

# Hire-MLP: Vision MLP via Hierarchical Rearrangement

Jianyuan Guo<sup>1,3\*</sup>, Yehui Tang<sup>1,2\*</sup>, Kai Han<sup>1</sup>, Xinghao Chen<sup>1</sup>,  
Han Wu<sup>3</sup>, Chao Xu<sup>2</sup>, Chang Xu<sup>3†</sup>, Yunhe Wang<sup>1†</sup>

<sup>1</sup>Noah’s Ark Lab, Huawei Technologies.

<sup>2</sup>Key Lab of Machine Perception (MOE), Dept. of Machine Intelligence, Peking University.

<sup>3</sup>School of Computer Science, Faculty of Engineering, University of Sydney.  
{jianyuan.guo, kai.han, yunhe.wang}@huawei.com, c.xu@sydney.edu.au.

## Abstract

Previous vision MLPs such as MLP-Mixer and ResMLP accept linearly flattened image patches as input, making them inflexible for different input sizes and hard to capture spatial information. Such approach withholds MLPs from getting comparable performance with their transformer-based counterparts and prevents them from becoming a general backbone for computer vision. This paper presents Hire-MLP, a simple yet competitive vision MLP architecture via **Hierarchical rearrangement**, which contains two levels of rearrangements. Specifically, the inner-region rearrangement is proposed to capture local information inside a spatial region, and the cross-region rearrangement is proposed to enable information communication between different regions and capture global context by circularly shifting all tokens along spatial directions. Extensive experiments demonstrate the effectiveness of Hire-MLP as a versatile backbone for various vision tasks. In particular, Hire-MLP achieves competitive results on image classification, object detection and semantic segmentation tasks, e.g., 83.8% top-1 accuracy on ImageNet, 51.7% box AP and 44.8% mask AP on COCO val2017, and 49.9% mIoU on ADE20K, surpassing previous transformer-based and MLP-based models with better trade-off for accuracy and throughput. Code is available at <https://github.com/ggry/Hire-Wave-MLP.pytorch>.

## 1. Introduction

Attention mechanism based transformers have shown great superiority in the realm of natural language processing in recent years. Several works such as ViT [11] and DeiT [43] propose to transfer the transformers into visual recognition tasks [14], and have achieved awesome results

which are comparable with conventional convolutional neural networks (CNNs). However, the heavy computational burdens caused by the self-attention modules in transformers withhold the models from better trade-off between accuracy and latency. Recently, models composed of only multi-layer perceptrons (MLPs) have become a new trend in vision community [41, 42]. These MLP-based models can achieve comparable results with CNNs while discarding the heavy self-attention module. For example, MLP-Mixer [41] extracts per-location information through MLPs that are applied to every image patch, and captures long-range information through MLPs that are applied across patches.

Although MLP-Mixer can obtain the global receptive field, there are two intractable flaws that prevent the model from becoming a more general backbone for vision tasks: (i) The number of the patches (tokens) will change as the input size changes, which means it cannot be directly fine-tuned at other resolutions that are different from those used in pre-training phase, making MLP-Mixer infeasible to be transferred into downstream vision tasks such as detection and segmentation. (ii) MLP-Mixer rarely explores the local information, which is demonstrated as an useful inductive bias in both CNNs and transformer-based architectures [17, 49]. The above challenges naturally motivate us to explore an efficient MLP-based architecture which can encode both local and global information while being compatible with flexible input resolutions at the same time.

To address the two aforementioned challenges, we propose the Hire-MLP, which innovates the existing MLP-based models by using hierarchical rearrangement operations. Taking the first challenge into account, the sequence of tokens in MLP-Mixer [41] are denoted as  $X \in \mathbb{R}^{HW \times C}$ , where  $HW$  and  $C$  denote the number of tokens and channels, respectively. MLP-Mixer first uses a token-mixing MLP which acts on the columns of  $X$  to map  $\mathbb{R}^{HW} \mapsto \mathbb{R}^{HW}$ , and then uses a channel-mixing MLP which acts on the rows of  $X$  to map  $\mathbb{R}^C \mapsto \mathbb{R}^C$ . The parameters of the token-mixing MLP are configured by the number of tokens

\*Equal contribution.

†Corresponding author.

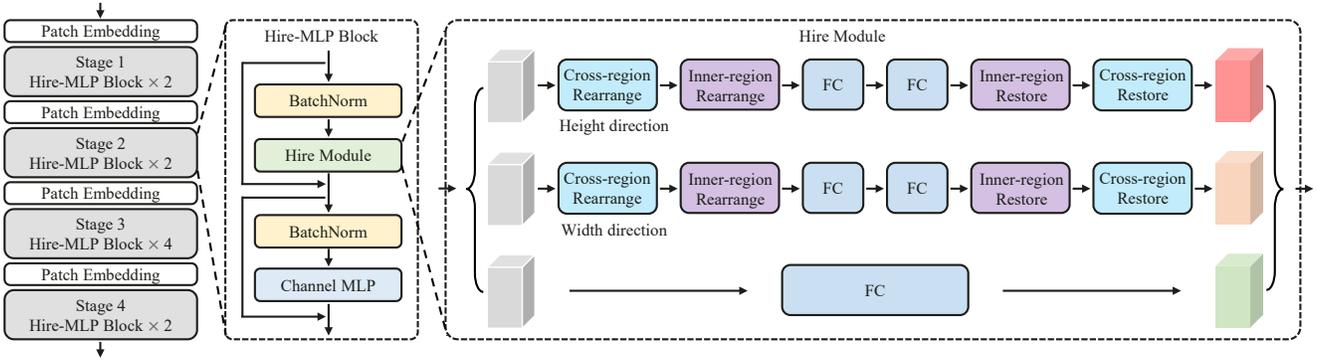


Figure 1. The overall architecture of the proposed Hire-MLP-Tiny. More details and other variants of Hire-MLP can be found in Table A.1 in supplementary materials. Rearrangement layer and restoration layer in hire module are illustrated in Figure 2.

$HW$ , which depends on the resolution of input images and results in the first challenge. To this end, we construct our Hire-MLP merely by channel-mixing MLPs applied on the channel dimension. As for the second challenge, we build the blocks of Hire-MLP based on hierarchical rearrangements and channel-mixing MLPs. The hierarchical rearrangement operation consists of the inner-region rearrangement and the cross-region rearrangement, in which both the local and global information can be easily captured in both height and width directions. We first split the input tokens into multiple regions along the height/width directions, and leverage the inner-region rearrangement operation to shuffle all adjacent tokens belonging to the same region into a one-dimensional vector, followed by two fully connected layers to capture local information within these features. After that, this one-dimensional vector is restored back to the initial arrangement, as illustrated in Figure 1. For the communication between tokens from different regions, a cross-region rearrangement operation is implemented by shifting all the tokens along a specific direction, as shown in Figure 2(c)(d). Such hierarchical rearrangement operation enables our model to obtain both local and global information, and can easily handle the flexible input resolutions.

To be specific, our Hire-MLP has a hierarchical architecture similar to conventional CNNs [17] and recently proposed transformers [32,46] to generate pyramid feature representations for downstream vision tasks. The overall architecture is shown in Figure 1. After the first projection layer, the resulting feature  $X \in \mathbb{R}^{H \times W \times C}$  is then fed into a sequence of Hire-MLP blocks. Hire module is a key component in Hire-MLP block, which consists of three independent branches. The first two branches consist of a cross-region rearrangement layer, an inner-region rearrangement layer, two channel-mixing fully connected (FC) layers, an inner-region restore layer and a cross-region restore layer to capture local and global information along specific direction, *i.e.*, the height and the width direction. The last branch is built upon a simple channel-mixing FC layer to

capture channel information. Compared to existing MLP-based models that spatially shift features in different directions [28,51] or leverage a new cycle fully connected operator [5], our Hire-MLP needs only the channel-mixing MLPs and rearrangement operations. Furthermore, the rearrangement operations can be easily realized by commonly used reshape and padding operations in Pytorch/Tensorflow. And our Hire-MLP is completely capable to serve as a versatile backbone for various computer vision tasks.

Experiments show that Hire-MLP can largely improve the performances of existing MLP-based models on various tasks, including image classification, object detection, instance segmentation, and semantic segmentation. For example, the Hire-MLP-Small attains an 82.1% top-1 accuracy on ImageNet, outperforming Swin-T [32] significantly with a higher throughput. Scaling up the model to larger sizes, we can further obtain 83.2% and 83.8% top-1 accuracy. Using Hire-MLP-Small as backbone, Cascade Mask R-CNN achieves 50.7% box AP and 44.2% mask AP on COCO val2017. In addition, Hire-MLP-Small obtains 46.1% single-scale mIoU on ADE20K, which has an improvement of +1.6% mIoU over Swin-T, demonstrating that Hire-MLP can achieve a better accuracy-latency trade-off than prior MLP-based and transformer-based architectures.

## 2. Related Work

**CNN-based Models.** LeCun *et al.* proposed the classical LeNet [27] in 1990s, which contained most of the basic components of modern CNNs (*e.g.*, convolution and pooling). In ILSVRC 2012 contest, AlexNet [26] achieved far higher performance than others and drew much attention to CNNs. VGGNet [37] constructed a plain model by stacking only convolutional layers with kernel size of  $3 \times 3$ . GoogLeNet [38] designed an inception module containing multiple branches to fuse features from diverse receptive fields. To train an extremely deep model for better performance, ResNet [17, 18] skipped multiple layers with an identity projection to alleviate gradient vanishing or explod-

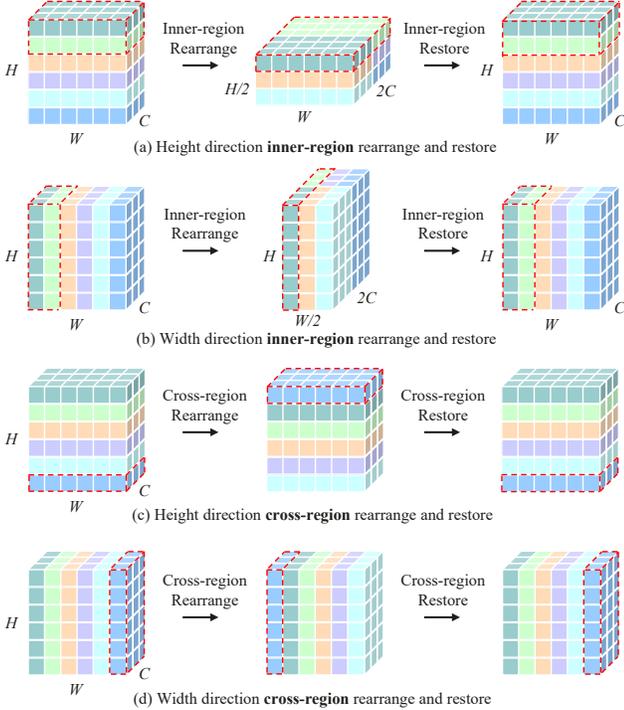


Figure 2. Illustration of the inner-region and cross-region rearrangement operations in hire module.

ing. In addition to accuracy, efficiency also plays a crucial part in the practical implementation of CNN-based models, especially on resource-limited devices such as mobile phones. MobileNet [21] adopted depth-wise convolutions to aggregate spatial information. ShuffleNet [57] introduced the shuffle operation to complement the information loss caused by group convolutions. Such operation can exchange the information across different groups. These well-designed CNNs have been widely used in various tasks such as image recognition [17], object detection [36], semantic segmentation [4] and video analysis [24].

**Transformer-based Models.** The classical transformer model [44] was originally designed to tackle natural language processing (NLP) tasks such as machine translation and English constituency parsing. Recently, Dosovitskiy *et al.* [11] introduced it to vision community by splitting an image into multiple patches and taking each patch as a token in NLP. Vision transformers can accommodate more training data and achieve higher performance compared to CNNs when the dataset is large enough. Touvron *et al.* [43] explored how to train data-efficient vision transformers and proposed a new distillation strategy. Extensive works [6, 12, 15, 32, 46, 47, 52, 53] were proposed to design the architecture of transformers. For example, PVT [46] designed a pyramid-like structure, where the spatial sizes of feature maps are reduced stage-by-stage, and validated the efficiency of transformers on dense prediction

tasks such as object detection and semantic segmentation. TNT [15] embedded small transformer blocks in original modules to capture local information. T2T-ViT [54] improved the tokenization process of input images, and proposed a layer-wise Tokens-to-Token transformation module by recursively aggregating neighboring tokens. The information of images can be preserved more sufficiently compared to the simple tokenization with a single layer. Considering the high computational cost of self-attention mechanism, Swin Transformer [32] calculated the attention between different tokens in shifted local windows, reducing the computational cost from quadratic to linear complexity. However, the self-attention mechanism is still computationally expensive and relatively slow on devices like GPUs.

**MLP-based Models.** Considering the large computational cost of attention modules in transformers, simple and efficient models that consist of only multi-layer perceptrons (MLPs) are proposed [41, 42]. For example, MLP-Mixer [41] used token-mixing MLP and channel-mixing MLP to capture the relationship between tokens and between channels, respectively. Concurrently, the performance of MLP-based models are further improved by designing new architectures [5, 20, 28, 50, 51]. CycleMLP [5] introduced a cycle fully connected layer to capture the spatial information, which replaces token-mixing MLP in [41]. AS-MLP [28] shifted tokens along vertical and horizontal directions to get an axial receptive field.  $S^2$ -MLP [51] also used the shift operation to achieve cross-patch communications. Different from them, our method can simultaneously capture both local and global spatial information by a hierarchical rearrangement operation, *i.e.*, rearranging tokens in/cross local regions, which also achieves a better trade-off between high performance and computational efficiency.

## 3. Method

### 3.1. Hire-MLP Block

The proposed Hire-MLP architecture is constructed by stacking multiple Hire-MLP blocks, as detailed in Figure 1. Similar to ViT [11] and MLP-Mixer [41], each Hire-MLP block consists of two sub-blocks, *i.e.*, the proposed hire module and channel MLP in [41], aggregating spatial information and channel information, respectively. Given the input feature  $X \in \mathbb{R}^{H \times W \times C}$  with height  $H$ , width  $W$ , and channel number  $C$ , a Hire-MLP block can be formulated as:

$$\begin{aligned} Y &= \text{Hire-Module}(\text{BN}(X)) + X, \\ Z &= \text{Channel-MLP}(\text{BN}(Y)) + Y, \end{aligned} \quad (1)$$

where  $Y$  and  $Z$  are intermediate feature and output feature of the block, respectively. BN denotes the batch normalization [23]. The whole Hire-MLP architecture is constructed by iteratively stacking the Hire-MLP block (Eq. 1). Compared with MLP-Mixer [41], the major difference is that

we replace token-mixing MLP in MLP-Mixer with the proposed hire module and have successfully managed to capture the relationship between different tokens effectively.

### 3.2. Hierarchical Rearrangement Module

In MLP-Mixer [41], the token-mixing MLP takes linearly flattened tokens as input, and uses fully connected layers to capture the cross-location information. As the dimension of fully connected layers is fixed, it is not compatible with sequences of variable lengths on dense prediction tasks such as object detection and semantic segmentation. Besides, each token-mixing operation captures and aggregates the global information, while some crucial local information might be neglected. In this section, we propose the **hierarchical rearrangement** (hire) module to replace the token-mixing MLP in [41] and address these challenges accordingly. Briefly, the *inner-region rearrangement* operation in hire module can help capture the local information of tokens in a pre-defined region, while the global information can be captured through *cross-region rearrangement* operation. And credited to the proposed *region partition*, the size of each region remains the same when taking inputs of different sizes. Therefore, our hire module can naturally tackle with sequences of variable lengths and has linear computational complexity with respect to input size. In the following, we will introduce the *region partition*, *inner-region rearrangement*, and *cross-region rearrangement* in details.

**Region Partition.** We first split the input features into multiple regions, and perform the *inner-region rearrangement* on tokens in each region. The feature can be split along both width and height directions. Taking the height direction *inner-region rearrangement* as an example, the input feature  $X$  of shape  $H \times W \times C$  will be divided into  $g$  regions, *i.e.*,  $X = [X_1, X_2, \dots, X_g]$ . Each region  $X_i \in \mathbb{R}^{h \times W \times C}$  contains  $h$  tokens along the height direction, where  $h = H/g$ .

**Inner-region Rearrangement.** Given an input feature  $X_i \in \mathbb{R}^{h \times W \times C}$  of the  $i$ -th region along the height direction, different tokens will exchange the information adequately through *inner-region rearrangement* operation. Specifically, we concatenate all tokens in  $X_i$  along the channel dimension, and get the rearranged feature  $X_i^c$  with the shape of  $W \times hC$  ( $h = 2$  in Figure 2(a)). Then  $X_i^c$  is sent to an MLP module  $\mathcal{F}$  to mix information along the last dimension and produce output feature  $X_i^o \in \mathbb{R}^{W \times hC}$ . For efficiency, the MLP  $\mathcal{F}$  is implemented by two linear projections with bottleneck, *i.e.*, the feature is first reduced to  $W \times \frac{C}{2}$  and then restored to  $W \times hC$ . A non-linear activation function (*e.g.*, ReLU [13] and GeLU [19]) and normalization layer (*e.g.*, BN [23] and LN [1]) can also be inserted into linear projections to enhance representation ability and stabilize training. At last, the output feature  $X_i^o \in \mathbb{R}^{W \times hC}$  is restored to the original shape for next module, *i.e.*, it is split into multiple tokens along the last dimension to get feature

$X'_i \in \mathbb{R}^{h \times W \times C}$ . In this way, different tokens in each region can be mixed adequately for generating output features.

**Cross-region Rearrangement.** Although the *inner-region rearrangement* enables the communication among tokens in a local region, the receptive field of output feature is limited by the size of each region. Here we introduce the *cross-region rearrangement* operation that exchanges information across different regions by shifting tokens along the height/width direction, and in return enables the model to aggregate global spatial information.

The *cross-region rearrangement* is implemented by recurrently shifting all the tokens along a specific direction with a given step size  $s$ , as illustrated in Figure 2(c) ( $s = 1$  along the height direction) and Figure 2(d) ( $s = 1$  along the width direction). After shifting, tokens included in the local region split by *region partition* will change. It is worth noting that this operation can be easily accomplished by the ‘‘circular padding’’ in Pytorch/Tensorflow. To get a global receptive field, the *cross-region rearrangement* operations are inserted before the *inner-region rearrangement* operation every two blocks. The positions of shifted tokens are also restored after the *inner-region restoration* operation to preserve the relative position between different tokens. And this restoration can further boost the accuracy of our Hire-MLP, as shown in Table 5.

Note that Zhang *et al.* [57] uses the channel shuffle operation to communicate across different groups, which disorganizes channels totally. In the contrast, our proposed *cross-region rearrangement* preserves the relative position between different tokens. We argue that the relative position is vital to achieve high representation ability, and related ablation study for these two strategies is investigated in Table 6. We also visualize the feature maps after two *cross-region rearrangement* manners (ShuffleNet [57] manner vs. our shifted manner) in supplementary materials.

**Hire Module.** Considering an input feature  $X$  of size  $H \times W \times C$ , the spatial information communication is conducted within two branches, *i.e.*, along the height direction and the width direction. Inspired by shortcut connections in ResNet [17] and ViP [20], an extra branch without spatial communication is also added, where only a fully connected layer is leveraged to encode information along the channel dimension. The input  $X$  is sent to above three branches to get features  $X'_W$ ,  $X'_H$ , and  $X'_C$ , respectively. Then the output feature  $X'$  is obtained by summing up these features, *i.e.*,  $X' = X'_W + X'_H + X'_C$ , as depicted in Figure 1.

**Complexity Analysis.** In hire module, the fully connected layer (FC) consumes the major memory and computational cost. Consider the height direction branch in Figure 1, given an input feature  $X \in \mathbb{R}^{H \times W \times C}$ , we first split it into  $H/h$  regions with the shape of  $h \times W \times C$ . And the shape of the feature after *inner-region rearrangement* is  $H/h \times W \times hC$ .

We empirically set the channel dimension in bottleneck to  $C/2$ , thereby this branch occupies  $hC \times \frac{C}{2} \times 2 = hC^2$  parameters and  $\frac{H}{h} \times W \times hC \times \frac{C}{2} \times 2 = HWC^2$  FLOPs. Considering a hire module with three branches (height, width, and channel), the total parameters and FLOPs are  $(2hC^2 + C^2)$  and  $3HWC^2$ , respectively.

### 3.3. Overall Architecture

An overview of the Hire-MLP-Tiny architecture is shown in Figure 1, more details and other variants of Hire-MLP are presented in Table A.1 in supplementary materials. We adopt a pyramid-like architecture for Hire-MLP following the commonly used design of CNNs [17, 37] and vision transformers [32, 46]. It first splits the input image into patches (tokens) by a patch embedding layer [45]. Then two Hire-MLP blocks referred to as ‘‘Stage 1’’ are applied on the tokens above. As the network gets deeper, the number of tokens is reduced by another patch embedding layer and output channels are doubled at the same time. Especially, the whole architecture contains four stages, where the feature resolution reduces from  $\frac{H}{4} \times \frac{W}{4}$  to  $\frac{H}{32} \times \frac{W}{32}$  and the output dimension increases accordingly. The pyramid architecture aggregates the spatial feature for extracting semantic information, which can be applied to image classification, object detection, and semantic segmentation.

We develop diverse variants of Hire-MLP architectures with different memory and computational cost. The ‘‘Base’’ model (Hire-MLP-Base) contains {4, 6, 24, 3} layers for each stage. ‘‘Tiny’’ and ‘‘Small’’ variants have fewer layers to realize efficient implementation, while the ‘‘Large’’ variant has larger representation capacity to achieve higher performance. Detailed configurations can also be found in supplementary materials.

## 4. Experiments

In this section, we investigate the effectiveness of Hire-MLP architectures by conducting experiments on several vision tasks. We first compare the proposed Hire-MLP with previous state-of-the-art models for image classification on ImageNet-1K [10], and then we ablate the important design elements of Hire-MLP. We also present the results of object detection and semantic segmentation on COCO [30] and ADE20K [59], respectively.

### 4.1. Image Classification on ImageNet

**Experimental Settings.** We conduct experiments on the challenging ImageNet-1K [10], which is a image classification benchmark containing 1.28M training images and 50K validation images of 1000 classes. ImageNet-1K is also

<sup>1</sup>AS-MLP [28] reported the throughput under the mixed precision training mode, here we reproduce it and report the throughput under the pure precision training mode for a fair comparison with other methods.

Network	Params	FLOPs	Throughput (image / s)	Top-1
CNN-based				
RegNetY-4GF [34]	39M	4.0G	1156.7	81.0
RegNetY-16GF [34]	84M	16.0G	334.7	82.9
EfficientNet-B4 [40]	19M	4.2G	349.4	82.9
EfficientNet-B5 [40]	30M	9.9G	169.1	83.6
EfficientNet-B6 [40]	43M	19.0G	96.9	84.0
Transformer-based				
DeiT-S [43]	22M	4.6G	940.4	79.8
T2T-ViT <sub>t</sub> -14 [53]	22M	5.2G	-	80.7
Swin-T [32]	29M	4.5G	755.2	81.3
CPVT-S-GAP [7]	22M	4.6G	942.3	81.5
PVT-M [46]	44M	6.7G	528.1	81.2
PVT-L [46]	61M	9.8G	358.8	81.7
T2T-ViT <sub>t</sub> -24 [53]	64M	15.0G	-	82.6
TNT-B [15]	66M	14.1G	-	82.9
Swin-B [32]	88M	15.4G	278.1	83.5
MLP-based				
gMLP-Ti [31]	6M	1.4G	-	72.3
CycleMLP-B1 [5]	15M	2.1G	1038.4 <sup>‡</sup>	78.9
Hire-MLP-Tiny (ours)	18M	2.1G	1561.7	<b>79.7</b>
ResMLP-S12 [42]	15M	3.0G	1415.1	76.6
ViP-Small/7 [20]	25M	-	719.0	81.5
AS-MLP-T* [28]	28M	4.4G	863.6 <sup>‡</sup>	81.3
CycleMLP-B2 [5]	27M	3.9G	640.6 <sup>‡</sup>	81.6
Hire-MLP-Small (ours)	33M	4.2G	807.6	<b>82.1</b>
Mixer-B/16 [41]	59M	12.7G	-	76.4
S <sup>2</sup> -MLP-deep [51]	51M	10.5G	-	80.7
ResMLP-B24 [42]	116M	23.0G	231.3	81.0
ViP-Medium/7 [20]	55M	-	418.0	82.7
CycleMLP-B4 [5]	52M	10.1G	320.8 <sup>‡</sup>	83.0
AS-MLP-S* [28]	50M	8.5G	478.4 <sup>‡</sup>	83.1
Hire-MLP-Base (ours)	58M	8.1G	440.6	<b>83.2</b>
S <sup>2</sup> -MLP-wide [51]	71M	14.0G	-	80.0
CycleMLP-B5 [5]	76M	12.3G	246.9 <sup>‡</sup>	83.2
gMLP-B [31]	73M	15.8G	-	81.6
ViP-Large/7 [20]	88M	-	298.0	83.2
AS-MLP-B* [28]	88M	15.2G	312.4 <sup>‡</sup>	83.3
Hire-MLP-Large (ours)	96M	13.4G	290.1	<b>83.8</b>

Table 1. Experimental results of different networks on ImageNet-1K. Throughput is measured as the number of images that we can process per second on a single V100 GPU following [32, 43]. \* means AS-MLP [28] accelerates the AS operation by CUDA implementation. ‡ means the throughput result is reproduced by us<sup>1</sup>.

utilized to conduct the ablation studies. For fair comparisons with recent works, we adopt the same training and augmentation strategy as those in DeiT [43], *i.e.*, models are trained for 300 epochs using the AdamW [33] optimizer with weight decay 0.05 and the batch size of 1024. We use a linear warmup for early 20 epochs, the initial learning rate is set to 1e-3 and gradually drops to 1e-5. The data augmentation methods include Rand-Augment [9], MixUp [56], CutMix [55], Label Smoothing [39], Random Erasing [58], and DropPath [22]. All models are trained on 8 NVIDIA Tesla

Num. of $h$ and $w$	Top-1 (%)	Num. of $h$ and $w$	Top-1 (%)
(2, 2, 2, 2)	81.62	(2, 2, 3, 3)	81.73
(3, 2, 2, 2)	81.82	(3, 3, 2, 2)	81.78
(3, 3, 3, 2)	81.87	(3, 3, 3, 3)	81.79
(4, 3, 3, 2)	<b>82.07</b>	(4, 3, 3, 3)	81.86
(4, 4, 3, 3)	81.81	(4, 4, 4, 4)	81.72
(5, 4, 3, 3)	81.74	(6, 4, 3, 3)	81.49

Table 2. Ablation study on the number of tokens in each region in *Region Partition*. Given an input feature of size  $H \times W \times C$ , we split it into  $H/h$  ( $W/w$ ) regions along the height (width) direction, and the size of each region is  $h \times W \times C$ . We set  $h=w$  as default for  $224 \times 224$  input resolution. For example, (4, 3, 3, 2) indicates  $h$  and  $w$  are set to 4, 3, 3, and 2 for stage 1, stage 2, stage 3, and stage 4, respectively. The step size  $s$  here is set to (2, 2, 1, 1).

Num. of $s$	Top-1 (%)	Num. of $s$	Top-1 (%)
(0, 0, 0, 0)	81.18	(1, 1, 1, 1)	81.88
(2, 2, 1, 1)	<b>82.07</b>	(2, 2, 2, 2)	81.71

Table 3. Ablation study about the step size of shifted tokens ( $s$ ) in *cross-region rearrangement*. For example, (2, 2, 1, 1) means  $s$  is set to 2, 2, 1, and 1 for stage 1, stage 2, stage 3, and stage 4, respectively. (0, 0, 0, 0) indicates there is no *cross-region rearrangement* in Hire-MLP. The  $h$  and  $w$  here are set to (4, 3, 3, 2).

Padding mode	Top-1 (%)	Padding mode	Top-1 (%)
Zero padding	81.62	Circular padding	<b>82.07</b>
Reflect padding	81.48	Replicated padding	81.60

Table 4. Different padding modes for *inner-region rearrangement*.

Model	Top-1 (%)
Hire-MLP-Small	<b>82.07</b>
w/o cross-region restore	81.70
w/o cross-region rearrange and restore	81.18
w/o inner-region rearrange and restore	80.17
w/o extra FC branch	81.32

Table 5. Impact of different components in hire module.

Model	Top-1 (%)
Shifted manner	<b>82.07</b>
ShuffleNet manner [57]	80.90

Table 6. Different manners for cross-region communication.

V100 GPUs, we report the experimental results with single-crop top-1 accuracy, parameters, FLOPs and throughput.

**Main Results.** We compare the proposed Hire-MLP with previous CNN-based, transformer-based, and MLP-based models on Imagenet as shown in Table 1. The resolution of input image is set to  $224 \times 224$ . For example, our

Num. of FC layer	# Params	# FLOPs	Top-1 (%)
1	49.65M	5.65G	82.15
2	33.11M	4.24G	82.07
3	32.98M	4.23G	81.81
4	33.26M	4.24G	81.85

Table 7. Ablation study about the number of intermediate FC layers in first two branches in hire module.

Hire-MLP-Small achieves 82.1% top-1 accuracy with only 4.2G FLOPs, which is better than all other existing MLP-based models. When compared to recently proposed AS-MLP [28] and CycleMLP [5], our Hire-MLP can obtain better performances (+0.5~0.8) without any complicated shift operations or variants of fully connected layer. Scaling up our model to 8.1G and 13.1G can achieve 83.2% and 83.8% top-1 accuracy, respectively. The superiority of Hire-MLP demonstrates that the proposed hire module can better capture both local and global information, which is crucial for classification. In addition, we show the comparison with conventional CNN-based and transformer-based models. When compared to transformer-based models such as DeiT [43], Swin Transformer [32], and PVT [15], our model can get better results with a faster inference speed. When compared to CNN-based architectures such as RegNetY [34], our Hire-MLP can achieve better results with smaller model size and lower computational cost. However, there is still a small gap between our model and the state-of-the-art EfficientNet-B6. We argue that MLP-based architectures have their unique advantages of simplicity and faster inference speed (290.1 vs. 96.9), and there are still opportunities for further enhancements for MLP-based models.

## 4.2. Ablation Study

The core component in Hire-MLP is the hierarchical rearrangement module (Sec. 3.2). We conduct the ablation studies about the number of tokens in each region in region partition, the number of shifted regions and different rearrangement manners for cross-region rearrangement, the padding mode in inner-region rearrangement, and the number of FC layers in hire module. All ablation experiments are conducted based on the Hire-MLP-Small.

### The number of tokens in each region in region partition.

Table 2 investigates how region partition affects the final performance based on Hire-MLP-Small, where  $h$  and  $w$  denote the size of each region. Consider that the resolution of input image is  $224 \times 224$  in ImageNet, we set  $h = w$  if not specified. A small region size implies few adjacent tokens are mixed via the inner-region rearrangement operation, which emphasizes more on local information. We empirically find that a larger region size is required in lower layer to tackle the feature maps with more tokens and obtain

Backbone	RetinaNet 1×					Mask R-CNN 1×						
	Param / FLOPs	AP	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>	Param / FLOPs	AP <sup>b</sup>	AP <sup>b</sup> <sub>50</sub>	AP <sup>b</sup> <sub>75</sub>	AP <sup>m</sup>	AP <sup>m</sup> <sub>50</sub>	AP <sup>m</sup> <sub>75</sub>
ResNet18 [17]	21.3M / 188.7G	31.8	16.3	34.3	43.2	31.2M / 207.3G	34.0	54.0	36.7	31.2	51.0	32.7
PVT-Tiny [46]	23.0M / 189.5G	36.7	22.6	38.8	50.0	32.9M / 208.1G	36.7	59.2	39.3	35.1	56.7	37.3
CycleMLP-B1 [5]	24.9M / 195.0G	38.6	21.9	41.8	<b>50.7</b>	34.8M / 213.6G	39.4	61.4	43.0	36.8	58.6	39.1
Hire-MLP-Tiny	27.8M / 195.3G	<b>38.9</b>	<b>24.9</b>	<b>42.7</b>	<b>50.7</b>	37.7M / 213.8G	<b>39.6</b>	<b>61.7</b>	<b>43.1</b>	<b>37.0</b>	<b>59.1</b>	<b>39.6</b>
ResNet50 [17]	37.7M / 239.3G	36.3	19.3	40.0	48.8	44.2M / 260.1G	38.0	58.6	41.4	34.4	55.1	36.7
PVT-Small [46]	34.2M / 226.5G	40.4	25.0	42.9	55.7	44.1M / 245.1G	40.4	62.9	43.8	37.8	60.1	40.3
CycleMLP-B2 [5]	36.5M / 230.9G	40.9	23.4	44.7	53.4	46.5M / 249.5G	41.7	63.6	45.8	38.2	60.4	41.0
Swin-T [32]	38.5M / 244.8G	41.5	25.1	44.9	<b>55.5</b>	47.8M / 264.0G	42.2	64.6	46.2	39.1	61.6	42.0
Hire-MLP-Small	42.8M / 237.6G	<b>41.7</b>	<b>25.3</b>	<b>45.4</b>	54.6	52.7M / 256.2G	<b>42.8</b>	<b>65.0</b>	<b>46.7</b>	<b>39.3</b>	<b>62.0</b>	<b>42.1</b>
ResNet101 [17]	56.7M / 315.4G	38.5	21.4	42.6	51.1	63.2M / 336.4G	40.4	61.1	44.2	36.4	57.7	38.8
PVT-Medium [46]	53.9M / 283.1G	41.9	25.0	44.9	57.6	63.9M / 301.7G	42.0	64.4	45.6	39.0	61.6	42.1
CycleMLP-B3 [5]	48.1M / 291.3G	42.5	25.2	45.5	56.2	58.0M / 309.9G	43.4	65.0	47.7	39.5	62.0	42.4
CycleMLP-B4 [5]	61.5M / 356.6G	43.2	26.6	46.5	57.4	71.5M / 375.2G	44.1	65.7	48.1	40.2	62.7	43.5
Swin-S [32]	59.8M / 334.8G	<b>44.5</b>	27.4	48.0	<b>59.9</b>	69.1M / 353.8G	44.8	66.6	48.9	40.9	63.4	<b>44.2</b>
Hire-MLP-Base	68.0M / 316.5G	44.3	<b>28.0</b>	<b>48.4</b>	58.0	77.8M / 334.9G	<b>45.2</b>	<b>66.9</b>	<b>49.3</b>	<b>41.0</b>	<b>64.0</b>	<b>44.2</b>
PVT-Large [46]	71.1M / 345.7G	42.6	25.8	46.0	58.4	81.0M / 364.3G	42.9	65.0	46.6	39.5	61.9	42.5
CycleMLP-B5 [5]	85.9M / 402.2G	42.7	24.1	46.3	57.4	95.3M / 421.1G	44.1	65.5	48.4	40.1	62.8	43.0
Hire-MLP-Large	105.8M / 424.5G	<b>44.9</b>	<b>28.9</b>	<b>48.9</b>	<b>57.5</b>	115.2M / 443.5G	<b>45.9</b>	<b>67.2</b>	<b>50.4</b>	<b>41.7</b>	<b>64.7</b>	<b>45.3</b>

Table 8. Object detection and instance segmentation results on COCO val2017. We compare Hire-MLP with other backbones based on RetinaNet and Mask R-CNN frameworks, all models are trained in “1x” schedule. FLOPs is calculated on 1280×800 input.

Backbone	Mask R-CNN 3×							Cascade Mask R-CNN 3×						
	FLOPs	AP <sup>b</sup>	AP <sup>b</sup> <sub>50</sub>	AP <sup>b</sup> <sub>75</sub>	AP <sup>m</sup>	AP <sup>m</sup> <sub>50</sub>	AP <sup>m</sup> <sub>75</sub>	FLOPs	AP <sup>b</sup>	AP <sup>b</sup> <sub>50</sub>	AP <sup>b</sup> <sub>75</sub>	AP <sup>m</sup>	AP <sup>m</sup> <sub>50</sub>	AP <sup>m</sup> <sub>75</sub>
ResNet50 [17]	260.1G	41.0	61.7	44.9	37.1	58.4	40.1	738.7G	46.3	64.3	50.5	40.1	61.7	43.4
AS-MLP-T [28]	260.1G	46.0	67.5	50.7	41.5	64.6	44.5	739.0G	50.1	68.8	54.3	43.5	66.3	46.9
Swin-T [32]	264.0G	46.0	<b>68.2</b>	50.2	41.6	65.1	44.8	742.4G	50.5	69.3	54.9	43.7	66.6	47.1
Hire-MLP-Small	256.2G	<b>46.2</b>	<b>68.2</b>	<b>50.9</b>	<b>42.0</b>	<b>65.6</b>	<b>45.3</b>	734.6G	<b>50.7</b>	<b>69.4</b>	<b>55.1</b>	<b>44.2</b>	<b>66.9</b>	<b>48.1</b>
Swin-S [32]	353.8G	<b>48.5</b>	<b>70.2</b>	<b>53.5</b>	<b>43.3</b>	<b>67.3</b>	46.6	832.4G	<b>51.8</b>	<b>70.4</b>	<b>56.3</b>	44.7	<b>67.9</b>	<b>48.5</b>
AS-MLP-S [28]	346.0G	47.8	68.9	52.5	42.9	66.4	46.3	823.8G	51.1	69.8	55.6	44.2	67.3	48.1
Hire-MLP-Base	334.9G	48.1	69.6	52.7	43.1	66.8	<b>46.7</b>	813.2G	51.7	70.2	56.1	<b>44.8</b>	67.8	<b>48.5</b>

Table 9. Instance segmentation results on COCO val2017. Mask R-CNN and Cascade Mask R-CNN are trained in “3x” schedule.

larger receptive fields. When the region size is further increased, the performance will drop slightly. We conjecture that there might be some information loss in the bottleneck structure with the increasing region size.

**The step size  $s$  of shifted token in cross-region rearrangement.** The cross-region rearrangement is implemented by shifting tokens with a given step size  $s$ , whose impact is investigated in Table 3. When the tokens are not shifted, *i.e.*,  $s = (0, 0, 0, 0)$ , there is no communication between different regions (without cross-region rearrangement operation). Obviously, the lack of global information leads to a bad performance.

**The impacts of different padding methods.** The resolution of the input image from ImageNet [10] is of size  $224 \times 224$ , therefore the shape of output feature in stage 4 is  $7 \times 7$ , which is not divisible by any  $h$  and  $w$ . In consequence, we need to pad the feature map. Table 4 evaluates the influence of different padding methods. And we find

that the “Circular padding” is the most suitable for the design of hire module.

**The impacts of different components in hire module.** Table 5 ablates the impacts of different components in hire module (Sec. 3.2). We can find that the inner-region rearrangement is the most important component to capture local information. The cross-region restoration operation can bring about 0.3% improvement on top-1 accuracy. If we discard the cross-region rearrangement (including the restoration), the model cannot exchange information across different regions, and the performance would drop to 81.18%. And removing the third branch in Figure 1 would harm the top-1 accuracy by 0.7%.

**Different strategies for cross-region communication.** We compare two different strategies for cross-region communication in Table 6. The shifted manner achieves better result compared to ShuffleNet manner, indicating shifted manner can preserve more relative position information for model.

Semantic FPN					UperNet					
Backbone	Param	FLOPs	FPS	SS mIoU	Backbone	Param	FLOPs	FPS	SS mIoU	MS mIoU
PVT-Small [46]	28M	163G	43.9 <sup>‡</sup>	39.8	Swin-T [32]	60M	945G	18.5	44.5	46.1
CycleMLP-B2 [5]	31M	167G	44.5 <sup>‡</sup>	42.4	AS-MLP-T [28]	60M	937G	17.7 <sup>‡</sup>	-	46.5
Hire-MLP-Small	37M	174G	47.3	<b>44.3</b>	Hire-MLP-Small	63M	930G	19.3	<b>46.1</b>	<b>47.1</b>
CycleMLP-B3 [5]	42M	229G	31.0 <sup>‡</sup>	44.5	ResNet-101 [17]	86M	1029G	20.1	43.8	44.9
GFNet-Base [35]	75M	261G	-	44.8	Swin-S [32]	81M	1038G	15.2	47.6	49.5
CycleMLP-B4 [5]	56M	296G	23.6 <sup>‡</sup>	45.1	AS-MLP-S [28]	81M	1024G	14.4 <sup>‡</sup>	-	49.2
Hire-MLP-Base	62M	255G	31.8	<b>46.2</b>	Hire-MLP-Base	88M	1011G	16.0	<b>48.3</b>	<b>49.6</b>
Swin-B [5]	53M	274G	23.4 <sup>‡</sup>	45.2 <sup>†</sup>	Swin-B [32]	121M	1188G	13.3 <sup>‡</sup>	48.1	49.7
CycleMLP-B5 [5]	79M	343G	22.9 <sup>‡</sup>	45.6	AS-MLP-B [28]	121M	1166G	11.0 <sup>‡</sup>	-	49.5
Hire-MLP-Large	99M	366G	24.5	<b>46.6</b>	Hire-MLP-Large	127M	1125G	13.7	<b>48.8</b>	<b>49.9</b>

Table 10. Results of semantic segmentation on ADE20K validation set. FLOPs is calculated with the input size of 2048×512. FPS is measured by using a 32G Tesla V100 GPU. † indicates the results are from GFNet [35]. ‡ indicates the results are measured by us.

More details and corresponding visualization of these two strategies can be found in supplementary materials.

**The number of FC layers in hire module.** The bottleneck design of MLP  $\mathcal{F}$  in hire module (Sec. 3.2) can help eliminate the heavy burden of FLOPs brought by the increase in channels. Ablation studies about the number of FC layers are reported in Table 7. Although using one FC layer achieves the best performance, the parameter and FLOPs are larger than other counterparts. A bottleneck with two FC layers can obtain a better trade-off between accuracy and computational cost. Furthermore, adding more FC layers cannot bring more benefits, demonstrating that the improvements come from our hierarchical rearrangement operation rather than the increase in the number of FC layers.

### 4.3. Object Detection on COCO

**Experimental Settings.** We conduct the object detection and instance segmentation experiments on COCO 2017 benchmark [30], which contains 118K training images and 5K validation images. Following PVT [46] and Swin Transformer [32], we consider three typical object detection frameworks: RetinaNet [29], Mask R-CNN [16] and Cascade Mask R-CNN [2] in mmdetection [3]. We utilize the single-scale training and multi-scale training for the “1x” and “3x” schedules, respectively. More details are introduced in the supplementary materials.

**Results.** We report the results of object detection and instance segmentation under different frameworks and training schedules in Table 8 and Table 9, respectively. As shown in Table 8, Hire-MLP based RetinaNet and Mask R-CNN consistently surpasses the CNN-based ResNet [17], transformer-based PVT [46] and MLP-based CycleMLP [5] under similar FLOPs constraints. Consider RetinaNet [29] as the basic framework, our Hire-MLPs bring consistent +5.8~7.1 AP gains over ResNets [17] and bring +0.3~2.2 AP gains over CycleMLPs [5] with slightly larger model

size and FLOPs. The results indicate that Hire-MLP can serve as an excellent backbone for object detection. Furthermore, Hire-MLP based Cascade Mask R-CNN surpasses the AS-MLP counterpart by 0.6~0.7 in both box AP and mask AP with less FLOPs, as shown in Table 9.

### 4.4. Semantic Segmentation on ADE20K

**Experimental Settings.** We conduct semantic segmentation experiments on ADE20K benchmark [59], which contains 20,210 training images and 2,000 validation images. Following [5, 32, 46], we consider two typical frameworks: Semantic FPN [25] and UperNet [48] in mmsegmentation [8]. See supplementary materials for more details.

**Results.** Table 10 lists the parameters, FLOPs, FPS, single-scale (SS) and multi-scale (MS) mIoU for different backbones based on two typical frameworks. We first choose Semantic FPN [25] as the basic framework following [5, 46]. It can be seen that Hire-MLP outperforms CycleMLP [5] and PVT [46] by a large margin (44.3 vs. 42.4) with similar FLOPs and higher FPS, indicating the superiority of hierarchical rearrangement operation to model at various input scales. In addition, we follow [28, 32] to validate our Hire-MLP based on another commonly used framework UperNet [48]. The proposed Hire-MLP achieves better MS mIoU compared to the state-of-the-art Swin Transformer [32], and is +1.6 mIoU higher than Swin-T on SS mIoU. It seems that Swin Transformer can obtain a larger improvement during multi-scale testing. We speculate one main reason is that the self-attention mechanism in Swin Transformer can capture scale information easier than our hire module. The related ablation studies can be found in supplementary materials.

## 5. Conclusion

This paper proposes a novel variant of MLP-based architecture via hierarchically rearranging tokens to aggregate

both local and global spatial information. Input features are first split into multiple regions along the height/width directions. Different tokens in each region can communicate adequately via inner-region rearrangement operation, which mixes channels from different tokens to extract local information. Then tokens from different regions are rearranged by token shifting. This cross-region rearrangement operation not only exchanges the information between regions, but also preserves the relative position. Based on hierarchical rearrangement operations above, an effective and efficient Hire-MLP is constructed and has achieved significant performance improvements in various vision tasks.

## References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 4
- [2] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018. 8
- [3] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. 8
- [4] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, 2018. 3
- [5] Shoufa Chen, Enze Xie, Chongjian Ge, Ding Liang, and Ping Luo. Cyclemlp: A mlp-like architecture for dense prediction. *arXiv preprint arXiv:2107.10224*, 2021. 2, 3, 5, 6, 7, 8
- [6] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting spatial attention design in vision transformers. *arXiv preprint arXiv:2104.13840*, 2021. 3
- [7] Xiangxiang Chu, Bo Zhang, Zhi Tian, Xiaolin Wei, and Huaxia Xia. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021. 5
- [8] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. <https://github.com/open-mmlab/mmdetection>, 2020. 8
- [9] 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*, 2020. 5
- [10] 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*, 2009. 5, 7
- [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. *arXiv preprint arXiv:2010.11929*, 2020. 1, 3
- [12] Yuxin Fang, Xinggang Wang, Rui Wu, and Wenyu Liu. What makes for hierarchical vision transformer? *arXiv preprint arXiv:2107.02174*, 2021. 3
- [13] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics. JMLR Workshop and Conference Proceedings*, 2011. 4
- [14] Kai Han, Yunhe Wang, Hanqing Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on visual transformer. *arXiv preprint arXiv:2012.12556*, 2020. 1
- [15] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *arXiv preprint arXiv:2103.00112*, 2021. 3, 5, 6
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, 2017. 8
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 2, 3, 4, 5, 7, 8
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645, 2016. 2
- [19] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 4
- [20] Qibin Hou, Zihang Jiang, Li Yuan, Ming-Ming Cheng, Shuicheng Yan, and Jiashi Feng. Vision permutator: A permutable mlp-like architecture for visual recognition. *arXiv preprint arXiv:2106.12368*, 2021. 3, 4, 5
- [21] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 3
- [22] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *European conference on computer vision*, 2016. 5
- [23] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, 2015. 3, 4
- [24] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014. 3
- [25] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 8

- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, pages 1097–1105, 2012. [2](#)
- [27] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. [2](#)
- [28] Dongze Lian, Zehao Yu, Xing Sun, and Shenghua Gao. As-mlp: An axial shifted mlp architecture for vision. *arXiv preprint arXiv:2107.08391*, 2021. [2](#), [3](#), [5](#), [6](#), [7](#), [8](#)
- [29] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, 2017. [8](#)
- [30] 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. [5](#), [8](#)
- [31] Hanxiao Liu, Zihang Dai, David R So, and Quoc V Le. Pay attention to mlps. *arXiv preprint arXiv:2105.08050*, 2021. [5](#)
- [32] 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. *arXiv preprint arXiv:2103.14030*, 2021. [2](#), [3](#), [5](#), [6](#), [7](#), [8](#)
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. [5](#)
- [34] Ilija Radosavovic, Raj Prateek Kosaraju, Ross Girshick, Kaiming He, and Piotr Dollár. Designing network design spaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. [5](#), [6](#)
- [35] Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. Global filter networks for image classification. *arXiv preprint arXiv:2107.00645*, 2021. [8](#)
- [36] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015. [3](#)
- [37] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [2](#), [5](#)
- [38] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015. [2](#)
- [39] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016. [5](#)
- [40] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, 2019. [5](#)
- [41] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, et al. Mlp-mixer: An all-mlp architecture for vision. *arXiv preprint arXiv:2105.01601*, 2021. [1](#), [3](#), [4](#), [5](#)
- [42] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby, Edouard Grave, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, and Herve Jegou. Resmlp: Feedforward networks for image classification with data-efficient training. *arXiv preprint arXiv:2105.03404*, 2021. [1](#), [3](#), [5](#)
- [43] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020. [1](#), [3](#), [5](#), [6](#)
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. [3](#)
- [45] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvtv2: Improved baselines with pyramid vision transformer. *arXiv preprint arXiv:2106.13797*, 2021. [5](#)
- [46] 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. *arXiv preprint arXiv:2102.12122*, 2021. [2](#), [3](#), [5](#), [7](#), [8](#)
- [47] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808*, 2021. [3](#)
- [48] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. [8](#)
- [49] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions help transformers see better. *arXiv preprint arXiv:2106.14881*, 2021. [1](#)
- [50] Tan Yu, Xu Li, Yunfeng Cai, Mingming Sun, and Ping Li. S<sup>2</sup>-mlpv2: Improved spatial-shift mlp architecture for vision. *arXiv preprint arXiv:2108.01072*, 2021. [3](#)
- [51] Tan Yu, Xu Li, Yunfeng Cai, Mingming Sun, and Ping Li. S2-mlp: Spatial-shift mlp architecture for vision. *arXiv preprint arXiv:2106.07477*, 2021. [2](#), [3](#), [5](#)
- [52] Kun Yuan, Shaopeng Guo, Ziwei Liu, Aojun Zhou, Fengwei Yu, and Wei Wu. Incorporating convolution designs into visual transformers. *arXiv preprint arXiv:2103.11816*, 2021. [3](#)
- [53] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986*, 2021. [3](#), [5](#)
- [54] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986*, 2021. [3](#)

- [55] 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*, 2019. 5
- [56] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 5
- [57] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6848–6856, 2018. 3, 4, 6
- [58] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. 5
- [59] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 2019. 5, 8