

Mobile-Former: Bridging MobileNet and Transformer

Yinpeng Chen¹Xiyang Dai¹Dongdong Chen¹Mengchen Liu¹Xiaoyi Dong²Lu Yuan¹Zicheng Liu¹¹ Microsoft² University of Science and Technology of China

{yiche, xidai, dochen, mengcliu, luyuan, zliu}@microsoft.com,

dlight@mail.ustc.edu.cn

Abstract

We present *Mobile-Former*, a parallel design of *MobileNet* and *transformer* with a two-way bridge in between. This structure leverages the advantages of *MobileNet* at local processing and *transformer* at global interaction. And the bridge enables bidirectional fusion of local and global features. Different from recent works on vision transformer, the *transformer* in *Mobile-Former* contains very few tokens (e.g. 6 or fewer tokens) that are randomly initialized to learn global priors, resulting in low computational cost. Combining with the proposed light-weight cross attention to model the bridge, *Mobile-Former* is not only computationally efficient, but also has more representation power. It outperforms *MobileNetV3* at low FLOP regime from 25M to 500M FLOPs on ImageNet classification. For instance, *Mobile-Former* achieves 77.9% top-1 accuracy at 294M FLOPs, gaining 1.3% over *MobileNetV3* but saving 17% of computations. When transferring to object detection, *Mobile-Former* outperforms *MobileNetV3* by 8.6 AP in *RetinaNet* framework. Furthermore, we build an efficient end-to-end detector by replacing backbone, encoder and decoder in *DETR* with *Mobile-Former*, which outperforms *DETR* by 1.1 AP but saves 52% of computational cost and 36% of parameters.

1. Introduction

Recently, vision transformer (ViT) [9, 32] demonstrates the advantage of global processing and achieves significant performance boost over CNNs. However, when constraining the computational budget within 1G FLOPs, the gain of ViT diminishes. If we further challenge the computational cost, *MobileNet* [15, 16, 26] and its extensions [11, 18] still dominate their backyard (e.g. fewer than 300M FLOPs for ImageNet classification) due to their efficiency in local processing filters via decomposition of depthwise and pointwise convolution. This in turn naturally raises a question:

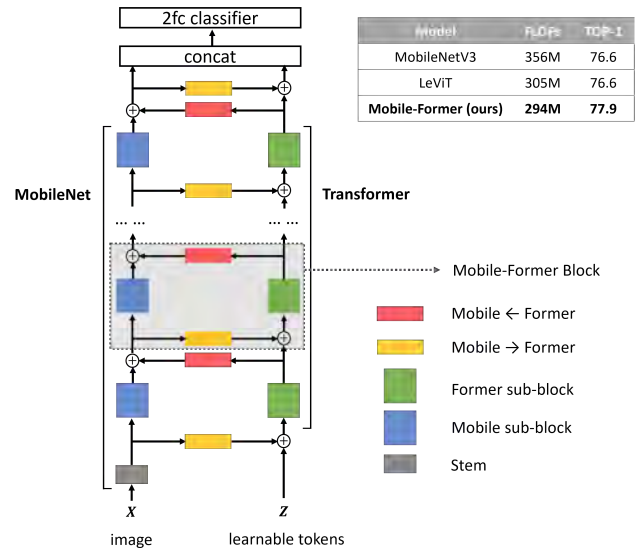


Figure 1. **Overview of Mobile-Former**, which parallelizes *MobileNet* [26] on the left side and *Transformer* [36] on the right side. Different from vision transformer [9] that uses image patches to form tokens, the *transformer* in *Mobile-Former* takes *very few learnable tokens* as input that are randomly initialized. *Mobile* (refers to *MobileNet*) and *Former* (refers to *transformer*) communicate through a bidirectional bridge, which is modeled by the proposed light-weight cross attention. Best viewed in color.

How to design efficient networks to effectively encode both local processing and global interaction?

A straightforward idea is to combine convolution and vision transformer. Recent works [10, 39, 40] show the benefit of combining convolution and vision transformer in *series*, either using convolution at the beginning or intertwining convolution into each transformer block.

In this paper, we shift the design paradigm from *series* to *parallel*, and propose a new network that parallelizes *MobileNet* and *transformer* with a two-way bridge in between (see Figure 1). We name it *Mobile-Former*, where *Mobile* refers to *MobileNet* and *Former* stands for *transformer*. *Mo-*

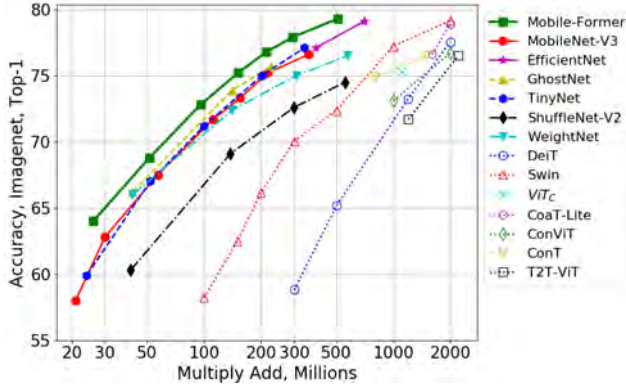


Figure 2. **Comparison among Mobile-Former, efficient CNNs and vision transformers**, in terms of accuracy over FLOPs. The comparison is performed on ImageNet classification. Mobile-Former consistently outperforms both efficient CNNs and vision transformers in low FLOP regime (from 25M to 500M MAAdd). Note that we implement Swin [22] and DeiT [32] at low computational budget from 100M to 2G FLOPs. Best viewed in color.

mobile takes an image as input and stacks mobile (or inverted bottleneck) blocks [26]. It leverages the efficient depthwise and pointwise convolution to extract local features. *Former* takes a few learnable tokens as input and stacks multi-head attention and feed-forward networks (FFN). These tokens are used to encode global features of the image.

Mobile and *Former* communicate through a two-way bridge to fuse local and global features. This is crucial since it feeds local features to *Former*'s tokens as well as introduces global views to every pixel of feature map in *Mobile*. We propose a light-weight cross attention to model this bidirectional bridge by (a) performing the cross attention at the bottleneck of *Mobile* where the number of channels is low, and (b) removing projections on query, key and value (W^Q , W^K , W^V) from *Mobile* side.

This parallel structure with a bidirectional bridge leverages the advantages of both MobileNet and transformer. Decoupling of local and global features in parallel leverages MobileNet's efficiency in extracting local features as well as transformer's power in modeling global interaction. More importantly, this is achieved in an efficient way via a thin transformer with very few tokens and a light-weight bridge to exchange local and global features between *Mobile* and *Former*. The bridge and *Former* consume less than 20% of the total computational cost, but significantly improve the representation capability. This showcases an efficient and effective implementation of part-whole hierarchy [14].

Mobile-Former achieves solid performance on both image classification and object detection. For example, it achieves 77.9% top-1 accuracy on ImageNet classification at 294M FLOPs, outperforming MobileNetV3 [15] and LeViT [10] by a clear margin (see Figure 1). More importantly, Mobile-Former consistently outperforms both ef-

ficient CNNs and vision transformers from 25M to 500M FLOPs (see Figure 2), showcasing the usage of transformer at the low FLOP regime where efficient CNNs dominate.

When transferring from image classification to object detection, Mobile-Former significantly outperforms MobileNetV3 as backbone in RetinaNet [20], gaining 8.6 AP (35.8 vs. 27.2 AP) with even less computational cost. In addition, we build an efficient end-to-end detector by using Mobile-Former to replace backbone and encoder/decoder in DETR [1]. Using the same number of object queries (100), it gains 1.1 AP over DETR (43.1 vs. 42.0 AP) but has significantly fewer FLOPs (41G vs. 86G) and smaller model size (26.6M vs. 41.3M).

Finally we note that *exploring the optimal network parameters (e.g. width, height) in Mobile-Former is not a goal of this work*, rather we demonstrate that the parallel design provides an efficient and effective network architecture.

2. Related Work

Light-weight convolutional neural networks (CNNs):

MobileNets [15, 16, 26] efficiently encode local features by stacking depthwise and pointwise convolutions. ShuffleNet [25, 46] uses group convolution and channel shuffle to simplify pointwise convolution. MicroNet [18] presents micro-factorized convolution to handle extremely low FLOPs by lowering node connectivity to enlarge network width. Dynamic operators [3, 4, 17, 43] have been studied to boost performance for MobileNet with negligible computational cost. Other efficient operators include butterfly transform [34], cheap linear transformations in GhostNet [11], and using additions to trade multiplications in AdderNet [2]. In addition, MixConv [29] explores mixing up multiple kernel sizes, and Sandglass [48] flips the structure of inverted residual block. EfficientNet [28, 30] and TinyNet [12] study the compound scaling of depth, width and resolution.

Vision transformers (ViTs):

Recently, ViT [9] and its follow-ups [8, 22, 32, 35, 44] achieve impressive performance on multiple vision tasks. The original ViT requires training on large dataset such as JFT-300M. DeiT [32] introduces training strategies on the smaller ImageNet-1K dataset. Later, hierarchical transformers are proposed to handle high resolution images. Swin [22] computes self-attention within shifted local windows and CSWin [8] further improves it by introducing cross-shaped window. T2T-ViT [44] progressively converts an image to tokens by recursively aggregating neighboring tokens. HaloNet [35] improves speed, memory usage and accuracy by two extensions (blocked local attention and attention downsampling).

Combination of CNN and ViT:

Recent works [6, 10, 27, 39, 40] show advantages of combining convolution and transformer. BoTNet [27] improves both instance segmentation and object detection by just using self-attention in the last

three blocks of ResNet [13]. ConViT [6] presents a gated positional self-attention for soft convolutional inductive biases. CvT [39] introduces depthwise/pointwise convolution before multi-head attention. LeViT [10] and ViT_C [40] use convolutional stem to replace the patchify stem and achieve clear improvement. In this paper, we propose a different design that *parallelizes* MobileNet and transformer with bidirectional cross attention in between. Our approach is efficient and effective, outperforming both efficient CNNs and ViT variants at low FLOP regime.

3. Our Method: Mobile-Former

In this section, we first overview the design of Mobile-Former (see Figure 1), and then discuss details within a Mobile-Former block (see Figure 3). Finally, we show the network specification and variants with different FLOPs.

3.1. Overview

Parallel structure: Mobile-Former parallelizes MobileNet and transformer, and connects them by bidirectional cross attention (see Figure 1). *Mobile* (refers to MobileNet) takes an image as input ($\mathbf{X} \in \mathbb{R}^{HW \times 3}$) and applies inverted bottleneck blocks [26] to extract local features. *Former* (refers to transformer) takes learnable parameters (or tokens) as input, denoted as $\mathbf{Z} \in \mathbb{R}^{M \times d}$ where M and d are the number and dimension of tokens, respectively. These tokens are randomly initialized. Different from vision transformer (ViT) [9], where tokens project the local image patch linearly, *Former* has significantly fewer tokens ($M \leq 6$ in this paper), each represents a global prior of the image. This results in much less computational cost.

Low cost two-way bridge: *Mobile* and *Former* communicate through a two-way bridge where local and global features are fused bidirectionally. The two directions are denoted as *Mobile*→*Former* and *Mobile*←*Former*, respectively. We propose a light-weight cross attention to model them, in which the projections (\mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V) are removed from *Mobile* side to save computations, but kept at *Former* side. The cross attention is computed at the bottleneck of *Mobile* where the number of channels is low. Specifically, the light-weight cross attention from local feature map \mathbf{X} to global tokens \mathbf{Z} is computed as:

$$\mathcal{A}_{\mathbf{X} \rightarrow \mathbf{Z}} = \left[\text{Attn}(\tilde{\mathbf{z}}_i \mathbf{W}_i^Q, \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_i) \right]_{i=1:h} \mathbf{W}^O, \quad (1)$$

where the local feature \mathbf{X} and global tokens \mathbf{Z} are split into h heads as $\mathbf{X} = [\tilde{\mathbf{x}}_1 \cdots \tilde{\mathbf{x}}_h]$, $\mathbf{Z} = [\tilde{\mathbf{z}}_1 \cdots \tilde{\mathbf{z}}_h]$ for multi-head attention. The split for the i^{th} head $\tilde{\mathbf{z}}_i \in \mathbb{R}^{M \times \frac{d}{h}}$ is different to the i^{th} token $\mathbf{z}_i \in \mathbb{R}^d$. \mathbf{W}_i^Q is the query projection matrix for the i^{th} head. \mathbf{W}^O is used to combine multiple heads together. $\text{Attn}(Q, K, V)$ is the standard attention function [36] over query Q , key K , and value V

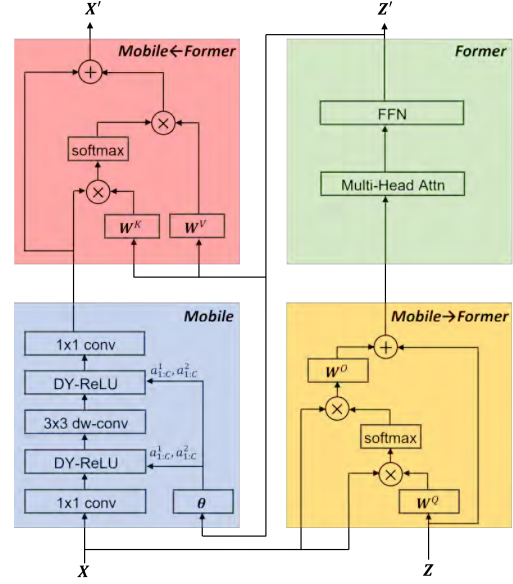


Figure 3. **Mobile-Former block** that includes four modules: *Mobile* sub-block modifies inverted bottleneck block in [26] by replacing ReLU with dynamic ReLU [4]. *Mobile*→*Former* uses light-weight cross attention to fuse local features into global features. *Former* sub-block is a standard transformer block including multi-head attention and FFN. Note that the output of *Former* is used to generate parameters for dynamic ReLU in *Mobile* sub-block. *Mobile*←*Former* bridges from global to local features.

as $\text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$. $[\cdot]_{1:h}$ denotes the concatenation of h elements. Note that the projection matrices for the key and value are removed from *Mobile* side, while the project matrix \mathbf{W}_i^Q for the query is kept at *Former* side. Similarly, the cross attention from global to local is computed as:

$$\mathcal{A}_{\mathbf{Z} \rightarrow \mathbf{X}} = \left[\text{Attn}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{z}}_i \mathbf{W}_i^K, \tilde{\mathbf{z}}_i \mathbf{W}_i^V) \right]_{i=1:h}, \quad (2)$$

where \mathbf{W}_i^K and \mathbf{W}_i^V are the projection matrices for the key and value at *Former* side. The projection matrix of the query is removed from *Mobile* side.

3.2. Mobile-Former Block

Mobile-Former consists of stacked *Mobile-Former blocks* (see Figure 1). Each block has four pillars: a *Mobile* sub-block, a *Former* sub-block, and two-way cross attention *Mobile*←*Former* and *Mobile*→*Former* (shown in Figure 3).

Input and output: Mobile-Former block has two inputs: (a) local feature map $\mathbf{X} \in \mathbb{R}^{HW \times C}$, which has C channels over height H and width W , and (b) global tokens $\mathbf{Z} \in \mathbb{R}^{M \times d}$, where M and d are the number and dimension of tokens, respectively. Note that M and d are identical across all blocks. Mobile-Former block outputs the updated local feature map \mathbf{X}' and global tokens \mathbf{Z}' , which are used as input for the next block.

Mobile sub-block: As shown in Figure 3, *Mobile* sub-block takes the feature map \mathbf{X} as input and its output is taken as the input for *Mobile* \leftarrow *Former*. It is slightly different to the inverted bottleneck block in [26] by replacing ReLU with dynamic ReLU [4] as the activation function. Different from the original dynamic ReLU, in which the parameters are generated by applying two MLP layers on the average pooled feature, we save the average pooling by applying the two MLP layers (θ in Figure 3) on the first global token output z'_1 from *Former*. Note that the kernel size of depthwise convolution is 3×3 for all blocks.

Former sub-block: *Former* sub-block is a standard transformer block including a multi-head attention (MHA) and a feed-forward network (FFN). Expansion ratio 2 (instead of 4) is used in FFN. We follow [36] to use post layer normalization. *Former* is processed between *Mobile* \rightarrow *Former* and *Mobile* \leftarrow *Former* (see Figure 3).

Mobile \rightarrow **Former:** The proposed light-weight cross attention (Equation 1) is used to fuse local features \mathbf{X} to global tokens \mathbf{Z} . Compared to the standard attention, the projection matrices for key \mathbf{W}^K and value \mathbf{W}^V (on the local features \mathbf{X}) are removed to save computations (see Figure 3).

Mobile \leftarrow **Former:** Here, the cross attention (Equation 2) is on the opposite direction to *Mobile* \rightarrow *Former*. It fuses global tokens to local features. The local features are the query and global tokens are the key and value. Therefore, we keep the projection matrices for the key \mathbf{W}^K and value \mathbf{W}^V , but remove the projection matrix for the query \mathbf{W}^Q to save computations, as shown in Figure 3.

Computational complexity: The four pillars of a Mobile-Former block have different computational costs. Given an input feature map with size $HW \times C$, and M global tokens with d dimensions, *Mobile* consumes the most computations $O(HWC^2)$. *Former* and the two-way bridge have light weight, consuming less than 20% of the total computational cost. Specifically, *Former* has complexity $O(M^2d + Md^2)$ for self attention and FFN. *Mobile* \rightarrow *Former* and *Mobile* \leftarrow *Former* share the complexity $O(MHWC + MdC)$ for cross attention.

3.3. Network Specification

Architecture: Table 1 shows a Mobile-Former architecture with 294M FLOPs for image size 224×224 , which stacks 11 Mobile-Former blocks at different input resolutions. All blocks have six global tokens with dimension 192. It starts with a 3×3 convolution as stem and a lite bottleneck block [18] at stage 1, which expands and then squeezes the number of channels by stacking a 3×3 depthwise and a pointwise convolution. Stage 2–5 consists of Mobile-Former blocks. Each stage handles downsampling by a downsample variant denoted as Mobile-Former $^\downarrow$ (see appendix A for details). The classification head applies av-

Stage	Input	Operator	exp size	#out	Stride
tokens	6×192	–	–	–	–
stem	$224^2\times 3$	conv2d, 3×3	–	16	2
1	$112^2\times 16$	bneck-lite	32	16	1
2	$112^2\times 16$	Mobile-Former $^\downarrow$	96	24	2
	$56^2\times 24$	Mobile-Former	96	24	1
3	$56^2\times 24$	Mobile-Former $^\downarrow$	144	48	2
	$28^2\times 48$	Mobile-Former	192	48	1
4	$28^2\times 48$	Mobile-Former $^\downarrow$	288	96	2
	$14^2\times 96$	Mobile-Former	384	96	1
	$14^2\times 96$	Mobile-Former	576	128	1
	$14^2\times 128$	Mobile-Former	768	128	1
5	$14^2\times 128$	Mobile-Former $^\downarrow$	768	192	2
	$7^2\times 192$	Mobile-Former	1152	192	1
	$7^2\times 192$	Mobile-Former	1152	192	1
	$7^2\times 192$	conv2d, 1×1	–	1152	1
head	$7^2\times 1152$	pool, 7×7	–	1152	–
	$1^2\times 1152$	concat w/ cls token	–	1344	–
	$1^2\times 1344$	FC	–	1920	–
	$1^2\times 1920$	FC	–	1000	–

Table 1. **Specification for Mobile-Former-294M.** “bneck-lite” denotes the lite bottleneck block [18]. “Mobile-Former $^\downarrow$ ” denotes the variant of downsample block.

erage pooling on the local features, concatenates with the first global token, and then passes through two fully connected layers with h-swish [15] in between.

Mobile-Former variants: Mobile-Former has seven models with different computational costs from 26M to 508M FLOPs. They share the similar architecture, but have different width and height. We follow [40] to refer our models by their FLOPs, e.g. Mobile-Former-294M, Mobile-Former-96M. The details of network architecture for these models are listed in appendix A (see Table 11).

4. Efficient End-to-End Object Detection

Mobile-Former can be easily used for object detection in both backbone and head, providing an *efficient end-to-end* detector. Using the same number of object queries (100), it outperforms DETR [1], but uses much lower FLOPs.

Backbone–Head architecture: We use Mobile-Former blocks in both backbone and head (see Figure 4), which have separate tokens. The backbone has six global tokens while the head has 100 object queries generated similarly to DETR [1]. Different from DETR that has a single scale ($\frac{1}{32}$ or $\frac{1}{16}$) in the head, Mobile-Former head employs multi-scales ($\frac{1}{32}$, $\frac{1}{16}$, $\frac{1}{8}$) in low FLOPs due to its computational efficiency. The upsampling is achieved via bilinear interpolation followed by adding the feature output from the backbone (with the same resolution). All object queries progressively refine their representation across scales from coarse to fine, saving the manual process in FPN [19] to allocate objects across scales by size. We follow DETR to use prediction FFNs and auxiliary losses in the head during the training. The head is trained from scratch and the backbone

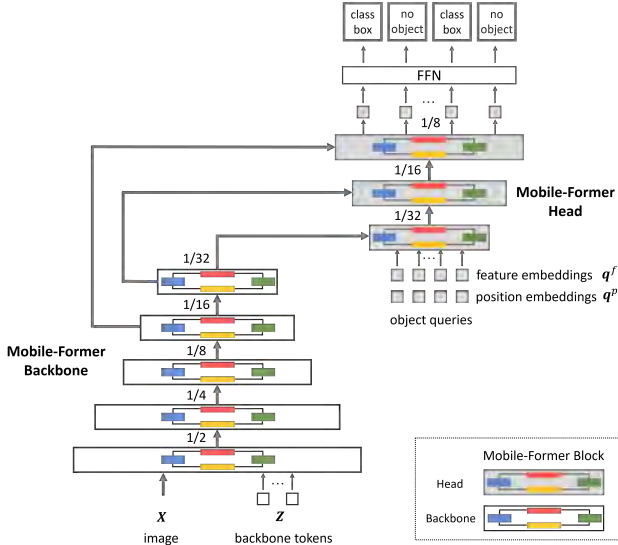


Figure 4. **Mobile-Former for object detection.** Both backbone and head use Mobile-Former blocks (see Figure 1, 3). The backbone has 6 global tokens while the head has 100 object queries. All object queries pass through multiple resolutions ($\frac{1}{32}$, $\frac{1}{16}$, $\frac{1}{8}$) in the head. Similar to DETR [1], feed forward network (FFN) is used to predict class label and bounding box. Best viewed in color.

is pretrained on ImageNet. Our end-to-end Mobile-Former detector is computationally efficient. The total cost of E2E-MF-508M that uses Mobile-Former-508M as backbone and nine Mobile-Former blocks in the head is 41.4G FLOPs, significantly less than DETR (86G FLOPs). But it outperforms DETR by 1.1 AP (43.1 vs. 42.0 AP). The details of head structure is listed in appendix A (see Table 13).

Spatial-aware dynamic ReLU in backbone: We extend dynamic ReLU in the backbone from *spatial-shared* to *spatial-aware* by involving all global tokens to generate parameters, rather than just using the first one, as these tokens have different spatial focuses. Let us denote the parameter generation of the spatial-shared dynamic ReLU as $\theta = f(z_1)$, where z_1 is the first global token and $f(\cdot)$ is modeled by two MLP layers with ReLU in between. By contrast, the spatial-aware dynamic ReLU generates parameters θ_i per spatial position i in a feature map, by using all global tokens $\{z_j\}$ as follows:

$$\theta_i = \sum_j \alpha_{i,j} f(z_j), \text{ s.t. } \sum_j \alpha_{i,j} = 1, \quad (3)$$

where $\alpha_{i,j}$ is the attention between the feature at position i and token z_j . Its calculation is cheap by just normalizing the cross attention obtained in *Mobile*→*Former* along tokens $\{z_j\}$. Spatial-aware dynamic ReLU is on par with its spatial-shared counterpart on image classification, but gains 1.1 AP on COCO object detection (see Table 10).

Adapting position embedding in head: Different from

DETR [1] that shared the position embedding of object queries across all decoder layers, we refine the position embedding after each block in the head as the feature map changes per block. Let us denote the feature and position embedding of a query at the k^{th} block as q_k^f and q_k^p , respectively. The sum of them ($q_k^f + q_k^p$) is used to compute cross attention between object queries and feature map as well as self attention among object queries, after which the feature embedding is updated as the input for the next block q_{k+1}^f . Here, we adapt the position embedding based upon the feature embedding as:

$$q_{k+1}^p = q_k^p + g(q_{k+1}^f), \quad (4)$$

where the adaptation function $g(\cdot)$ is implemented by two MLP layers with ReLU in between. Thus, object queries can adapt their positions across blocks based on the content.

5. Experimental Results

We evaluate the proposed Mobile-Former on both ImageNet classification [7], and COCO object detection [21].

5.1. ImageNet Classification

Image classification experiments are conducted on ImageNet [7] that has 1000 classes, including 1,281,167 images for training and 50,000 images for validation.

Training setup: The image resolution is 224×224 . All models are trained from scratch using AdamW [23] optimizer for 450 epochs with cosine learning rate decay. A batch size of 1024 is used. Data augmentation includes Mixup [45], auto-augmentation [5], and random erasing [47]. Different combinations of initial learning rate, weight decay and dropout are used for models with different complexities, which are listed in appendix A (see Table 12).

Comparison with efficient CNNs: Table 2 shows the comparison between Mobile-Former and classic efficient CNNs: (a) MobileNetV3 [15], (b) EfficientNet [28], and (c) ShuffleNetV2 [25] and its extension WeightNet [24]. The comparison covers the FLOP range from 26M to 508M, organized in seven groups based on similar FLOPs. Mobile-Former consistently outperforms efficient CNNs with even less computational cost except the group around 150M FLOPs, where Mobile-Former costs slightly more FLOPs than ShuffleNet/WeightNet (151M vs. 138M/141M), but achieves significantly higher top-1 accuracy (75.2% vs. 69.1%/72.4%). This demonstrates that our parallel design improves the representation capability efficiently.

Comparison with ViTs: In Table 3, we compare Mobile-Former with multiple variants (DeiT [32], T2T-ViT [44], PVT [37], ConViT [6], CoaT [41], ViT_C [40], Swin [22]) of vision transformer. All variants use image resolution 224×224 and are trained *without* distillation from a teacher

Model	Input	#Params	MAdds	Top-1
MobileNetV3 Small 1.0× [15]	160 ²	2.5M	30M	62.8
Mobile-Former-26M	224 ²	3.2M	26M	64.0
MobileNetV3 Small 1.0× [15]	224 ²	2.5M	57M	67.5
Mobile-Former-52M	224 ²	3.5M	52M	68.7
MobileNetV3 1.0× [15]	160 ²	5.4M	112M	71.7
Mobile-Former-96M	224 ²	4.6M	96M	72.8
ShuffleNetV2 1.0× [25]	224 ²	2.2M	138M	69.1
ShuffleNetV2 1.0×+WeightNet 4× [24]	224 ²	5.1M	141M	72.4
MobileNetV3 0.75× [15]	224 ²	4.0M	155M	73.3
Mobile-Former-151M	224 ²	7.6M	151M	75.2
MobileNetV3 1.0× [15]	224 ²	5.4M	217M	75.2
Mobile-Former-214M	224 ²	9.4M	214M	76.7
ShuffleNetV2 1.5× [25]	224 ²	3.5M	299M	72.6
ShuffleNetV2 1.5×+WeightNet 4× [24]	224 ²	9.6M	307M	75.0
MobileNetV3 1.25× [15]	224 ²	7.5M	356M	76.6
EfficientNet-B0 [28]	224 ²	5.3M	390M	77.1
Mobile-Former-294M	224 ²	11.4M	294M	77.9
ShuffleNetV2 2× [25]	224 ²	5.5M	557M	74.5
ShuffleNetV2 2×+WeightNet 4× [24]	224 ²	18.1M	573M	76.5
Mobile-Former-508M	224 ²	14.0M	508M	79.3

Table 2. **Comparing Mobile-Former with efficient CNNs** evaluated on ImageNet [7] classification.

Model	Input	#Params	MAdds	Top-1
T2T-ViT-7 [44]	224 ²	4.3M	1.2G	71.7
DeiT-Tiny [32]	224 ²	5.7M	1.2G	72.2
ConViT-Tiny [6]	224 ²	6.0M	1.0G	73.1
ConT-Ti [42]	224 ²	5.8M	0.8G	74.9
ViT _C [40]	224 ²	4.6M	1.1G	75.3
ConT-S [42]	224 ²	10.1M	1.5G	76.5
Swin-1G [22] ‡	224 ²	7.3M	1.0G	77.3
Mobile-Former-294M	224 ²	11.4M	294M	77.9
PVT-Tiny [37]	224 ²	13.2M	1.9G	75.1
T2T-ViT-12 [44]	224 ²	6.9M	2.2G	76.5
CoaT-Lite Tiny [41]	224 ²	5.7M	1.6G	76.6
ConViT-Tiny+ [6]	224 ²	10.0M	2G	76.7
DeiT-2G [32] ‡	224 ²	9.5M	2.0G	77.6
CoaT-Lite Mini [41]	224 ²	11.0M	2.0G	78.9
BoT-S1-50 [27]	224 ²	20.8M	4.3G	79.1
Swin-2G [22] ‡	224 ²	12.8M	2.0G	79.2
Mobile-Former-508M	224 ²	14.0M	508M	79.3

Table 3. **Comparing Mobile-Former with vision transformer variants** evaluated on ImageNet [7] classification. Here, we choose ViT variants that use image resolution 224×224 and are trained *without* distillation from a teacher network. We group ViT models based on FLOPs (using 1.5G as threshold) and rank them based on top-1 accuracy. ‡ indicates our implementation.

network. Mobile-Former achieves higher accuracy but uses 3~4 times less computational cost. This is because that Mobile-Former uses significantly fewer tokens to model global interaction and leverages MobileNet to extract local features efficiently. Note that our Mobile-Former (trained in 450 epochs without distillation) even outperforms LeViT [10] which leverages the distillation of a teacher network and much longer training (1000 epochs). Our method achieves higher top-1 accuracy (77.9% vs. 76.6%) but uses

Model	#Params	MAdds	Top-1	Top-5
<i>Mobile</i> (using ReLU)	6.1M	259M	74.2	91.8
+ <i>Former</i> and Bridge	10.1M	290M	76.8 _(+2.6)	93.2 _(+1.4)
+ DY-ReLU in <i>Mobile</i>	11.4M	294M	77.8 _(+1.0)	93.7 _(+0.5)

Table 4. **Ablation of *Former*+bridge and dynamic ReLU** evaluated on ImageNet classification. Mobile-Former-294M is used.

fewer FLOPs (294M vs. 305M) than LeViT.

Accuracy–FLOP plot: Figure 2 compares Mobile-Former with more CNN models (e.g. GhostNet [11]) and vision transformer variants (e.g. Swin [22] and DeiT [32]) in one plot. We implement Swin and DeiT from 100M to 2G FLOPs, by carefully reducing network width and height. Mobile-Former clearly outperforms both CNNs and ViT variants, showcasing the parallel design to integrate MobileNet and transformer. Although vision transformers are inferior to efficient CNNs by a large margin, our work demonstrates that the transformer can also contribute to the low FLOP regime with proper architecture design.

5.2. Ablations

In this subsection, we show Mobile-Former is effective and efficient via several ablations performed on ImageNet classification. Here, Mobile-Former-294M is used and all models are trained for 300 epochs. Moreover, we summarize interesting observations on the visualization of the two-way bridge between *Mobile* and *Former*.

Mobile-Former is effective: Mobile-Former is more effective than MobileNet as it encodes global interaction via *Former*, resulting in more accurate prediction. As shown in Table 4, adding *Former* and bridge (*Mobile*→*Former* and *Mobile*←*Former*) only costs 12% of the computational cost, but gains 2.6% top-1 accuracy over the baseline that uses *Mobile* alone. In addition, using dynamic ReLU [4] in *Mobile* sub-block (see Figure 3) gains additional 1.0% top-1 accuracy. Note that the parameters in dynamic ReLU is generated by using the first global token. This validates our parallel design in Mobile-Former. We also find that increasing kernel size ($3\times 3 \rightarrow 5\times 5$) of the depthwise convolution in *Mobile* only introduces negligible gain (see Table 14 in appendix B), as the reception field of *Mobile* is enlarged by fusing global features from *Former*.

Mobile-Former is efficient: Mobile-Former is not only effective in encoding both local processing and global interaction, but achieves this *efficiently*. Ablations below show that *Former* only requires a few global tokens with low dimension. In addition, the efficient parallel design of Mobile-Former is stable when removing FFN in *Former* or replacing multi-head attention with position mixing MLP [31].

Number of tokens in Former: Table 5 shows the ImageNet classification results for using different number of global

#Tokens	#Params	MAdds	Top-1	Top-5
1	11.4M	269M	77.1	93.2
3	11.4M	279M	77.6	93.6
6	11.4M	294M	77.8	93.7
9	11.4M	309M	77.7	93.8

Table 5. **Ablation of the number of tokens** on ImageNet classification. Mobile-Former-294M is used.

Token Dimension	#Params	MAdds	Top-1	Top-5
64	7.3M	277M	76.8	93.1
128	9.1M	284M	77.3	93.5
192	11.4M	294M	77.8	93.7
256	14.3M	308M	77.8	93.7
320	17.9M	325M	77.6	93.6

Table 6. **Ablation of token dimension** on ImageNet classification. Mobile-Former-294M is used.

Attention	FFN	#Params	MAdds	Top-1	Top-5
MHA	✓	11.4M	294M	77.8	93.7
MHA	✗	9.8M	284M	77.5	93.6
Pos-Mix-MLP	✓	10.5M	284M	77.3	93.5

Table 7. **Ablation of multi-head attention (MHA) and FFN** on ImageNet classification. Mobile-Former-294M is used.

tokens in *Former*. The token dimension is 192. Interestingly, even a single global token achieves a good performance (77.1% top-1 accuracy). Additional improvement (0.5% and 0.7% top-1 accuracy) is achieved when using 3 and 6 tokens. The improvement stops when more than 6 tokens are used. Such compactness of global tokens is a key contributor to the efficiency of Mobile-Former.

Token dimension: Table 6 shows the results for different token lengths (or dimensions). Here, six global tokens are used in *Former*. The accuracy improves from 76.8% to 77.8% when token dimension increases from 64 to 192, but converges when higher dimension is used. This further supports the efficiency of *Former*. With six tokens of dimension 192, the total computational cost of *Former* and the bridge only consumes 12% of the overall budget (35M/294M).

FFN in Former: As shown in Table 7, removing FFN introduces a small drop in top-1 accuracy (−0.3%). Compared to the important role of FFN in the original vision transformer, FFN has limited contribution in Mobile-Former. We believe this is because FFN is not the only module for channel fusion in Mobile-Former. The 1×1 convolution in *Mobile* helps the channel fusion of local features, while the projection matrix W^O in *Mobile*→*Former* (see Equation 1) contributes to the fusion between local and global features.

Multi-head attention (MHA) vs. position-mixing MLP: Table 7 shows the results of replacing multi-head attention (MHA) with token/position mixing MLP [31] in both *Former* and bridge (*Mobile*→*Former* and *Mobile*←*Former*).

Model	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	MAdds (G)	#Params (M)
Shuffle-V2 [25]	25.9	41.9	26.9	12.4	28.0	36.4	2.6 (161)	0.8 (10.4)
MF-151M	34.2	53.4	36.0	19.9	36.8	45.3	2.6 (161)	4.9 (14.4)
Mobile-V3 [15]	27.2	43.9	28.3	13.5	30.2	37.2	4.7 (162)	2.8 (12.3)
MF-214M	35.8	55.4	38.0	21.8	38.5	46.8	3.9 (162)	5.7 (15.2)
ResNet18 [13]	31.8	49.6	33.6	16.3	34.3	43.2	29 (181)	11.2 (21.3)
MF-294M	36.6	56.6	38.6	21.9	39.5	47.9	5.5 (164)	6.5 (16.1)
ResNet50 [13]	36.5	55.4	39.1	20.4	40.3	48.1	84 (239)	23.3 (37.7)
PVT-Tiny [37]	36.7	56.9	38.9	22.6	38.8	50.0	70 (221)	12.3 (23.0)
ConT-M [42]	37.9	58.1	40.2	23.0	40.6	50.4	65 (217)	16.8 (27.0)
MF-508M	38.0	58.3	40.3	22.9	41.2	49.7	9.8 (168)	8.4 (17.9)

Table 8. **COCO object detection results in RetinaNet framework.** All models are trained on `train2017` for 12 epochs ($1\times$) from ImageNet pretrained weights, and tested on `val2017`. We use initial **MF** (e.g. MF-508M) to refer Mobile-Former. MAdds and #Params are in the format of “backbone (total)”. MAdds is based on the image size 800×1333 .

The top-1 accuracy drops from 77.8% to 77.3%. The implementation of MLP is more efficient by a single matrix multiplication, but it is static (i.e. not adaptive to different input images).

Mobile-Former is explainable: We observe three interesting patterns in the two-way bridge (*Mobile*→*Former* and *Mobile*←*Former*). First, from low to high levels, the global tokens change their focus from edges/corners, to foreground/background, and finally on the most discriminative region. Second, the cross attention has more diversity across tokens at lower levels than high levels. Thirdly, the separation between foreground and background is found at middle layers of *Mobile*←*Former*. The detailed visualization is shown in appendix C.

5.3. Object Detection

Object detection experiments are conducted on COCO 2017 [21], which contains 118K training and 5K validation images. We evaluate Mobile-Former in two detection frameworks: (a) comparing with other backbone networks in RetinaNet [20] framework that has dense proposals, and (b) end-to-end comparison with DETR [1] where sparse proposals are used.

RetinaNet training setup: We follow the standard settings of RetinaNet [20] and replace the backbone with our Mobile-Former to generate multi-scale feature maps. All models are trained for 12 epochs ($1\times$) from ImageNet pretrained weights.

DETR training setup: All Mobile-Former models are trained for 300 epochs on 8 GPUs with 2 images per GPU. AdamW optimizer is used with initial learning rate $1e-5$ for the backbone and $1e-4$ for the head. The learning rate drops by a factor of 10 after 200 epochs. The weight decay is $1e-4$ and dropout rate is 0.1. BatchNorm layers in the ImageNet pretrained Mobile-Former backbone are frozen. The head includes 100 object queries with 256 channels.

Model	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	MAdds (G)	#Params (M)
DETR [1]	42.0	62.4	44.2	20.5	45.8	61.1	86	41.3
DETR-DC5 [1]	43.3	63.1	45.9	22.5	47.3	61.1	187	41.3
E2E-MF-508M	43.1	61.9	46.8	23.8	46.5	60.4	41.4	26.6
E2E-MF-294M	40.5	58.8	43.5	20.6	44.0	56.9	24.1	25.1
E2E-MF-214M	39.3	57.3	42.1	19.9	42.4	56.6	17.8	20.1
E2E-MF-151M	37.2	54.5	39.9	17.4	39.8	54.9	12.7	14.8

Table 9. **End-to-end object detection results on COCO.** All models are trained on `train2017` and tested on `val2017`. DETR baselines are trained for 500 epochs, while our MobileFormer models are trained for 300 epochs. We use initial **E2E-MF** (e.g. E2E-MF-508M) to refer end-to-end MobileFormer detectors. MAdds is based on image size 800×1333 .

Efficient and effective backbone in RetinaNet: In Table 8, we compare MobileFormer with both CNNs (ResNet [13], MobileNetV3 [15], ShuffleNetV2 [25]) and vision transformers (PVT [37] and ConT [42]). MobileFormer significantly outperforms MobileNetV3 and ShuffleNetV2 by 8.3+ AP under similar computational cost. Compared to ResNet and transformer variants, MobileFormer achieves higher AP with significantly less FLOPs in the backbone. Specifically, MobileFormer-508M only takes 9.8G FLOPs in backbone but achieves 38.0 AP, outperforming ResNet-50, PVT-Tiny, and ConT-M which consume 7 times more computation (65G to 84G FLOPs) in the backbone. This showcases MobileFormer as an effective and efficient backbone in object detection.

Efficient and effective end-to-end detector: Table 9 compares end-to-end MobileFormer detectors (denoted by prefix E2E-MF) with DETR [1]. All models use 100 object queries. Our E2E-MF-508M gains 1.1 AP over DETR but consumes fewer FLOPs (48%), fewer model parameters (64%) and fewer training epochs (300 vs. 500). It is just a little behind DETR-DC5 that has four times more FLOPs. The other three MobileFormer variants achieve 40.5, 39.3 and 37.2 AP with 24.1G, 17.8G and 12.7G FLOPs respectively, providing more compact end-to-end detectors.

Ablations of key components: Table 10 shows the effects of three proposed components in the end-to-end MobileFormer detector: (a) spatial-aware dynamic ReLU in backbone, (b) multi-scale MobileFormer head, and (c) adapting position embedding. E2E-MF-508M is used and all models are trained for 300 epochs. Compared to DETR trained in 300 epochs, replacing ResNet-50 backbone with MF-508M saves computations, but results in 1.2 AP drop. Adding spatial-aware dynamic ReLU gains 1.1 AP ($39.4 \rightarrow 40.5$). Then using multi-scale MobileFormer head to replace DETR’s encoder/decoder gains another 0.9 AP ($40.5 \rightarrow 41.4$). The detection of small and medium objects is improved, while a slight degradation is found in large objects. Finally, adapting position embedding provides additional 1.7 AP gain. The three proposed components are

Backbone	SP-DY ReLU	Head	Adapt PE	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
ResNet-50		DETR		40.6	61.6	42.7	19.9	44.3	60.2
MF-508M		DETR		39.4	59.4	41.3	18.0	42.7	58.8
MF-508M	✓	DETR		40.5	60.7	42.4	19.0	44.0	60.0
MF-508M	✓	MF-Head		41.4	60.4	43.9	21.1	45.0	59.3
MF-508M	✓	MF-Head	✓	43.1	61.9	46.8	23.8	46.5	60.4

Table 10. **Ablation of end-to-end MobileFormer on COCO object detection.** All models are trained on `train2017` for 300 epochs and tested on `val2017`. The first line is the baseline of DETR [1] trained in 300 epochs. MF-508M refers to MobileFormer-508M backbone. SP-DY-ReLU refers to spatial-aware dynamic ReLU in backbone. DETR head includes 6 encoder and 6 decoder layers with resolution $\frac{1}{32}$. MF head includes 9 MobileFormer blocks (5, 2, 2 at resolution $\frac{1}{32}$, $\frac{1}{16}$, and $\frac{1}{8}$ respectively).

complementary, gaining 3.7 AP ($39.4 \rightarrow 43.1$) in total.

6. Limitations and Discussions

Although MobileFormer has faster inference than MobileNetV3 [15] for large images, it is slower as the image becomes smaller. Please see Figure 5 in appendix B for comparison on inference latency between MobileFormer-214M and MobileNetV3 Large. They have similar FLOPs (214M vs 217M), but MobileFormer is more accurate. The comparison is performed on multiple image sizes due to the resolution variation across tasks (e.g. classification, detection). As image resolution decreases, MobileFormer loses its leading position to MobileNetV3. This is because *Former* and embedding projections in *Mobile* \rightarrow *Former* and *Mobile* \leftarrow *Former* are resolution independent, and their PyTorch implementations are not as efficient as convolution. Thus, the overhead is relative large when image is small, but becomes negligible as image size grows. The runtime performance of MobileFormer can be further improved by optimizing the implementation of these components. We will investigate these in the future work.

Another limitation is that MobileFormer is not efficient in parameters especially when performing image classification, due to the parameter-heavy classification head. For instance, the head of MobileFormer-294M consumes 4.6M of total 11.4M parameters (40%). This problem is mitigated when switching to object detection, due to the removal of image classification head. In addition, *Former* and two-way bridge are computationally but not parametrically efficient.

7. Conclusion

This paper presents MobileFormer, a new parallel design of MobileNet and Transformer with two-way bridge in between to communicate. It leverages the efficiency of MobileNet in local processing and the advantage of Transformer in encoding global interaction. This design is not only effective to boost accuracy, but also efficient to save

computational cost. It outperforms both efficient CNNs and vision transformer variants with a clear margin on image classification and object detection in the low FLOP regime. Furthermore, we build an end-to-end Mobile-Former detector that outperforms DETR but consumes significantly less computations and parameters. We hope Mobile-Former encourage new design of efficient CNNs and transformers.

References

- [1] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. [2](#), [4](#), [5](#), [7](#), [8](#)
- [2] Hanqing Chen, Yunhe Wang, Chunjing Xu, Boxin Shi, Chao Xu, Qi Tian, and Chang Xu. Addernet: Do we really need multiplications in deep learning? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [2](#)
- [3] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#)
- [4] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic relu. In *ECCV*, 2020. [2](#), [3](#), [4](#), [6](#)
- [5] Ekin D. Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [5](#)
- [6] Stéphane d’Ascoli, Hugo Touvron, Matthew Leavitt, Ari Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. *arXiv preprint arXiv:2103.10697*, 2021. [2](#), [3](#), [5](#), [6](#)
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [5](#), [6](#), [12](#)
- [8] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. *arXiv preprint arXiv:2107.00652*, 2021. [2](#)
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. [1](#), [2](#), [3](#)
- [10] Benjamin Graham, Alaaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Hervé Jégou, and Matthijs Douze. Levit: a vision transformer in convnet’s clothing for faster inference. *arXiv preprint arXiv:22104.01136*, 2021. [1](#), [2](#), [3](#), [6](#)
- [11] Kai Han, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. Ghostnet: More features from cheap operations. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [1](#), [2](#), [6](#)
- [12] Kai Han, Yunhe Wang, Qiulin Zhang, Wei Zhang, Chunjing XU, and Tong Zhang. Model rubiks cube: Twisting resolution, depth and width for tinynets. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 19353–19364. Curran Associates, Inc., 2020. [2](#)
- [13] 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. [3](#), [7](#), [8](#)
- [14] Geoffrey E. Hinton. How to represent part-whole hierarchies in a neural network. *CoRR*, abs/2102.12627, 2021. [2](#)
- [15] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, Quoc V. Le, and Hartwig Adam. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. [1](#), [2](#), [4](#), [5](#), [6](#), [7](#), [8](#), [11](#), [12](#)
- [16] 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. [1](#), [2](#)
- [17] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [2](#)
- [18] Yunsheng Li, Yinpeng Chen, Xiyang Dai, Dongdong Chen, Mengchen Liu, Lu Yuan, Zicheng Liu, Lei Zhang, and Nuno Vasconcelos. Micronet: Improving image recognition with extremely low flops. In *International Conference on Computer Vision*, 2021. [1](#), [2](#), [4](#), [11](#), [12](#)
- [19] T. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 936–944, July 2017. [4](#)
- [20] 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 (ICCV)*, Oct 2017. [2](#), [7](#)
- [21] 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*, pages 740–755. Springer, 2014. [5](#), [7](#), [12](#)
- [22] 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](#), [5](#), [6](#)
- [23] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. [5](#)
- [24] Ningning Ma, X. Zhang, J. Huang, and J. Sun. Weightnet: Revisiting the design space of weight networks. volume abs/2007.11823, 2020. [5](#), [6](#)

- [25] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *The European Conference on Computer Vision (ECCV)*, September 2018. 2, 5, 6, 7, 8
- [26] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018. 1, 2, 3, 4, 11, 12
- [27] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16519–16529, June 2021. 2, 6
- [28] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, pages 6105–6114, Long Beach, California, USA, 09–15 Jun 2019. 2, 5, 6
- [29] Mingxing Tan and Quoc V. Le. Mixconv: Mixed depthwise convolutional kernels. In *30th British Machine Vision Conference 2019*, 2019. 2
- [30] Mingxing Tan, Ruoming Pang, and Quoc V. Le. Efficientdet: Scalable and efficient object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2
- [31] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. Mlp-mixer: An all-mlp architecture for vision, 2021. 6, 7
- [32] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers and distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020. 1, 2, 5, 6
- [33] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers, 2021. 13
- [34] Keivan Alizadeh vahid, Anish Prabhu, Ali Farhadi, and Mohammad Rastegari. Butterfly transform: An efficient fft based neural architecture design. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2
- [35] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12894–12904, June 2021. 2
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 1, 3, 4
- [37] 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, 2021. 5, 6, 7, 8
- [38] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. 11
- [39] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers, 2021. 1, 2, 3
- [40] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross B. Girshick. Early convolutions help transformers see better. *CoRR*, abs/2106.14881, 2021. 1, 2, 3, 4, 5, 6
- [41] Weijian Xu, Yifan Xu, Tyler Chang, and Zhuowen Tu. Co-scale conv-attentional image transformers, 2021. 5, 6
- [42] Haotian Yan, Zhe Li, Weijian Li, Changhu Wang, Ming Wu, and Chuang Zhang. Contnet: Why not use convolution and transformer at the same time? *CoRR*, abs/2104.13497, 2021. 6, 7, 8
- [43] Brandon Yang, Gabriel Bender, Quoc V. Le, and Jiquan Ngiam. Condconv: Conditionally parameterized convolutions for efficient inference. In *NeurIPS*, 2019. 2
- [44] 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. 2, 5, 6
- [45] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. 5
- [46] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [47] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2020. 5
- [48] Daquan Zhou, Qi-Bin Hou, Y. Chen, Jiashi Feng, and S. Yan. Rethinking bottleneck structure for efficient mobile network design. In *ECCV*, August 2020. 2
- [49] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer, 2021. 13
- [50] Daquan Zhou, Yujun Shi, Bingyi Kang, Weihao Yu, Zihang Jiang, Yuan Li, Xiaojie Jin, Qibin Hou, and Jiashi Feng. Refiner: Refining self-attention for vision transformers, 2021. 13

A. Mobile-Former Architecture

Seven model variants: Table 11 shows six Mobile-Former models (508M–52M). The smallest model Mobile-Former-26M has similar architecture to Mobile-Former-52M except replacing all 1×1 convolutions with group convolution (group=4). They are used either in image classification or as the backbone of object detectors. These models are manually designed without searching for the optimal architecture parameters (e.g. width or depth). We follow the well known rules used in MobileNet [15, 26]: (a) the number of channels increases across stages, and (b) the channel expansion rate starts from three at low levels and increases to six at high levels. For the four bigger models (508M–151M), we use six global tokens with dimension 192 and eleven Mobile-Former blocks. But these four models have different widths. Mobile-Former-96M and Mobile-Former-52M are shallower (with only eight Mobile-Former blocks) to meet the low computational budget.

Downsample Mobile-Former block: Note that stage 2–5 has a downsample variant of Mobile-Former block (denoted as $M-F^\downarrow$ in Table 11) to handle the spatial downsampling. $M-F^\downarrow$ has a slightly different *Mobile* sub-block that includes four (instead of three) convolutional layers (depthwise→pointwise→depthwise→pointwise), where the first depthwise convolution layer has stride two. The number of channels expands in each depthwise convolution, and squeezes in the following pointwise convolution. This saves computations as the two costly pointwise

convolutions are performed at the lower resolution after downsampling.

Training hyper-parameters: Table 12 lists three hyper-parameters (initial learning rate, weight decay and dropout rate) used for training Mobile-Former models in ImageNet classification. Their values increase as the model becomes bigger to prevent overfitting. Our implementation is based on timm framework [38].

Head model variants in end-to-end object detection: Table 13 shows the head structures for four end-to-end Mobile-Former detectors. All share similar structure and have 100 object queries with dimension 256. The largest model (E2E-MF-508M) has the heaviest head with 9 Mobile-Former blocks over three scales, while the other three smaller models have 9, 7, 5 blocks respectively over two scales to save computations.

All models start by projecting the input feature map linearly to 256 channels through a 1×1 convolution. Then multiple Mobile-Former blocks are stacked with upsampling block in between to move upscale. The upsampling block (denoted as “up-conv”) includes three steps: (a) increasing feature resolution by two using bilinear interpolation, (b) adding the feature output from the backbone, and (c) applying a 3×3 depthwise and a pointwise convolution. To handle the computational boost due to resolution increasing, we use lite bottleneck [18] in *Mobile*. Moreover, we find that the performance can be further improved at a small additional cost by adding multi-head attention in *Mo-*

Stage	Mobile-Former-508M			Mobile-Former-294M			Mobile-Former-214M			Mobile-Former-151M			Mobile-Former-96M			Mobile-Former-52M		
	Block	#exp	#out	Block	#exp	#out	Block	#exp	#out	Block	#exp	#out	Block	#exp	#out	Block	#exp	#out
token	6×192			6×192			6×192			6×192			4×128			3×128		
stem	conv 3×3	–	24	conv 3×3	–	16	conv 3×3	–	12	conv 3×3	–	12	conv 3×3	–	12	conv 3×3	–	8
1	bneck-lite	48	24	bneck-lite	32	16	bneck-lite	24	12	bneck-lite	24	12	bneck-lite	24	12			
2	$M-F^\downarrow$	144	40	$M-F^\downarrow$	96	24	$M-F^\downarrow$	72	20	$M-F^\downarrow$	72	16	$M-F^\downarrow$	72	16	bneck-lite [↓]	24	12
	M-F	120	40	M-F	96	24	M-F	60	20	M-F	48	16	M-F	48	16	M-F	36	12
3	$M-F^\downarrow$	240	72	$M-F^\downarrow$	144	48	$M-F^\downarrow$	120	40	$M-F^\downarrow$	96	32	$M-F^\downarrow$	96	32	$M-F^\downarrow$	72	24
	M-F	216	72	M-F	192	48	M-F	160	40	M-F	96	32	M-F	96	32	M-F	72	24
4	$M-F^\downarrow$	432	128	$M-F^\downarrow$	288	96	$M-F^\downarrow$	240	80	$M-F^\downarrow$	192	64	$M-F^\downarrow$	192	64	$M-F^\downarrow$	144	48
	M-F	512	128	M-F	384	96	M-F	320	80	M-F	256	64	M-F	256	64	M-F	192	48
	M-F	768	176	M-F	576	128	M-F	480	112	M-F	384	88	M-F	384	88	M-F	288	64
	M-F	1056	176	M-F	768	128	M-F	672	112	M-F	528	88	M-F	528	88			
5	$M-F^\downarrow$	1056	240	$M-F^\downarrow$	768	192	$M-F^\downarrow$	672	160	$M-F^\downarrow$	528	128	$M-F^\downarrow$	528	128	$M-F^\downarrow$	384	96
	M-F	1440	240	M-F	1152	192	M-F	960	160	M-F	768	128	M-F	768	128	M-F	576	96
	M-F	1440	240	M-F	1152	192	M-F	960	160	M-F	768	128	M-F	768	128	conv 1×1	–	768
	conv 1×1	–	1440	conv 1×1	–	1152	conv 1×1	–	960	conv 1×1	–	768	conv 1×1	–	768	conv 1×1	–	576
pool concat	–	–	1632	–	–	1344	–	–	1152	–	–	960	–	–	896	–	–	704
FC1	–	–	1920	–	–	1920	–	–	1600	–	–	1280	–	–	1280	–	–	1024
FC2	–	–	1000	–	–	1000	–	–	1000	–	–	1000	–	–	1000	–	–	1000

Table 11. **Specification of Mobile-Former models.** “bneck-lite” denotes the lite bottleneck block [18]. “bneck-lite[↓]” denotes the downsample variant of lite bottleneck, in which the depthwise convolution has stride 2. “M-F” denotes the Mobile-Former block and “ $M-F^\downarrow$ ” denotes the Mobile-Former block for downsampling. Mobile-Former-26M has a similar architecture to Mobile-Former-52M except replacing all 1×1 convolutions with group convolution (group=4).

Model	Learning Rate	Weight Decay	Dropout
Mobile-Former-26M	8e-4	0.08	0.1
Mobile-Former-52M	8e-4	0.10	0.2
Mobile-Former-96M	8e-4	0.10	0.2
Mobile-Former-151M	9e-4	0.10	0.2
Mobile-Former-214M	9e-4	0.15	0.2
Mobile-Former-294M	1e-3	0.20	0.3
Mobile-Former-508M	1e-3	0.20	0.3

Table 12. **Hyper-parameters** of seven Mobile-Former models for ImageNet [7] classification.

Stage	E2E-MF 508M	E2E-MF 294M	E2E-MF 214M	E2E-MF 151M
query	100×256	100×256	100×256	100×256
$\frac{1}{32}$	projection †M-F ×5	projection †M-F ×6	projection †M-F ×5	projection †M-F ×3
$\frac{1}{16}$	up-conv M-F ×2	up-conv M-F ×3	up-conv M-F ×2	up-conv M-F ×2
$\frac{1}{8}$	up-conv M-F ×2	–	–	–

Table 13. **Specification of head variants** in end-to-end Mobile-Former object detectors. 100 object queries with dimension 256 are used. “projection” denotes projecting an input feature map linearly to 256 channels (through a 1×1 convolution). “up-conv” denotes a convolutional block for upsampling that includes bilinear interpolation followed by a 3×3 depthwise convolution. “M-F $\times 2$ ” refers to stacking two Mobile-Former blocks. In the detection head, we use lite bottleneck [18] in *Mobile* sub-block to reduce the computational cost. At the lowest resolution $\frac{1}{32}$, multi-head attention is added into *Mobile*, which is denoted as †M-F.

Mobile sub-block at the lowest scale ($\frac{1}{32}$) of the head (denoted as †M-F). It is especially helpful for detecting large objects.

B. More Experimental Results

Inference latency: Figure 5 compares between Mobile-Former-214M and MobileNetV3 Large [15] on inference latency, as they have similar FLOPs (214M vs. 217M). The latency is measured on an Intel(R) Xeon(R) CPU E5-2650 v3 (2.3GHz), following the common settings (single-thread with batch size 1) in [15, 26]. The comparison is performed on multiple image sizes due to the resolution variation across tasks (e.g. classification, detection). Mobile-Former is behind MobileNetV3 at low resolution (224×224). As the image resolution increases, the gap shrinks until resolution 750×750 , after which Mobile-Former has faster inference.

This is because *Former* and embedding projections in *Mobile*→*Former* and *Mobile*←*Former* are resolution independent, and their PyTorch implementations are not as efficient as convolution. Thus, the overhead is relative large when image is small, but becomes negligible as image size

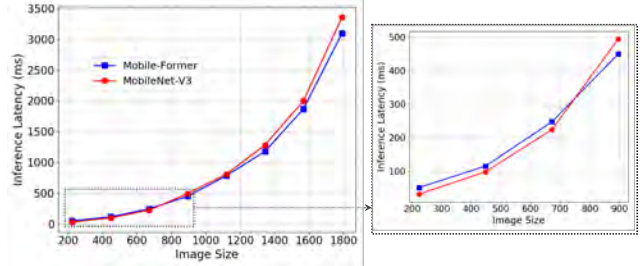


Figure 5. **Inference latency** over different image sizes. The latency is measured on an Intel(R) Xeon(R) CPU E5-2650 v3 (2.3GHz), following the common settings (single-thread with batch size 1) in [15, 26]. Mobile-Former-214M is compared with MobileNetV3 Large [15] as they have similar FLOPs (214M vs. 217M). Mobile-Former is slower when image size is small, but has faster inference than MobileNetV3 as image size grows above 750×750 . Best viewed in color.

Kernel Size in Mobile	#Param	MAdds	Top-1	Top-5
3×3	11.4M	294M	77.8	93.7
5×5	11.5M	332M	77.9	93.9

Table 14. **Ablation of the kernel size in the depthwise convolution** (in *Mobile* sub-block). The evaluation is performed on ImageNet [7] classification. Mobile-Former-294M is used.

MHA in <i>Mobile</i> at scale $\frac{1}{32}$	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	MAdds (G)	#Params (M)
	42.5	61.0	46.0	23.2	46.3	58.7	36.0	23.7
✓	43.1	61.9	46.8	23.8	46.5	60.4	41.4	26.6

Table 15. **Ablation of multi-head attention (MHA) in *Mobile*** at resolution $\frac{1}{32}$ of the detection head. The evaluation is performed on COCO [21] object detection. Both models are trained on train2017 for 300 epochs and tested on val2017. E2E-MF-508M is used. MAdds is based on image size 800×1333 .

grows. The runtime performance of Mobile-Former can be further improved by optimizing the implementation of these components. We will investigate this in the future work.

Ablation of the kernel size in *Mobile*: We perform an ablation on the kernel size of the depthwise convolution in *Mobile*, to validate the contribution of *Former* and bridge on global interaction. Table 14 shows that the gain of increasing kernel size (from 3×3 to 5×5) is negligible. We believe this is because *Former* and the bridge enlarge the reception field for *Mobile* via fusing global features. Therefore, using larger kernel size is not necessary in Mobile-Former.

Ablation of multi-head attention in *Mobile* at resolution $\frac{1}{32}$ of the detection head: Table 15 shows the effect of using multi-head attention (MHA) in the five blocks at the lowest resolution $\frac{1}{32}$ (†M-F in Table 13). Without MHA, a solid performance (42.5 AP) is achieved at low FLOPs (36.0G). Adding MHA gains 0.6 AP with 15% additional



Figure 6. **Cross attention over the entire featuremap for the first token in $Mobile \rightarrow Former$** across all Mobile-Former blocks. Attention is normalized over pixels, showing the focused region. The focused region changes from low to high levels. The token starts paying more attention to edges/corners at block 2–4. Then it focuses more on a large region rather than scattered small pieces at block 5–12. The focused region shifts between the foreground (person and horse) and background (grass). Finally, it locks the most discriminative part (horse body and head) for classification. Best viewed in color.

computational cost. It is especially helpful for detecting large objects (58.7 \rightarrow 60.4 AP_L).

C. Visualization

In order to understand the collaboration between *Mobile* and *Former*, we visualize the cross attention on the two-way bridge (i.e. $Mobile \rightarrow Former$ and $Mobile \leftarrow Former$) in Figure 6, 7, and 8. The ImageNet pretrained Mobile-Former-294M is used, which includes six global tokens and eleven Mobile-Former blocks. We observe three interesting patterns as follows:

Patten 1 – global tokens shift focus over levels: The focused regions of global tokens change progressively from low to high levels. Figure 6 shows the cross attention over pixels for the first token in $Mobile \rightarrow Former$. This token begins focusing on local features, e.g. edges/corners (at block 2-4). Then it pays more attention to regions with connected pixels. Interestingly, the focused region shifts between foreground (person and horse) and background (grass) across blocks. Finally, it locates the most discriminative region (horse body and head) for classification.

Pattern 2 – foreground and background are separated

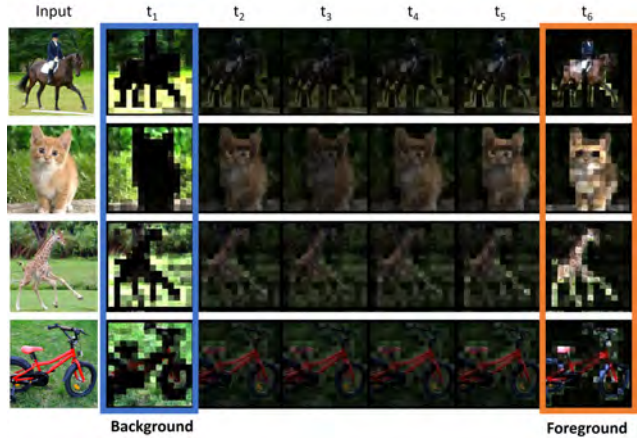


Figure 7. **Cross attention in $Mobile \leftarrow Former$ separates foreground and background at middle layers.** Attention is normalized over tokens showing the contribution of different tokens at each pixel. Block 8 is chosen where background pixels pay more attention to the first token and foreground pixels pay more attention to the last token. Best viewed in color.

in middle layers: The separation between foreground and background is surprisingly found in $Mobile \leftarrow Former$ at middle layers (e.g. block 8). Figure 7 shows the cross attention over six tokens for each pixel in the featuremap. Clearly, the foreground and background are separated in the first and last tokens. This shows that global tokens learn meaningful prototypes that cluster pixels with similar semantics.

Pattern 3 – attention diversity across tokens diminishes:

The attention has more diversity across tokens at low levels than high levels. As shown in Figure 8, each column corresponds to a token, and each row corresponds to a head in the corresponding multi-head cross attention. Note that the attention is normalized over pixels in $Mobile \rightarrow Former$ (left half), showing the focused region per token. In contrast, the attention in $Mobile \leftarrow Former$ is normalized over tokens, comparing the contribution of different tokens at each pixel. Clearly, the six tokens at block 3 and 5 have different cross attention patterns in both $Mobile \rightarrow Former$ and $Mobile \leftarrow Former$. Similar attention maps over tokens are clearly observed at block 8. At block 12, the last five tokens share a similar attention pattern. Note that the first token is the classification token fed into the classifier. The similar observation on token diversity has been identified in recent studies on ViT [33, 49, 50]. The full visualization of two-way cross attention for all blocks is shown in Figure 9.

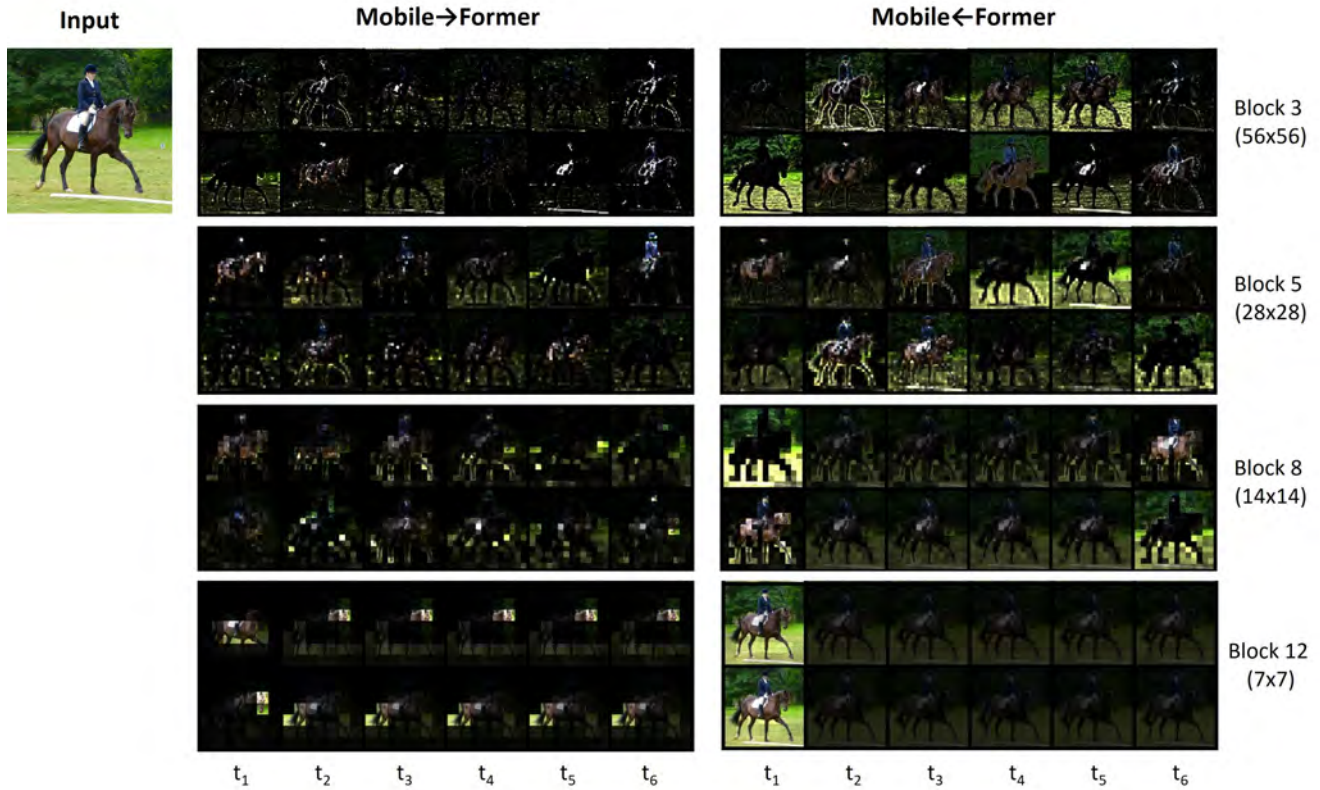


Figure 8. **Visualization of the two-way cross attention:** *Mobile*→*Former* and *Mobile*←*Former*. Mobile-Former-294M is used, which includes six tokens (each corresponds to a column). Four blocks with different input resolutions are selected and each has two attention heads that are visualized in two rows. Attention in *Mobile*→*Former* (left half) is normalized over pixels, showing the focused region per token. Attention in *Mobile*←*Former* (right half) is normalized over tokens showing the contribution of different tokens at each pixel. The cross attention has less variation across tokens at high levels than low levels. Specifically, token 2–5 in the last block have very similar cross attention. Best viewed in color.

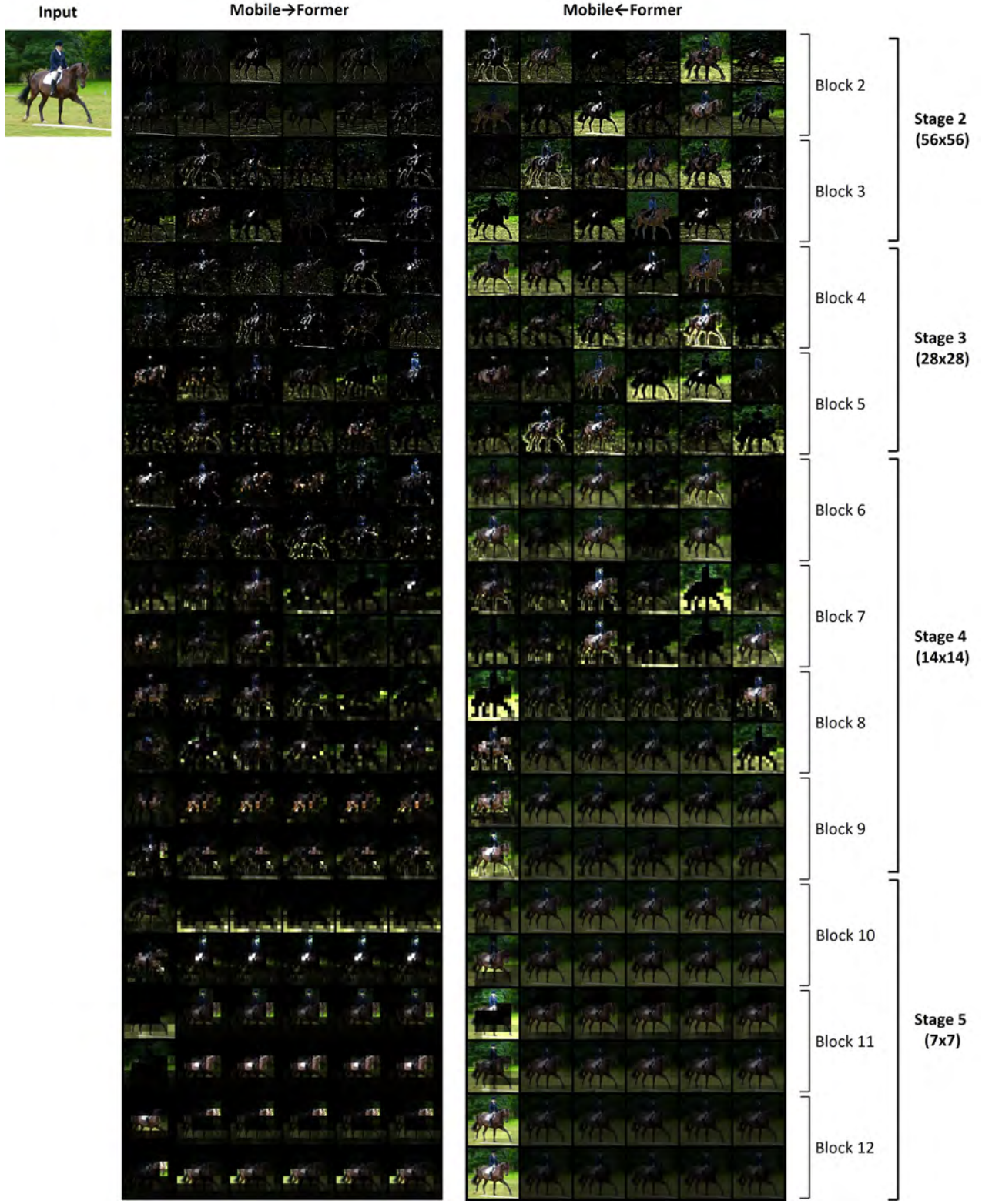


Figure 9. **Visualization of the two-way cross attention:** *Mobile→Former* and *Mobile←Former*. MobileFormer-294M is used, which includes six tokens (each corresponds to a column) and eleven MobileFormer blocks (block 2–12) across four stages. Each block has two attention heads that are visualized in two rows. Attention in *Mobile→Former* (left) is normalized over pixels, showing the focused region per token. Attention in *Mobile←Former* (right) is normalized over tokens showing the contribution of different tokens at each pixel.