

# **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

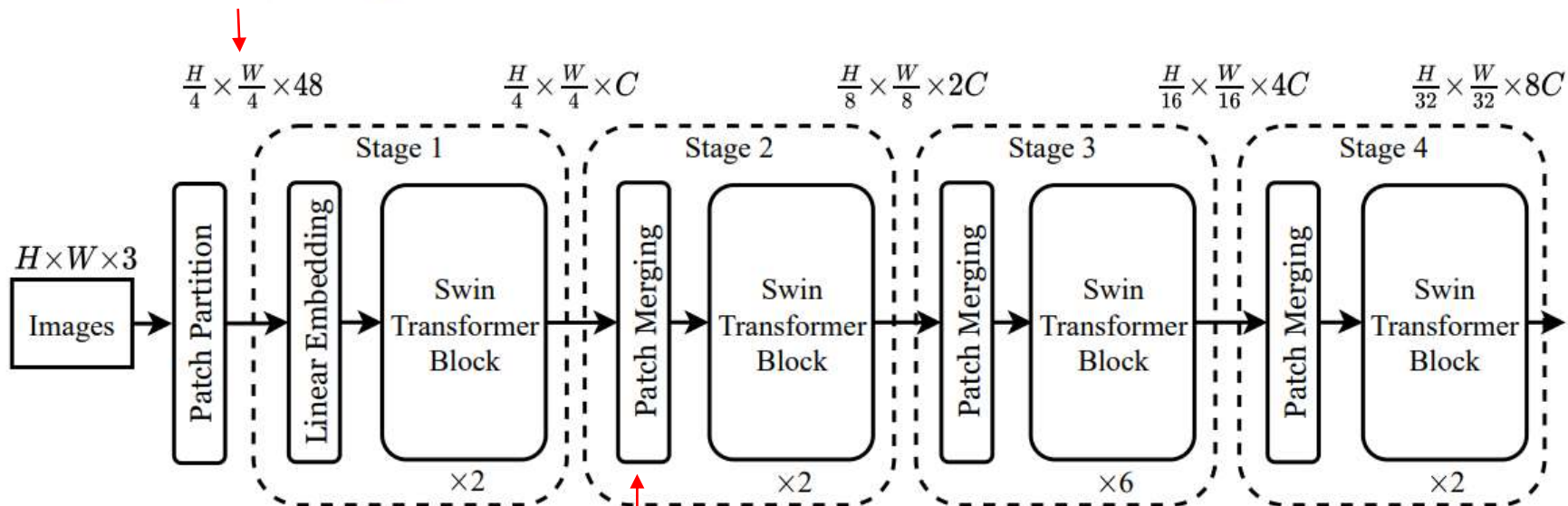
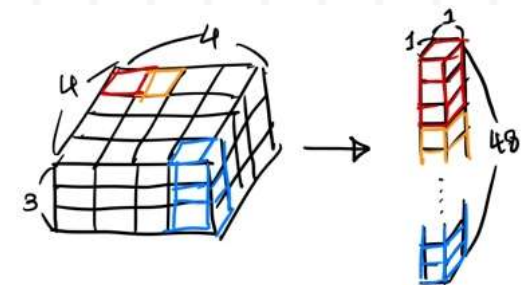
`{v-zeliu1, v-yutlin, yuecao, hanhu, v-yixwe, zhez, stevelin, bainguo}@microsoft.com`

# Swin Transformer

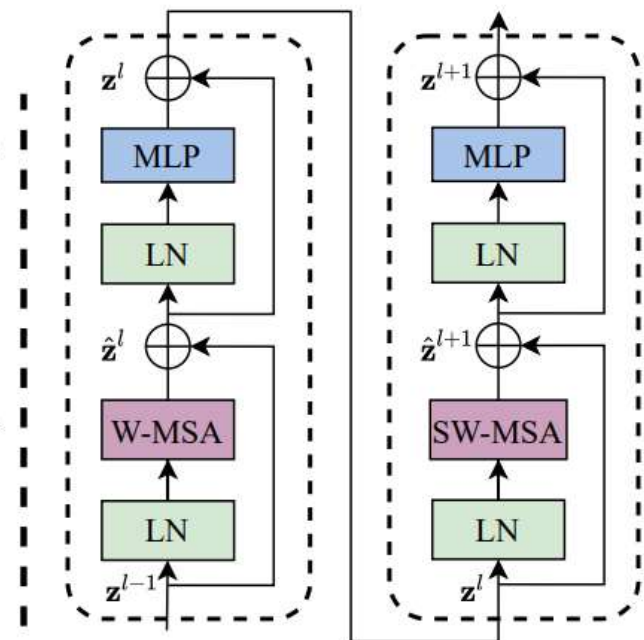
---

- General-purpose **backbone** for computer vision.
- **Hierarchical** Transformer whose representation is computed with **shifted windows**.

# Architecture

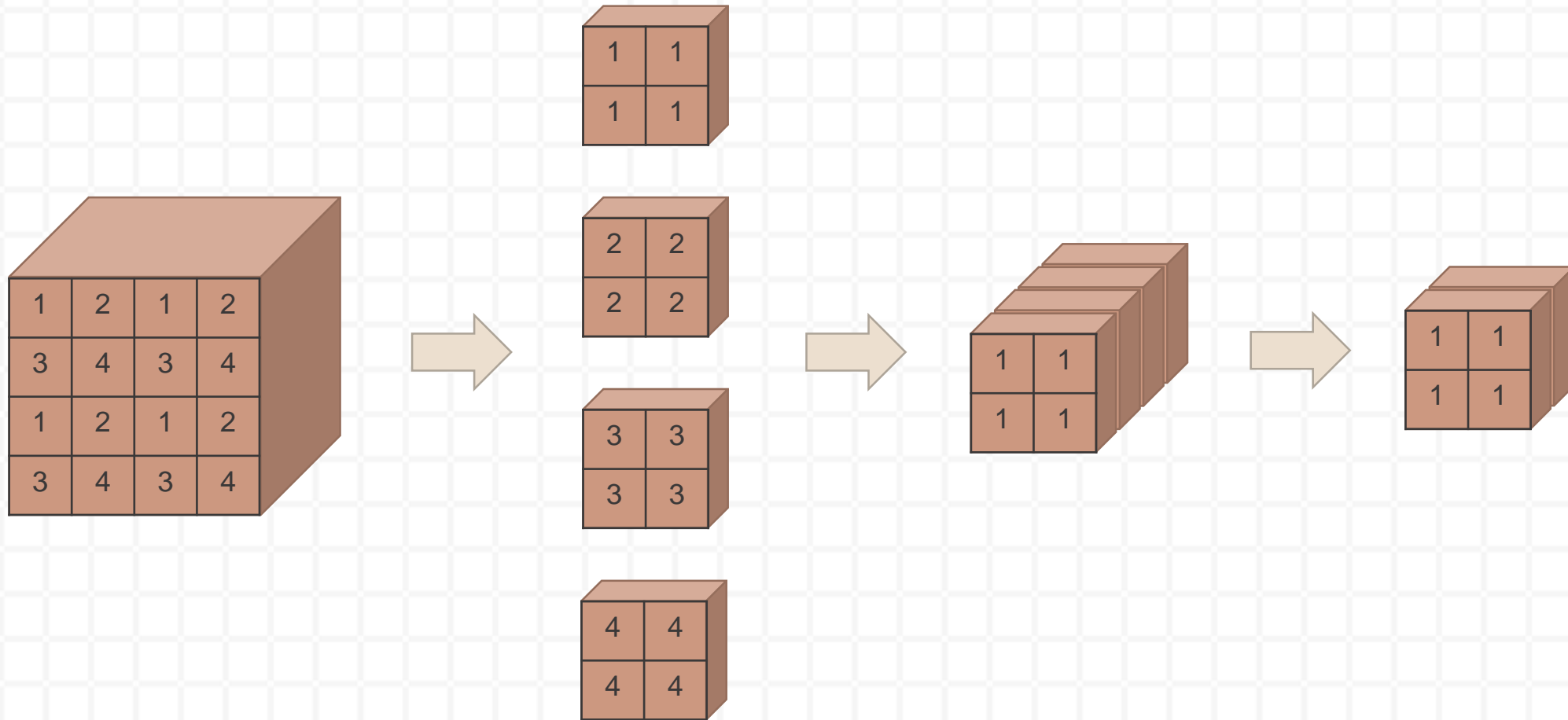


(a) Architecture



(b) Two Successive Swin Transformer Blocks

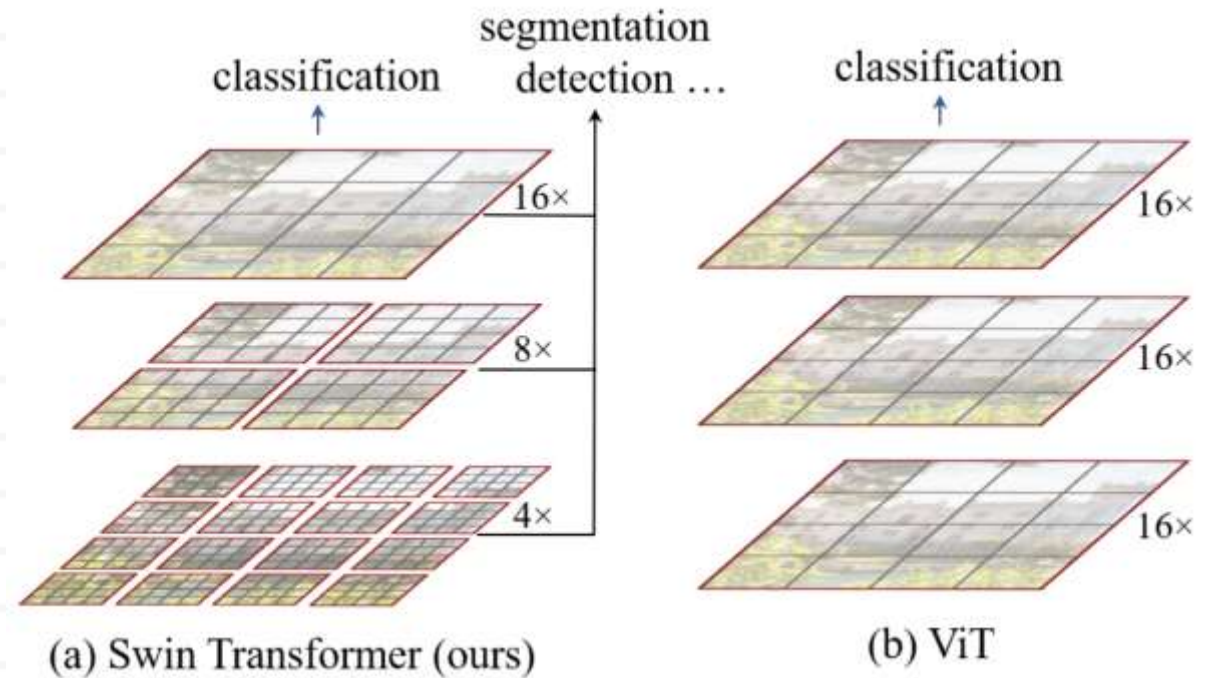
# Patch Merging



# Self-attention in non-overlapped windows

- $M \times M$  ( $7 \times 7$ )
- Propose to compute self-attention within local windows.

$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C,$$
$$\Omega(\text{W-MSA}) = 4hwC^2 + 2M^2hwC,$$



# Shifted window partitioning in successive blocks

- The window-based self-attention module lacks connections across windows, which limits its modeling power.
- W-MSA
- SW-MSA

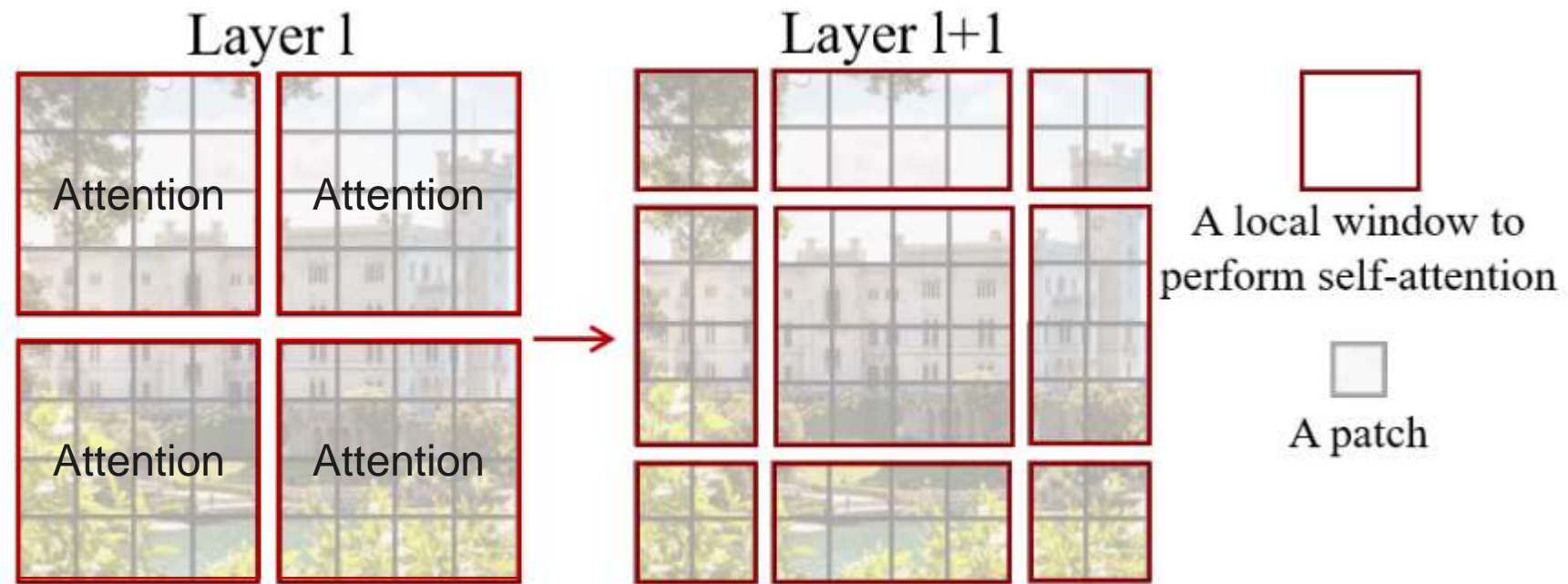
$$\hat{\mathbf{z}}^l = \text{W-MSA}(\text{LN}(\mathbf{z}^{l-1})) + \mathbf{z}^{l-1},$$

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

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

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

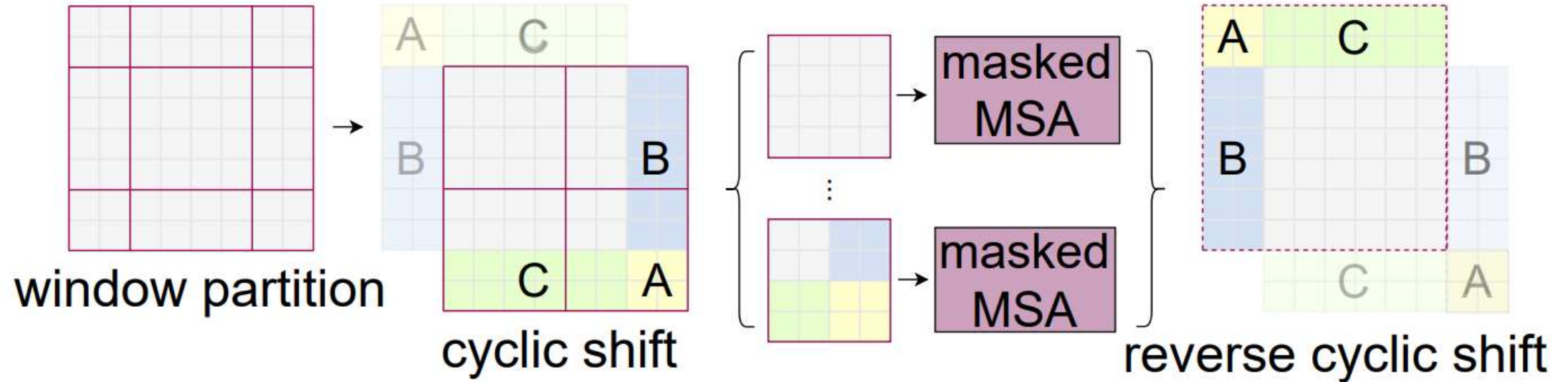
# Swin Transformer





# Efficient batch computation for shifted configuration

- if.. smaller than  $M \times M$  ??
  - Add padding -> increased computation
  - **Cyclic-Shift**





0	1	2
3	4	5
6	7	8

4	5	3
7	8	6
1	2	0

4	5	3
7	8	6
1	2	0

Q

4  
4  
4  
4

5  
3  
5  
3

7  
7  
1  
1

8  
6  
2  
0

K(transposed)

4 4 4 4

5 3 5 3

7 7 1 1

8 6 2 0

4	4	4	4
4	4	4	4
4	4	4	4
4	4	4	4

5		5	
	3		3
5		5	
	3		3

7	7		
7	7		
		1	1
		1	1

8			
	6		
		2	
			0

# Relative position bias

---

- Including a **relative position** bias :  $B \in \mathbb{R}^{M^2 \times M^2}$

$$Attention(Q, K, V) = Softmax\left(\frac{QK^T}{\sqrt{d}} + \mathbf{B}\right)V$$

# Swin Transformer

---

- We also introduce **Swin-T**, **Swin-S** and **Swin-L**, which are versions of about **0.25x**, **0.5x** and **2x** the model size and computational complexity, respectively.
- $M = 7 / d = 32$

Swin-T:  $C = 96$ , layer numbers =  $\{2, 2, 6, 2\}$

Swin-S:  $C = 96$ , layer numbers =  $\{2, 2, 18, 2\}$

Swin-B:  $C = 128$ , layer numbers =  $\{2, 2, 18, 2\}$

Swin-L:  $C = 192$ , layer numbers =  $\{2, 2, 18, 2\}$

# Experiments - ImageNet

**(a) Regular ImageNet-1K trained models**

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
RegNetY-4G [47]	224 <sup>2</sup>	21M	4.0G	1156.7	80.0
RegNetY-8G [47]	224 <sup>2</sup>	39M	8.0G	591.6	81.7
RegNetY-16G [47]	224 <sup>2</sup>	84M	16.0G	334.7	82.9
EffNet-B3 [57]	300 <sup>2</sup>	12M	1.8G	732.1	81.6
EffNet-B4 [57]	380 <sup>2</sup>	19M	4.2G	349.4	82.9
EffNet-B5 [57]	456 <sup>2</sup>	30M	9.9G	169.1	83.6
EffNet-B6 [57]	528 <sup>2</sup>	43M	19.0G	96.9	84.0
EffNet-B7 [57]	600 <sup>2</sup>	66M	37.0G	55.1	84.3
ViT-B/16 [19]	384 <sup>2</sup>	86M	55.4G	85.9	77.9
ViT-L/16 [19]	384 <sup>2</sup>	307M	190.7G	27.3	76.5
DeiT-S [60]	224 <sup>2</sup>	22M	4.6G	940.4	79.8
DeiT-B [60]	224 <sup>2</sup>	86M	17.5G	292.3	81.8
DeiT-B [60]	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	83.3
Swin-B	384 <sup>2</sup>	88M	47.0G	84.7	84.2

14.2 million images and 22K classes.

**(b) ImageNet-22K pre-trained models**

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
R-101x3 [37]	384 <sup>2</sup>	388M	204.6G	-	84.4
R-152x4 [37]	480 <sup>2</sup>	937M	840.5G	-	85.4
ViT-B/16 [19]	384 <sup>2</sup>	86M	55.4G	85.9	84.0
ViT-L/16 [19]	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	86.0
Swin-L	384 <sup>2</sup>	197M	103.9G	42.1	86.4



# Experiments - 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>	#param	FLOPs	FPS
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* [11]	-	-	52.1	-	-	-
GCNet* [6]	51.8	44.7	52.3	45.4	-	1041G
RelationNet++* [12]	-	-	52.7	-	-	-
SpineNet-190 [20]	52.6	-	52.8	-	164M	1885G
ResNeSt-200* [75]	52.5	-	53.3	47.1	-	-
EfficientDet-D7 [58]	54.4	-	55.1	-	77M	410G
DetectoRS* [45]	-	-	55.7	48.5	-	-
YOLOv4 P7* [3]	-	-	55.8	-	-	-
Copy-paste [25]	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	-

# Semantic Segmentation

ADE20K		val	test	#param.	FLOPs	FPS
Method	Backbone	mIoU	score			
DANet [22]	ResNet-101	45.2	-	69M	1119G	15.2
DLab.v3+ [10]	ResNet-101	44.1	-	63M	1021G	16.0
ACNet [23]	ResNet-101	45.9	38.5	-		
DNL [68]	ResNet-101	46.0	56.2	69M	1249G	14.8
OCRNet [70]	ResNet-101	45.3	56.0	56M	923G	19.3
UperNet [66]	ResNet-101	44.9	-	86M	1029G	20.1
OCRNet [70]	HRNet-w48	45.7	-	71M	664G	12.5
DLab.v3+ [10]	ResNeSt-101	46.9	55.1	66M	1051G	11.9
DLab.v3+ [10]	ResNeSt-200	48.4	-	88M	1381G	8.1
SETR [78]	T-Large <sup>‡</sup>	50.3	61.7	308M	-	-
UperNet	DeiT-S <sup>†</sup>	44.0	-	52M	1099G	16.2
UperNet	Swin-T	46.1	-	60M	945G	18.5
UperNet	Swin-S	49.3	-	81M	1038G	15.2
UperNet	Swin-B <sup>‡</sup>	51.6	-	121M	1841G	8.7
UperNet	Swin-L <sup>‡</sup>	<b>53.5</b>	<b>62.8</b>	234M	3230G	6.2



# Ablation Study

	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>

method	MSA in a stage (ms)				Arch. (FPS)		
	S1	S2	S3	S4	T	S	B
sliding window (naive)	122.5	38.3	12.1	7.6	183	109	77
sliding window (kernel)	7.6	4.7	2.7	1.8	488	283	187
Performer [13]	4.8	2.8	1.8	1.5	638	370	241
window (w/o shifting)	2.8	1.7	1.2	0.9	770	444	280
shifted window (padding)	3.3	2.3	1.9	2.2	670	371	236
shifted window (cyclic)	3.0	1.9	1.3	1.0	755	437	278

	Backbone	ImageNet		COCO		ADE20k
		top-1	top-5	AP <sup>box</sup>	AP <sup>mask</sup>	mIoU
sliding window	Swin-T	81.4	95.6	50.2	43.5	45.8
Performer [13]	Swin-T	79.0	94.2	-	-	-
shifted window	Swin-T	81.3	95.6	50.5	43.7	46.1



# Conclusion

---

- Proposed Swin Transformer which a new vision Transformer.
- Hierarchical feature representation and has linear computational complexity with respect to input image size.
- Shifted window based self-attention is shown to be effective and efficient on vision problems.