

Demystifying Local Vision Transformer: Sparse Connectivity, Weight Sharing, and Dynamic Weight

Qi Han^{1*} Zejia Fan^{2*} Qi Dai³ Lei Sun³

Ming-Ming Cheng¹ Jiaying Liu² Jingdong Wang^{3†}

¹TKLNDST, Nankai University ²Peking University ³Microsoft Research Asia

Abstract

Vision Transformer (ViT) attains state-of-the-art performance in visual recognition, and the variant, Local Vision Transformer, makes further improvements. The major component in Local Vision Transformer, local attention, performs the attention separately over small local windows.

We rephrase local attention as a channel-wise locally-connected layer and analyze it from two network regularization manners, sparse connectivity and weight sharing, as well as weight computation. *Sparse connectivity*: there is no connection across channels, and each position is connected to the positions within a small local window. *Weight sharing*: the connection weights for one position are shared across channels or within each group of channels. *Dynamic weight*: the connection weights are dynamically predicted according to each image instance.

We point out that local attention resembles depth-wise convolution and its dynamic version in sparse connectivity. The main difference lies in weight sharing - depth-wise convolution shares connection weights (kernel weights) across spatial positions. We empirically observe that the models based on depth-wise convolution and the dynamic variant with lower computation complexity perform on-par with or sometimes slightly better than Swin Transformer, an instance of Local Vision Transformer, for ImageNet classification, COCO object detection and ADE semantic segmentation. These observations suggest that Local Vision Transformer takes advantage of two regularization forms and dynamic weight to increase the network capacity. The source code is available in <https://github.com/Atten4Vis/DemystifyLocalViT>.

1 Introduction

Vision Transformer [8, 12, 14, 17, 18, 33, 53, 57, 59, 62, 65] has shown promising performance in ImageNet classification. The improved variants, Local Vision Transformer [7, 37, 54], adopt the local attention mechanism, which partitions the image space into a set of small windows, and conducts the attention over the windows simultaneously. Local attention leads to great improvement in memory and computation efficiency and makes the extension to downstream tasks easier and more efficient, such as object detection and semantic segmentation.

We exploit the conventional network regularization schemes [16], sparse connectivity that controls the model complexity, and weight sharing that relaxes the requirement of increasing the training data scale, as well as dynamic weight prediction that increases the model capability, to study the local attention mechanism. We rephrase local attention as a channel-wise spatially-locally connected layer with dynamic connection weights. The main properties are summarized as follows. (i) Sparse connectivity: there is no connection across channels, and each output position is only connected to

*Equal contribution. Work done as interns at Microsoft Research, Beijing, P.R. China.

†jingdw@microsoft.com.

the input positions within a local window. (ii) Weight sharing: the connection weights are shared across channels or within each group of channels. (iii) Dynamic weight: the connection weights are dynamically predicted according to each image instance.

We compare local attention to depth-wise convolution [6, 23] that is also a channel-wise spatially-locally connected layer. They are similar in sparse connectivity. The major difference lies in the weight sharing pattern: depth-wise convolution shares the weights across spatial positions other than across channels. Other than learning the weights as static model parameters, depth-wise convolution also benefits from dynamic connection weights (convolutional kernel weights) [20].

We take the recently-developed Local Vision Transformer, Swin Transformer [37], as an example, and study the empirical performance of local attention and (dynamic) depth-wise convolution in the training setting same as Swin Transformer. We replace the local attention layer with the (dynamic) depth-wise convolution layer, keeping the overall structure unchanged. The results show that the (dynamic) depth-wise convolution-based approaches achieve comparable or slightly higher performance for ImageNet classification and two downstream tasks, COCO object detection and ADE semantic segmentation, and (dynamic) depth-wise convolution takes lower computation complexity.

We summarize the main findings in the following.

- Local attention adopted by local Vision Transformer takes advantage of existing regularization schemes, sparse connectivity and weight sharing, as well as dynamic weight prediction, for increasing the capability without requiring a corresponding increase in model complexity and training data.
- Local attention and (dynamic) depth-wise convolution are similar in sparse connectivity and differ in weight sharing and dynamic weight prediction forms. The empirical results on visual recognition imply that the regularization forms and the dynamic weight prediction scheme exploited by local attention and (dynamic) depth-wise convolution perform similarly.
- In addition, we present a relation graph to connect convolution and attention, as well as the concurrently-developed MLP-based methods, e.g., ResMLP [52] and MLP-Mixer [51]. The relation graph shows that these methods essentially take advantage of different sparse connectivity and weight sharing patterns for model regularization optionally with dynamic weight prediction.

2 Understanding Local Attention

2.1 Sparse Connectivity, Weight Sharing, and Dynamic Weight

We give a brief introduction of two regularization forms, sparse connectivity and weight sharing, and dynamic weight, and their benefits. We will use the three forms to analyze local attention and connect it to depth-wise convolution.

Sparse connectivity means that there are no connections between some output neurons (variables) and some input neurons in a layer. It reduces the model complexity without decreasing the number of neurons, e.g., the size of the (hidden) representations.

Weight sharing indicates that some connection weights are equal. It lowers the number of model parameters and increases the network size without requiring a corresponding increase in training data [16].

Dynamic weight refers to learning specialized connection weights for each instance. It generally aims to increase the model capacity. If regarding the learned connection weights as hidden variables, dynamic weight can be viewed as introducing second-order operations that increase the capability of the network. The connection to Hopfield networks is discussed in [44].

2.2 Local Attention

Vision Transformer [14] forms a network by repeating the attention layer and the subsequent point-wise MLP (point-wise convolution). The local Vision Transformer, such as Swin Transformer [37] and HaloNet [54], adopts the local attention layer, which partitions the space into a set of small windows and performs the attention operation within each window simultaneously, to improve the memory and computation efficiency.

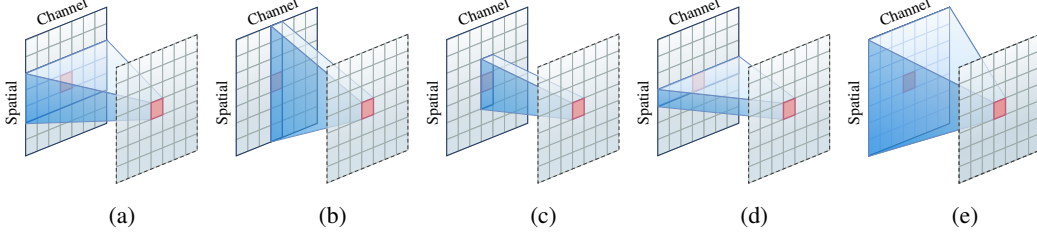


Figure 1: Illustration of connectivity for (a) convolution, (b) global attention and spatial mixing MLP, (c) local attention and depth-wise convolution, (d) point-wise MLP or 1×1 convolution, and (e) MLP (fully-connected layer). In the spatial dimension, we use 1D to illustrate the local-connectivity pattern for clarity.

The local attention mechanism forms the keys and values in a window that the query lies in. The attention output for the query $\mathbf{x}_i \in \mathbb{R}^D$ at the position i is the aggregation of the corresponding values in the local window, $\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iN_k}\}$, weighted by the corresponding attention weights $\{a_{i1}, a_{i2}, \dots, a_{iN_k}\}$ ³:

$$\mathbf{y}_i = \sum_{j=1}^{N_k} a_{ij} \mathbf{x}_{ij}, \quad (1)$$

where $N_k = K_w \times K_h$ is the size of the local window. The attention weight a_{ij} is computed as the softmax normalization of the dot-product between the query \mathbf{x}_i and the key \mathbf{x}_{ij} :

$$a_{ij} = \frac{e^{\frac{1}{\sqrt{D}} \mathbf{x}_i^\top \mathbf{x}_{ij}}}{Z_i} \quad \text{where} \quad Z_i = \sum_{j=1}^{N_k} e^{\frac{1}{\sqrt{D}} \mathbf{x}_i^\top \mathbf{x}_{ij}}. \quad (2)$$

The multi-head version partitions the D -dimensional query, key and value vectors into M subvectors (each with $\frac{D}{M}$ dimensions), and conducts the attention process M times, each over the corresponding subvector. The whole output is the concatenation of M outputs, $\mathbf{y}_i = [\mathbf{y}_{i1}^\top \mathbf{y}_{i2}^\top \dots \mathbf{y}_{iM}^\top]^\top$. The m th output \mathbf{y}_{im} is calculated by

$$\mathbf{y}_{im} = \sum_{j=1}^{N_k} a_{ijm} \mathbf{x}_{ijm}, \quad (3)$$

where \mathbf{x}_{ijm} is the m th value subvector and a_{ijm} is the attention weight computed from the m th head in the same way as Equation 2.

2.3 Properties

We show that local attention is a channel-wise spatially-locally connected layer with dynamic weight computation, and discuss its properties. Figure 1 (c) illustrates the connectivity pattern.

The aggregation processes (Equation 1 and Equation 3) for local attention can be rewritten equivalently in a form of element-wise multiplication:

$$\mathbf{y}_i = \sum_{j=1}^{N_k} \mathbf{w}_{ij} \odot \mathbf{x}_{ij}, \quad (4)$$

where \odot is the element-wise multiplication operator, and $\mathbf{w}_{ij} \in \mathbb{R}^D$ is the weight vector formed from the attention weight a_{ij} or $\{a_{ij1}, a_{ij2}, \dots, a_{ijM}\}$.

Sparse connectivity. The local attention layer is spatially sparse: each position is connected to the N_k positions in a small local window. There are also no connections across channels. The element-wise multiplication in Equation 4 indicates that given the attention weights, each output element, e.g., y_{id} (the i th position for the d th channel), is only dependent on the corresponding input elements from the same channel in the window, $\{x_{i1d}, x_{i2d}, \dots, x_{iN_kd}\}$, and not related to other channels.

Weight sharing. The weights are shared with respect to channels. In the single-head attention case, all the elements $\{w_{ij1}, w_{ij2}, \dots, w_{ijD}\}$ in the weight vector \mathbf{w}_{ij} are the same: $w_{ijd} = a_{ij}$, $1 \leq d \leq D$. In the multi-head attention case, the weight vector \mathbf{w}_{ij} is group-wise same: \mathbf{w}_{ij} is partitioned to M subvectors each corresponding to one attention head, $\{\mathbf{w}_{ij1}, \mathbf{w}_{ij2}, \dots, \mathbf{w}_{ijM}\}$, and the elements in each subvector \mathbf{w}_{ijm} are the same and are equal to the m th attention weight, a_{ijm} .

³For presentation convenience, we ignore the linear projections conducted to the queries, the keys and the values. In vision applications, the value and the corresponding key are from the same feature possibly with different linear projections, and we denote them using the same symbol \mathbf{x}_{ij} .

Dynamic weight. The weights, $\{\mathbf{w}_{i1}, \mathbf{w}_{i2}, \dots, \mathbf{w}_{iN_k}\}$, are dynamically predicted from the query \mathbf{x}_i and the keys $\{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iN_k}\}$ in the local window as shown in Equation 2. We rewrite it as:

$$\{\mathbf{w}_{i1}, \mathbf{w}_{i2}, \dots, \mathbf{w}_{iN_k}\} = f(\mathbf{x}_i; \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iN_k}). \quad (5)$$

Each weight may obtain the information across all the channels, and serves as a bridge to deliver the across-channel information to each output channel.

Set representation. The keys/values for one query are collected as a set with the spatial-order information lost. This leads to that the spatial correspondence between the keys/values across windows is not exploited. The order information loss is partially remedied by encoding the positions as embeddings [14, 53], or learning a so-called relative position embedding (e.g., [37]) in which the spatial-order information is preserved as the keys/values in a local window are collected as a vector.

2.4 Connection to Depth-Wise Convolution

Depth-wise convolution is a type of convolution that applies a single convolutional filter for each channel: $\bar{\mathbf{X}}_d = \mathbf{C}_d \otimes \mathbf{X}_d$, where \mathbf{X}_d and $\bar{\mathbf{X}}_d$ are the d th input and output channel maps, $\mathbf{C}_d \in \mathbb{R}^{N_k}$ is the corresponding kernel weight, and \otimes is the convolution operation. It can be equivalently written in the form of element-wise multiplication for each position:

$$\mathbf{y}_i = \sum_{j=1}^{N_k} \mathbf{w}_{\text{offset}(i,j)} \odot \mathbf{x}_{ij}. \quad (6)$$

Here, $\text{offset}(i, j)$ is the relative offset, $\text{offset}(i, j) = 2D(j) - 2D(i)$, from the 2D coordinate of the position j to the 2D coordinate of the central position i . The weights $\{\mathbf{w}_{\text{offset}(i,j)} \in \mathbb{R}^D; j = 1, 2, \dots, N_k\}$ are reshaped from $\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_D$. The N_k weight vectors are model parameters and shared for all the positions.

We describe the similarities and differences between (dynamic) depth-wise convolution and local attention. Figure 1 (c) illustrates the connectivity patterns.

Similarities. Depth-wise convolution resembles local attention in *sparse connectivity*. There are no connections across channels. Each position is only connected to the positions in a small local window for each channel.

Differences. One main difference lies in *weight sharing*: depth-wise convolution shares the connection weights across spatial positions, while local attention shares the weights across channels or within each group of channels.

The second difference is that the connection weights for depth-wise convolution are *static* and learned as model parameters, while the connection weights for local attention are *dynamic* and predicted from each instance. Depth-wise convolution can also benefit from dynamic weight prediction, e.g., using the weight prediction scheme in SENet [26] to predict the convolutional kernel weights for each instance. More discussions about the dynamic weights are given in Appendix B.

One more difference lies in window representation. Local attention represents the positions in a window by utilizing a *set* form with spatial-order information lost. It explores the spatial-order information implicitly using the positional embedding or explicitly using the learned so-called relative positional embedding. Depth-wise convolution exploits a *vector* form: aggregate the representations within a local window with the weights indexed by the relative position (see Equation 6); keep spatial correspondence between the positions for different windows, thus exploring the spatial-order information explicitly.

2.5 Relation Graph

We present the connectivity patterns in Figure 1, and the relation graph in Figure 2 with the summarization in Table 1 to describe the relation between convolution, depth-wise separable convolution (depth-wise convolution + 1×1 convolution) [23, 6], Vision Transformer [14, 53], Local Vision Transformer [37, 54], as well as multilayer perceptron (MLP), Separable MLP (Sep. MLP, e.g., MLP-Mixer [51], ResMLP [52] and [40]) in terms of sparse connectivity, weight sharing, and dynamic weight. We discuss their relation in the matrix forms in Appendix A.

Multilayer perceptron (MLP) is a fully-connected layer: each neuron (an element at each position and each channel) in one layer is connected with all the neurons in the previous layer⁴. Convolution

⁴We use the widely-used definition for the term MLP: fully-connected layer. There might be other definitions.

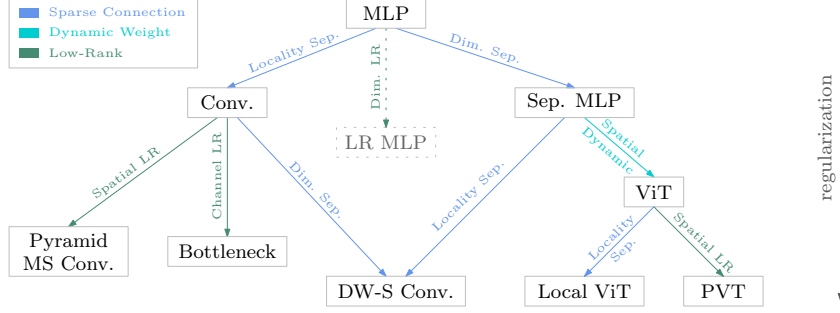


Figure 2: Relation graph for convolution (Conv.), depth-wise separable convolution (DW-S Conv.), Vision Transformer (ViT) building block, local ViT building block, as well as Sep. MLP (e.g., MLP-Mixer and ResMLP) in terms of sparse connectivity and dynamic weight. We also include the low-rank regularization studied for convolutions and ViT and potentially for MLP, and the explanation for pyramid as low rank and other details (not our focus) are given in the appendix. The weight sharing patterns are discussed in Section 2.5. Here, ViT and Local ViT refer to the corresponding building blocks, and PVT means the pyramid way for spatial low-rank. Dim. = dimension including spatial and channel, Sep. = separable, LR = low rank, MS Conv. = multi-scale convolution, PVT = pyramid vision transformer.

and separable MLP are sparse versions of MLP. The connection weights can be formulated as a tensor (e.g., 3D tensor, two dimension for space and one dimension for channel) and the low-rank approximation of the tensor can be used to regularize the MLP (LR MLP, details in Appendix A).

Convolution is a locally-connected layer, formed by connecting each neuron to the neurons in a small local window with the weights shared across the spatial positions. Depth-wise separable convolution is formed by decomposing the convolution into two components: one is point-wise 1×1 convolution, mixing the information across channels, and the other is depth-wise convolution, mixing the spatial information. Other variants of convolution, such as bottleneck, multi-scale convolution or pyramid, can be regarded as low-rank variants (details in Appendix A).

Separable MLP (e.g., MLP-Mixer and ResMLP) reshapes the 3D tensor into a 2D format with the spatial dimension and channel dimension. Separable MLP consists of two sparse MLP along the two dimensions separately, which are formed by separating the input neurons into groups. Regarding channel sparsity, the neurons in the same channel form a group, and an MLP is performed over each group with the MLP parameters shared across groups, forming the first sparse MLP (spatial/token mixing). A similar process is done by viewing the neurons at the same position into a group, forming the second sparse MLP (channel mixing).

Vision Transformer is a dynamic version of separable MLP. The weights in the first sparse MLP (spatial/token mixing) are dynamically predicted from each instance. Local Vision Transformer is a spatially-sparsier version of Vision Transformer: each output neuron is connected to the input neurons in a local window. PVT [57] is a pyramid (spatial sampling/ low-rank) variant of Vision Transformer.

Depth-wise separable convolution can also be regarded as a spatially-sparsier version of separable MLP. In the first sparse MLP (spatial/token mixing), each output neuron is only dependent on the input neurons in a local window, forming depth-wise convolution. In addition, the connection weights are shared across spatial positions, instead of across channels.

3 Experimental Study

We conduct empirical comparisons between local attention and depth-wise convolutions on three visual recognition tasks (studied on Swin Transformer [37]): ImageNet classification, COCO object detection, and ADE semantic segmentation. We follow the structure of Swin Transformer to build the depth-wise convolution-based networks. We apply the same training and evaluation settings from Swin Transformer to our models. In addition, we study the effects of weight sharing and dynamic weight in the two methods.

3.1 Architectures

We use the recently-developed Swin Transformer as the example of local attention-based networks and study the performance over the tiny and base networks: Swin-T and Swin-B, provided by the

Table 1: The comparison of attention, local attention, convolution, depth-wise convolution (DW-Conv.) and the dynamic variant (D-DW-Conv.), as well as MLP and MLP variants in terms of the patterns of sparse connectivity, weight sharing, and dynamic weight. [†]Channel Sep. MLP corresponds to token-mixer MLP. [‡] 1×1 Conv. is also called point-wise MLP. ^bThe weights might be shared within each group of channels.

	Sparse between positions		Sparse between channels	Weight sharing across		Dynamic weight
	non-local	full		position	channel	
Local attention	✓		✓		✓ ^b	✓
Attention			✓		✓	✓
DW-Conv.	✓		✓	✓		
D-DW-Conv.	✓		✓	✓		✓
Conv.	✓			✓		
MLP						
Channel Sep. MLP [†]			✓		✓	
1×1 Conv. [‡]		✓		✓		

authors [37]⁵. We follow the tiny and base networks to build two depth-wise convolution-based networks, DW-Conv.-T and DW-Conv.-B so that the overall architectures are the same, making the comparison fair. We also build the dynamic versions, D-DW-Conv.-T and D-DW-Conv.-B, by predicting the dynamic weights using the similar technique as SENet [26]. We simply replace local attention in Swin Transformer by depth-wise convolution of the same window size, where the pre- and post- linear projections over the values are replaced by 1×1 convolutions. We adopt the convolutional network design pattern to append BN [31] and ReLU [41] to the convolution. The details are available in Appendix D. In terms of parameter and computation complexity, the depth-wise convolution-based networks are lower (Table 2) because there are linear projections for keys and values in local attention.

3.2 Datasets and Implementation Details

ImageNet classification. The ImageNet-1K recognition dataset [13] contains 1.28M training images and 50K validation images with totally 1,000 classes. We use the exactly-same training setting as Swin Transformer [37]. The AdamW [38] optimizer for 300 epochs is adopted, with a cosine decay learning rate scheduler and 20 epochs of linear warm-up. The weight decay is 0.05, and the initial learning rate is 0.001. The augmentation and regularization strategies include RandAugment [11], Mixup [68], CutMix [67], stochastic depth [28], etc.

COCO object detection. The COCO 2017 dataset [35] contains 118K training and 5K validation images. We follow Swin Transformer to adopt Cascade Mask R-CNN [4] for comparing backbones. We use the training and test settings from Swin Transformer: multi-scale training - resizing the input such that the shorter side is between 480 and 800 and the longer side is at most 1333; AdamW optimizer with the initial learning rate 0.0001; weight decay - 0.05; batch size - 16; and epochs - 36.

ADE semantic segmentation. The ADE20K [77] dataset contains 25K images, 20K for training, 2K for validation, and 3K for testing, with 150 semantic categories. The same setting as Swin Transformer [37] is adopted. UPerNet [60] is used as the segmentation framework. Details are provided in Appendix E.

3.3 Main Results

ImageNet classification. The comparison for ImageNet classification is given in Table 2. One can see that the local attention-based networks, Swin Transformer, and the depth-wise convolution-based networks, perform on par (with a slight difference of 0.1) in terms of top-1 accuracy and real accuracy [3] for both tiny and base models. In the tiny model case, the dynamic depth-wise convolution-based network performs higher. In particular, the depth-wise convolution-based networks are more efficient in parameters and computation complexities. In the tiny model case, the parameters and computation complexities are reduced by 14.2% and 15.5%, respectively. Similarly, in the base model case, the two costs are reduced by 15.9% and 16.2%, respectively. The dynamic variant takes more parameters but with almost the same complexity efficiency.

In addition, we report the results for other models: ResNet - with normal convolutions and bottleneck forming residual units; channel and spatial separable MLP - MLP-Mixer [51] and ResMLP [52]; and ViT and DeiT - global attention, viewed as dynamic separable MLP. The reason that the results of

⁵<https://github.com/microsoft/Swin-Transformer> (MIT License)

Table 2: ImageNet classification comparison for ResNet, Mixer and ResMLP, ViT and DeiT, Swin (Swin Transformer), DW-Conv. (depth-wise convolution), and D-DW-Conv. (dynamic depth-wise convolution).

method	img. size	#param.	FLOPs	throughput (img. / s)	top-1 acc.	real acc.
<i>Bottleneck: convolution with low rank</i>						
ResNet-50 [21]	224 ²	26M	4.1G	1128.3	76.2	82.5
ResNet-101 [21]	224 ²	45M	7.9G	652.0	77.4	83.7
ResNet-152 [21]	224 ²	60M	11.6G	456.7	78.3	84.1
<i>Channel and spatial separable MLP, spatial separable MLP = point-wise 1×1 convolution</i>						
Mixer-B/16 [51]	224 ²	46M	-	-	76.4	82.4
Mixer-L/16 [51]	224 ²	189M	-	-	71.8	77.1
ResMLP-12 [52]	224 ²	15M	3.0G	-	76.6	83.3
ResMLP-24 [52]	224 ²	30M	6.0G	-	79.4	85.3
ResMLP-36 [52]	224 ²	45M	8.9G	-	79.7	85.6
<i>Global attention: dynamic channel separable MLP + spatial separable MLP</i>						
ViT-B/16 [14]	384 ²	86M	55.4G	83.4	77.9	83.6
ViT-L/16 [14]	384 ²	307M	190.7G	26.5	76.5	82.2
DeiT-S [53]	224 ²	22M	4.6G	947.3	79.8	85.7
DeiT-B [53]	224 ²	86M	17.5G	298.2	81.8	86.7
DeiT-B [53]	384 ²	86M	55.4G	82.7	83.1	87.7
<i>Local attention: perform attention in local small windows</i>						
Swin-T [37]	224 ²	28M	4.5G	713.5	81.3	86.6
Swin-B [37]	224 ²	88M	15.4G	263.0	83.3	87.9
<i>Depth-wise convolution + point-wise 1×1 convolution</i>						
DW-Conv.-T	224 ²	24M	3.8G	928.7	81.3	86.8
DW-Conv.-B	224 ²	74M	12.9G	327.6	83.2	87.9
D-DW-Conv.-T	224 ²	51M	3.8G	897.0	81.9	87.3
D-DW-Conv.-B	224 ²	162M	13.0G	322.4	83.2	87.9

ResNets are lower than ResMLP might be the strong training setting used in MLP based methods. The overall conclusion seems to be that the locality-based sparsity pattern (adopted in depth-wise convolution and local attention) besides sparsity between channels/spatial positions still facilitates the network training for ImageNet-1K, though separable MLP achieves promising performance.

COCO object detection. The comparisons between local attention (Swin Transformer), depth-wise convolution, and dynamic depth-wise convolution are shown in Table 3. In the tiny model case, depth-wise convolution performs a little lower than local attention, and dynamic depth-wise convolution performs better than the static version and on par with local attention. In the base model case, (dynamic) depth-wise convolution performs a little worse than local attention.

ADE semantic Segmentation. The comparisons of single scale testing on ADE semantic segmentation are shown in Table 3. In the tiny model case, (dynamic) depth-wise convolution is $\sim 1.0\%$ higher than local attention. In the base model case, the performances are similar⁶.

Summary. In ImageNet classification, depth-wise convolution and its dynamic variant are superior over local attention: almost the same accuracy with higher computation efficiency. Dynamic depth-wise convolution is more advantageous in the tiny model case.

In COCO object detection, dynamic depth-wise convolution performs the same with local attention for the tiny model, and local attention is superior for the base model. The reasons might be: (i) the training setting for local attention [37] might not be suitable for depth-wise convolution, or (ii) it is helpful for detection that each position in local attention has its own dynamic weights encoding the information of the corresponding object. We will conduct a further study by predicting the weights for each position in dynamic depth-wise convolution as done [58].

In ADE semantic segmentation, depth-wise convolution and its dynamic variant are superior over local attention for the tiny model, and the performance is similar for the base model.

⁶We conducted an additional experiment by changing the ending learning rate from 0 to $1e - 6$. The base model with depth-wise convolutions achieves a higher mIoU score: 48.9.

Table 3: Comparison results on COCO object detection and ADE semantic segmentation.

	COCO Object Detection						ADE20K Semantic Segmentation		
	#param.	FLOPs	AP ^{box}	AP ^{box} ₅₀	AP ^{box} ₇₅	AP ^{mask}	#param.	FLOPs	mIoU
Swin-T	86M	747G	50.5	69.3	54.9	43.7	60M	947G	44.5
DW Conv.-T	82M	730G	49.9	68.6	54.3	43.4	56M	928G	45.5
D-DW Conv.-T	108M	730G	50.5	69.5	54.6	43.7	83M	928G	45.7
Swin-B	145M	986G	51.9	70.9	56.5	45.0	121M	1192G	48.1
DW Conv.-B	132M	924G	51.1	69.6	55.4	44.2	108M	1129G	48.3
D-DW Conv.-B	219M	924G	51.2	70.0	55.4	44.4	195M	1129G	48.0

3.4 Additional Studies

Weight sharing. We study how the performance is affected by the number of channels in each group across which the weights are shared (the numbers of attention heads at each stage are accordingly changed). We use the tiny Swin Transformer model for this study and the subsequent studies. The results from Swin Transformer shown in Figure 3 imply that in the case of too many channels and too few channels in each group, the accuracy is not the best.

In addition, we study how sharing weights across channels for depth-wise convolution affects the performance. We use the same weight sharing pattern across channels in Swin Transformer for sharing weights across channels in depth-wise convolution. The ImageNet top-1 accuracy is slightly reduced: from 81.3 to 81.1, implying that proper weight sharing across channels does not have big impact for depth-wise convolution.

Dynamic weight. We study how dynamic weight in local attention affects the performance. We study the static variant: learn the weights in each window as model parameters (the weights are not shared across windows). The static version achieves the ImageNet top-1 accuracy 80.3%, lower than the dynamic version 81.3% for the tiny model, implying that dynamic weight is helpful. We point out that the static variant is a locally-connected version of separable MLP (ResMLP): the MLP over each channel (spatial/token mixing) is done over each window, other than the whole image space. The results are shown in Table 4 (DW = depth-wise conv.). As a comparison, we also show the results of dynamic depth-wise convolution.

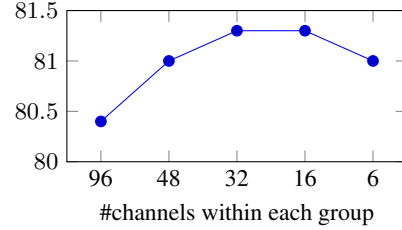


Figure 3: **Effect of #channels sharing the weights** on ImageNet classification. Too many and too few channels sharing the weights do not lead to the best ImageNet classification accuracy.

Set representation. Local attention represents the positions in a window as a set with the spatial-order information lost. Swin Transformer learns relative positional embeddings where the positions in a window are actually described as a vector keeping the spatial-order information. It is reported in [37] that removing the relative positional embeddings leads to a 1.2% accuracy drop, indicating the spatial-order information is important.

Retraining on 384×384 images. Similar to [37], we study the performance of fine-tuning the models: first learn with 224×224 images, then fine-tune on large images of 384×384 . We study two cases: (1) keep the window size 7×7 unchanged; and (2) upsample the kernel weights from 7×7 to 12×12 as done in [37] for upsampling the relative positional embeddings.

The results are in Table 5⁷. In the case of keeping the window size 7×7 unchanged, depth-wise convolution (DW) performs better. When using a larger window size 12×12 , depth-wise convolution performs worse than 7×7 . We suspect that this is because upsampling the kernel weights is not a good starting for fine-tuning. In Swin Transformer, using a larger window size improves the performance. We believe that this is because the local attention mechanism is suitable for variable window sizes.

Cooperating with SE. Squeeze-and-excitation [26] (SE) is a parameter- and computation-efficient dynamic module, initially designed for improving the ResNet performance. The results in Table 6 show that depth-wise convolution (DW), a static module, benefits from the SE module, while Swin

⁷Swin Transformer takes slightly higher FLOPs for 7×7 than 12×12 . The higher computation cost comes from larger padding than 12×12 .

Table 4: Dynamic weight.					Table 5: Retrain on larger images.					Table 6: Cooperate with SE.				
	dynamic	#params	FLOPs	Acc.	model	ws.	#param.	FLOPs	Acc.	model	SE	#param.	FLOPs	Acc.
Swin	✗	26M	3.8G	80.3	Swin	7×7	28M	14.4G	81.8	Swin	✓	28M	4.5G	81.3
		28M	4.5G	81.3		12×12	28M	14.2G	82.4			29M	4.5G	81.2
DW	✓	24M	3.8G	81.3	DW	7×7	24M	11.1G	82.2	DW	✓	24M	3.8G	81.3
		51M	3.8G	81.9		12×12	25M	11.5G	82.1			24M	3.8G	81.7

Transformer, already a dynamic module, does not benefit from dynamic module SE. The reason is still unclear, and might lie in the optimization.

4 Related Work

Sparse connectivity. Sparse connection across channels is widely explored for removing redundancy in the channel domain. The typical schemes are depth-wise convolution adopted by MobileNet [23, 45], ShuffleNetV2 [39] and IGCv3 [46], and group convolution adopted by ResNeXt [61], merge-and-run [75], ShuffleNetV1 [72], and IGC [71].

The self-attention unit⁸ in Vision Transformer, its variants [5, 8, 14, 19, 22, 34, 37, 42, 53, 54, 57, 59, 64, 65, 69, 74, 78], and the spatial information fusion unit (e.g., token-mixer in MLP-Mixer [51] and ResMLP [52]) have no connections across channels.

1×1 (point-wise) convolution (in ShuffleNetV2 [39], MobileNet [23, 45], IGC [71], ViT [14], local ViT [37, 54], MLP-Mixer [51], ResMLP [52]) has no connections across spatial positions. The convolutions with other kernel sizes and local attention [74, 37, 54] have connections between each position and the positions within a small local window, respectively.

In addition to hand-crafted sparse connections, various methods are developed for learning sparse connections, e.g., CondenseNet [27] and dynamic grouping [73].

Weight sharing. Weight sharing across spatial positions is mainly used in convolution, including normal convolution, depth-wise convolution and point-wise convolution. Weight sharing across channels is adopted in the attention unit [55], its variants [7, 8, 14, 34, 37, 53, 54, 57, 59, 65], and token-mixer MLP in MLP-mixer [51] and ResMLP [52].

Dynamic weight. Predicting the connection weights is widely studied in convolutional networks. There are basically two types. One is to learn homogeneous connection weights, e.g., SENet [26], dynamic convolution [32]. The other is to learn the weights for each region or each position (GENet [25], Lite-HRNet [63], Involution [34]). The attention unit in ViT or local ViT learns dynamic connection weights for each position.

Networks built with depth-wise separable convolutions. There are many networks built upon depth-wise separable convolution or its variants, such as MobileNet [23, 45], ShuffleNet [39], IGC [71], Xception [6], and EfficientNet [48, 49]. In this paper, instead of proposing new convolutional modules or improving depth-wise separable convolution, our goal is to compare depth-wise convolution with local attention.

Convolution vs Transformer. The study in [10] shows that a multi-head self-attention layer can simulate a convolutional layer by taking into consideration the linear projection conducted on values, and with specific conditions, e.g., well-designed relative positional embeddings and losing the dynamic weight scheme. Differently, our analysis and comparison do not need the linear projection conducted on values, and the connections are discussed for local attention with depth-wise convolution other than normal convolution. In [1], the mathematical connection (in terms of the tensor form) between convolution and attention is presented. The opinion that convolution and attention are essentially about the model complexity control is similar to ours, and we make the detailed analysis and report empirical studies.

The concurrently-developed work in NLP [50] empirically compares lightweight depth-wise convolution [58] to Transformer for NLP tasks, and reaches a conclusion similar to ours for vision tasks: convolution and Transformer obtain on-par results. Differently, we attempt to understand why they perform on par from three perspectives: sparse connectivity, weight sharing and dynamic weight, and discuss their similarities and differences.

⁸The pre- and post- linear projections for values can be regarded as 1×1 convolutions. The attention weights generated from keys and values with linear projections in some sense mix the information across channels.

5 Conclusion

We aim to understand local attention through the connection to depth-wise convolution. The experiments imply that the performance of local attention is on par with (dynamic) depth-wise convolution, suggesting that the good performance of local attention essentially stems from two regularization forms, sparse connectivity and weight sharing, and dynamic weight. In addition, we also discuss how the concurrently-developed works, e.g., ResMLP and MLP-Mixer, are related to ViT and depth-wise convolution. As future works, we will study if the training settings and the architecture design for depth-wise convolution can be improved over the current settings adopted from Swin Transformer.

References

- [1] Jean-Marc Andreoli. Convolution, attention and structure embedding. *arXiv preprint arXiv:1905.01289*, 2019.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord. Are we done with imagenet? *arXiv preprint arXiv:2006.07159*, 2020.
- [4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: high quality object detection and instance segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019.
- [5] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. *arXiv preprint arXiv:2012.00364*, 2020.
- [6] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1251–1258, 2017.
- [7] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting spatial attention design in vision transformers. *arXiv preprint arXiv:2104.13840*, 2021.
- [8] Xiangxiang Chu, Bo Zhang, Zhi Tian, Xiaolin Wei, and Huaxia Xia. Do we really need explicit position encodings for vision transformers? *arXiv preprint arXiv:2102.10882*, 2021.
- [9] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. <https://github.com/open-mmlab/mms Segmentation>, 2020.
- [10] Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between self-attention and convolutional layers. In *Int. Conf. Learn. Represent.*, 2020.
- [11] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 702–703, 2020.
- [12] Stéphane d’Ascoli, Hugo Touvron, Matthew Leavitt, Ari Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. *arXiv preprint arXiv:2103.10697*, 2021.
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 248–255. Ieee, 2009.
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Int. Conf. Learn. Represent.*, 2021.
- [15] Shang-Hua Gao, Qi Han, Duo Li, Pai Peng, Ming-Ming Cheng, and Pai Peng. Representative batch normalization with feature calibration. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.

- [16] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- [17] Meng-Hao Guo, Zheng-Ning Liu, Tai-Jiang Mu, and Shi-Min Hu. Beyond self-attention: External attention using two linear layers for visual tasks. *arXiv preprint arXiv:2105.02358*, 2021.
- [18] Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on visual transformer. *arXiv preprint arXiv:2012.12556*, 2020.
- [19] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and Yunhe Wang. Transformer in transformer. *arXiv preprint arXiv:2103.00112*, 2021.
- [20] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic neural networks: A survey. *arXiv preprint arXiv:2102.04906*, 2021.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 770–778, 2016.
- [22] Byeongho Heo, Sangdoo Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon Oh. Rethinking spatial dimensions of vision transformers. *arXiv preprint arXiv:2103.16302*, 2021.
- [23] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [24] Han Hu, Zheng Zhang, Zhenda Xie, and Stephen Lin. Local relation networks for image recognition. In *Int. Conf. Comput. Vis.*, pages 3464–3473, 2019.
- [25] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Andrea Vedaldi. Gather-excite: Exploiting feature context in convolutional neural networks. In *Adv. Neural Inform. Process. Syst.*, 2018.
- [26] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7132–7141, 2018.
- [27] Gao Huang, Shichen Liu, Laurens Van der Maaten, and Kilian Q Weinberger. Condensenet: An efficient densenet using learned group convolutions. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2752–2761, 2018.
- [28] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *Eur. Conf. Comput. Vis.*, pages 646–661. Springer, 2016.
- [29] Lang Huang, Yuhui Yuan, Jianyuan Guo, Chao Zhang, Xilin Chen, and Jingdong Wang. Interlaced sparse self-attention for semantic segmentation. *CoRR*, abs/1907.12273, 2019.
- [30] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Int. Conf. Comput. Vis.*, pages 603–612, 2019.
- [31] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Int. Conf. Mach. Learn.*, pages 448–456. PMLR, 2015.
- [32] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc Van Gool. Dynamic filter networks. In *Adv. Neural Inform. Process. Syst.*, 2016.
- [33] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *arXiv preprint arXiv:2101.01169*, 2021.
- [34] Duo Li, Jie Hu, Changhu Wang, Xiangtai Li, Qi She, Lei Zhu, Tong Zhang, and Qifeng Chen. Involution: Inverting the inherence of convolution for visual recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.

- [35] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Eur. Conf. Comput. Vis.*, pages 740–755. Springer, 2014.
- [36] Hanxiao Liu, Zihang Dai, David R So, and Quoc V Le. Pay attention to mlps. *arXiv preprint arXiv:2105.08050*, 2021.
- [37] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*, 2021.
- [38] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Int. Conf. Learn. Represent.* OpenReview.net, 2019.
- [39] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Eur. Conf. Comput. Vis.*, pages 116–131, 2018.
- [40] Luke Melas-Kyriazi. Do you even need attention? a stack of feed-forward layers does surprisingly well on imagenet. *arXiv preprint arXiv:2105.02723*, 2021.
- [41] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Int. Conf. Mach. Learn.*, 2010.
- [42] Zizheng Pan, Bohan Zhuang, Jing Liu, Haoyu He, and Jianfei Cai. Scalable visual transformers with hierarchical pooling. *arXiv preprint arXiv:2103.10619*, 2021.
- [43] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. In *Adv. Neural Inform. Process. Syst.*, pages 68–80, 2019.
- [44] Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Thomas Adler, Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Sandve, et al. Hopfield networks is all you need. *arXiv preprint arXiv:2008.02217*, 2020.
- [45] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4510–4520, 2018.
- [46] Ke Sun, Mingjie Li, Dong Liu, and Jingdong Wang. Igcv3: Interleaved low-rank group convolutions for efficient deep neural networks. In *Brit. Mach. Vis. Conf.*, 2018.
- [47] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5693–5703, 2019.
- [48] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *Int. Conf. Mach. Learn.*, pages 6105–6114. PMLR, 2019.
- [49] Mingxing Tan and Quoc V Le. Efficientnetv2: Smaller models and faster training. *arXiv preprint arXiv:2104.00298*, 2021.
- [50] Yi Tay, Mostafa Dehghani, Jai Gupta, Dara Bahri, Vamsi Aribandi, Zhen Qin, and Donald Metzler. Are pre-trained convolutions better than pre-trained transformers? *arXiv preprint arXiv:2105.03322*, 2021.
- [51] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, et al. Mlp-mixer: An all-mlp architecture for vision. *arXiv preprint arXiv:2105.01601*, 2021.
- [52] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby, Edouard Grave, Armand Joulin, Gabriel Synnaeve, Jakob Verbeek, and Hervé Jégou. Resmlp: Feedforward networks for image classification with data-efficient training. *arXiv preprint arXiv:2105.03404*, 2021.

- [53] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020.
- [54] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.
- [55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Adv. Neural Inform. Process. Syst.*, pages 5998–6008, 2017.
- [56] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.
- [57] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. *arXiv preprint arXiv:2102.12122*, 2021.
- [58] Felix Wu, Angela Fan, Alexei Baevski, Yann N. Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. In *Int. Conf. Learn. Represent.*, 2019.
- [59] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808*, 2021.
- [60] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Eur. Conf. Comput. Vis.*, pages 418–434, 2018.
- [61] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1492–1500, 2017.
- [62] Weijian Xu, Yifan Xu, Tyler Chang, and Zhuowen Tu. Co-scale conv-attentional image transformers. *arXiv preprint arXiv:2104.06399*, 2021.
- [63] Changqian Yu, Bin Xiao, Changxin Gao, Lu Yuan, Lei Zhang, Nong Sang, and Jingdong Wang. Lite-hrnet: A lightweight high-resolution network. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.
- [64] Kun Yuan, Shaopeng Guo, Ziwei Liu, Aojun Zhou, Fengwei Yu, and Wei Wu. Incorporating convolution designs into visual transformers. *arXiv preprint arXiv:2103.11816*, 2021.
- [65] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986*, 2021.
- [66] Li Yuan, Qibin Hou, Zihang Jiang, Jiashi Feng, and Shuicheng Yan. Volo: Vision outlooker for visual recognition, 2021.
- [67] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Int. Conf. Comput. Vis.*, pages 6023–6032, 2019.
- [68] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *Int. Conf. Learn. Represent.*, 2018.
- [69] Pengchuan Zhang, Xiyang Dai, Jianwei Yang, Bin Xiao, Lu Yuan, Lei Zhang, and Jianfeng Gao. Multi-scale vision longformer: A new vision transformer for high-resolution image encoding. *arXiv preprint arXiv:2103.15358*, 2021.
- [70] Qinglong Zhang and Yubin Yang. Rest: An efficient transformer for visual recognition. *arXiv preprint arXiv:2105.13677*, 2021.

- [71] Ting Zhang, Guo-Jun Qi, Bin Xiao, and Jingdong Wang. Interleaved group convolutions. In *Int. Conf. Comput. Vis.*, pages 4373–4382, 2017.
- [72] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6848–6856, 2018.
- [73] Zhaoyang Zhang, Jingyu Li, Wenqi Shao, Zhanglin Peng, Ruimao Zhang, Xiaogang Wang, and Ping Luo. Differentiable learning-to-group channels via groupable convolutional neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3542–3551, 2019.
- [74] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, June 2020.
- [75] Liming Zhao, Mingjie Li, Depu Meng, Xi Li, Zhaoxiang Zhang, Yueting Zhuang, Zhuowen Tu, and Jingdong Wang. Deep convolutional neural networks with merge-and-run mappings. In Jérôme Lang, editor, *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pages 3170–3176. ijcai.org, 2018.
- [76] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Assoc. Adv. Artif. Intell.*, volume 34, pages 13001–13008, 2020.
- [77] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 633–641, 2017.
- [78] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*, 2021.

APPENDIX

A Matrix Form for Explaining Relation Graph

We use the matrix form to explain sparsity connectivity in various layers and how they are obtained by modifying the MLP. We reshaw the relation graph in Figure 4.

MLP. The term MLP, Multilayer Perceptron, is used ambiguously, sometimes loosely to any feedforward neural network. We adopt one of the common definitions, and use it to refer to fully-connected layers. Our discussion is based on a single fully-connected layer, and can be easily generalized to two or more fully-connected layers. One major component, except the nonlinear units and others, is a linear transformation:

$$\mathbf{y} = \mathbf{W}\mathbf{x}, \quad (7)$$

where \mathbf{x} represents the input neurons, \mathbf{y} represents the output neurons, and \mathbf{W} represents the connection weights, e.g., $\mathbf{W} \in \mathbb{R}^{N_C \times N_C}$, where N is the number of positions, and C is the number of channels.

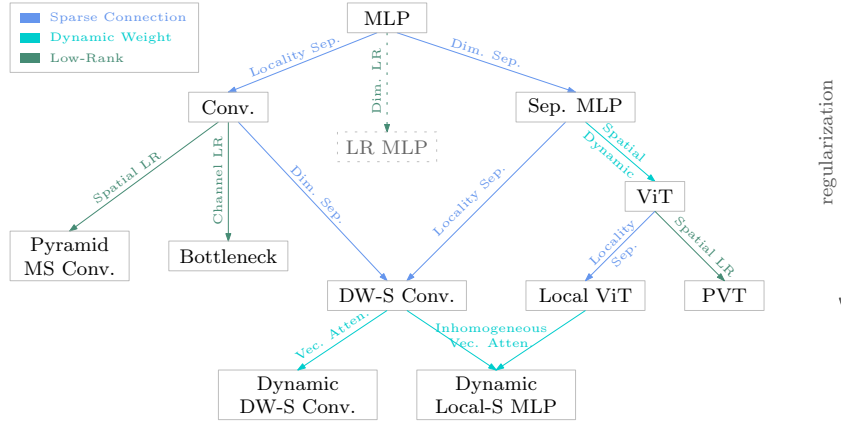


Figure 4: Relation graph for convolution (Conv.), depth-wise separable convolution (DW-S Conv.), Vision Transformer (ViT) building block, local ViT building block, Sep. MLP (e.g., MLP-Mixer and ResMLP), dynamic depth-wise separable convolution (Dynamic DW-S Conv.), as well as dynamic local separable MLP (e.g., involution [34] and inhomogeneous dynamic depth-wise convolution) in terms of sparse connectivity and dynamic weight. Dim. = dimension including spatial and channel, Sep. = separable, LR = low rank, MS Conv. = multi-scale convolution, PVT = pyramid vision transformer.

Convolution. Considering the 1D case with a single channel (the 2D case is similar), the connection weight matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$ is in the following sparse form, also known as the Toeplitz matrix (We use the window size 3 as an example):

$$\mathbf{W} = \begin{bmatrix} a_2 & a_3 & 0 & 0 & \cdots & 0 & a_1 \\ a_1 & a_2 & a_3 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ a_3 & 0 & 0 & 0 & \cdots & a_1 & a_2 \end{bmatrix}. \quad (8)$$

For the C -channel case, we organize the input into a vector channel by channel: $[\mathbf{x}_1^\top \mathbf{x}_2^\top \cdots \mathbf{x}_C^\top]^\top$, and accordingly the connection weight matrix channel by channel for the c_o th output channel, $\mathbf{W}_{c_o} = [\mathbf{W}_{c_o1} \mathbf{W}_{c_o2} \cdots \mathbf{W}_{c_oC}]$ (the form of $\mathbf{W}_{c_o i}$ is the same as Equation 8). The whole form could be written as

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_C \end{bmatrix} = \begin{bmatrix} \mathbf{W}_1 \\ \mathbf{W}_2 \\ \vdots \\ \mathbf{W}_C \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_C \end{bmatrix}. \quad (9)$$

Sep. MLP. Sep. MLP, e.g., ResMLP and MLP-Mixer, is formed with two kinds of block-sparse matrices: one for channel-mixing and the other for spatial-mixing. In the case that the input is organized channel by channel (the neurons in each channel form a group), $\mathbf{x} = [\mathbf{x}_1^\top \mathbf{x}_2^\top \dots \mathbf{x}_C^\top]^\top$, the connection weight matrix is in a block-sparse form:

$$\mathbf{W} = \begin{bmatrix} \mathbf{W}_c & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{W}_c & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{W}_c \end{bmatrix}, \quad (10)$$

where the block matrices $\mathbf{W}_c \in \mathbb{R}^{N \times N}$ are shared across all the channels, and the sharing pattern can be modified to share weights within each group of channels.

The input can be reshaped position by position (the neurons at each position forms a group): $\mathbf{x} = [\mathbf{x}_1^\top \mathbf{x}_2^\top \dots \mathbf{x}_N^\top]^\top$, and similarly one more connection weight matrix can be formulated in a block-sparse form (it is essentially a 1×1 convolution, $\mathbf{W}_p \in \mathbb{R}^{C \times C}$):

$$\mathbf{W}' = \begin{bmatrix} \mathbf{W}_p & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{W}_p & \dots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{W}_p \end{bmatrix}. \quad (11)$$

The forms of block-sparsity are studied in interleaved group convolutions [71] without sharing the weights across groups.

Sep. MLP can also be regarded as using Kronecker product to approximate the connection matrix,

$$\mathbf{W}\mathbf{x} = \text{vec}(\mathbf{A} \text{mat}(\mathbf{x})\mathbf{B}). \quad (12)$$

Here, $\mathbf{W} = \mathbf{B}^\top \otimes \mathbf{A} = \mathbf{W}_c^\top \otimes \mathbf{W}_p$, and \otimes is the Kronecker product operator. $\text{mat}(\mathbf{x})$ reshapes the vector \mathbf{x} in a 2D matrix form, while $\text{vec}(\mathbf{x})$ reshapes the 2D matrix into a vector form. In Sep. MLP, the 2D matrix, $\text{mat}(\mathbf{x}) \in \mathbb{R}^{C \times N}$, is organized so that each row corresponds to one channel and each column corresponds to one spatial position. CCNet [30] and interlaced self-attention [29] use Kronecker product to approximate the spatial connection: the former reshapes the vector in a 2D matrix form along the x and y axes, and the latter reshapes the vector windows by windows.

Vision Transformer (ViT). The matrix form is similar to Sep. MLP. The difference is that the matrix \mathbf{W}_c is predicted from each image instance. The weight prediction manner in ViT has a benefit: handle an arbitrary number of input neurons.

Depth-wise separable convolution. There are two basic components: depth-wise convolution, and 1×1 convolution that is the same as channel-mixing MLP in Sep. MLP. Depth-wise convolution can be written in the matrix form:

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_C \end{bmatrix} = \begin{bmatrix} \mathbf{W}_{11} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{W}_{22} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{W}_{CC} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_C \end{bmatrix}, \quad (13)$$

where the form of \mathbf{W}_{cc} is the same as Equation 8.

Local ViT. In the non-overlapping window partition case, local ViT simply repeats ViT over each window separately with the linear projections, applied to keys, values, and queries, shared across windows. In the overlapping case, the form is a little complicated, but the intuition is the same. In the extreme case, the partition is the same as convolution, and the form is as the following:

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_C \end{bmatrix} = \begin{bmatrix} \mathbf{W}^d & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{W}^d & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{W}^d \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_C \end{bmatrix}, \quad (14)$$

where the dynamic weight matrix \mathbf{W}^d is like the form below:

$$\mathbf{W}^d = \begin{bmatrix} a_{12} & a_{13} & 0 & 0 & \cdots & 0 & a_{11} \\ a_{21} & a_{22} & a_{23} & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{N3} & 0 & 0 & 0 & \cdots & a_{N1} & a_{N2} \end{bmatrix}. \quad (15)$$

Low-rank MLP. Low-rank MLP approximates the connection weight matrix $\mathbf{W} \in \mathbb{R}^{D_o \times D_i}$ in Equation 7 using the product of two low-rank matrix:

$$\mathbf{W} \leftarrow \mathbf{W}_{D_o r} \mathbf{W}_{r D_i}, \quad (16)$$

where r is a number smaller than D_i and D_o

Pyramid. The downsampling process in the pyramid networks can be regarded as spatial low rank: $\mathbf{W}(\in \mathbb{R}^{N^C \times N^C}) \rightarrow \mathbf{W}'(\in \mathbb{R}^{N'^C \times N'^C})$, where N' is equal to $\frac{N}{4}$ in the case that the resolution is reduced by $\frac{1}{2}$. If the numbers of input and output channels are different, it becomes $\mathbf{W}(\in \mathbb{R}^{N^C' \times N^C}) \rightarrow \mathbf{W}'(\in \mathbb{R}^{N'^C' \times N'^C})$.

Multi-scale parallel convolution. Multi-scale parallel convolution used in HRNet [56, 47] can also be regarded as spatial low rank. Consider the case with four scales, multi-scale parallel convolution can be formed as the following,

$$\mathbf{W} \rightarrow \begin{bmatrix} \mathbf{W}_1 \in \mathbb{R}^{N C_1} \\ \mathbf{W}_2 \in \mathbb{R}^{N C_2} \\ \mathbf{W}_3 \in \mathbb{R}^{N C_3} \\ \mathbf{W}_4 \in \mathbb{R}^{N C_4} \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{W}'_1 \in \mathbb{R}^{N C_1} \\ \mathbf{W}'_2 \in \mathbb{R}^{\frac{N}{4} C_2} \\ \mathbf{W}'_3 \in \mathbb{R}^{\frac{N}{16} C_3} \\ \mathbf{W}'_4 \in \mathbb{R}^{\frac{N}{64} C_4} \end{bmatrix}, \quad (17)$$

where C_1, C_2, C_3 , and C_4 are the numbers of the channels in four resolutions.

B Local Attention vs Convolution: Dynamic Weights

We take the 1D case with the window size $2K + 1$ as an example to illustrate the dynamic weight prediction manner. Let $\{\mathbf{x}_{i-K}, \dots, \mathbf{x}_i, \dots, \mathbf{x}_{i+K}\}$ correspond to the $(2K + 1)$ positions in the i th window, and $\{w_{i-K}, \dots, w_i, \dots, w_{i+K}\}$ be the corresponding dynamic weights for updating the representation of the i th (center) position. The discussion can be easily extended to multiple weights for each positions, like the M -head attention and updating the representations for other positions.

Inhomogeneous dynamic convolution. We use the case using only a single linear projection to illustrate inhomogeneous dynamic convolution. The properties we will discuss are similar for more linear projections. The dynamic weights are predicted as the following:

$$\begin{bmatrix} w_{i-K} \\ \vdots \\ w_i \\ \vdots \\ w_{i+K} \end{bmatrix} = \Theta \mathbf{x}_i = \begin{bmatrix} \boldsymbol{\theta}_{-K}^\top \\ \vdots \\ \boldsymbol{\theta}_0^\top \\ \vdots \\ \boldsymbol{\theta}_K^\top \end{bmatrix} \mathbf{x}_i. \quad (18)$$

It can be seen that dynamic convolution learns the weights for each position through the parameters that are different for different positions, e.g., $\boldsymbol{\theta}_k$ corresponds to w_{i+k} . It regards the positions in the window as the vector form, keeping the spatial order information.

Dot-product attention. The dot-product attention mechanism in the single-head case predicts the weights as the following⁹:

$$\begin{bmatrix} w_{i-K} \\ \vdots \\ w_i \\ \vdots \\ w_{i+K} \end{bmatrix} = \begin{bmatrix} (\mathbf{x}_{i-K})^\top \\ \vdots \\ (\mathbf{x}_i)^\top \\ \vdots \\ (\mathbf{x}_{i+K})^\top \end{bmatrix} \mathbf{P}_k^\top \mathbf{P}_q \mathbf{x}_i. \quad (19)$$

Dot-product attention uses the same parameters $\mathbf{P}_k^\top \mathbf{P}_q$ for all the positions. The weight depends on the features at the same position, e.g., w_{i-k} corresponds to \mathbf{x}_{i-k} . It in some sense regards the positions in the window as a set form, losing the spatial order information.

We rewrite it as the following

$$\Theta_d = \begin{bmatrix} (\mathbf{x}_{i-K})^\top \\ \vdots \\ (\mathbf{x}_i)^\top \\ \vdots \\ (\mathbf{x}_{i+K})^\top \end{bmatrix} \mathbf{P}_k^\top \mathbf{P}_q, \quad (20)$$

from which we can see that the parameters Θ_d is dynamically predicted. In other words, dot-product attention can be regarded as a two-level dynamic scheme.

Relative position embeddings is equivalent to adding static weights that keeps the spatial order information:

$$\begin{bmatrix} w_{i-K} \\ \vdots \\ w_i \\ \vdots \\ w_{i+K} \end{bmatrix} = \Theta_d \mathbf{x}_i + \begin{bmatrix} \beta_{-K} \\ \vdots \\ \beta_0 \\ \vdots \\ \beta_K \end{bmatrix}. \quad (21)$$

A straightforward variant is a combination of the static Θ and the dynamic Θ_d :

$$\begin{bmatrix} w_{i-K} \\ \vdots \\ w_i \\ \vdots \\ w_{i+K} \end{bmatrix} = (\Theta_d + \Theta) \mathbf{x}_i. \quad (22)$$

Convolutional attention. We introduce a convolutional attention framework so that it enjoys the benefits of dynamic convolution and dot-product attention: keep the spatial order information and two-level dynamic weight prediction.

⁹For presentation clarity, we omit the softmax normalization and the scale in dot-product. What we discuss still holds if softmax and scale are included.

The post-convolutional attention mechanism left-multiplies a matrix (with the kernel size being 3):

$$\Theta_d = \begin{bmatrix} a_2 & a_3 & 0 & 0 & \cdots & 0 & a_1 \\ a_1 & a_2 & a_3 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ a_3 & 0 & 0 & 0 & \cdots & a_1 & a_2 \end{bmatrix} \begin{bmatrix} (\mathbf{x}_{i-K})^\top \\ \vdots \\ (\mathbf{x}_i)^\top \\ \vdots \\ (\mathbf{x}_{i+K})^\top \end{bmatrix} \mathbf{P}_k^\top \mathbf{P}_q. \quad (23)$$

This can be reviewed as a variant of relative positional embeddings (Equation 21). In the simplified case that the left matrix is diagonal, it can be regarded as the product version of relative positional embeddings (Equation 21 is an addition version).

We can perform a convolution with the kernel size being 3, the kernel weights shared across channels (it is also fine not to share weights), and then do dot-product attention. This is called pre-convolutional attention: perform convolutions on the representations. The two processes are can be written as follows (omit BN and ReLU that follow the convolution),

$$\begin{bmatrix} w_{i-K} \\ \vdots \\ w_i \\ \vdots \\ w_{i+K} \end{bmatrix} = \begin{bmatrix} a_1 & a_2 & a_3 & \cdots & 0 & 0 & 0 \\ 0 & a_1 & a_1 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_2 & a_3 & 0 \\ 0 & 0 & 0 & \cdots & a_1 & a_2 & a_3 \end{bmatrix} \begin{bmatrix} (\mathbf{x}_{i-K-1})^\top \\ (\mathbf{x}_{i-K})^\top \\ \vdots \\ (\mathbf{x}_i)^\top \\ \vdots \\ (\mathbf{x}_{i+K})^\top \\ (\mathbf{x}_{i+K+1})^\top \end{bmatrix} \mathbf{P}_k^\top \mathbf{P}_q [\mathbf{x}_{i-1} \quad \mathbf{x}_i \quad \mathbf{x}_{i+1}] \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}. \quad (24)$$

It can be generalized to using normal convolution:

$$\begin{bmatrix} w_{i-K} \\ \vdots \\ w_i \\ \vdots \\ w_{i+K} \end{bmatrix} = \mathbf{C}' \begin{bmatrix} \mathbf{x}_{i-K-1} & \mathbf{x}_{i-K-1} & \cdots & \mathbf{x}_{i-K-1} \\ \mathbf{x}_{i-K} & \mathbf{x}_{i-K} & \cdots & \mathbf{x}_{i-K} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_i & \mathbf{x}_i & \cdots & \mathbf{x}_i \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_{i+K} & \mathbf{x}_{i+K} & \cdots & \mathbf{x}_{i+K} \\ \mathbf{x}_{i+K+1} & \mathbf{x}_{i+K+1} & \cdots & \mathbf{x}_{i+K+1} \end{bmatrix} \mathbf{P}_k^\top \mathbf{P}_q \mathbf{C}_3 \begin{bmatrix} \mathbf{x}_{i-1} \\ \mathbf{x}_i \\ \mathbf{x}_{i+1} \end{bmatrix}. \quad (25)$$

Here, \mathbf{C}' is a $(2K+1)$ -row matrix and can be easily derived from the convolutional kernel \mathbf{C}_3 . The $(2K+1)$ weights, $\{w_{i-1}, w_i, w_{i+1}\}$, correspond to the $(2K+1)$ rows in \mathbf{C} , respectively. This means that the three positions are differentiated and the same position in each window corresponds to the same row. This explains why the positional embeddings are not necessary when convolutions are adopted [59]. Using different pairs $(\mathbf{W}_q, \mathbf{W}_k)$ leads to more weights for each position, e.g., M pairs correspond to M -head attention.

C Local Attention vs Convolution: Equivalence to Translation

In local attention, the equivalence to translation depends if the keys/values are changed, i.e., if the query lies in the same window, when the feature map is translated. In the case of sparsely-sampled window, e.g., [24, 37, 43, 54], (for efficient implementation), local attention is equivalent to block-wise translation, i.e., the translation is a block with the size same as the window size $K_w \times K_h$ or multiple blocks. In the case that the windows are densely sampled (e.g., [74]), local attention is equivalent to translation.

Table 7: Architectures details of Swin Transformer and depth-wise convolution-based network (DW Conv.) for the tiny model. The architectures for the base model can be easily obtained.

	downsp. rate (output size)	Swin	DW Conv.
stage 1	$4\times$ (56×56)	concat 4×4 , linear 96-d, LN	concat 4×4 , linear 96-d, LN
		LN, linear 96x3-d local sa. 7×7 , head 3 linear 96-d LN, linear 384-d GELU, linear 96-d $\times 2$	linear 96-d, BN, ReLU depthwise conv. 7×7 , BN, ReLU linear 96-d, BN BN, linear 384-d GELU, linear 96-d $\times 2$
stage 2	$8\times$ (28×28)	concat 2×2 , linear 192-d, LN	concat 2×2 , linear 192-d, LN
		LN, linear 192x3-d local sa. 7×7 , head 6 linear 192-d LN, linear 768-d GELU, linear 192-d $\times 2$	linear 192-d, BN, ReLU depthwise conv. 7×7 , BN, ReLU linear 192-d, BN BN, linear 768-d GELU, linear 192-d $\times 2$
stage 3	$16\times$ (14×14)	concat 2×2 , linear 384-d, LN	concat 2×2 , linear 384-d, LN
		LN, linear 384x3-d local sa. 7×7 , head 12 linear 384-d LN, linear 1536-d GELU, linear 384-d $\times 6$	linear 384-d, BN, ReLU depthwise conv. 7×7 , BN, ReLU linear 384-d, BN BN, linear 1536-d GELU, linear 384-d $\times 6$
stage 4	$32\times$ (7×7)	concat 2×2 , linear 768-d, LN	concat 2×2 , linear 768-d, LN
		LN, linear 768x3-d local sa. 7×7 , head 24 linear 768-d LN, linear 3072-d GELU, linear 768-d $\times 2$	linear 768-d, BN, ReLU depthwise conv. 7×7 , BN, ReLU linear 768-d, BN BN, linear 3072-d GELU, linear 768-d $\times 2$
stage 4	1×1	LN, AvgPool. 1×1 linear classifier	LN, AvgPool. 1×1 linear classifier

Depth-wise convolution is similar to local attention in equivalence to translation. Depth-wise convolution is equivalence to any translation and not limited in block translation in local attention. This is because of weight sharing across spatial positions¹⁰ [16].

D Architecture Details

Overall structures. Following local vision transformer, Swin Transformer [37], we build two depth-wise convolution-based networks, namely DW-Conv.-T and DW-Conv.-B. The corresponding dynamic versions are D-DW-Conv.-T and D-DW-Conv.-B. The depth-wise convolution-based networks follow the overall structure of Swin Transformer. We replace local self attention by depth-wise convolution with the same window size. We use batch normalization [31] and ReLU [41] instead of layer normalization [2] in the convolution blocks.

Table 7 shows the architecture details of Swin Transformer and depth-wise convolution-based networks for the tiny model. Normalizations are performed within the residual block, same as Swin Transformer. The base model is similarly built by following Swin Transformer to change the number of channels and the depth of the third stage.

¹⁰The boundary positions are often taken into no consideration when talking about equivalence to translation.

Table 8: ImageNet classification comparison for ResNet, HRNet, Mixer and ResMLP and gMLP, ViT and DeiT, Swin (Swin Transformer), DW-Conv. (depth-wise convolution), and D-DW-Conv. (dynamic depth-wise convolution). † means that ResNet is built by using two 3×3 convolutions to form the residual units.

method	img. size	#param.	FLOPs	throughput (img. / s)	top-1 acc.	real acc.
<i>Convolution: local connection</i>						
ResNet-38 † [56]	224^2	28M	3.8G	2123.7	75.4	-
ResNet-72 † [56]	224^2	48M	7.5G	623.0	76.7	-
ResNet-106 † [56]	224^2	65M	11.1G	452.8	77.3	-
<i>Bottleneck: convolution with low rank</i>						
ResNet-50 [21]	224^2	26M	4.1G	1128.3	76.2	82.5
ResNet-101 [21]	224^2	45M	7.9G	652.0	77.4	83.7
ResNet-152 [21]	224^2	60M	11.6G	456.7	78.3	84.1
<i>Pyramid: convolution with pyramid (spatial low rank) features.</i>						
HRNet-W18 [56]	224^2	21M	4.0G	-	76.8	-
HRNet-W32 [56]	224^2	41M	8.3G	-	78.5	-
HRNet-W48 [56]	224^2	78M	16.1G	-	79.3	-
<i>Channel and spatial separable MLP, spatial separable MLP = point-wise 1×1 convolution</i>						
Mixer-B/16 [51]	224^2	46M	-	-	76.4	82.4
Mixer-L/16 [51]	224^2	189M	-	-	71.8	77.1
ResMLP-12 [52]	224^2	15M	3.0G	-	76.6	83.3
ResMLP-24 [52]	224^2	30M	6.0G	-	79.4	85.3
ResMLP-36 [52]	224^2	45M	8.9G	-	79.7	85.6
gMLP-Ti [36]	224^2	6M	1.4G	-	72.0	-
gMLP-S [36]	224^2	20M	4.5G	-	79.4	-
gMLP-B [36]	224^2	73M	15.8G	-	81.6	-
<i>Global attention: dynamic channel separable MLP + spatial separable MLP</i>						
ViT-B/16 [14]	384^2	86M	55.4G	83.4	77.9	83.6
ViT-L/16 [14]	384^2	307M	190.7G	26.5	76.5	82.2
DeiT-S [53]	224^2	22M	4.6G	947.3	79.8	85.7
DeiT-B [53]	224^2	86M	17.5G	298.2	81.8	86.7
DeiT-B [53]	384^2	86M	55.4G	82.7	83.1	87.7
<i>Pyramid attention: perform attention with spatial low rank</i>						
PVT-S [57]	224^2	25M	3.8G	-	79.8	-
PVT-M [57]	224^2	44M	6.7G	-	81.2	-
PVT-L [57]	224^2	61M	9.8G	-	81.7	-
<i>Local attention: perform attention in local small windows</i>						
Swin-T [37]	224^2	28M	4.5G	713.5	81.3	86.6
Swin-B [37]	224^2	88M	15.4G	263.0	83.3	87.9
<i>Depth-wise convolution + point-wise 1×1 convolution</i>						
DW-Conv.-T	224^2	24M	3.8G	928.7	81.3	86.8
DW-Conv.-B	224^2	74M	12.9G	327.6	83.2	87.9
D-DW-Conv.-T	224^2	51M	3.8G	897.0	81.9	87.3
D-DW-Conv.-B	224^2	162M	13.0G	322.4	83.2	87.9

Dynamic depth-wise convolution. Dynamic depth-wise convolution generates the connection weights according to the instance. We conduct the global average pooling operation to get a vector, and adopt two linear projections: the first one reduces the dimension by $1/4$, followed by BN and ReLU, and then generate the kernel weights. Unlike SENet [26], we currently do not use the Sigmoid activation function for generating the weights.

E Setting Details

ImageNet pretraining. We use the identical training setting with Swin Transformer in ImageNet pretraining for fair comparison. The default input size is 224×224 . The AdamW optimizer [38], with the initial learning rate 0.001 and the weight decay 0.05, is used for 300 epochs. The learning rate is scheduled by a cosine decay schema and warm-up with linear schema for the first 20 epochs. We train

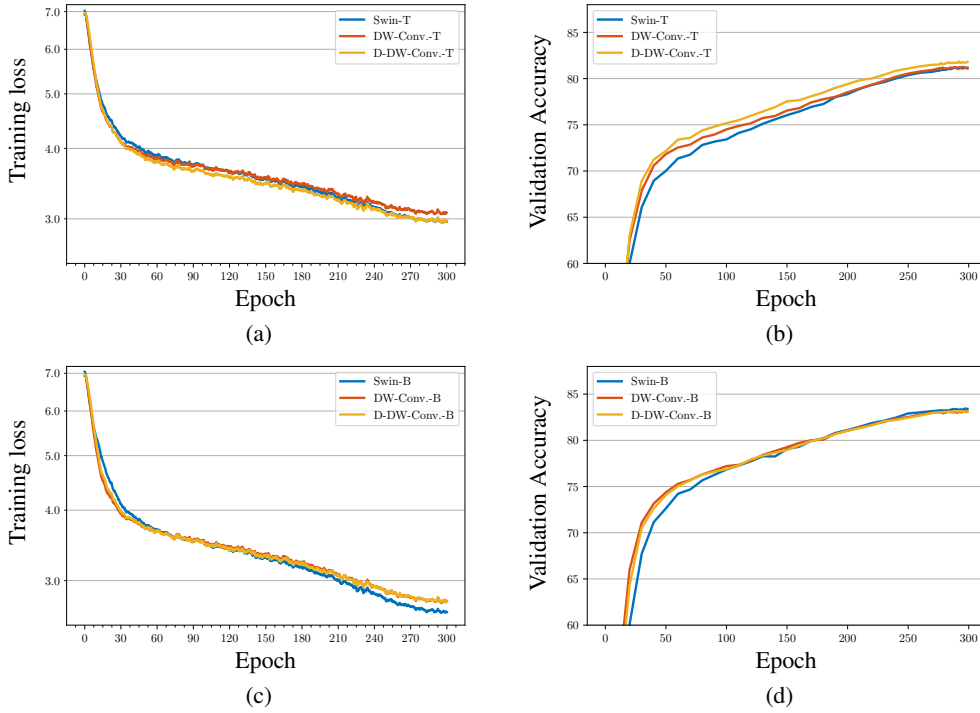


Figure 5: Training and validation curves for ImageNet classification. (a) and (b) are the training loss and validation top-1 accuracy curves for the tiny model, and (c) and (d) are for the base model.

the model on 8 GPUs with the total batch size 1024. The augmentation and regularization strategies are same as Swin Transformer, which includes RandAugment [11], Mixup [68], CutMix [67], random erasing [76] and stochastic depth [28]. The stochastic depth rate is employed as 0.2 and 0.5 for the tiny and base models, respectively, the same as Swin Transformer.

COCO object detection. We follow Swin Transformer to adopt Cascade Mask R-CNN [4] for comparing backbones. We use the training and test settings from Swin Transformer: multi-scale training - resizing the input such that the shorter side is between 480 and 800 and the longer side is at most 1333; AdamW optimizer with the initial learning rate 0.0001; weight decay - 0.05; batch size - 16; and epochs - 36.

ADE semantic segmentation. Following Swin Transformer, we use UPerNet [60] as the segmentation framework. We use the same setting as the Swin Transformer: the AdamW optimizer with initial learning rate 0.00006; weight decay 0.01; linear learning rate decay; 160,000 iterations with warm-up for 1500 iterations; 8 GPUs with mini-batch 2 per GPU. We use the same data augmentation as Swin Transformer based on MMSegmentation [9]. The experimental results are reported as single scale testing.

Static version of Swin Transformer. We remove the linear projections applied to keys and queries, accordingly dot production and softmax normalization. The connection weights (corresponding to attention weights in the dynamic version) are set as static model parameters which are learnt during the training and shared for all the images.

Retraining on 384×384 . We retrain the depth-wise convolution-based network on the ImageNet dataset with 384×384 input images from the model trained with 224×224 images. We use learning rate 10^{-5} , weight decay 10^{-8} and stochastic depth ratio 0.1 for 30 epochs for both 7×7 and 12×12 windows.

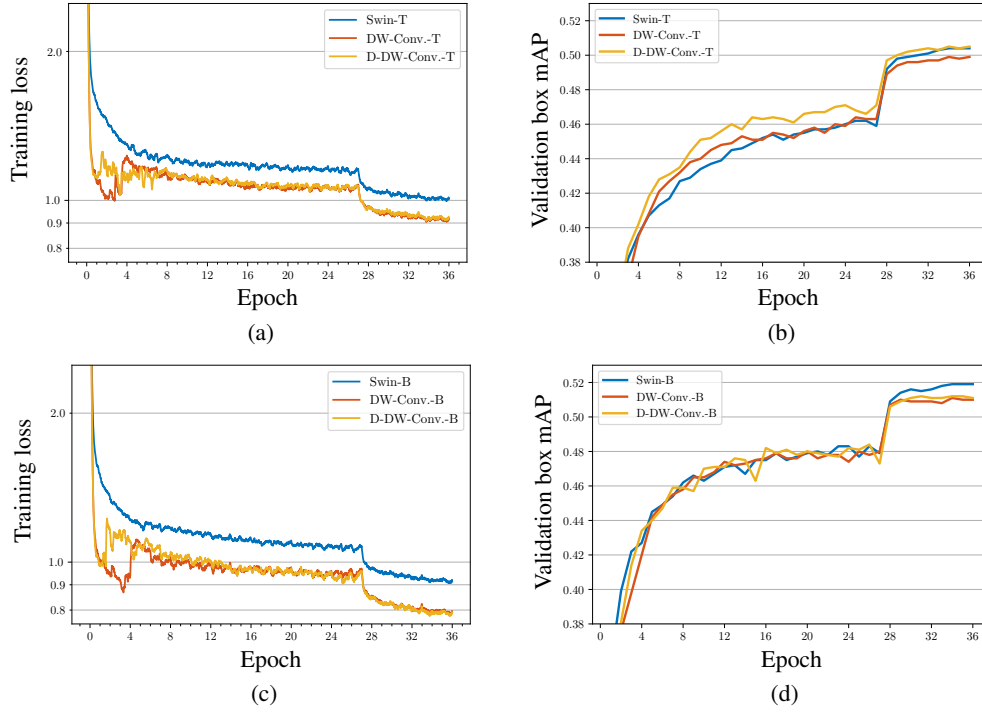


Figure 6: Training and validation curves for COCO object detection. (a) and (b) are the training loss and validation box AP curves for the tiny model, and (c) and (d) are for the base model. It is not expected that depth-wise convolution-based models have lower training errors, but lower detection scores.

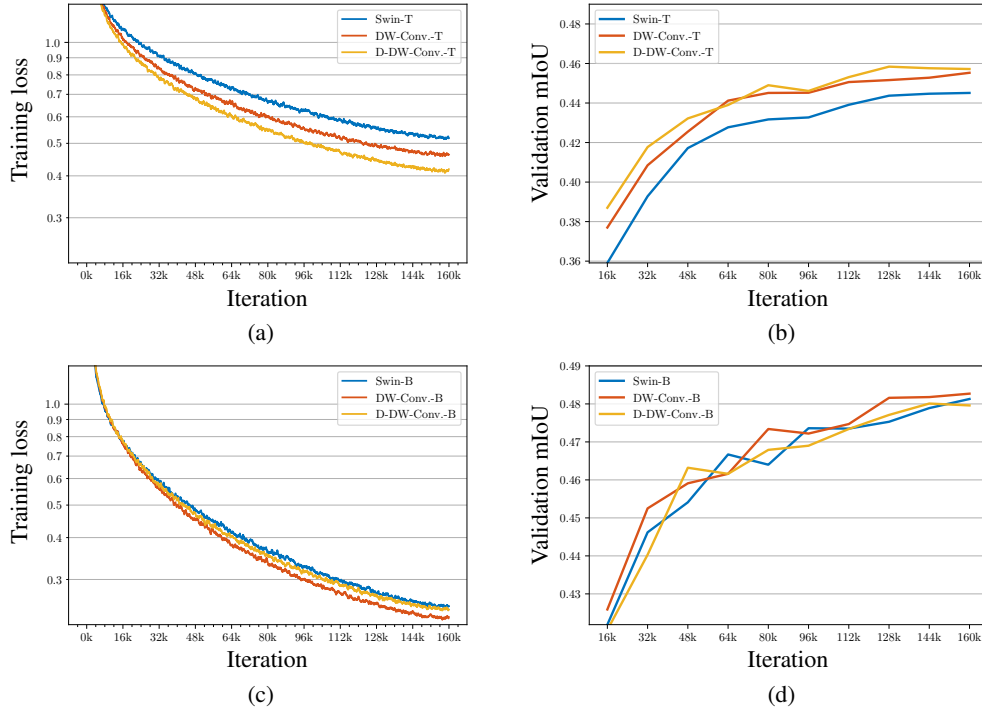


Figure 7: Training and validation curves for ADE semantic segmentation. (a) and (b) are the training loss and validation curves mIoU for the tiny model, and (c) and (d) are for the base model.

Table 9: Exploring normalization schemes of Swin Transformer and depth-wise convolution based networks (DW Conv.) for the tiny model. The results are reported on the ImageNet top-1 accuracy.

	Layer Norm.	Batch Norm.	Centering calibrated Batch Norm.	Top-1 Acc.
Swin	✓			81.3
Swin		✓		80.9
Swin			✓	81.2
DW Conv.	✓			81.2
DW Conv.		✓		81.3
DW Conv.			✓	81.7

F Additional Experiments and Analysis

More results on ImageNet classification. We give more experimental results with different sparse connection strategies, as shown in Table 8. These results also verify that locality-based sparsity pattern (adopted in depth-wise convolution and local attention) besides sparsity between channels/spatial positions still facilitates the network training for ImageNet-1K.

Training curves on ImageNet, COCO and ADE. Figures 5, 6 and 7 show the training and validation curves for Swin Transformer and depth-wise convolution-based methods on ImageNet classification, COCO object detection and ADE20K semantic segmentation.

The curves for ImageNet classification and ADE20K semantic segmentation are normal, but the curves for object detection shown in Figure 6 are not normal: depth-wise convolutions get lower training errors, but lower validation scores. The reason is not clear, and might be the training setting (same as Swin Transformer on COCO object detection) or other issues.

Cooperating with different normalization functions. Transformers usually use the layer normalization to stabilize the training, while convolutional architectures adopt batch normalization. We verify different combinations of backbones (Swin and DW Conv.) and normalization functions. The popular used layer normalization (LN), batch normalization (BN), and the dynamic version of batch normalization - centering calibrated batch normalization [15] (CC. BN) are verified in the experiments. Table 9 shows the results on ImageNet classification.

Combining weight sharing across positions and channels. Depth-wise convolution shares weights across positions, while local transformer shares weights across channels or within each group of channels. In static Swin Transformer, we study a further variant, the weight parameters are shared across windows. In depth-wise convolution-based networks, we additionally share the weights across channels in the same way as Swin Transformer. The results are reported in Table 10.

Spatial inhomogeneous dynamic convolutional weights. In our experiment, we use weights shared across positions for the dynamic version of depth-wise convolution-based networks. This may be enhanced by using weights not shared across positions, such as GENet [25], Involution [34], and Lite-HRNet [63].

We made an initial investigation (inhomogeneous dynamic): generate local weights for each position using two 1×1 convolutions to predict the weights shared across each group of channels, which is a generalization of homogeneous dynamic weight prediction and similar to [34, 58, 63], and share the weights within each group of channels. We use 32 groups by default. The results are shown in Table 11. The higher performance from our new dynamic weight prediction way may stem from that the weights using the attention mechanism are predicted by regarding the keys as a set and our approach generates the kernel weights as a feature vector.

Cooperating with additional depth-wise convolutions. Compared with local self attention, depth-wise convolution slightly reduces the computational cost and model parameters. Thus we simply add additional depth-wise convolutions between the original local fusion block and FFN block, with residual connection. Optionally, we add an SE layer following the additional depth-wise convolution to introduce dynamic characteristic among channel. Results are shown in Tab. 12.

Depth-wise convolution with other architectures. We conduct experiments on other local attention designs, such as SVT [7] and VOLO [66] whose implementations are publicly available. SVT uses local self attention as a basic spatial feature fusion operation, while VOLO proposes a new attention

Table 10: Combination of weight sharing across channels and positions. The results are reported on the ImageNet top-1 accuracy.

	sharing across channels	sharing across positions	Acc.
Swin	✓	✗	80.3
	✓	✓	80.3
DW Conv.	✗	✓	81.3
	✓	✓	81.1

Table 11: Generate local weights for each position and share the weights with each group of channels (inhomogeneous dynamic, I-Dynamic).

	ImageNet				COCO				ADE20K		
	#param.	FLOPs	top-1 acc.	real acc.	#param.	FLOPs	AP ^{box}	AP ^{mask}	#param.	FLOPs	mIoU
Swin-T	28M	4.5G	81.3	86.6	86M	747G	50.5	43.7	60M	947G	44.5
DW Conv.-T	24M	3.8G	81.3	86.8	82M	730G	49.9	43.4	56M	928G	45.5
D-DW Conv.-T	51M	3.8G	81.9	87.3	108M	730G	50.5	43.7	83M	928G	45.7
I-Dynamic-T	26M	4.4G	81.8	87.1	84M	741G	50.8	44.0	58M	939G	46.2
Swin-B	88M	15.4G	83.3	87.9	145M	986G	51.9	45.0	121M	1192G	48.1
DW Conv.-B	74M	12.9G	83.2	87.9	132M	924G	51.1	44.2	108M	1129G	48.3
D-DW Conv.-B	162M	13.0G	83.2	87.9	219M	924G	51.2	44.4	195M	1129G	48.0
I-Dynamic-B	80M	14.3G	83.4	88.0	137M	948G	51.8	44.8	114M	1153G	47.8

module named Vision Outlooker. We replace the local self attention with depth-wise convolution in SVT same as the paper, and replace Vision Outlooker with 7×7 local self attention and 7×7 depth-wise convolution, respectively. The remaining structures are unchanged and the same training setting is used as the original papers. The experimental results are shown in Tab 13 and the observations are the same as the Swin Transformer design.

G Potential Studies

Complexity balance between point-wise (1×1) convolution and depth-wise (spatial) convolution. Depth-wise convolution takes only about 2% computation in the depth-wise convolution-based architecture. The major computation complexity comes from 1×1 convolutions. The solutions to this issue could be: group 1×1 convolution studied in IGC [71, 46], and channel-wise weighting (like SENet) studied in Lite-HRNet [63] and EfficientNet [48, 49], or simply add more depth-wise (spatial) convolutions.

Attention weights as channel maps. Attention weights in attention can be regarded as channel maps. The operations, such as convolution or simple weighting, can be applied to the attention weights. The resT approach [70] performs 1×1 convolutions over the attention weight maps.

Dynamic weights. In Swin Transformer and our developed dynamic depth-wise convolution networks, only the spatial part, attention and depth-wise convolution, explores dynamic weights. Lite-HRNet instead studies dynamic weight for point-wise (1×1) convolution. It is interesting to explore dynamic weight for both parts.

Convolution-style MLP weights. The weights of the spatial-mixing MLP in MLP-Mixer and ResMLP could be modified in the convolution-like style with more weights (some like the relative

Table 12: Co-operations with additional depth-wise convolution before FFN (DW) and optionally with SE layer (SE). Results are reported on ImageNet classification with tiny model.

	#param.	FLOPs	top-1 acc.
DW Conv.	24M	3.8G	81.3
DW Conv. + DW	24M	3.8G	81.9
DW Conv. + DW + SE	25M	3.8G	82.1

Table 13: Comparison between local attention and depth-wise convolution in VOLO [66] and SVT [7] architecture. Results are reported on ImageNet classification with tiny model.

	#param.	FLOPs	top-1 acc.
VOLO-d1 [66]	27M	7.0G	84.1
VOLO (Local SA)-d1	27M	7.2G	84.2
DW Conv.-d1	26M	6.9G	84.2
SVT-S [7]	24M	2.8G	81.7
DW Conv.-S	22M	2.7G	81.9

position embeddings used in local attention, larger than the image window size) so that it could be extended to larger images and downstream tasks with different image sizes.