

Vision Transformers with Hierarchical Attention

Yun Liu¹ Yu-Huan Wu² Guolei Sun³ Le Zhang⁴
Ajad Chhatkuli³ Luc Van Gool³

¹Institute for Infocomm Research (I2R), A*STAR, Singapore 138632, Singapore

²Institute of High Performance Computing (IHPC), A*STAR, Singapore 138632, Singapore

³Computer Vision Lab, ETH Zürich, Zürich 8092, Switzerland

⁴School of Information and Communication Engineering, University of Electronic Science and Technology of China (UESTC), Chengdu 611731, China

Abstract: This paper tackles the high computational/space complexity associated with multi-head self-attention (MHSA) in vanilla vision transformers. To this end, we propose hierarchical MHSA (H-MHSA), a novel approach that computes self-attention in a hierarchical fashion. Specifically, we first divide the input image into patches as commonly done, and each patch is viewed as a token. Then, the proposed H-MHSA learns token relationships within local patches, serving as local relationship modeling. Then, the small patches are merged into larger ones, and H-MHSA models the global dependencies for the small number of the merged tokens. At last, the local and global attentive features are aggregated to obtain features with powerful representation capacity. Since we only calculate attention for a limited number of tokens at each step, the computational load is reduced dramatically. Hence, H-MHSA can efficiently model global relationships among tokens without sacrificing fine-grained information. With the H-MHSA module incorporated, we build a family of hierarchical-attention-based transformer networks, namely HAT-Net. To demonstrate the superiority of HAT-Net in scene understanding, we conduct extensive experiments on fundamental vision tasks, including image classification, semantic segmentation, object detection and instance segmentation. Therefore, HAT-Net provides a new perspective for vision transformers. Code and pretrained models are available at <https://github.com/yun-liu/HAT-Net>.

Keywords: Vision transformer, hierarchical attention, global attention, local attention, scene understanding.

Citation: Y. Liu, Y. H. Wu, G. Sun, L. Zhang, A. Chhatkuli, L. V. Gool. Vision transformers with hierarchical attention. *Machine Intelligence Research*, vol.21, no.4, pp.670–683, 2024. <http://doi.org/10.1007/s11633-024-1393-8>

1 Introduction

In the last decade, convolutional neural networks (CNNs) have been the go-to architecture in computer vision, owing to their powerful capability in learning representations from images/videos^[1–12]. Meanwhile, in another field of natural language processing (NLP), the transformer architecture^[13] has been the de-facto standard to handle long-range dependencies^[14, 15]. Transformers rely heavily on self-attention to model global relationships of sequence data. Although global modeling is also essential for vision tasks, the 2D/3D structures of vision data make it less straightforward to apply transformers therein. This predicament has been recently broken by Dosovitskiy et al.^[16], by applying a pure transformer to sequences of image patches.

Motivated by ^[16], a large amount of literature on vision transformers has emerged to resolve the problems

caused by the domain gap between computer vision and NLP^[17–21]. From our point of view, one major problem of vision transformers is that the sequence length of image patches is much longer than that of tokens (words) in an NLP application, thus leading to high computational/space complexity when computing the multi-head self-attention (MHSA). Some efforts have been dedicated to resolving this problem. ToMe^[22] improves the throughput of existing ViT models^[16] by systematically merging similar tokens through the utilization of a general and light-weight matching algorithm. Pyramid vision transformer (PVT)^[19] and multiscale vision transformer (MViT)^[21] downsample the feature to compute attention in a reduced length of tokens but at the cost of losing fine-grained details. Swin transformer^[18] computes attention within small windows to model local relationships, and it gradually enlarges the receptive field by shifting windows and stacking more layers. From this point of view, Swin transformer^[18] may still be suboptimal because it works in a similar manner to CNNs and needs many layers to model long-range dependencies^[16].

Building upon the discussed strengths of downsampling-based transformers^[19, 21] and window-based transformers^[18], each with its distinctive merits, we aim to

Research Article
Special Issue on Multi-modal Representation Learning
Manuscript received on September 3, 2023; accepted on January 8, 2024; published online on April 19, 2024
Recommended by Associate Editor Wanli Ouyang
Colored figures are available in the online version at <https://link.springer.com/journal/11633>
© The Author(s) 2024

harness their complementary advantages. Downsampling-based transformers excel at directly modeling global dependencies but may sacrifice fine-grained details, while window-based transformers effectively capture local dependencies but may fall short in global dependency modeling. As widely accepted, both global and local information is essential for visual scene understanding. Motivated by this insight, our approach seeks to amalgamate the strengths of both paradigms, enabling the direct modeling of both global and local dependencies.

To achieve this, we introduce the hierarchical multi-head self-attention (H-MHSA), a novel mechanism that enhances the flexibility and efficiency of self-attention computation in transformers. Our methodology begins by segmenting an image into patches, treating each patch akin to a token^[16]. Rather than computing attention across all patches, we further organize these patches into small grids, performing attention computation within each grid. This step is instrumental in capturing local relationships and generating more discriminative local representations. Subsequently, we amalgamate these smaller patches into larger ones and treat the merged patches as new tokens, resulting in a substantial reduction in their number. This enables the direct modeling of global dependencies by calculating self-attention for the new tokens. Ultimately, the attentive features from both local and global hierarchies are aggregated to yield potent features with rich granularities. Notably, as the attention calculation at each step is confined to a small number of tokens, our hierarchical strategy mitigates the computational and space complexity of vanilla transformers. Empirical observations underscore the efficacy of this hierarchical self-attention mechanism, revealing improved generalization results in our experiments.

By simply incorporating H-MHSA, we build a family of hierarchical-attention-based transformer networks (HAT-Net). To evaluate the efficacy of HAT-Net in scene understanding, we experiment HAT-Net for fundamental vision tasks, including image classification, semantic segmentation, object detection and instance segmentation. Experimental results demonstrate that HAT-Net performs favorably against previous backbone networks. Note that H-MHSA is based on a very simple and intuitive idea, so H-MHSA is expected to provide a new perspective for the future design of vision transformers.

2 Related work

Convolutional neural networks. More than two decades ago, LeCun et al.^[23] built the first deep CNN, i.e., LeNet, for document recognition. About ten years ago, AlexNet^[1] introduced pooling layers into CNNs and pushed forward the state of the art of ImageNet classification^[24] significantly. Since then, CNNs have become the de-facto standard of computer vision owing to its powerful ability in representation learning. Brilliant achieve-

ments have been seen in this direction. VGGNet^[2] investigates networks of increasing depth using small (3×3) convolution filters. ResNet^[3] manages to build very deep networks by resolving the gradient vanishing/exploding problem with residual connections^[25]. GoogLeNet^[26] presents the inception architecture^[27, 28] using multiple branches with different convolution kernels. ResNeXt^[29] improves ResNet^[3] by replacing the 3×3 convolution in the bottleneck with a grouped convolution. DenseNets^[30] present dense connections, i.e., using the feature maps of all preceding layers as inputs for each layer. MobileNets^[31, 32] decompose the traditional convolution into a pointwise convolution and a depthwise separable convolution for acceleration, and an inverted bottleneck is proposed for ensuring accuracy. ShuffleNets^[33, 34] further decompose the pointwise convolution into pointwise group convolution and channel shuffle to reduce computational cost. MansNet^[35] proposes an automated mobile neural architecture search approach to search for a model with a good trade-off between accuracy and latency. EfficientNet^[36] introduces a scaling method to uniformly scale depth/width/resolution dimensions of the architecture searched by MansNet^[35]. The above advanced techniques are the engines driving the development of computer vision in the last decade. This paper aims at improving feature representation learning by designing new transformers.

Self-attention mechanism. Inspired by the human visual system, the self-attention mechanism is usually adopted to enhance essential information and suppress noisy information. Spatial transformer network (STN)^[37] presents the spatial attention mechanism through learning an appropriate spatial transformation for each input. Chen et al.^[38] propose the channel attention model and achieve promising results on the image captioning task. Wang et al.^[39] explore self-attention in well-known residual networks^[3]. SENet^[40] applies channel attention to backbone network design and boosts the accuracy of ImageNet classification^[24]. Convolutional block attention module (CBAM)^[41] sequentially applies channel and spatial attention for adaptive feature refinement in deep networks. Bottleneck attention module (BAM)^[42] produces a 3D attention map by combining channel and spatial attention. Selective kernel network (SK-Net)^[43] uses channel attention to selectively fuse multiple branches with different kernel sizes. Non-local network^[44] presents non-local attention for capturing long-range dependencies. ResNeSt^[45] is a milestone in this direction. It applies channel attention on different network branches to capture cross-feature interactions and learn diverse representations. Our work shares some similarities with these works by applying self-attention for learning feature representations. The difference is that we propose H-MHSA to learn global relationships rather than a simple feature recalibration using spatial or channel attention in these works.

Vision transformer. Transformer^[13] entirely relies

on self-attention to handle long-range dependencies of sequence data. It was first proposed for NLP tasks^[14, 15]. In order to apply transformers on image data, Dosovitskiy et al.^[16] split an image into patches and treat them as tokens. Then, a pure transformer^[13] can be adopted. Such a vision transformer (ViT) attains competitive accuracy for ImageNet classification^[24]. More recently, lots of efforts have been dedicated to improving ViT. Tokens-to-token ViT (T2T-ViT)^[46] proposes to split an image into tokens of overlapping patches so as to represent local structures by surrounding tokens. CaiT^[47] builds a deeper transformer network by introducing a per-channel weighting and specific class attention. DeepViT^[48] proposes Re-attention to re-generate attention maps to increase their diversity at different layers. DeiT^[49] presents a knowledge distillation strategy for improving the training of ViT^[16]. Srinivas et al.^[50] try to add the bottleneck structure to vision transformers. Some works build pyramid transformer networks to generate multi-scale features^[17–21]. PVT^[19] adopts convolution operation to down-sample the feature map in order to reduce the sequence length in MHSA, thus reducing the computational load. Similar to PVT^[19], MViT^[21] utilizes pooling to compute attention on a reduced sequence length. Swin transformer^[18] computes attention within small windows and shifts windows to gradually enlarge the receptive field. CoaT^[20] computes attention in the channel dimension rather than in the traditional spatial dimension. Token merging (ToMe)^[22] enhances the throughput of existing ViT models^[16] without requiring retraining, which is achieved by gradually combining similar tokens in a transformer using a matching algorithm. In this paper, we introduce a novel design to reduce the computational complexity of MHSA and learn both the global and local relationship modeling through vision transformers.

Vision MLP networks. While CNNs and vision transformers have been widely adopted for computer vision tasks, Tolstikhin et al.^[51] challenge the necessity of convolutions and attention mechanisms. They introduce the MLP-Mixer architecture, which relies solely on multi-layer perceptrons (MLPs). MLP-Mixer incorporates two types of layers: One applies MLPs independently to image patches, facilitating the mixing of per-location features, and the other applies MLPs across patches, enabling the mixing of spatial information. Despite lacking convolutions and attention, MLP-Mixer demonstrates competitive performance in image classification compared to state-of-the-art models. Liu et al.^[52] introduce gMLP, an MLP-based model with gating, showcasing its comparable performance to transformers in crucial language and vision applications. In contrast to other MLP-like models that encode spatial information along flattened spatial dimensions, Vision permutator^[53] uniquely encodes feature representations along height and width dimensions using linear projections. Wang et al.^[54] propose a novel positional spatial gating unit, leveraging classical

relative positional encoding to efficiently capture cross-token relations for token mixing. Despite these advancements, the performance of vision MLP networks still lags behind that of vision transformers. In this paper, we focus on the design of a new vision transformer network.

3 Methodology

In this section, we first provide a brief review of vision transformers^[16] in Section 3.1. Then, we present the proposed H-MHSA and analyze its computational complexity in Section 3.2. Finally, we describe the configuration details of the proposed HAT-Net in Section 3.3.

3.1 Review of vision transformers

Transformer^[13, 16] heavily relies on MHSA to model long-range relationships. Suppose $\mathbf{X} \in \mathbf{R}^{H \times W \times C}$ denotes the input, where H , W and C are the height, width and the feature dimension, respectively. We reshape \mathbf{X} and define the query \mathbf{Q} , key \mathbf{K} , value \mathbf{V} as

$$\begin{aligned} \mathbf{X} \in \mathbf{R}^{H \times W \times C} &\rightarrow \mathbf{X} \in \mathbf{R}^{(H \times W) \times C}, \\ \mathbf{Q} &= \mathbf{X} \mathbf{W}^q, \quad \mathbf{K} = \mathbf{X} \mathbf{W}^k, \quad \mathbf{V} = \mathbf{X} \mathbf{W}^v \end{aligned} \quad (1)$$

where $\mathbf{W}^q \in \mathbf{R}^{C \times C}$, $\mathbf{W}^k \in \mathbf{R}^{C \times C}$, and $\mathbf{W}^v \in \mathbf{R}^{C \times C}$ are the trainable weight matrices of linear transformations. With a mild assumption that the input and output have the same dimension, the traditional MHSA can be formulated as

$$\mathbf{A} = \text{Softmax}(\mathbf{Q} \mathbf{K}^T / \sqrt{d}) \mathbf{V} \quad (2)$$

in which \sqrt{d} means an approximate normalization, and the softmax function is applied to the rows of the matrix. Note that we omit the concept of multiple heads here for simplicity. In (2), the matrix product of $\mathbf{Q} \mathbf{K}^T$ first computes the similarity between each pair of tokens. Each new token is then derived over the combination of all tokens according to the similarity. After the computation of MHSA, a residual connection is further added to ease the optimization, like

$$\begin{aligned} \mathbf{X} \in \mathbf{R}^{(H \times W) \times C} &\rightarrow \mathbf{X} \in \mathbf{R}^{H \times W \times C}, \\ \mathbf{A}' &= \mathbf{A} \mathbf{W}^p + \mathbf{X} \end{aligned} \quad (3)$$

in which $\mathbf{W}^p \in \mathbf{R}^{C \times C}$ is a trainable weight matrix for feature projection. At last, a multilayer perceptron (MLP) is adopted to enhance the representation, like

$$\mathbf{Y} = \text{MLP}(\mathbf{A}') + \mathbf{A}' \quad (4)$$

where \mathbf{Y} denotes the output of a transformer block.

It is easy to infer that the computational complexity of MHSA (Equation (2)) is

$$\Omega(\text{MHSA}) = 3HWC^2 + 2H^2W^2C. \quad (5)$$

Similarly, the space complexity (memory consumption) also includes the term of $O(H^2W^2)$. As commonly known, $O(H^2W^2)$ could become very large for high-resolution inputs. This limits the applicability of transformers for vision tasks. Motivated by this, we aim at improving MHSA to reduce such complexity and maintain the capacity of global relationship modeling without the risk of sacrificing performances.

3.2 Hierarchical multi-head self-attention

In this section, we present an approach to alleviate the computational and space demands associated with (2) through the utilization of our proposed H-MHSA mechanism. Rather than computing attention over the entire input, we adopt a hierarchical strategy, allowing each step to process only a limited number of tokens.

The initial step concentrates on local attention computation. Assuming the input feature map is denoted as $\mathbf{X} \in \mathbf{R}^{H \times W \times C}$, we partition the feature map into small grids of size $G_1 \times G_1$ and reshape it as follows:

$$\begin{aligned} \mathbf{X} \in \mathbf{R}^{H \times W \times C} &\rightarrow \mathbf{X}_1 \in \mathbf{R}^{(\frac{H}{G_1} \times G_1) \times (\frac{W}{G_1} \times G_1) \times C} \\ &\rightarrow \mathbf{X}_1 \in \mathbf{R}^{(\frac{H}{G_1} \times \frac{W}{G_1}) \times (G_1 \times G_1) \times C}. \end{aligned} \quad (6)$$

The query, key and value are then calculated by

$$\mathbf{Q}_1 = \mathbf{X}_1 \mathbf{W}_1^q, \quad \mathbf{K}_1 = \mathbf{X}_1 \mathbf{W}_1^k, \quad \mathbf{V}_1 = \mathbf{X}_1 \mathbf{W}_1^v \quad (7)$$

where $\mathbf{W}_1^q, \mathbf{W}_1^k, \mathbf{W}_1^v \in \mathbf{R}^{C \times C}$ are trainable weight matrices. Subsequently, (2) is applied to generate the local attentive feature \mathbf{A}_1 . To ease network optimization, we reshape \mathbf{A}_1 back to the shape of \mathbf{X} through

$$\begin{aligned} \mathbf{A}_1 &\in \mathbf{R}^{(\frac{H}{G_1} \times \frac{W}{G_1}) \times (G_1 \times G_1) \times C} \\ &\rightarrow \mathbf{A}_1 \in \mathbf{R}^{(\frac{H}{G_1} \times G_1) \times (\frac{W}{G_1} \times G_1) \times C} \\ &\rightarrow \mathbf{A}_1 \in \mathbf{R}^{H \times W \times C} \end{aligned} \quad (8)$$

and incorporate a residual connection:

$$\mathbf{A}_1 = \mathbf{A}_1 + \mathbf{X}. \quad (9)$$

As the local attentive feature \mathbf{A}_1 is computed within each small $G_1 \times G_1$ grid, a substantial reduction in computational and space complexity is achieved.

The second step focuses on global attention calculation. Here, we downsample \mathbf{A}_1 by a factor of G_2 during the computation of key and value matrices. This downsampling enables efficient global attention calculation, treating each $G_2 \times G_2$ grid as a token. This process can be expressed as

$$\hat{\mathbf{A}}_1 = \text{AvePool}_{G_2}(\mathbf{A}_1) \quad (10)$$

where $\text{AvePool}_{G_2}(\cdot)$ denotes downsampling a feature map by G_2 times using average pooling with both the kernel size and stride set to G_2 . Consequently, we have $\hat{\mathbf{A}}_1 \in \mathbf{R}^{\frac{H}{G_2} \times \frac{W}{G_2} \times C}$. We then reshape \mathbf{A}_1 and $\hat{\mathbf{A}}_1$ as follows:

$$\begin{aligned} \mathbf{A}_1 &\in \mathbf{R}^{H \times W \times C} \rightarrow \mathbf{A}_1 \in \mathbf{R}^{(H \times W) \times C}, \\ \hat{\mathbf{A}}_1 &\in \mathbf{R}^{\frac{H}{G_2} \times \frac{W}{G_2} \times C} \rightarrow \hat{\mathbf{A}}_1 \in \mathbf{R}^{(\frac{H}{G_2} \times \frac{W}{G_2}) \times C}. \end{aligned} \quad (11)$$

Following this, we compute the query, key and value as

$$\mathbf{Q}_2 = \mathbf{A}_1 \mathbf{W}_2^q, \quad \mathbf{K}_2 = \hat{\mathbf{A}}_1 \mathbf{W}_2^k, \quad \mathbf{V}_2 = \hat{\mathbf{A}}_1 \mathbf{W}_2^v \quad (12)$$

where $\mathbf{W}_2^q, \mathbf{W}_2^k, \mathbf{W}_2^v \in \mathbf{R}^{C \times C}$ are trainable weight matrices. It is easy to derive that we have $\mathbf{Q}_2 \in \mathbf{R}^{(H \times W) \times C}$, $\mathbf{K}_2 \in \mathbf{R}^{(\frac{H}{G_2} \times \frac{W}{G_2}) \times C}$, and $\mathbf{V}_2 \in \mathbf{R}^{(\frac{H}{G_2} \times \frac{W}{G_2}) \times C}$. Subsequently, (2) is called to obtain the global attentive feature $\mathbf{A}_2 \in \mathbf{R}^{(H \times W) \times C}$, followed by a reshaping operation:

$$\mathbf{A}_2 \in \mathbf{R}^{(H \times W) \times C} \rightarrow \mathbf{A}_2 \in \mathbf{R}^{H \times W \times C}. \quad (13)$$

The final output of H-MHSA is given by

$$\text{H-MHSA}(\mathbf{X}) = (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{W}^p + \mathbf{X} \quad (14)$$

where \mathbf{W}^p has the same meaning as in (3). In this way, H-MHSA effectively models both local and global relationships, akin to vanilla MHSA.

The computational complexity of H-MHSA can be expressed as

$$\Omega(\text{H-MHSA}) = HWC(4C + 2G_1^2) + 2\frac{HW}{G_2^2}C(C + HW). \quad (15)$$

Compared to (5), this represents a reduction in computational complexity from $O(H^2W^2)$ to $O(HWG_1^2 + \frac{H^2W^2}{G_2^2})$. The same conclusion can be easily derived for space complexity.

We continue by comparing H-MHSA with existing vision transformers, highlighting distinctive features. Swin transformer^[18] focuses on modeling local relationships, progressively expanding the receptive field through shifted windows and additional layers. Conversely, PVT^[19] prioritizes global relationships through downsampling key and value matrices but overlooks local information. In contrast, our proposed H-MHSA excels by concurrently capturing both local and global relationships. While Swin transformer employs a fixed window size (i.e., a fixed-size bias matrix), and PVT uses a constant downsampling ratio (i.e., a convolution with the kernel size equal to the stride), these approaches necessitate retraining on the ImageNet dataset^[24] for any re-parameterization. In con-

trast, the parameter-free nature of G_1 and G_2 in H-MHSA allows flexible configuration adjustments for downstream vision tasks without the need for retraining on ImageNet.

In computer vision, achieving a comprehensive understanding of scenes relies on the simultaneous consideration of both global and local information. Within the framework of our proposed H-MHSA, global self-attention calculation (Equations (10)–(13)) is instrumental in establishing the foundation for scene interpretation, enabling the recognition of overarching patterns and aiding in high-level decision-making processes. Concurrently, local self-attention calculation (Equations (6)–(9)) is crucial for refining the understanding of individual components within the larger context, facilitating more detailed and nuanced scene analysis. H-MHSA excels in striking the delicate balance between global and local information, thereby facilitating a nuanced and accurate comprehension of diverse scenes. In essence, the seamless integration of global and local self-attention within the H-MHSA framework empowers transformers to navigate the intricacies of scene understanding, facilitating context-aware decision-making.

3.3 Network architecture

This part introduces the network architecture of HAT-Net. We follow the common practice in CNNs to use a global average pooling layer and a fully connected layer to predict image classes^[18]. This is different from existing transformers which rely on another 1 D class token to make predictions^[16, 17, 19–21, 46–49, 55–57]. We also observe that existing transformers^[16–21, 46–49] usually adopt the Gaussian error linear unit (GELU) function^[58] for nonlinear activation. However, GELU is memory-hungry during network training. We empirically found that the sigmoid-weighted linear unit (SiLU) function^[59], originally coined in [58], performs on-par with GELU and is more memory-friendly. Hence, HAT-Net uses SiLU^[59] for nonlinear activation. Besides, we add a depthwise separable convolution (DW-Conv)^[31] inside the MLP as widely done.

The overall architecture of HAT-Net is illustrated in Fig. 1. At the beginning of HAT-Net, instead of flatten-

ing image patches^[16], we apply two sequential vanilla 3×3 convolutions, each of which has a stride of 2, to downsample the input image into $\frac{1}{4}$ scale. Then, we stack H-MHSA and MLP alternatively, which can be divided into four stages with pyramid feature scales of $\frac{1}{4}$, $\frac{1}{8}$, $\frac{1}{16}$, and $\frac{1}{32}$, respectively. For feature downsampling at the end of each stage, a vanilla 3×3 convolution with a stride of 2 is used. The configuration details of HAT-Net are summarized in Table 1. We provide four versions of HAT-Net: HAT-Net-Tiny, HAT-Net-Small, HAT-Net-Medium and HAT-Net-Large, whose number of parameters is similar to ResNet18, ResNet50, ResNet101 and ResNet152^[3], respectively. We only adopt simple parameter settings without careful tuning to demonstrate the effectiveness and generality of HAT-Net. The dimension of each head in the multi-head setting is set to 48 for HAT-Net-Tiny and 64 for other versions.

To enhance the applicability of HAT-Net across diverse vision tasks, we present guidelines for configuring the parameter-free G_1 and G_2 . While established models like Swin transformer^[18] adhere to a fixed window size of 7, and PVT^[19] employs a set of constant downsampling ratios 8, 4, 2 for the t -th stage ($t = 2, 3, 4$), we advocate for certain adjustments. Practically, we find that a window size of 8 is more pragmatic than 7, given that input resolutions often align with multiples of 8. Moreover, augmenting downsampling ratios serves to mitigate computational complexity. Consequently, for image classification on the ImageNet dataset^[24], where the standard input resolution is 224×224 pixels, we designate $G_1 = 8, 7, 7$ and $G_2 = 8, 4, 2$ for the t -th stage ($t = 2, 3, 4$). Here, a window size of 7 is necessitated by the chosen resolution and small downsampling rates are in line with the approach taken by PVT^[19]. In scenarios involving downstream tasks like semantic segmentation, object detection and instance segmentation, where input resolutions tend to be larger, we opt for $G_1 = 8, 8, 8$ for convenience and $G_2 = 16, 8, 4$ to curtail computational expenses. For a comprehensive analysis of the impact of different G_1 and G_2 settings, we conduct an ablation study in Section 4.4.

4 Experiments

To show the superiority of HAT-Net in feature repres-

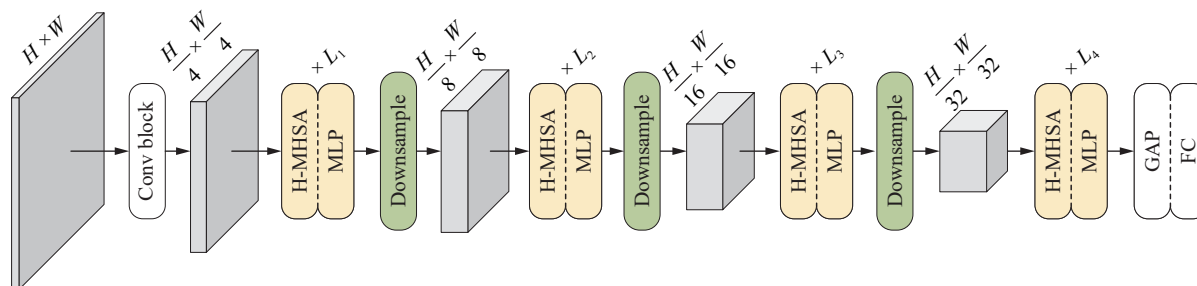


Fig. 1 Illustration of the proposed HAT-Net. GAP: Global average pooling; FC: Fully-connected layer. $\times L_i$ means that the transformer block is repeated for L_i times. H and W denote the height and width of the input image, respectively.

Table 1 Network configurations of HAT-Net. The settings of building blocks are shown in brackets, with the number of blocks stacked. For the first stage, each convolution has C channels and a stride of S . For the other four stages, each MLP uses a $K \times K$ DW-Conv and an expansion ratio of E . Note that we omit the downsampling operation after the t -th stage ($t = \{2, 3, 4\}$) for simplicity. “#Param” refers to the number of parameters.

Stage	Input size	Operator	HAT-Net-Tiny	HAT-Net-Small	HAT-Net-Medium	HAT-Net-Large
1	224×224	3×3 conv.	$C = 16, S = 2$ $C = 48, S = 2$	$C = 16, S = 2$ $C = 64, S = 2$	$C = 16, S = 2$ $C = 64, S = 2$	$C = 16, S = 2$ $C = 64, S = 2$
2	56×56	H-MHSA MLP	$\begin{bmatrix} C = 48 \\ K = 3 \\ E = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} C = 64 \\ K = 3 \\ E = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} C = 64 \\ K = 5 \\ E = 8 \end{bmatrix} \times 3$	$\begin{bmatrix} C = 64 \\ K = 3 \\ E = 8 \end{bmatrix} \times 3$
3	28×28	H-MHSA MLP	$\begin{bmatrix} C = 96 \\ K = 3 \\ E = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} C = 128 \\ K = 3 \\ E = 8 \end{bmatrix} \times 3$	$\begin{bmatrix} C = 128 \\ K = 3 \\ E = 8 \end{bmatrix} \times 6$	$\begin{bmatrix} C = 128 \\ K = 3 \\ E = 8 \end{bmatrix} \times 8$
4	14×14	H-MHSA MLP	$\begin{bmatrix} C = 240 \\ K = 3 \\ E = 4 \end{bmatrix} \times 6$	$\begin{bmatrix} C = 320 \\ K = 3 \\ E = 4 \end{bmatrix} \times 8$	$\begin{bmatrix} C = 320 \\ K = 5 \\ E = 4 \end{bmatrix} \times 18$	$\begin{bmatrix} C = 320 \\ K = 3 \\ E = 4 \end{bmatrix} \times 27$
5	7×7	H-MHSA MLP	$\begin{bmatrix} C = 384 \\ K = 3 \\ E = 4 \end{bmatrix} \times 3$	$\begin{bmatrix} C = 512 \\ K = 3 \\ E = 4 \end{bmatrix} \times 3$	$\begin{bmatrix} C = 512 \\ K = 3 \\ E = 4 \end{bmatrix} \times 3$	$\begin{bmatrix} C = 640 \\ K = 3 \\ E = 4 \end{bmatrix} \times 3$
	1×1	–	Global average pooling, 1 000-d FC, softmax			
	#Param		12.7 M	25.7 M	42.9 M	63.1 M

entation learning, this section evaluates HAT-Net for image classification, semantic segmentation, object detection and instance segmentation.

4.1 Image classification

Experimental setup. The ImageNet dataset^[24] consists of 1.28 M training images and 50 K validation images from 1 000 categories. We adopt the training set to train our networks and the validation set to test the performance. We implement HAT-Net using the popular PyTorch framework^[60]. For a fair comparison, we follow the same training protocol as DeiT^[49], which is the standard protocol for training transformer networks nowadays. Specifically, the input images are randomly cropped to 224×224 pixels, followed by random horizontal flipping and mixup^[61] for data augmentation. Label smoothing^[27] is used to avoid overfitting. The AdamW optimizer^[62] is adopted with the momentum of 0.9, the weight decay of 0.05, and a mini-batch size of 128 per GPU by default. The initial learning rate is set to 1×10^{-3} , which decreases following the cosine learning rate schedule^[63]. The training process lasts for 300 epochs on eight NVIDIA Tesla V100 GPUs. Note that for ablation studies, we utilize a mini-batch size of 64 and 100 training epochs to save training time. Moreover, we set $G_1 = \{8, 7, 7\}$ and $G_2 = \{8, 4, 2\}$ for the t -th stage ($t = \{2, 3, 4\}$), respectively. The fifth stage can be processed directly using the vanilla MHSA mechanism. For model evaluation, we apply a center crop of 224×224 pixels on validation images to evaluate the recognition accuracy. We report the top-1 classification accuracy on the ImageNet validation set^[24] as well as the number of parameters and the number of FLOPs for each model.

Experimental results. We compare HAT-Net with state-of-the-art network architectures, including CNN-based ones like ResNet^[3], ResNeXt^[29], RegNetY^[64], ResNeSt^[45], and transformer-based ones like ViT^[16], DeiT^[49], T2T-ViT^[46], TNT^[55], CvT^[65], MViT^[21], PVT^[19], Swin transformer^[18], Twins^[66], ToMe^[22]. The results are summarized in Table 2. We can observe that HAT-Net achieves state-of-the-art performance. Specifically, with similar numbers of parameters and FLOPs, HAT-Net-Tiny, HAT-Net-Small, HAT-Net-Medium and HAT-Net-Large outperform the second best results by 1.1%, 0.6%, 0.8% and 0.6% in terms of the top-1 accuracy, respectively. Since the performance for image classification implies the ability of a network for learning feature representations, the above comparison suggests that the proposed HAT-Net has great potential for generic scene understanding.

4.2 Semantic segmentation

Experimental setup. We continue by applying HAT-Net to a fundamental downstream vision task, semantic segmentation, which aims at predicting a class label for each pixel in an image. We follow [19, 66] to replace the backbone of the well-known segmentor, Semantic FPN^[68], with HAT-Net or other backbone networks for a fair comparison. Experiments are conducted on the challenging ADE20K dataset^[69]. This dataset has 20 000 training images, 2 000 validation images, and 3 302 testing images. We train Semantic FPN^[68] using the training set and evaluate it on the validation set. The training optimizer is AdamW^[62] with weight decay of 1×10^{-4} . We apply the poly learning rate schedule with $\gamma = 0.9$ and the initial learning rate of 1×10^{-4} . During training, the

batch size is 16, and each image has a resolution of 512×512 through resizing and cropping. During testing, each image is resized to the shorter side of 512 pixels,

Table 2 Comparison to state-of-the-art methods on the ImageNet validation set^[24]. “*” indicates the performance of a method using the default training setting in the original paper. “#Param” and “#FLOPs” refer to the number of parameters and the number of FLOPs, respectively. “†” marks models that use the input size of 384×384 ; Otherwise, models use the input size of 224×224 .

Arch.	Models	#Param	#FLOPs	Top-1 Acc. (%)
CNN	ResNet18*[3]	11.7 M	1.8 G	69.8
	ResNet18[3]	11.7 M	1.8 G	68.5
	DeiT-Ti/16 ^[49]	5.7 M	1.3 G	72.2
Trans	PVT-Tiny ^[19]	13.2 M	1.9 G	75.1
	PVTv2-B1 ^[67]	13.1 M	2.1 G	78.7
	HAT-Net-Tiny	12.7 M	2.0 G	79.8
CNN	ResNet50*[3]	25.6 M	4.1 G	76.1
	ResNet50[3]	25.6 M	4.1 G	78.5
	ResNeXt50-32x4d*[29]	25.0 M	4.3 G	77.6
	ResNeXt50-32x4d ^[29]	25.0 M	4.3 G	79.5
	RegNetY-4G ^[64]	20.6 M	4.0 G	80.0
	ResNeSt-50 ^[45]	27.5 M	5.4 G	81.1
	ToMe-ViT-S/16 ^[22]	22.1 M	2.7 G	79.4
	DeiT-S/16 ^[49]	22.1 M	4.6 G	79.8
	T2T-ViT _t -14 ^[46]	21.5 M	5.2 G	80.7
	TNT-S ^[55]	23.8 M	5.2 G	81.3
Trans	CvT-13 ^[65]	20.0 M	4.5 G	81.6
	PVT-Small ^[19]	24.5 M	3.8 G	79.8
	PVTv2-B2 ^[67]	25.4 M	4.0 G	82.0
	Swin-T ^[18]	28.3 M	4.5 G	81.3
	Twins-SVT-S ^[66]	24.0 M	2.8 G	81.7
	HAT-Net-Small	25.7 M	4.3 G	82.6
CNN	ResNet101*[3]	44.7 M	7.9 G	77.4
	ResNet101[3]	44.7 M	7.9 G	79.8
	ResNeXt101-32x4d*[29]	44.2 M	8.0 G	78.8
	ResNeXt101-32x4d ^[29]	44.2 M	8.0 G	80.6
	RegNetY-8G ^[64]	39.2 M	8.0 G	81.7
	ResNeSt-101 ^[45]	48.3 M	10.3 G	83.0
	T2T-ViT _t -19 ^[46]	39.2 M	8.4 G	81.4
	CvT-21 ^[65]	31.5 M	7.1 G	82.5
	MViT-B-16 ^[21]	37.0 M	7.8 G	82.5
	PVT-Medium ^[19]	44.2 M	6.7 G	81.2
Trans	PVTv2-B3 ^[67]	45.2 M	6.9 G	83.2
	Swin-S ^[18]	49.6 M	8.7 G	83.0
	HAT-Net-Medium	42.9 M	8.3 G	84.0

Table 2 (continued) Comparison to state-of-the-art methods on the ImageNet validation set^[24]. “*” indicates the performance of a method using the default training setting in the original paper. “#Param” and “#FLOPs” refer to the number of parameters and the number of FLOPs, respectively. “†” marks models that use the input size of 384×384 ; Otherwise, models use the input size of 224×224 .

Arch.	Models	#Param	#FLOPs	Top-1 Acc. (%)
CNN	ResNet152*[3]	60.2 M	11.6 G	78.3
	ResNeXt101-64x4d*[29]	83.5 M	15.6 G	79.6
	ResNeXt101-64x4d ^[29]	83.5 M	15.6 G	81.5
Trans	ViT-B/16† ^[16]	86.6 M	55.4 G	77.9
	ViT-L/16† ^[16]	304.3 M	190.7 G	76.5
	ToMe-ViT-L/16 ^[22]	304.3 M	22.3 G	84.2
	DeiT-B/16 ^[49]	86.6 M	17.6 G	81.8
	MViT-B-24 ^[21]	53.5 M	10.9 G	83.1
	TNT-B ^[55]	65.6 M	14.1 G	82.8
Trans	PVT-Large ^[19]	61.4 M	9.8 G	81.7
	PVTv2-B4 ^[67]	62.6 M	10.1 G	83.6
	Swin-B ^[18]	87.8 M	15.4 G	83.3
	Twins-SVT-B ^[66]	56.0 M	8.3 G	83.2
	HAT-Net-Large	63.1 M	11.5 G	84.2

without multi-scale testing or flipping. We adopt the well-known MMSegmentation toolbox^[70] for the above experiments. We set $G_1 = \{8, 8, 8\}$ and $G_2 = \{16, 8, 4\}$ for the t -th stage ($t = \{2, 3, 4\}$), respectively.

Experimental results. The results are depicted in Table 3. We compare with typical CNN networks, i.e., ResNets^[3] and ResNeXts^[29], and transformer networks, i.e., Swin transformer^[18], PVT^[19], PVTv2^[67] and Twins-SVT^[66]. As can be observed, the proposed HAT-Net achieves significantly better performance than previous competitors. Specifically, HAT-Net-Tiny, HAT-Net-Small, HAT-Net-Medium and HAT-Net-Large attain 1.9%, 0.4%, 1.9% and 0.7% higher mIoU than the second better results with similar number of parameters and FLOPs. This demonstrates the superiority of HAT-Net in learning effective feature representations for dense prediction tasks.

4.3 Object detection and instance segmentation

Experimental setup. Since object detection and instance segmentation are also fundamental downstream vision tasks, we apply HAT-Net to both tasks for further evaluating its effectiveness. Specifically, we utilize two well-known detectors, i.e., RetinaNet^[71] for object detection and Mask R-CNN^[5] for instance segmentation. HAT-Net is compared to some well-known CNN and transformer networks by only replacing the backbone of the above two detectors. Experiments are conducted on the

Table 3 Experimental results on the ADE20K validation dataset^[69] for semantic segmentation. We replace the backbone of Semantic FPN^[68] with various network architectures. The number of FLOPs is calculated with the input size of 512×512 .

Backbone	Semantic FPN ^[68]		
	#Param (M) ↓	FLOPs (G) ↓	mIoU (%) ↑
ResNet-18 ^[3]	15.5	31.9	32.9
PVT-Tiny ^[19]	17.0	32.1	35.7
PVTv2-B1 ^[67]	17.8	33.1	41.5
HAT-Net-Tiny	15.9	33.2	43.6
ResNet-50 ^[3]	28.5	45.4	36.7
PVT-Small ^[19]	28.2	42.9	39.8
Swin-T ^[18]	31.9	46.0	41.5
Twins-SVT-S ^[66]	28.3	37.0	43.2
PVTv2-B2 ^[67]	29.1	44.1	46.1
HAT-Net-Small	29.5	49.6	46.6
ResNet-101 ^[3]	47.5	64.8	38.8
ResNeXt-101-32x4d ^[29]	47.1	64.6	39.7
PVT-Medium ^[19]	48.0	59.4	41.6
Swin-S ^[18]	53.2	70.0	45.2
Twins-SVT-B ^[66]	60.4	67.0	45.3
PVTv2-B3 ^[67]	49.0	60.7	47.3
HAT-Net-Medium	46.7	74.7	49.3
ResNeXt-101-64x4d ^[29]	86.4	104.2	40.2
PVT-Large ^[19]	65.1	78.0	42.1
Swin-B ^[18]	91.2	107.0	46.0
Twins-SVT-L ^[66]	102.0	103.7	46.7
PVTv2-B4 ^[67]	66.3	79.6	48.6
HAT-Net-Large	66.8	96.4	49.5

large-scale MS-COCO dataset^[72] by training on the train2017 set (~118K images) and evaluating on the val2017 set (5K images). We adopt MMDetection toolbox^[73] for experiments and follow the experimental settings of PVT^[19] for a fair comparison. During training, we initialize the backbone weights with the ImageNet-pretrained models. The detectors are fine-tuned using the AdamW optimizer^[62] with an initial learning rate of 1×10^{-4} that is decreased by 10 times after the 8-th and 11-th epochs, respectively. The whole training lasts for 12 epochs with a batch size of 16. Each image is resized to a shorter side of 800 pixels, but the longer side is not allowed to exceed 1333 pixels. We set $G_1 = \{8, 8, 8\}$ and $G_2 = \{16, 8, 4\}$ for the t -th stage ($t = \{2, 3, 4\}$), respectively.

Experimental results. The results are displayed in Table 4. As can be seen, HAT-Net substantially improves the accuracy over other network architectures with a similar number of parameters. Twins-SVT^[66] combines the advantages of PVT^[19] and Swin transformer^[18]

by alternatively stacking their basic blocks. When RetinaNet^[71] is adopted as the detector, HAT-Net-Small attains 1.8%, 1.6% and 1.8% higher results than Twins-SVT-S^[66] in terms of AP, AP₅₀ and AP₇₅, respectively. Correspondingly, HAT-Net-Large gets 1.0%, 0.5% and 1.5% higher results than Twins-SVT-B^[66]. With Mask R-CNN^[5] as the detector, HAT-Net-Large achieves 2.2%, 1.7% and 2.8% higher results than Twins-SVT-B^[66] in terms of bounding box metrics AP_b, AP₅₀^b and AP₇₅^b, respectively. HAT-Net-Large achieves 1.6%, 2.0% and 1.8% higher results than Twins-SVT-B^[66] in terms of mask metrics AP_m, AP₅₀^m and AP₇₅^m, respectively. Such significant improvement in object detection and instance segmentation shows the superiority of HAT-Net in learning effective feature representations.

4.4 Ablation studies

In this part, we evaluate various design choices of the proposed HAT-Net. As discussed above, we only train all ablation models for 100 epochs to save training time. The batch size and learning rate are also reduced by half accordingly. HAT-Net-Small is adopted for these ablation studies.

Effect of the proposed H-MHSA. Starting from the window attention^[18] based transformer network, we gradually replace the window attention with our proposed H-MHSA at different stages. The results are summarized in Table 5. Since the feature map at the fifth stage is small enough for directly computing MHSA, the fifth stage is excluded from Table 5. Note that the first stage of HAT-Net only consists of convolutions so that it is also excluded. From Table 5, we can observe that the performance for both image classification and semantic segmentation is improved when more stages adopt H-MHSA. This verifies the effectiveness of the proposed H-MHSA in feature presentation learning. It is interesting to find that the usage of H-MHSA at the fourth stage leads to more significant improvement than other stages. Intuitively, the fourth stage has the most transformer blocks, so the changes at this stage would lead to more significant effects.

A pure transformer version of HAT-Net VS. PVT^[19]. When we remove all depthwise separable convolutions from HAT-Net and train the resulting transformer network for 100 epochs, it achieves 77.7% top-1 accuracy on the ImageNet validation set^[24]. In contrast, the well-known transformer network, PVT^[19], attains 75.8% top-1 accuracy under the same condition. This suggests that our proposed H-MHSA is very effective in feature representation learning.

SiLU^[59] VS. GELU^[58]. We use SiLU function^[59] for nonlinearization rather than the widely-used GELU function^[58] in transformers^[13, 16]. Here, we evaluate the effect of this choice. HAT-Net with SiLU^[59] attains 82.6% top-1 accuracy on the ImageNet validation set^[24] when trained for 300 epochs. HAT-Net with GELU^[58] gets 82.7% top-1 accuracy, slightly higher than SiLU^[59]. However, HAT-

Table 4 Object detection results with RetinaNet^[71] and instance segmentation results with Mask R-CNN^[5] on the MS-COCO val2017 set^[72]. “R” and “X” represent ResNet^[3] and ResNeXt^[29], respectively. The number of FLOPs is computed with the input size of 800×1280 .

Backbone	Object detection								Instance segmentation							
	#Param (M) ↓	#FLOPs (G) ↓	RetinaNet ^[71]						#Param (M) ↓	#FLOPs (G) ↓	Mask R-CNN ^[5]					
			AP (%)	AP ₅₀ (%)	AP ₇₅ (%)	AP _S (%)	AP _M (%)	AP _L (%)			AP ^b (%)	AP ^b ₅₀ (%)	AP ^b ₇₅ (%)	AP ^m (%)	AP ^m ₅₀ (%)	AP ^m ₇₅ (%)
R-18 ^[3]	21.3	190.0	31.8	49.6	33.6	16.3	34.3	43.2	31.2	209.0	34.0	54.0	36.7	31.2	51.0	32.7
ViL-Tiny ^[74]	16.6	204.0	40.8	61.3	43.6	26.7	44.9	53.6	26.9	223.0	41.4	63.5	45.0	38.1	60.3	40.8
PVT-Tiny ^[19]	23.0	205.0	36.7	56.9	38.9	22.6	38.8	50.0	32.9	223.0	36.7	59.2	39.3	35.1	56.7	37.3
HAT-Net-Tiny	21.6	212.0	42.5	63.3	45.8	26.9	46.1	56.6	31.8	231.0	43.1	65.4	47.4	39.7	62.5	42.4
R-50 ^[3]	37.7	239.0	36.3	55.3	38.6	19.3	40.0	48.8	44.2	260.0	38.0	58.6	41.4	34.4	55.1	36.7
PVT-Small ^[19]	34.2	261.0	40.4	61.3	43.0	25.0	42.9	55.7	44.1	280.0	40.4	62.9	43.8	37.8	60.1	40.3
Swin-T ^[18]	38.5	248.0	41.5	62.1	44.2	25.1	44.9	55.5	47.8	264.0	42.2	64.6	46.2	39.1	61.6	42.0
ViL-Small ^[74]	35.7	292.0	44.2	65.2	47.6	28.8	48.0	57.8	45.0	310.0	44.9	67.1	49.3	41.0	64.2	44.1
Twins-SVT-S ^[66]	34.3	236.0	43.0	64.2	46.3	28.0	46.4	57.5	44.0	254.0	43.4	66.0	47.3	40.3	63.2	43.4
HAT-Net-Small	35.5	286.0	44.8	65.8	48.1	28.8	48.6	59.5	45.4	303.0	45.2	67.6	49.9	41.6	64.6	44.7
R-101 ^[3]	56.7	315.0	38.5	57.8	41.2	21.4	42.6	51.1	63.2	336.0	40.4	61.1	44.2	36.4	57.7	38.8
X-101-32x4d ^[29]	56.4	319.0	39.9	59.6	42.7	22.3	44.2	52.5	62.8	340.0	41.9	62.5	45.9	37.5	59.4	40.2
PVT-Medium ^[19]	53.9	349.0	41.9	63.1	44.3	25.0	44.9	57.6	63.9	367.0	42.0	64.4	45.6	39.0	61.6	42.1
Swin-S ^[18]	59.8	336.0	44.5	65.7	47.5	27.4	48.0	59.9	69.1	354.0	44.8	66.6	48.9	40.9	63.4	44.2
HAT-Net-Medium	52.7	405.0	45.9	66.9	49.2	29.7	50.0	61.6	62.6	424.0	47.0	69.0	51.5	42.7	66.0	46.0
X-101-64x4d ^[29]	95.5	473.0	41.0	60.9	44.0	23.9	45.2	54.0	101.9	493.0	42.8	63.8	47.3	38.4	60.6	41.3
PVT-Large ^[19]	71.1	450.0	42.6	63.7	45.4	25.8	46.0	58.4	81.0	469.0	42.9	65.0	46.6	39.5	61.9	42.5
Twins-SVT-B ^[66]	67.0	376.0	45.3	66.7	48.1	28.5	48.9	60.6	76.3	395.0	45.2	67.6	49.3	41.5	64.5	44.8
HAT-Net-Large	73.1	519.0	46.3	67.2	49.6	30.0	50.6	62.4	82.7	537.0	47.4	69.3	52.1	43.1	66.5	46.6

Table 5 Ablation studies for the hierarchical attention in HAT-Net. The configuration of HAT-Net-Small is adopted for all experiments. “√” indicates that we replace the window attention^[18] with the hierarchical attention at the i -th stage. “Top-1 Acc” is the top-1 accuracy on the ImageNet validation dataset^[24]. “mIoU” is the mean IoU for semantic segmentation on the ADE20K dataset^[69].

Design	#Stage			Top-1 Acc. (%)	mIoU (%)
	2	3	4		
1				78.2	42.1
2	√			78.2	42.4
3	√	√		78.4	42.5
4	√	√	√	79.3	43.4

Net with GELU^[58] only obtains 45.7% mIoU on the ADE20K dataset, 0.8% lower than HAT-Net with SiLU. When using a batch size of 128 per GPU, HAT-Net with SiLU^[59] occupies 20.2GB GPU memory during training, while HAT-Net with GELU^[58] occupies 23.8GB GPU memory. Hence, HAT-Net with SiLU^[59] can achieve slightly better performance with less GPU memory consumption.

Settings of G_1 and G_2 . In HAT-Net, the paramet-

ers G_1 and G_2 play pivotal roles, controlling grid sizes for local attention calculation and downsampling rates for global attention calculation, respectively. In this evaluation, we assess the model’s performance under various configurations of G_1 and G_2 . By default, for tasks such as object detection and instance segmentation, we employ $G_1 = \{8, 8, 8\}$ and $G_2 = \{16, 8, 4\}$ for the t -th stage ($t = \{2, 3, 4\}$), respectively. Subsequently, we systematically vary G_1 and G_2 , evaluating the performance of Mask R-CNN^[5] with HAT-Net-Small as the backbone. The evaluation results, conducted on the MS-COCO val2017 dataset^[72], are presented in Table 6. The findings indicate that HAT-Net demonstrates robustness across different G_1 and G_2 settings. Notably, altering G_1 from its default $\{8, 8, 8\}$ configuration has a marginal impact on performance, resulting in slight performance reduction. Similarly, adjusting the values of G_2 yields a trade-off: Decreasing values enhances performance at the expense of increased computational cost, while increasing values reduces computational cost at the cost of slightly degraded performance. Our default choice of $G_1 = \{8, 8, 8\}$ and $G_2 = \{16, 8, 4\}$ strikes a favorable balance between accuracy and efficiency, offering a practical configuration for general use.

Table 6 Ablation studies for the settings of G_1 and G_2 in HAT-Net. The performance assessment is conducted using Mask R-CNN^[5] with HAT-Net-Small as the backbone. Evaluation results are reported on the MS-COCO val2017 dataset^[72].

Settings		#FLOPs (G)↓	Mask R-CNN ^[5]					
G_1	G_2		AP ^b (%)	AP ^b ₅₀ (%)	AP ^b ₇₅ (%)	AP ^m (%)	AP ^m ₅₀ (%)	AP ^m ₇₅ (%)
{8, 8, 8}	{12, 4, 2}	338.0	45.7	67.8	50.4	41.7	64.8	44.7
{8, 8, 8}	{12, 6, 3}	313.0	45.3	67.6	49.7	41.6	64.8	44.9
{8, 8, 8}	{16, 8, 4}	303.0	45.2	67.6	49.9	41.6	64.6	44.7
{8, 8, 8}	{32, 16, 8}	291.0	44.6	66.8	49.1	41.0	63.8	44.3
{16, 16, 8}	{16, 8, 4}	309.0	45.1	67.5	49.5	41.3	64.5	44.4
{4, 4, 4}	{16, 8, 4}	300.0	45.1	67.1	49.5	41.3	64.3	44.5

5 Conclusions

This paper addresses the inefficiency inherent in vanilla vision transformers due to the elevated computational and space complexity associated with MHSA. In response to this challenge, we introduce a novel hierarchical framework for MHSA computation, denoted as H-MHSA, aiming to alleviate the computational and space demands. Compared to existing approaches in this domain, such as PVT^[19] and Swin transformer^[18], H-MHSA distinguishes itself by directly capturing both global dependencies and local relationships. Integrating the proposed H-MHSA, we formulate the HAT-Net family, showcasing its prowess through comprehensive experiments spanning image classification, semantic segmentation, object detection, and instance segmentation. Our results affirm the efficacy and untapped potential of HAT-Net in advancing representation learning.

Applications of HAT-Net. The versatility of HAT-Net extends its utility across diverse real-world scenarios and downstream vision tasks. As a robust backbone network for feature extraction, HAT-Net seamlessly integrates with existing prediction heads and decoder networks, enabling proficient execution of various scene understanding tasks. Furthermore, HAT-Net's adaptability to different input resolutions and computational resource constraints is facilitated by the flexible adjustment of parameters, specifically G_1 and G_2 . Users can tailor HAT-Net to their specific requirements, selecting from different HAT-Net versions to align with their objectives.

In conclusion, HAT-Net not only presents a pragmatic solution to the limitations of vanilla vision transformers but also opens avenues for innovation in the future design of such architectures. The simplicity of the proposed H-MHSA underscores its potential as a transformative element in the evolving landscape of vision transformer development.

Acknowledgements

This work was supported by A*STAR Career Development Fund, Singapore (No. C233312006). Open access

funding provided by Swiss Federal Institute of Technology Zurich.

Declarations of conflict of interest

The authors declared that they have no conflicts of interest to this work.

Open Access

This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made.

The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

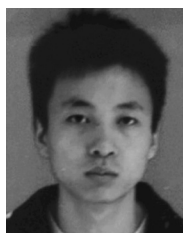
References

- [1] A. Krizhevsky, I. Sutskever, G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems*, Lake Tahoe, USA, pp.1097–1105, 2012. DOI: [10.5555/2999134.2999257](https://doi.org/10.5555/2999134.2999257).
- [2] K. Simonyan, A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proceedings of the 3rd International Conference on Learning Representations*, San Diego, USA, 2015.
- [3] K. M. He, X. Y. Zhang, S. Q. Ren, J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Las Vegas, USA, pp.770–778, 2016. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).

- [4] S. Q. Ren, K. M. He, R. Girshick, J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol.39, no.6, pp.1137–1149, 2017. DOI: [10.1109/TPAMI.2016.2577031](https://doi.org/10.1109/TPAMI.2016.2577031).
- [5] K. M. He, G. Gkioxari, P. Dollár, R. Girshick. Mask R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision*, Venice, Italy, pp.2980–2988, 2017. DOI: [10.1109/ICCV.2017.322](https://doi.org/10.1109/ICCV.2017.322).
- [6] H. S. Zhao, J. P. Shi, X. J. Qi, X. G. Wang, J. Y. Jia. Pyramid scene parsing network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, USA, pp.6230–6239, 2017. DOI: [10.1109/CVPR.2017.660](https://doi.org/10.1109/CVPR.2017.660).
- [7] Y. Liu, M. M. Cheng, X. W. Hu, J. W. Bian, L. Zhang, X. Bai, J. H. Tang. Richer convolutional features for edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol.41, no.8, pp.1939–1946, 2019. DOI: [10.1109/TPAMI.2018.2878849](https://doi.org/10.1109/TPAMI.2018.2878849).
- [8] Y. Liu, M. M. Cheng, D. P. Fan, L. Zhang, J. W. Bian, D. C. Tao. Semantic edge detection with diverse deep supervision. *International Journal of Computer Vision*, vol.130, no.1, pp.179–198, 2022. DOI: [10.1007/s11263-021-01539-8](https://doi.org/10.1007/s11263-021-01539-8).
- [9] Y. Liu, M. M. Cheng, X. Y. Zhang, G. Y. Nie, M. Wang. DNA: Deeply supervised nonlinear aggregation for salient object detection. *IEEE Transactions on Cybernetics*, vol.52, no.7, pp.6131–6142, 2022. DOI: [10.1109/TCYB.2021.3051350](https://doi.org/10.1109/TCYB.2021.3051350).
- [10] Y. Liu, Y. C. Gu, X. Y. Zhang, W. W. Wang, M. M. Cheng. Lightweight salient object detection via hierarchical visual perception learning. *IEEE Transactions on Cybernetics*, vol.51, no.9, pp.4439–4449, 2021. DOI: [10.1109/TCYB.2020.3035613](https://doi.org/10.1109/TCYB.2020.3035613).
- [11] Y. Liu, Y. H. Wu, Y. F. Ban, H. F. Wang, M. M. Cheng. Rethinking computer-aided tuberculosis diagnosis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Seattle, USA, pp.2643–2652, 2020. DOI: [10.1109/CVPR42600.2020.00272](https://doi.org/10.1109/CVPR42600.2020.00272).
- [12] Y. Liu, Y. H. Wu, P. S. Wen, Y. J. Shi, Y. Qiu, M. M. Cheng. Leveraging instance-, image- and dataset-level information for weakly supervised instance segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol.44, no.3, pp.1415–1428, 2022. DOI: [10.1109/TPAMI.2020.3023152](https://doi.org/10.1109/TPAMI.2020.3023152).
- [13] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, Long Beach, USA, pp.6000–6010, 2017. DOI: [10.5555/3295222.3295349](https://doi.org/10.5555/3295222.3295349).
- [14] J. Devlin, M. W. Chang, K. Lee, K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Minneapolis, USA, pp.4171–4186, 2019. DOI: [10.18653/v1/N19-1423](https://doi.org/10.18653/v1/N19-1423).
- [15] Z. H. Dai, Z. L. Yang, Y. M. Yang, J. Carbonell, Q. Le, R. Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy, pp.2978–2988, 2019. DOI: [10.18653/v1/P19-1285](https://doi.org/10.18653/v1/P19-1285).
- [16] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. H. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, N. Houlsby. An image is worth 16×16 words: Transformers for image recognition at scale. In *Proceedings of the 9th International Conference on Learning Representations*, 2021.
- [17] B. Heo, S. Yun, D. Han, S. Chun, J. Choe, S. J. Oh. Rethinking spatial dimensions of vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp.11916–11925, 2021. DOI: [10.1109/ICCV48922.2021.01172](https://doi.org/10.1109/ICCV48922.2021.01172).
- [18] Z. Liu, Y. T. Lin, Y. Cao, H. Hu, Y. X. Wei, Z. Zhang, S. Lin, B. N. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp.9992–10002, 2021. DOI: [10.1109/ICCV48922.2021.00986](https://doi.org/10.1109/ICCV48922.2021.00986).
- [19] W. H. Wang, E. Z. Xie, X. Li, D. P. Fan, K. T. Song, D. Liang, T. Lu, P. Luo, L. Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp.548–558, 2021. DOI: [10.1109/ICCV48922.2021.00061](https://doi.org/10.1109/ICCV48922.2021.00061).
- [20] W. J. Xu, Y. F. Xu, T. Chang, Z. W. Tu. Co-scale convolutional image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp.9961–9970, 2021. DOI: [10.1109/ICCV48922.2021.00983](https://doi.org/10.1109/ICCV48922.2021.00983).
- [21] H. Q. Fan, B. Xiong, K. Mangalam, Y. H. Li, Z. C. Yan, J. Malik, C. Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp.6804–6815, 2021. DOI: [10.1109/ICCV48922.2021.00675](https://doi.org/10.1109/ICCV48922.2021.00675).
- [22] D. Bolya, C. Y. Fu, X. L. Dai, P. Z. Zhang, C. Feichtenhofer, J. Hoffman. Token merging: Your ViT but faster. In *Proceedings of the 11th International Conference on Learning Representations*, Kigali, Rwanda, 2023.
- [23] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, vol.86, no.11, pp.2278–2324, 1998. DOI: [10.1109/5.726791](https://doi.org/10.1109/5.726791).
- [24] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. A. Ma, Z. H. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, L. Fei-Fei. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, vol.115, no.3, pp.211–252, 2015. DOI: [10.1007/s11263-015-0816-y](https://doi.org/10.1007/s11263-015-0816-y).
- [25] R. K. Srivastava, K. Greff, J. Schmidhuber. Highway networks, [Online], Available: <https://arxiv.org/abs/1505.00387>, 2015.
- [26] C. Szegedy, W. Liu, Y. Q. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Boston, USA, pp.1–9, 2015. DOI: [10.1109/CVPR.2015.7298594](https://doi.org/10.1109/CVPR.2015.7298594).
- [27] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Las Vegas, USA, pp.2818–2826, 2016. DOI: [10.1109/CVPR.2016.308](https://doi.org/10.1109/CVPR.2016.308).

- [28] C. Szegedy, S. Ioffe, V. Vanhoucke, A. A. Alemi. Inception-v4, Inception-ResNet and the impact of residual connections on learning. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, San Francisco, USA, pp. 4278–4284, 2017. DOI: [10.5555/3298023.3298188](https://doi.org/10.5555/3298023.3298188).
- [29] S. N. Xie, R. Girshick, P. Dollár, Z. W. Tu, K. M. He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, USA, pp. 5987–5995, 2017. DOI: [10.1109/CVPR.2017.634](https://doi.org/10.1109/CVPR.2017.634).
- [30] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, USA, pp. 2261–2269, 2017. DOI: [10.1109/CVPR.2017.243](https://doi.org/10.1109/CVPR.2017.243).
- [31] A. G. Howard, M. L. Zhu, B. Chen, D. Kalenichenko, W. J. Wang, T. Weyand, M. Andreetto, H. Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications, [Online], Available: <https://arxiv.org/abs/1704.04861>, 2017.
- [32] M. Sandler, A. Howard, M. L. Zhu, A. Zhmoginov, L. C. Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Salt Lake City, USA, pp. 4510–4520, 2018. DOI: [10.1109/CVPR.2018.00474](https://doi.org/10.1109/CVPR.2018.00474).
- [33] X. Y. Zhang, X. Y. Zhou, M. X. Lin, J. Sun. ShuffleNet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Salt Lake City, USA, pp. 6848–6856, 2018. DOI: [10.1109/CVPR.2018.00716](https://doi.org/10.1109/CVPR.2018.00716).
- [34] N. N. Ma, X. Y. Zhang, H. T. Zheng, J. Sun. ShuffleNet V2: Practical guidelines for efficient CNN architecture design. In *Proceedings of the 15th European Conference on Computer Vision*, Munich, Germany, pp. 122–138, 2018. DOI: [10.1007/978-3-030-01264-9_8](https://doi.org/10.1007/978-3-030-01264-9_8).
- [35] M. X. Tan, B. Chen, R. M. Pang, V. Vasudevan, M. Sandler, A. Howard, Q. V. Le. MnasNet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Long Beach, USA, pp. 2815–2823, 2019. DOI: [10.1109/CVPR.2019.00293](https://doi.org/10.1109/CVPR.2019.00293).
- [36] M. X. Tan, Q. V. Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, Long Beach, USA, pp. 6105–6114, 2019.
- [37] M. Jaderberg, K. Simonyan, A. Zisserman, K. Kavukcuoglu. Spatial transformer networks. In *Proceedings of the 28th International Conference on Neural Information Processing Systems*, Montreal, Canada, pp. 2017–2025, 2015. DOI: [10.5555/2969442.2969465](https://doi.org/10.5555/2969442.2969465).
- [38] L. Chen, H. W. Zhang, J. Xiao, L. Q. Nie, J. Shao, W. Liu, T. S. Chua. SCA-CNN: Spatial and channel-wise attention in convolutional networks for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, USA, pp. 6298–6306, 2017. DOI: [10.1109/CVPR.2017.667](https://doi.org/10.1109/CVPR.2017.667).
- [39] F. Wang, M. Q. Jiang, C. Qian, S. Yang, C. Li, H. G. Zhang, X. G. Wang, X. O. Tang. Residual attention network for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, USA, pp. 6450–6458, 2017. DOI: [10.1109/CVPR.2017.683](https://doi.org/10.1109/CVPR.2017.683).
- [40] J. Hu, L. Shen, S. Albanie, G. Sun, E. H. Wu. Squeeze-and-excitation networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 8, pp. 2011–2023, 2020. DOI: [10.1109/TPAMI.2019.2913372](https://doi.org/10.1109/TPAMI.2019.2913372).
- [41] S. Woo, J. Park, J. Y. Lee, I. S. Kweon. CBAM: Convolutional block attention module. In *Proceedings of the 15th European Conference on Computer Vision*, Munich, Germany, pp. 3–19, 2018. DOI: [10.1007/978-3-030-01234-2_1](https://doi.org/10.1007/978-3-030-01234-2_1).
- [42] J. Park, S. Woo, J. Y. Lee, I. S. Kweon. BAM: Bottleneck attention module. In *Proceedings of the British Machine Vision Conference*, Newcastle, UK, Article number 147, 2018.
- [43] X. Li, W. H. Wang, X. L. Hu, J. Yang. Selective kernel networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Long Beach, USA, pp. 510–519, 2019. DOI: [10.1109/CVPR.2019.00060](https://doi.org/10.1109/CVPR.2019.00060).
- [44] X. L. Wang, R. Girshick, A. Gupta, K. M. He. Non-local neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Salt Lake City, USA, pp. 7794–7803, 2018. DOI: [10.1109/CVPR.2018.00813](https://doi.org/10.1109/CVPR.2018.00813).
- [45] H. Zhang, C. R. Wu, Z. Y. Zhang, Y. Zhu, H. B. Lin, Z. Zhang, Y. Sun, T. He, J. Mueller, R. Manmatha, M. Li, A. Smola. ResNeSt: Split-attention networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, New Orleans, USA, pp. 2735–2745, 2022. DOI: [10.1109/CVPRW56347.2022.00309](https://doi.org/10.1109/CVPRW56347.2022.00309).
- [46] L. Yuan, Y. P. Chen, T. Wang, W. H. Yu, Y. J. Shi, Z. H. Jiang, F. E. H. Tay, J. S. Feng, S. C. Yan. Tokens-to-token ViT: Training vision transformers from scratch on ImageNet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp. 538–547, 2021. DOI: [10.1109/ICCV48922.2021.00060](https://doi.org/10.1109/ICCV48922.2021.00060).
- [47] H. Touvron, M. Cord, A. Sablayrolles, G. Synnaeve, H. Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp. 32–42, 2021. DOI: [10.1109/ICCV48922.2021.00010](https://doi.org/10.1109/ICCV48922.2021.00010).
- [48] D. Q. Zhou, B. Y. Kang, X. J. Jin, L. J. Yang, X. C. Lian, Z. H. Jiang, Q. B. Hou, J. S. Feng. DeepViT: Towards deeper vision transformer, [Online], Available: <https://arxiv.org/abs/2103.11886>, 2021.
- [49] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, H. Jégou. Training data-efficient image transformers & distillation through attention. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 10347–10357, 2021.
- [50] A. Srinivas, T. Y. Lin, N. Parmar, J. Shlens, P. Abbeel, A. Vaswani. Bottleneck transformers for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Nashville, USA, pp. 16514–16524, 2021. DOI: [10.1109/CVPR46437.2021.01625](https://doi.org/10.1109/CVPR46437.2021.01625).
- [51] I. O. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. H. Zhai, T. Unterthiner, J. Yung, A. Steiner, D. Keysers, J. Uszkoreit, M. Lucic, A. Dosovitskiy. MLP-Mixer: An all-MLP architecture for vision. In *Proceedings of the 34th Advances in Neural Information Processing Systems*, pp. 24261–24272, 2021.
- [52] H. X. Liu, Z. H. Dai, D. R. So, Q. V. Le. Pay attention to MLPs. In *Proceedings of the 34th Advances in Neural In-*

- formation Processing Systems, pp. 9204–9215, 2021.
- [53] Q. B. Hou, Z. H. Jiang, L. Yuan, M. M. Cheng, S. C. Yan, J. S. Feng. Vision Permutator: A permutable MLP-like architecture for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 1, pp. 1328–1334, 2023. DOI: [10.1109/TPAMI.2022.3145427](https://doi.org/10.1109/TPAMI.2022.3145427).
- [54] Z. C. Wang, Y. B. Hao, X. Y. Gao, H. Zhang, S. Wang, T. T. Mu, X. N. He. Parameterization of cross-token relations with relative positional encoding for vision MLP. In *Proceedings of the 30th ACM International Conference on Multimedia*, Lisboa, Portugal, pp. 6288–6299, 2022. DOI: [10.1145/3503161.3547953](https://doi.org/10.1145/3503161.3547953).
- [55] K. Han, A. Xiao, E. H. Wu, J. Y. Guo, C. J. Xu, Y. H. Wang. Transformer in transformer. In *Proceedings of the 35th Conference on Neural Information Processing Systems*, pp. 15908–15919, 2021.
- [56] Y. W. Li, K. Zhang, J. Z. Cao, R. Timofte, L. Van Gool. LocalViT: Bringing locality to vision transformers, [Online], Available: <https://arxiv.org/abs/2104.05707>, 2021.
- [57] K. Yuan, S. P. Guo, Z. W. Liu, A. J. Zhou, F. W. Yu, W. Wu. Incorporating convolution designs into visual transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp. 559–568, 2021. DOI: [10.1109/ICCV48922.2021.00062](https://doi.org/10.1109/ICCV48922.2021.00062).
- [58] D. Hendrycks, K. Gimpel. Gaussian error linear units (GELUs), [Online], Available: <https://arxiv.org/abs/1606.08415>, 2016.
- [59] S. Elfving, E. Uchibe, K. Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural Networks*, vol. 107, pp. 3–11, 2018. DOI: [10.1016/j.neunet.2017.12.012](https://doi.org/10.1016/j.neunet.2017.12.012).
- [60] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. M. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. J. Bai, S. Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Red Hook, USA, Article number 721, 2019. DOI: [10.5555/3454287.3455008](https://doi.org/10.5555/3454287.3455008).
- [61] H. Y. Zhang, M. Cissé, Y. N. Dauphin, D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *Proceedings of the 6th International Conference on Learning Representations*, Vancouver, Canada, 2018.
- [62] I. Loshchilov, F. Hutter. Decoupled weight decay regularization. In *Proceedings of the 7th International Conference on Learning Representations*, New Orleans, USA, 2019.
- [63] I. Loshchilov, F. Hutter. SGDR: Stochastic gradient descent with warm restarts. In *Proceedings of the 5th International Conference on Learning Representations*, Toulon, France, 2017.
- [64] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. M. He, P. Dollár. Designing network design spaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Seattle, USA, pp. 10425–10433, 2020. DOI: [10.1109/CVPR42600.2020.01044](https://doi.org/10.1109/CVPR42600.2020.01044).
- [65] H. P. Wu, B. Xiao, N. Codella, M. C. Liu, X. Y. Dai, L. Yuan, L. Zhang. CvT: Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp. 22–31, 2021. DOI: [10.1109/ICCV48922.2021.00009](https://doi.org/10.1109/ICCV48922.2021.00009).
- [66] X. X. Chu, Z. Tian, Y. Q. Wang, B. Zhang, H. B. Ren, X. L. Wei, H. X. Xia, C. H. Shen. Twins: Revisiting the design of spatial attention in vision transformers. In *Proceedings of the 34th Advances in Neural Information Processing Systems*, pp. 9355–9366, 2021.
- [67] W. H. Wang, E. Z. Xie, X. Li, D. P. Fan, K. T. Song, D. Liang, T. Lu, P. Luo, L. Shao. PVT v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, vol. 8, no. 3, pp. 415–424, 2022. DOI: [10.1007/s41095-022-0274-8](https://doi.org/10.1007/s41095-022-0274-8).
- [68] A. Kirillov, R. Girshick, K. M. He, P. Dollár. Panoptic feature pyramid networks. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*, Long Beach, USA, pp. 6392–6401, 2019. DOI: [10.1109/CVPR.2019.00656](https://doi.org/10.1109/CVPR.2019.00656).
- [69] B. L. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, A. Torralba. Scene parsing through ADE20K dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, USA, pp. 5122–5130, 2017. DOI: [10.1109/CVPR.2017.544](https://doi.org/10.1109/CVPR.2017.544).
- [70] MMSegmentation Contributors. MMSegmentation: OpenMMLab semantic segmentation toolbox and benchmark, [Online], Available: <https://github.com/open-mmlab/mmdetection>, 2020.
- [71] T. Y. Lin, P. Goyal, R. Girshick, K. M. He, P. Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, Venice, Italy, pp. 2999–3007, 2017. DOI: [10.1109/ICCV.2017.324](https://doi.org/10.1109/ICCV.2017.324).
- [72] T. Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick. Microsoft COCO: Common objects in context. In *Proceedings of the 13th European Conference on Computer Vision*, Zürich, Switzerland, pp. 740–755, 2014. DOI: [10.1007/978-3-319-10602-1](https://doi.org/10.1007/978-3-319-10602-1).
- [73] K. Chen, J. Q. Wang, J. M. Pang, Y. H. Cao, Y. Xiong, X. X. Li, S. Y. Sun, W. S. Feng, Z. W. Liu, J. R. Xu, Z. Zhang, D. Z. Cheng, C. C. Zhu, T. H. Cheng, Q. J. Zhao, B. Y. Li, X. Lu, R. Zhu, Y. Wu, J. F. Dai, J. D. Wang, J. P. Shi, W. L. Ouyang, C. C. Loy, D. H. Lin. MMDetection: Open MMLab detection toolbox and benchmark, [Online], Available: <https://arxiv.org/abs/1906.07155>, 2019.
- [74] P. C. Zhang, X. Y. Dai, J. W. Yang, B. Xiao, L. Yuan, L. Zhang, J. F. Gao. Multi-scale vision Longformer: A new vision transformer for high-resolution image encoding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, Montreal, Canada, pp. 2978–2988, 2021. DOI: [10.1109/ICCV48922.2021.00299](https://doi.org/10.1109/ICCV48922.2021.00299).



Yun Liu received the B.Eng. and Ph.D. degrees in computer science from Nankai University, China in 2016 and 2020, respectively. Then, he worked with Prof. Luc Van Gool for one and a half years as a postdoctoral scholar at Computer Vision Lab, ETH Zürich, Switzerland. Currently, he is a senior scientist at Institute for Infocomm Research (I2R), A*STAR, Singa-

pore.

His research interests include computer vision and machine learning (especially deep learning).

E-mail: vagrantlyun@gmail.com

ORCID iD: 0000-0001-6143-0264



Yu-Huan Wu received the Ph.D. degree in computer science from Nankai University, China in 2022, advised by Prof. Ming-Ming Cheng. He is a scientist at the Institute of High Performance Computing (IHPC), A*STAR, Singapore. He has published 10+ papers on top-tier conferences and journals such as IEEE TPAMI/TIP/CVPR/ICCV.

His research interests include computer vision and medical imaging.

E-mail: wu_yuhuan@ihpc.a-star.edu.sg

ORCID iD: 0000-0001-8666-3435



Guolei Sun received the M.Sc. degree in computer science from King Abdullah University of Science and Technology, Saudi Arabia in 2018. From 2018 to 2019, he worked as a research engineer at the Inception Institute of Artificial Intelligence, UAE. Currently, he is a Ph.D. degree candidate at ETH Zürich, Switzerland under supervision of Prof. Luc Van Gool. He has

published more than 20 papers in top journals and conferences such as TPAMI, CVPR, ICCV, and ECCV.

His research interests include computer vision and deep learning for tasks such as semantic segmentation, video understanding, and object counting.

E-mail: sunguolei.kaust@gmail.com (Corresponding author)

ORCID iD: 0000-0001-8667-9656



Le Zhang received the Ph.D. degree in electrical and electronic engineering from Nanyang Technological University (NTU), Singapore in 2016. He is a professor with the School of Information and Communication Engineering, University of Electronic Science and Technology of China (UESTC), China. From 2016 to 2018, he was a postdoc fellow at Advanced Digital

Sciences Center (ADSC), Singapore. From 2018 to 2021, he was

a research scientist at the Institute for Infocomm Research (I2R), A*STAR, Singapore. He is an Associate Editor of *Neural Networks*, *Neurocomputing*, and *IET Biometrics*.

His research interests include computer vision and machine learning.

E-mail: lezhang@uestc.edu.cn (Corresponding author)

ORCID iD: 0000-0002-6930-8674



Ajad Chhatkuli received the M.Sc. degree in computer vision from the University of Burgundy, France in 2013, and the Ph.D. degree in computer vision from the University of Clermont Auvergne, France in 2017 under the supervision of Prof. Adrien Bartoli and Dr. Daniel Pizarro. He is currently a postdoctoral researcher supervised by Prof. Luc Van Gool

at ETH Zürich, Switzerland.

His research interests include template-based and template-free non-rigid 3D reconstruction.

E-mail: ajad.chhatkuli@vision.ee.ethz.ch

ORCID iD: 0000-0003-2051-2209



Luc Van Gool received the B.Eng. degree in electromechanical engineering from the Katholieke Universiteit Leuven, Belgium in 1981. Currently, he is a professor at the Katholieke Universiteit Leuven, Belgium, and the ETH in Zürich, Switzerland. He leads computer vision research at both places, and also teaches at both. He has been a program committee member of sev-

eral major computer vision conferences. He received several Best Paper awards, won a David Marr Prize and a Koenderink Award, and was nominated Distinguished Researcher by the IEEE Computer Science Committee. He is a co-founder of 10 spin-off companies.

His research interests include 3D reconstruction and modeling, object recognition, tracking, gesture analysis, and the combination of those.

E-mail: vangool@vision.ee.ethz.ch

ORCID iD: 0000-0002-3445-5711