# Visformer: The Vision-friendly Transformer

Zhengsu Chen, Lingxi Xie, Jianwei Niu, *Senior Member*, *IEEE*,
Xuefeng Liu, Longhui Wei, Qi Tian, *Fellow*, *IEEE*

*Abstract*—The past few years have witnessed the rapid development of applying the Transformer module to vision problems. While some researchers have demonstrated that Transformer-based models enjoy a favorable ability of fitting data, there are still growing number of evidences showing that these models suffer over-fitting especially when the training data is limited. This paper offers an empirical study by performing step-by-step operations to gradually transit a Transformer-based model to a convolution-based model. The results we obtain during the transition process deliver useful messages for improving visual recognition. Based on these observations, we propose a new architecture named Visformer, which is abbreviated from the 'Vision-friendly Transformer'. With the same computational complexity, Visformer outperforms both the Transformer-based and convolution-based models in terms of ImageNet classification and object detection performance, and the advantage becomes more significant when the model complexity is lower or the training set is smaller. The code is available at https://github.com/danczs/Visformer.

*Index Terms*—vision-friendly Transformer, vision Transformer, convolutional neural network, image recognition.

## I. INTRODUCTION

IN the past decade, convolution used to play a central role in the deep learning models [1]–[4] for visual recognition. This situation starts to change when the Transformer [5], a module that originates from natural language processing [5]–[7], is transplanted to the vision scenarios. It was shown in the ViT model [8] that an image can be partitioned into a grid of patches and the Transformer is directly applied upon the grid as if each patch is a visual word. ViT requires a large amount of training data (*e.g.*, the ImageNet-21K [9] or the JFT-300M dataset), arguably because the Transformer is equipped with long-range attention and interaction, and is prone to over-fitting. The follow-up efforts [10] improved ViT to some extent, but these models still perform badly especially under limited training data or moderate data augmentation compared with convolution-based models.

On the other hand, vision Transformers can achieve much better performance than convolution-based models when

Zhengsu Chen and Xuefeng Liu are with the School of Computer Science and Engineering, Beihang University, 100191, Beijing, China (e-mail: danczs@buaa.edu.cn; liu_xuefeng@buaa.edu.cn)

Lingxi Xie is with Johns Hopkins University, Baltimore, USA (e-mail: 198808xc@gmail.com)

Jianwei Niu is with Hangzhou Innovation Institute of Beihang University, Zhengzhou University, and Beihang University, 100191, Beijing, China (e-mail: @niujianwei@buaa.edu.cn)

Longhui wei is with the University of Science and Technology of China, 230026, Hefei, China (e-mail:longhuiwei@pku.edu.cn)

Qi Tian is with Xidian University, 710126, Xi'an, China (e-mail: wywq-tian@gmail.com)

TABLE I
THE COMPARISON AMONG RESNET-50, DEIT-S, AND THE PROPOSED VISFORMER-S MODEL ON IMAGENET.

| Network | | ResNet-50 | DeiT-S | Visformer-S |
|---|---|---|---|---|
| FLOPs (G) | | 4.1 | 4.6 | 4.9 |
| Parameters (M) | | 25.6 | 21.8 | 40.2 |
| Full data | base setting | 77.43 | 63.12 | 77.20 |
| | elite setting | 78.73 | 80.07 | 82.19 |
| Part of data | 10% labels | 58.37 | 40.41 | 58.74 |
| | 10% classes | 89.90 | 80.06 | 90.06 |

trained with large amount of data. Namely, vision Transformers have higher 'upper-bound' while convolution-based models are better in 'lower-bound'. Both upper-bound and lower-bound are important properties for neural networks. Upper-bound is the potential to achieve higher performance and lower-bound enables networks to perform better when trained with limited data or scaled to different complexity.

Based on the observation of lower-bound and upper-bound on Transformer-based and convolution-based networks, the main goal of this paper is to identify the reasons behind the difference, by which we can design networks with higher lower-bound and upper-bound. The gap between Transformer-based and convolution-based networks can be revealed with two different training settings on ImageNet. The first one is the base setting. It is the standard setting for convolution-based models, *i.e.*, the training schedule is shorter and the data augmentation only contains basic operators such as random-size cropping [11] and flipping. The performance under this setting is called **base performance** in this paper. The other one is the training setting used in [10]. It is carefully tuned for Transformer-based models, *i.e.*, the training schedule is longer and the data augmentation is stronger (*e.g.*, RandAugment [12], CutMix [13], *etc.*, have been added). We use the **elite performance** to refer to the accuracy produced by it.

We take DeiT-S [10] and ResNet-50 [4] as the examples of Transformer-based and convolution-based models. As shown in Table I, Deit-S and ResNet-50 employ comparable FLOPs and parameters. However, they behave very differently trained on the full data under these two settings. Deit-S has higher elite performance, but changing the setting from elite to base can cause a $10\%+$ accuracy drop for DeiT-S. ResNet-50 performs much better under the base setting, yet the improvement for the elite setting is merely $1.3\%$. This motivates us to study the difference between these models. With these two settings, we can roughly estimate the lower-bound and upper-bound of the models. The methodology we use is to perform step-by-step operations to gradually transit one model into another, by

which we can identify the properties of modules and designs in these two networks. The entire transition process, taking a total of 8 steps, is illustrated in Figure 1.

Specifically, from DeiT-S to ResNet-50, one should (i) use global average pooling (not the classification token), (ii) introduce step-wise patch embeddings (not large patch flattening), (iii) adopt the stage-wise backbone design, (iv) use batch normalization [14] (not layer normalization [15]), (v) leverage $3 \times 3$ convolutions, (vi) discard the position embedding scheme, (vii) replace self-attention with convolution, and finally (viii) adjust the network shape (*e.g.*, depth, width, *etc.*). After a thorough analysis on the reasons behind the results, we absorb all the factors that are helpful to visual recognition and derive the **Visformer**, *i.e.*, the Vision-friendly Transformer.

Evaluated on ImageNet classification, Visformer claims better performance than the competitors, DeiT and ResNet, as shown in Table I. With the elite setting, the Visformer-S model outperforms DeiT-S and ResNet-50 by $2.12\%$ and $3.46\%$, respectively, under a comparable model complexity. Different from Deit-S, Visformer-S also survives two extra challenges, namely, when the model is trained with 10% labels (images) and 10% classes. Visformer-S even performs better than ResNet-50, which reveals the high lower-bound of Visformer-S. Additionally, for tiny models, Visformer-Ti significantly outperforms Deit-Ti by more than 6%.

The contribution of this paper is three-fold. **First**, for the first time, we introduce the lower-bound and upper-bound to investigate the performance of Transformer-based vision models. **Second**, we close the gap between the Transformer-based and convolution-based models by a gradual transition process and thus identify the properties of the designs in the Transformer-based and convolution-based models. **Third**, we propose the Visformer as the final model that achieves satisfying lower-bound and upper-bound.

The preliminary version of this paper appeared as [16]. In the extended version, we further explore the recently proposed work and provide more experiments and analysis. The main improvements over the preliminary version are summarized as follows:

- We optimize the architecture of Visformer according to the experimental observations and propose VisformerV2 which substantially outperforms the old version.
- We analyze the overflow problem when utilizing half-precision in Transformers and propose an efficient method to avoid overflow without degrading the performance.
- We generalize Visformer to downstream vision tasks and observe consistent improvements.

## II. RELATED WORK

Image classification is a fundamental task in computer vision. In the deep learning era, the most popular method is to use deep neural networks [2], [4], [17]. One of the fundamental units to build such networks is convolution, where a number of convolutional kernels are used to capture repeatable local patterns in the input image and intermediate data. To reduce the computational costs as well as alleviate the risk of over-fitting, it was believed that the convolutional kernels should be of a small size, *e.g.*, $3 \times 3$. However, this brings the difficulty for faraway contexts in the image to communicate with each other – this is partly the reason that the number of layers has been increasing. Despite stacking more and more layers, researchers consider another path which is to use attention-based approaches to ease the propagation of visual information.

Since Transformers achieved remarkable success in natural language processing (NLP) [5]–[7], many efforts have been made to introduce Transformers to vision tasks. These works mainly fall into two categories. The first category consists of pure attention models [8], [10], [18]–[22]. These models usually only utilize self-attention and attempt to build vision models without convolutions. However, it is computationally expensive to relate all pixels with self-attention for realistic full-sized images. Thus, there has some interest in forcing self-attention to only concentrate on the pixels in local neighborhoods (*e.g.*, SASA [18], LRNet [19], SANet [20]). These methods replace convolutions with local self-attentions to learn local relations and achieve promising results. However, it requires complex engineering to efficiently apply self-attention to every local region in an image. Another way to solve the complexity problem is to apply self-attention to reduced resolution. These methods either reduce the resolution and color space first [21] or regard image patches rather pixels as tokens (*i.e.*, words) [8], [10]. However, resolution reduction and patch flattening usually make it more difficult to utilize the local prior in natural images. Thus, these methods usually obtain suboptimal results [21] or require huge dataset [8] and heavy augmentation [10].

The second category contains the networks built with not only self-attentions but also convolutions. Self-attention was first introduced to CNNs by non-local neural networks [23]. These networks aim to capture global dependencies in images and videos. Note that non-local neural networks are inspired by the classical non-local method in vision tasks [24] and unlike those in Transformers, the self-attentions in non-local networks are usually not equipped with multi-heads and position embedding [23], [25], [26]. Afterwards, Transformers achieve remarkable success in NLP tasks [6], [7] and, therefore, self-attentions that inherits NLP settings (*e.g.*, multi-heads, position encodings, classification token, *etc.*) are combined with convolutions to improve vision tasks [18], [27], [28]. A common combination is to utilize convolutions first and apply self-attention afterwards [8], [29]. [8] builds hybrids of self-attention and convolution by adding a ResNet backbone before Transformers. Afterwards, more and more methods are proposed to combine self-attention and convolution. Besides utilizing convolution in early layers [30], BotNet [29] designs bottleneck cells for self-attention. Conformer [31] fuses the feature of a convolution neural network and a Transformer with the feature coupling unit to combine the global and local representations. CvT [32] introduces convolution to the feature projection of self-attention, by which the query, key and value in the self-attention can capture the local information. CoAt-Net [33] unifies depthwise convolution with self-attention and vertically stacks convolution layers and self-attention layers to improve the generalization, capacity and efficiency. However,

these methods usually combine convolution and self-attention empirically or heuristically. Our method, in contrast, explores the full process of converting a Transformer to a convolution neural network.

There is also work that studies scaling or training vision Transformer. CaiT [34] find that it is very efficient to scale up vision Transformer in depth dimension with LayerScale. AutoFormer [35] builds a super Transformer network by which they can evaluate the different designs and search efficient Transformer architecture. Zhai *et al.* [36] observe that most vision Transformers can benefit from increasing compute resources and larger dataset. They suggest scaling up compute, data and model together. Steiner *et al.* [37] further study data, augmentation and regularization in vision Transformer and finds that augmentation can yield the same performance as that trained on an order of magnitude more data.

Additionally, self-attention has been used in many downstream vision tasks (detection [38], segmentation [39]) and low vision tasks [40]. These tasks usually utilize much larger input resolution than classification. For example, the frameworks in COCO [41] usually utilize $1280 \times 800$ inputs. This is a critical problem for vision Transformer, since the complexity increases quadratically with pixel numbers. The widely used solution is adopting sliding windows to capture local patterns and building extra pipelines for information exchange among the windows. Swin Transformer [42] shifts the windows alternately in different layers, by which the tokens can build long-distance relations as the depth increases. CSWin Transformer [43] further develops the cross-shaped self-attention mechanism to ensure that the windows can access the global feature in one dimension. MSG-Transformer [44], by contrast, exchanges the local information by messenger tokens.

## III. METHODOLOGY

### A. *Transformer-based and convolution-based visual recognition models*

Recognition is the fundamental task in computer vision. This work mainly considers image classification, where the input image is propagated through a deep network to derive the output class label. Most deep networks are designed in a hierarchical manner and composed of a series of layers.

We consider two popular layers named convolution and Transformer. Convolution originates from the intuition to capture local patterns which are believed more repeatable than global patterns. It uses a number of learnable kernels to compute the responses of the input to different patterns, for which a sliding window is moved along both axes of the input data and the inner-product between the data and kernel is calculated. In this paper, we constrain our study in the scope of residual blocks, a combination of 2 or 3 convolutional layers and a skip-connection. Non-linearities such as activation and normalization are inserted between the neighboring convolutional layers.

On the other hand, Transformer originates from natural language processing and aims to frequently formulate the relationship between *any* two elements (called tokens) even when they are far from each other. This is achieved by generating three features for each token, named the query, key, and value, respectively. Then, the response of each token is calculated as a weighted sum over all the values, where the weights are determined by the similarity between its query and the corresponding keys. This is often referred to as multi-head self-attention (MHSA), followed by other operations including normalization and linear mapping.

Throughout the remaining part, we consider DeiT-S [10] and ResNet-50 [4] as the representative of Transformer-based and convolution-based models, respectively. Besides the basic building block, there are also differences in design, *e.g.*, ResNet-50 has a few down-sampling layers that partition the model into stages, but the number of tokens remains unchanged throughout DeiT-S. The impact of these details will be elaborated in Section III-C.

### B. *Settings: The base and elite performance*

Although DeiT-S reports a $80.1\%$ accuracy which is higher than $78.7\%$ of ResNet-50, we notice that DeiT-S has changed the training strategy significantly, *e.g.*, the number of epochs is enlarged by more than $3\times$ and the data augmentation becomes much stronger. Interestingly, DeiT-S seems to heavily rely on the carefully-tuned training strategy, and other Transformer-based models including ViT [8] and PIT [40] also reported their dependency on other factors, *e.g.*, a large-scale training set. In what follows, we provide a comprehensive study on this phenomenon.

We evaluate all classification models on the ImageNet dataset [45] which has 1K classes, 1.28M training images and 50K testing images. Each class has roughly the same number of training images. This is one of the most popular datasets for visual recognition.

There are two settings to optimize each recognition model. The first one is named the **base setting** which is widely adopted by convolution-based networks. Specifically, the model is trained for 90 epochs with the SGD optimizer. The learning rate starts with $0.2$ for batch size $512$ and gradually decays to $0.00001$ following the cosine annealing function. A moderate data augmentation strategy with random-size cropping [11] and flipping is used. The second one is named the **elite setting** which has been verified effective to improve the Transformer-based models. The Adamw optimizer with an initial learning rate of $0.0005$ for batch size $512$ is used. The data augmentation and regularization strategy is made much stronger to avoid over-fitting, for which intensive operations including RandAugment [12], Mixup [46], Cut-Mix [13], Random Erasing [47], Repeated Augmentation [48], [49] and Stochastic Depth [50] are used. Correspondingly, the training lasts 300 epochs, much longer than that of the base setting.

Throughout the remaining part of this paper, we refer to the classification accuracy under the base and elite settings as **base performance** and **elite performance**, respectively. We expect the numbers to provide complementary views for us to understand the studied models.
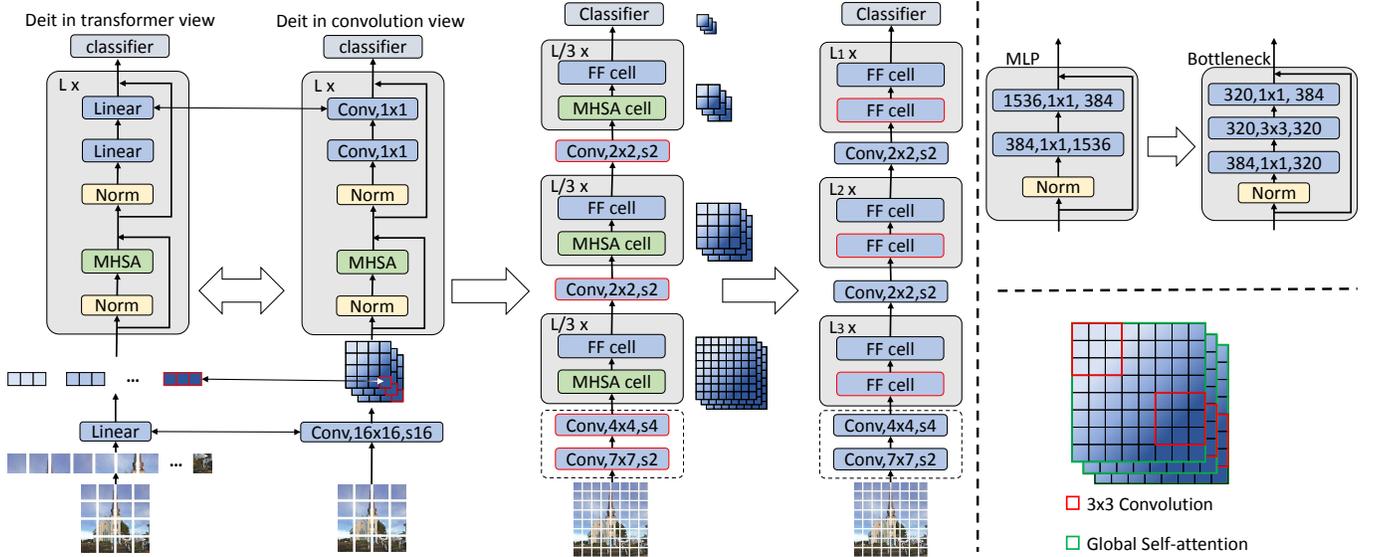
Fig. 1. The transition process that starts with DeiT and ends with ResNet-50. To save space, we only show three important movements. The first movement converts DeiT from the Transformer to convolution view (Section III-C1). The second movement replaces the patch flattening module with step-wise patch embedding (elaborated in Section III-C2) and introduces the stage-wise design (Section III-C3) . The third movement replaces the self-attention module with convolution (Section III-C7). The upper-right area shows a relatively minor modifications, inserting $3 \times 3$ convolution (Section III-C5). The lower-right area compares the receptive fields of a $3 \times 3$ convolution and self-attention. This figure is best viewed in color.

## C. The transition from DeiT-S to ResNet-50

This subsection displays a step-by-step process in which we gradually transit a model from DeiT-S to ResNet-50. There are eight steps in total. The key steps are illustrated in Figure 1, and the results, including the base and elite performance and the model statistics, are summarized in Table II.

*1) Using global average pooling to replace the classification token:* The first step of the transition is to remove the classification token and add global average pooling to the Transformer-based models. Unlike the convolution-based models, Transformers usually add a classification token to the inputs and utilize the corresponding output token to perform classification, which is inherited from NLP tasks [6]. As a contrast, the classification features in convolution-based models are obtained by conducting global average pooling in the space dimension.

By removing the classification token, the Transformer can be equivalent translated to the convolutional version as shown in Figure 1. Specifically, the patch embedding operation is equivalent to a convolution whose kernel size and stride is the patch size [8]. The shape of the intermediate features can be naturally converted from a sequence of tokens (*i.e.*, words) to a bundle feature maps and the tokens become the vector in channel dimension (illustrated in Figure 1). The linear layers in MHSA and MLP blocks are equivalent to $1 \times 1$ convolutions.

The performance of the obtained network (Net1) is shown in Table II. As can be seen, this transition can substantially improve the base performance. Our further experiments show that adding global pooling itself can improve the base performance from 64.17% to 69.44%. In other words, the global average pooling operation which is widely used in convolution-based models since NIN [51], enables the network to learn more

efficiently under moderate augmentation. Furthermore, this transition can slightly improve the elite performance.

*2) Replacing patch flattening with step-wise patch embedding:* DeiT and ViT models directly encode the image pixels with a patch embedding layer which is equivalent to a convolution with a large kernel size and stride (*e.g.*, 16). This operation flattens the image patches to a sequence of tokens so that Transformers can handle images. However, patch flattening impairs the position information within each patch and makes it more difficult to extract the patterns within patches. To solve this problem, existing methods usually attach a preprocessing module before patch embedding. The preprocessing module can be a feature extraction convnet [8] or a specially designed Transformer [52].

We found that there is a rather simple solution, which is factorizing the large patch embedding to step-wise small patch embeddings. Specifically, We first add the stem layer in ResNet to the Transformer, which is a $7 \times 7$ convolution layer with a stride of two. The stem layer can be seen as a $2 \times 2$ patching embedding operation with pixel overlap (*i.e.*, $7 \times 7$ kernel size). Since the patch size in the original DeiT model is 16, we still need to embed $8 \times 8$ patches after the stem. We further factorize the $8 \times 8$ patch embedding to a $4 \times 4$ embedding and a $2 \times 2$ embedding, which are $4 \times 4$ and $2 \times 2$ convolution layers with stride 4 and 2 in the perspective of convolution. Additionally, we add an extra $2 \times 2$ convolution to further upgrade the patch size from $16 \times 16$ to $32 \times 32$ before classification. These patch embedding layers can also be seen as the down-sampling layers and we double the channel numbers after embedding following the practice in convolution-based models.

By utilizing step-wise embeddings, the position prior within patches is encoded into features. As a result, the model can

TABLE II

THE CLASSIFICATION ACCURACY ON IMAGENET DURING THE TRANSITION PROCEDURE FROM DEIT-S TO RESNET-50. BOTH THE BASE SETTING AND THE ELITE SETTING ARE CONSIDERED (FOR THE DETAILS, SEE SECTION III-B), AND WE MARK THE POSITIVE MODIFICATIONS IN RED AND THE NEGATIVE MODIFICATIONS IN BLUE. NOTE THAT A MODIFICATION CAN IMPACT THE BASE AND ELITE PERFORMANCE DIFFERENTLY. THOUGH THE NUMBER OF PARAMETERS INCREASES CONSIDERABLY AT THE INTERMEDIATE STATUS, THE COMPUTATIONAL COSTS MEASURED BY FLOPS DOES NOT CHANGE SIGNIFICANTLY.

| Model Name | **added** | **removed** | **base** perf. | **elite** perf. | FLOPs (G) | Params (M) |
|---|---|---|---|---|---|---|
| DeiT-S | | – | 64.17 | 80.07 | 4.60 | 22.1 |
| Net1 | global average pooling | classification token | 69.81 (+5.64) | 80.16 (+0.09) | 4.57 | 22.0 |
| Net2 | step-wise embeddings | large patch embedding | 73.01 (+3.20) | 81.35 (+1.19) | 4.77 | 23.9 |
| Net3 | stage-wise design | – | 75.76 (+2.75) | 80.19 (-1.14) | 4.79 | 39.5 |
| Net4 | batch norm | layer norm | 76.49 (+0.73) | 80.97 (+0.78) | 4.79 | 39.5 |
| Net5 | $3 \times 3$ convolution | – | 77.37 (+0.88) | 80.15 (-0.82) | 4.76 | 39.2 |
| Net6 | – | position embedding | 77.31 (-0.06) | 79.86 (-0.29) | 4.76 | 39.0 |
| Net7 | convolution | self-attention | 76.24 (-1.07) | 79.01 (-0.85) | 4.83 | 45.0 |
| ResNet-50 | network shape adjustment | | 77.43 (+1.19) | 78.73 (-0.28) | 4.09 | 25.6 |

learn patterns more efficiently. As can be seen in Table II, this transition can significantly improve the base performance and elite performance of the network. It indicates that step-wise embedding is a better choice than larger patch embedding in Transformer-based models. Additionally, this transition is computationally efficient and only introduces about $4\%$ extra FLOPs.

*3) Stage-wise design:* In this section, we split networks into stages like ResNets. The blocks in the same stage share the same feature resolution. Since step-wise embeddings in the last transition have split the network into different stages, the transition in this section is to reassign the blocks to different stages as shown in Figure 1. However, unlike convolution blocks, the complexity of self-attention blocks increases by $O(N^4)$ with respect to the feature size. Thus we only insert blocks to the $8 \times 8$, $16 \times 16$ and $32 \times 32$ patch embedding stages, which correspond to $28 \times 28$, $14 \times 14$ and $7 \times 7$ feature resolutions respectively for $224 \times 224$ inputs. Additionally, we halve the head dimension and feature dimension before self-attention in $28 \times 28$ stage to ensure that the blocks in different stages utilize similar FLOPs.

This transition leads to interesting results. The base performance is further improved. It is conjectured that the stage-wise design leverages the image local priors and thus can perform better under moderate augmentation. However, the elite performance of the network decreases markedly. To study reasons, we conduct ablation experiments and find that self-attention does not work well in very large resolutions. We conjecture that large resolution contains too many tokens and it is much more difficult for self-attention to learn relations among them. We will detail it in section III-D.

*4) Replacing LayerNorm with BatchNorm:* Transformer-based models usually normalize the features with Layer-Norm [15], which is inherited from NLP tasks [5], [6]. As a contrast, convolution-based models like ResNets usually utilize BatchNorm [14] to stabilize the training process. LayerNorm is independent of batch size and more friendly for specific tasks compared with BatchNorm, while BatchNorm usually can achieve better performance given appropriate batch size [53]. We replace all the LayerNorm layers with BatchNorm layers and the results show that BatchNorm performs better than LayerNorm. It can improve both the base

performance and elite performance of the network.

In addition, we also try to add BatchNorm to Net2 to further improve the elite performance. However, this Net2-BN network suffers from convergence problems. This may explain why BatchNorm is not widely used in the pure self-attention models. But for our mixed model, BatchNorm is a reliable method to advance performance.

*5) Introducing $3 \times 3$ convolutions:* Since the tokens of the network are present as feature maps, it is natural to introduce convolutions with kernel sizes larger than $1 \times 1$. The specific meaning of large kernel convolution is illustrated at the bottom right of Figure 1. When global self-attentions attempt to build the relations among all the tokens (*i.e.*, pixels), convolutions focus on relating the tokens within local neighborhoods. We chose to insert $3 \times 3$ convolutions between the $1 \times 1$ convolutions in feed-forward blocks, which transforms the MLP blocks into bottleneck blocks as exhibited at the top right of Figure 1. Note that the channel numbers of the $3 \times 3$ convolution layers are tuned to ensure that the FLOPs of the feed-forward blocks are nearly unchanged. The obtained bottleneck blocks are similar to the bottleneck blocks in ResNet-50, although they have different bottleneck ratios (*i.e.*, the factor of reducing the channel numbers before the $3 \times 3$ convolution). We replace the MLP blocks with bottleneck blocks in all three stages.

Not surprisingly, $3 \times 3$ convolutions which can leverage the local priors in images further improve the network base performance. The base performance (77.37%) becomes comparable with ResNet-50 (77.43%). However, the elite performance decreases by 0.82%. We conduct more experiments to study the reasons. Instead of adding $3 \times 3$ convolutions to all stages, we insert $3 \times 3$ convolutions to different stages separately. We observe that $3 \times 3$ convolutions only work well on the high-resolution features. We conjecture that leveraging local relations is important for the high-resolution features in natural images. For the low-resolution features, however, local convolutions become unimportant when equipped with global self-attention. We will detail it in section III-D.

*6) Removing position embedding:* In Transformer-based models, position embedding is proposed to encode the position information inter tokens. In the transition network, we utilize learnable position embedding as in [6] and add them to

features after patch embeddings. To approaching ResNet-50, position embedding should be removed.

The results are exhibited in Table II. The base performance is almost unchanged and the elite performance declines slightly (0.29%). As a comparison, We test to remove the position embedding of DeiT-S and elite performance decreases significantly by 3.95%. It reveals that position embedding is less important in the transition model than that in the pure Transformer-based models. It is because that the position prior inter tokens is preserved by the feature maps and convolutions with spatial kernels can encode and leverage it. Consequently, the harm of removing position embedding is remarkably reduced in the transition network. It also explains why convolution-based models do not need position embedding.

*7) Replacing self-attention with feed-forward:* In this section, we remove the self-attention blocks in each stage and utilize a feed-forward layer instead, so that the network becomes a pure convolution-based network. To keep the FLOPs unchanged, several bottleneck blocks are added to each stage. After the replacement, the obtained network consists of bottleneck blocks like ResNet-50.

The performance of the obtained network (Net7) is shown in Table II. The pure convolution-based network performs much worse both in base performance and elite performance. *It indicates that self-attentions do drive neural networks to higher elite performance and is not responsible for the poor base performance in ViT or DeiT. It is possible to design a self-attention network with high base performance and elite performance.*

*8) Adjusting the shape of network:* There are still many differences between Net7 and ResNet-50. First, the shape of Net7 is different from ResNet-50. Their depths, widths, bottleneck ratios and block numbers in network stages are different. Second, they normalize the features in different positions. Net7 only normalizes input features in a block, while ResNet-50 normalizes features after each convolutional layer. Third, ResNet-50 down-samples the features with bottleneck blocks but Net7 utilizes a single convolution layer (*i.e.*, patch embedding layer). In addition, Net7 employs a few more FLOPs. Nevertheless, both these two networks are convolution-based networks. The performance gap between these two networks can be attributed to architecture design strategy.

As shown in Table II, the base performance is improved after transition. It demonstrates that ResNet-50 has better network architecture and can perform better with fewer FLOPs. However, ResNet-50 obtains worse elite performance. It indicates that the inconsistencies between base performance and elite performance exist not only in self-attention models but also in pure convolution-based networks.

### D. Summary: the Visformer model

We aim to build a network with high base performance and elite performance. The transition study has shown that there are some inconsistencies between base performance and elite performance. The first problem is the stage-wise design,

#### TABLE III
IMPACT OF REPLACING THE SELF-ATTENTION BLOCKS WITH THE BOTTLENECK BLOCKS IN EACH STAGE OF NET5. THESE EXPERIMENTS ARE PERFORMED INDIVIDUALLY.

| Network | base perf.(%) | elite perf.(%) |
| --- | --- | --- |
| Net5 | 77.37 | 80.15 |
| Net5-DS1 | 77.29 (-0.08) | 80.13 (-0.02) |
| Net5-DS2 | 77.34 (-0.02) | 79.75 (-0.40) |
| Net5-DS3 | 77.05 (-0.32) | 79.59 (-0.56) |

#### TABLE IV
IMPACT OF REPLACING THE MLP LAYERS WITH THE BOTTLENECK BLOCKS IN EACH STAGE OF NET4. THESE EXPERIMENTS ARE PERFORMED INDIVIDUALLY.

| Network | base perf.(%) | elite perf.(%) |
| --- | --- | --- |
| Net4 | 76.49 | 80.97 |
| Net4-S1 | 77.02 (+0.53) | 81.10 (+0.13) |
| Net4-S2 | 76.55 (+0.06) | 80.50 (-0.47) |
| Net4-S3 | 76.82 (+0.33) | 80.44 (-0.53) |
| Net5 | 77.37 (+0.88) | 80.15 (-0.82) |

which increases the base performance but decreases the elite performance. Stage-wise design re-arrange the blocks from one stage to three stages. Thus, for elite performance, some blocks in the new two stages must work less efficiently than those in the original stage. We replace the self-attention blocks with bottleneck blocks in each stage separately for Net5, by which we can estimate the importance of self-attention in different stages. The results are shown in Table III. The replacement of self-attention in all three stages reduces both the base performance and the elite performance. There is a trend that self-attentions in lower resolutions play more important roles than those in higher resolutions. Additionally, replacing the self-attentions in the first stage almost has no effect on the network performance. Larger resolutions contain much more tokens and we conjecture that it is more difficult for self-attentions to learn relations among them.

The second problem is adding $3 \times 3$ convolutions to the feed-forward blocks, which decreases the elite performance by 0.82%. Based on Net4, we replace MLP blocks with bottleneck blocks in each stage separately. As can be seen in Table III-D, although all stages obtain improvements in base performance, only the first stage benefits from bottleneck blocks in elite performance. The $3 \times 3$ convolutions are not necessary for the other two low-resolution stages when self-attentions already have a global view in these positions. On the high-resolution stage, for which self-attentions have difficulty in handling all tokens, the $3 \times 3$ convolutions can provide improvement.

Integrating the observation above, we propose the **Visformer** as vision-friendly, Transformer-based models. The detailed architectures are shown in Table V. Besides the positive transitions, Visformer adopts the stage-wise design for higher base performance. But self-attentions are only utilized in the last two stages, considered that self-attention in the high-resolution stage is relatively inefficient. Visformer employs bottleneck blocks in the first stage and utilizes group $3 \times 3$ convolutions in bottleneck blocks inspired by ResNeXt [54].

| | output size | Visformer-Ti | Visformer-S | VisformerV2-Ti | VisformerV2-S |
|---|---|---|---|---|---|
| stem | $112 \times 112$ | $7 \times 7$, 16, stride 2 | $7 \times 7$, 32, stride 2 | $7 \times 7$, 24, stride 2 | $7 \times 7$, 32, stride 2 |
| emb. | $56 \times 56$ | - | - | $2 \times 2$, 48, stride 2 | $2 \times 2$, 64, stride 2 |
| s0 | $56 \times 56$ | - | - | $\begin{bmatrix} 1 \times 1,\ 96 \\ 3 \times 3,\ 96 \\ (\text{group} = 8) \\ 1 \times 1,\ 48 \end{bmatrix} \times 1$ | $\begin{bmatrix} 1 \times 1,\ 128 \\ 3 \times 3,\ 128 \\ (\text{group} = 8) \\ 1 \times 1,\ 64 \end{bmatrix} \times 1$ |
| emb. | $28 \times 28$ | $4 \times 4$, 96, stride 4 | $4 \times 4$, 192, stride 4 | $2 \times 2$, 96, stride 2 | $2 \times 2$, 128, stride 2 |
| s1 | $28 \times 28$ | $\begin{bmatrix} 1 \times 1,\ 192 \\ 3 \times 3,\ 192 \\ (\text{group} = 8) \\ 1 \times 1,\ 96 \end{bmatrix} \times 7$ | $\begin{bmatrix} 1 \times 1,\ 384 \\ 3 \times 3,\ 384 \\ (\text{group} = 8) \\ 1 \times 1,\ 192 \end{bmatrix} \times 7$ | $\begin{bmatrix} 1 \times 1,\ 192 \\ 3 \times 3,\ 192 \\ (\text{group} = 8) \\ 1 \times 1,\ 96 \end{bmatrix} \times 4$ | $\begin{bmatrix} 1 \times 1,\ 256 \\ 3 \times 3,\ 256 \\ (\text{group} = 8) \\ 1 \times 1,\ 128 \end{bmatrix} \times 10$ |
| emb. | $14 \times 14$ | $2 \times 2$, 192, stride 2 | $2 \times 2$, 384, stride 2 | | |
| s2 | $14 \times 14$ | $\begin{bmatrix} \text{MHSA},\ 192 \\ 1 \times 1,\ 768 \\ 1 \times 1,\ 192 \end{bmatrix} \times 4$ | $\begin{bmatrix} \text{MHSA},\ 384 \\ 1 \times 1,\ 1536 \\ 1 \times 1,\ 384 \end{bmatrix} \times 4$ | $\begin{bmatrix} \text{MHSA},\ 192 \\ 1 \times 1,\ 768 \\ 1 \times 1,\ 192 \end{bmatrix} \times 6$ | $\begin{bmatrix} \text{MHSA},\ 256 \\ 1 \times 1,\ 1024 \\ 1 \times 1,\ 256 \end{bmatrix} \times 14$ |
| emb. | $7 \times 7$ | $2 \times 2$, 384, stride 2 | $2 \times 2$, 768, stride 2 | | |
| s3 | $7 \times 7$ | $\begin{bmatrix} \text{MHSA},\ 384 \\ 1 \times 1,\ 1536 \\ 1 \times 1,\ 384 \end{bmatrix} \times 4$ | $\begin{bmatrix} \text{MHSA},\ 768 \\ 1 \times 1,\ 3072 \\ 1 \times 1,\ 768 \end{bmatrix} \times 4$ | $\begin{bmatrix} \text{MHSA},\ 384 \\ 1 \times 1,\ 1536 \\ 1 \times 1,\ 384 \end{bmatrix} \times 2$ | $\begin{bmatrix} \text{MHSA},\ 512 \\ 1 \times 1,\ 2048 \\ 1 \times 1,\ 512 \end{bmatrix} \times 3$ |
| | $1 \times 1$ | global average pool, 1000-d fc, softmax | | | |
| FLOPs | | $1.3 \times 10^9$ | $4.9 \times 10^9$ | $1.3 \times 10^9$ | $4.3 \times 10^9$ |

We also introduce BatchNorm to patch embedding modules as in CNNs. We name Visformer-S to denote the model that directly comes from DeiT-S. In addition, we can adjust the complexity by changing the output dimensionality of multi-head attentions. Here, we shrink the dimensionality by half and derive the Visformer-Ti model, which requires around $1/4$ computational costs of the Visformer-S model.

### E. VisformerV2: optimizing the architecture configuration

Some architecture configurations of Visformer are not carefully tuned. For example, when splitting the network into different stages, we averagely assign the 12 blocks to the three stages, expect that we utilize 3 more blocks in the first stage to compensate for the removal of self-attention. In other words, the stage configuration ([7, 4, 4]) is not carefully designed. Furthermore, depth and width, which are also not polished in Visformer, have been demonstrated to be very important configurations for network performance. Therefore, we conduct many experiments to explore the architecture of Visformer and propose VisformerV2. VisformerV2 is much better than the original Visformer and the architecture is shown in Table V. The detailed analysis and experiments are shown in Section IV-C.

### F. Transformer with Half-precision

Recently, quantization has been widely used to accelerate the training process and save GPU memory. Specifically, half-precision floating-point (FP16), the lowest precision that can preserve the network performance, has been adopted by many researchers. However, some works [10], [55], [56] have shown that half-precision can lead to overflows in Transformers and we also observe this problem in Visformer.

Based on our experimental analysis, we find that attention score generation can cause overflow. With the queries ($Q$) and keys ($K$), the standard self-attention scores can be computed as:

$$A_{score} = \text{softmax}(\frac{QK^T}{\sqrt{d}}) \qquad (1)$$

However, $Q$ and $K$ can be very large matrices and the elements in $QK^T$ will be the dot-product of two very long vectors. As a result, the scores can overflow easily while utilizing 16-bit precision. To solve this problem, we first try to pre-normalize $Q$ and $K$:

$$A_{score} = \text{softmax}((\frac{Q}{\sqrt[4]{d}})(\frac{K^T}{\sqrt[4]{d}})) \qquad (2)$$

Where $d$ is the length of the vector. Nevertheless, the attention scores still overflow sometimes. This is because the scores are only normalized with $\sqrt{d}$ overall. As the dot product of two vectors with length $d$, the scores are still under the risk of overflow. Consequently, we try to normalize the score with $d$:

$$A_{score} = \text{softmax}((\frac{Q}{\sqrt{d}})(\frac{K^T}{\sqrt{d}})) \qquad (3)$$

TABLE VI
COMPARISON OF INFERENCE TIME FOR VISFORMERV2-S WITH
DIFFERENT SCORE GENERATING METHODS. THE TESTED GPU IS V100
AND THE BATCH SIZE IS 32.

| Method | original | PB-Relax | ours |
|---|---|---|---|
| Batch Time (ms) | 42.8 | 46.6 | 43.0 |

TABLE VII
THE COMPARISON OF BASE AND ELITE PERFORMANCE AS WELL AS THE
FLOPS BETWEEN VISFORMER AND DEIT, THE DIRECT BASELINE.

| Network | base perf. (%) | elite perf. (%) | FLOPs (G) |
|---|---|---|---|
| Visformer-Ti | 74.34 | 78.62 | 1.3 |
| DeiT-Ti | 63.87 | 72.21 | 1.3 |
| Visformer-S | 77.20 | 82.19 | 4.9 |
| DeiT-S | 63.12 | 80.07 | 4.6 |

TABLE VIII
COMPARISON AMONG VISFORMER, DEIT, AND RESNET, IN TERMS OF
CLASSIFICATION ACCURACY (%) USING LIMITED TRAINING DATA. THE
ELITE SETTING WITH 300 EPOCHS IS USED FOR ALL MODELS.

| Network | 100% classes | 10% classes | 1% classes | 10% images | 1% images |
|---|---|---|---|---|---|
| DeiT-S | 80.07 | 80.06 | 73.40 | 40.41 | 6.94 |
| ResNet-50 | 78.73 | 89.90 | 93.20 | 58.37 | 13.59 |
| Visformer-S | 82.19 | 90.06 | 91.60 | 58.74 | 16.56 |
| Deit-Ti | 72.33 | 78.72 | 74.40 | 38.44 | 6.53 |
| ResNet-50-55% | 72.84 | 87.10 | 91.40 | 51.48 | 10.68 |
| Visformer-Ti | 78.62 | 89.48 | 90.60 | 55.14 | 11.79 |

In our experiments, we observed that it can effectively avoid overflow during computing scores and will not degrade the network performance.

Note that CogView [55] also proposes PB-Relax to eliminate overflow in attention scores. PB-Relax pre-minuses the maximum of the attention scores. However, this method needs to tune a hyper-parameter and usually considerably increases the network runtime, as shown in Table VI. As a contrast, our method nearly does not introduce extra runtime.

## IV. MORE EXPERIMENTS ON VISFORMER

### A. The improvements on the upper-bound and lower-bound

We first compare Visformer against DeiT, the direct baseline. Results are summarized in Table VII. Using comparable computational costs, the Visformer models outperform the corresponding DeiT models significantly. Specifically, the advantages of Visformer-S and Visformer-Ti over DeiT-S and DeiT-Ti under the elite setting are 2.12% and 6.41%, while under the base setting, the numbers grow to 14.08% and 10.47%, respectively. In other words, the advantage becomes more significant under the base setting, which is more frequently used for visual recognition.

### B. Training with limited data

We evaluate the performance of Visformer in the scenario with limited training data, which we consider is an important ability of being vision-friendly, while prior Transformer-based models mostly required abundant training data [8].

Four subsets of ImageNet are used, with 10% and 1% randomly chosen classes (all data), and with 10% and 1% randomly chosen images (all classes), respectively. To challenge the models, we still use the elite setting with 300 epochs (not extended). As shown in Table VIII, it is observed that the DeiT-S model reports dramatic accuracy drops in all the four tests (note that the accuracy of using only 10% and 1% classes should be much higher if epochs are extended). In comparison, Visformer remains robust in these scenarios, showing its potential of being used for visual recognition with limited data.

In tiny level, ResNet-50-55% is obtained by reducing the channel numbers (like other tiny models) to 55% (so that the FLOPs, 1.3G, is similar to Visformer-Ti and Deit-Ti). The conclusion is similar: Visformer-Ti is still the best overall model, and the advantage is slightly enlarged because the risk of over-fitting has been reduced.

### C. Designing VisformerV2

We design VisformerV2 by polishing the original Visformer. First, we apply relative position bias [42] to Visformer, which improves the results to 82.39% as shown in Table IV-C. Then we test whether we need to utilize an extra early stage. As shown in Table IX, assigning a block to the new stage does not improve the performance and furthermore, the performance decreases when more blocks are assigned to it. However, we find that this stage can improve the detection and segmentation results, which will be detailed in Section IV-F. Therefore, we decide to utilize one block in the new stage. Next, we test to utilize deep and narrow architecture [43]. We first narrow down the network and directly assign blocks to the self-attention stages (the last two stages). The tested stage configurations are {1, 3, 11, 11} and {1, 6, 10, 10}. These settings degrade the performance. Then we try to assign more blocks to the third stage ( i.e., {1, 3, 18, 3} and {1, 6, 16, 3}), which is the default setting for many convolution [4] and Transformer networks [42], [43]. It improves the networks significantly. We also find that the second pure convolution stage is very important. Moving the blocks from this stage to the other stages will substantially degrade the network. Therefore we assign more blocks to this stage and obtain VisformerV2-S. With a similar study, we design VisformerV2-Ti. The detailed architecture is shown in Table V.

Note that the deep and narrow architecture significantly increases the runtime on GPU. This is because that the wide and shallow architecture has a better parallelization property. To compensate for the loss in runtime, we utilize fewer FLOPs for 'deep-narrow' networks. More importantly, we find that when the input resolution is enlarged (detection and segmentation tasks in Section IV-F) or the model is scaled up, the parallelization property will be improved and the runtime on GPU becomes more consistent with FLOPs.

### D. Comparison to the state-of-the-arts

We then compare Visformer and VisformerV2 to other Transformer-based approaches in Table X. At the tiny level, Visformer-Ti and VisformerV2-Ti outperform other vision Transformers that with similar FLOPs. For larger models, Visformer-S performs much better than most of the models with similar FLOPs. VisformerV2-S further improves the performance and outperforms other vision Transformer models. Note that VisformerV2-S utilize fewer FLOPs and parameters than Visformer-S.

TABLE IX
THE ELITE PERFORMANCE AND INFERENCE TIME OF DIFFERENT VISFORMER MODELS. THE BATCH TIME IS TESTED ON A V100 GPU WITH A BATCH SIZE OF 32.

| Network | block numbers | channel numbers | elite perf.(%) | FLOPs (G) | Params (M) | Batch Time (ms) |
|---|---|---|---|---|---|---|
| Visformer-S | {0, 7, 4, 4} | {96, 192, 384, 768} | 82.39 | 4.9 | 40.2 | 36.9 |
| | {1, 6, 4, 4} | {96, 192, 384, 768} | 82.37 | 4.9 | 40.3 | 37.3 |
| | {3, 4, 4, 4} | | 81.70 | 4.9 | 39.3 | 38.4 |
| | {1, 3, 11, 11} | | 81.73 | 4.2 | 45.4 | 42.1 |
| | {1, 6, 10, 10} | {64, 128, 256, 512} | 82.20 | 4.2 | 41.9 | 41.7 |
| | {1, 3, 18, 3} | | 82.51 | 4.2 | 25.8 | 43.4 |
| | {1, 6, 16, 3} | | 82.89 | 4.2 | 24.6 | 42.8 |
| VisformerV2-S | {1, 10, 14, 3} | {64, 128, 256, 512} | 82.97 | 4.3 | 23.6 | 43.0 |

TABLE X
COMPARISON AMONG OUR METHOD AND OTHER TRANSFORMER-BASED VISION MODELS. '*' INDICATES THAT WE RE-RUN THE MODEL USING THE ELITE SETTING. 'KD' STANDS FOR KNOWLEDGE DISTILLATION [57].

| Methods | Top-1(%) | FLOPs (G) | Params (M) |
|---|---|---|---|
| ResNet-18 [4] | 69.8 | 1.8 | 11.7 |
| DeiT-Ti [10] | 72.2 | 1.3 | 5.7 |
| DeiT-Ti (KD) [10] | 74.6 | 1.3 | 5.7 |
| AutoFormer-Ti [35] | 74.7 | 1.3 | 5.7 |
| PVT-Ti [22] | 75.1 | 1.9 | 13.2 |
| PVTv2-B1 [58] | 78.7 | 2.1 | 13.1 |
| **Visformer-Ti (ours)** | 78.6 | 1.3 | 10.3 |
| **VisformerV2-Ti (ours)** | 79.6 | 1.3 | 9.4 |
| ResNet-50 [4] | 76.2 | 4.1 | 25.6 |
| ResNet-50* [4] | 78.7 | 4.1 | 25.6 |
| RegNetY-4GF [59] | 79.4 | 4.0 | 20.6 |
| RegNetY-8GF [59] | 79.9 | 8.0 | 39.2 |
| RegNetY-4GF* [59] | 80.0 | 4.0 | 20.6 |
| DeiT-S [10] | 79.8 | 4.6 | 21.8 |
| DeiT-S* [10] | 80.1 | 4.6 | 21.8 |
| DeiT-B [10] | 81.8 | 17.4 | 86.3 |
| PVT-S [22] | 79.8 | 3.8 | 24.5 |
| PVT-Medium [22] | 81.2 | 6.7 | 44.2 |
| PVTv2-B2-Li [58] | 82.1 | 3.9 | 22.6 |
| PVTv2-B2 [58] | 82.0 | 4.0 | 25.4 |
| Swin-T [42] | 81.3 | 4.5 | 29 |
| CvT-13 [32] | 81.6 | 4.5 | 20 |
| CvT-13-NAS [32] | 82.2 | 4.1 | 18 |
| CvT-13(384) [32] | 83.0 | 16.3 | 20 |
| Conformer-Ti [31] | 81.3 | 5.2 | 23.5 |
| T2T-ViT$_t$-14 [52] | 80.7 | 5.2 | 21.5 |
| T2T-ViT$_t$-19 [52] | 81.4 | 8.4 | 39.0 |
| BoTNet-S1-59 [29] | 81.7 | 7.3 | 33.5 |
| CSWin-T [43] | 82.7 | 4.3 | 23 |
| AutoFormer-S [35] | 81.7 | 5.1 | 22.9 |
| **Visformer-S (ours)** | 82.2 | 4.9 | 40.2 |
| **VisformerV2-S (ours)** | 83.0 | 4.3 | 23.6 |

TABLE XI
COMPARISON OF INFERENCE EFFICIENCY AMONG VISFORMER AND OTHER MODELS ON A 32G-V100. A BATCH SIZE OF 32 IS USED FOR TESTING. '*' INDICATES THAT THE MODEL IS RE-TRAINED WITH THE ELITE SETTING.

| Methods | Top-1 (%) | FLOPs (G) | Batch Time (ms) |
|---|---|---|---|
| ResNet-50* | 78.7 | 4.1 | 34.2 |
| DeiT-S* | 80.1 | 4.6 | 36.9 |
| RegNetY-4GF* | 80.0 | 4.0 | 40.2 |
| Swin-T | 81.3 | 4.5 | 47.6 |
| CSwin-T | 82.7 | 4.3 | 57.5 |
| PVT-S | 79.8 | 3.8 | 47.6 |
| PVTv2-B2 | 82.0 | 4.0 | 57.1 |
| PVTv2-B2-Li | 82.1 | 3.9 | 56.8 |
| EfficientNet-B3 [60] | 81.6 | 1.8 | 48.3 |
| EfficientNet-B4 [60] | 82.9 | 4.2 | 81.7 |
| Visformer-S (ours) | 82.2 | 4.9 | 36.7 |
| VisformerV2-S (ours) | 83.0 | 4.3 | 43.0 |

## F. COCO Object Detection

Last but not least, we evaluate our models on the COCO object detection task. Since the standard self-attention in Visformer models is not efficient for high-resolution inputs, we simply replace self-attention with the shifted window (Swin) self-attention [42] to apply our models to the detection task. Therefore, Swin Transformers are our important baseline models. It should be emphasized that the self-attention in Visformer can also be replaced with other resolution-friendly self-attentions like CSWin self-attention and MSG self-attention. We just utilize the widely used Swin self-attention to show the superiority of Visformer architecture.

The models are evaluated with two frameworks: Mask-RCNN [61] and Cascade Mask-RCNN [62]. We train the models on COCO 2017 dataset and report the results on COCO val2017. We inherit the training settings in [42]: the AdamW optimizer with a learning rate of 0.0001 and the weight decay of 0.05. The batch size is 16 and we show the results of $1\times$ (12 epochs) and $3\times$ (36 epochs) schedule. The FPS is measured on a V100 GPU with a batch size of 1. The FLOPs are computed with $1280 \times 800$ resolution.

We first test different methods with the Mask R-CNN $1\times$ schedule. As shown at the top of Table XII, Visformer-S slightly outperform Swin-T. When we assign a block to the first stage (Visformer-S-F), although the classification performance is not improved (illustrated in Section IV-C), the detection result becomes better. We conjecture that the

## E. Inference efficiency

Although VisformerV2-S is not as efficient as Visformer-S in runtime, it is still much faster than most vision Transformer models as shown in Table XI. As for the state-of-the-art EfficientNet convnets, Visformer-S are below the EfficientNets with similar FLOPs. However, EfficientNets are computing inefficient on GPUs. It is shown that Visformer-S is significantly faster than EfficientNet-B3 which performance is slightly worse than our model. VisformerV2-S and EfficientNet-B4 have similar FLOPs and performance, but VisformerV2-S is significantly faster than EfficientNet-B4.

| Method | Backbone | $AP^{box}$ | $AP^{box}_{50}$ | $AP^{box}_{75}$ | $AP^{mask}$ | $AP^{mask}_{50}$ | $AP^{mask}_{75}$ | FLOPs | Params | FPS |
|---|---|---|---|---|---|---|---|---|---|---|
| Mask R-CNN 1× schedule | R-50 | 38.0 | 58.6 | 41.4 | 34.4 | 55.1 | 36.7 | 260 | 44 | 18.6 |
| | Swin-T | 42.6 | 65.1 | 46.2 | 39.3 | 62.0 | 42.1 | 267 | 48 | 14.8 |
| | Visformer-S | 43.0 | 65.3 | 47.2 | 39.6 | 62.4 | 42.4 | 275 | 60 | 13.1 |
| | Visformer-S-F | 43.5 | 65.9 | 47.7 | 39.8 | 62.5 | 42.6 | 275 | 60 | 13.0 |
| | VisformerV2-S | 44.8 | 66.8 | 49.4 | 40.7 | 63.9 | 43.7 | 262 | 43 | 15.2 |
| Mask R-CNN 3× + MS schedule | R-50 | 41.0 | 61.7 | 44.9 | 37.1 | 58.4 | 40.1 | 260 | 44 | 18.6 |
| | Swin-T | 46.0 | 68.2 | 50.2 | 41.6 | 65.1 | 44.8 | 267 | 48 | 14.8 |
| | VisformerV2-S | 47.8 | 69.5 | 52.6 | 42.5 | 66.4 | 45.8 | 262 | 43 | 15.2 |
| Cascade Mask R-CNN 1× + MS schedule | R-50 | 43.7 | 61.7 | 47.5 | 38.0 | 58.8 | 41.0 | 739 | 82 | 10.6 |
| | Swin-T | 48.1 | 67.1 | 52.2 | 41.7 | 64.4 | 45.0 | 745 | 86 | 9.5 |
| | VisformerV2-S | 49.3 | 68.1 | 53.6 | 42.3 | 65.1 | 45.7 | 740 | 81 | 9.6 |
| Cascade Mask R-CNN 3× + MS schedule | R-50 | 46.3 | 64.3 | 50.5 | 40.1 | 61.7 | 43.4 | 739 | 82 | 10.6 |
| | DeiT-S | 48.0 | 67.2 | 51.7 | 41.4 | 64.2 | 44.3 | 889 | 80 | - |
| | Swin-T | 50.5 | 69.3 | 54.9 | 43.7 | 66.6 | 47.1 | 745 | 86 | 9.5 |
| | MSG-T [44] | 50.3 | 69.0 | 54.7 | 43.6 | 66.5 | 47.5 | 758 | 86 | 9.5 |
| | PVTv2-B2-Li [58] | 50.9 | 69.5 | 55.2 | - | - | - | 725 | 80 | 8.2 |
| | PVTv2-B2 [58] | 51.1 | 69.8 | 55.3 | - | - | - | 788 | 83 | 7.1 |
| | VisformerV2-S | 51.6 | 70.1 | 56.4 | 44.1 | 67.5 | 47.8 | 740 | 81 | 9.6 |

TABLE XII

OBJECT DETECTION AND INSTANCE SEGMENTATION PERFORMANCE ON COCO 2017. THE FPS IS MEASURED ON A V100 GPU WITH A BATCH SIZE OF 1. THE FLOPS ARE COMPUTED WITH 1280 × 800 RESOLUTION. 'MS' INDICATES MULTI-SCALE TRAINING [38], [64]

.

block in the first stage can help the FPN [63] to explore the low-level features. The VisformerV2-S further improves the performance and outperforms Swin-T by 2.2%. Additionally, because of the improvement on parallelization property, the FPS becomes consistent with FLOPs and VisformerV2-S is faster than Visformer-S.

For the Cascade Mask R-CNN framework, VisformerV2-S still outperforms Swin-T by a large margin. We compare VisformerV2-S with more methods for 3× and multi-scale schedule [38], [64], and our model still performs better than the other methods. As for FPS, our method is as efficient as Swin-T and MSG-T, and is faster than other vision Transformer Methods.

## V. CONCLUSIONS

This paper presents Visformer, a Transformer-based model that is friendly to visual recognition. We propose to use two protocols, the base and elite setting, to evaluate the performance of each model. To study the reason why Transformer-based models and convolution-based models behave differently, we decompose the gap between these models and design an eight-step transition procedure that bridges the gap between DeiT-S and ResNet-50. By absorbing the advantages and discarding the disadvantages, we obtain the Visformer-S model that outperforms both DeiT-S and ResNet-50. Visformer also shows a promising ability when it is transferred to a compact model and when it is evaluated on small datasets.

## VI. ACKNOWLEDGMENTS

## REFERENCES

[1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.

[4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[5] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.

[6] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[7] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," 2018.

[8] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2009, pp. 248–255.

[10] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, "Training data-efficient image transformers & distillation through attention," *arXiv preprint arXiv:2012.12877*, 2020.

[11] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826.

[12] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, "Randaugment: Practical automated data augmentation with a reduced search space," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 702–703.

[13] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6023–6032.

[14] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

[15] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.

[16] Z. Chen, L. Xie, J. Niu, X. Liu, L. Wei, and Q. Tian, "Visformer: The vision-friendly transformer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 589–598.

[17] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[18] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens, "Stand-alone self-attention in vision models," *arXiv preprint arXiv:1906.05909*, 2019.

[19] H. Hu, Z. Zhang, Z. Xie, and S. Lin, "Local relation networks for image recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 3464–3473.

[20] H. Zhao, J. Jia, and V. Koltun, "Exploring self-attention for image recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 076–10 085.

[21] M. Chen, A. Radford, R. Child, J. Wu, H. Jun, D. Luan, and I. Sutskever, "Generative pretraining from pixels," in *International Conference on Machine Learning*. PMLR, 2020, pp. 1691–1703.

[22] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," *arXiv preprint arXiv:2102.12122*, 2021.

[23] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7794–7803.

[24] A. Buades, B. Coll, and J.-M. Morel, "A non-local algorithm for image denoising," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 2. IEEE, 2005, pp. 60–65.

[25] Y. Cao, J. Xu, S. Lin, F. Wei, and H. Hu, "Gcnet: Non-local networks meet squeeze-excitation networks and beyond," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.

[26] Y. Li, X. Jin, J. Mei, X. Lian, L. Yang, C. Xie, Q. Yu, Y. Zhou, S. Bai, and A. L. Yuille, "Neural architecture search for lightweight non-local networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 297–10 306.

[27] I. Bello, B. Zoph, A. Vaswani, J. Shlens, and Q. V. Le, "Attention augmented convolutional networks," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 3286–3295.

[28] I. Bello, "Lambdanetworks: Modeling long-range interactions without attention," *arXiv preprint arXiv:2102.08602*, 2021.

[29] A. Srinivas, T.-Y. Lin, N. Parmar, J. Shlens, P. Abbeel, and A. Vaswani, "Bottleneck transformers for visual recognition," *arXiv preprint arXiv:2101.11605*, 2021.

[30] A. Vaswani, P. Ramachandran, A. Srinivas, N. Parmar, B. Hechtman, and J. Shlens, "Scaling local self-attention for parameter efficient visual backbones," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 894–12 904.

[31] Z. Peng, W. Huang, S. Gu, L. Xie, Y. Wang, J. Jiao, and Q. Ye, "Conformer: Local features coupling global representations for visual recognition," *arXiv preprint arXiv:2105.03889*, 2021.

[32] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, "Cvt: Introducing convolutions to vision transformers," *arXiv preprint arXiv:2103.15808*, 2021.

[33] Z. Dai, H. Liu, Q. V. Le, and M. Tan, "Coatnet: Marrying convolution and attention for all data sizes," *arXiv preprint arXiv:2106.04803*, 2021.

[34] H. Touvron, M. Cord, A. Sablayrolles, G. Synnaeve, and H. Jégou, "Going deeper with image transformers," *arXiv preprint arXiv:2103.17239*, 2021.

[35] M. Chen, H. Peng, J. Fu, and H. Ling, "Autoformer: Searching transformers for visual recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12 270–12 280.

[36] X. Zhai, A. Kolesnikov, N. Houlsby, and L. Beyer, "Scaling vision transformers," *arXiv preprint arXiv:2106.04560*, 2021.

[37] A. Steiner, A. Kolesnikov, X. Zhai, R. Wightman, J. Uszkoreit, and L. Beyer, "How to train your vit? data, augmentation, and regularization in vision transformers," *arXiv preprint arXiv:2106.10270*, 2021.

[38] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision*. Springer, 2020, pp. 213–229.

[39] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, and Y. Zhou, "Transunet: Transformers make strong encoders for medical image segmentation," *arXiv preprint arXiv:2102.04306*, 2021.

[40] H. Chen, Y. Wang, T. Guo, C. Xu, Y. Deng, Z. Liu, S. Ma, C. Xu, C. Xu, and W. Gao, "Pre-trained image processing transformer," *arXiv preprint arXiv:2012.00364*, 2020.

[41] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.

[42] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," *arXiv preprint arXiv:2103.14030*, 2021.

[43] X. Dong, J. Bao, D. Chen, W. Zhang, N. Yu, L. Yuan, D. Chen, and B. Guo, "Cswin transformer: A general vision transformer backbone with cross-shaped windows," *arXiv preprint arXiv:2107.00652*, 2021.

[44] J. Fang, L. Xie, X. Wang, X. Zhang, W. Liu, and Q. Tian, "Msg-transformer: Exchanging local spatial information by manipulating messenger tokens," *arXiv preprint arXiv:2105.15168*, 2021.

[45] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

[46] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," 2017.

[47] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 13 001–13 008.

[48] M. Berman, H. Jégou, A. Vedaldi, I. Kokkinos, and M. Douze, "Multi-grain: a unified image embedding for classes and instances," *arXiv preprint arXiv:1902.05509*, 2019.

[49] E. Hoffer, T. Ben-Nun, I. Hubara, N. Giladi, T. Hoefler, and D. Soudry, "Augment your batch: Improving generalization through instance repetition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8129–8138.

[50] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, "Deep networks with stochastic depth," in *European Conference on Computer Vision*. Springer, 2016, pp. 646–661.

[51] M. Lin, Q. Chen, and S. Yan, "Network in network," *arXiv preprint arXiv:1312.4400*, 2013.

[52] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, F. E. Tay, J. Feng, and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," *arXiv preprint arXiv:2101.11986*, 2021.

[53] Y. Wu and K. He, "Group normalization," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.

[54] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," *arXiv preprint arXiv:1611.05431*, 2016.

[55] M. Ding, Z. Yang, W. Hong, W. Zheng, C. Zhou, D. Yin, J. Lin, X. Zou, Z. Shao, H. Yang *et al.*, "Cogview: Mastering text-to-image generation via transformers," *arXiv preprint arXiv:2105.13290*, 2021.

[56] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong *et al.*, "Swin transformer v2: Scaling up capacity and resolution," *arXiv preprint arXiv:2111.09883*, 2021.

[57] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.

[58] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pvtv2: Improved baselines with pyramid vision transformer," *arXiv preprint arXiv:2106.13797*, 2021.

[59] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. He, and P. Dollár, "Designing network design spaces," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 428–10 436.

[60] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6105–6114.

[61] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *International Conference on Computer Vision*. IEEE, 2017.

[62] Z. Cai and N. Vasconcelos, "Cascade r-cnn: Delving into high quality object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6154–6162.

[63] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.

[64] P. Sun, R. Zhang, Y. Jiang, T. Kong, C. Xu, W. Zhan, M. Tomizuka, L. Li, Z. Yuan, C. Wang *et al.*, "Sparse r-cnn: End-to-end object detection with learnable proposals," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 454–14 463.