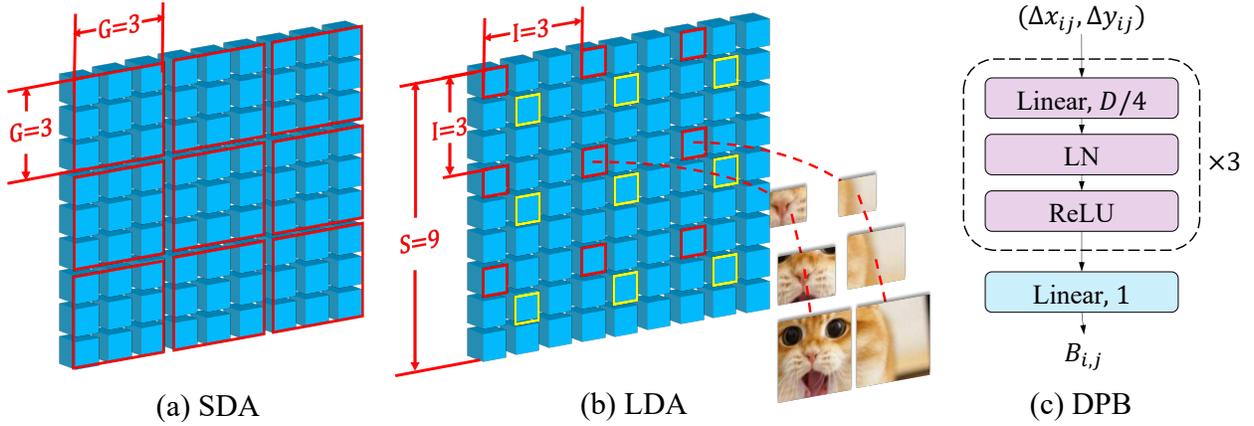


Fig. 3: (a) Short distance attention (SDA). Embeddings (blue cubes) are grouped by red boxes. (b) Long distance attention (LDA). Embeddings with the same color borders belong to the same group. Large patches of embeddings in the same group are adjacent. (c) Dynamic position bias (DPB). The dimensions of intermediate layers are $D/4$, and the output is a scalar.

the classification head in Fig. 1) follows after the final stage accounting for a specific task.

3.1 Cross-scale Embedding Layer (CEL)

Cross-scale embedding layer (CEL) is leveraged to generate input tokens for each stage. Fig. 2 takes the first CEL, which is ahead of *Stage-1*, as an example. It receives an image as input, then sampling patches using four kernels of different sizes. The stride of four kernels is kept the same so that they generate the same number of tokens². As we can observe in Fig. 2, every four corresponding patches have the same center but different scales, and all these four patches will be embedded and concatenated as one embedding. In practice, the process of sampling and embedding can be fulfilled through four convolutional layers.

For a cross-scale token, one problem is how to set the embedded dimension of each scale. The computational budget of a convolutional layer is proportional to $K^2 D^2$, where K and D represent kernel size and input/output dimension, respectively (assuming that the input dimension equals to the output dimension). Therefore, given the same dimension, a large kernel consumes a greater budget than a smaller one. To control the total budget of the CEL, we use a lower dimension for large kernels while a higher dimension for small kernels. Fig. 2 provides the specific allocation rule in its subtable, and a 128-dimensional example is given. Compared with allocating the dimension equally, our scheme saves much computational cost but does not explicitly affect the model’s performance. The cross-scale embedding layers in other stages work in a similar way. As shown in Fig. 1, CELs for *Stage-2/3/4* use two different kernels (2×2 and 4×4). Further, to form a pyramid structure, the strides of CELs for *Stage-2/3/4* are set as 2×2 to reduce the number of embeddings to a quarter.

3.2 CrossFormer Block

Each CrossFormer block consists of a long short distance attention module (*i.e.*, LSDA, which involves a short distance

attention (SDA) module or a long distance attention (LDA) module) and a multi-layer perceptron (MLP). As shown in Fig. 1b, SDA and LDA appear alternately in different blocks, and the dynamic position bias (DPB) module works in both SDA and LDA for obtaining embeddings’ position representations. Following the prior vision transformers, residual connections are used in each block.

3.2.1 Long Short Distance Attention (LSDA)

We split the self-attention module into two parts: short distance attention (SDA) and long distance attention (LDA). For SDA, all $G \times G$ adjacent embeddings are grouped together. Fig. 3a gives an example where $G = 3$. For LDA with input of size $S \times S$, the embeddings are sampled with a fixed interval I . For example in Fig. 3b ($I = 3$), all embeddings with a red border belong to a group, and those with a yellow border comprise another group. The group’s height or width for LDA is computed as $G = S/I$ (*i.e.*, $G = 3$ in this example). After grouping embeddings, both SDA and LDA employ the vanilla self-attention within each group. As a result, the memory/computational cost of the self-attention module is reduced from $O(S^4)$ to $O(S^2 G^2)$.

It is worth noting that the effectiveness of LDA also benefits from cross-scale embeddings. Specifically, we draw all the patches comprising two embeddings in Fig. 3b. As we can see, the small-scale patches of two embeddings are non-adjacent, so it is difficult to judge their relationship without the help of the context. In other words, it will be hard to build the dependency between these two embeddings if they are only constructed by small-scale patches (*i.e.*, single-scale feature). On the contrary, adjacent large-scale patches provide sufficient context to link these two embeddings, which makes long-distance cross-scale attention easier to compute and more meaningful.

3.2.2 Dynamic Position Bias (DPB)

Relative position bias (RPB) indicates embeddings’ relative positions by adding a bias to their attentions. Formally, the LSDA’s attention map with RPB becomes:

$$\text{Attn} = \text{Softmax}(QK^T / \sqrt{d} + B)V, \quad (1)$$

2. The image will be padded if necessary.

TABLE 1: Variants of CrossFormer and CrossFormer++ for image classification. The classification head just contains a global average pooling and a linear layer, which are not shown in the table below for simplicity. D , G , and I are token dimension, group size, and interval for SDA and LDA, respectively.

Version	Stage	Layer Name	Tiny (-T)	Small (-S)	Base (-B)	Large (-L)	Huge (-H)	
CrossFormer		Cross Embed.	Kernel size: $4 \times 4, 8 \times 8, 16 \times 16, 32 \times 32$, Stride=4					
	Stage-1	SDA/LDA MLP	$\begin{bmatrix} D_1 = 64 \\ G_1 = 7 \\ I_1 = 8 \end{bmatrix} \times 1$	$\begin{bmatrix} D_1 = 96 \\ G_1 = 7 \\ I_1 = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} D_1 = 96 \\ G_1 = 7 \\ I_1 = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} D_1 = 128 \\ G_1 = 7 \\ I_1 = 8 \end{bmatrix} \times 2$	–	
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2					
	Stage-2	SDA/LDA MLP	$\begin{bmatrix} D_2 = 128 \\ G_2 = 7 \\ I_2 = 4 \end{bmatrix} \times 1$	$\begin{bmatrix} D_2 = 192 \\ G_2 = 7 \\ I_2 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} D_2 = 192 \\ G_2 = 7 \\ I_2 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} D_2 = 256 \\ G_2 = 7 \\ I_2 = 4 \end{bmatrix} \times 2$	–	
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2					
	Stage-3	SDA/LDA MLP	$\begin{bmatrix} D_3 = 256 \\ G_3 = 7 \\ I_3 = 2 \end{bmatrix} \times 8$	$\begin{bmatrix} D_3 = 384 \\ G_3 = 7 \\ I_3 = 2 \end{bmatrix} \times 6$	$\begin{bmatrix} D_3 = 384 \\ G_3 = 7 \\ I_3 = 2 \end{bmatrix} \times 18$	$\begin{bmatrix} D_3 = 512 \\ G_3 = 7 \\ I_3 = 2 \end{bmatrix} \times 18$	–	
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2					
	Stage-4	SDA/LDA MLP	$\begin{bmatrix} D_4 = 512 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 6$	$\begin{bmatrix} D_4 = 768 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	$\begin{bmatrix} D_4 = 768 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	$\begin{bmatrix} D_4 = 1024 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	–	
CrossFormer++		Cross Embed.	Kernel size: $4 \times 4, 8 \times 8, 16 \times 16, 32 \times 32$, Stride=4					
	Stage-1	SDA/LDA MLP/ACL	–	$\begin{bmatrix} D_1 = 64 \\ G_1 = 4 \\ I_1 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} D_1 = 96 \\ G_1 = 4 \\ I_1 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} D_1 = 128 \\ G_1 = 4 \\ I_1 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} D_1 = 128 \\ G_1 = 4 \\ I_1 = 4 \end{bmatrix} \times 6$	
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2					
	Stage-2	SDA/LDA MLP/ACL	–	$\begin{bmatrix} D_2 = 128 \\ G_2 = 4 \\ I_2 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} D_2 = 192 \\ G_2 = 4 \\ I_2 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} D_2 = 256 \\ G_2 = 4 \\ I_2 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} D_2 = 256 \\ G_2 = 4 \\ I_2 = 4 \end{bmatrix} \times 6$	
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2					
	Stage-3	SDA/LDA MLP/ACL	–	$\begin{bmatrix} D_3 = 256 \\ G_3 = 14 \\ I_3 = 1 \end{bmatrix} \times 18$	$\begin{bmatrix} D_3 = 384 \\ G_3 = 14 \\ I_3 = 1 \end{bmatrix} \times 18$	$\begin{bmatrix} D_3 = 512 \\ G_3 = 14 \\ I_3 = 1 \end{bmatrix} \times 18$	$\begin{bmatrix} D_3 = 512 \\ G_3 = 14 \\ I_3 = 1 \end{bmatrix} \times 18$	
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2					
	Stage-4	SDA/LDA MLP/ACL	–	$\begin{bmatrix} D_4 = 512 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	$\begin{bmatrix} D_4 = 768 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	$\begin{bmatrix} D_4 = 1024 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	$\begin{bmatrix} D_4 = 1024 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	

where $Q, K, V \in \mathbb{R}^{G^2 \times D}$ represent *query, key, value* in the self-attention module, respectively. \sqrt{d} is a constant normalizer, and $B \in \mathbb{R}^{G^2 \times G^2}$ is the RPB matrix. In previous work [4], $B_{i,j} = \hat{B}_{\Delta x_{ij}, \Delta y_{ij}}$, where \hat{B} is a fixed-sized matrix, and $(\Delta x_{ij}, \Delta y_{ij})$ is the coordinate distance between the i_{th} and j_{th} embeddings. It is obvious that the image/group size is restricted in case that $(\Delta x_{ij}, \Delta y_{ij})$ exceeds the size of \hat{B} . In contrast, we propose an MLP-based module called DPB to generate the relative position bias dynamically, *i.e.*,

$$B_{i,j} = \text{DPB}(\Delta x_{ij}, \Delta y_{ij}). \quad (2)$$

The structure of DPB is displayed in Fig. 3c. Its non-linear transformation consists of three fully-connected layers with layer normalization [16] and ReLU [17]. The input dimension of DPB is 2, *i.e.*, $(\Delta x_{ij}, \Delta y_{ij})$, and intermediate layers' dimension is set as $D/4$, where D is the dimension of embeddings. The output B_{ij} is a scalar, encoding the relative position feature between the i_{th} and j_{th} embeddings. DPB is a trainable module optimized along with the whole transformer model. It can deal with any image/group size without worrying about the bound of $(\Delta x_{ij}, \Delta y_{ij})$.

In the appendix, we prove that DPB is equivalent to RPB if the image/group size is fixed. In this case, we can transform a trained DPB to RPB in the testing phase. We also provide an efficient $O(G^2)$ implementation of DPB when

image/group size is variable (the complexity is $O(G^4)$ in a normal case because $B \in \mathbb{R}^{G^2 \times G^2}$).

3.3 Variants of CrossFormer

TABLE 1 lists the detailed configurations of CrossFormer's four variants for image classification. To re-use the pre-trained weights, the models for other tasks (*e.g.*, object detection) employ the same backbones as classification except that they may use different G and I . Specifically, besides the configurations same to classification, we also test with $G_1 = G_2 = 14, I_1 = 16$, and $I_2 = 8$ for the detection (and segmentation) models' first two stages to adapt to larger images. Notably, the group size or the interval (*i.e.*, G or I) does not affect the shape of weight tensors, so the backbones pre-trained on ImageNet can be readily fine-tuned on other tasks even though they use different G or I .

4 CROSSFORMER++

In this section, we propose a progressive group size (PGS) paradigm to adapt to vision transformers' gradually expanding group size. Moreover, an amplitude cooling layer (ACL) is proposed to alleviate the amplitude explosion issue. The PGS and ACL are plugged into CrossFormer, yielding an improved version, dubbed CrossFormer++.

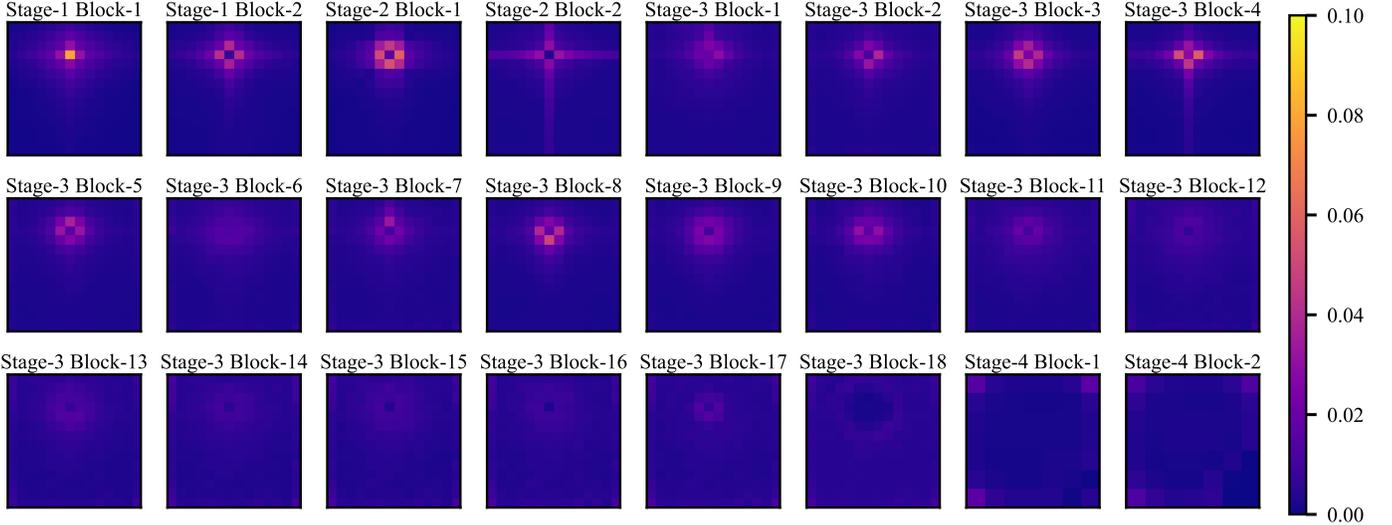


Fig. 4: The attention maps of a random token in CrossFormer-B’s blocks. The attention map size is 14×14 (except 7×7 for Stage-4). The attention concentrates around each token itself at shallow blocks but gradually disperses and distributes evenly at deep blocks.

4.1 Progressive Group Size (PGS)

Existing work [18] has explored the vanilla ViTs’ mechanisms and found that ViTs prefer global attentions even from early layers, which work in a different way from CNNs. However, vision transformers with a pyramid structure take smaller patches as input, resorting to group self-attention, and may perform differently from the vanilla ViTs. We take CrossFormer as an example and compute its average attention maps. Specifically, the attention maps of a certain group can be represented as:

$$Attn \in \mathbb{R}^{B \times H \times G \times G \times G \times G}, \quad (3)$$

where B, H, G represent batch size, number of heads, and group size, respectively. It means that there are $G \times G$ tokens in all, and that the attention map for each token is of size $G \times G$. The attention map of a token at image b , head h , and position (i, j) is represented as:

$$Attn_{b,h,i,j} \in \mathbb{R}^{G \times G}. \quad (4)$$

For the token at position (i, j) , the attention map averaged over batches and multi-heads is:

$$\overline{Attn}_{i,j} = \frac{1}{BH} \sum_{b=0}^B \sum_{h=0}^H Attn_{b,h,i,j} \in \mathbb{R}^{G \times G}. \quad (5)$$

We train a CrossFormer-B with a large group size of 14×14 and compute the average attention \overline{Attn} of each token. The results of a random token’s \overline{Attn} are shown in Fig. 4. As we can see, the attention regions gradually expand from shallow to deep layers. For example, tokens in the first two stages mainly attend to regions of size 4×4 around themselves, while the attention maps of deep layers from stage-3 are evenly distributed. The results indicate that tokens from shallow layers prefer local dependencies, while those from deep layers prefer global attentions.

To this end, we propose a PGS paradigm, *i.e.*, adopting a smaller group size in shallow layers to lower the computational budget and a larger group size in deep layers for

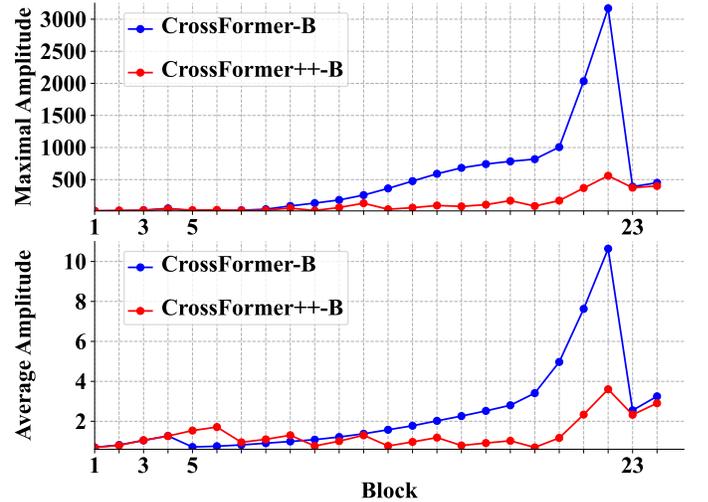


Fig. 5: Maximal/Average amplitude of each block’s output for CrossFormer-B and CrossFormer++-B. Both models use the same group size for all stages. CrossFormer++-B contains ACLs while CrossFormer-B does not.

global attentions. Under this guide, we first empirically set the group size to $\{4, 4, 14, 7\}$ ³ for four stages, respectively. Then, a linear scaling group size policy is proposed, *i.e.*, expanding the group size linearly from shallow to deep layers.

4.2 Amplitude Cooling Layer (ACL)

In addition to attention maps, we also explore the output’s amplitude of each block. As shown in Fig. 5, for a CrossFormer-B model, the amplitude increases greatly as the block goes deeper. In particular, the maximal output for the 22_{nd} block becomes over 3000, about 300 times larger

3. The last stage’s group size decreases because its feature maps size is 7×7 , and a group size of 7 already means a global attention.

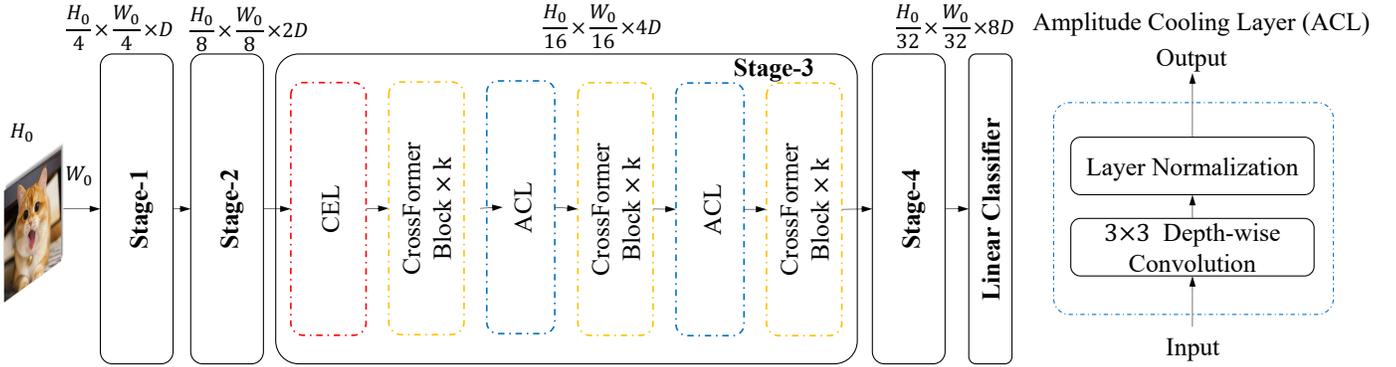


Fig. 6: The architecture of CrossFormer++ and its amplitude cooling layer (ACL).

than the value for the 1_{st} block. Besides, the average amplitude of the 22_{nd} also becomes about 15 times larger than that of the 1_{st} block. The extreme value makes the training process unstable and hinders the model from converging.

Fortunately, we also observe that the amplitude shrinks back to a small value at the start of each stage (e.g., the 5_{th} and 23_{th} blocks). We think that all block’s outputs are gradually accumulated through the residual connections in the model. The CEL at the beginning of each stage does not have a residual connection and cuts off the accumulation process, so it can effectively cool down the amplitude. While CEL still contains cumbersome normal convolution layers, we propose a more lightweight counterpart, dubbed amplitude cooling layer (ACL). As shown in Fig. 6, similar to CEL, an ACL does not use any residual connection, either. It only consists of a depth-wise convolution layer and a normalization layer. The comparison in Fig. 5 shows that ACL can also cool down the amplitude, but it introduces fewer parameters and a less computational budget than CEL because a depth-wise convolution with a small kernel (instead of a normal convolutional layer) is used.

However, ACL without a residual connection will prolong the back-propagation path and aggravate the vanishing gradient issue. To prevent this, we put an ACL layer after each k blocks with $k > 1$, as shown in Fig. 6. Empirically, $k = 3$ achieves a satisfying trade-off between amplitude cooling and back-propagation.

4.3 Variants of CrossFormer++

Armed with PGS and ACL, we further improve CrossFormer and propose CrossFormer++. The architectures are listed as TABLE 1. Wherein, CrossFormer++-B and CrossFormer++-L inherit each stage’s depth and channel from CrossFormer-B and CrossFormer-L, respectively. Besides, CrossFormer++ focuses more on larger models, so the “Tiny (-T)” version is ignored, and a new “Huge (-H)” version is constructed. Particularly, CrossFormer++-H puts more layers in the first two stages for refined low-level features. Moreover, we find that for the the “Small (-S)” version, a deep slim model performs similar to a shallow wide one but introduces fewer parameters, so CrossFormer++-S resorts to a deeper model than CrossFormer-S while using fewer channels for each block.

5 EXPERIMENTS

The experiments are carried out on four challenging tasks: image classification, object detection, instance segmentation, and semantic segmentation. To entail a fair comparison, we keep the same data augmentation and training settings as the other vision transformers as far as possible. The competitors are all competitive vision transformers, including DeiT [2], PVT [3], T2T-ViT [19], TNT [20], CViT [21], Twins [22], Swin [4], NesT [23], CvT [9], ViL [24], CAT [5], ResT [25], TransCNN [26], Shuffle [27], BoTNet [28], RegionViT [29], ViTAEv2 [30], MPViT [31], ScalableViT [32], DaViT [33], and CoAtNet [34].

5.1 Image Classification

Experimental Settings. The experiments on image classification are conducted on the ImageNet dataset. It contains 1.28M natural images for training and 50,000 images for evaluation. The images are resized to 224×224 for both training and evaluation by default. The same training settings as the other vision transformers are adopted. In particular, we use an AdamW [35] optimizer training for 300 epochs with a cosine decay learning rate scheduler, and 20 epochs of linear warm-up are used. The batch size is 1,024 split on 8 V100 GPUs. An initial learning rate of 0.001 and a weight decay of 0.05 are used. Besides, we use drop path rate of 0.1, 0.2, 0.3, 0.5 for CrossFormer-T, CrossFormer-S, CrossFormer-B, and CrossFormer-L, respectively. As well, the drop path rates of 0.2, 0.3, 0.5, 0.7 are used for CrossFormer++-S, CrossFormer++-B, CrossFormer++-L, and CrossFormer++-H, respectively. Further, similar to Swin [4], RandAugment [36], Mixup [37], Cutmix [38], random erasing [39], and stochastic depth [40] are used for data augmentation.

Results. The results are shown in TABLE 2. CrossFormer outperforms all other contemporaneous vision transformers. In specific, compared against strong baselines DeiT, PVT, and Swin, our CrossFormer outperforms them at least absolute 1.2% in accuracy on small models. Compared with CrossFormer, CrossFormer++ brings about 0.8% accuracy improvement in average. For example, CrossFormer-B achieves 83.4% in accuracy, while CrossFormer++-B reaches 84.2% with a negligible extra computational budget. Moreover, CrossFormer++ outperforms all existing vision transformers with similar parameters and a comparable computational budget.

TABLE 2: Results on the ImageNet validation set. The input size is 224×224 for most models, while it is 384×384 for the model with a \dagger . Results of other architectures are drawn from original papers. The blue numbers in brackets represent the improvements of CrossFormer++ over CrossFormer.

Architectures	#Params	FLOPs	Accuracy	Architectures	#Params	FLOPs	Accuracy
PVT-S	24.5M	3.8G	79.8%	BoTNet-S1-59	33.5M	7.3G	81.7%
RegionViT-T	13.8M	2.4G	80.4%	PVT-L	61.4M	9.8G	81.7%
Twins-SVT-S	24.0M	2.8G	81.3%	CvT-21	32.0M	7.1G	82.5%
CrossFormer-T	27.8M	2.9G	81.5%	CAT-B	52.0M	8.9G	82.8%
DeiT-S	22.1M	4.6G	79.8%	Swin-S	50.0M	8.7G	83.0%
T2T-ViT	21.5M	5.2G	80.7%	RegionViT-M	41.2M	7.4G	83.1%
CViT-S	26.7M	5.6G	81.0%	Twins-SVT-B	56.0M	8.3G	83.1%
PVT-M	44.2M	6.7G	81.2%	NesT-S	38.0M	10.4G	83.3%
TNT-S	23.8M	5.2G	81.3%	CrossFormer-B	52.0M	9.2G	83.4%
Swin-T	29.0M	4.5G	81.3%	CrossFormer+++B	52.0M	9.5G	84.2% (+0.8%)
NesT-T	17.0M	5.8G	81.5%	DeiT-B	86.0M	17.5G	81.8%
CvT-13	20.0M	4.5G	81.6%	DeiT-B \dagger	86.0M	55.4G	83.1%
ResT	30.2M	4.3G	81.6%	ViL-B	55.7M	13.4G	83.2%
CAT-S	37.0M	5.9G	81.8%	RegionViT-B	72.0M	13.3G	83.3%
ViL-S	24.6M	4.9G	81.8%	Twins-SVT-L	99.2M	14.8G	83.3%
RegionViT-S	30.6M	5.3G	82.5%	Swin-B	88.0M	15.4G	83.3%
Shuffle-T	29.0M	4.6G	82.5%	NesT-B	68.0M	17.9G	83.8%
ViTAEv2-S	19.2M	5.7G	82.6%	DaViT-Base	87.9M	15.5G	84.6%
Shunted-S	22.4M	4.9G	82.9%	CrossFormer-L	92.0M	16.1G	84.0%
MPViT-S	22.8M	4.7G	83.0%	CrossFormer+++L	92.0M	16.6G	84.7% (+0.7%)
ScalableViT-S	32.0M	4.2G	83.1%	ScalableViT-L	104.0M	14.7G	84.4%
CrossFormer-S	30.7M	4.9G	82.5%	CoAtNet	168.0M	34.7G	84.5%
CrossFormer+++S	23.3M	4.4G	83.2% (+0.8%)	CrossFormer+++H	96.0M	21.8G	84.9%

TABLE 3: Object detection results on the COCO 2017 *val* set with RetinaNets as detectors. Numbers in blue fonts represent the improvement of CrossFormer++ over CrossFormer. CrossFormers with \ddagger use different group sizes from classification models.

Method	Backbone	#Params	FLOPs	AP ^b	AP ₅₀ ^b	AP ₇₅ ^b	AP _S ^b	AP _M ^b	AP _L ^b
RetinaNet 1× schedule	PoolFormer-S36	40.6M	—	39.5	60.5	41.8	22.5	42.9	52.4
	CAT-B	62.0M	337.0G	41.4	62.9	43.8	24.9	44.6	55.2
	Swin-T	38.5M	245.0G	41.5	62.1	44.2	25.1	44.9	55.5
	ViL-M	50.8M	338.9G	42.9	64.0	45.4	27.0	46.1	57.2
	RegionViT-B	83.4M	308.9G	43.3	65.2	46.4	29.2	46.4	57.0
	TransCNN-B	36.5M	—	43.4	64.2	46.5	27.0	47.4	56.7
	DaViT-Tiny	—	244.0G	44.0	—	—	—	—	—
	PVTv2-B3	35.1M	—	44.6	65.6	47.6	27.4	48.8	58.6
	CrossFormer-S	40.8M	282.0G	44.4	65.8	47.4	28.2	48.4	59.4
	CrossFormer-S\ddagger	40.8M	272.1G	44.2	65.7	47.2	28.0	48.0	59.1
	CrossFormer+++S\ddagger	40.8M	272.1G	45.1 (+0.9)	66.6	48.5	28.7	49.4	60.3
	Twins-SVT-B	67.0M	322.0G	44.4	66.7	48.1	28.5	48.9	60.6
	RegionViT-B+	84.5M	328.2G	44.6	66.4	47.6	29.6	47.6	59.0
	Swin-B	98.4M	477.0G	44.7	65.9	49.2	—	—	—
	Twins-SVT-L	110.9M	455.0G	44.8	66.1	48.1	28.4	48.3	60.1
	ScalableViT-B	85.0M	330.0G	45.8	67.3	49.2	29.9	49.5	61.0
	DaViT-Small	—	332.0G	46.0	—	—	—	—	—
	PVTv2-B4	72.3M	—	46.1	66.9	49.2	28.4	50.0	62.2
	CrossFormer-B	62.1M	389.0G	46.2	67.8	49.5	30.1	49.9	61.8
	CrossFormer-B\ddagger	62.1M	379.1G	46.1	67.7	49.0	29.5	49.9	61.5
CrossFormer+++B\ddagger	62.2M	389.0G	46.6 (+0.5)	68.4	50.1	31.3	50.8	61.9	

5.2 Object Detection and Instance Segmentation

Experimental Settings. The experiments on object detection and instance segmentation are both done on the COCO 2017 dataset [41], which contains 118K training and 5K val images. We use MMDetection-based [42] RetinaNet [43] and Mask R-CNN [44] as the object detection and instance segmentation head, respectively. For both tasks, the backbones are initialized with the weights pre-trained on ImageNet. Then the whole models are trained with batch size 16 on 8 V100 GPUs, and an AdamW optimizer with an initial learning rate of 1×10^{-4} is used. Following previous works,

we adopt 1× training schedule (*i.e.*, the models are trained for 12 epochs) when taking RetinaNets as detectors, and images are resized to 800 pixels for the short side. While for Mask R-CNN, both 1× and 3× training schedules are used. It is noted that multi-scale training [45] is also employed when taking 3× training schedules.

Results. The results on RetinaNet and Mask R-CNN are shown in TABLE 3 and TABLE 4, respectively. As we can see, CrossFormer outperforms most existing vision transformers and shows higher superiority than on image classification task. We think it is because CrossFormer utilizes cross-scale

TABLE 4: Object detection and instance segmentation results on COCO *val* 2017 with Mask R-CNNs as detectors. AP^b and AP^m are box average precision and mask average precision, respectively.

Method	Backbone	#Params	FLOPs	AP^b	AP^b_{50}	AP^b_{75}	AP^m	AP^m_{50}	AP^m_{75}
Mask R-CNN 1× schedule	PVT-M	63.9M	–	42.0	64.4	45.6	39.0	61.6	42.0
	Swin-T	47.8M	264.0G	42.2	64.6	46.2	39.1	61.6	42.0
	Twins-PCPVT-S	44.3M	245.0G	42.9	65.8	47.1	40.0	62.7	42.9
	TransCNN-B	46.4M	–	44.0	66.4	48.5	40.2	63.3	43.2
	ViL-M	60.1M	261.1G	43.3	65.9	47.0	39.7	62.8	42.0
	RegionViT-B	92.2M	287.9G	43.5	66.7	47.4	40.1	63.4	43.0
	RegionViT-B+	93.2M	307.1G	44.5	67.6	48.7	41.0	64.4	43.9
	PVTv2-B2	45.0M	–	45.3	67.1	49.6	41.2	64.2	44.4
	ELSA-T	49.0M	269.0G	45.6	67.9	50.3	41.1	64.8	44.0
	ScalableViT-S	46.0M	256.0G	45.8	67.6	50.0	41.7	64.7	44.8
	CrossFormer-S	50.2M	301.0G	45.4	68.0	49.7	41.4	64.8	44.6
	CrossFormer-S[‡]	50.2M	291.1G	45.0	67.9	49.1	41.2	64.6	44.3
	CrossFormer++-S[‡]	43.0M	287.4G	46.4 (+1.4)	68.8	51.3	42.1 (+0.9)	65.7	45.4
	CAT-B	71.0M	356.0G	41.8	65.4	45.2	38.7	62.3	41.4
	PVT-L	81.0M	364.0G	42.9	65.0	46.6	39.5	61.9	42.5
	Twins-SVT-B	76.3M	340.0G	45.1	67.0	49.4	41.1	64.1	44.4
	ViL-B	76.1M	365.1G	45.1	67.2	49.3	41.0	64.3	44.2
	Twins-SVT-L	119.7M	474.0G	45.2	67.5	49.4	41.2	64.5	44.5
	Swin-S	69.1M	354.0G	44.8	66.6	48.9	40.9	63.4	44.2
	Swin-B	107.2M	496.0G	45.5	–	–	41.3	–	–
PVTv2-B4	82.2M	–	47.5	68.7	52.0	42.7	66.1	46.1	
ScalableViT-B	95.0M	349.0G	46.8	68.7	51.5	42.5	65.8	45.9	
CrossFormer-B	71.5M	407.9G	47.2	69.9	51.8	42.7	66.6	46.2	
CrossFormer-B[‡]	71.5M	398.1G	47.1	69.9	52.0	42.7	66.5	46.1	
CrossFormer++-B[‡]	71.5M	408.0G	47.7 (+0.6)	70.2	52.7	43.2 (+0.5)	67.3	46.7	
Mask R-CNN 3× schedule	PVT-M	63.9M	–	44.2	66.0	48.2	45.0	63.1	43.5
	ViL-M	60.1M	261.1G	44.6	66.3	48.5	40.7	63.8	43.7
	Swin-T	47.8M	264.0G	46.0	68.2	50.2	41.6	65.1	44.8
	Shuffle-T	48.0M	268.0G	46.8	68.9	51.5	42.3	66.0	45.6
	ViTAEv2-S	37.0M	–	47.8	69.4	52.2	42.6	66.6	45.8
	ScalableViT-S	46.0M	256.0G	48.7	70.1	53.6	43.6	67.2	47.2
	CrossFormer-S[‡]	50.2M	291.1G	48.7	70.7	53.7	43.9	67.9	47.3
	CrossFormer++-S[‡]	43.0M	287.4G	49.5 (+0.8)	71.6	54.1	44.3 (+0.4)	68.5	47.6
	PVT-L	81.0M	364.0G	44.5	66.0	48.3	40.7	63.4	43.7
	ViL-B	76.1M	365.1G	45.7	67.2	49.9	41.3	64.4	44.5
	Shuffle-S	69.0M	359.0G	48.4	70.1	53.5	43.3	67.3	46.7
	Swin-S	69.1M	354.0G	48.5	70.2	53.5	43.3	67.3	46.6
	ScalableViT-B	95.0M	349.0G	49.0	70.3	53.6	43.8	67.4	47.5
	Shunted-B	59.0M	–	50.1	70.9	54.1	45.2	68.0	48.0
	CrossFormer-B[‡]	71.5M	398.1G	49.8	71.6	54.9	44.5	68.8	47.9
	CrossFormer++-B[‡]	71.5M	408.0G	50.2 (+0.4)	71.8	54.9	44.6 (+0.1)	68.7	48.1

features explicitly, and cross-scale features are particularly vital for dense prediction tasks like objection detection and instance segmentation. However, there are still some other architectures that perform better than CrossFormer, *e.g.*, PVTv2-B3 achieves 0.2% AP higher than CrossFormer-S. Instead, CrossFormer++ surpasses CrossFormer by at least 0.5% AP and outperforms all existing methods. Further, its performance gain over the other architectures gets sharper when enlarging the model, indicating that CrossFormer++ enjoys greater potentials. For example, CrossFormer++-S outperforms ScalableViT-S by 0.8 in AP (48.7 vs. 49.5) when using Mask R-CNN as the detection head, while the CrossFormer++-B outperforms ScalableViT-B by 1.2 in AP.

5.3 Semantic Segmentation

Experimental Settings. ADE20K [46] is used as the benchmark for semantic segmentation. It covers a broad range of 150 semantic categories, including 20K images for training and 2K for validation. Similar to models for detection, we initialize the backbones with weights pre-trained on ImageNet, and MMSegmentation-based [47] semantic FPN

and UPerNet [48] are taken as the segmentation head. For FPN [49], we use an AdamW optimizer with learning rate and weight decay of 1×10^{-4} . Models are trained for 80K iterations with batch size 16. For UPerNet, an AdamW optimizer with an initial learning rate of 6×10^{-5} and a weight decay of 0.01 is used, and models are trained for 160K iterations.

Results. All results are shown in TABLE 5. Compared with architectures released earlier, CrossFormer exhibits a greater performance gain over the others when enlarging the model, which is similar to object detection. For example, CrossFormer-T achieves 1.4% absolutely higher in IOU than Twins-SVT-B, but CrossFormer-B achieves 3.1% absolutely higher in IOU than Twins-SVT-L. Totally, CrossFormer shows a more significant advantage over the others on dense prediction tasks (*e.g.*, detection and segmentation) than on classification, implying that cross-scale interactions in the attention module are more important for dense prediction tasks than for classification.

Moreover, CrossFormer++ shows a more significant advantage than CrossFormer when using a more powerful

TABLE 5: Semantic segmentation results on the ADE20K validation set. “MS IOU” means testing with variable input size.

Models	Semantic FPN (80K iterations)		
	Parameters	FLOPs	IOU
PoolFormer-S36	34.6M	–	42.0
Twins-SVT-B	60.4M	261.0G	45.0
Swin-S	53.2M	274.0G	45.2
ScalableViT-S	30.0M	174.0G	44.9
Shunted-S	26.1M	183.0G	48.2
CrossFormer-S	34.3M	220.7G	46.0
CrossFormer-S[‡]	34.3M	209.8G	46.4
CrossFormer++-S[‡]	27.1M	199.5G	47.4 (+1.0)
PoolFormer-M36	59.8M	–	42.4
PVTv2-B3	49.0M	249.6G	47.3
ScalableViT-B	79.0M	270.0G	48.4
CrossFormer-B	55.6M	331.0G	47.7
CrossFormer-B[‡]	55.6M	320.1G	48.0
CrossFormer++-B[‡]	55.6M	331.1G	48.6 (+0.6)
Twins-SVT-L	103.7M	397.0G	45.8
PVTv2-B5	85.7M	364.4G	48.7
ScalableViT-L	105.0M	402.0G	49.4
CrossFormer-L	95.4M	497.0G	48.7
CrossFormer-L[‡]	95.4M	482.7G	49.1
CrossFormer++-L[‡]	95.5M	482.8G	49.5 (+0.4)
CrossFormer++-H[‡]	99.4M	612.3G	49.7

TABLE 6: Ablation studies on different kernel sizes of CELs. The results are tested on the ImageNet validation set and the COCO instance segmentation task, and the baseline model is CrossFormer-S (82.5%) with Mask R-CNN being the segmentation head (41.4 AP).

Stage-1	CEL’s Kernel Size			#Params/FLOPs	Acc.	AP ^m
	Stage-2	Stage-3	Stage-4			
4	2	2	2	28.3M / 4.5G	81.5%	39.7
8	2	2	2	28.3M / 4.5G	81.9%	40.2
4, 8	2, 4	2, 4	2, 4	30.6M / 4.8G	82.3%	–
4, 8, 16, 32	2, 4	2, 4	2, 4	30.7M / 4.9G	82.5%	41.4
4, 8, 16, 32	2, 4, 8	2, 4	2	29.4M / 5.0G	82.4%	–

segmentation head. As we can see in TABLE 5, UPerNet is a more powerful model than semantic FPN. CrossFormer++S outperforms CrossFormer-S by 1.0 AP when using semantic FPN (47.4 vs. 46.4), while outperforms 2.0 AP when using UPerNet as the segmentation head.

5.4 Ablation Studies

5.4.1 Cross-scale Tokens vs. Single-scale Tokens

We conduct the experiments by replacing cross-scale embedding layers with single-scale ones. As we can see in TABLE 6, when using single-scale embeddings, the 8×8 kernel in Stage-1 brings 0.4% (81.9% vs. 81.5%) absolute improvement compared with the 4×4 kernel. It tells us that overlapping receptive fields help improve the model’s performance. Besides, all models with cross-scale embeddings perform better than those with single-scale embeddings. In particular, our CrossFormer achieves 1% (82.5% vs. 81.5%) absolute performance gain compared with using single-scale embeddings for all stages. For cross-scale embeddings, we also try several different combinations of kernel sizes, and they all show similar performance (82.3% \sim 82.5%). In summary, cross-scale embeddings can bring a large performance gain, yet the model is relatively robust with respect to different choices of kernel size.

Models	UPerNet (160K iterations)			
	Parameters	FLOPs	IOU	MS IOU
ELSA-Swin-T	61.0M	946.0G	–	47.7
MPViT-S	52.0M	943.0G	–	48.3
ScalableViT-S	57.0M	931.0G	48.5	49.4
Shunted-S	52.0M	940.0G	48.9	49.9
CrossFormer-S	62.3M	979.5G	47.6	48.4
CrossFormer-S[‡]	62.3M	968.5G	47.4	48.2
CrossFormer++-S[‡]	53.1M	963.5G	49.4 (+2.0)	50.0 (+1.8)
Swin-S	81.0M	1038.0G	47.6	49.5
ELSA-Swin-S	85.0M	1046.0G	–	50.3
ViTAEv2-S	49.0M	–	45.0	48.0
MPViT-B	105.0M	1186.0G	–	50.3
ScalableViT-B	107.0M	1029.0G	49.5	50.4
CrossFormer-B	83.6M	1089.7G	49.7	50.6
CrossFormer-B[‡]	83.6M	1078.8G	49.2	50.1
CrossFormer++-B[‡]	83.7M	1089.8G	50.7 (+1.5)	51.0 (+0.9)
Swin-B	121.0M	1088.0G	48.1	49.7
ScalableViT-L	135.0M	1162.0G	49.8	50.7
CrossFormer-L	125.5M	1257.8G	50.4	51.4
CrossFormer-L[‡]	125.5M	1243.5G	50.5	51.4
CrossFormer++-L[‡]	125.5M	1257.9G	51.0 (+0.5)	51.9 (+0.5)
CrossFormer++-H[‡]	129.5M	1416.1G	51.2	51.8

TABLE 7: Ablation studies on different self-attention mechanisms. The results are tested on the ImageNet validation set, and the baseline model is CrossFormer-S (82.5%), *i.e.*, LSDA with CEL and DPB.

Self-attentions	CEL	DPB	Accuracy
CAT	✓	✓	81.1%
PVT	✓	–	81.3%
LSDA	–	–	81.5%
Swin	✓	✓	81.9%
DaViT	✓	✓	82.1%
CSwin	✓	✓	82.5%
ScalableViT	✓	✓	82.5%
LSDA	✓	✓	82.5%

5.4.2 LSDA vs. Other Self-attentions

Self-attention mechanisms used in other vision transformers are also compared. Specifically, we replace LSDA in CrossFormer-S with self-attention modules proposed in previous work, and CEL and DPB are retained (except for PVT because DPB cannot be applied to PVT). As shown in TABLE 7, the self-attention mechanisms have great influences on the models’ accuracies. In particular, LSDA outperforms CAT-like self-attention by 1.4% (82.5% vs. 81.1%). Besides, some self-attention mechanisms can achieve similar accuracy to LSDA (CSwin and ScalableViT), which indicates that the most proper self-attention mechanism is not unique.

5.4.3 Ablation Studies about PGS and ACL

Manually Designed Group Size. We adopt manually designed group size for the CrossFormer++, *i.e.*, [4, 4, 14, 7] for four stages, respectively. We also test other group size, and the results are shown in TABLE 8. Compare CrossFormer++-S with CrossFormer++-S₃ and CrossFormer++-S₄, and the results show that a 7×7 group size for the first two stages is inessential, which leads to a greater computational budget but no accuracy improvement. The comparison between CrossFormer++-S and CrossFormer++-S₁ (83.2% vs 82.6%) shows that a large group size for Stage-3 is pivotal. These

TABLE 8: Results on the ImageNet validation set. We test with different group sizes and ACL Layers. $x \rightarrow y$ indicates linearly scaling from x to y .

Models	ACL	Group Size				Acc.
		Stage-1	Stage-2	Stage-3	Stage-4	
CrossFormer-S ₁		4 × 4	4 × 4	7 × 7	7 × 7	82.4%
CrossFormer-S ₂		7 × 7	7 × 7	7 × 7	7 × 7	82.5%
CrossFormer-S ₃		4 × 4	4 × 4	14 × 14	7 × 7	82.9%
CrossFormer++-S ₁	✓	4 × 4	4 × 4	7 × 7	7 × 7	82.6%
CrossFormer++-S ₂	✓	4 × 4 → 14 × 14		7 × 7		82.9%
CrossFormer++-S ₃	✓	7 × 7	7 × 7	14 × 14	7 × 7	83.2%
CrossFormer++-S ₄	✓	4 × 4	7 × 7	14 × 14	7 × 7	83.2%
CrossFormer++-S	✓	4 × 4	4 × 4	14 × 14	7 × 7	83.2%
CrossFormer-B		4 × 4	4 × 4	14 × 14	7 × 7	83.9%
CrossFormer++-B	✓	4 × 4	4 × 4	14 × 14	7 × 7	84.2%
CrossFormer-L		4 × 4	4 × 4	14 × 14	7 × 7	84.3%
CrossFormer++-L	✓	4 × 4	4 × 4	14 × 14	7 × 7	84.7%
Swin-T		7 × 7	7 × 7	7 × 7	7 × 7	81.3%
Swin-T ₁		4 × 4	4 × 4	14 × 14	7 × 7	82.8%
Swin-T ₂	✓	4 × 4	4 × 4	14 × 14	7 × 7	83.2%

TABLE 9: Comparisons between different position representations. The base model is CrossFormer-S. Flexibility indicates whether the representation applies to dynamic group size.

Method	#Params/FLOPs	Throughput	Flexibility	Acc.
APE	30.9M/4.9G	686 imgs/sec		82.1%
RPB	30.7M/4.9G	684 imgs/sec		82.5%
DPB	30.7M/4.9G	672 imgs/sec	✓	82.5%
DPB-residual	30.7M/4.9G	672 imgs/sec	✓	82.4%

conclusions also coincide with our observation that the self-attention at shallow layers concentrates on a small region around each token, while the attention gradually disperses at deep layers.

Linearly Scaling Group Size. Additionally, we also try a simple non-manually designed group size, *i.e.*, we expand group size from 4×4 to 14×14 linearly from Stage-1 to Stage-3. However, the linear group size does not work as well as a manually designed group size.

Attention Mechanisms vs. Group Size. It is also worth emphasizing that the conclusions about group size also apply to other vision transformers. As shown in TABLE 8, an appropriate group size for Swin-T also brings a significant accuracy improvement, from 81.3% to 82.8%. In comparison, replacing shifted window in Swin with LSDA only brings 0.6% improvement (as shown in TABLE 7), which indicates that an appropriate group size may be more important than the choice of the self-attention mechanism.

Ablation Studies about ACL. Regarding ACL, the ACL layer brings about 0.3%~0.4% accuracy improvement. Concretely, CrossFormer++-B improves from 83.9% to 84.2% after plugging an ACL, and CrossFormer++-L improves from 84.3% to 84.7%. Moreover, ACL is a universal layer that also applies to other vision transformers. For example, plugging an ACL into Swin-T can bring 0.4% accuracy improvement (82.8% vs. 83.2%).

5.4.4 DPB vs. Other Position Representations

We compare the parameters, FLOPs, throughputs, and accuracies of the models among absolute position embedding (APE), relative position bias (RPB), and DPB. The results

are shown in TABLE 9. DPB-residual means DPB with residual connections. Both DPB and RPB outperform APE with absolute 0.4% accuracy, which indicates that relative position representations are more beneficial than absolute ones.

Further, DPB achieves the same accuracy (82.5%) as RPB with an ignorable extra cost; however, as we described in Sec. 3.2.2, it is more flexible than RPB and applies to variable image size or group size. Besides, the results also show that residual connection in DPB does not help improve or even degrades the model’s performance (82.5% vs. 82.4%).

6 CONCLUSIONS AND FUTURE WORKS

In this paper, we first proposed a universal visual backbone, dubbed CrossFormer. It utilizes cross-scale features through our proposed CEL and LSDA. The experimental results show that utilizing cross-scale features explicitly can significantly improve the vision transformers’ performance on image classification and other downstream tasks. Besides, a more flexible relation position representation, DPB, is also proposed to make the CrossFormer apply to dynamic group size. Based on CrossFormer, we further analyzed the self-attention module’s and MLPs’ outputs in the CrossFormer and proposed progressive group size (PGS) and amplitude cooling layer (ACL). Based on PGS and ACL, CrossFormer++ achieves better performance than CrossFormer and outperforms all the existing state-of-the-art vision transformers on representative visual tasks. Besides, extensive experiments show that CEL, PGS, and ACL are all universal and can consistently bring performance gains when plugged into other vision transformers.

Despite these contributions, the proposed modules still have some limitations. Concretely, though we demonstrated the importance of a progressive appropriate group size, an empirically designed group size is used finally. An adaptive and automated group size policy is still expected. Then, to prevent the gradient vanishing issue, the ACL layer cannot be used in large quantities in vision transformers. So, we will also explore how to cool down the vision transformers’ amplitude without cutting the residual connection.

ACKNOWLEDGMENTS

This work was supported in part by The National Nature Science Foundation of China (Grant Nos.: U1909203, 62273303), in part by Innovation Capability Support Program of Shaanxi (Program No. 2021TD-05).

REFERENCES

- [1] 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 *International Conference on Learning Representations, ICLR*, 2021.
- [2] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, “Training data-efficient image transformers & distillation through attention,” in *International Conference on Machine Learning, ICML*, vol. 139, 2021, pp. 10 347–10 357.
- [3] W. Wang, E. Xie, X. Li, D. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, “Pyramid vision transformer: A versatile backbone for dense prediction without convolutions,” *CoRR*, vol. abs/2102.12122, 2021.

- [4] 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," *CoRR*, vol. abs/2103.14030, 2021.
- [5] H. Lin, X. Cheng, X. Wu, F. Yang, D. Shen, Z. Wang, Q. Song, and W. Yuan, "CAT: cross attention in vision transformer," *CoRR*, vol. abs/2106.05786, 2021.
- [6] P. Shaw, J. Uszkoreit, and A. Vaswani, "Self-attention with relative position representations," in *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL*, 2018, pp. 464–468.
- [7] W. Wang, L. Yao, L. Chen, B. Lin, D. Cai, X. He, and W. Liu, "Crossformer: A versatile vision transformer hinging on cross-scale attention," in *International Conference on Learning Representations, ICLR*, 2022.
- [8] Y. Li, C. Wu, H. Fan, K. Mangalam, B. Xiong, J. Malik, and C. Feichtenhofer, "Mvity2: Improved multiscale vision transformers for classification and detection," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4794–4804, 2021.
- [9] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, "Cvt: Introducing convolutions to vision transformers," *CoRR*, vol. abs/2103.15808, 2021.
- [10] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 124–12 134.
- [11] Z. Pan, B. Zhuang, J. Liu, H. He, and J. Cai, "Scalable visual transformers with hierarchical pooling," *CoRR*, vol. abs/2103.10619, 2021.
- [12] X. Chu, Z. Tian, B. Zhang, X. Wang, X. Wei, H. Xia, and C. Shen, "Conditional positional encodings for vision transformers," *arXiv preprint arXiv:2102.10882*, 2021.
- [13] C. Feichtenhofer, "X3d: Expanding architectures for efficient video recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 203–213.
- [14] M. Chen, H. Peng, J. Fu, and H. Ling, "Autoformer: Searching transformers for visual recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 12 270–12 280.
- [15] B. Chen, P. Li, C. Li, B. Li, L. Bai, C. Lin, M. Sun, J. Yan, and W. Ouyang, "Glit: Neural architecture search for global and local image transformer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12–21.
- [16] L. J. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *CoRR*, vol. abs/1607.06450, 2016.
- [17] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *International Conference on Machine Learning, ICML*, J. Fürnkranz and T. Joachims, Eds., 2010, pp. 807–814.
- [18] M. Raghu, T. Unterthiner, S. Kornblith, C. Zhang, and A. Dosovitskiy, "Do vision transformers see like convolutional neural networks?" in *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021, pp. 12 116–12 128. [Online]. Available: <https://proceedings.neurips.cc/paper/2021/hash/652cf38361a209088302ba2b8b7f51e0-Abstract.html>
- [19] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, F. E. H. Tay, J. Feng, and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," *CoRR*, vol. abs/2101.11986, 2021.
- [20] K. Han, A. Xiao, E. Wu, J. Guo, C. Xu, and Y. Wang, "Transformer in transformer," *CoRR*, vol. abs/2103.00112, 2021.
- [21] C. Chen, Q. Fan, and R. Panda, "Crossvit: Cross-attention multi-scale vision transformer for image classification," *CoRR*, vol. abs/2103.14899, 2021.
- [22] X. Chu, Z. Tian, Y. Wang, B. Zhang, H. Ren, X. Wei, H. Xia, and C. Shen, "Twins: Revisiting spatial attention design in vision transformers," *CoRR*, vol. abs/2104.13840, 2021.
- [23] Z. Zhang, H. Zhang, L. Zhao, T. Chen, and T. Pfister, "Aggregating nested transformers," *CoRR*, vol. abs/2105.12723, 2021.
- [24] P. Zhang, X. Dai, J. Yang, B. Xiao, L. Yuan, L. Zhang, and J. Gao, "Multi-scale vision longformer: A new vision transformer for high-resolution image encoding," *CoRR*, vol. abs/2103.15358, 2021.
- [25] Q. Zhang and Y. Yang, "Rest: An efficient transformer for visual recognition," *CoRR*, vol. abs/2105.13677, 2021.
- [26] Y. Liu, G. Sun, Y. Qiu, L. Zhang, A. Chhatkuli, and L. V. Gool, "Transformer in convolutional neural networks," *CoRR*, vol. abs/2106.03180, 2021.
- [27] Z. Huang, Y. Ben, G. Luo, P. Cheng, G. Yu, and B. Fu, "Shuffle transformer: Rethinking spatial shuffle for vision transformer," *CoRR*, vol. abs/2106.03650, 2021.
- [28] A. Srinivas, T. Lin, N. Parmar, J. Shlens, P. Abbeel, and A. Vaswani, "Bottleneck transformers for visual recognition," in *Conference on Computer Vision and Pattern Recognition, CVPR*, 2021.
- [29] C. Chen, R. Panda, and Q. Fan, "Regionvit: Regional-to-local attention for vision transformers," *CoRR*, vol. abs/2106.02689, 2021.
- [30] Q. Zhang, Y. Xu, J. Zhang, and D. Tao, "Vitaev2: Vision transformer advanced by exploring inductive bias for image recognition and beyond," *CoRR*, vol. abs/2202.10108, 2022.
- [31] Y. Lee, J. Kim, J. Willette, and S. J. Hwang, "Mpvit: Multi-path vision transformer for dense prediction," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. IEEE, 2022, pp. 7277–7286.
- [32] R. Yang, H. Ma, J. Wu, Y. Tang, X. Xiao, M. Zheng, and X. Li, "Scalablevit: Rethinking the context-oriented generalization of vision transformer," in *Computer Vision - ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXIV*, ser. Lecture Notes in Computer Science, S. Avidan, G. J. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, Eds., vol. 13684. Springer, 2022, pp. 480–496.
- [33] M. Ding, B. Xiao, N. Codella, P. Luo, J. Wang, and L. Yuan, "Davvit: Dual attention vision transformers," in *Computer Vision - ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXIV*, ser. Lecture Notes in Computer Science, S. Avidan, G. J. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, Eds., vol. 13684. Springer, 2022, pp. 74–92.
- [34] Z. Dai, H. Liu, Q. V. Le, and M. Tan, "Coatnet: Marrying convolution and attention for all data sizes," in *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021, pp. 3965–3977.
- [35] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations, ICLR*, 2015.
- [36] E. D. Cubuk, B. Zoph, J. Shlens, and Q. Le, "Randaugment: Practical automated data augmentation with a reduced search space," in *Neural Information Processing Systems, NeurIPS*, 2020.
- [37] H. Zhang, M. Cissé, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," in *International Conference on Learning Representations, ICLR*, 2018.
- [38] S. Yun, D. Han, S. Chun, S. J. Oh, Y. Yoo, and J. Choe, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *International Conference on Computer Vision, ICCV*, 2019, pp. 6022–6031.
- [39] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," in *Association for the Advancement of Artificial Intelligence, AAAI*, 2020, pp. 13 001–13 008.
- [40] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, "Deep networks with stochastic depth," in *European Conference on Computer Vision, ECCV*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., vol. 9908, 2016, pp. 646–661.
- [41] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," in *European Conference on Computer Vision, ECCV*, vol. 8693, 2014, pp. 740–755.
- [42] K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu, Z. Zhang, D. Cheng, C. Zhu, T. Cheng, Q. Zhao, B. Li, X. Lu, R. Zhu, Y. Wu, J. Dai, J. Wang, J. Shi, W. Ouyang, C. C. Loy, and D. Lin, "MMDetection: Open mmlab detection toolbox and benchmark," *arXiv preprint arXiv:1906.07155*, 2019.
- [43] T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *Transactions on Pattern Analysis and Machine Intelligence, PAMI*, vol. 42, no. 2, pp. 318–327, 2020.
- [44] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask R-CNN," in *International Conference on Computer Vision, ICCV*, 2017, pp. 2980–2988.
- [45] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Computer Vision European Conference, ECCV*, vol. 12346, 2020, pp. 213–229.

- [46] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Scene parsing through ADE20K dataset," in *Conference on Computer Vision and Pattern Recognition, CVPR, 2017*, pp. 5122–5130.
- [47] M. Contributors, "MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark," 2020.
- [48] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun, "Unified perceptual parsing for scene understanding," in *European Conference on Computer Vision, ECCV*, vol. 11209, 2018, pp. 432–448.
- [49] A. Kirillov, R. B. Girshick, K. He, and P. Dollár, "Panoptic feature pyramid networks," in *Conference on Computer Vision and Pattern Recognition, CVPR, 2019*, pp. 6399–6408.



Boxi Wu is a distinguished research fellow at the School of Software Technology, Zhejiang University. He received his bachelor and Ph.D. degree in computer science and technology from Zhejiang University in 2016 and 2022. His research interests include machine learning and computer vision.



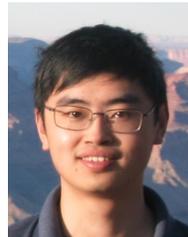
Wenxiao Wang is a distinguished research fellow, School of Software Technology of Zhejiang University, China. He received the Ph.D. degree in computer science and technology from Zhejiang University in 2022. His research interests include deep learning and computer vision.



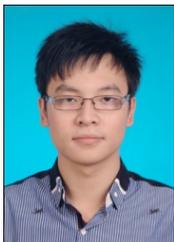
Binbin Lin is an assistant professor in the School of Software Technology at Zhejiang University, China. He received a Ph.D. degree in computer science from Zhejiang University in 2012. His research interests include machine learning and decision making.



Wei Chen is a senior student from College of Computer Science and Technology at Zhejiang University, China. He majors in Artificial Intelligence and is currently an intern in the State Key Lab of CAD&CG at Zhejiang University. His research interests include deep learning and computer vision.



Xiaofei He received the B.S. degree in Computer Science from Zhejiang University, China, in 2000 and the Ph.D. degree in Computer Science from the University of Chicago, in 2005. He is a Professor in the State Key Lab of CAD&CG at Zhejiang University, China. Prior to joining Zhejiang University, he was a Research Scientist at Yahoo! Research Labs, Burbank, CA. His research interests include machine learning, information retrieval, and computer vision. He is a senior member of IEEE.



Qibo Qiu received the M.S. degree in computer science at Zhejiang University, Hangzhou, China, in 2017. He is currently a senior algorithm engineer with Zhejiang Laboratory, his research interests include road marker detection, 2D/3D lane detection, general place recognition, localization and perception for intelligent vehicles and mobile robots.



Wei Liu (M'14-SM'19-F'23) received the Ph.D. degree from Columbia University in 2012 in Electrical Engineering and Computer Science. He is currently a Distinguished Scientist of Tencent and the Director of Ads Multimedia AI at Tencent Data Platform. Prior to that, he has been a research staff member of IBM T. J. Watson Research Center, USA from 2012 to 2015. Dr. Liu has long been devoted to fundamental research and technology development in core fields of AI, including deep learning, machine learning, computer vision, pattern recognition, information retrieval, big data, etc. To date, he has published extensively in these fields with more than 280 peer-reviewed technical papers, and also issued 32 US patents. He currently serves on the editorial boards of IEEE TPAMI, TNNLS, and IEEE Intelligent Systems. He is an Area Chair of top-tier computer science and AI conferences, e.g., NeurIPS, ICML, IEEE CVPR, IEEE ICCV, IJCAI, and AAAI. Dr. Liu is a Fellow of the IEEE, IAPR, and IMA.



Long Chen received the Ph.D. degree in Computer Science from Zhejiang University in 2020, and the B.Eng. degree in Electrical Information Engineering from Dalian University of Technology in 2015. He is currently a postdoctoral research scientist at Columbia University. He is also an incoming assistant professor at Hong Kong University of Science and Technology. He was a senior researcher at Tencent AI Lab. His research interests are computer vision, machine learning, and multimedia.

Appendices for “CrossFormer++: A Versatile Vision Transformer Hinging on Cross-scale Attention”



APPENDIX A EFFICIENT DYNAMIC POSITION BIAS (DPB)

Figure 1 gives an example of computing $(\Delta x_{ij}, \Delta y_{ij})$ with $G = 5$ in the DPB module. For a group of size $G \times G$, it is easy to deduce that:

$$\begin{aligned} 0 &\leq x, y < G \\ 1 - G &\leq \Delta x_{ij} \leq G - 1 \\ 1 - G &\leq \Delta y_{ij} \leq G - 1. \end{aligned} \quad (1)$$

Thus, motivated by the relative position bias, we construct a matrix $\hat{B} \in \mathbb{R}^{(2G-1) \times (2G-1)}$, where

$$\hat{B}_{i,j} = DPB(1 - G + i, 1 - G + j), \quad 0 \leq i, j < 2G - 1. \quad (2)$$

The complexity of computing \hat{B} is $O(G^2)$. Then, the bias matrix B in DPB can be drawn from \hat{B} , i.e.,

$$B_{i,j} = \hat{B}_{\Delta x_{ij}, \Delta y_{ij}}. \quad (3)$$

When the image/group size (i.e., G) is fixed, both \hat{B} and B will be also unchanged in the test phase. Therefore, we only need to compute \hat{B} and B once, and DPB is equivalent to relative position bias in this case.

APPENDIX B VARIANTS OF CROSSFORMER FOR DETECTION AND SEGMENTATION

We test two different backbones for dense prediction tasks. The variants of CrossFormer for dense prediction (object detection, instance segmentation, and semantic segmentation) are in Table 1. The architectures are the same as those for image classification except that different G and I in the first two stages are used. Notably, group size (i.e., G and I) does not affect the shape of weight tensors, so backbones pre-trained on ImageNet can be fine-tuned directly on other tasks even if they use different G and I .

We test two different backbones for dense prediction tasks. The variants of CrossFormer for dense prediction (object detection, instance segmentation, and semantic segmentation) are in Table 1. The architectures are the same as those for image classification except that different G and I in the first two stages are used. Notably, group size (i.e., G and I) does not affect the shape of weight tensors, so backbones pre-trained on ImageNet can be fine-tuned directly on other tasks even if they use different G and I .

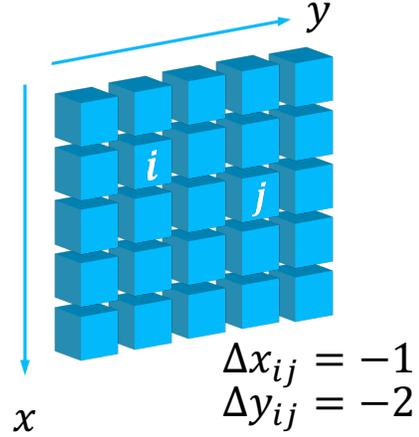


Fig. 1: An example of computing $(\Delta x_{ij}, \Delta y_{ij})$ for DPB.

APPENDIX C EXPERIMENTS

C.1 Object Detection

Table 2 provides more results on object detection with RetinaNet and Mask-RCNN as detection heads. As we can see, a smaller (G, I) achieves a higher AP than a larger one, but the performance gain is marginal. Considering that a larger (G, I) can save more memory cost, we think $(G_1 = 14, I_1 = 16, G_2 = 14, I_2 = 8)$, which accords with configurations in Table 1, achieves a better trade-off between the performance and cost.

C.2 Semantic Segmentation

Similar to object detection, we test two different configurations of (G, I) for semantic segmentation’s backbones. The results are shown in Table 3. As we can see, the memory costs of the two configurations are almost the same, which is different from experiments on object detection. Further, when taking semantic FPN as the detection head, CrossFormers[‡] show advantages over CrossFormers on both IOU (e.g., 46.4 vs. 46.0) and FLOPs (e.g., 209.8G vs. 220.7G). When taking UPerNet as the segmentation head, a smaller (G, I) achieves higher performance like object detection.

TABLE 1: CrossFormer-based backbones for object detection and semantic/instance segmentation. The example input size is 1280×800 . D and H mean embedding dimension and the number of heads in the multi-head self-attention module, respectively. G and I are group size and interval for SDA and LDA, respectively.

	Output Size	Layer Name	CrossFormer-T	CrossFormer-S	CrossFormer-B	CrossFormer-L
		Cross Embed.	Kernel size: $4 \times 4, 8 \times 8, 16 \times 16, 32 \times 32$, Stride=4			
Stage-1	320×200	SDA/LDA MLP	$\begin{bmatrix} D_1 = 64 \\ H_1 = 2 \\ G_1 = 14 \\ I_1 = 16 \end{bmatrix} \times 1$	$\begin{bmatrix} D_1 = 96 \\ H_1 = 3 \\ G_1 = 14 \\ I_1 = 16 \end{bmatrix} \times 2$	$\begin{bmatrix} D_1 = 96 \\ H_1 = 3 \\ G_1 = 14 \\ I_1 = 16 \end{bmatrix} \times 2$	$\begin{bmatrix} D_1 = 128 \\ H_1 = 4 \\ G_1 = 14 \\ I_1 = 16 \end{bmatrix} \times 2$
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2			
Stage-2	160×100	SDA/LDA MLP	$\begin{bmatrix} D_2 = 128 \\ H_2 = 4 \\ G_2 = 14 \\ I_2 = 8 \end{bmatrix} \times 1$	$\begin{bmatrix} D_2 = 192 \\ H_2 = 6 \\ G_2 = 14 \\ I_2 = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} D_2 = 192 \\ H_2 = 6 \\ G_2 = 14 \\ I_2 = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} D_2 = 256 \\ H_2 = 8 \\ G_2 = 14 \\ I_2 = 8 \end{bmatrix} \times 2$
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2			
Stage-3	80×50	SDA/LDA MLP	$\begin{bmatrix} D_3 = 256 \\ H_3 = 8 \\ G_3 = 7 \\ I_3 = 2 \end{bmatrix} \times 8$	$\begin{bmatrix} D_3 = 384 \\ H_3 = 12 \\ G_3 = 7 \\ I_3 = 2 \end{bmatrix} \times 6$	$\begin{bmatrix} D_3 = 384 \\ H_3 = 12 \\ G_3 = 7 \\ I_3 = 2 \end{bmatrix} \times 18$	$\begin{bmatrix} D_3 = 512 \\ H_3 = 16 \\ G_3 = 7 \\ I_3 = 2 \end{bmatrix} \times 18$
		Cross Embed.	Kernel size: $2 \times 2, 4 \times 4$, Stride=2			
Stage-4	40×25	SDA/LDA MLP	$\begin{bmatrix} D_4 = 512 \\ H_4 = 16 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 6$	$\begin{bmatrix} D_4 = 768 \\ H_4 = 24 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	$\begin{bmatrix} D_4 = 768 \\ H_4 = 24 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$	$\begin{bmatrix} D_4 = 1024 \\ H_4 = 32 \\ G_4 = 7 \\ I_4 = 1 \end{bmatrix} \times 2$

TABLE 2: Object detection results on COCO *val* 2017. “Memory” means the allocated memory per GPU reported by *torch.cuda.max_memory_allocated()*. \ddagger indicates that models use different (G, I) from classification models.

Method	Backbone	G_1	I_1	G_2	I_2	Memory	#Params	FLOPs	AP ^b	AP ^b ₅₀	AP ^b ₇₅
RetinaNet 1× schedule	CrossFormer-S	7	8	7	4	14.7G	40.8M	282.0G	44.4	65.8	47.4
	CrossFormer-S \ddagger	14	16	14	8	11.9G	40.8M	272.1G	44.2	65.7	47.2
	CrossFormer-B	7	8	7	4	22.8G	62.1M	389.0G	46.2	67.8	49.5
	CrossFormer-B \ddagger	14	16	14	8	20.2G	62.1M	379.0G	46.1	67.7	49.0
Mask-RCNN 1× schedule	CrossFormer-S	7	8	7	4	15.5G	50.2M	301.0G	45.4	68.0	49.7
	CrossFormer-S \ddagger	14	16	14	8	12.7G	50.2M	291.1G	45.0	67.9	49.1
	CrossFormer-B	7	8	7	4	23.8G	71.5M	407.9G	47.2	69.9	51.8
	CrossFormer-B \ddagger	14	16	14	8	21.0G	71.5M	398.1G	47.1	69.9	52.0

TABLE 3: Semantic segmentation results on ADE20K validation set with semantic FPN or UPerNet as heads.

Backbone	G_1	I_1	G_2	I_2	Semantic FPN (80K iterations)				UPerNet (160K iterations)				
					Memory	#Params	FLOPs	IOU	Memory	#Params	FLOP	IOU	MS IOU
CrossFormer-S	7	8	7	4	20.9G	34.3M	220.7G	46.0	–	62.3M	979.5G	47.6	48.4
CrossFormer-S \ddagger	14	16	14	8	20.9G	34.3M	209.8G	46.4	14.6G	62.3M	968.5G	47.4	48.2
CrossFormer-B	7	8	7	4	14.6G	55.6M	331.0G	47.7	15.8G	83.6M	1089.7G	49.7	50.6
CrossFormer-B \ddagger	14	16	14	8	14.6G	55.6M	320.1G	48.0	15.8G	83.6M	1078.8G	49.2	50.1
CrossFormer-L	7	8	7	4	25.3G	95.4M	497.0G	48.7	18.1G	125.5M	1257.8G	50.4	51.4
CrossFormer-L \ddagger	14	16	14	8	25.3G	95.4M	482.7G	49.1	18.1G	125.5M	1243.5G	50.5	51.4