

# CMT: Convolutional Neural Networks Meet Vision Transformers

Jianyuan Guo<sup>1,2</sup>, Kai Han<sup>1</sup>, Han Wu<sup>2</sup>, Chang Xu<sup>2,\*</sup>,  
Yehui Tang<sup>1</sup>, Chunjing Xu<sup>1</sup>, Yunhe Wang<sup>1\*</sup>

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

<sup>2</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

Vision transformers have been successfully applied to image recognition tasks due to their ability to capture long-range dependencies within an image. However, there are still gaps in both performance and computational cost between transformers and existing convolutional neural networks (CNNs). In this paper, we aim to address this issue and develop a network that can outperform not only the canonical transformers, but also the high-performance convolutional models. We propose a new transformer based hybrid network by taking advantage of transformers to capture long-range dependencies, and of CNNs to model local features. Furthermore, we scale it to obtain a family of models, called CMTs, obtaining much better accuracy and efficiency than previous convolution and transformer based models. In particular, our CMT-S achieves 83.5% top-1 accuracy on ImageNet, while being 14x and 2x smaller on FLOPs than the existing DeiT and EfficientNet, respectively. The proposed CMT-S also generalizes well on CIFAR10 (99.2%), CIFAR100 (91.7%), Flowers (98.7%), and other challenging vision datasets such as COCO (44.3% mAP), with considerably less computational cost.

## 1 Introduction

The past decades have witnessed the extraordinary contribution of CNNs [15, 44, 42, 49, 50] in the field of computer vision due to its ability of extracting deep discriminative features. Meanwhile, self-attention based transformers [9, 52] has become the de facto most popular models for natural language processing (NLP) tasks, and shown excellent capability of capturing long-distance relationships. Recently, many researchers attempt to apply the transformer based architectures to vision domains, and achieve promising results in various tasks such as image classification [10, 51], object detection [2, 64], and semantic segmentation [63]. Vision transformer (ViT) [10] is the first work to replace the conventional CNN backbone with a pure transformer. Input images with the size of  $224 \times 224 \times 3$  are firstly split into  $14 \times 14 = 196$  discrete non-overlapping patches (with a fixed size of  $16 \times 16 \times 3$  per patch), which are analogous to the word tokens in NLP. The patches are then fed into stacked standard transformer blocks to model global relations and extract feature for image classification. The design paradigm of ViT has heavily inspired the following transformer based models for computer vision, such as IPT [3] for low-level vision and SETR [63] for semantic segmentation.

Despite that transformers have demonstrated excellent capabilities when migrated to vision tasks, their performances are still far inferior to similar-sized convolutional neural network counterparts, *e.g.*, EfficientNets [50]. We believe the reason of such weakness is threefold. Firstly, as mentioned above, images are split into patches in ViT [10] and other transformer based models such as IPT [3] and SETR [63] in order to adapt to the structure designed to take sequences as input in the first

\*Corresponding author.

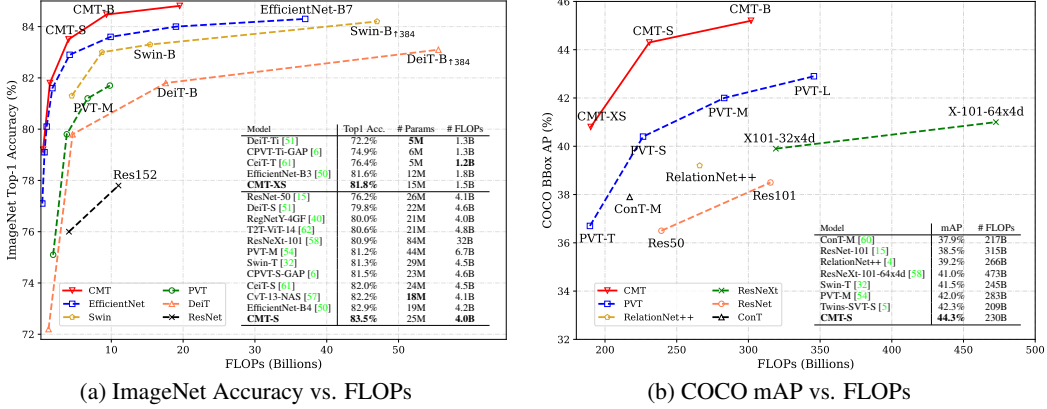


Figure 1: **Performance comparison between CMT and other state-of-the-art models.** (a) Top-1 accuracy on ImageNet, all numbers are for single-crop, single-model. (b) Object detection results on COCO val2017 of different backbones using RetinaNet framework, all numbers are for single-scale, “1x” training schedule.

place. Doing so can greatly simplify the process of applying transformer to image-based tasks. The sequence of patches can be directly fed into a standard transformer and the long-range dependencies between patches can be well captured. However, it ignores the fundamental difference between sequence-based NLP tasks and image-based vision tasks, *e.g.*, the 2D structure and spatial local information within each patch. Secondly, transformer is difficult to explicitly extract low-resolution and multi-scale feature maps due to the fixed patch size, which poses a big challenge to dense prediction tasks such as detection and segmentation. Thirdly, the computational and memory cost of self-attention based transformers are quadratic ( $\mathcal{O}(N^2C)$ ) to the resolution of input images, compared to  $\mathcal{O}(NC^2)$  of convolution-based CNNs. High resolution images are prevalent and common, *e.g.*,  $1333 \times 800$  in COCO [31],  $2048 \times 1024$  in Cityscapes [7]. Using transformers to process such images would inevitably cause the problem of insufficient GPU memory and low computation efficiency.

In this paper, we stand upon the intersection of CNNs and transformers, and propose a novel CMT (CNNs meet transformers) architecture for visual recognition. The proposed CMT takes the advantages of CNNs to compensate for the aforementioned limitations when utilizing pure transformers. As shown in Figure 2(c), input images first go through the convolution stem for fine-grained feature extraction, and are then fed into a stack of CMT blocks for representation learning. Specifically, the introduced CMT block is an improved variant of transformer block whose local information is enhanced by depth-wise convolution. Compared to ViT [10], the features generated from the first stage of CMT can maintain higher resolution, *i.e.*,  $H/4 \times W/4$  against  $H/16 \times W/16$  in ViT, which are essential for other dense prediction tasks. Furthermore, we adopt the stage-wise architecture design similar to CNNs [15, 50, 42] by using four convolutional layer with stride 2, to reduce the resolution (sequence length) gradually and meanwhile increase the dimension flexibly. The stage-wise design helps to extract multi-scale features and alleviate the computation burden caused by high resolution. The local perception unit (LPU) and inverted residual feed-forward network (IRFFN) in CMT block can help capture both local and global structure information within the intermediate features and promote the representation ability of the network. Finally, the average pooling is used to replace the class token in ViT for better classification results. In addition, we propose a simple scaling strategy to obtain a family of CMT variants. Extensive experiments on ImageNet and other downstream tasks demonstrate the superiority of our CMTs in terms of accuracy and FLOPs. For example, our CMT-S achieves 83.5% ImageNet top-1 with only 4.04B FLOPs, while being 14x and 2x smaller than the best existing DeiT [51] and EfficientNet [50], respectively. Besides image classification, CMT can also be easily transferred to other vision tasks to replace the conventional backbone. Our CMT-S based RetinaNet [30] achieves 44.3% mAP on COCO val2017 detection task, outperforming the PVT-based RetinaNet [54] by 3.9% with less computational cost.

## 2 Related Work

The computer vision community prospered in past decades riding the wave of deep learning, and the most popular deep neural networks are often built upon basic blocks, in which a series of convolutional layers are stacked sequentially to capture local information within intermediate features.

However, the limited receptive field of small convolutional kernels makes it difficult to obtain global information, withholding the networks of high performance on challenging tasks such as classification, detection, and segmentation. Thus, many researchers start to dig deeper into the self-attention based transformers which have the ability to capture long-range information. Here we briefly review the conventional CNNs and recently proposed vision transformers.

**Convolutional neural networks.** The first standard CNN was proposed by Yann *et al.* [29] for handwritten character recognition, and the past decades have witnessed that many powerful networks [28, 48, 44, 15, 20] achieved unprecedented success on large scale image classification task [8]. AlexNet [28] and VGGNet [44] showed that a deep neural network composed of convolutional layers and pooling layers can obtain adequate results in recognition. GoogleNet [48] and InceptionNet [49] demonstrated the effectiveness of multiple paths within a basic block. ResNet [15] showed better generalization by adding shortcut connections every two layers to the base network. To alleviate the limited receptive fields in prior research, some researches [53, 56, 19, 18, 37, 41] incorporated attention mechanisms as an operator for adaptation between modalities. Wang *et al.* [53] proposed to stack attention modules sequentially between the intermediate stages of deep residual networks. SENet [19] and GENet [18] adaptively recalibrated channel-wise feature responses by modelling interdependencies between channels. NLNet [55] incorporated the self-attention mechanism into neural networks, providing pairwise interactions across all spatial positions to augment the long-range dependencies. In addition to above architectural advances, there has also been works [21, 35, 12] on improving over-parameterized deep neural networks by trading accuracy for efficiency. MobileNets [42, 17] and EfficientNets [50] leveraged neural architecture search (NAS) to design efficient mobile-size network and achieved new state-of-the-art results.

**Vision transformers.** Since transformers achieved remarkable success in natural language processing (NLP) [52, 9], many attempts [10, 51, 62, 54, 57, 6, 13, 61, 32] have been made to introduce transformer-like architectures to vision tasks. The pioneering work ViT [10] directly applied the transformer architecture inherited from NLP to image classification with image patches as input. While ViT required a large private labelled image dataset JFT-300M [46] to achieve promising result, DeiT [51] introduced a new training paradigm to extend ViT to a data-efficient vision transformer directly trained on ImageNet-1K. T2T-ViT [62] proposed a T2T transformation to embed visual tokens by recursively aggregating neighboring tokens into one token. TNT [62] proposed to model both patch-level and pixel-level representation by the inner and outer transformer block, respectively. PVT [54] introduced the pyramid structure into ViT, which can generate multi-scale feature maps for various pixel-level dense prediction tasks. CPVT [6] and CvT [57] are the most related to our work which leverage a convolutional projection into conventional transformer block, but we carefully investigate how to maximize the advantage of utilizing both CNNs and transformers by studying the different components including shortcut and normalization functions and successfully obtain a more superior result. Besides, transformers are also used to solve other vision tasks such as object detection [2, 64], semantic segmentation [63], image retrieval [11], and low-level vision task [3].

Although there are many works successfully applying transformers for vision tasks, they have not shown satisfactory results compared to conventional CNNs, which are still the primary architectures for vision applications. Transformers are especially good at modelling long-range dependencies necessary for downstream vision tasks. However, locality should also be maintained for visual perception. In this paper, we demonstrate the potential of combining the transformer based network together with convolutional layer, the overall architecture follows the elaborated prior convolutional neural networks such as ResNet [15] and EfficientNet [50].

### 3 Approach

#### 3.1 Overall Architecture

Our intention is to build a hybrid network taking the advantages of both convolutional neural networks and transformers. An overview of ResNet-50 [15], DeiT [51], and the proposed small version (CMT-S) of CMT architectures are presented in Figure 2. As shown in Figure 2(b), DeiT directly splits an input image into non-overlapping patches, the in-patch structure information can only be poorly modeled with linear projections. To overcome this limitation, we utilize the stem architecture which has a  $3 \times 3$  convolution with a stride of 2 and an output channel of 32 to reduce the size of input images, followed by another two  $3 \times 3$  convolutions with stride 1 for better local information

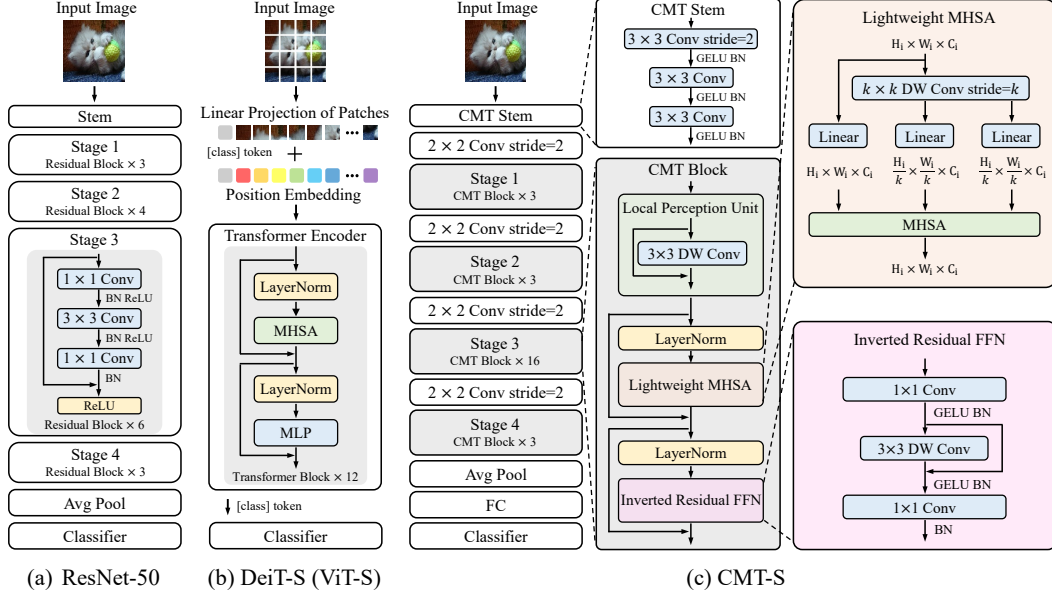


Figure 2: **Example CMT architecture for ImageNet.** (a) ResNet-50 model [15] as a reference. (b) DeiT-S [51] (ViT-S [10]) architecture. MHSA denotes the multi-head self-attention module. MLP is the multi-layer perceptron. (c) Our proposed CMT-S, described in Sec. 3. Table 1 shows more details and other variants.

extraction. Following the design in modern CNNs (*e.g.*, ResNet [15]), our model has four stages to generate feature maps of different scales which are important for dense prediction tasks. To produce the hierarchical representation, a patch aggregation layer consisting of a convolution and a layer normalization (LN) [1] is applied before each stage to reduce the size of intermediate feature (2x downsampling of resolution) and project it to a larger dimension (2x enlargement of dimension). In each stage, several CMT blocks are stacked sequentially for feature transformation while retaining the same resolution of the input. For example, the “Stage 3” of CMT-S contains 16 CMT blocks as illustrated in Figure 2(c). The CMT block is able to capture both local and long-range dependencies, and we will describe it in details in Sec. 3.2. The model ends with a global average pooling layer, a fully-connected (FC) layer and a 1000-way classification layer with softmax.

Given an input image, we can obtain four hierarchical feature maps with different resolutions, similar to those of typical convolutional neural networks, *e.g.*, ResNet [15] and EfficientNet [50]. With the above feature maps whose strides are 4, 8, 16, and 32 with respect to the input, our CMT can obtain multi-scale representations of input images and can be easily applied to most downstream tasks such as detection and segmentation.

### 3.2 CMT Block

The proposed CMT block consists of a local perception unit (LPU), a lightweight multi-head self-attention (LMHSA) module, and an inverted residual feed-forward network (IRFFN), as illustrated in Figure 2(c). We describe these three parts in the following.

**Local Perception Unit.** Rotation and shift are two commonly used data augmentation manners in vision tasks, and these operations should not alter the final results of the model. In other words, we expect translation-invariance [24] in those tasks. However, the absolute positional encoding used in previous transformers, initially designed to leverage the order of tokens, damages such invariance because it adds unique positional encoding to each patch [6]. Besides, vision transformers ignore the local relation [34] and the structure information [23] inside the patch. To alleviate the limitations, we propose the local perception unit (LPU) to extract local information, which is defined as:

$$\text{LPU}(\mathbf{X}) = \text{DWConv}(\mathbf{X}) + \mathbf{X}. \quad (1)$$

where  $\mathbf{X} \in \mathbb{R}^{H \times W \times d}$ ,  $H \times W$  is the resolution of the input of current stage,  $d$  indicates the dimension of features.  $\text{DWConv}(\cdot)$  denotes the depth-wise convolution.

**Lightweight Multi-head Self-attention.** In original self-attention module, the input  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is linearly transformed into query  $\mathbf{Q} \in \mathbb{R}^{n \times d_k}$ , key  $\mathbf{K} \in \mathbb{R}^{n \times d_k}$ , and value  $\mathbf{V} \in \mathbb{R}^{n \times d_v}$ , where  $n = H \times W$  is the number of patches. And we omit the reshape operation from  $H \times W \times d$  to  $n \times d$  of input tensor for simplicity. The notation  $d$ ,  $d_k$  and  $d_v$  are the dimensions of input, key (query) and value, respectively. Then the self-attention module is applied as:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}. \quad (2)$$

To mitigate the computation overhead, we use a  $k \times k$  depth-wise convolution with stride  $k$  to reduce the spatial size of  $\mathbf{K}$  and  $\mathbf{V}$  before the attention operation, *i.e.*,  $\mathbf{K}' = \text{DWConv}(\mathbf{K}) \in \mathbb{R}^{\frac{n}{k^2} \times d_k}$  and  $\mathbf{V}' = \text{DWConv}(\mathbf{V}) \in \mathbb{R}^{\frac{n}{k^2} \times d_v}$  as shown in Figure 2(c). In addition, we follow [43, 32] by adding a relative position bias  $\mathbf{B}$  to each self-attention module:

$$\text{LightweightAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}'^T}{\sqrt{d_k}} + \mathbf{B}\right)\mathbf{V}'. \quad (3)$$

where  $\mathbf{B} \in \mathbb{R}^{n \times \frac{n}{k^2}}$  is randomly initialized and learnable. The learnt relative position bias can also be easily transferred to  $\mathbf{B}' \in \mathbb{R}^{m_1 \times m_2}$  with a different size  $m_1 \times m_2$  through bicubic interpolation, *i.e.*,  $\mathbf{B}' = \text{Bicubic}(\mathbf{B})$ . Thus it is convenient to fine-tune the proposed CMT for other downstream vision tasks. Finally, the lightweight multi-head self-attention (LMHSA) module is defined by considering  $h$  “heads”, *i.e.*,  $h$  LightweightAttention functions are applied to the input. Each head outputs a sequence of size  $n \times \frac{d}{h}$ . These  $h$  sequences are then concatenated into a  $n \times d$  sequence.

**Inverted Residual Feed-forward Network.** The original FFN proposed in ViT [10] is composed of two linear layers separated by a GELU activation [16]. The first layer expands the dimension by a factor of 4, and the second layer reduces the dimension by the same ratio:

$$\text{FFN}(\mathbf{X}) = \text{GELU}(\mathbf{X}\mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2. \quad (4)$$

where  $\mathbf{W}_1 \in \mathbb{R}^{d \times 4d}$  and  $\mathbf{W}_2 \in \mathbb{R}^{4d \times d}$  indicate weights of the two linear layers, respectively. The notation  $b_1$  and  $b_2$  are the bias terms. Figure 2(c) provides a schematic visualization of our design. The proposed inverted residual feed-forward network (IRFFN) appears similar to inverted residual block [42] consisting of an expansion layer followed by a depth-wise convolution and a projection layer. Specifically, we change the location of shortcut connection for better performance:

$$\begin{aligned} \text{IRFFN}(\mathbf{X}) &= \text{Conv}(\mathcal{F}(\text{Conv}(\mathbf{X}))), \\ \mathcal{F}(\mathbf{X}) &= \text{DWConv}(\mathbf{X}) + \mathbf{X}. \end{aligned} \quad (5)$$

where the activation layer is omitted. We also include the batch normalization after the activation layer and the last linear layer according to [42]. The depth-wise convolution is used to extract local information with negligible extra computational cost. The motivation for inserting shortcut is similar to that of classic residual networks, which can promote the propagation ability of gradient across layers. We show that such shortcut helps the network achieve better results in our experiments.

With the aforementioned three components, the CMT block can be formulated as:

$$\mathbf{X}'_i = \text{LPU}(\mathbf{X}_{i-1}), \quad (6)$$

$$\mathbf{X}''_i = \text{LMHSA}(\text{LN}(\mathbf{X}'_i)) + \mathbf{X}'_i, \quad (7)$$

$$\mathbf{X}_i = \text{IRFFN}(\text{LN}(\mathbf{X}''_i)) + \mathbf{X}''_i. \quad (8)$$

where  $\mathbf{X}'_i$  and  $\mathbf{X}''_i$  denote the output features of the LPU and the LMHSA module for block  $i$ , respectively. The term LN denotes the layer normalization [1]. We stack several CMT blocks in each stage for feature transformation and aggregation.

### 3.3 Complexity Analysis

We analyze the computational cost between standard ViT and our CMT in this section. A standard transformer block consists of the MHSA and FFN. Supposing the input feature is of size  $n \times d$ , the computational complexity (FLOPs) can be calculated as:

$$\mathcal{O}(\text{MHSA}) = 2nd(d_k + d_v) + n^2(d_k + d_v), \quad (9)$$

$$\mathcal{O}(\text{FFN}) = 2nd^2r, \quad (10)$$



Table 1: **Architectures for ImageNet classification.** The output size is for CMT-S. Convolution and CMT blocks are shown in brackets (see also Figure 2(c)), with the numbers of blocks stacked.  $H_i$  and  $k_i$  are the number of heads and reduction rates in LMHSA of stage  $i$ , respectively.  $R_i$  denotes the expansion ratio in IRFFN of stage  $i$ . Patch Aggr. is the short for patch aggregation layer.

Output Size	Layer Name	CMT-Ti	CMT-XS	CMT-S	CMT-B
$112 \times 112$	Stem	$3 \times 3, 16, \text{stride } 2$ $[3 \times 3, 16] \times 2$	$3 \times 3, 16, \text{stride } 2$ $[3 \times 3, 16] \times 2$	$3 \times 3, 32, \text{stride } 2$ $[3 \times 3, 32] \times 2$	$3 \times 3, 38, \text{stride } 2$ $[3 \times 3, 38] \times 2$
$56 \times 56$	Patch Aggr.	$2 \times 2, 46, \text{stride } 2$	$2 \times 2, 52, \text{stride } 2$	$2 \times 2, 64, \text{stride } 2$	$2 \times 2, 76, \text{stride } 2$
Stage 1	LPU LMHSA IRFFN	$\begin{bmatrix} 3 \times 3, 46 \\ H_1=1, k_1=8 \\ R_1=3.6 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 52 \\ H_1=1, k_1=8 \\ R_1=3.8 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 64 \\ H_1=1, k_1=8 \\ R_1=4 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 76 \\ H_1=1, k_1=8 \\ R_1=4 \end{bmatrix} \times 4$
$28 \times 28$	Patch Aggr.	$2 \times 2, 92, \text{stride } 2$	$2 \times 2, 104, \text{stride } 2$	$2 \times 2, 128, \text{stride } 2$	$2 \times 2, 152, \text{stride } 2$
Stage 2	LPU LMHSA IRFFN	$\begin{bmatrix} 3 \times 3, 92 \\ H_2=2, k_2=4 \\ R_2=3.6 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 104 \\ H_2=2, k_2=4 \\ R_2=3.8 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 128 \\ H_2=2, k_2=4 \\ R_2=4 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 152 \\ H_2=2, k_2=4 \\ R_2=4 \end{bmatrix} \times 4$
$14 \times 14$	Patch Aggr.	$2 \times 2, 184, \text{stride } 2$	$2 \times 2, 208, \text{stride } 2$	$2 \times 2, 256, \text{stride } 2$	$2 \times 2, 304, \text{stride } 2$
Stage 3	LPU LMHSA IRFFN	$\begin{bmatrix} 3 \times 3, 184 \\ H_3=4, k_3=2 \\ R_3=3.6 \end{bmatrix} \times 10$	$\begin{bmatrix} 3 \times 3, 208 \\ H_3=4, k_3=2 \\ R_3=3.8 \end{bmatrix} \times 12$	$\begin{bmatrix} 3 \times 3, 256 \\ H_3=4, k_3=2 \\ R_3=4 \end{bmatrix} \times 16$	$\begin{bmatrix} 3 \times 3, 304 \\ H_3=4, k_3=2 \\ R_3=4 \end{bmatrix} \times 20$
$7 \times 7$	Patch Aggr.	$2 \times 2, 368, \text{stride } 2$	$2 \times 2, 416, \text{stride } 2$	$2 \times 2, 512, \text{stride } 2$	$2 \times 2, 608, \text{stride } 2$
Stage 4	LPU LMHSA IRFFN	$\begin{bmatrix} 3 \times 3, 368 \\ H_4=8, k_4=1 \\ R_4=3.6 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 416 \\ H_4=8, k_4=1 \\ R_4=3.8 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 512 \\ H_4=8, k_4=1 \\ R_4=4 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 608 \\ H_4=8, k_4=1 \\ R_4=4 \end{bmatrix} \times 4$
$1 \times 1$	FC	$1 \times 1, 1280$			
$1 \times 1$	Classifier	$1 \times 1, 1000$			
# Params		9.49 M	15.24 M	25.14 M	45.72 M
# FLOPs		0.64 B	1.54 B	4.04 B	9.33 B

where  $r$  is the expansion ratio of FFN,  $d_k$  and  $d_v$  are dimensions of key and value, respectively. More specifically, ViT sets  $d = d_k = d_v$  and  $r = 4$ , the cost can be simplified as:

$$\mathcal{O}(\text{Transformer block}) = \mathcal{O}(\text{MHSA}) + \mathcal{O}(\text{FFN}) = 12nd^2 + 2n^2d \quad (11)$$

Under the setting of ViT, the FLOPs of CMT block can be formulated as:

$$\mathcal{O}(\text{LPU}) = 9nd, \quad (12)$$

$$\mathcal{O}(\text{LMHSA}) = 2nd^2(1 + 1/k^2) + 2n^2d/k^2, \quad (13)$$

$$\mathcal{O}(\text{IRFFN}) = 8nd^2 + 36nd, \quad (14)$$

$$\mathcal{O}(\text{CMT block}) = \mathcal{O}(\text{LPU}) + \mathcal{O}(\text{LMHSA}) + \mathcal{O}(\text{IRFFN}) \quad (15)$$

$$= 10nd^2(1 + 0.2/k^2) + 2n^2d/k^2 + 45nd, \quad (16)$$

where  $k \geq 1$  is the reduction ratio in LMHSA. Compared to standard transformer block, the CMT block is more friendly to computational cost, and is easier to process the feature map under higher resolution (larger  $n$ ).

### 3.4 Scaling Strategy

Inspired by [50], we propose a new compound scaling strategy suitable for transformer-based networks, which uses a compound coefficient  $\phi$  to uniformly scale the number of layers (depth), dimensions, and input resolution in a principled way:

$$\text{depth} : \alpha^\phi, \quad \text{dimension} : \beta^\phi, \quad \text{resolution} : \gamma^\phi, \quad (17)$$

$$\text{s.t. } \alpha \cdot \beta^{1.5} \cdot \gamma^2 \approx 2.5, \quad \alpha \geq 1, \beta \geq 1, \gamma \geq 1 \quad (18)$$

where  $\alpha$ ,  $\beta$ , and  $\gamma$  are constants determined by grid search to decide how to assign resources to network depth, dimension and input resolution, respectively. Intuitively,  $\phi$  is the coefficient that controls how many more ( $\phi \geq 1$ ) or less ( $\phi \leq -1$ ) resources are available for model scaling. Notably, the FLOPs of the proposed CMT block is approximately proportional<sup>2</sup> to  $\alpha$ ,  $\beta^{1.5}$ , and  $\gamma^2$  according to E.q. 16. And we constraint  $\alpha \cdot \beta^{1.5} \cdot \gamma^2 \approx 2.5$  so that for a given new  $\phi$ , the total FLOPs will approximately increase by  $2.5^\phi$ . This will strike a balance between the increase of computational cost and performance gain. In our experiments, we empirically set  $\alpha=1.2$ ,  $\beta=1.3$ , and  $\gamma=1.15$ .

<sup>2</sup>The precious proportion is associated with  $n$  and  $d$ . For example, CMT-S has  $n=3136 \gg d=64$  in ‘‘stage 1’’ and  $n=49 \ll d=512$  in ‘‘stage 4’’. We find the above proportion can already generate good variants for CMT.

Table 2: **ImageNet Results of CMT.** CNNs and transformers with similar accuracy are grouped together for comparison. The proposed CMTs consistently outperform other methods with less computational cost.

Model	Top-1 Acc.	Top-5 Acc.	# Params	Resolution	# FLOPs	Ratio
CPVT-Ti-GAP [6]	74.9%	-	6M	224 <sup>2</sup>	1.3B	2.2×
DenseNet-169 [20]	76.2%	93.2%	14M	224 <sup>2</sup>	3.5B	5.8×
EfficientNet-B1 [50]	79.1%	94.4%	7.8M	240 <sup>2</sup>	0.7B	1.2×
<b>CMT-Ti</b>	<b>79.2%</b>	<b>94.6%</b>	9.5M	160 <sup>2</sup>	<b>0.6B</b>	<b>1×</b>
ResNet-50 [15]	76.2%	92.9%	25.6M	224 <sup>2</sup>	4.1B	2.7×
Coat-Lite Mini [59]	78.9%	-	11M	224 <sup>2</sup>	2.0B	1.3×
DeiT-S [51]	79.8%	-	22M	224 <sup>2</sup>	4.6B	3.1×
EfficientNet-B3 [50]	81.6%	95.7%	12M	300 <sup>2</sup>	1.8B	1.2×
<b>CMT-XS</b>	<b>81.8%</b>	<b>95.8%</b>	15.2M	192 <sup>2</sup>	<b>1.5B</b>	<b>1×</b>
ResNeXt-101-64x4d [58]	80.9%	95.6%	84M	224 <sup>2</sup>	32B	8×
T2T-ViT-19 [62]	81.2%	-	39.0	224 <sup>2</sup>	8.0B	2×
PVT-M [54]	81.2%	-	44.2M	224 <sup>2</sup>	6.7B	1.7×
Swin-T [32]	81.3%	-	29M	224 <sup>2</sup>	4.5B	1.1×
CPVT-S-GAP [6]	81.5%	-	23M	224 <sup>2</sup>	4.6B	1.2×
RegNetY-8GF [40]	81.7%	-	39.2M	224 <sup>2</sup>	8.0B	2×
CeiT-S [61]	82.0%	95.9%	24.2M	224 <sup>2</sup>	4.5B	1.1×
CvT-13-NAS [57]	82.2%	-	18M	224 <sup>2</sup>	4.1B	1×
EfficientNet-B4 [50]	82.9%	96.4%	19M	380 <sup>2</sup>	4.2B	1×
Twins-SVT-B [5]	83.1%	-	56.0M	224 <sup>2</sup>	8.3B	2.1×
<b>CMT-S</b>	<b>83.5%</b>	<b>96.6%</b>	25.1M	224 <sup>2</sup>	<b>4.0B</b>	<b>1×</b>
ViT-B/16 <sub>↑384</sub> [10]	77.9%	-	55.5M	384 <sup>2</sup>	77.9B	8.4×
T2T-ViT-24 [62]	82.2%	-	63.9M	224 <sup>2</sup>	12.6B	1.4×
CPVT-B [6]	82.3%	-	88M	224 <sup>2</sup>	17.6B	1.9×
TNT-B [13]	82.8%	96.3%	65.6M	224 <sup>2</sup>	14.1B	1.5×
DeiT-B <sub>↑384</sub> [51]	83.1%	-	85.8M	384 <sup>2</sup>	55.6B	6.0×
CvT-21 <sub>↑384</sub> [57]	83.3%	-	31.5M	384 <sup>2</sup>	24.9B	2.7×
Swin-B [32]	83.3%	-	88M	224 <sup>2</sup>	15.4B	1.7×
Twins-SVT-L [5]	83.3%	-	99.2M	224 <sup>2</sup>	14.8B	1.6×
CeiT-S <sub>↑384</sub> [61]	83.3%	96.5%	24.2M	384 <sup>2</sup>	12.9B	1.4×
BoTNet-S1-128 [45]	83.5%	96.5%	75.1M	256 <sup>2</sup>	19.3B	2.1×
EfficientNet-B6 [50]	84.0%	96.8%	43M	528 <sup>2</sup>	19.2B	2.0×
<b>CMT-B</b>	<b>84.5%</b>	<b>96.9%</b>	45.7M	256 <sup>2</sup>	<b>9.3B</b>	<b>1×</b>
EfficientNet-B7 [50]	84.3%	97.0%	66M	600 <sup>2</sup>	37B	1.9×
<b>CMT-L</b>	<b>84.8%</b>	<b>97.1%</b>	74.7M	288 <sup>2</sup>	<b>19.5B</b>	<b>1×</b>

We build our model CMT-S to have similar model size and computation complexity with DeiT-S (ViT-S) and EfficientNet-B4. We also introduce CMT-Ti, CMT-XS, CMT-B and CMT-L according to the proposed scaling strategy. The input resolutions are 160<sup>2</sup>, 192<sup>2</sup>, 224<sup>2</sup>, 256<sup>2</sup> and 288<sup>2</sup> for all four models, respectively. The detailed architecture hyper-parameters are shown in Table 1.

## 4 Experiments

In this section, we investigate the effectiveness of CMT architecture by conducting experiments on several tasks including classification, object detection, and instance segmentation. We first compare the proposed CMT with previous state-of-the-art models on above tasks, and then we ablate the important design elements of CMT.

### 4.1 ImageNet Classification

**Experimental Settings.** ImageNet [8] is a image classification benchmark which contains 1.28M training images and 50K validation images of 1000 classes. For fair comparisons with recent works, we adopt the same training and augmentation strategy as that in DeiT [51], *i.e.*, models are trained

for 300 epochs using the AdamW [33] optimizer. All models are trained on 8 NVIDIA Tesla V100 GPUs.

**Results of CMT.** Table 2 shows the performance of the proposed CMTs that are scaled from the CMT-S according to E.q. 18. Our models achieve better accuracy with fewer parameters and FLOPs compared to other convolution-based and transformer-based counterparts. In particular, our CMT-S achieves 83.5% top-1 accuracy with 4.0B FLOPs, which is 3.7% higher than the baseline model DeiT-S [51] and 2.0% higher than CPVT [6], indicating the benefit of CMT block for capturing both local and global information. Note that all previous transformer-based models are still inferior to EfficientNet [50] which is obtained via a thorough architecture search, however, our CMT-S is 0.6% higher than EfficientNet-B4 with less computational cost, which demonstrates the efficacy of the proposed hybrid structure and show strong potential for further improvement. We also plot the accuracy-FLOPs curve in Figure 1(a) to have an intuitive comparison between these models. We can see that CMTs consistently outperform other models by a large margin.

## 4.2 Ablation Study

**Stage-wise architecture.** Transformer-based ViT/DeiT can only generate single-scale feature map, losing a lot of multi-scale information crucial for dense prediction tasks. We change the columnar DeiT-S to DeiT-S-4Stage, which has 4 stages like CMT-S in Table 1, but maintain the original FFN. We also change MHSA to LMHSA to reduce computational cost. As shown in Table 3, DeiT-S-4Stage outperforms DeiT-S by 1.6% with less FLOPs, demonstrating that the widely-adopted stage-wise design in CNNs is a better choice for promoting transformer-based architecture.

**CMT block.** Ablations on different modules in CMT are shown in Table 4. DeiT-S-4Stage has 4 patch aggregation layers (the first is a  $4 \times 4$  convolution with stride 4). “+ Stem” indicates that we

add the CMT stem into the network and replace the first patch aggregation layer with a  $2 \times 2$  convolution with stride 2. The improvement shows the benefit of the convolution-based stem. Besides, the proposed LPU and IRFFN can further boost the network by 0.8% and 0.6%, respectively. It is worth noticing that the shortcut connections in LPU and IRFFN are also crucial for the final performance.

**Normalization function.** Transformer-based models usually use layer normalization (LN) [1] inherited from NLP tasks. However, convolution-based models usually utilize batch normalization (BN) [22] to stabilize the training. CMT maintains the LN before LMHSA and IRFFN, and inserts BN after the convolutional layer. If all LNs are replaced by BNs, the model can not converge during training. If all BNs are replaced by LNs, the performance of CMT-S drops to 82.8%, indicating that proper application of different normalization functions can improve the performance of the network.

**Scaling strategy.** Table 5 shows the ImageNet results of CMT architecture under different scaling strategies. Unidimensional scaling strategies are significantly inferior to the proposed compound scaling strategy, especially for depth-only scaling strategy which leads to an even worse result of 83.3% against 83.5% of the original CMT-S, when the network is scaled up.

Table 3: Ablation study about stage-wise architecture on ImageNet.

Model	Params	FLOPs	Top-1
DeiT-S	22M	4.6B	79.8%
DeiT-S-4Stage	25M	3.7B	81.4%

Table 4: Ablation study of CMT block.

Model	Params	FLOPs	Top-1
DeiT-S-4Stage	25M	3.7B	81.4%
+ Stem	25M	3.9B	81.9%
+ LPU	25M	3.9B	82.7%
w/o shortcut	25M	3.9B	82.0%
+ IRFFN	25M	3.9B	83.3%
w/o shortcut	25M	3.9B	82.5%
+ FC (CMT-S)	25M	4.0B	83.5%

Table 5: Ablation study of scaling strategy.

Model (based on CMT-S)	FLOPs	Top-1	FLOPs	Top-1
Scale: $\alpha=2.2$ (depth only)	1.7B ( $\phi=-1$ )	80.8%	8.6B ( $\phi=1$ )	83.3%
Scale: $\beta=1.6$ (dimension only)	1.7B ( $\phi=-1$ )	81.3%	10B ( $\phi=1$ )	83.7%
Scale: $\alpha=1.2, \beta=1.3, \gamma=1.15$	1.5B ( $\phi=-1$ )	81.8%	9.3B ( $\phi=1$ )	84.5%

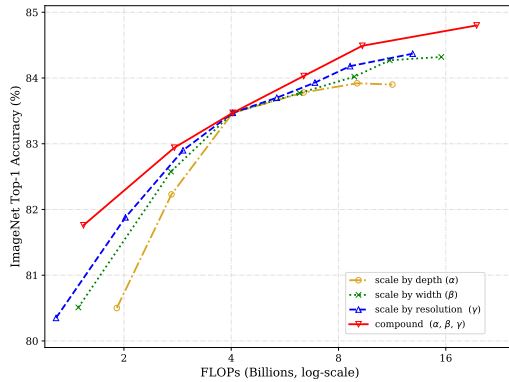


Figure 3: Results of different scaling strategies.



Table 6: **Transfer Learning Results of CMT**. All results are fine-tuned with the ImageNet pretrained checkpoint. † means the transfer results are from [25].

Model	# Params	# FLOPs	CIFAR10	CIFAR100	Cars	Flowers	Pets
ResNet-152† [15]	58.1M	11.3B	97.9%	87.6%	92.0%	97.4%	94.5%
Inception-v4† [47]	41.1M	16.1B	97.9%	87.5%	93.3%	98.5%	93.7%
EfficientNet-B7†600 [50]	64.0M	37.2B	98.9%	<b>91.7%</b>	<b>94.7%</b>	98.8%	<b>95.4%</b>
ViT-B/16†384 [10]	85.8M	17.6B	98.1%	87.1%	-	89.5%	93.8%
DeiT-B [51]	85.8M	17.6B	99.1%	90.8%	92.1%	98.4%	-
CeiT-S†384 [61]	24.2M	12.9B	99.1%	90.8%	94.1%	98.6%	94.9%
TNT-S†384 [13]	<b>23.8M</b>	17.3B	98.7%	90.1%	-	<b>98.8%</b>	94.7%
<b>CMT-S (ours)</b>	25.1M	<b>4.04B</b>	<b>99.2%</b>	<b>91.7%</b>	94.4%	98.7%	95.2%

### 4.3 Transfer Learning

In this section, we transfer CMT-S trained on ImageNet to other downstream vision tasks to further demonstrate the generalization ability of CMT.

#### 4.3.1 Object Detection and Instance Segmentation

**Experimental Settings.** The experiments are conducted on COCO [31], which contains 118K training images and 5K validation images of 80 classes. We evaluate the proposed CMT-S using two typical framework, RetinaNet [30] and Mask R-CNN [14], for object detection and instance segmentation, respectively. Specifically, we replace the original backbone with our CMT-S to build new detectors. All models are trained under standard “1x” schedule (12 epochs) and tested with single-scale following [54].

**Results of CMT.** We report the performance comparison results of object detection task and instance segmentation task in Table 7 and Table 8 respectively. For object detection with RetinaNet as basic framework, CMT-S outperforms Twins-PCPVT-S [5] with 1.3% mAP and Twins-SVT-S [5] with 2.0% mAP. For instance segmentation with Mask R-CNN as basic framework, CMT-S surpasses Twins-PCPVT-S [5] with 1.7% AP and Twins-SVT-S [5] with 1.9% AP.

#### 4.3.2 Other Vision Tasks

We also evaluate the proposed CMT on five commonly used transfer learning datasets, including CIFAR10 [27], CIFAR100 [27], Stanford Cars [26], Flowers [36], and Oxford-IIIT Pets [38] (see Appendix for more details). We fine-tune the ImageNet pretrained models on new datasets following [50, 13]. Table 6 shows the corresponding results. CMT-S outperforms other transformer-based models in all datasets with less FLOPs, and achieves comparable performance against EfficientNet-B7 [50] with 9x less FLOPs, which demonstrates the superiority of CMT architecture.

Table 7: **Object detection results on COCO val2017**. All models use RetinaNet as basic framework and are trained in “1x” schedule. FLOPs are calculated on 1280×800 input. † means the results are from [5].

Backbone	# Params	# FLOPs	mAP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
ConT-M [60]	217B	27.0M	37.9	58.1	40.2	23.0	40.6	50.4
ResNet-101 [15]	315B	56.7M	38.5	57.6	41.0	21.7	42.8	50.4
RelationNet++ [4]	266B	39.0M	39.4	58.2	42.5	-	-	-
ResNeXt-101-32x4d [58]	319B	56.4M	39.9	59.6	42.7	22.3	44.2	52.5
PVT-S [54]	226B	34.2M	40.4	61.3	43.0	25.0	42.9	55.7
Swin-T† [32]	245B	38.5M	41.5	62.1	44.2	25.1	44.9	55.5
Twins-SVT-S [5]	209B	34.3M	42.3	63.4	45.2	26.0	45.5	56.5
Twins-PCPVT-S [5]	226B	34.4M	43.0	64.1	46.0	27.5	46.3	57.3
<b>CMT-S (ours)</b>	231B	44.3M	<b>44.3</b>	<b>65.5</b>	<b>47.5</b>	<b>27.1</b>	<b>48.3</b>	<b>59.1</b>

## 5 Conclusion

This paper proposes a novel hybrid architecture named CMT for visual recognition and other downstream vision tasks, to address the limitations of utilizing transformers in a brutal force manner in the field of computer vision. The proposed CMT takes advantage of both CNNs and transformers

Table 8: **Instance segmentation results on COCO val2017**. All models use Mask R-CNN as basic framework and are trained in “1x” schedule. FLOPs are calculated on 1280×800 input. † means the results are from [5].

Backbone	# Params	# FLOPs	AP <sup>box</sup>	AP <sub>50</sub> <sup>box</sup>	AP <sub>75</sub> <sup>box</sup>	AP <sup>mask</sup>	AP <sub>50</sub> <sup>mask</sup>	AP <sub>75</sub> <sup>mask</sup>
ResNet-101 [15]	336B	63.2M	40.0	60.5	44.0	36.1	57.5	38.6
PVT-S [54]	245B	44.1M	40.4	62.9	43.8	37.8	60.1	40.3
ResNeXt-101-32x4d [58]	340B	62.8M	41.9	62.5	45.9	37.5	59.4	40.2
Swin-T† [32]	264B	47.8M	42.2	64.6	46.2	39.1	61.6	42.0
Twins-SVT-S [5]	228B	44.0M	42.7	65.6	46.7	39.6	62.5	42.6
Twins-PCPVT-S [5]	245B	44.3M	42.9	65.8	47.1	40.0	62.7	42.9
<b>CMT-S (ours)</b>	<b>249B</b>	<b>44.5M</b>	<b>44.6</b>	<b>66.8</b>	<b>48.9</b>	<b>40.7</b>	<b>63.9</b>	<b>43.4</b>

to capture local and global information, promoting the representation ability of the network. In addition, a scaling strategy is proposed to generate a family of CMT variants for different resource constraints. Extensive experiments on ImageNet and other downstream vision tasks demonstrate the effectiveness and superiority of the proposed CMT architecture.

## Appendix

In this supplementary material, we first compare the inference speed of our proposed CMT with other state-of-the-art networks. Then we list the training strategy for transfer learning in details. Finally we show more experiment results under different settings for object detection and instance segmentation on COCO benchmark.

### A Inference Speed

We evaluate the inference speed (throughput, images processed per second) of our proposed CMT-S and CMT-B on ImageNet as shown in Table 9. Noting that while EfficientNet is searched via NAS method, our CMT has stronger potential and better speed-accuracy trade-off. The proposed method also outperforms other transformer-based networks, demonstrating that combining convolution and transformer to capture both local and global information can produce impressive results.

Table 9: **Comparison of the inference speed between different models**. Throughput is measured on a single V100 GPU, following [32, 5].

Model	# FLOPs	Throughput (image/s)	ImageNet Top-1 Acc.	Model	# FLOPs	Throughput (image/s)	ImageNet Top-1 Acc.
DeiT-S/16 [51]	4.6B	940.4	79.8%	DeiT-B/16 [51]	17.6B	292.3	81.8%
RegNetY-4G [40]	4.0B	1156.7	80.0%	RegNetY-16G [40]	16.0B	334.7	82.9%
PVT-M [54]	6.7B	528.1	81.2%	PVT-L [54]	9.8B	358.8	81.7%
CPVT-S-GAP [6]	4.6B	942.3	81.5%	CPVT-B [6]	17.6B	285.5	82.3%
Swin-S [32]	8.7B	436.9	83.0%	Swin-B [32]	15.4B	278.1	83.3%
Twins-SVT-B [5]	8.3B	469.0	83.2%	Twins-SVT-L [5]	14.8B	288.0	83.7%
EfficientNet-B4 [50]	4.2B	349.4	82.9%	EfficientNet-B6 [50]	19.0B	96.9	84.0%
CMT-S (ours)	4.0B	562.5	83.5%	CMT-B (ours)	9.3B	285.4	84.5%

## B Transfer Learning

### B.1 Image Classification

**Datasets.** In addition to ImageNet, we also evaluate the proposed CMT on five commonly used transfer learning datasets, including CIFAR10 [27], CIFAR100 [27], Stanford Cars [26], Flowers [36], and Oxford-IIIT Pets [38]. The details of these datasets are listed in Table 10.

**Training details.** We describe our training strategy on the transfer learning datasets here. We build upon PyTorch [39], and adopt the same data augmentation strategy as that of ImageNet. We change the number of output units in the last classification layer to the number of classes in the target dataset and initialize the new classification layer randomly. The proposed CMT-S are fine-tuned with the

Table 10: **Datasets used for vision tasks.**

Dataset	Train size	Test size	# Classes
ImageNet [8]	1,281,167	50,000	1000
CIFAR10 [27]	50,000	10,000	10
CIFAR100 [27]	50,000	10,000	100
Stanford Cars [26]	8,144	8,041	196
Flowers [36]	2,040	6,149	102
Oxford-IIIT Pets [38]	3,680	3,669	37

Table 11: **More transfer learning results of CMT.** All results are fine-tuned with the ImageNet pretrained checkpoint. † means the transfer results are from [25].

Model	# Params	# FLOPs	CIFAR10	CIFAR100	Cars	Flowers	Pets
ResNet-152 <sup>†</sup> [15]	58.1M	11.3B	97.9%	87.6%	92.0%	97.4%	94.5%
Inception-v4 <sup>†</sup> [47]	41.1M	16.1B	97.9%	87.5%	93.3%	98.5%	93.7%
RegNetY-8GF [40]	39.2M	8.0B	-	-	94.0%	<b>99.0%</b>	-
CeiT-S [61]	<b>24.2M</b>	4.5B	99.0%	90.8%	93.2%	98.2%	94.6%
EfficientNet-B5 <sup>†</sup> <sub>456</sub> [50]	28.0M	9.9B	98.7%	<u>91.1%</u>	<u>93.9%</u>	98.5%	<u>94.9%</u>
<b>CMT-S (ours)</b>	<u>25.1M</u>	<b>4.0B</b>	<b>99.2%</b>	<b>91.7%</b>	<b>94.4%</b>	<u>98.7%</u>	<b>95.2%</b>
ViT-B/16 <sup>†</sup> <sub>384</sub> [10]	85.8M	17.6B	98.1%	87.1%	-	89.5%	93.8%
DeiT-B [51]	85.8M	17.6B	<u>99.1%</u>	90.8%	92.1%	98.4%	-
CeiT-S <sup>†</sup> <sub>384</sub> [61]	24.2M	12.9B	<u>99.1%</u>	90.8%	94.1%	98.6%	94.9%
TNT-S <sup>†</sup> <sub>384</sub> [13]	<b>23.8M</b>	17.3B	98.7%	90.1%	-	<u>98.8%</u>	94.7%
TNT-B <sup>†</sup> <sub>384</sub> [13]	65.6M	36.6B	<u>99.1%</u>	91.1%	-	<b>99.0%</b>	95.0%
EfficientNet-B7 <sup>†</sup> <sub>600</sub> [50]	64.0M	37.2B	98.9%	<u>91.7%</u>	<u>94.7%</u>	98.8%	95.4%
<b>CMT-B (ours)</b>	45.7M	<b>9.3B</b>	<b>99.3%</b>	<b>91.9%</b>	<b>94.9%</b>	<b>99.0%</b>	<b>95.5%</b>

image resolution of  $224 \times 224$  on all datasets. For **CIFAR10** and **CIFAR100**, the model is fine-tuned for 150 epochs with  $6e-5$  initial learning rate. For **Flowers** and **Pets**, the model is fine-tuned for 300 epochs with  $9e-5$  and  $6e-5$  initial learning rate, respectively. For **Cars**, the model is fine-tuned for 500 epochs with  $9e-5$  initial learning rate. Specifically, our proposed CMT achieves better result with smaller computational cost and less training epochs compared to EfficientNet [50].

**Results.** In addition to the CMT-S shown in our main paper, we present the transfer learning result of CMT-B in Table 11. We can find that CMT-B outperforms other previous models with less computational cost.

Table 12: **Object detection results on COCO val2017.** All models use RetinaNet as basic framework. “# P” means parameters, “# F” means FLOPs. FLOPs are calculated on  $1280 \times 800$  input. “1x” indicates 12 epochs, “3x” indicates 36 epochs, and “MS” indicates multi-scale training. † means the results are from [5].

Backbone	# P	# F	RetinaNet 1x						RetinaNet 3x + MS					
			mAP	AP <sub>50</sub>	AP <sub>75</sub>	mAP	AP <sub>50</sub>	AP <sub>75</sub>						
ConT-M [60]	217B	27.0M	37.9	58.1	40.2	23.0	40.6	50.4	-	-	-	-	-	-
ResNet-101 [15]	315B	56.7M	38.5	57.6	41.0	21.7	42.8	50.4	40.9	60.1	44.0	23.7	45.0	53.8
RelationNet++ [4]	266B	39.0M	39.4	58.2	42.5	-	-	-	-	-	-	-	-	-
ResNeXt-101-32x4d [58]	319B	56.4M	39.9	59.6	42.7	22.3	44.2	52.5	41.4	61.0	44.3	23.9	45.5	53.7
PVT-S [54]	226B	34.2M	40.4	61.3	43.0	25.0	42.9	55.7	42.2	62.7	45.0	26.2	45.2	57.2
Swin-T <sup>†</sup> [32]	245B	38.5M	41.5	62.1	44.2	25.1	44.9	55.5	43.9	64.8	47.1	28.4	47.2	57.8
Twins-SVT-S [5]	209B	34.3M	42.3	63.4	45.2	26.0	45.5	56.5	45.6	67.1	48.6	29.8	49.3	60.0
Twins-PCPVT-S [5]	226B	34.4M	43.0	64.1	46.0	<b>27.5</b>	46.3	57.3	45.2	66.5	48.6	30.0	48.8	58.9
<b>CMT-S (ours)</b>	231B	34.6M	<b>44.3</b>	<b>65.5</b>	<b>47.5</b>	27.1	<b>48.3</b>	<b>59.1</b>	<b>46.9</b>	<b>67.1</b>	<b>50.5</b>	<b>30.4</b>	<b>49.8</b>	<b>61.0</b>

## B.2 Object Detection

In addition to the “1x” setting presented in main paper, we also evaluate our CMT-S under the “3x” schedule. We follow the common multi-scale training strategy [54], *i.e.*, randomly resizing the input image so that its shorter side is between 640 and 800. The corresponding results are shown in Table 12.

### B.3 Instance Segmentation

Similar to object detection, we show the result of CMT-S based Mask R-CNN under “3x” schedule in Table 13. The proposed CMT architecture can surpass other transformer-based counterparts by a large margin with less FLOPs.

Table 13: **Instance segmentation results on COCO val2017.** All models use Mask R-CNN as basic framework. “# P” means parameters, “# F” means FLOPs. FLOPs are calculated on  $1280 \times 800$  input. “1x” indicates 12 epochs, “3x” indicates 36 epochs, and “MS” indicates multi-scale training. † means the results are from [5].

Backbone	# P	# F	Mask R-CNN 1x						Mask R-CNN 3x + MS					
			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>	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>
ResNet-101 [15]	336B	63.2M	40.0	60.5	44.0	36.1	57.5	38.6	42.8	63.2	47.1	38.5	60.1	41.3
PVT-S [54]	245B	44.1M	40.4	62.9	43.8	37.8	60.1	40.3	43.0	65.3	46.9	39.9	62.5	42.8
ConT-M [60]	237B	34.2M	40.5	-	-	38.1	-	-	-	-	-	-	-	-
ResNeXt-101-32x4d [58]	340B	62.8M	41.9	62.5	45.9	37.5	59.4	40.2	44.0	64.4	48.0	39.2	61.4	41.9
Swin-T† [32]	264B	47.8M	42.2	64.6	46.2	39.1	61.6	42.0	46.0	68.2	50.2	41.6	65.1	44.8
Twins-SVT-S [5]	228B	44.0M	42.7	65.6	46.7	39.6	62.5	42.6	46.8	69.2	51.2	42.6	66.3	45.8
Twins-PCPVT-S [5]	245B	44.3M	42.9	65.8	47.1	40.0	62.7	42.9	46.8	69.3	51.8	42.6	66.3	46.0
<b>CMT-S (ours)</b>	<b>249B</b>	<b>44.5M</b>	<b>44.6</b>	<b>66.8</b>	<b>48.9</b>	<b>40.7</b>	<b>63.9</b>	<b>43.4</b>	<b>48.3</b>	<b>70.4</b>	<b>52.3</b>	<b>43.7</b>	<b>67.7</b>	<b>47.1</b>

### References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, 2020.
- [3] Hanqing Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. *arXiv preprint arXiv:2012.00364*, 2020.
- [4] Cheng Chi, Fangyun Wei, and Han Hu. Relationnet++: Bridging visual representations for object detection via transformer decoder. *arXiv preprint arXiv:2010.15831*, 2020.
- [5] 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.
- [6] Xiangxiang Chu, Bo Zhang, Zhi Tian, Xiaolin Wei, and Huaxia Xia. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021.
- [7] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [8] 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.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *arXiv preprint arXiv:1810.04805*, 2018.
- [10] 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.
- [11] Alaaeldin El-Nouby, Natalia Neverova, Ivan Laptev, and Hervé Jégou. Training vision transformers for image retrieval. *arXiv preprint arXiv:2102.05644*, 2021.
- [12] Kai Han, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. Ghostnet: More features from cheap operations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [13] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *arXiv preprint arXiv:2103.00112*, 2021.
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, 2017.
- [15] 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*, 2016.
- [16] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [17] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019.
- [18] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Andrea Vedaldi. Gather-excite: Exploiting feature context in convolutional neural networks. *arXiv preprint arXiv:1810.12348*, 2018.

- [19] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- [20] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
- [21] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- [22] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, 2015.
- [23] Md Amirul Islam, Sen Jia, and Neil DB Bruce. How much position information do convolutional neural networks encode? *arXiv preprint arXiv:2001.08248*, 2020.
- [24] Osman Semih Kayhan and Jan C van Gemert. On translation invariance in cnns: Convolutional layers can exploit absolute spatial location. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [25] Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better imagenet models transfer better? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [26] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, 2013.
- [27] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [28] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 2012.
- [29] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- [30] 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.
- [31] 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*. Springer, 2014.
- [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.
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [34] David G Lowe. Object recognition from local scale-invariant features. In *Proceedings of the seventh IEEE international conference on computer vision*, 1999.
- [35] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European conference on computer vision (ECCV)*, 2018.
- [36] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, 2008.
- [37] Jongchan Park, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Bam: Bottleneck attention module. *arXiv preprint arXiv:1807.06514*, 2018.
- [38] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, 2012.
- [39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*, 2019.
- [40] 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.
- [41] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. *arXiv preprint arXiv:1906.05909*, 2019.
- [42] 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*, 2018.
- [43] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.
- [44] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [45] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. *arXiv preprint arXiv:2101.11605*, 2021.
- [46] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, 2017.
- [47] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2017.



- Intelligence*, 2017.
- [48] 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*, 2015.
  - [49] 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.
  - [50] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, 2019.
  - [51] 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.
  - [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *arXiv preprint arXiv:1706.03762*, 2017.
  - [53] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
  - [54] 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.
  - [55] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
  - [56] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, 2018.
  - [57] 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.
  - [58] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
  - [59] Weijian Xu, Yifan Xu, Tyler Chang, and Zhuowen Tu. Co-scale conv-attentional image transformers. *arXiv preprint arXiv:2104.06399*, 2021.
  - [60] Haotian Yan, Zhe Li, Weijian Li, Changhu Wang, Ming Wu, and Chuang Zhang. Contnet: Why not use convolution and transformer at the same time? *arXiv preprint arXiv:2104.13497*, 2021.
  - [61] 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.
  - [62] 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.
  - [63] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. *arXiv preprint arXiv:2012.15840*, 2020.
  - [64] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020.