

# An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale

Alexey Dosovitskiy\*,†, Lucas Beyer\*, Alexander Kolesnikov\*, Dirk  
Weissenborn\*, Xiaohua Zhai\*, Thomas Unterthiner, Mostafa Dehghani, Matthias  
Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby\*,†

\*equal technical contribution, †equal advising

Google Research, Brain Team

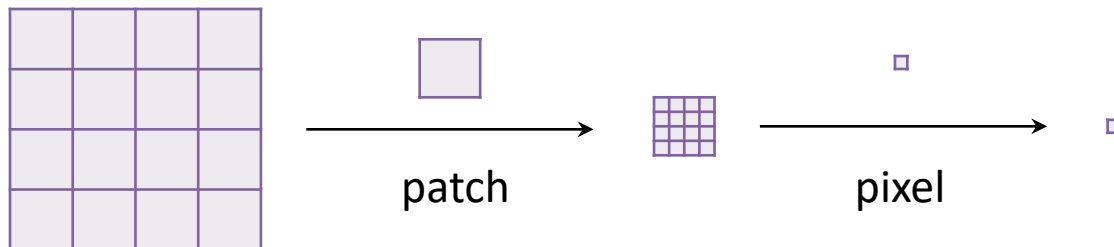
{adosovitskiy, neilhoulby}@google.com

# Motivation

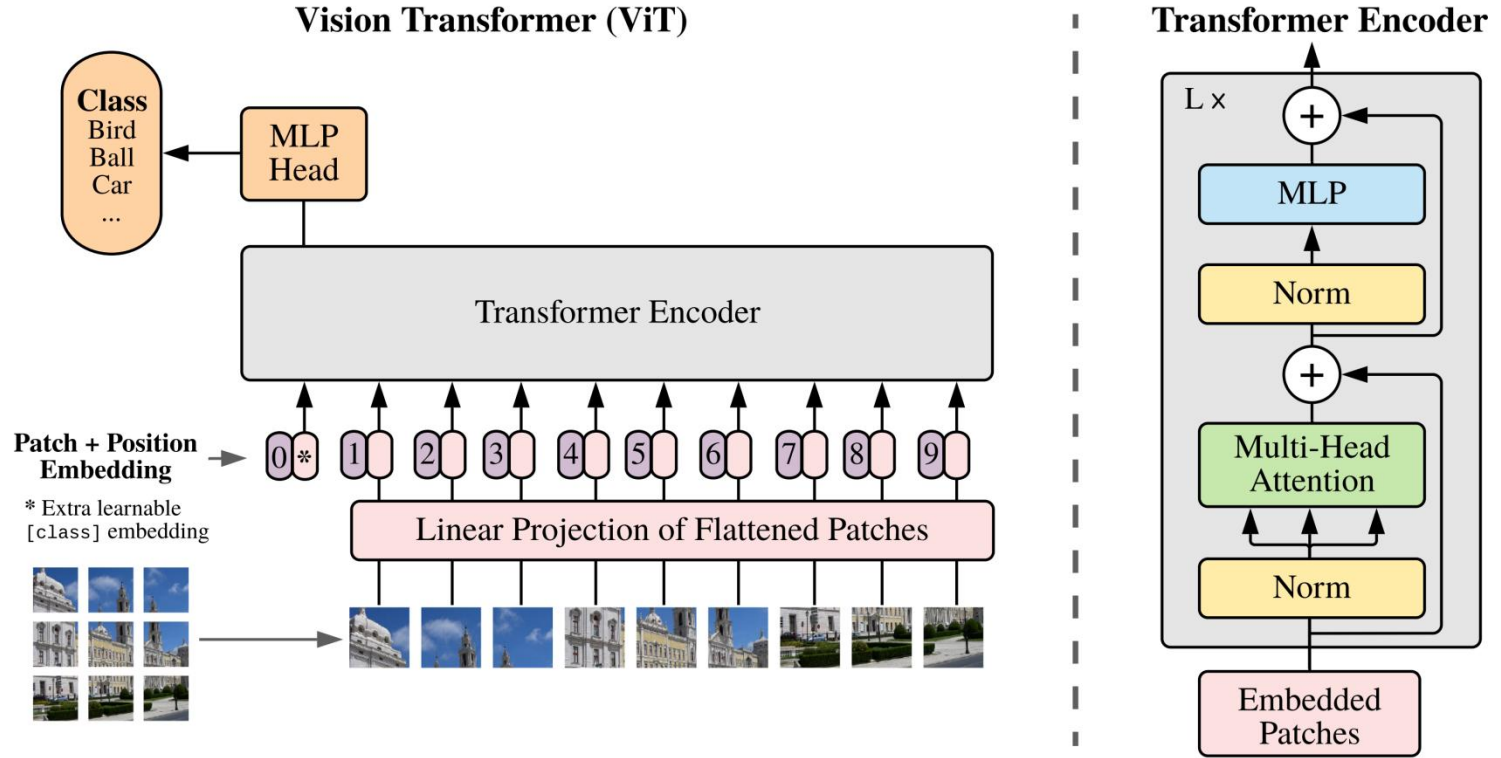
- Transformer architecture has achieved great success in multiple NLP tasks.
  - A larger CNN can lead to overfitting, but such phenomenon hasn't been observed in transformer.
  - Model consisted of pure transformer has not merged in CV (CNN+Transformer).
- ✓ Apply transformer *directly* in CV —————> ViT

# Challenge and Solution

- Transform 2D picture into 1D vector
  - Direct expand ( $224 \times 224 \rightarrow 50176$ )
  - Expand feature map ( $14 \times 14 \rightarrow 196$ )
  - Local window
  - Two self-attention step
- This paper: patch!



# Model



$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \cdots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}}, \quad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D} \quad (1)$$

$$\mathbf{z}'_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \quad \ell = 1 \dots L \quad (2)$$

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, \quad \ell = 1 \dots L \quad (3)$$

$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0) \quad (4)$$

# Finetuning

- a. Remove the pre-trained prediction head.
- b. Attach a feedforward layer.
- c. When feeding higher resolution images, we keep patch size the same.
- d. 2D interpolation of the pre-trained position embeddings (resolution adjustment is the only inductive bias manually injected into ViT).

# Experiments

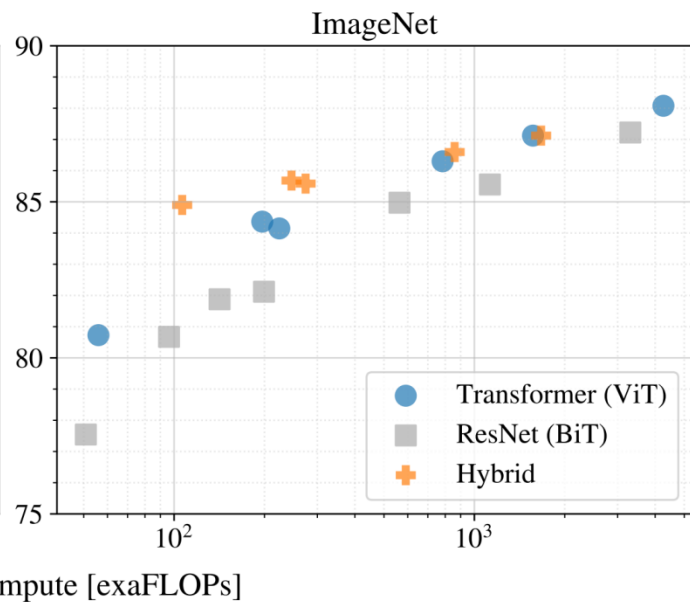
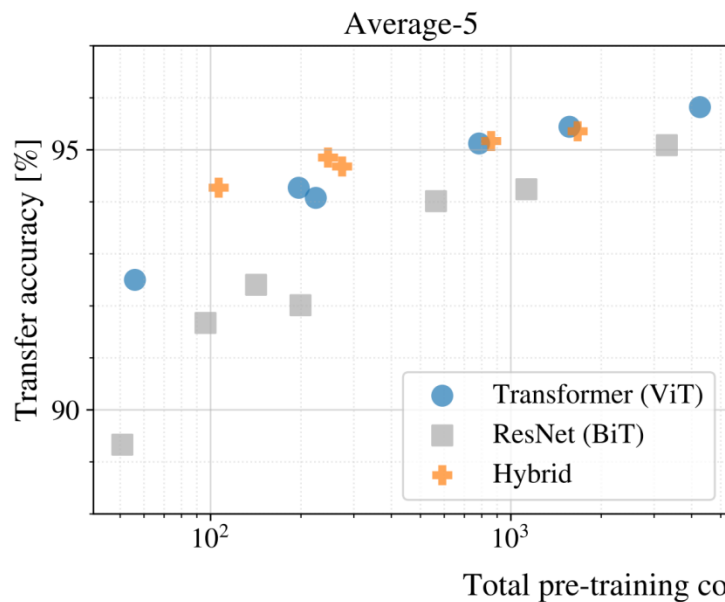
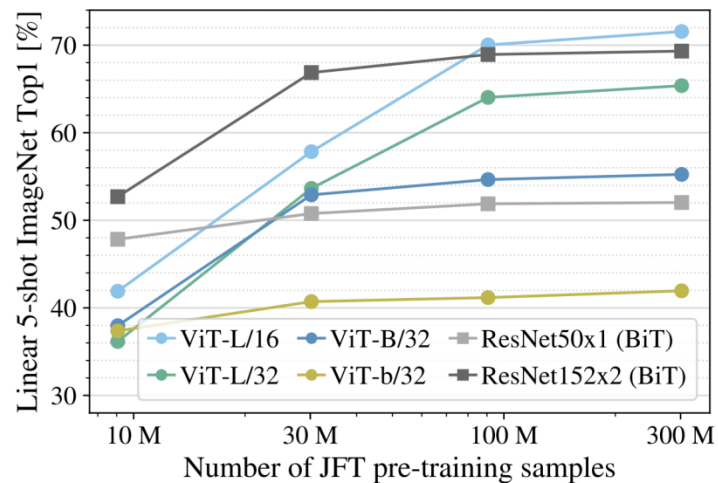
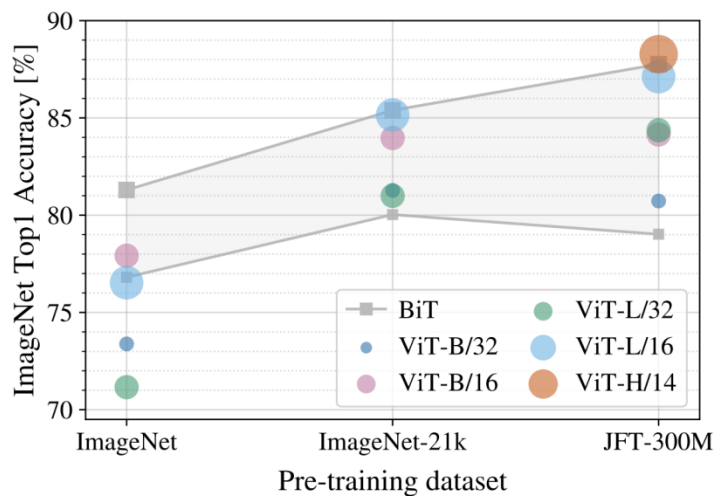
- Comparison to SOTA

Model	Layers	Hidden size $D$	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

Table 1: Details of Vision Transformer model variants.

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	<b>88.55</b> $\pm 0.04$	87.76 $\pm 0.03$	85.30 $\pm 0.02$	87.54 $\pm 0.02$	88.4/88.5*
ImageNet Real	<b>90.72</b> $\pm 0.05$	90.54 $\pm 0.03$	88.62 $\pm 0.05$	90.54	90.55
CIFAR-10	<b>99.50</b> $\pm 0.06$	99.42 $\pm 0.03$	99.15 $\pm 0.03$	99.37 $\pm 0.06$	—
CIFAR-100	<b>94.55</b> $\pm 0.04$	93.90 $\pm 0.05$	93.25 $\pm 0.05$	93.51 $\pm 0.08$	—
Oxford-IIIT Pets	<b>97.56</b> $\pm 0.03$	97.32 $\pm 0.11$	94.67 $\pm 0.15$	96.62 $\pm 0.23$	—
Oxford Flowers-102	99.68 $\pm 0.02$	<b>99.74</b> $\pm 0.00$	99.61 $\pm 0.02$	99.63 $\pm 0.03$	—
VTAB (19 tasks)	<b>77.63</b> $\pm 0.23$	76.28 $\pm 0.46$	72.72 $\pm 0.21$	76.29 $\pm 1.70$	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

- Pre-training setups



- Positional embedding

- Providing no positional information: Considering the inputs as *a bag of patches*.



- 1-dimensional positional embedding: Considering the inputs as *a sequence of patches in the raster order*.



- 2-dimensional positional embedding: Considering the inputs as *a grid of patches in two dimensions*.



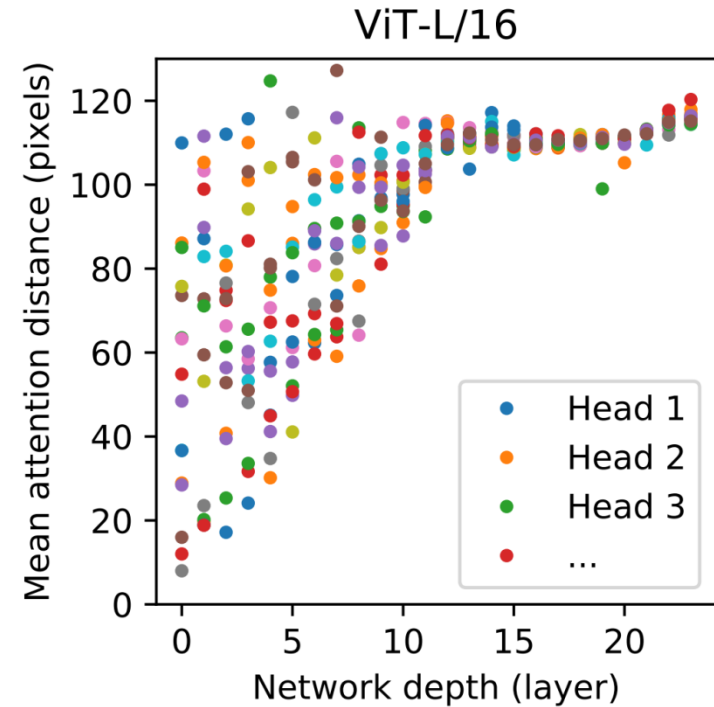
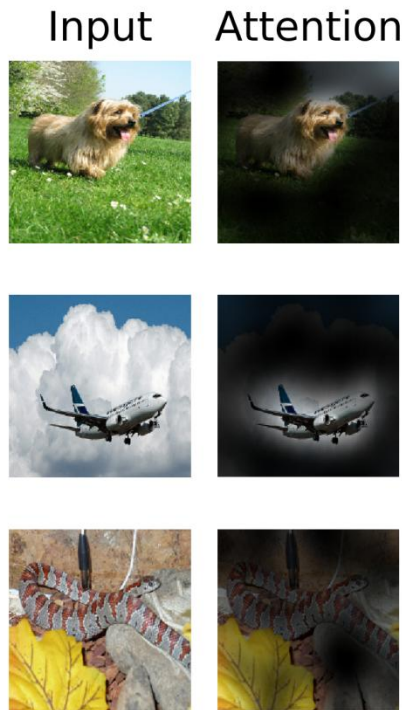
- Relative positional embeddings: Considering the *relative distance* between patches to encode the spatial information as instead of their absolute position.





Pos. Emb.	Default/Stem	Every Layer	Every Layer-Shared
No Pos. Emb.	0.61382	N/A	N/A
1-D Pos. Emb.	0.64206	0.63964	0.64292
2-D Pos. Emb.	0.64001	0.64046	0.64022
Rel. Pos. Emb.	0.64032	N/A	N/A

- Attention



$$\text{Mean attention distance} = d_{AB} \times W_{AB}$$

# Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

Ze Liu<sup>†\*</sup> Yutong Lin<sup>†\*</sup> Yue Cao<sup>\*</sup> Han Hu<sup>\*‡</sup> Yixuan Wei<sup>†</sup> Zheng Zhang Stephen  
Lin Baining Guo

Microsoft Research Asia

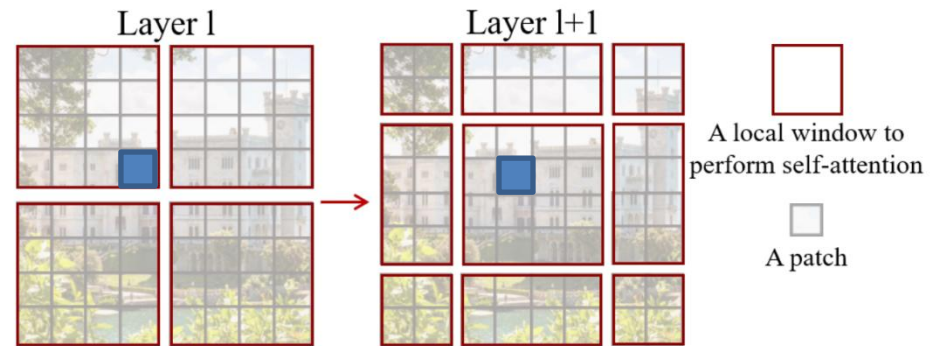
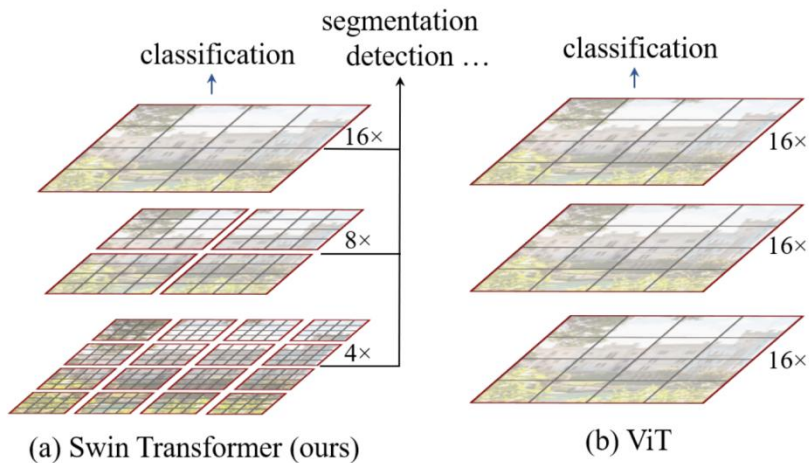
{v-zeliu1,v-yutlin,yuecao,hanhu,v-yixwe,zhez,stevelin,bainguo}@microsoft.com

# Motivations

- ViT is only designed for image classification.
- A transformer based network for various CV tasks, like image classification, object detection, semantic segmentation, etc.

# Challenges and Solutions

- Large variations in the scale of visual entities.
- High resolution of pixels in images.



1. Self-attention within local window
2. Patch merging
3. Linear computation complexity to image size

1. Shifted window (in NLP?)
2. *Global* modeling

# Model

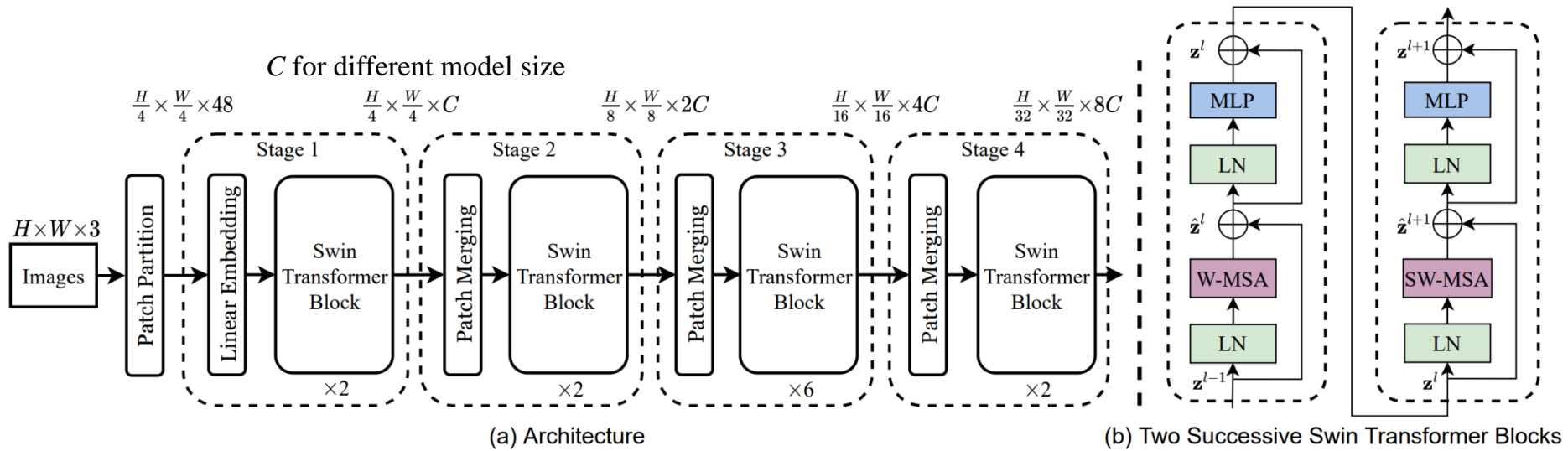
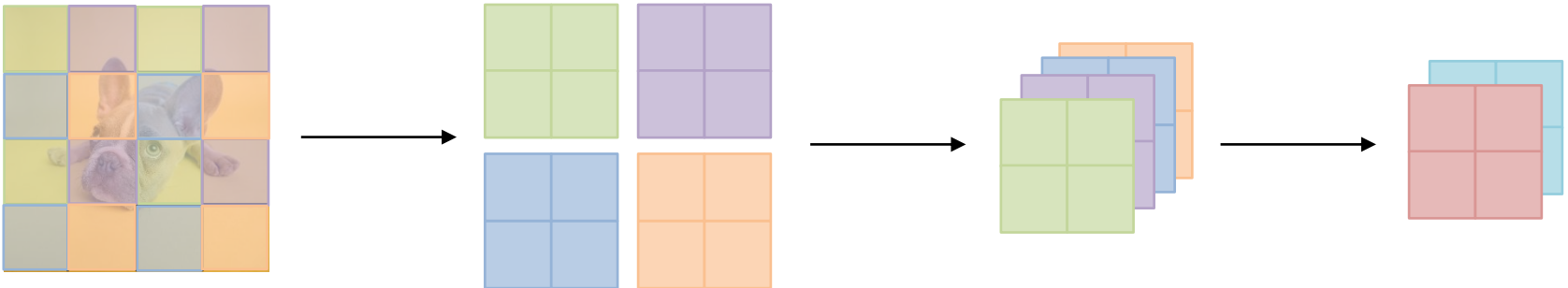
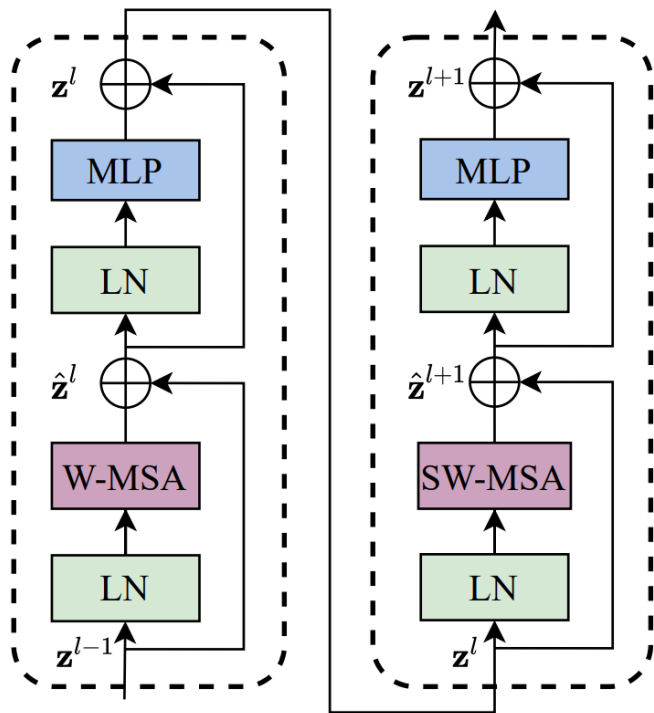


Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

## Patch Merging

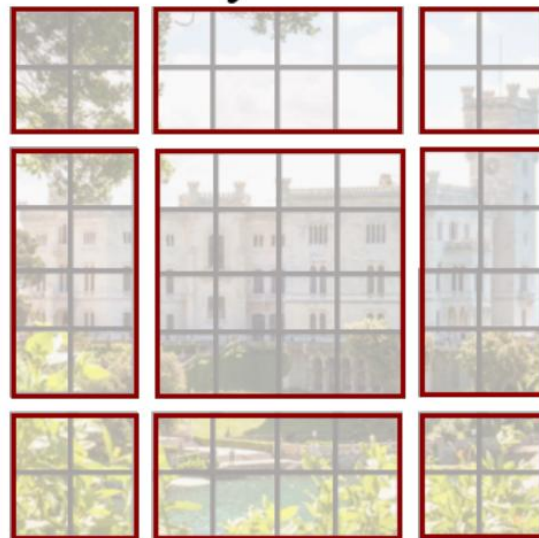


## Shifted window



Problem: Window sizes are different

Layer l+1



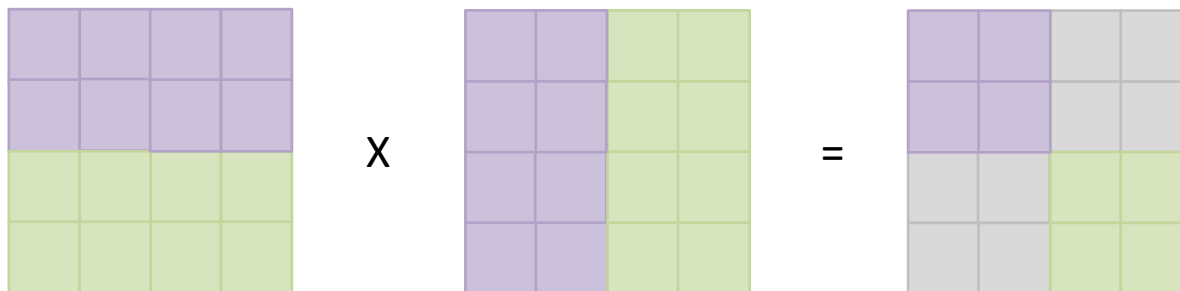
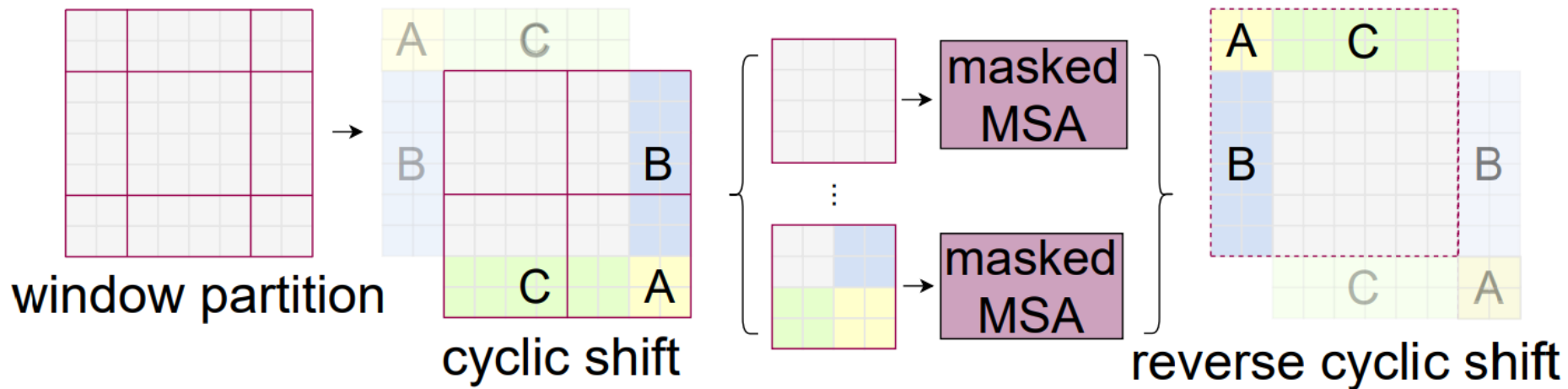
$$\hat{\mathbf{z}}^l = \text{W-MSA} (\text{LN} (\mathbf{z}^{l-1})) + \mathbf{z}^{l-1},$$

$$\mathbf{z}^l = \text{MLP} (\text{LN} (\hat{\mathbf{z}}^l)) + \hat{\mathbf{z}}^l,$$

$$\hat{\mathbf{z}}^{l+1} = \text{SW-MSA} (\text{LN} (\mathbf{z}^l)) + \mathbf{z}^l,$$

$$\mathbf{z}^{l+1} = \text{MLP} (\text{LN} (\hat{\mathbf{z}}^{l+1})) + \hat{\mathbf{z}}^{l+1},$$

# Partition and masking



# Partition and masking

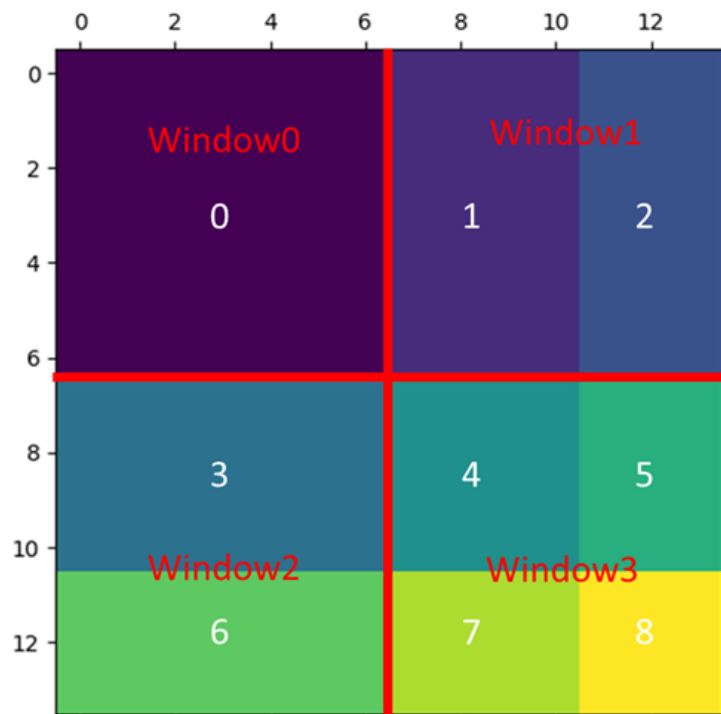
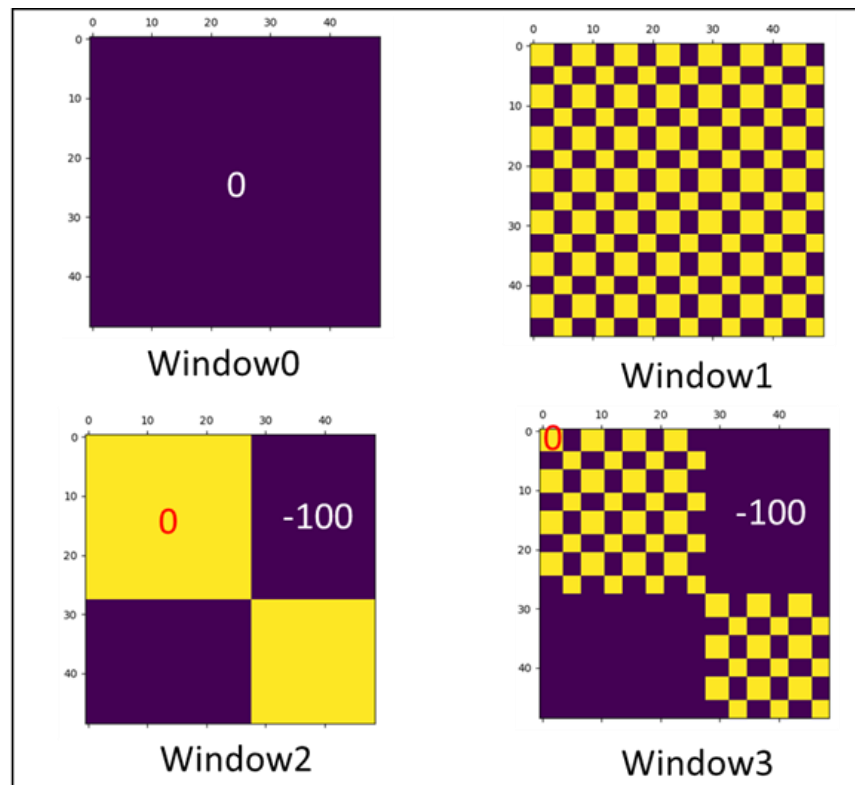


Image Mask  
(14x14, window 7x7, shift 3)



Attn Mask



# Experiments

## Image classification

(a) Regular ImageNet-1K trained models					
method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
RegNetY-4G [48]	224 <sup>2</sup>	21M	4.0G	1156.7	80.0
RegNetY-8G [48]	224 <sup>2</sup>	39M	8.0G	591.6	81.7
RegNetY-16G [48]	224 <sup>2</sup>	84M	16.0G	334.7	82.9
EffNet-B3 [58]	300 <sup>2</sup>	12M	1.8G	732.1	81.6
EffNet-B4 [58]	380 <sup>2</sup>	19M	4.2G	349.4	82.9
EffNet-B5 [58]	456 <sup>2</sup>	30M	9.9G	169.1	83.6
EffNet-B6 [58]	528 <sup>2</sup>	43M	19.0G	96.9	84.0
EffNet-B7 [58]	600 <sup>2</sup>	66M	37.0G	55.1	84.3
ViT-B/16 [20]	384 <sup>2</sup>	86M	55.4G	85.9	77.9
ViT-L/16 [20]	384 <sup>2</sup>	307M	190.7G	27.3	76.5
DeiT-S [63]	224 <sup>2</sup>	22M	4.6G	940.4	79.8
DeiT-B [63]	224 <sup>2</sup>	86M	17.5G	292.3	81.8
DeiT-B [63]	384 <sup>2</sup>	86M	55.4G	85.9	83.1
Swin-T	224 <sup>2</sup>	29M	4.5G	755.2	81.3
Swin-S	224 <sup>2</sup>	50M	8.7G	436.9	83.0
Swin-B	224 <sup>2</sup>	88M	15.4G	278.1	83.5
Swin-B	384 <sup>2</sup>	88M	47.0G	84.7	84.5
(b) ImageNet-22K pre-trained models					
method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
R-101x3 [38]	384 <sup>2</sup>	388M	204.6G	-	84.4
R-152x4 [38]	480 <sup>2</sup>	937M	840.5G	-	85.4
ViT-B/16 [20]	384 <sup>2</sup>	86M	55.4G	85.9	84.0
ViT-L/16 [20]	384 <sup>2</sup>	307M	190.7G	27.3	85.2
Swin-B	224 <sup>2</sup>	88M	15.4G	278.1	85.2
Swin-B	384 <sup>2</sup>	88M	47.0G	84.7	86.4
Swin-L	384 <sup>2</sup>	197M	103.9G	42.1	87.3

Table 1. Comparison of different backbones on ImageNet-1K classification. Throughput is measured using the GitHub repository of [68] and a V100 GPU, following [63].

## Object detection

(a) Various frameworks							
Method	Backbone	AP <sup>box</sup>	AP <sup>box</sup> <sub>50</sub>	AP <sup>box</sup> <sub>75</sub>	#param.	FLOPs	FPS
Cascade	R-50	46.3	64.3	50.5	82M	739G	18.0
Mask R-CNN	Swin-T	<b>50.5</b>	<b>69.3</b>	<b>54.9</b>	86M	745G	15.3
ATSS	R-50	43.5	61.9	47.0	32M	205G	28.3
	Swin-T	<b>47.2</b>	<b>66.5</b>	<b>51.3</b>	36M	215G	22.3
RepPointsV2	R-50	46.5	64.6	50.3	42M	274G	13.6
	Swin-T	<b>50.0</b>	<b>68.5</b>	<b>54.2</b>	45M	283G	12.0
Sparse	R-50	44.5	63.4	48.2	106M	166G	21.0
R-CNN	Swin-T	<b>47.9</b>	<b>67.3</b>	<b>52.3</b>	110M	172G	18.4
(b) Various backbones w. Cascade Mask R-CNN							
	AP <sup>box</sup>	AP <sup>box</sup> <sub>50</sub>	AP <sup>box</sup> <sub>75</sub>	AP <sup>mask</sup>	AP <sup>mask</sup> <sub>50</sub>	AP <sup>mask</sup> <sub>75</sub>	paramFLOPsFPS
DeiT-S <sup>†</sup>	48.0	67.2	51.7	41.4	64.2	44.3	80M 889G 10.4
R50	46.3	64.3	50.5	40.1	61.7	43.4	82M 739G 18.0
Swin-T	<b>50.5</b>	<b>69.3</b>	<b>54.9</b>	<b>43.7</b>	<b>66.6</b>	<b>47.1</b>	86M 745G 15.3
X101-32	48.1	66.5	52.4	41.6	63.9	45.2	101M 819G 12.8
Swin-S	<b>51.8</b>	<b>70.4</b>	<b>56.3</b>	<b>44.7</b>	<b>67.9</b>	<b>48.5</b>	107M 838G 12.0
X101-64	48.3	66.4	52.3	41.7	64.0	45.1	140M 972G 10.4
Swin-B	<b>51.9</b>	<b>70.9</b>	<b>56.5</b>	<b>45.0</b>	<b>68.4</b>	<b>48.7</b>	145M 982G 11.6
(c) System-level Comparison							
Method	mini-val		test-dev		#param.		FLOPs
	AP <sup>box</sup>	AP <sup>mask</sup>	AP <sup>box</sup>	AP <sup>mask</sup>			
RepPointsV2* [12]	-	-	52.1	-	-	-	-
GCNet* [7]	51.8	44.7	52.3	45.4	-	1041G	-
RelationNet++* [13]	-	-	52.7	-	-	-	-
SpineNet-190 [21]	52.6	-	52.8	-	164M	1885G	-
ResNeSt-200* [78]	52.5	-	53.3	47.1	-	-	-
EfficientDet-D7 [59]	54.4	-	55.1	-	77M	410G	-
DetectoRS* [46]	-	-	55.7	48.5	-	-	-
YOLOv4 P7* [4]	-	-	55.8	-	-	-	-
Copy-paste [26]	55.9	47.2	56.0	47.4	185M	1440G	-
X101-64 (HTC++)	52.3	46.0	-	-	155M	1033G	-
Swin-B (HTC++)	56.4	49.1	-	-	160M	1043G	-
Swin-L (HTC++)	57.1	49.5	57.7	50.2	284M	1470G	-
Swin-L (HTC++)*	<b>58.0</b>	<b>50.4</b>	<b>58.7</b>	<b>51.1</b>	284M	-	-

# Ablation studies

	ImageNet		COCO		ADE20k
	top-1	top-5	AP <sup>box</sup>	AP <sup>mask</sup>	mIoU
w/o shifting	80.2	95.1	47.7	41.5	43.3
shifted windows	<b>81.3</b>	<b>95.6</b>	<b>50.5</b>	<b>43.7</b>	<b>46.1</b>
no pos.	80.1	94.9	49.2	42.6	43.8
abs. pos.	80.5	95.2	49.0	42.4	43.2
abs.+rel. pos.	81.3	95.6	50.2	43.4	44.0
rel. pos. w/o app.	79.3	94.7	48.2	41.9	44.1
rel. pos.	<b>81.3</b>	<b>95.6</b>	<b>50.5</b>	<b>43.7</b>	<b>46.1</b>

Table 4. Ablation study on the *shifted windows* approach and different position embedding methods on three benchmarks, using the Swin-T architecture. w/o shifting: all self-attention modules adopt regular window partitioning, without *shifting*; abs. pos.: absolute position embedding term of ViT; rel. pos.: the default settings with an additional relative position bias term (see Eq. (4)); app.: the first scaled dot-product term in Eq. (4).