

# Swin Transformer

---

Ze Liu Yutong Lin Yue Cao Han Hu Yixuan Wei Zheng Zhang Stephen Lin Baining Guo

Swin Transformer : Hierarchical Vision Transformer using Shifted Windows. ICCV 2021

<https://arxiv.org/abs/2103.14030>

# Swin Transformer

---

 State of the Art Object Detection on COCO test-dev (using additional training data)

 State of the Art Instance Segmentation on COCO test-dev

 State of the Art Semantic Segmentation on ADE20K (using additional training data)

 Ranked #3 Action Classification on Kinetics-400 (using additional training data)

**Image Classification:** Included in this repo. See [get\\_started.md](#) for a quick start.

**Object Detection and Instance Segmentation:** See [Swin Transformer for Object Detection](#).

**Semantic Segmentation:** See [Swin Transformer for Semantic Segmentation](#).

**Video Action Recognition:** See [Video Swin Transformer](#).

**Semi-Supervised Object Detection:** See [Soft Teacher](#).

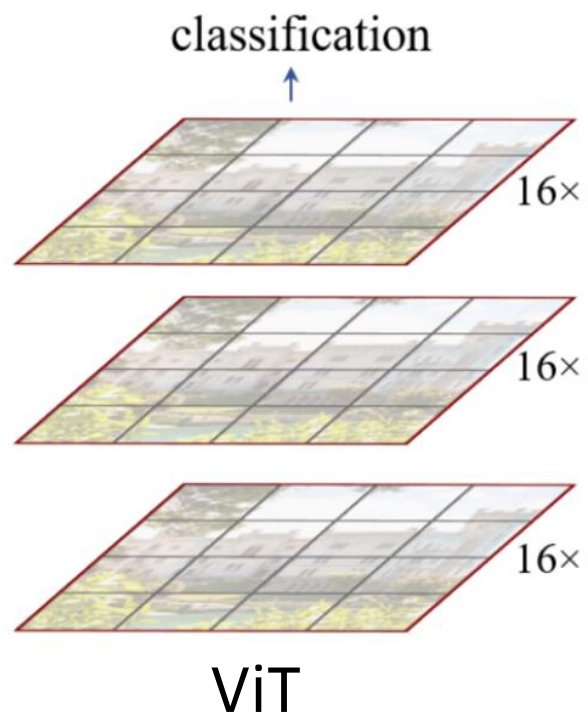
**SSL: Contrastive Learning:** See [Transformer-SSL](#).



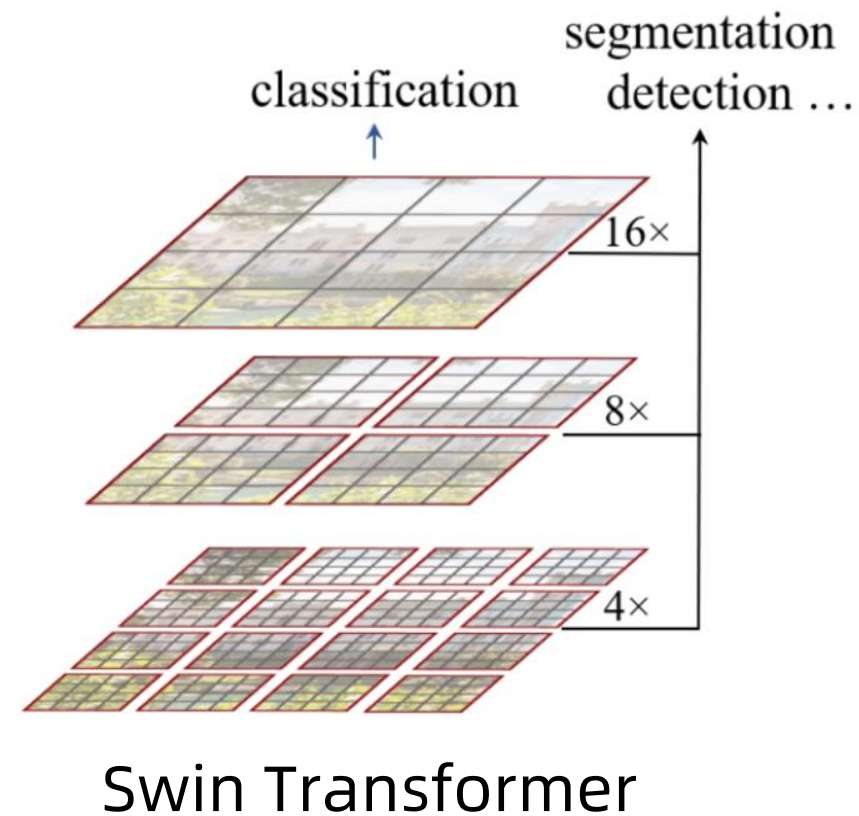
**SSL: Masked Image Modeling:** See [SimMIM](#).

# Problems of ViT

- Does not consider the difference between textual and visual signals
- Mainly for image classification

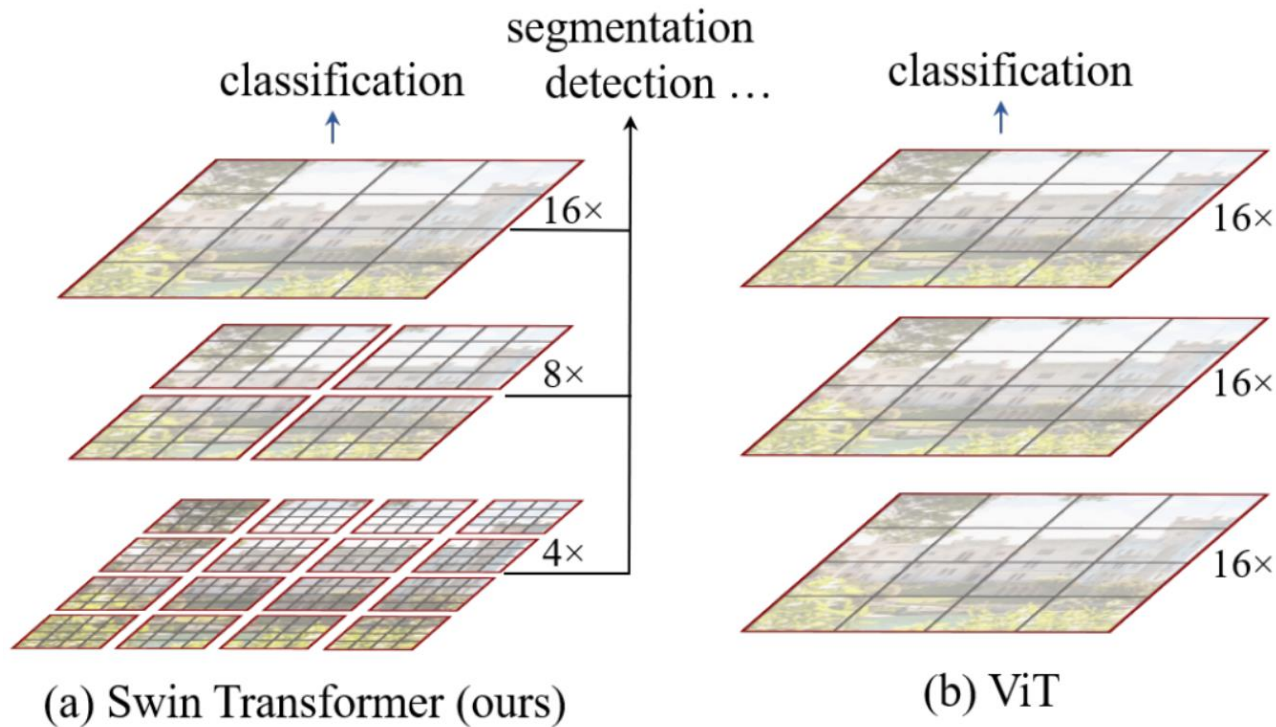


**good priors for visual signals**  
(hierarchy / locality / translation invariance)



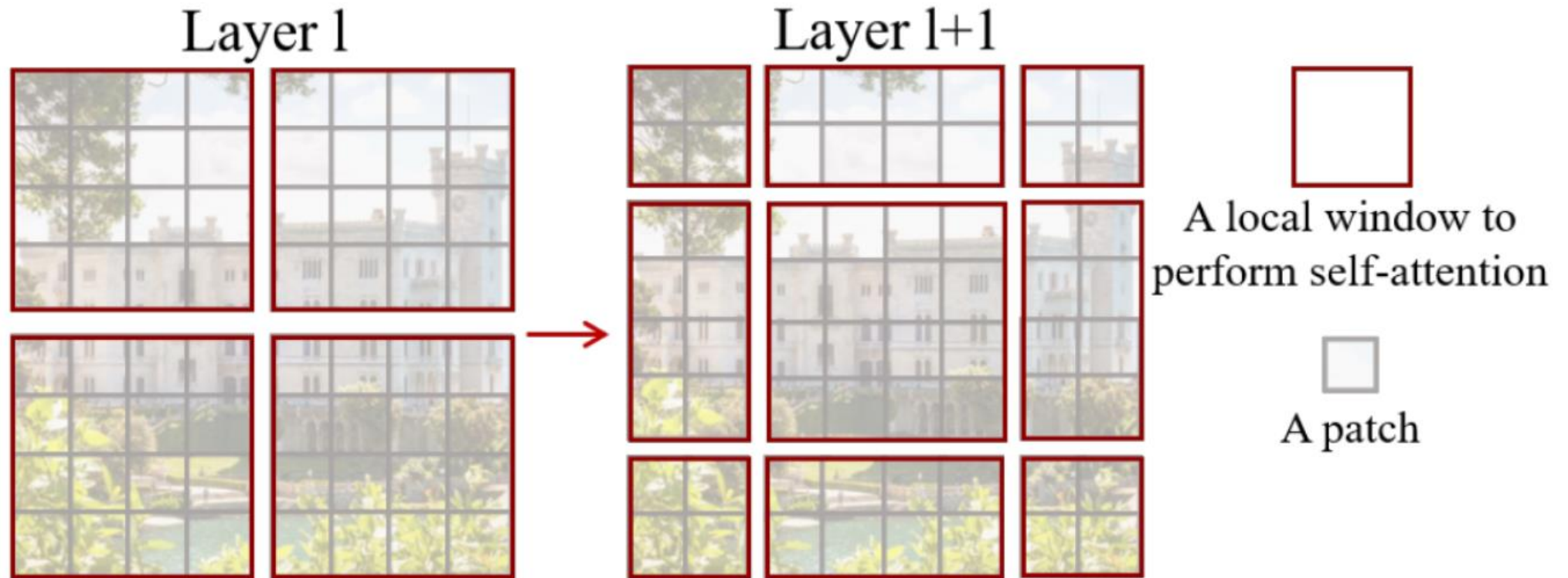
# Key tech innovation : locality by Shifted windows

- Non-overlapped windows (faster real speed than sliding windows)
- Windows are shifted in the next layer



- Hierarchical feature maps
- Windows Multi-Head Self-Attention (W-MSA)  
Shifted Windows Multi-Head Self-Attention (SW-MSA)

# Key tech innovation : locality by Shifted windows



- 1) 自注意的计算在局部的非重叠窗口内进行。不同query会共享同样的key集合，从而对硬件友好
- 2) 在前后两层的Transformer模块中，非重叠窗口的配置相比前一层做了半个窗口的移位，使得上一层中不同窗口的信息进行了交换。



# Self-attention in non-overlapped windows

For efficient modeling, we propose to compute self attention within local windows.

The windows are arranged to evenly partition the image in a non-overlapping manner.

Supposing each window contains  $M \times M$  patches, the computational complexity of a global MSA module and a window based one on an image of  $h \times w$  patches are :

$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C, \quad (1)$$

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

where the former is quadratic to patch number  $hw$ , and the latter is linear when  $M$  is fixed (set to 7 by default). Global self-attention computation is generally unaffordable for a large  $hw$ , while the window based self-attention is scalable.



# The architecture of a Swin Transformer (Swin-T)

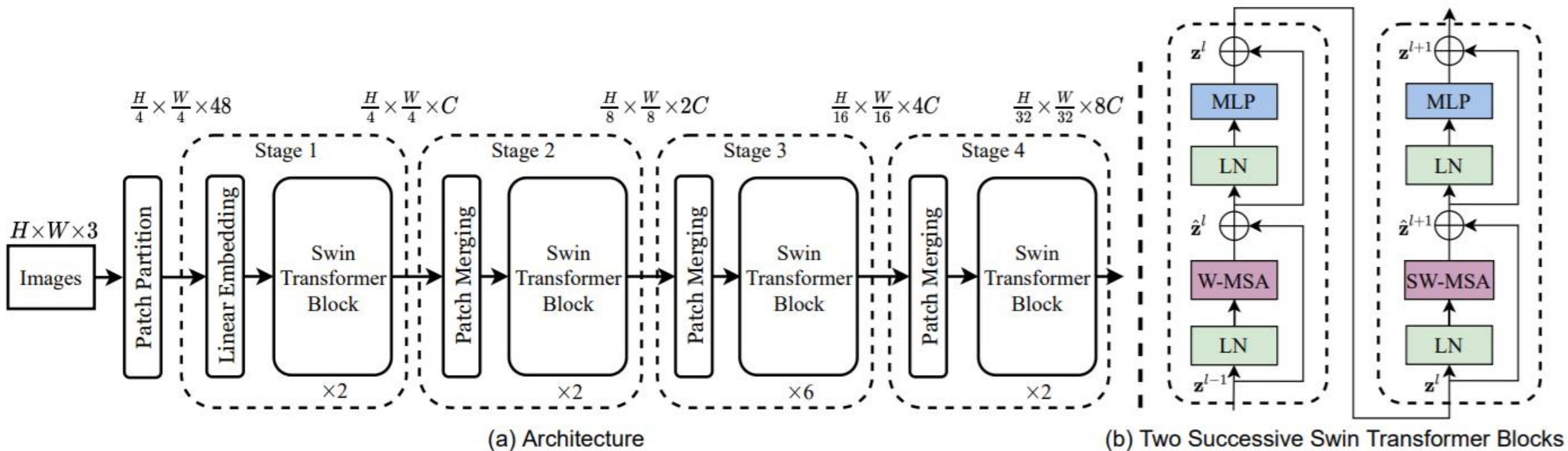
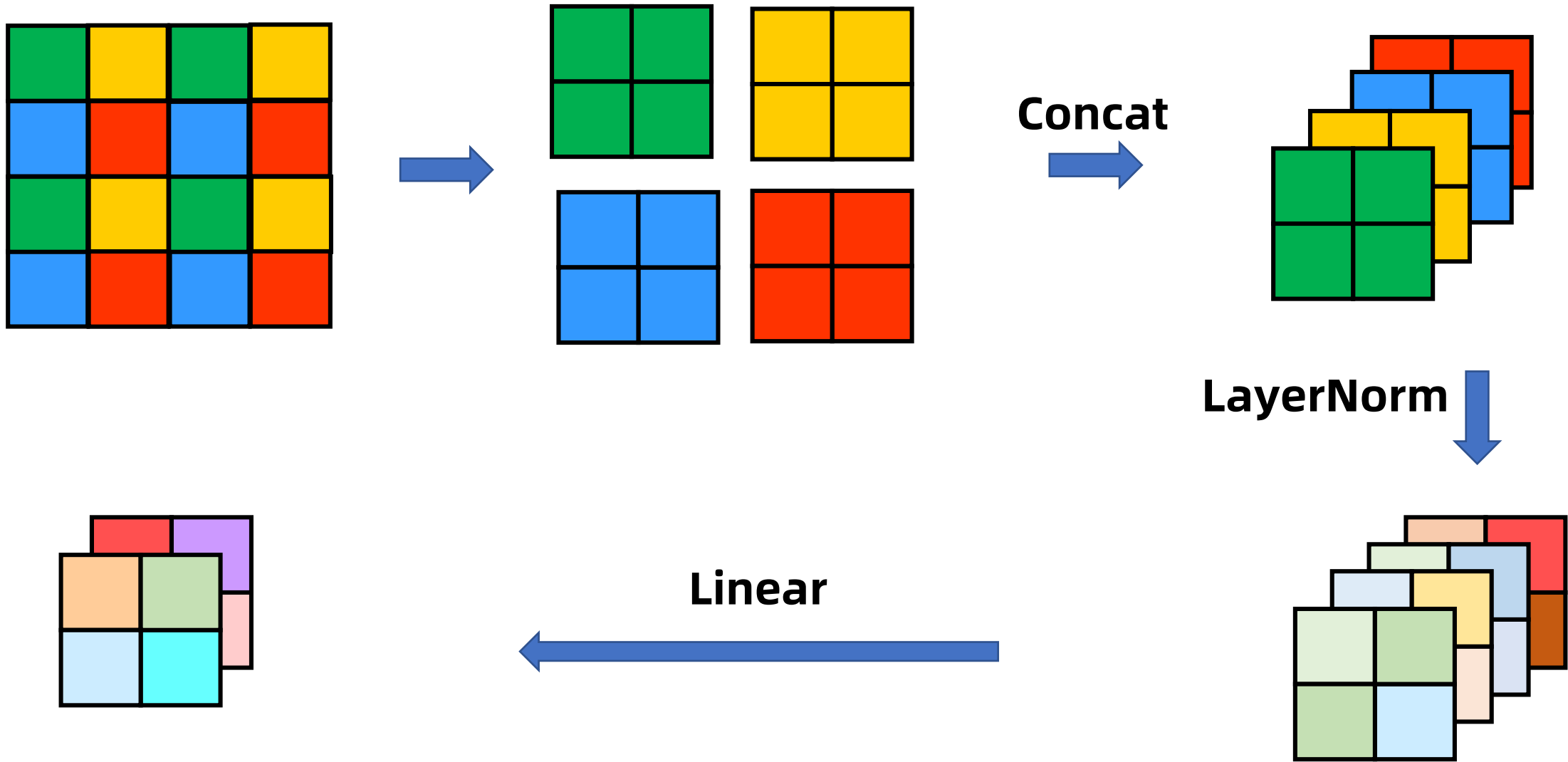


Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

# Patch Merging

经过Patch Merging后，feature map的高和宽会减半，深度会加倍

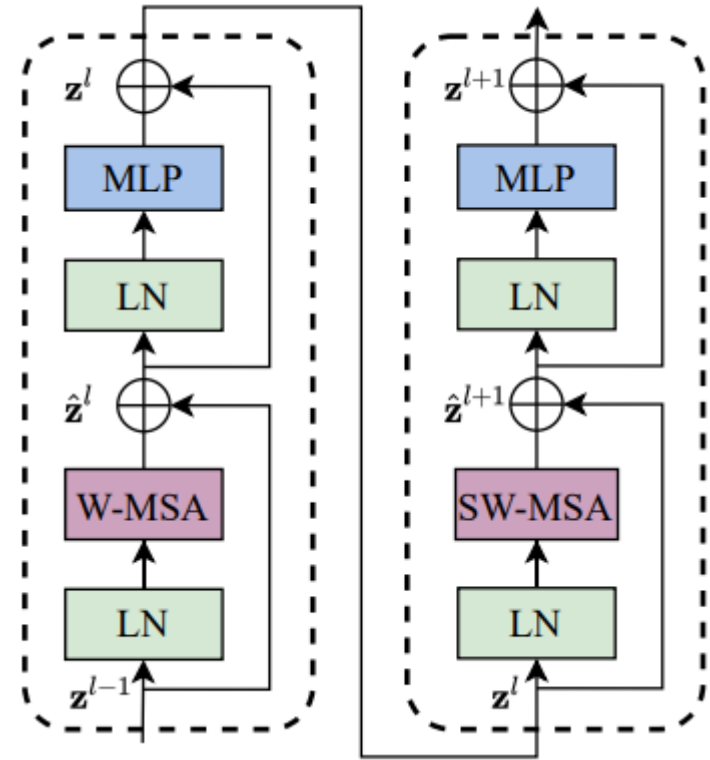


# Swin Transformer block

Swin Transformer is built by replacing the standard **multi-head self attention (MSA)** module in a Transformer block by a module based on shifted windows, with other layers kept the same.

A Swin Transformer block consists of a shifted window based MSA module, followed by a 2-layer MLP with GELU nonlinearity in between.

A LayerNorm (LN) layer is applied **before** each MSA module and each MLP, and a residual connection is applied after each module.



# Shifted window partitioning in successive blocks

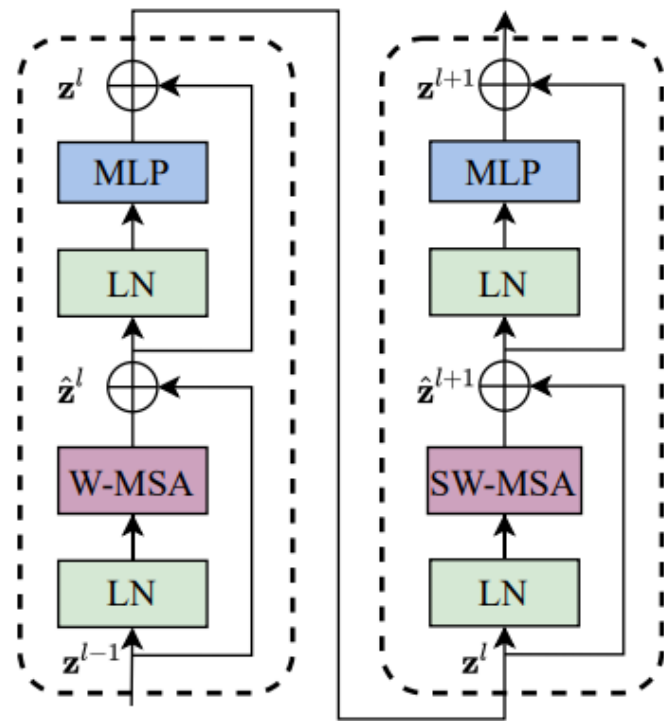
With the shifted window partitioning approach, consecutive Swin Transformer blocks are computed as

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

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

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

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



where  $\hat{\mathbf{z}}^l$  and  $\mathbf{z}^l$  denote the output features of the (S)W-MSA module and the MLP module for block  $l$ , respectively;

## Swin T(Tiny), S(Small), B(Base), L(Large)

	downsp. rate (output size)	Swin-T	Swin-S	Swin-B	Swin-L
stage 1	4× (56×56)	concat 4×4, 96-d, LN	concat 4×4, 96-d, LN	concat 4×4, 128-d, LN	concat 4×4, 192-d, LN
		$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 96, head 3} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 96, head 3} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 128, head 4} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 192, head 6} \end{array} \right] \times 2$
stage 2	8× (28×28)	concat 2×2, 192-d, LN	concat 2×2, 192-d, LN	concat 2×2, 256-d, LN	concat 2×2, 384-d, LN
		$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 192, head 6} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 192, head 6} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 256, head 8} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 384, head 12} \end{array} \right] \times 2$
stage 3	16× (14×14)	concat 2×2, 384-d, LN	concat 2×2, 384-d, LN	concat 2×2, 512-d, LN	concat 2×2, 768-d, LN
		$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 384, head 12} \end{array} \right] \times 6$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 384, head 12} \end{array} \right] \times 18$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 512, head 16} \end{array} \right] \times 18$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 768, head 24} \end{array} \right] \times 18$
stage 4	32× (7×7)	concat 2×2, 768-d, LN	concat 2×2, 768-d, LN	concat 2×2, 1024-d, LN	concat 2×2, 1536-d, LN
		$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 768, head 24} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 768, head 24} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 1024, head 32} \end{array} \right] \times 2$	$\left[ \begin{array}{l} \text{win. sz. } 7 \times 7, \\ \text{dim 1536, head 48} \end{array} \right] \times 2$

Table 7. Detailed architecture specifications.

- win. sz. 7x7表示使用的窗口 (Windows) 的大小
- dim表示feature map的channel深度)
- head表示多头注意力模块中head的个数

# Architecture Variants

- 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 =  $\{2, 2, 18, 2\}$

where  $C$  is the channel number of the hidden layers in the first stage.

The window size is set to  $M = 7$  by default. The query dimension of each head is  $d = 32$ , and the expansion layer of each MLP is  $\alpha = 4$ .



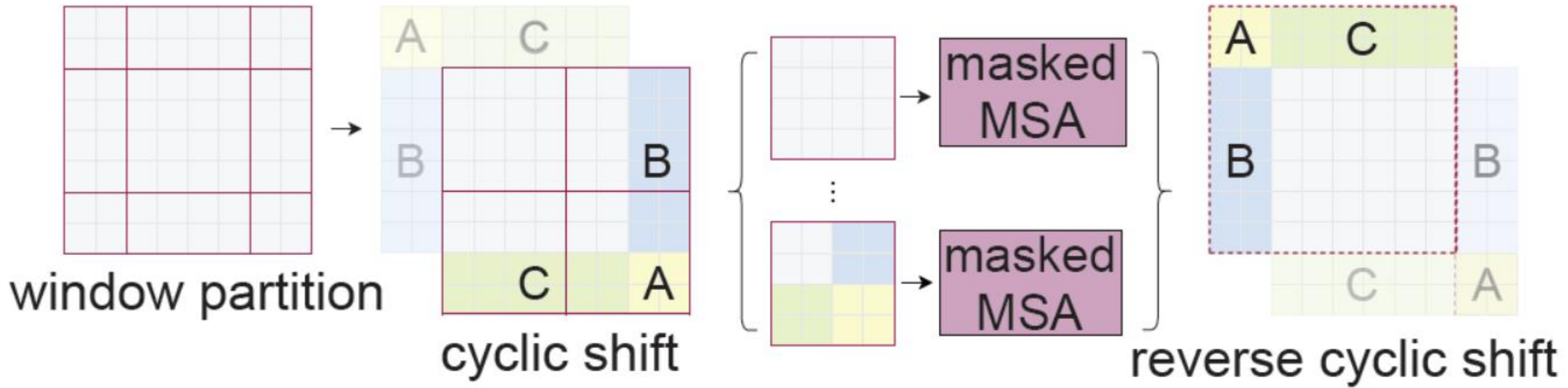
## Efficient batch computation for shifted configuration

An issue with shifted window partitioning is that it will result in more windows, from  $\lceil \frac{h}{M} \rceil \times \lceil \frac{w}{M} \rceil$  to  $(\lceil \frac{h}{M} \rceil + 1) \times (\lceil \frac{w}{M} \rceil + 1)$  in the shifted configuration, and some of the windows will be smaller than  $M \times M$ <sup>4</sup>. A naive solution is to

<sup>4</sup>To make the window size  $(M, M)$  divisible by the feature map size of  $(h, w)$ , bottom-right padding is employed on the feature map if needed.



# Illustration of an efficient batch computation approach for self-attention in shifted window partitioning



**Relative position bias** In computing self-attention, we follow [49, 1, 32, 33] by including a relative position bias  $B \in \mathbb{R}^{M^2 \times M^2}$  to each head in computing similarity:

$$\text{Attention}(Q, K, V) = \text{SoftMax}(QK^T / \sqrt{d} + B)V, \quad (4)$$

where  $Q, K, V \in \mathbb{R}^{M^2 \times d}$  are the *query*, *key* and *value* matrices;  $d$  is the *query/key* dimension, and  $M^2$  is the number of patches in a window. Since the relative position along each axis lies in the range  $[-M + 1, M - 1]$ , we parameterize a smaller-sized bias matrix  $\hat{B} \in \mathbb{R}^{(2M-1) \times (2M-1)}$ , and values in  $B$  are taken from  $\hat{B}$ .

# Gaussian Error Linear Unit (GELU)

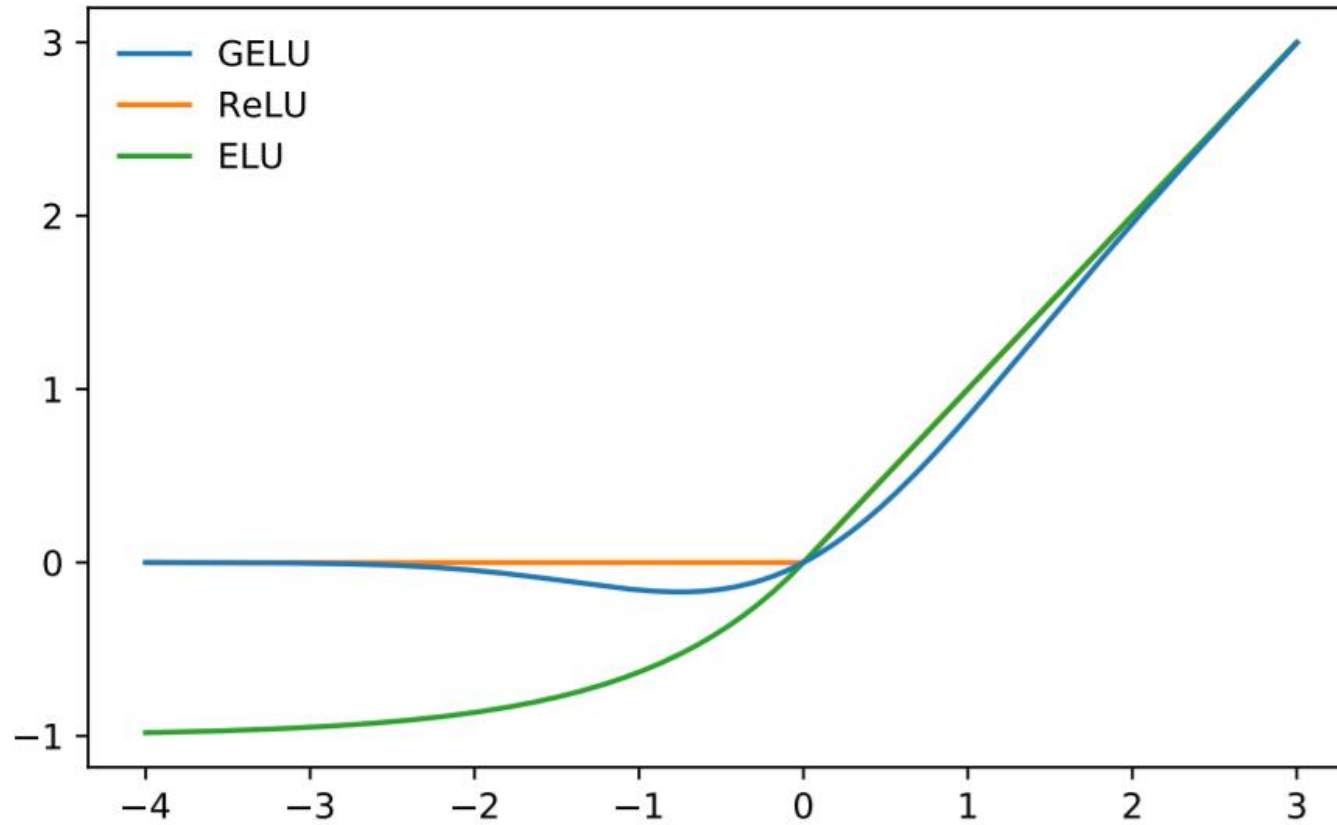


Figure 1: The GELU ( $\mu = 0, \sigma = 1$ ), ReLU, and ELU ( $\alpha = 1$ ).

# Experimental Designs

3 datasets to cover various recognition tasks of different granularities

- **Image-level** ImageNet-1K classification (1.28M images; 1000 classes).
- **Region-level** coco object detection (115K images; 80 classes).
- **Pixel-level** ADE20K semantic segmentation (20K images; 150 classes)

	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).

**(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



## Results on COCO object detection and instance segmentation

(a) Various frameworks							
Method	Backbone	$AP^{box}$	$AP_{50}^{box}$	$AP_{75}^{box}$	#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-CNN	R-50	44.5	63.4	48.2	106M	166G	21.0
	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								
	$AP^{box}$	$AP_{50}^{box}$	$AP_{75}^{box}$	$AP^{mask}$	$AP_{50}^{mask}$	$AP_{75}^{mask}$	param	FLOPsFPS
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



## Results of semantic segmentation on the ADE20K val and test set

ADE20K		val	test	#param.	FLOPs	FPS
Method	Backbone	mIoU	score			
DANet [23]	ResNet-101	45.2	-	69M	1119G	15.2
DLab.v3+ [11]	ResNet-101	44.1	-	63M	1021G	16.0
ACNet [24]	ResNet-101	45.9	38.5	-		
DNL [71]	ResNet-101	46.0	56.2	69M	1249G	14.8
OCRNet [73]	ResNet-101	45.3	56.0	56M	923G	19.3
UperNet [69]	ResNet-101	44.9	-	86M	1029G	20.1
OCRNet [73]	HRNet-w48	45.7	-	71M	664G	12.5
DLab.v3+ [11]	ResNeSt-101	46.9	55.1	66M	1051G	11.9
DLab.v3+ [11]	ResNeSt-200	48.4	-	88M	1381G	8.1
SETR [81]	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