

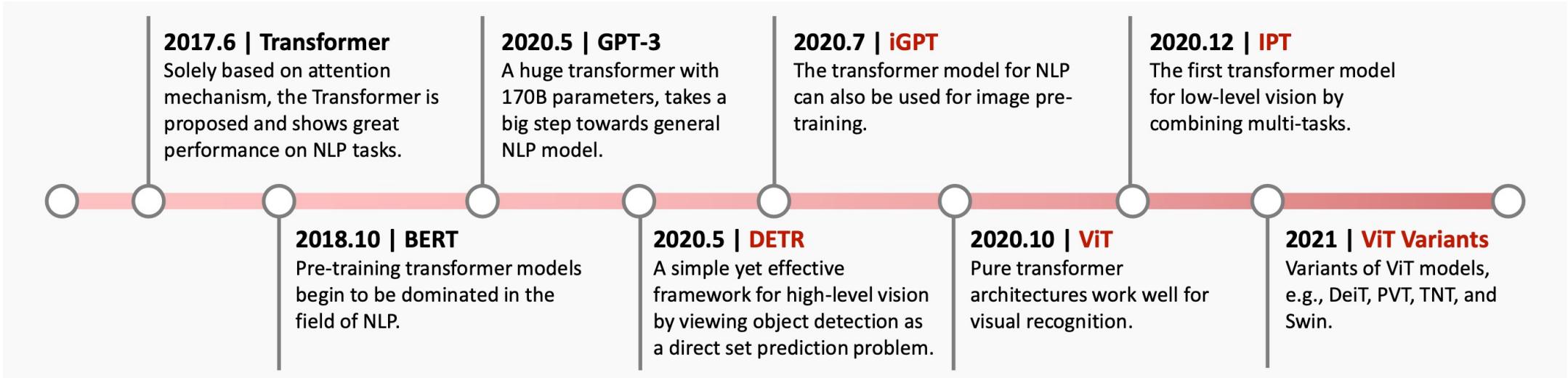
# ViT & Its Variants

10.30

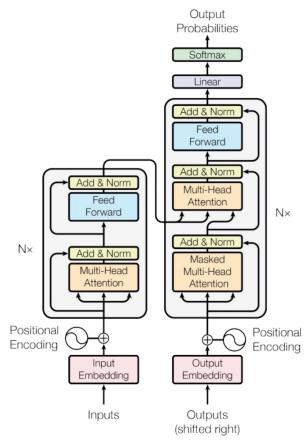
# Content

1. Transformer Development
2. Architecture Of Transformer
  - ① Formulation Of Self-attention
  - ② Formulation Of Multi-head attention
  - ③ Other Key Concepts In Transformer
3. ViT & ViT Variants
  - ① Training data-efficient image transformers & distillation through attention
  - ② Swin Transformer: Hierarchical Vision Transformer using Shifted Windows
  - ③ Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction without Convolutions

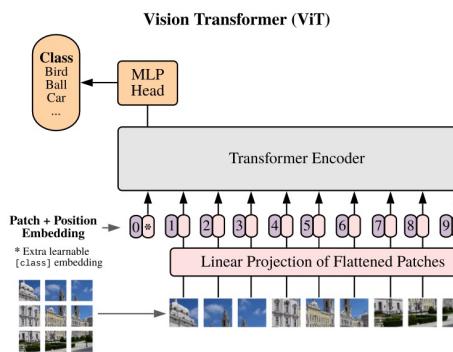
# 1. Development



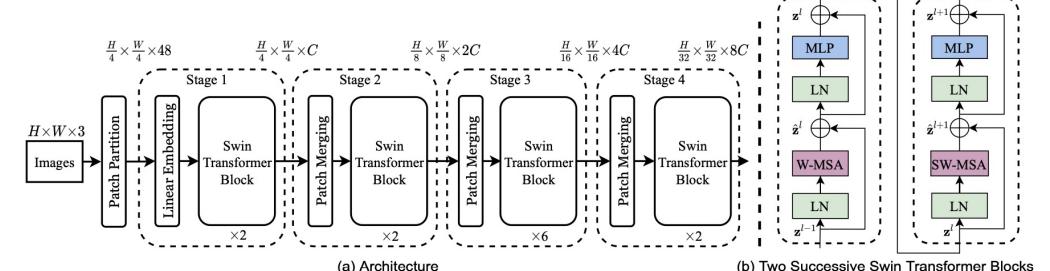
## ■ Transformer



## ■ ViT

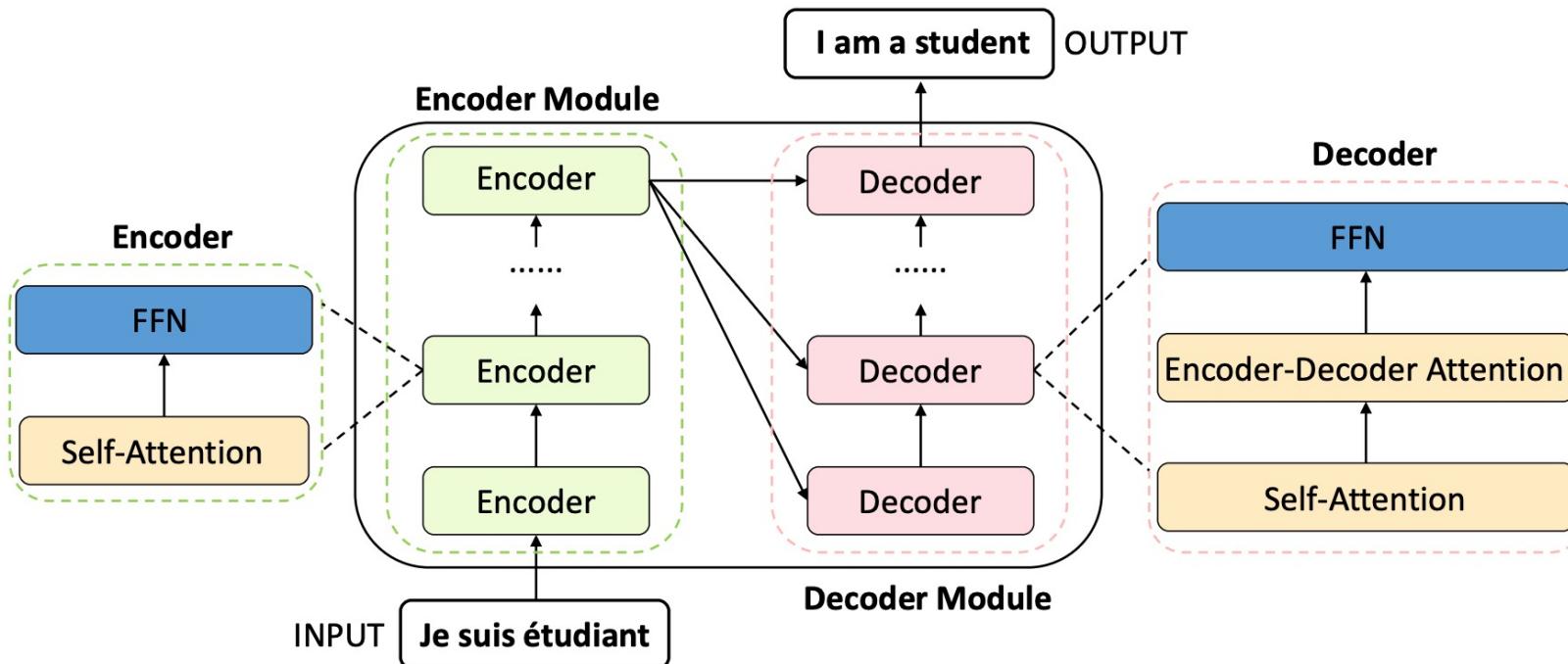


## ■ Swin Transformer



## 2. Architecture Of Transformer

- Transformer was first used in the field of NLP on machine translation tasks.
- Transformer consists if an encoder module and a decoder module.
- Each encoder is composed of a **Multi-head self-attention** layer and a **feed-forward** neural network. **Layernorm** is applied before every block, and **residual connections** after every block.



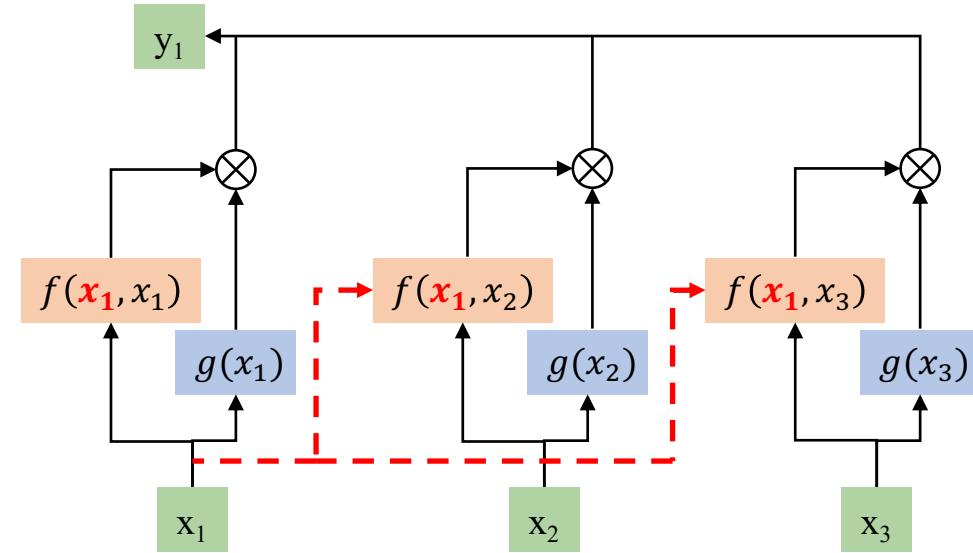
# 2.1 Formulation Of Self-attention

- Self-attention module computes the responses at each position in a sequence by estimating attention scores to all positions and gathering the corresponding embeddings based on the scores accordingly.

Given an input signal  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , where  $n$  is the number of sequences and  $d$  is the number of channels.

$$y_i = \sum_{\forall j} f(x_i, x_j)g(x_j)$$

- Where  $x_i \in \mathbb{R}^{1 \times d}$ ,  $y_i \in \mathbb{R}^{1 \times d}$  indicate the  $i^{\text{th}}$  position of the input signal  $\mathbf{X}$  and output signal  $\mathbf{Y}$ , respectively.
- Subscript  $j$  is the index that enumerates **all positions**.
- $f(\cdot)$  computes a relationship between  $i$  and all  $j$ .
- $g(\cdot)$  computes a representation of the input signal at position  $j$ .



# 2.1 Formulation Of Self-attention

There are many choices for the pairwise function  $\mathbf{f}(\cdot)$ .

For example,  $\mathbf{f}(\cdot)$  can be formulated as:

$$f(x_i, x_j) = \frac{e^{\theta(x_i)\emptyset(x_j)^T}}{\sum_j e^{\theta(x_i)\emptyset(x_j)^T}} = \alpha'_{i,j}$$

- Where  $\theta(\cdot)$  and  $\emptyset(\cdot)$  can be any embedding layers.

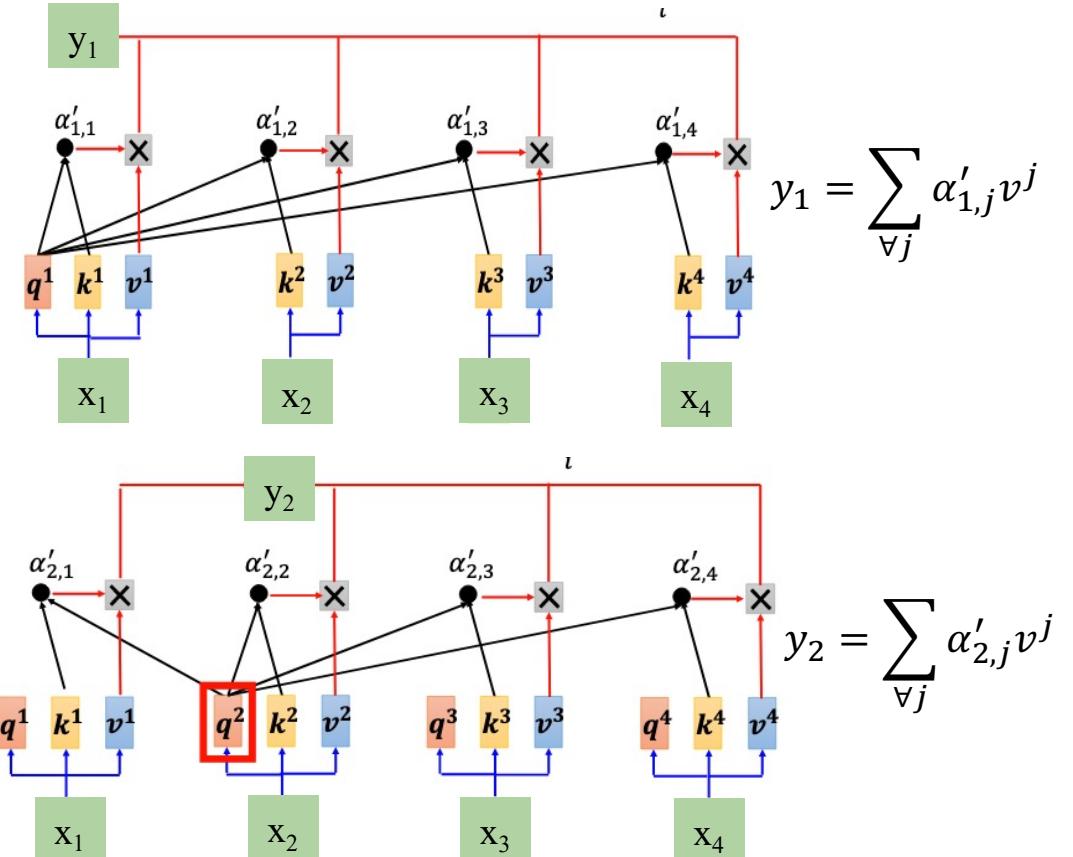
If we consider the  $\theta(\mathbf{X}) = \mathbf{XW}_q$ ,  $\emptyset(\mathbf{X}) = \mathbf{XW}_k$ ,  $\mathbf{g}(\mathbf{X}) = \mathbf{XW}_v$ .  $\mathbf{W}_\theta \in \mathbb{R}^{d \times d_{\text{model}}}$ ,  $\mathbf{W}_\emptyset \in \mathbb{R}^{d \times d_{\text{model}}}$ ,  $\mathbf{W}_g \in \mathbb{R}^{d \times d_{\text{model}}}$ .

$$y_i = \sum_{\forall j} \frac{e^{x_i W_q W_k^T x_j^T}}{\sum_j e^{x_i W_q W_k^T x_j^T}} x_j W_v = \sum_{\forall j} \frac{e^{q_i k_j}}{\sum_j e^{q_i k_j}} v_j$$

$$Y = \text{softmax}(XW_q W_k^T X^T)XW_v$$

- $\mathbf{Q} = \theta(\mathbf{X}) = \mathbf{XW}_q$ ,  $\mathbf{K} = \emptyset(\mathbf{X}) = \mathbf{XW}_k$ ,  $\mathbf{V} = \mathbf{g}(\mathbf{X}) = \mathbf{XW}_v$

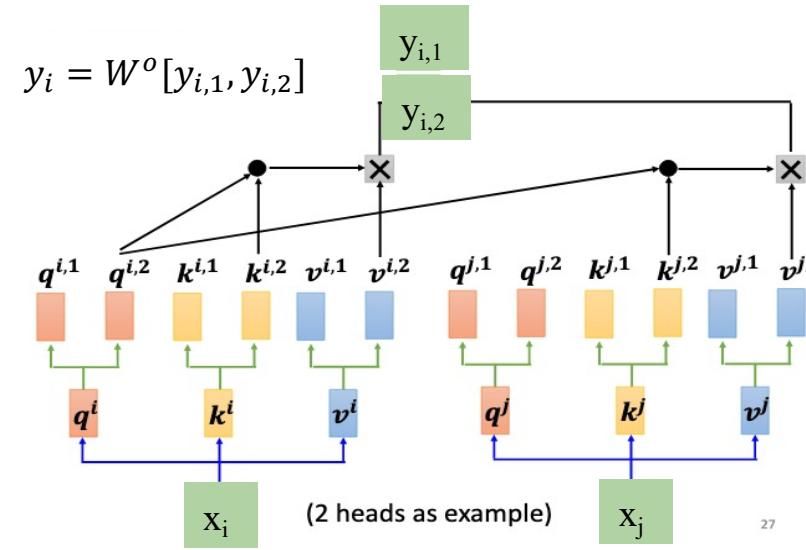
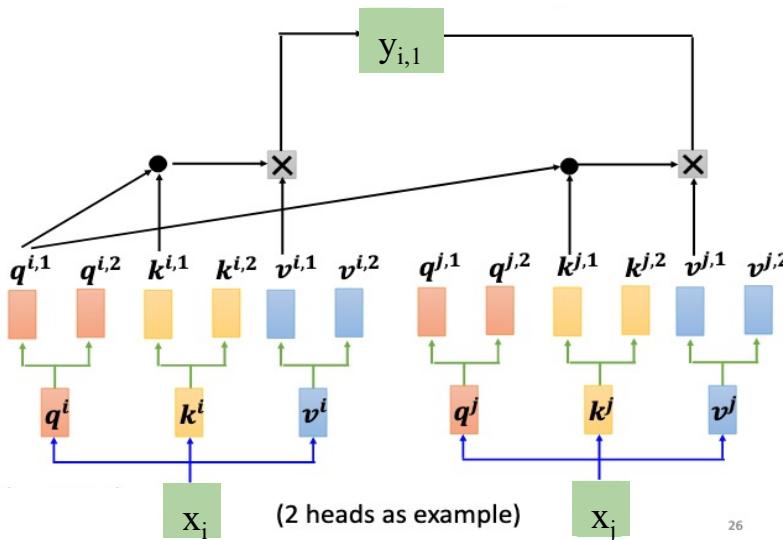
$$Y = \text{softmax}(QK^T)V$$



## 2.2 Formulation Of Multi-head Attention

- Given an input vector  $\mathbf{x}_i$  and the number of heads  $\mathbf{h}$ .
- The input vector  $\mathbf{x}_i$  is first transformed into three different groups of vector: query group  $\mathbf{q}^i$ , key group  $\mathbf{k}^i$ , value group  $\mathbf{v}^i$ .
- In each group, there are  $\mathbf{h}$  vectors with dimension  $\mathbf{d}_{\mathbf{q}'} = \mathbf{d}_{\mathbf{k}'} = \mathbf{d}_{\mathbf{v}'} = \mathbf{d}_{\text{model}}/\mathbf{h}$ .
- The vectors derived from different inputs and packed together into three groups of matrices:  $\{\mathbf{Q}_i\}_{i=1}^h$ ,  $\{\mathbf{K}_i\}_{i=1}^h$  and  $\{\mathbf{V}_i\}_{i=1}^h$ .

$$\text{MultiHead}(Q', K', V') = \text{Concat}(\text{head}_1, \dots, \text{head}_n)W^o, \text{ where } \text{head}_i = \text{Attention}(Q_i, K_i, V_i)$$



# 2.3 Other Key Concepts

## I. Positional Encoding

- Representing the **position** of each input sequence as a **vector**.
- Position vector  $\mathbf{X}_{pe}$  has the same dimension as input sequence  $\mathbf{X}$ .
- The input of transformers will be the sum of  $\mathbf{X}_{pe}$  and  $\mathbf{X}$ .

$$\mathbf{X}' = \mathbf{X} + \mathbf{X}_{pe}, \quad \mathbf{X} \in R^{B \times n \times d}, \mathbf{X}_{pe} \in R^{B \times n \times d}$$

- Where **B**: batch size, **n**: length of sequence, **d**: dimension of each vector.

### ① Sinusoidal Position Encoding

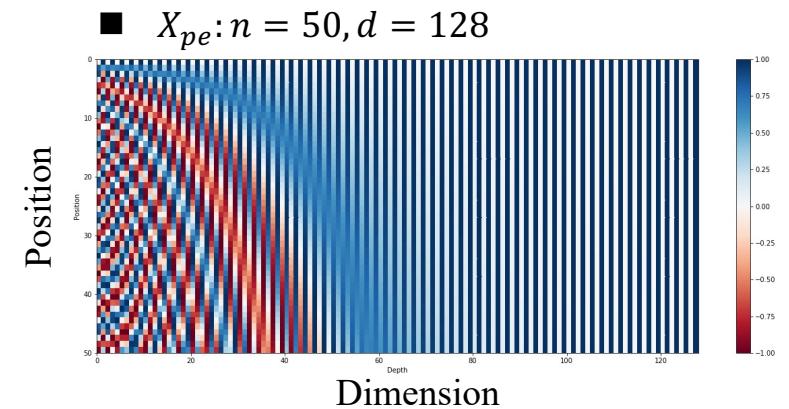
**pos**: the position in an input sequence, **i**: the dimension of a vector

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

### ② Learned Positional Embedding

Random initialization  $X_{pe}$  and  $\mathbf{X}_{pe}$  can be learned.



# 2.3 Other Key Concepts

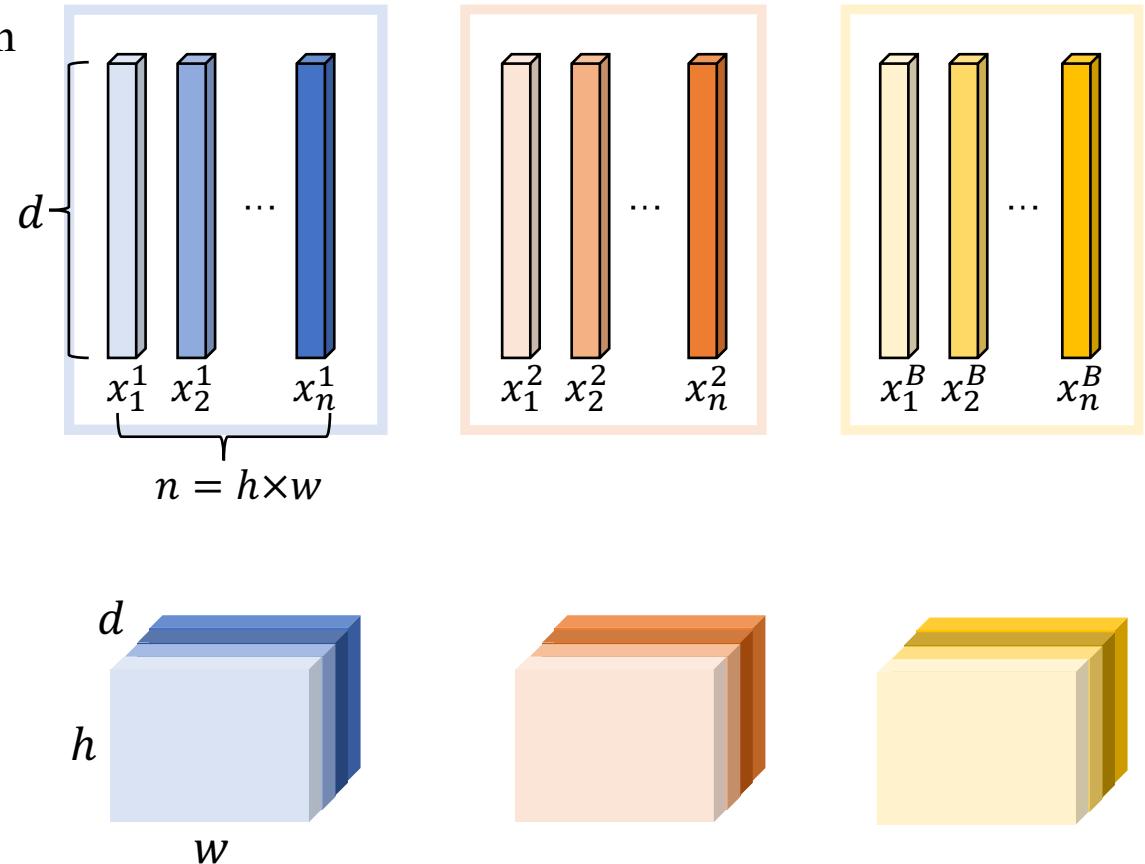
## II. Layer Norm

- Let **B**: batch size, **n**: number of sequence, **d**: dimension of vector.

$$\mu^b = \frac{1}{H} \sum_{i=1}^n \sum_{j=1}^d x_{i,j}^b, \quad H = n \times d$$

$$\sigma^b = \sqrt{\frac{1}{H} \sum_{i=1}^n \sum_{j=1}^d (x_{i,j}^b - \mu^b)^2} \quad \hat{x}_{i,j}^b = \frac{x_{i,j}^b - \mu^b}{\sqrt{(\sigma^b)^2 + \epsilon}}$$

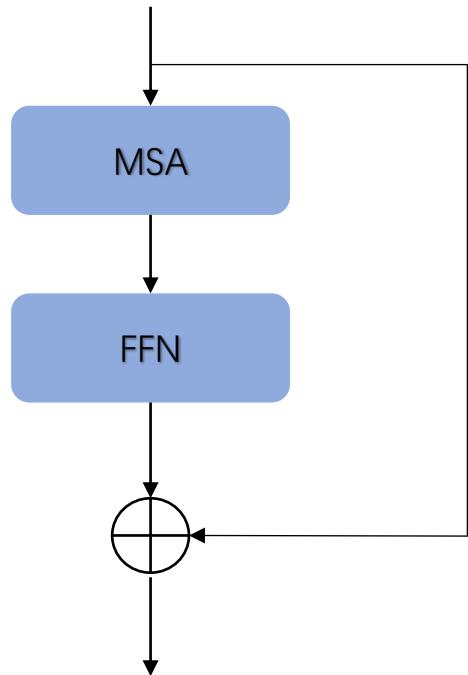
- $\mathbf{x}_i \in \mathbf{R}^{1 \times d}$  indicates the  $i^{\text{th}}$  position of the input sequence  $\mathbf{X}$ .
- $\mathbf{b}$  indicates the  $b^{\text{th}}$  sample in a batch.
- $\mathbf{j}$  indicates the  $j^{\text{th}}$  element in a vector.



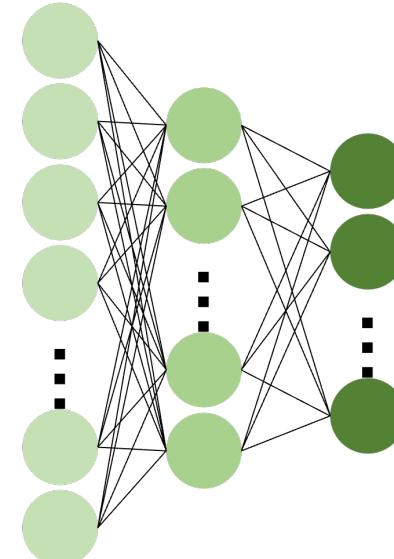
# 2.3 Other Key Concepts

## III. Residual Connection

$$X' = X + f(X)$$



## IV. FFN(Feed Forward Network) MLP(Multi layer Perception)

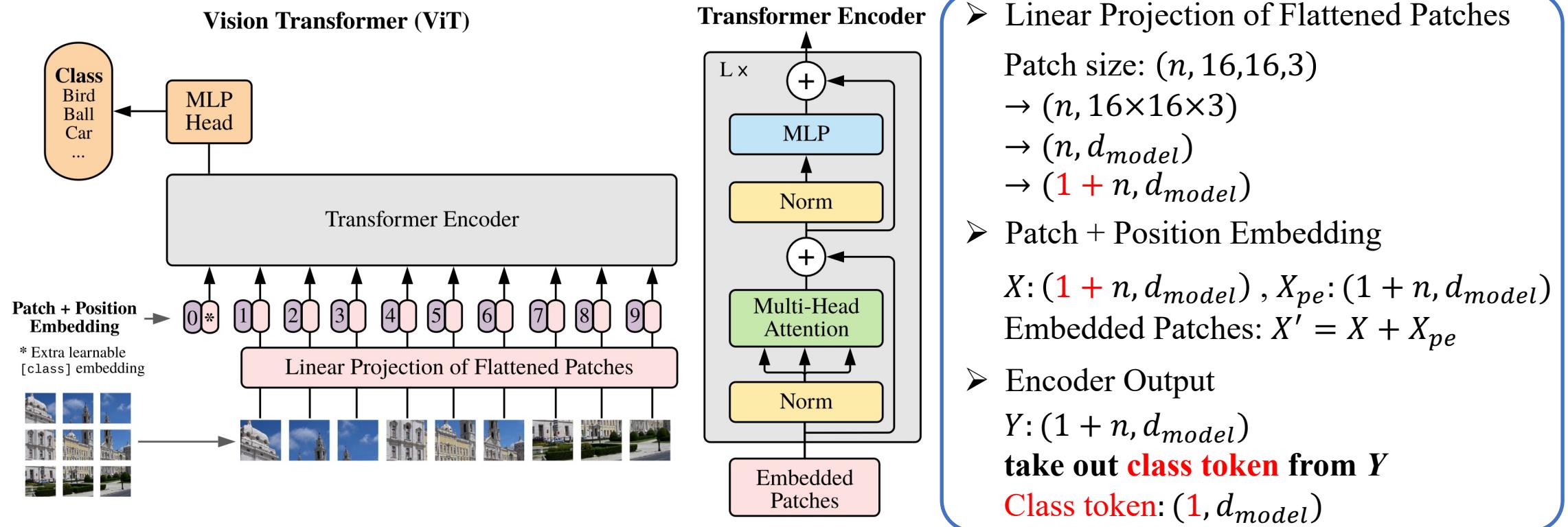


# AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

CVPR 2020

# 3.1 Architecture Of ViT

- ViT is a **pure transformer** that performs well on image classification task when applied directly to the sequences of image patches.
- ViT follows the original Transformer, but only contains encoder layers.
- Each encoder is composed of a **Multi-head self-attention** layer and a **feed-forward** neural network. **Layernorm** is applied before every block, and **residual connections** after every block.



### 3.1 ViT Experiments

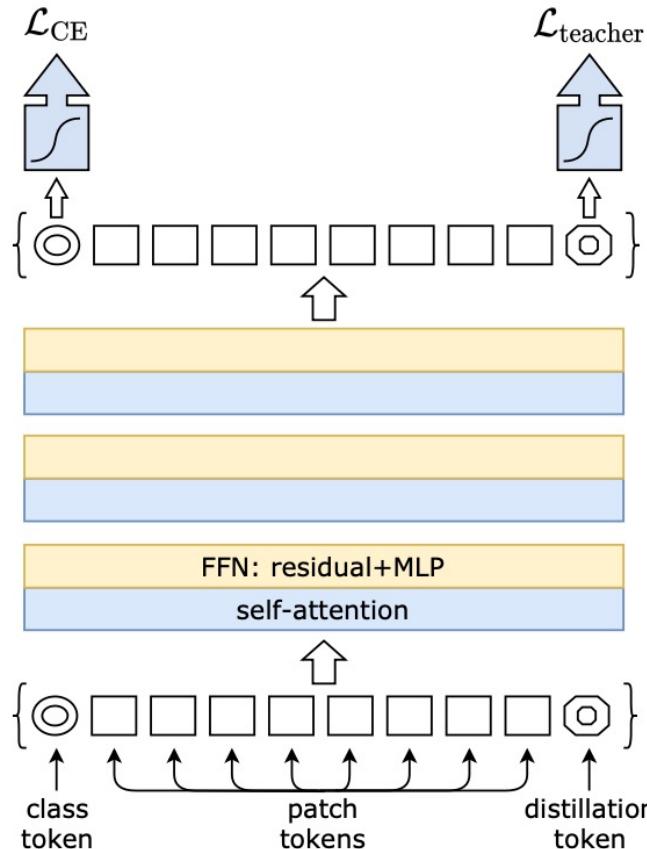
	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	<b>88.55</b> ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4/88.5*
ImageNet ReaL	<b>90.72</b> ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	<b>99.50</b> ± 0.06	99.42 ±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

Training **data-efficient** image transformers & distillation through  
attention

PMLR 2021

## 3.2 Architecture Of DeiT

- High-performing vision transformers are pre-trained with hundreds of millions of images using a **large infrastructure**, thereby limiting their adoption.
- DeiT introduces a **new distillation procedure** based on a distillation token, and the model pre-trained on Imagenet are competitive when transferred to downstream tasks.



Let  $\mathbf{Z}_t$  be the logits of the teacher model,  $\mathbf{Z}_s$  the logits of the student model. We denote by  $\tau$  the temperature for the distillation,  $\lambda$  the coefficient balancing the Kullback–Leibler divergence loss (**KL**) and the cross-entropy ( **$L_{CE}$** ) on ground truth labels  $y$ , and  $\psi$  the softmax function.

➤ Soft distillation

$$\mathcal{L}_{\text{global}} = (1 - \lambda)\mathcal{L}_{CE}(\psi(Z_s), y) + \lambda\tau^2\text{KL}(\psi(Z_s/\tau), \psi(Z_t/\tau)).$$

➤ Hard distillation ——  $y_t = \text{argmax}_c Z_t(c)$

$$\mathcal{L}_{\text{global}}^{\text{hardDistill}} = \frac{1}{2}\mathcal{L}_{CE}(\psi(Z_s), y) + \frac{1}{2}\mathcal{L}_{CE}(\psi(Z_s), y_t).$$

➤ Hard distillation + Label smoothing

$$\hat{y}_t(i) = \begin{cases} 1 - \epsilon & i = \text{target} \\ \epsilon/(K - 1) & i \neq \text{target} \end{cases}, \epsilon = 0.1$$

0	0.025
1	0.9
0	0.025
0	0.025
0	0.025

## 3.2 Training/Test procedure Of DeiT

- **Training**

Teacher Model is RegNetY-16GF and let its output :  $y$ . In training process, the model has two output: Output of class token:  $y_{ct}$  and output of distillation token:  $y_{dt}$ .

- Training Loss:  $L_{total} = \alpha L_{base} + (1 - \alpha) L_{dist}$
- If distillation procedure uses Soft-distillation :  $L_{total} = \alpha L_{CE}(y_{ct}, y) + (1 - \alpha) L_{KL}(y_{dt}, y_{teacher})$
- If distillation procedure uses Hard-distillation:  $L_{total} = \alpha L_{CE}(y_{ct}, y) + (1 - \alpha) L_{CE}(y_{dt}, \text{argmax}(y_{teacher}))$
- If using Label-smoothing Trick:  $L_{base}(\cdot) = L_{CE}(y_{ct}, \text{label\_smoothing}(y))$

$$\text{And : } L_{total} = \alpha L_{CE}(y_{ct}, \text{label\_smoothing}(y)) + (1 - \alpha) L_{dist}$$

- **Testing**

$$\text{Model Output: } y = \frac{1}{2} y_{ct} + \frac{1}{2} y_{dt}$$

## 3.2 DeiT Experiments

DeiT: method ↓	supervision		ImageNet top-1 (%)			
	label	teacher	Ti 224	S 224	B 224	B↑384
no distillation	✓	✗	72.2	79.8	81.8	83.1
usual distillation	✗	soft	72.2	79.8	81.8	83.2
hard distillation	✗	hard	74.3	80.9	83.0	84.0
class embedding	✓	hard	73.9	80.9	83.0	84.2
distil. embedding	✓	hard	74.6	81.1	83.1	84.4
DeiT $\natural$ : class+distil.	✓	hard	74.5	81.2	83.4	84.5

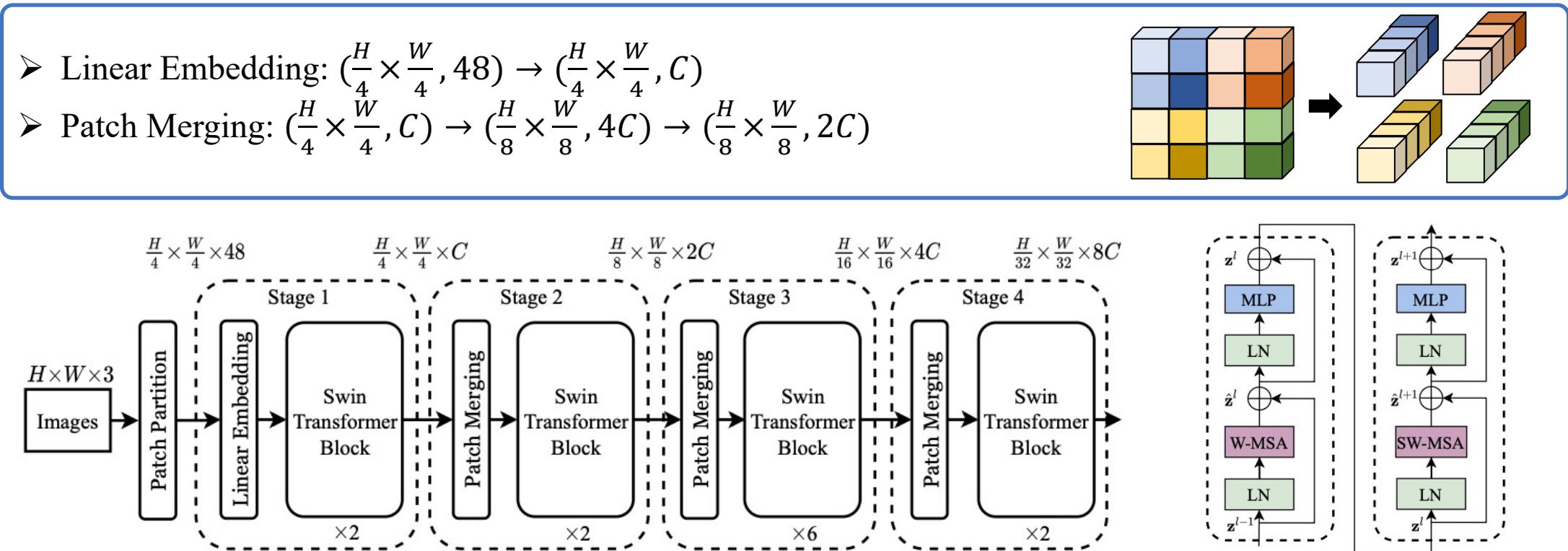
Model	ImageNet	CIFAR-10		CIFAR-100		Flowers	Cars	iNat-18	iNat-19	im/sec
		CIFAR-10	CIFAR-100	CIFAR-100	CIFAR-100					
EfficientNet-B7	84.3	98.9	91.7	98.8	94.7	-	-	-	-	55.1
ViT-B/32	73.4	97.8	86.3	85.4	-	-	-	-	-	394.5
ViT-B/16	77.9	98.1	87.1	89.5	-	-	-	-	-	85.9
ViT-L/32	71.2	97.9	87.1	86.4	-	-	-	-	-	124.1
ViT-L/16	76.5	97.9	86.4	89.7	-	-	-	-	-	27.3
DeiT-B	81.8	99.1	90.8	98.4	92.1	73.2	77.7	77.7	77.7	292.3
DeiT-B↑384	83.1	99.1	90.8	98.5	93.3	79.5	81.4	81.4	81.4	85.9
DeiT-B $\natural$	83.4	99.1	91.3	98.8	92.9	73.7	78.4	78.4	78.4	290.9
DeiT-B $\natural$ ↑384	84.4	99.2	91.4	98.9	93.9	80.1	83.0	83.0	83.0	85.9

# Swin Transformer: **Hierarchical** Vision Transformer using Shifted Windows

ICCV 2021 Best Paper

### 3.3 Architecture Of Swin Transformer

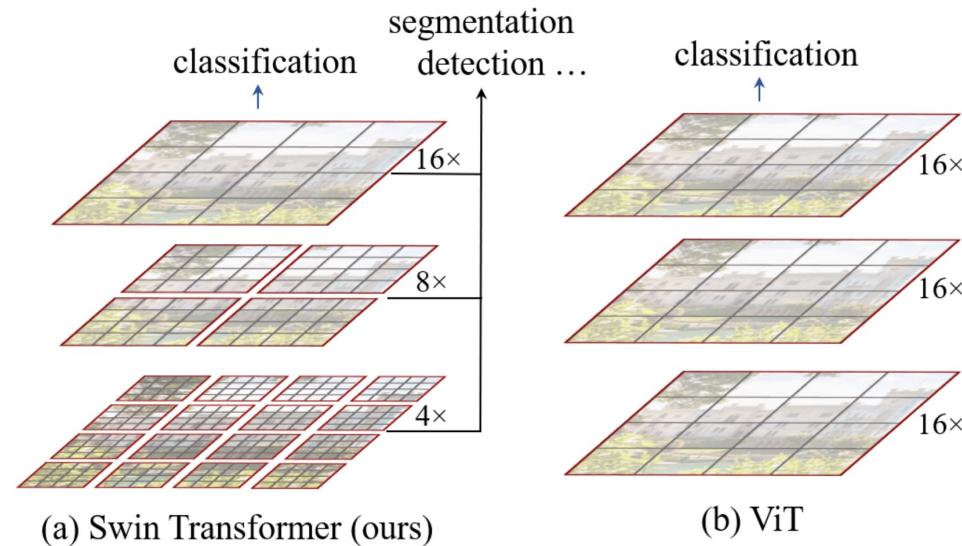
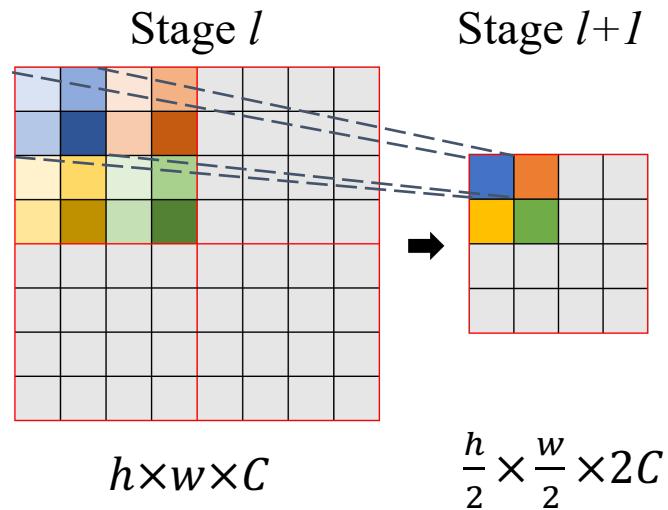
- Swin Transformer, that capably serves as a general-purpose **backbone** for computer vision, has a **hierarchical** architecture.
- In Swin Transformer , representation is computed with **Shifted windows**.



- **W-MSA:** Regular Window Multi-head Self-attention
- **SW-MSA:** Shifted Window Multi-head Self-attention

### 3.3 Why Hierarchical?

- The architecture of ViT is **unsuitable** for use as a general-purpose backbone network on **dense** vision tasks when input image **resolution is high**.
- As the structure deepens, the **receptive field** of patch in each layer gradually **increases**, eg., VGG and ResNet.

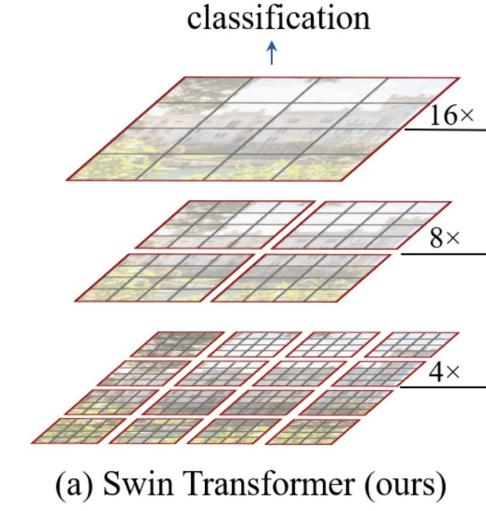
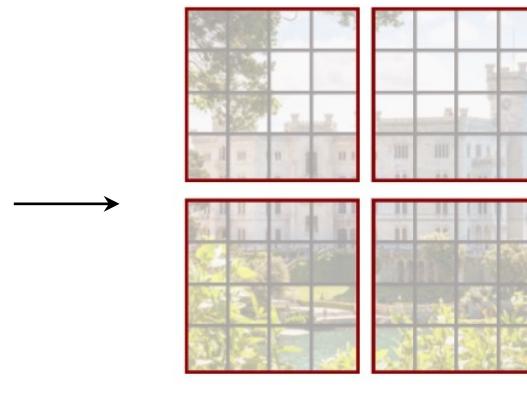
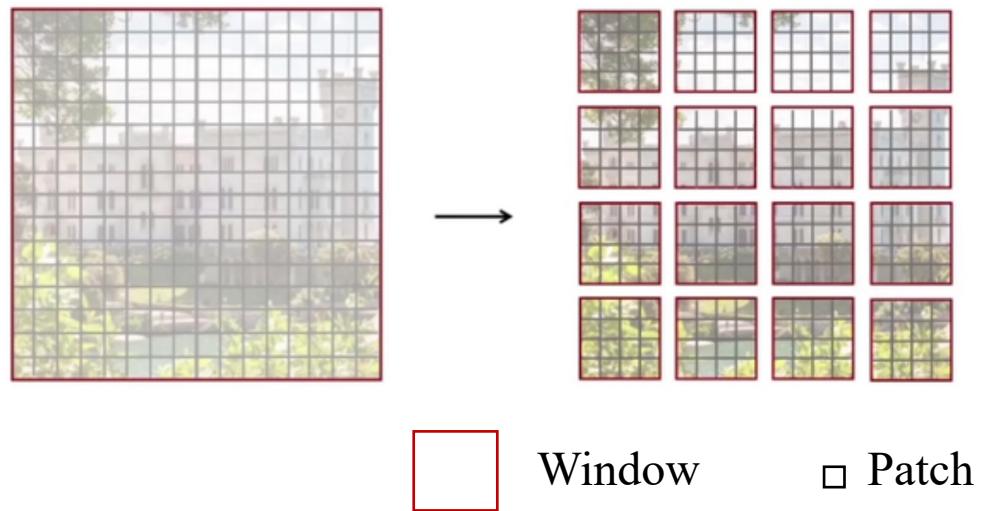


### 3.3 Window Multihead Self-attention

- Windows are arranged to evenly partition the image in a non-overlapping manner.
- The number of patches in each window is fixed, each window contains  $M \times M$  patches.
- The computational complexity of global self-attention is quadratic to image size . Different from global self-attention, **window is the unit for self-attention** . So, it has linear complexity to image size.
- Supposing a global MSA module and a window based one on an image of  $h \times w$  patches are:

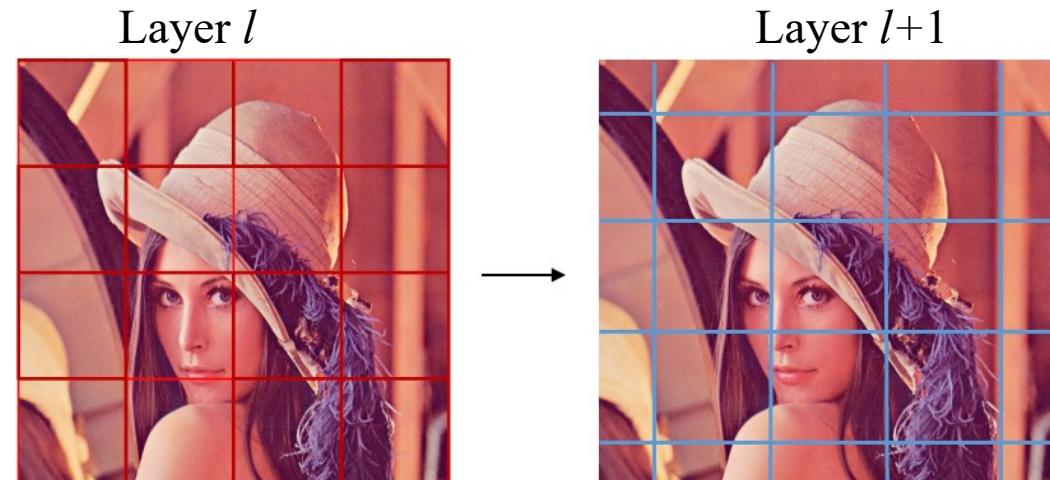
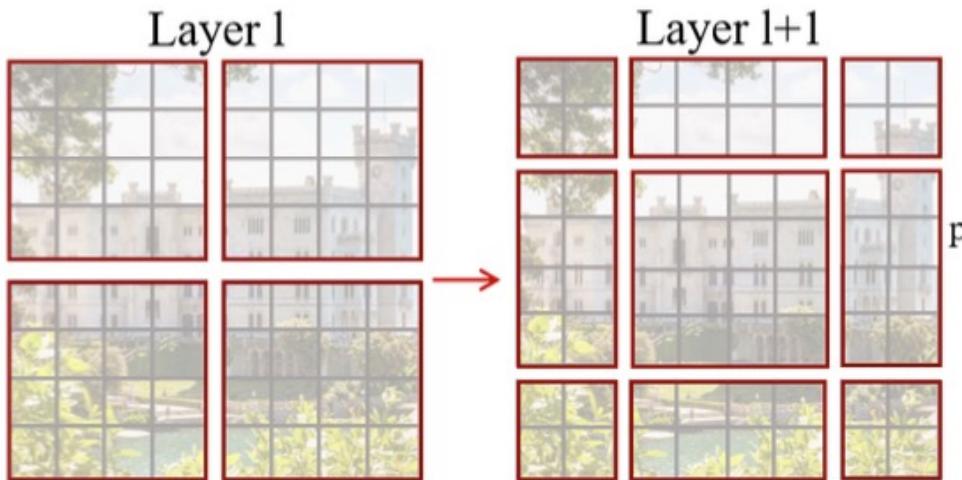
$$\Omega(\text{MSA}) = 4hwC^2 + 2(hw)^2C, \quad (1)$$

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

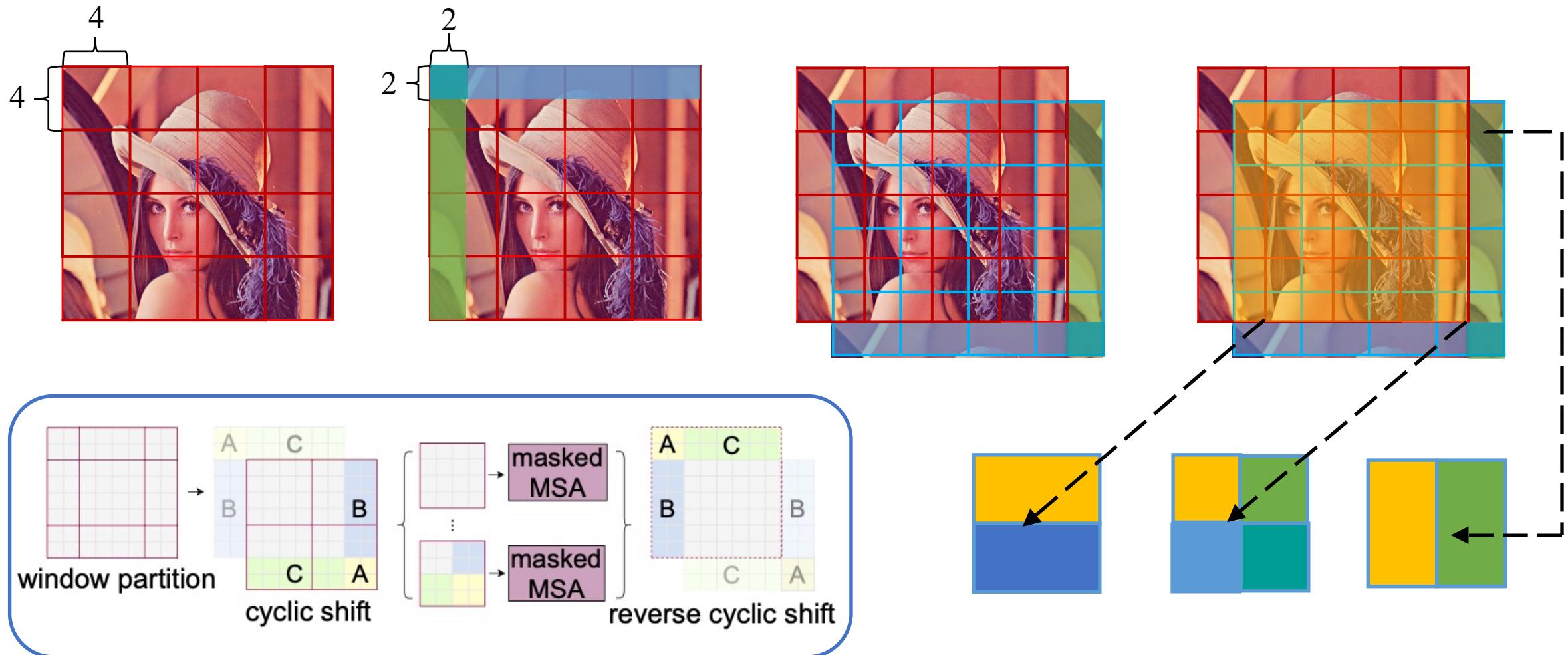


### 3.3 Shifted Window Multihead Self-attention

- The self-attention computation in the new windows **crosses the boundaries of the previous windows** in layer  $l$ , providing connections among them.



### 3.3 Shifted Window Multihead Self-attention



# 3.3 Experiment Of SwinT

Classification on ImageNet-1K

(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

Object detection on COCO

(a) Various frameworks							
Method	Backbone	AP <sup>box</sup>	AP <sub>50</sub> <sup>box</sup>	AP <sub>75</sub> <sup>box</sup>	#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
	R-50	43.5	61.9	47.0	32M	205G	28.3
ATSS	Swin-T	<b>47.2</b>	<b>66.5</b>	<b>51.3</b>	36M	215G	22.3
	R-50	46.5	64.6	50.3	42M	274G	13.6
RepPointsV2	Swin-T	<b>50.0</b>	<b>68.5</b>	<b>54.2</b>	45M	283G	12.0
	R-50	44.5	63.4	48.2	106M	166G	21.0
Sparse 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 <sub>50</sub> <sup>box</sup>	AP <sub>75</sub> <sup>box</sup>	AP <sup>mask</sup>	AP <sub>50</sub> <sup>mask</sup>	AP <sub>75</sub> <sup>mask</sup>	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

### 3.3 Experiment Of SwinT



	ImageNet top-1	ImageNet top-5	COCO AP <sup>box</sup>	COCO AP <sup>mask</sup>	ADE20k 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>

**END**

10.30