

# Transformer 系列模型

Presenter: 王浩

Advisor: 林嘉文

Date: 2023/07/26

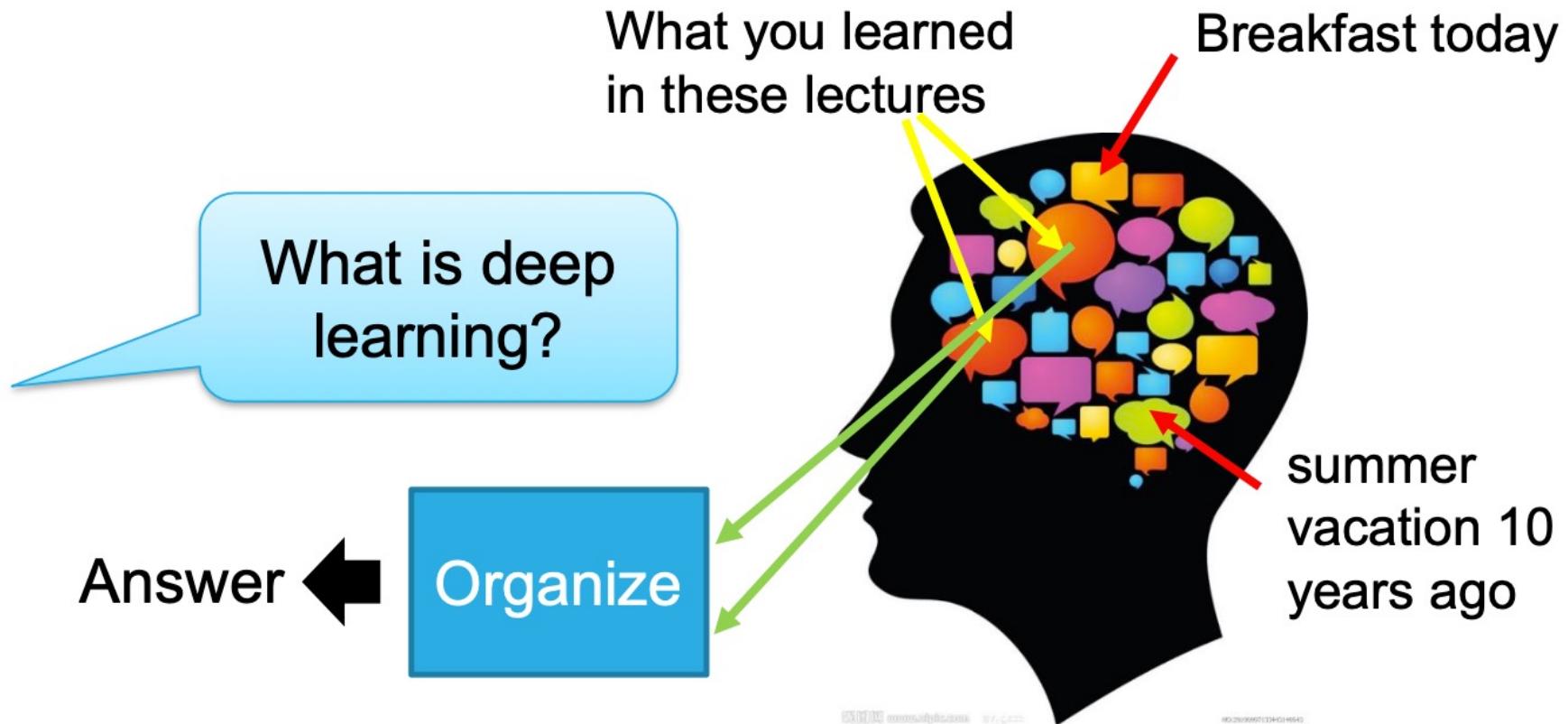
# Outlines

- Attention Mechanism
- Transformer
- Vision Transformer
- Swin Transformer

# Outlines

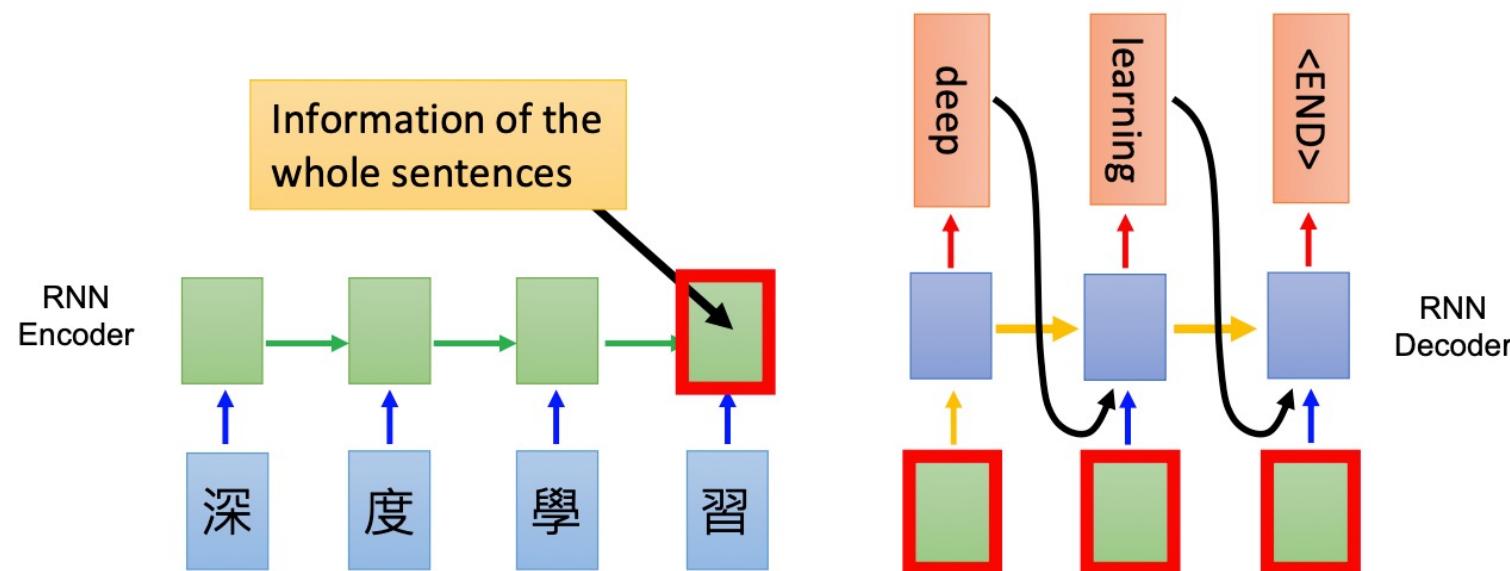
- Attention Mechanism
- Transformer
- Vision Transformer
- Swin Transformer

# Attention Mechanism

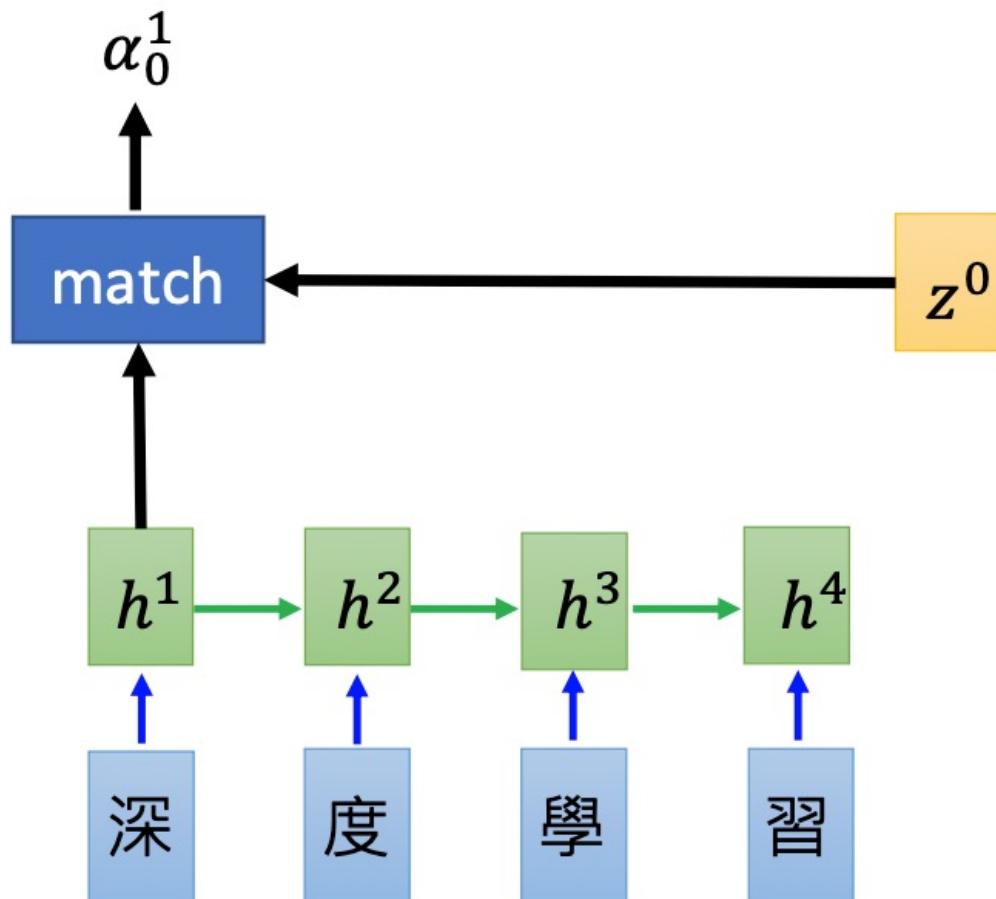


# Machine Translation

- Sequence-to-sequence learning: both input and output are both sequences with different lengths.
- E.g. 深度學習 → deep learning



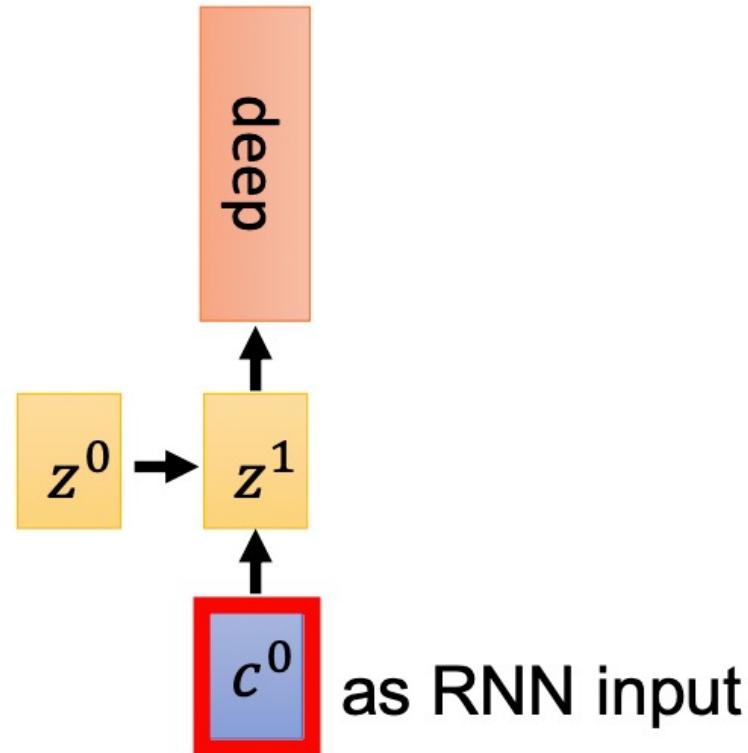
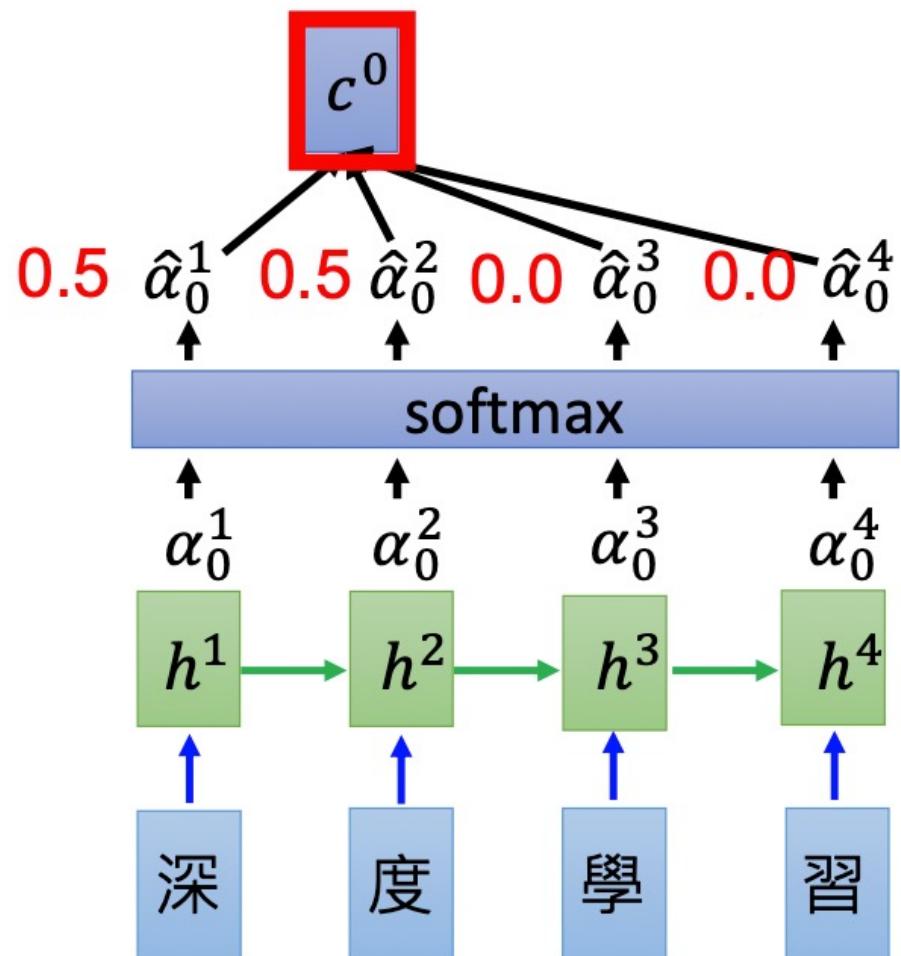
# Machine Translation with Attention



What is **match** ?

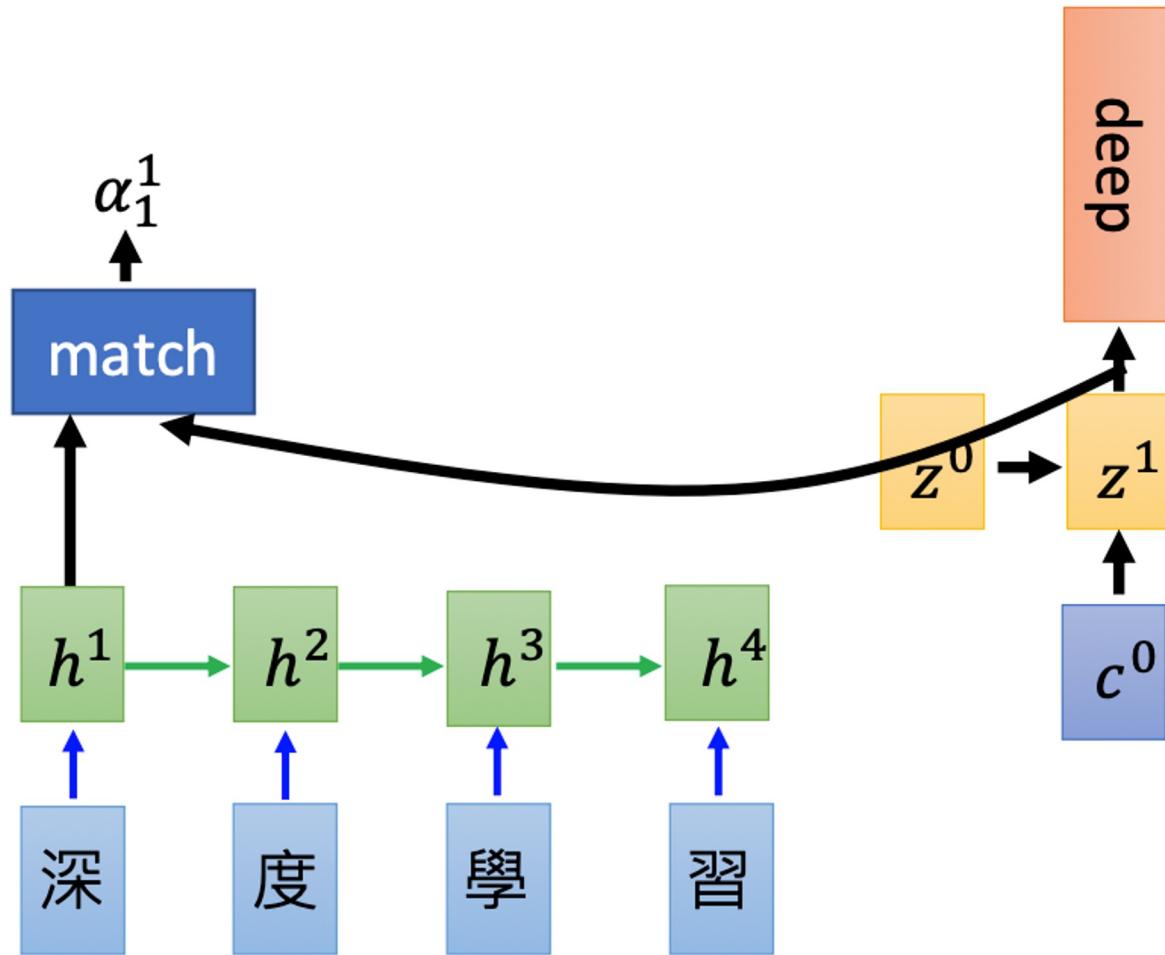
- Cosine similarity of  $z$  and  $h$
- Small NN whose input is  $z$  and  $h$ , output a scalar
- $\alpha = h^T W z$

# Machine Translation with Attention

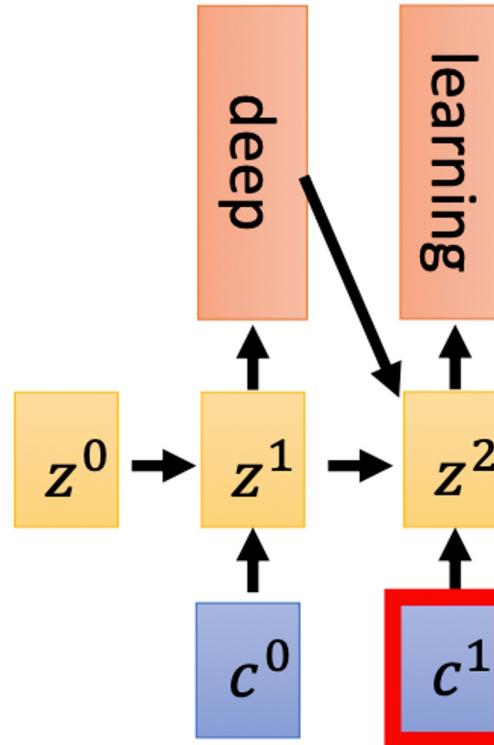
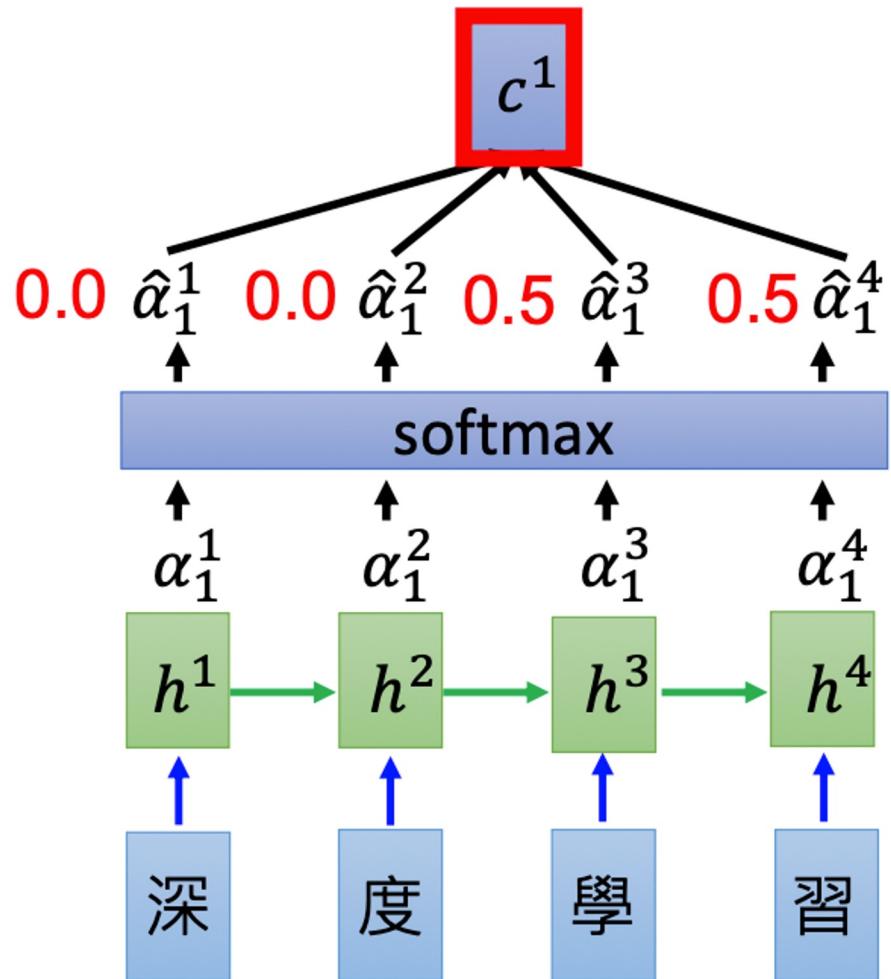


$$c^0 = \sum \hat{\alpha}_0^i h^i = 0.5h^1 + 0.5h^2$$

# Machine Translation with Attention

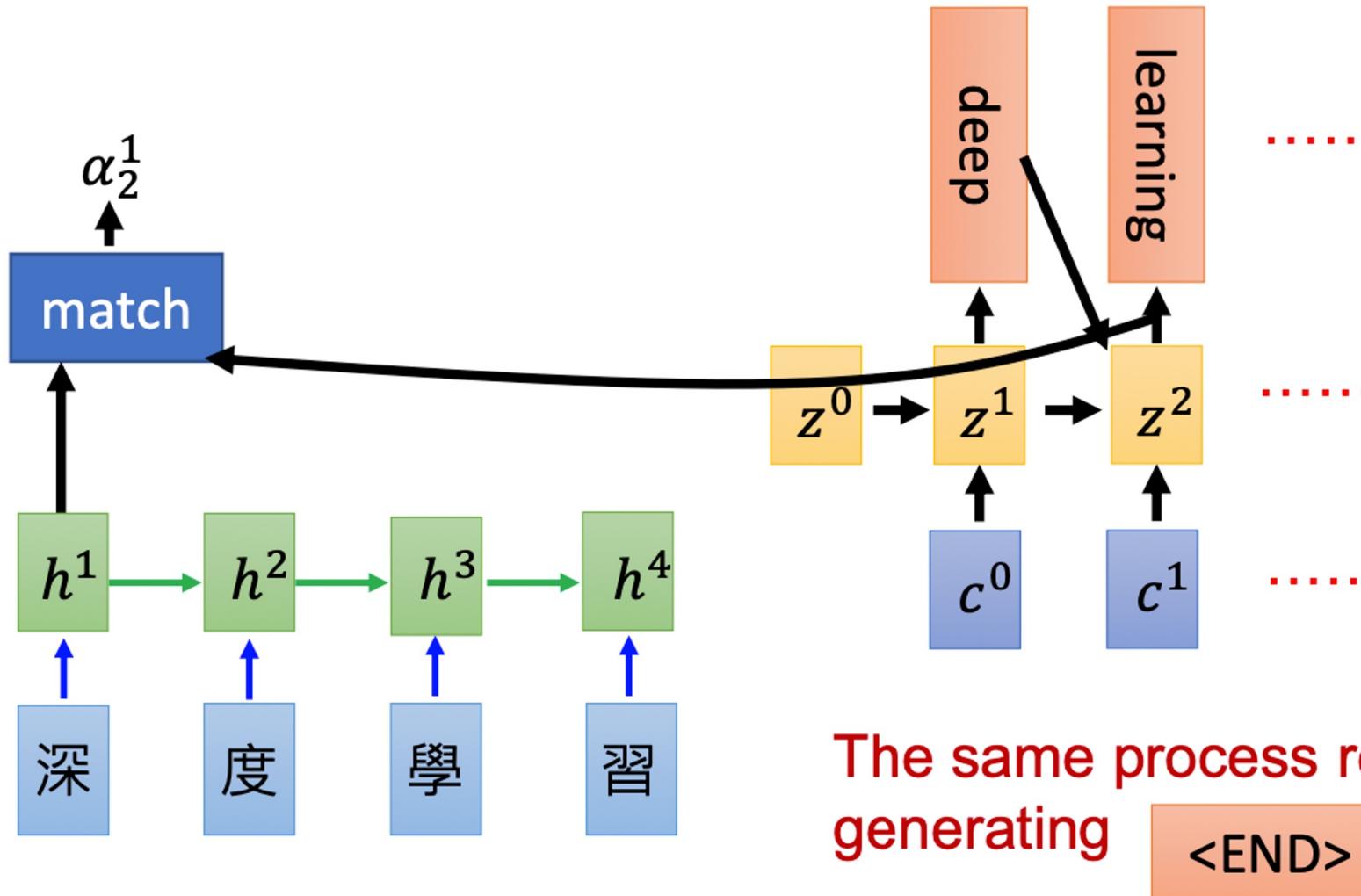


# Machine Translation with Attention

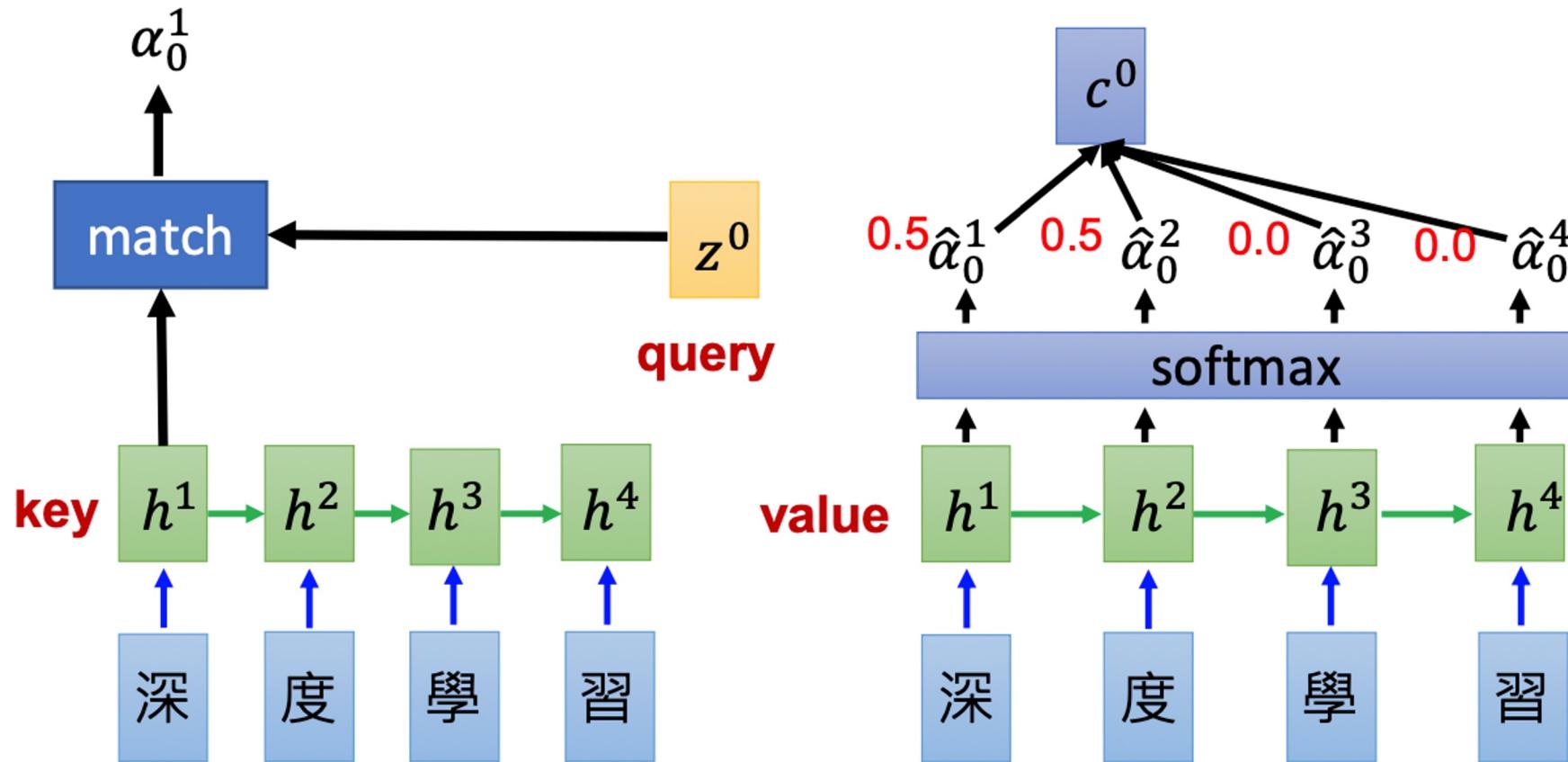


$$c^1 = \sum \hat{\alpha}_1^i h^i = 0.5h^3 + 0.5h^4$$

# Machine Translation with Attention



# Machine Translation with Attention



# Dot-Product Attention

- Input: a query  $q$  and a set of key-value ( $k$ - $v$ ) pairs to an output
- Output: weighted sum of values

Inner product of  
query and corresponding key

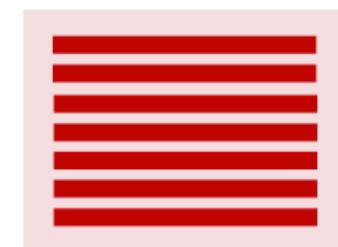
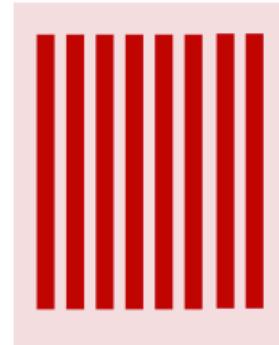
$$A(q, K, V) = \sum_i \frac{\exp(q \cdot k_i)}{\sum_j \exp(q \cdot k_j)} v_i$$

# Dot-Product Attention in Matrix

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

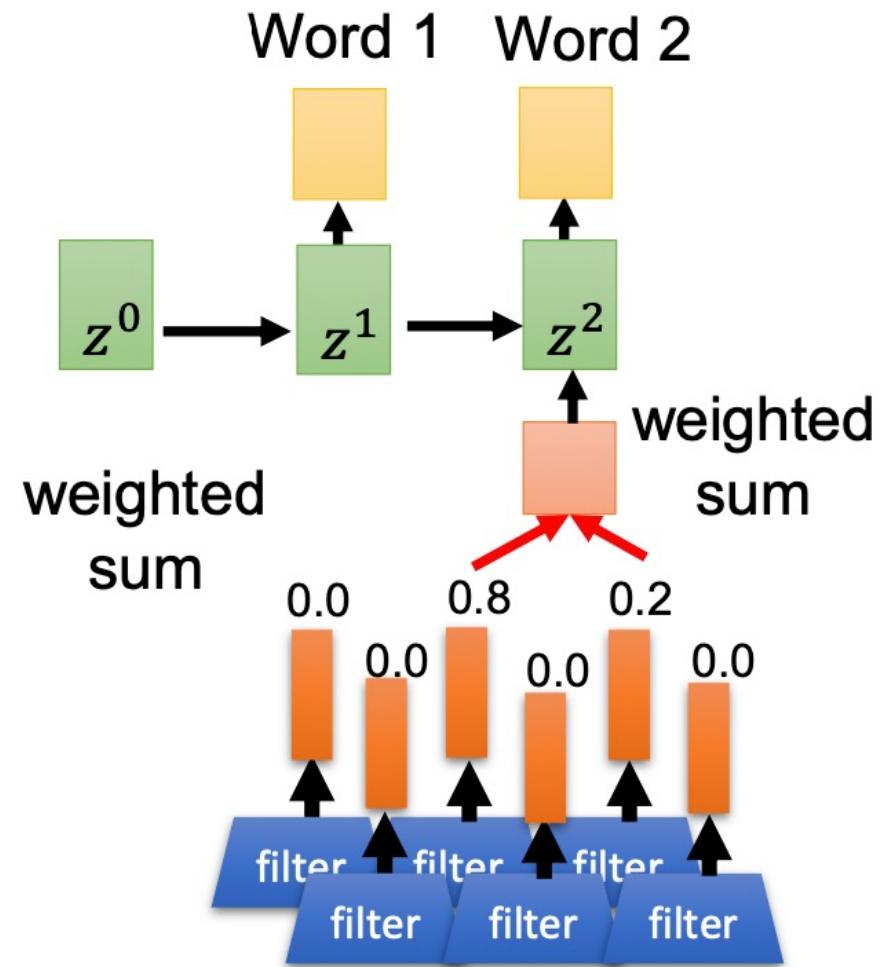
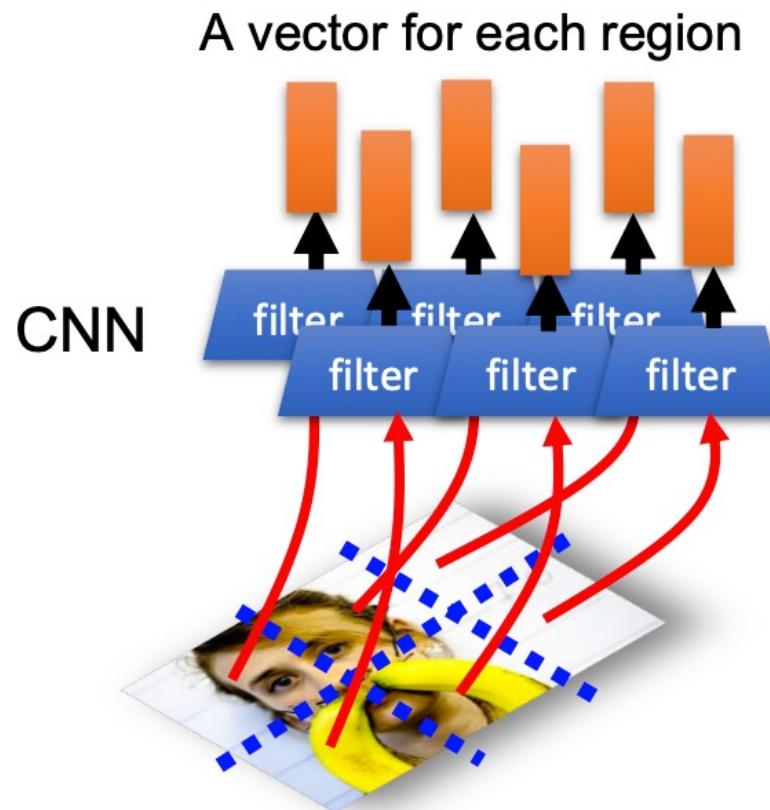
$$[|Q| \times d_k] \times [d_k \times |K|] \times [|K| \times d_v]$$

softmax  
row-wise



$$= [|Q| \times d_v]$$

# Image Captioning with Attention



# Image Captioning

- Good examples



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



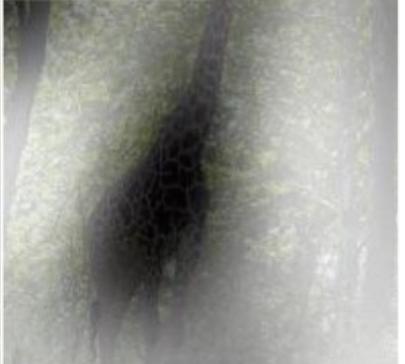
A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

# Image Captioning

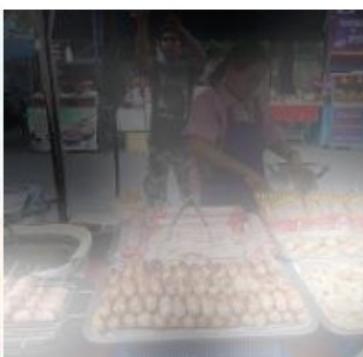
- Bad examples



A large white bird standing in a forest.

A woman holding a clock in her hand.

A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.

A woman is sitting at a table with a large pizza.

A man is talking on his cell phone while another man watches.

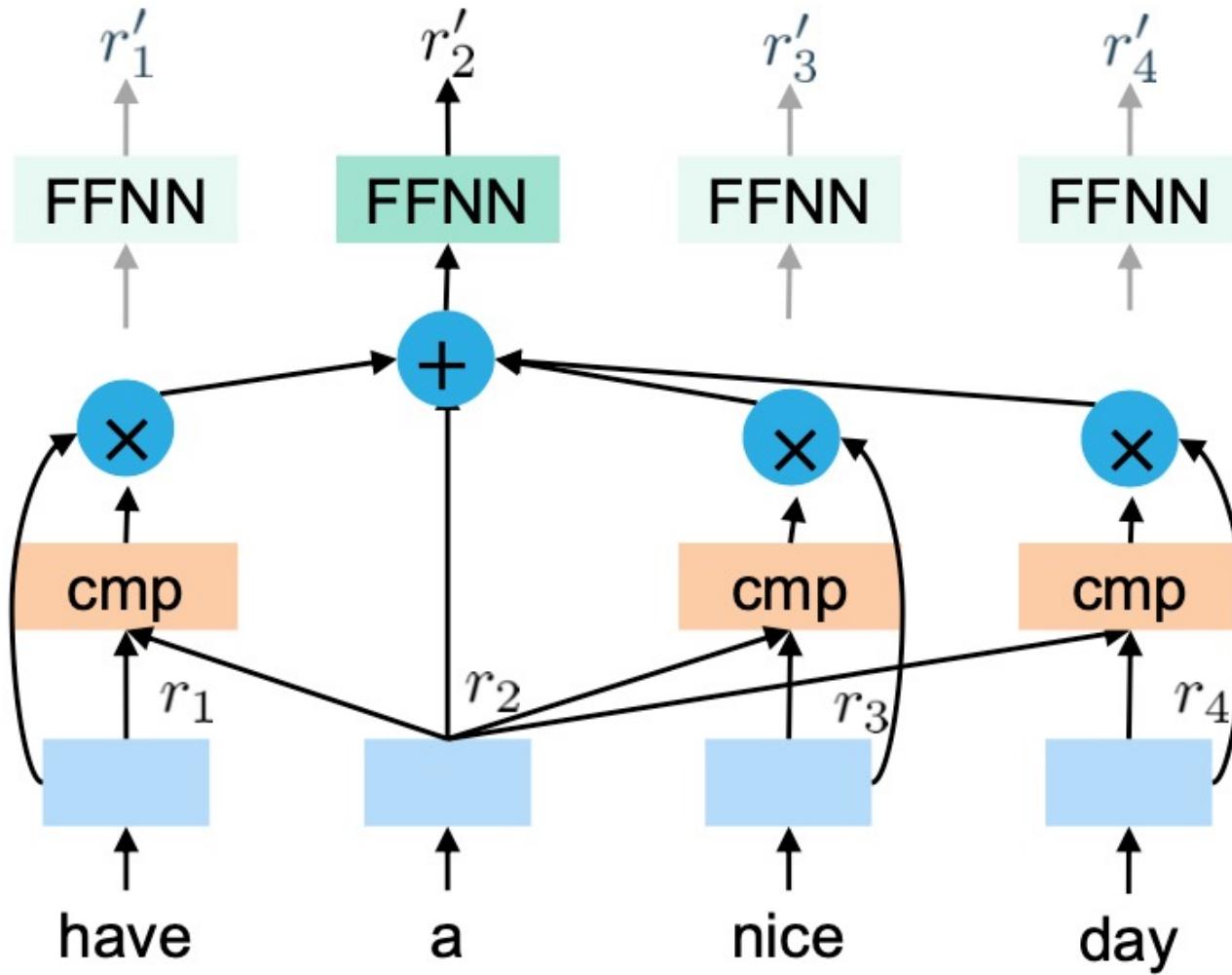
# Outlines

- Attention Mechanism
- Transformer
- Vision Transformer
- Swin Transformer

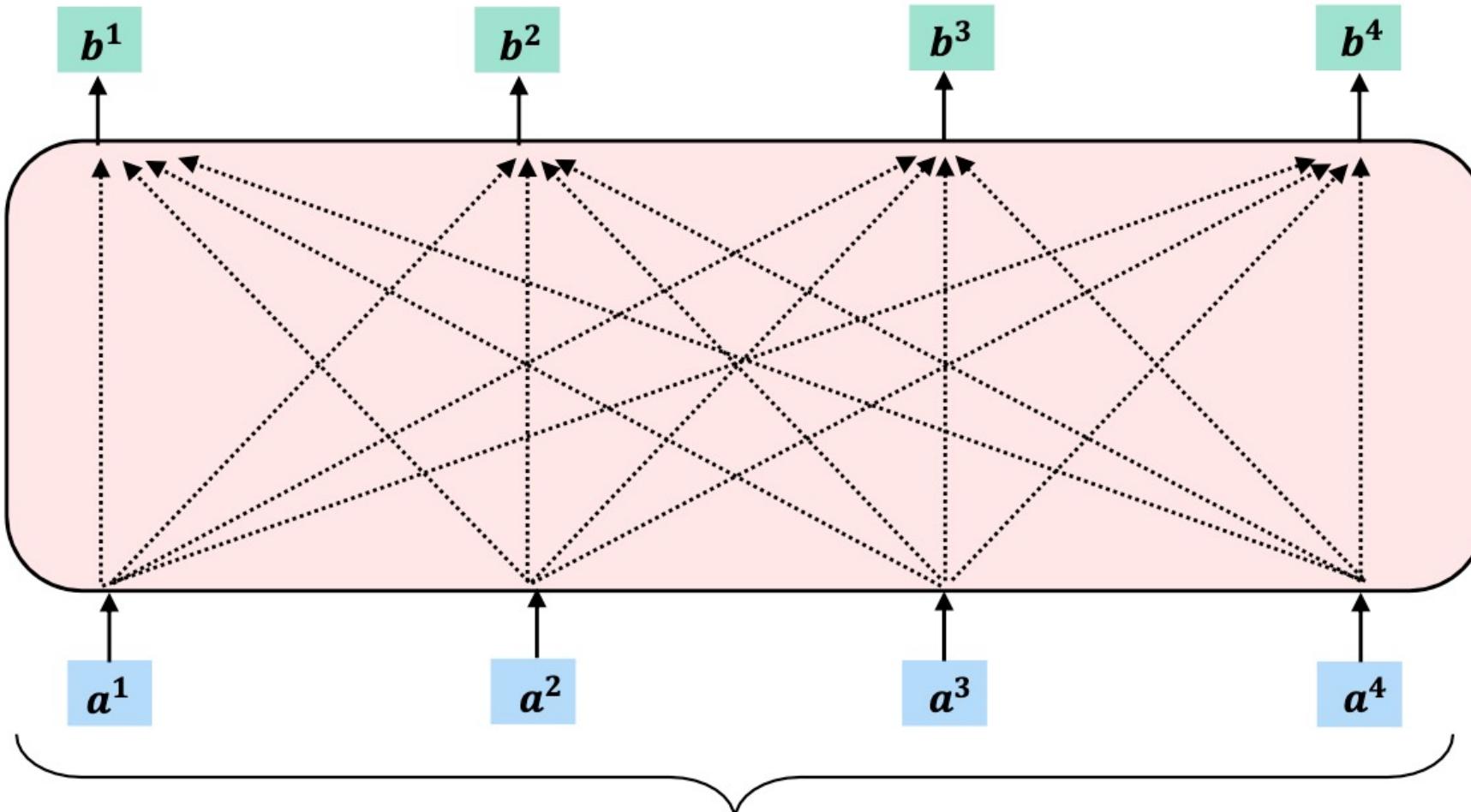
# Comparison

- RNN
  - No explicit modeling of long and short range dependencies
  - LSTM?
- CNN
  - Long-distance dependencies require many layers
- Attention allows us to access any state

# Self-Attention



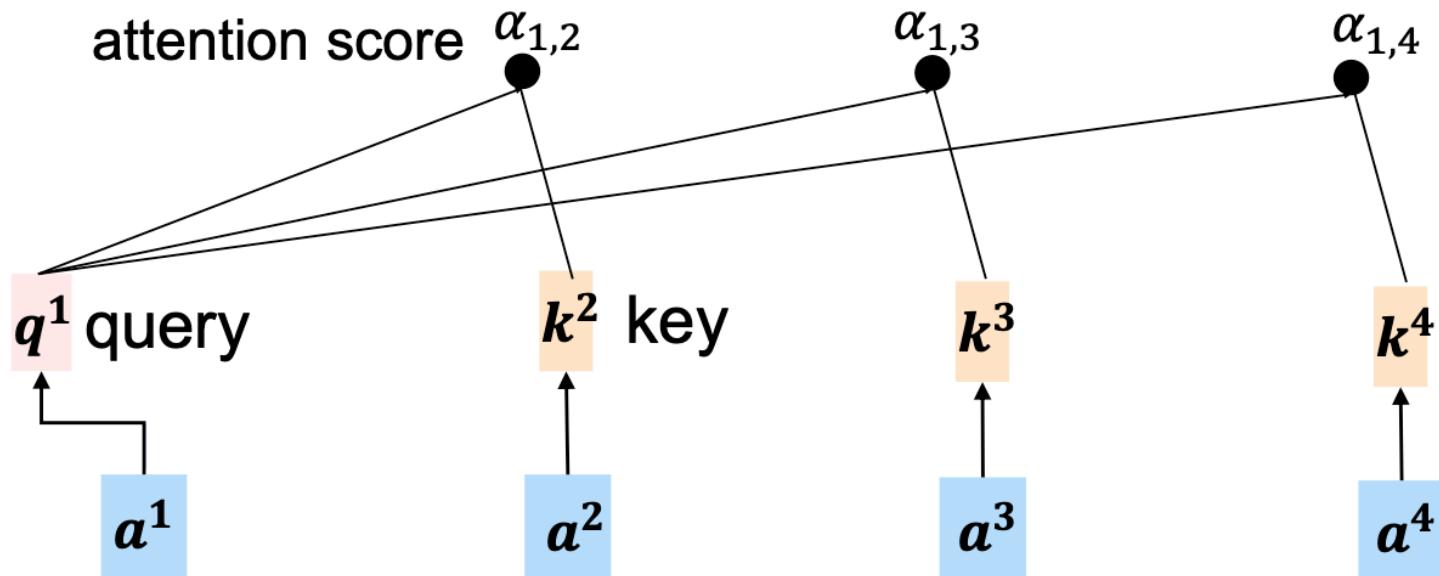
# Self-Attention



Can be either **input or a hidden layer**

# Self-Attention

$$\alpha_{1,2} = q^1 \cdot k^2 \quad \alpha_{1,3} = q^1 \cdot k^3 \quad \alpha_{1,4} = q^1 \cdot k^4$$



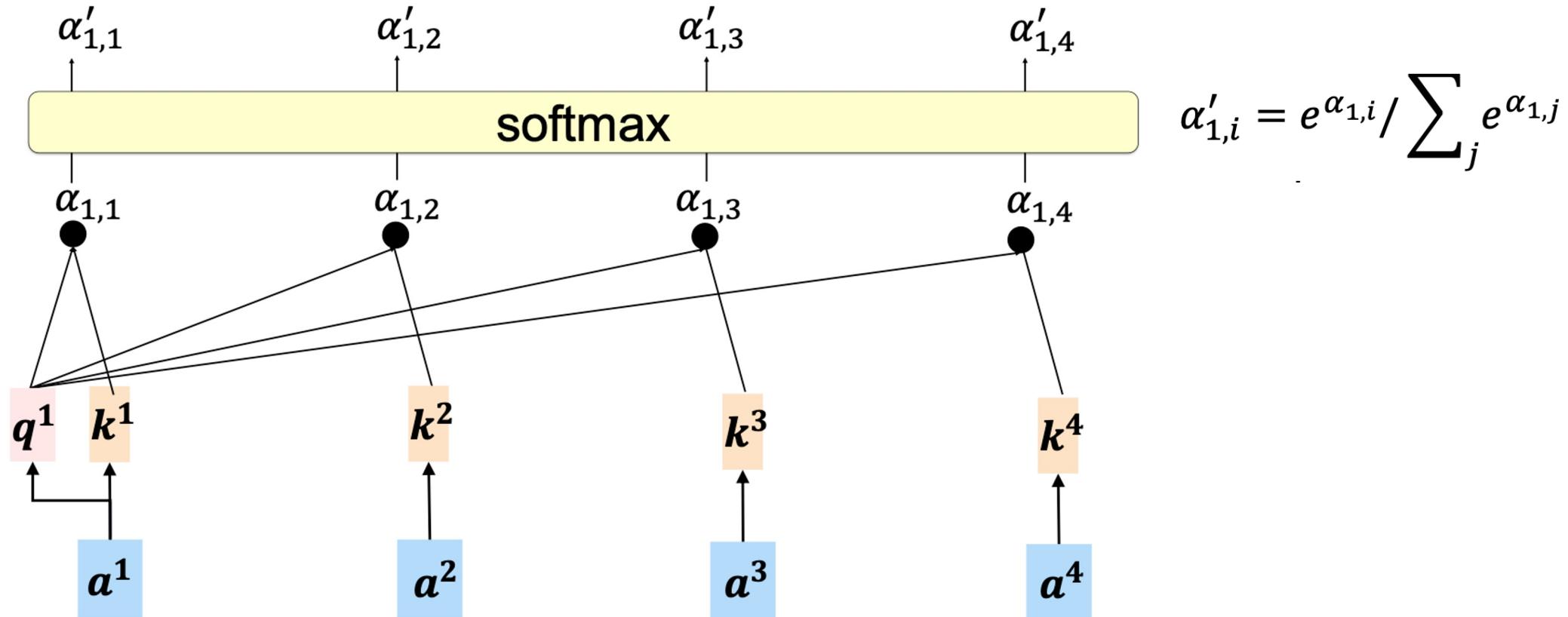
$$q^1 = W^q a^1$$

$$k^2 = W^k a^2$$

$$k^3 = W^k a^3$$

$$k^4 = W^k a^4$$

# Self-Attention



$$q^1 = W^q a^1$$

$$k^2 = W^k a^2$$

$$k^3 = W^k a^3$$

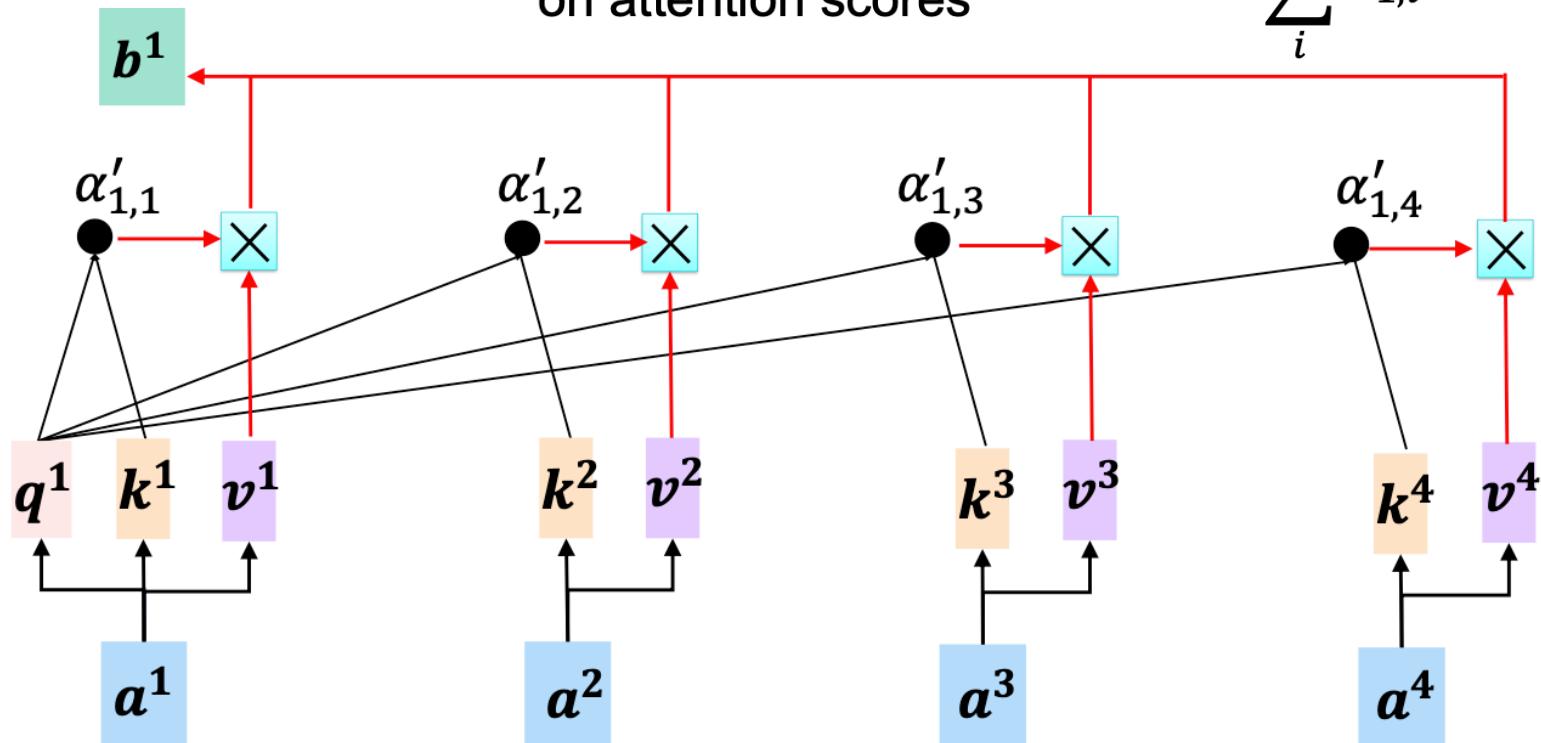
$$k^4 = W^k a^4$$

$$k^1 = W^k a^1$$

# Self-Attention

extract information based  
on attention scores

$$b^1 = \sum_i \alpha'_{1,i} v^i$$



$$v^1 = W^v a^1$$

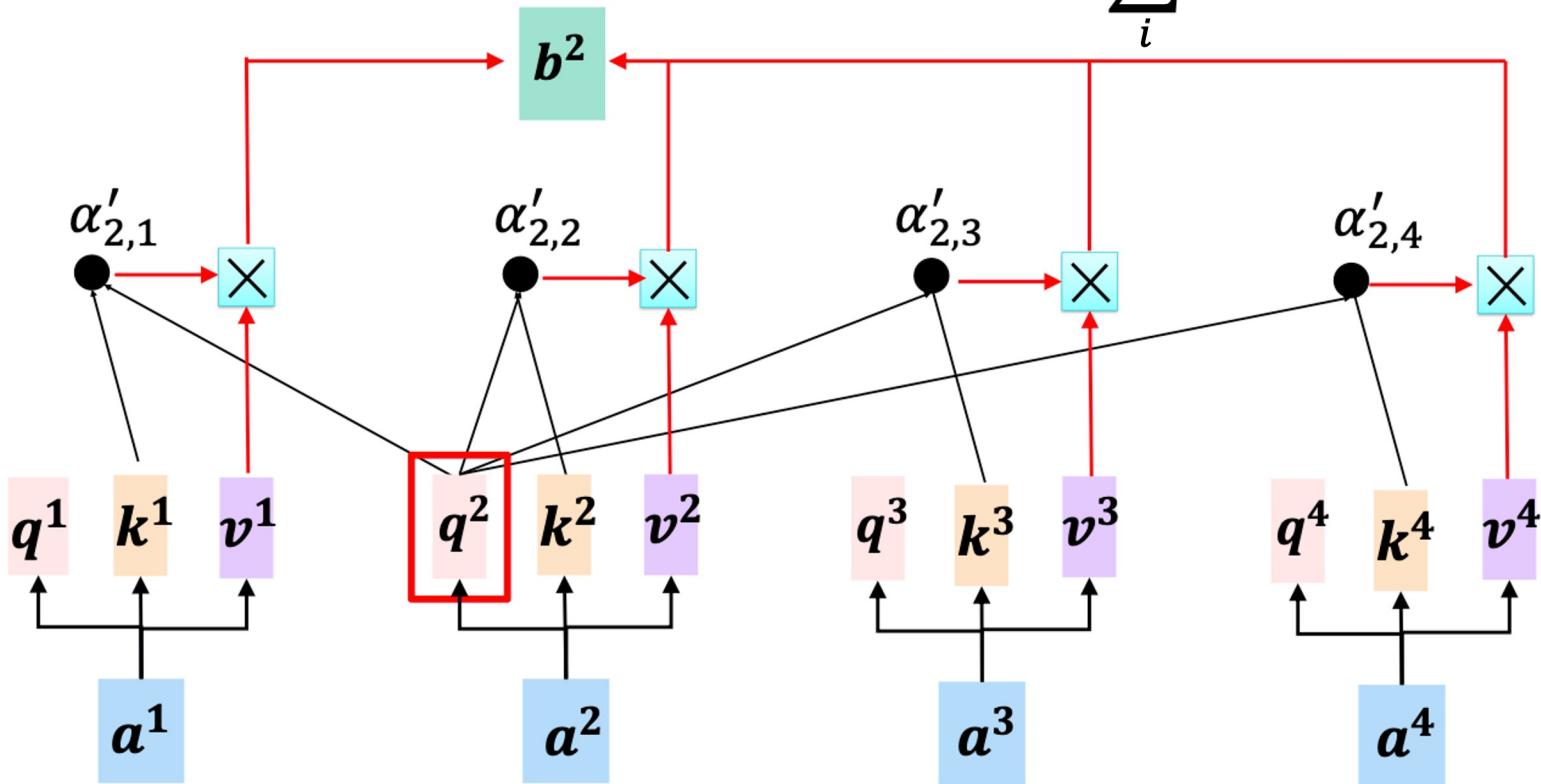
$$v^2 = W^v a^2$$

$$v^3 = W^v a^3$$

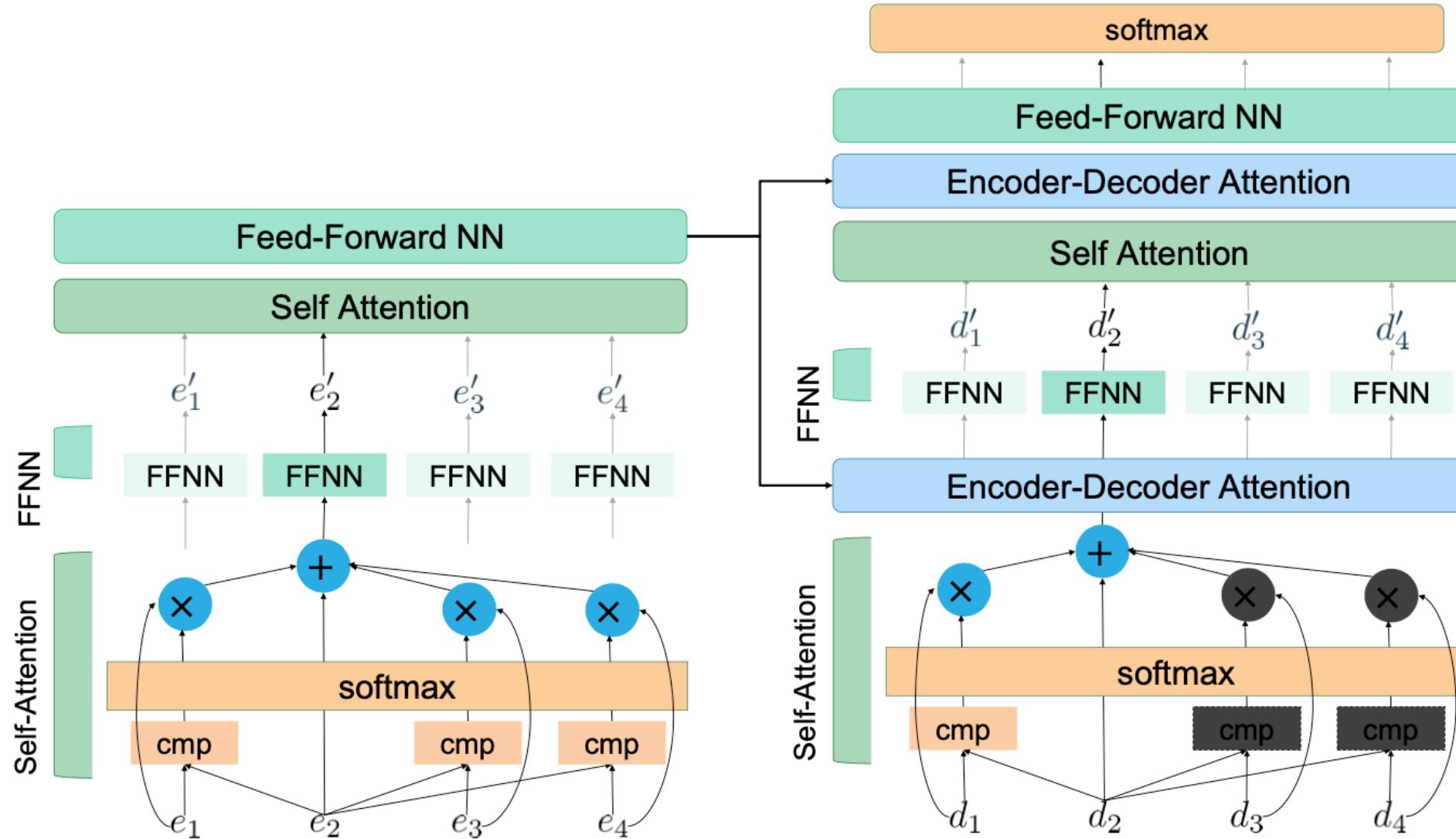
$$v^4 = W^v a^4$$

# Self-Attention

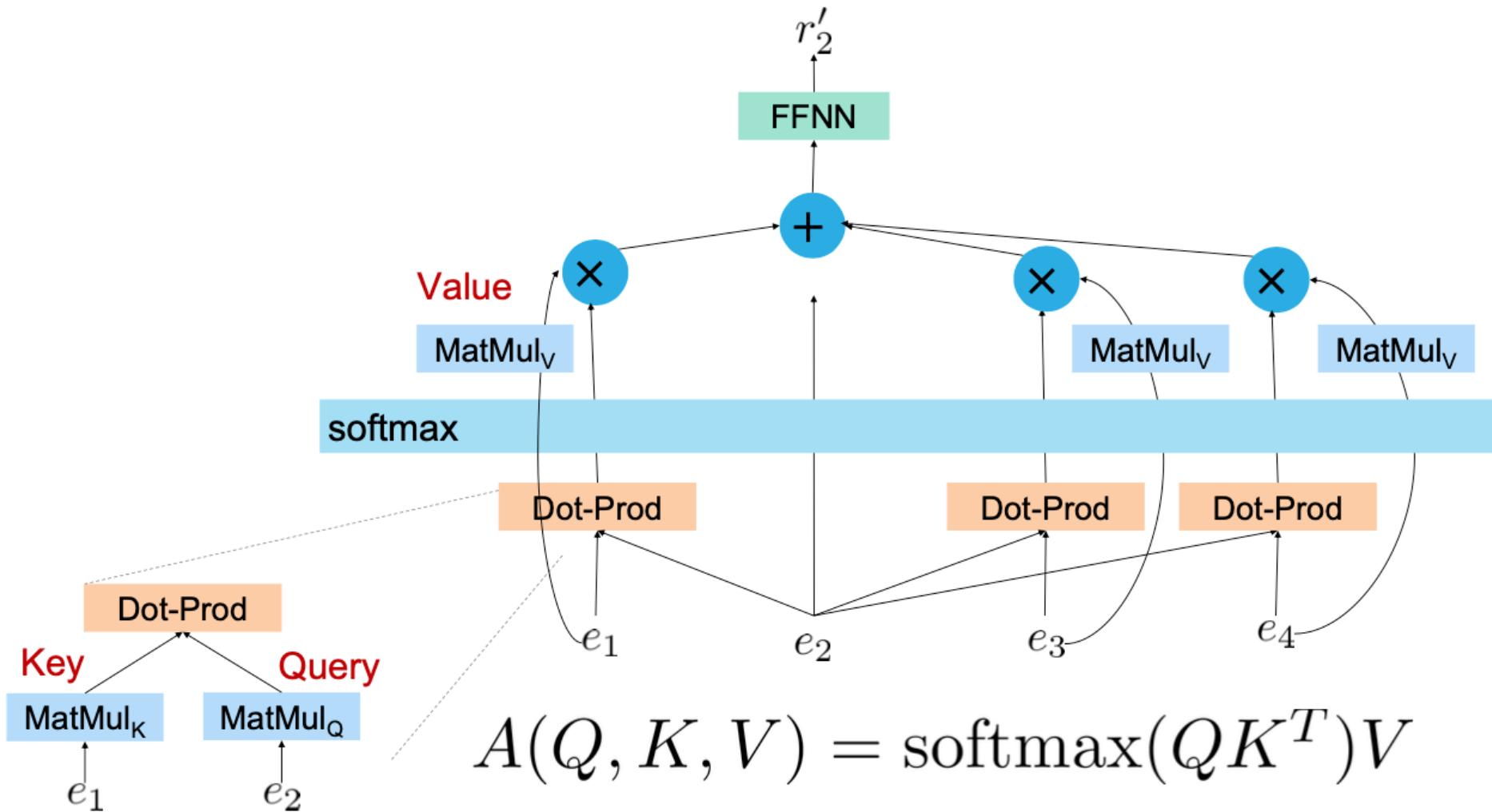
$$b^2 = \sum_i \alpha'_{2,i} v^i$$



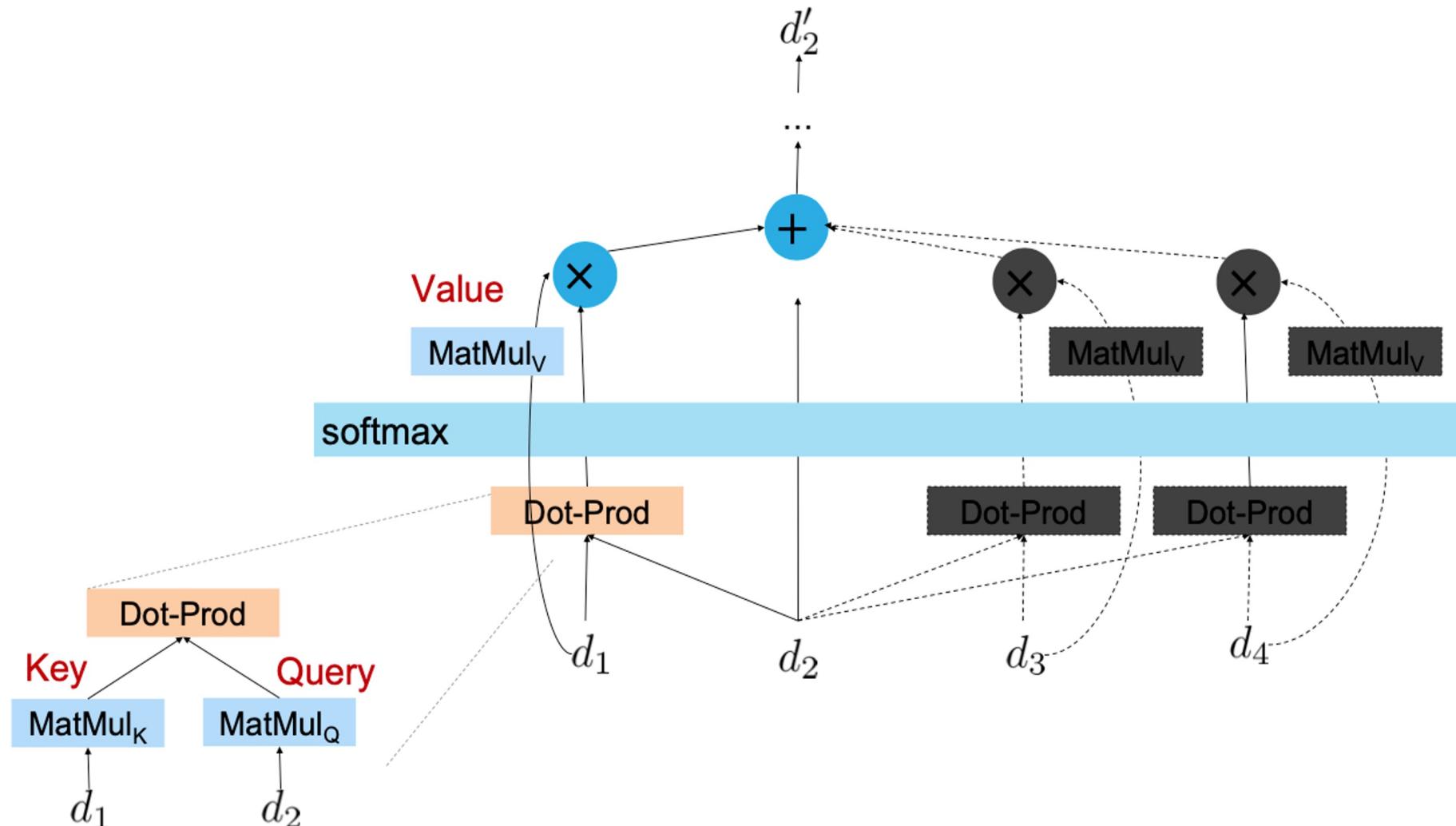
# Transformer Idea



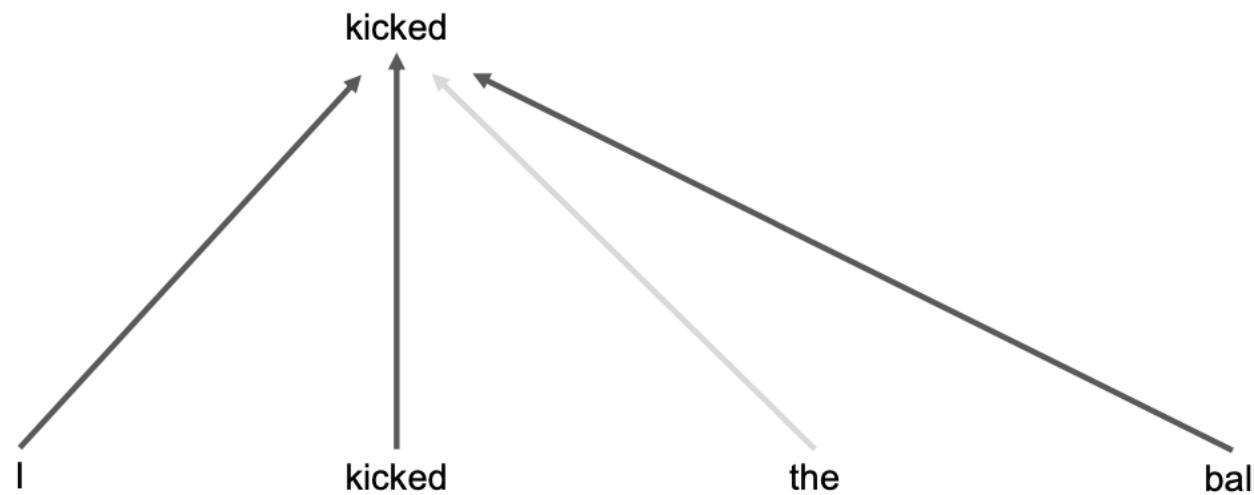
# Encoder Self-Attention



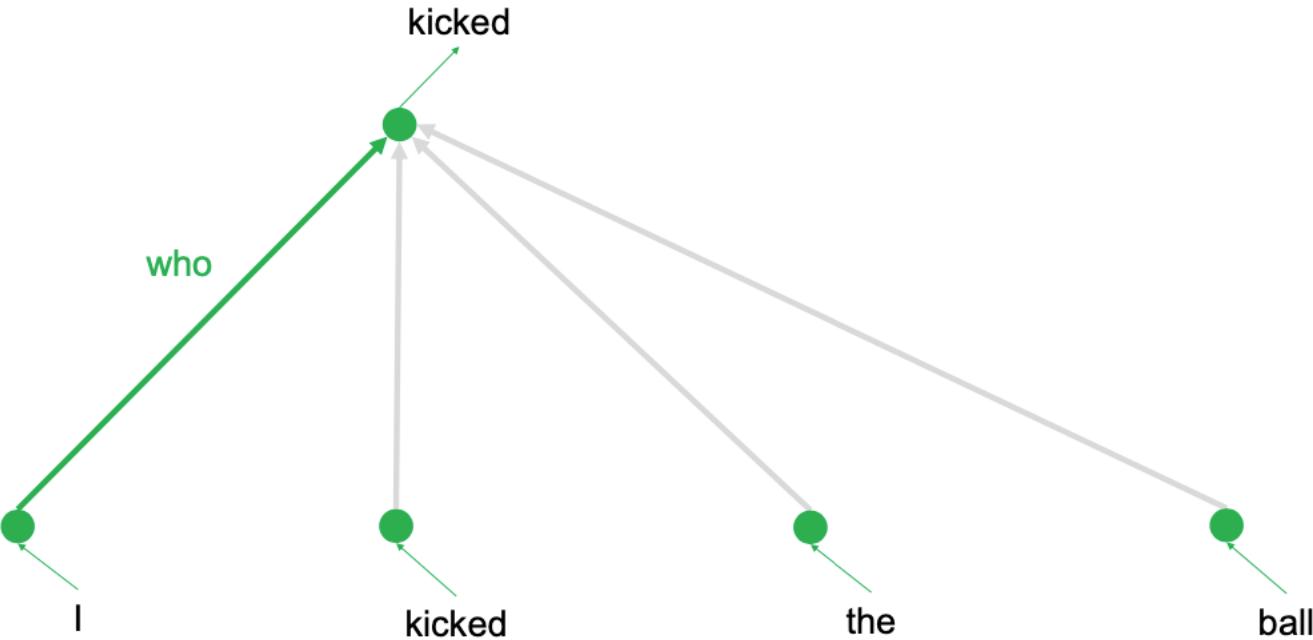
# Decoder Self-Attention



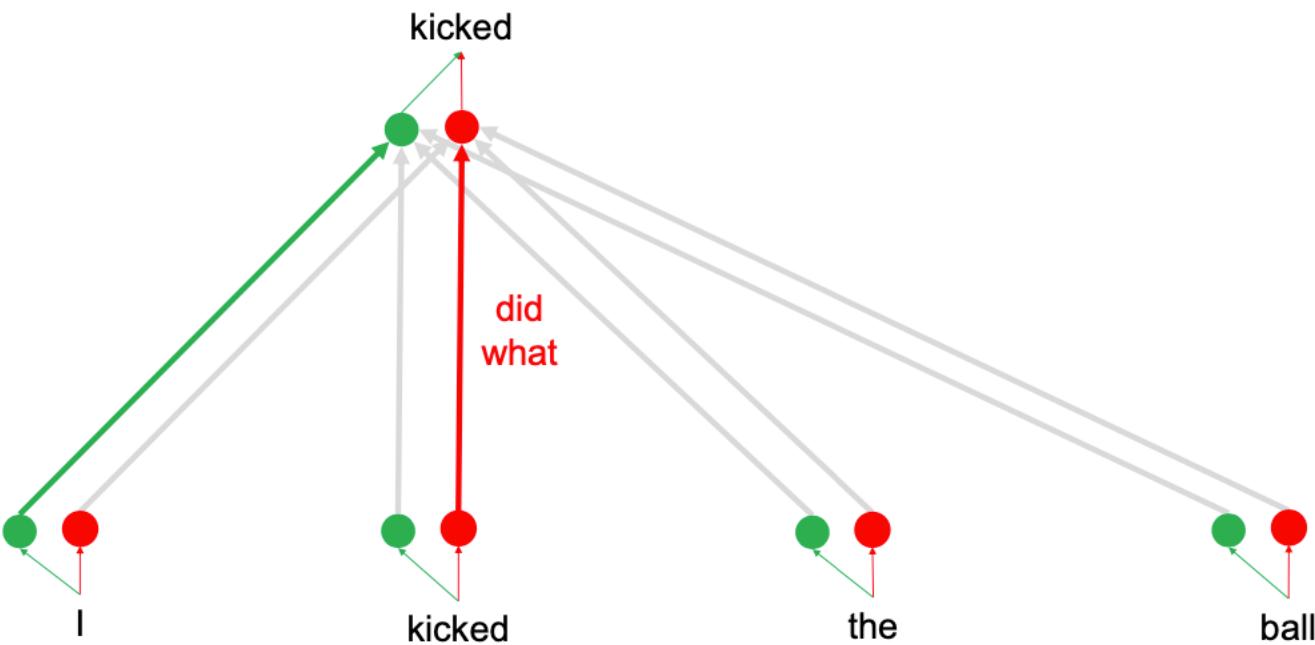
# Self-Attention



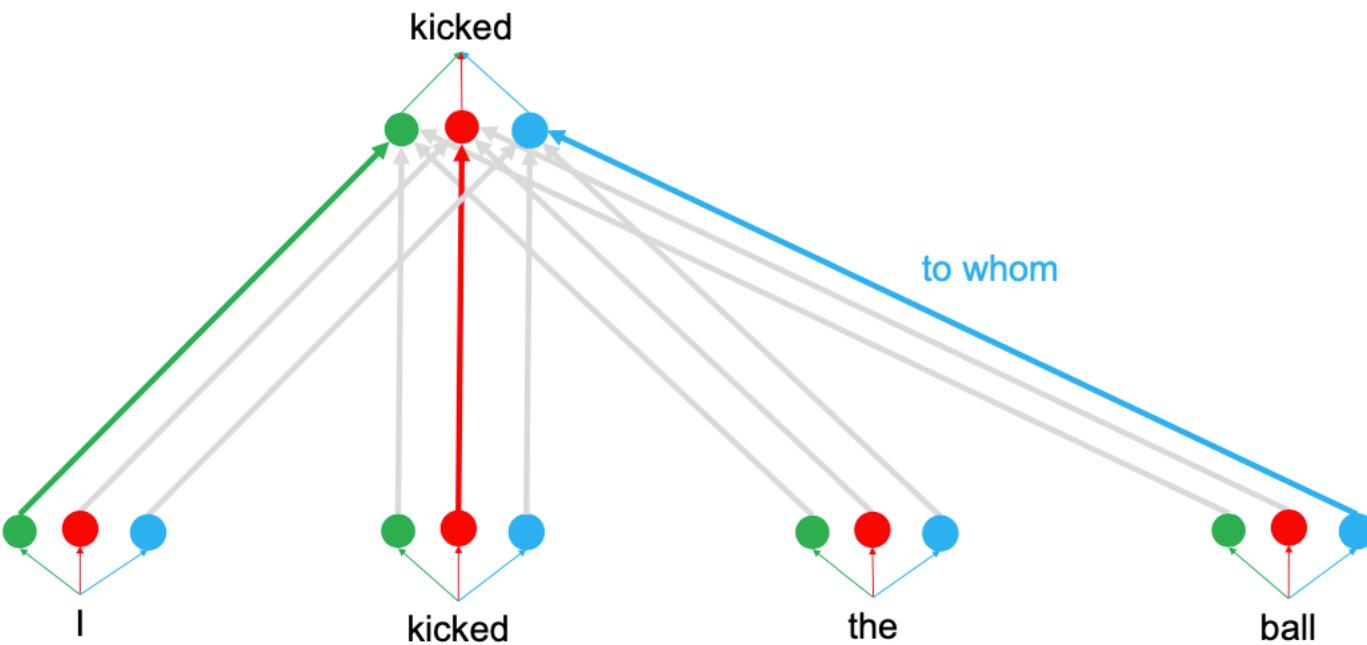
# Self-Attention



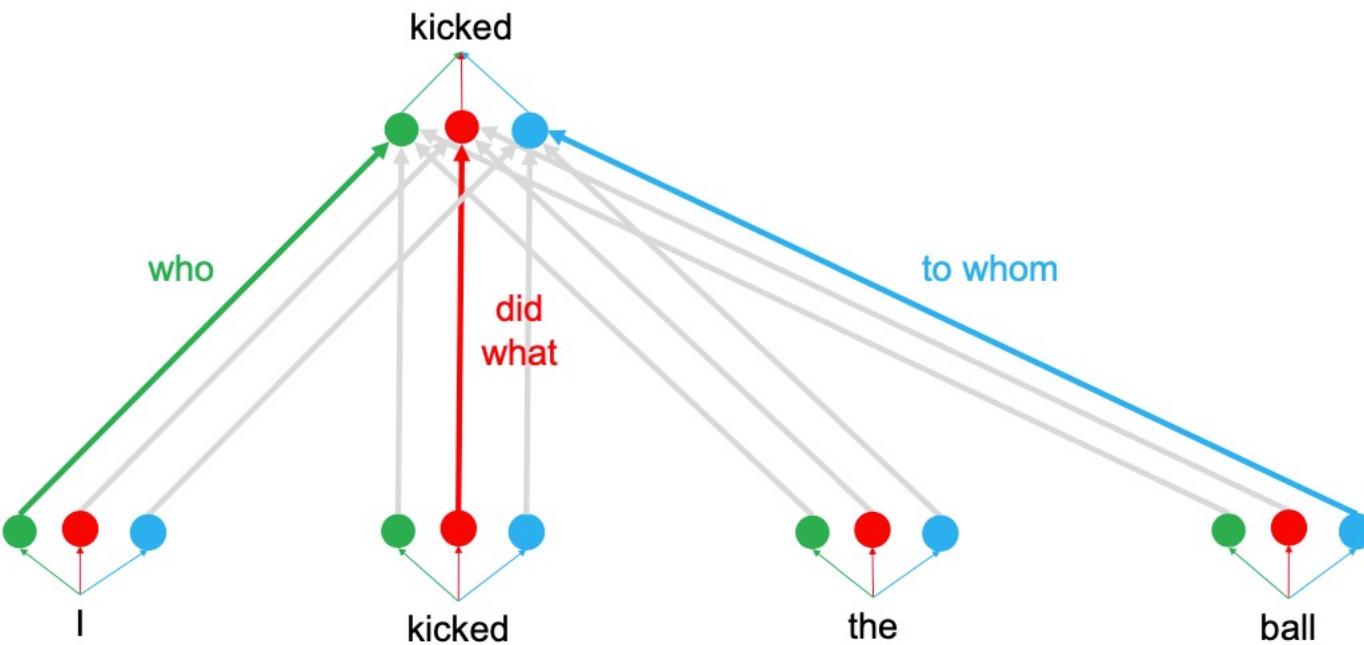
# Self-Attention



# Self-Attention



# Self-Attention



# Comparison

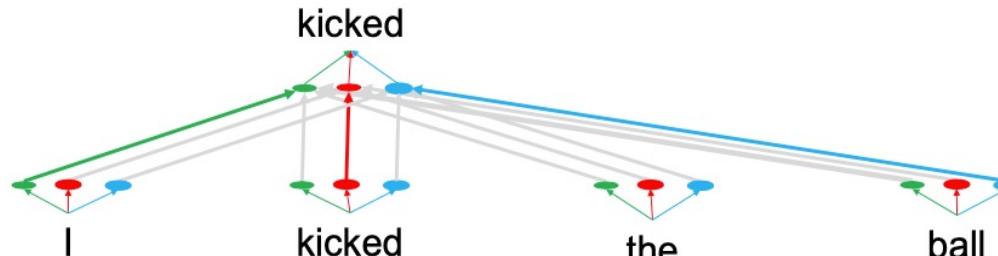
- Convolution: different linear transformations by relative positions



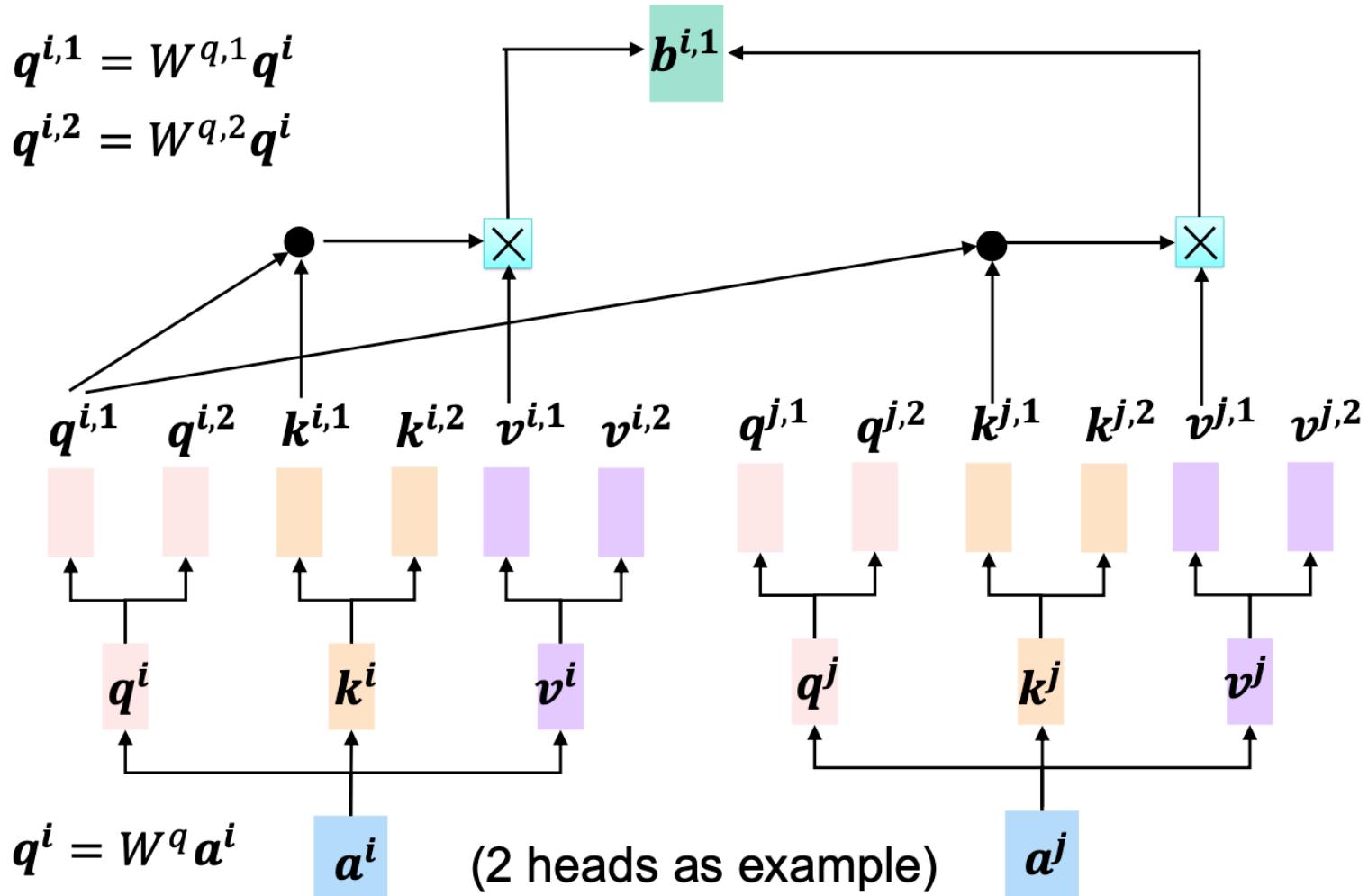
- Attention: a weighted average



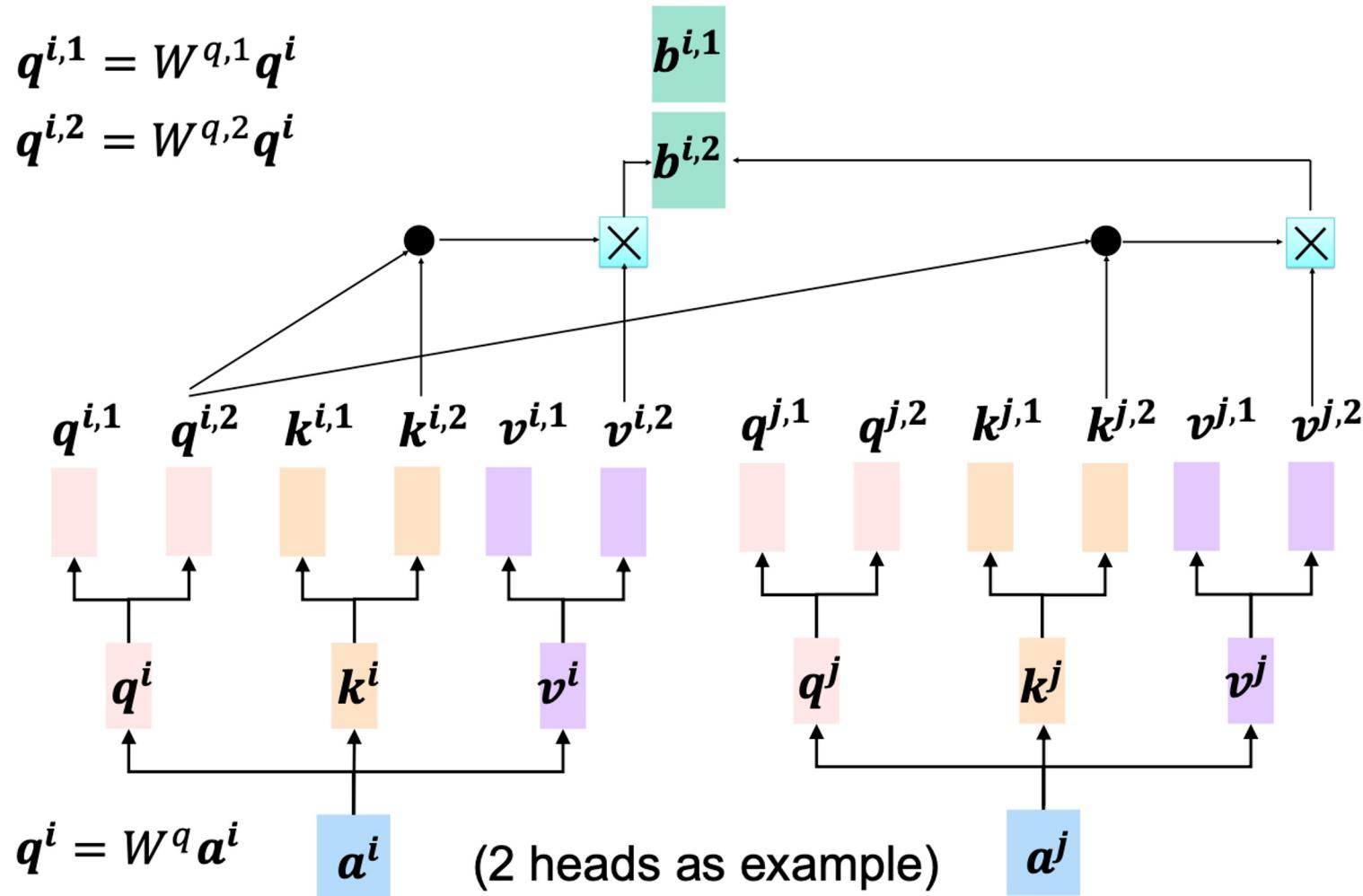
- Multi-Head Attention: parallel attention layers with different linear transformations on input/output



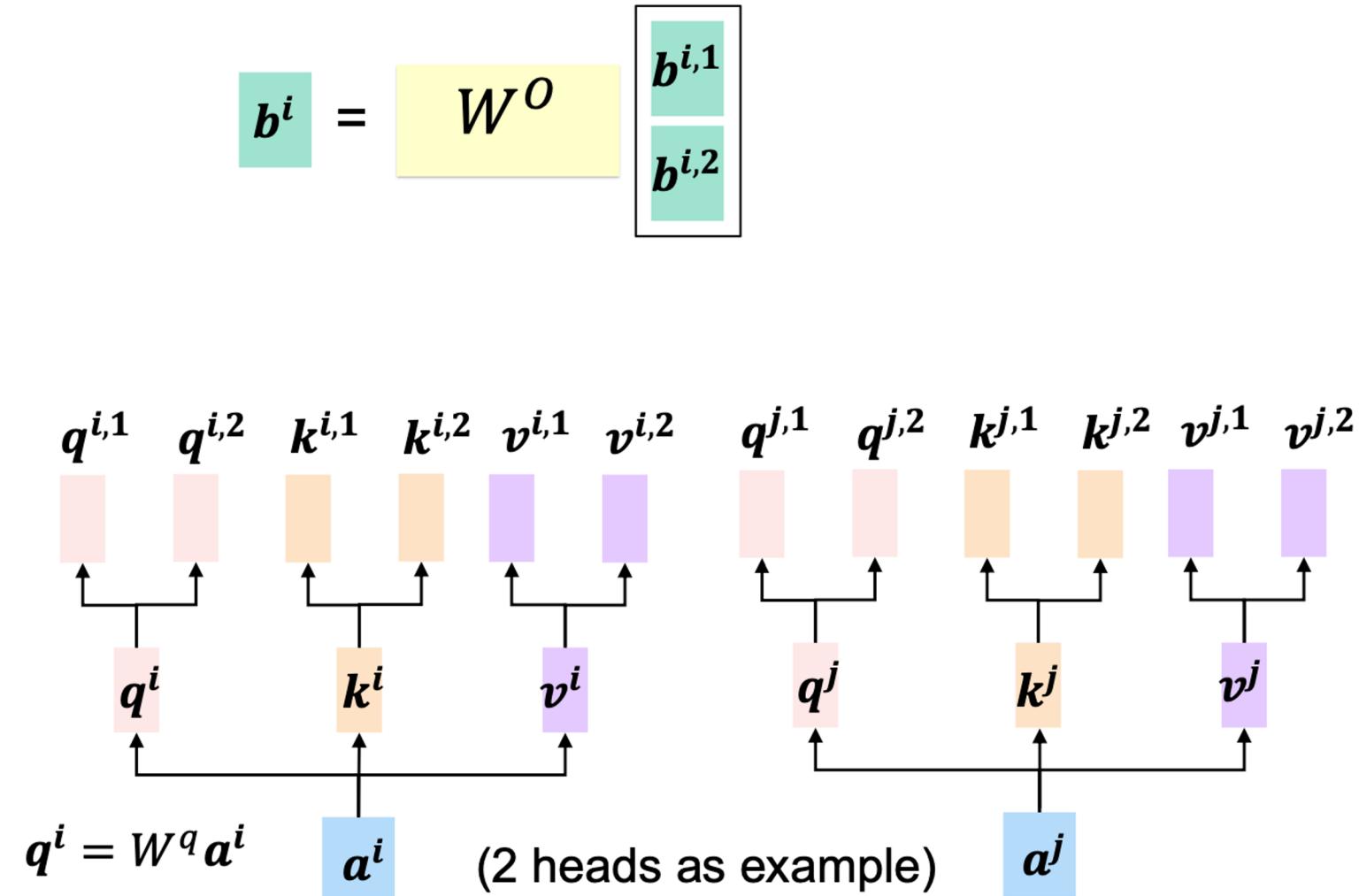
# Multi-Head Attention



# Multi-Head Attention



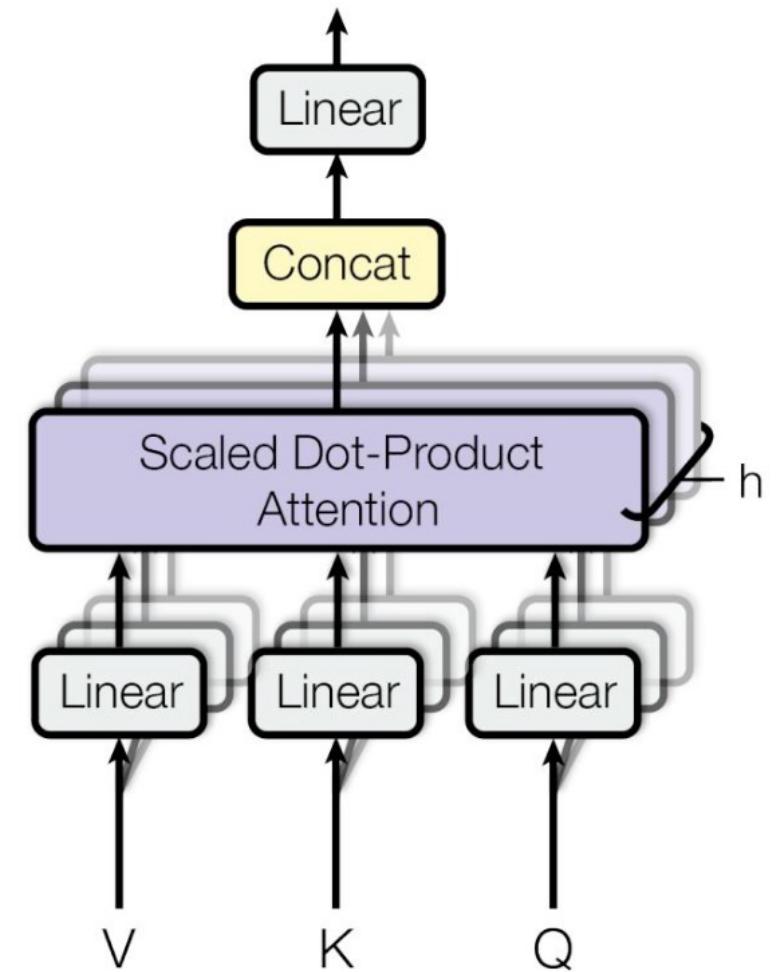
# Multi-Head Attention



# Multi-Head Attention

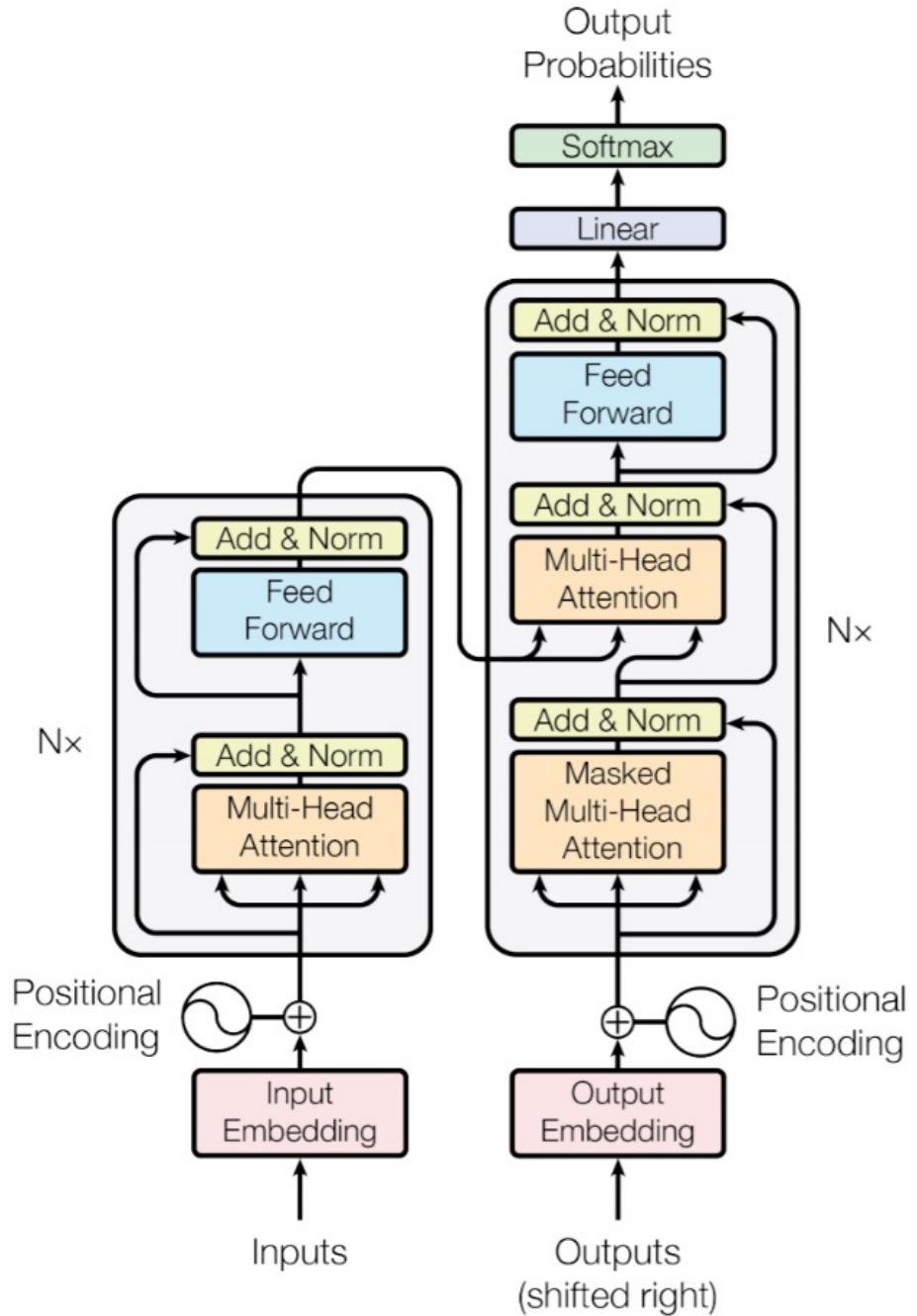
$$\begin{aligned}\text{MultiHead}(Q, K, V) \\ = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O\end{aligned}$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$



# Transformer Overview

$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



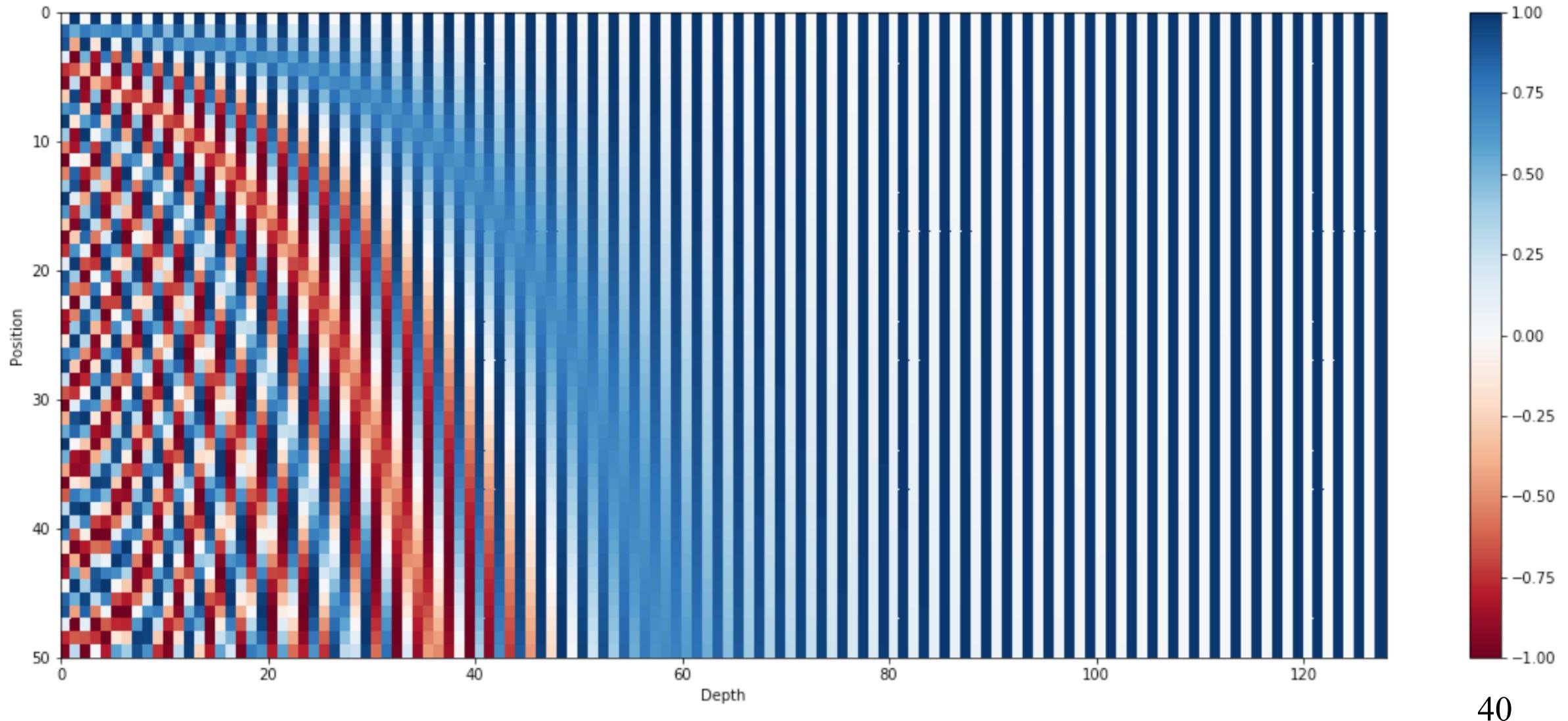
# Sinusoidal Positional Encoding

- Criteria for positional embeddings
  - Unique encoding for each position
  - Deterministic
  - Distance between neighboring positions should be the same
  - Model can easily generalize to longer sentences

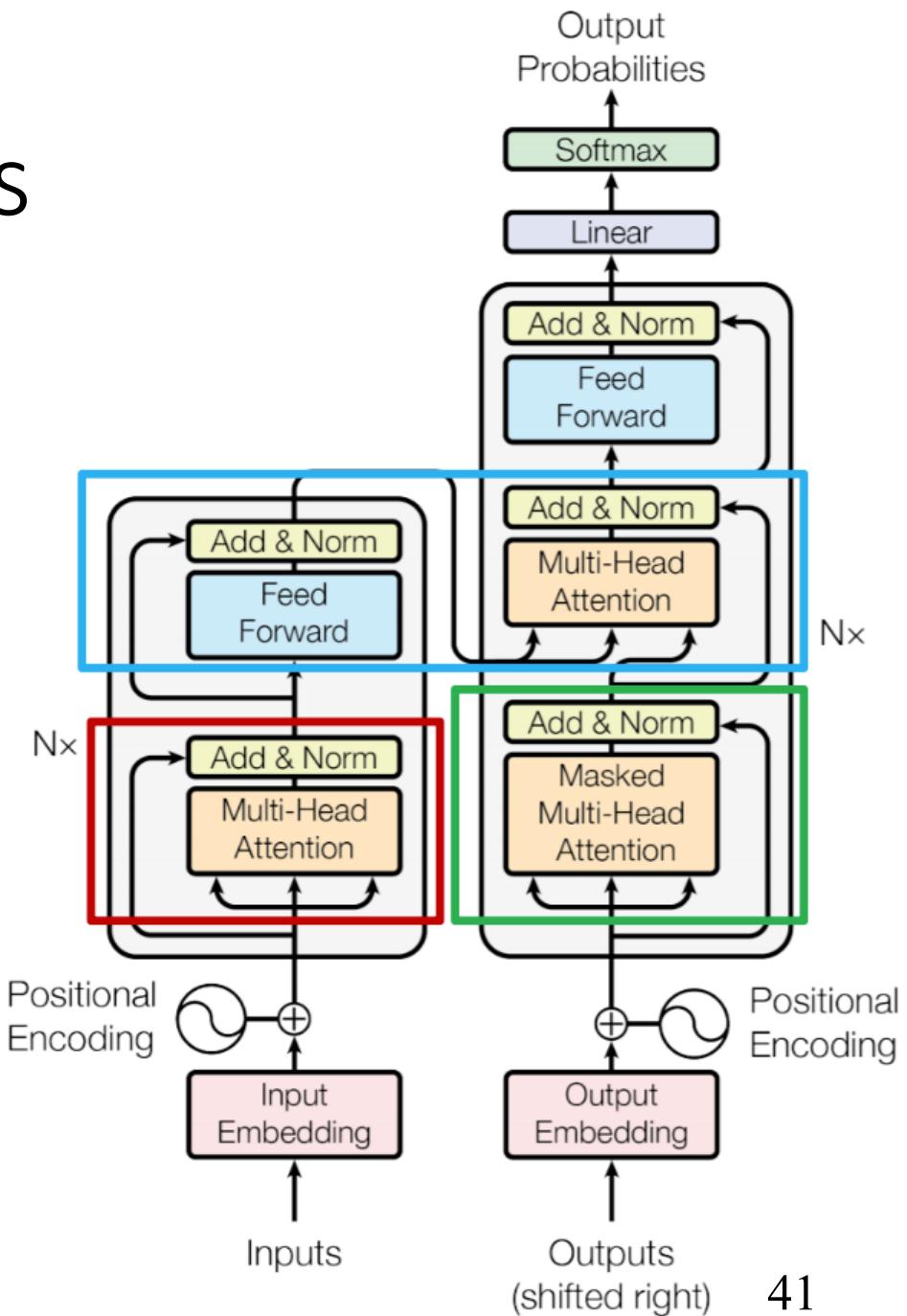
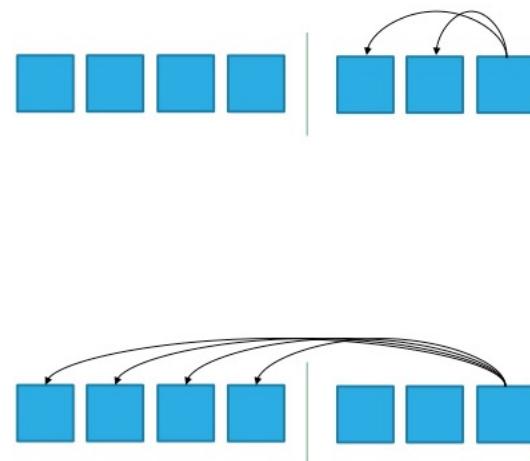
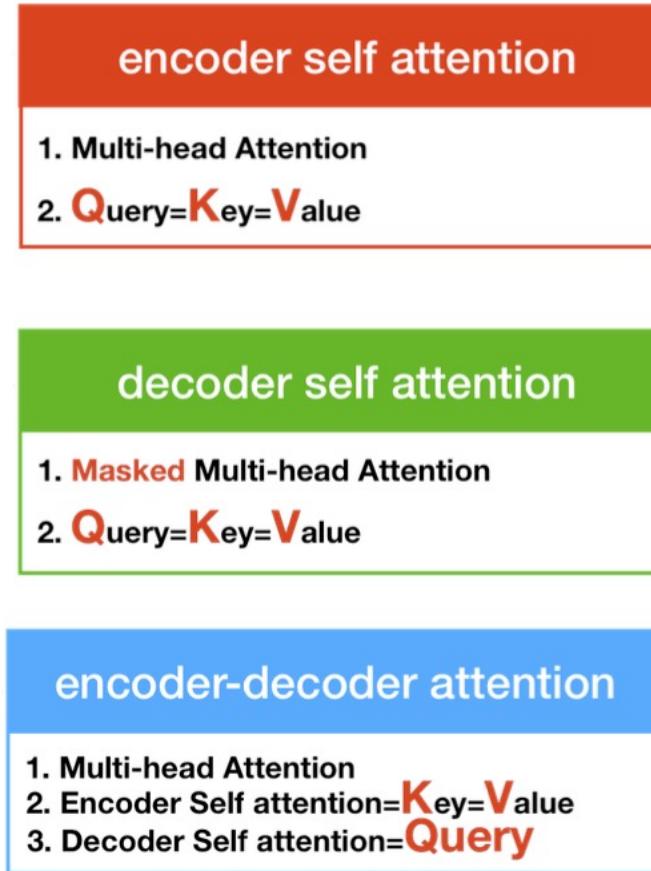
$$\text{PE}(pos, 2i) = \sin\left(\frac{pos}{100000^{2i/d}}\right)$$

$$\text{PE}(pos, 2i + 1) = \cos\left(\frac{pos}{100000^{2i/d}}\right)$$

# Sinusoidal Positional Encoding



# Multi-Head Attention Details



# MT Experiments

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1		<b><math>3.3 \cdot 10^{18}</math></b>
Transformer (big)	<b>28.4</b>	<b>41.8</b>		$2.3 \cdot 10^{19}$

# Concluding Remarks

- **Non-recurrence** model is easy to parallelize
- **Multi-head attention** captures different aspects by interacting between words
- **Positional encoding** captures location information
- Each transformer block can be **applied to diverse tasks**

# Outlines

- Attention Mechanism
- Transformer
- Vision Transformer
- Swin Transformer

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

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

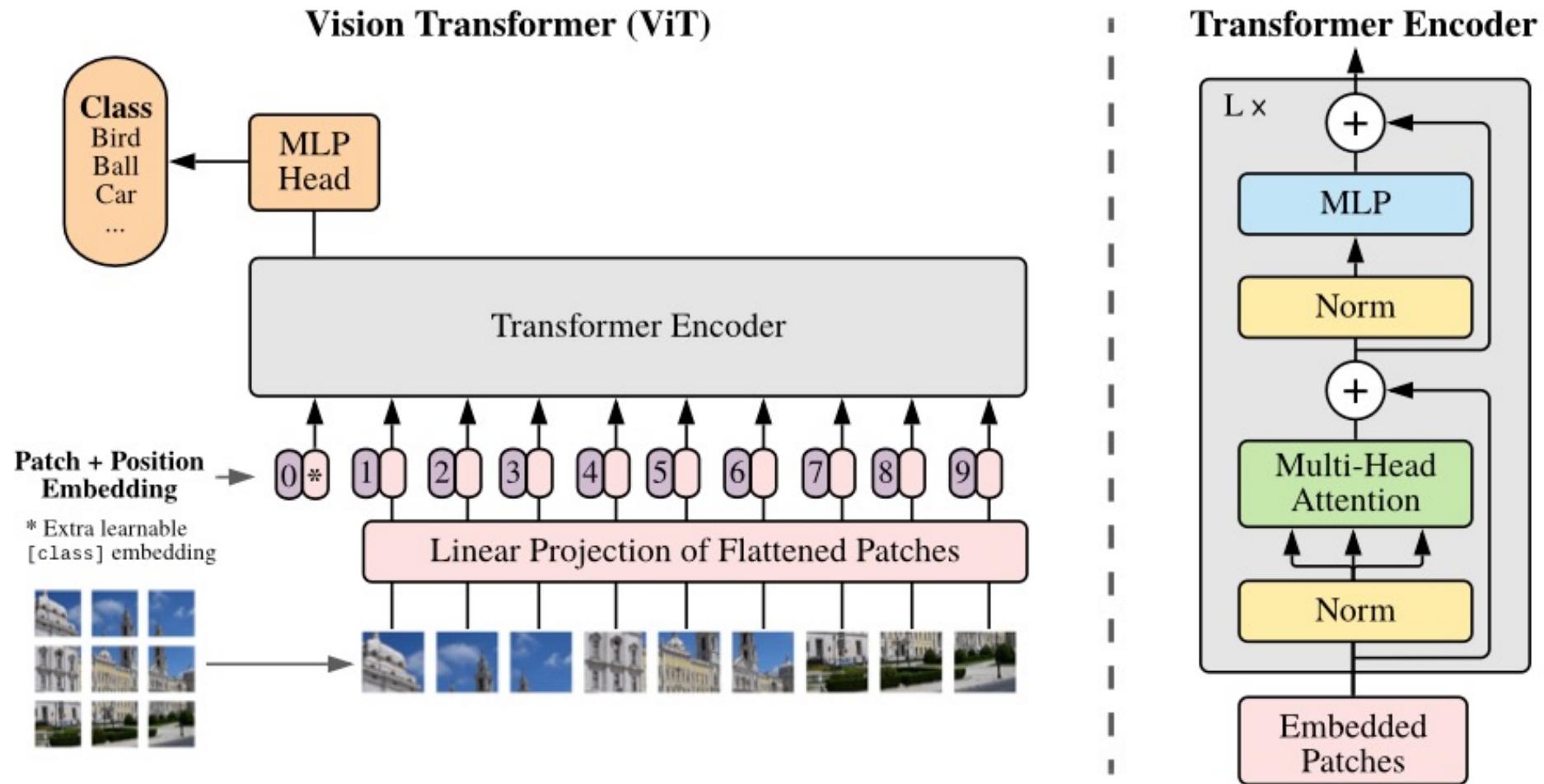
\*equal technical contribution, †equal advising

Google Research, Brain Team

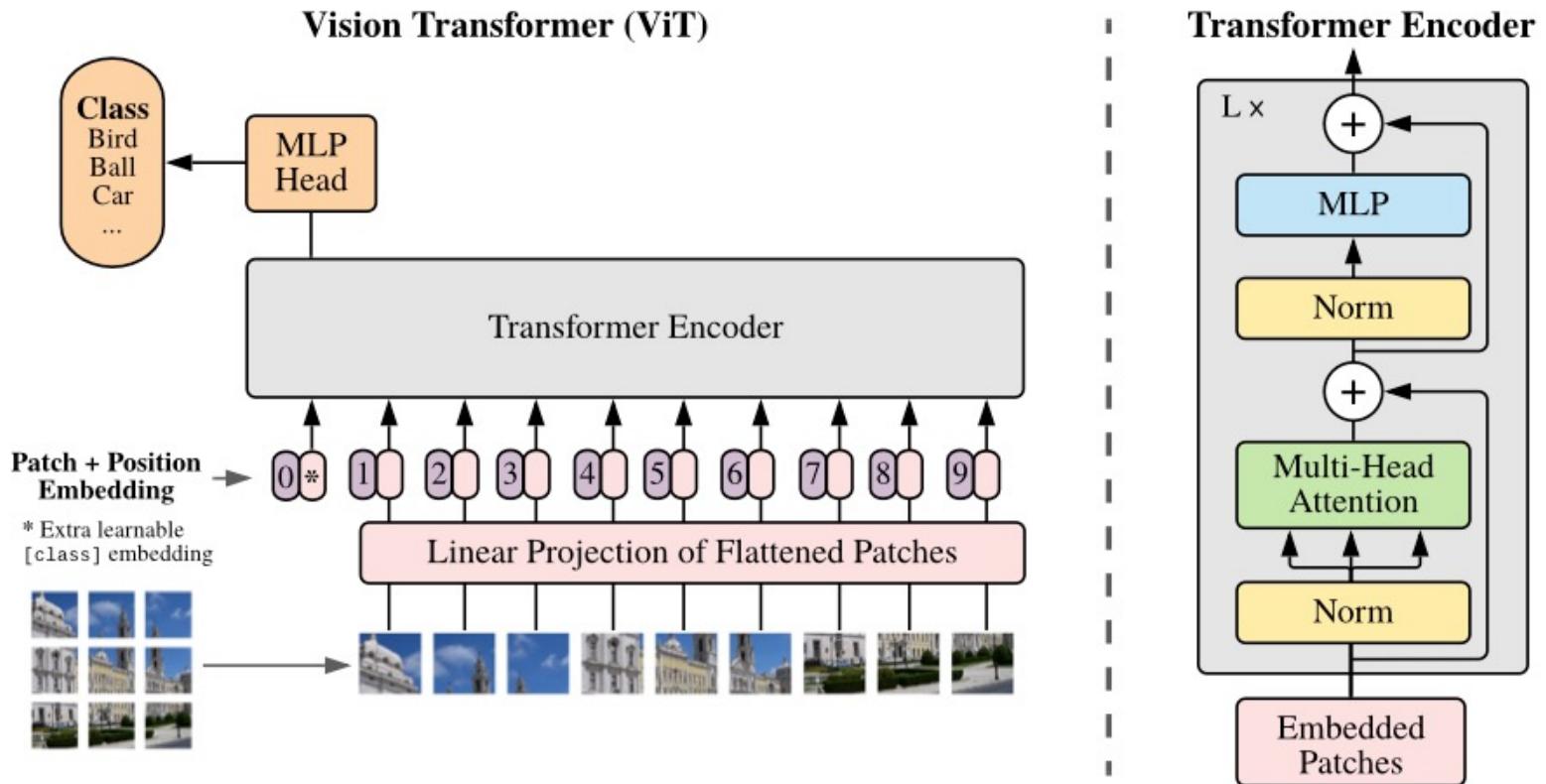
{adosovitskiy, neilhoulsby}@google.com

ICLR 2021

# Vision Transformer (ViT)



# Vision Transformer (ViT)



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

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

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

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

# Vision Transformer (ViT)

- Reshape the image  $x \in \mathbb{R}^{H*W*C}$  into a sequence of flattened 2D patches  $x_p \in \mathbb{R}^{N*(P^2*C)}$ ,  $N = HW/P^2$
- Maps to  $D$  dimensions
- Similar to BERT's [class] token, prepend a learnable patch ( $z_0^0 = x_{class}$ )
- Add positional embedding
- Alternating layers of MSA and MLP blocks
- Layernorm (LN) is applied before every block, and residual connections after every block

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

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

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

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

# Details of Vision Transformer model variants

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.

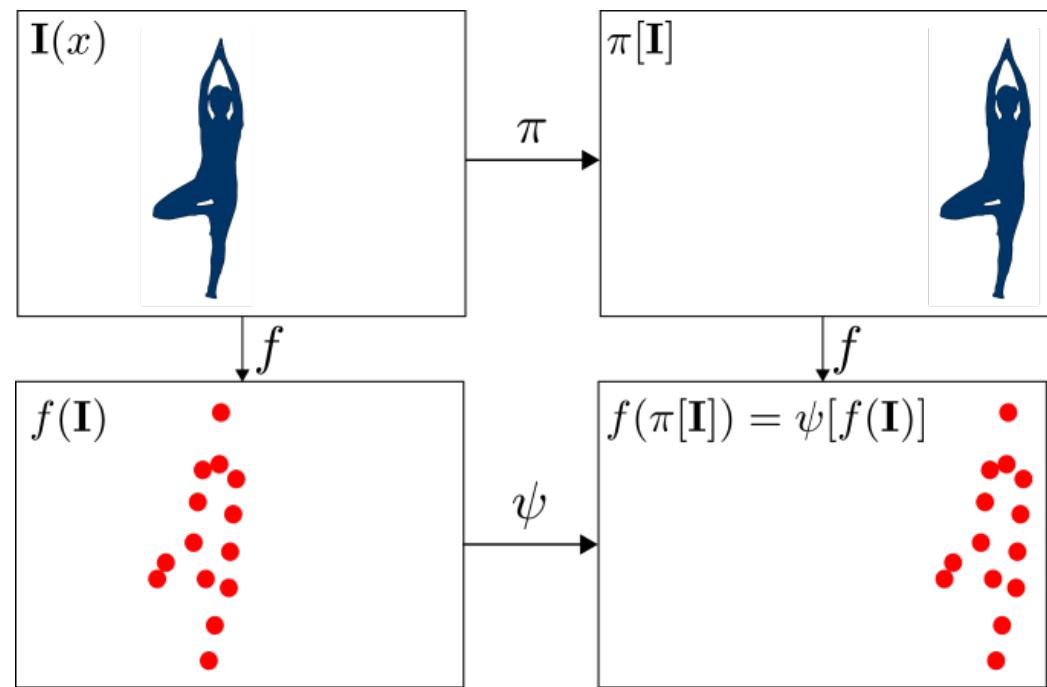
# Scores

名稱	模型	pre-train dataset
<b>Ous-JFT</b>	ViT-H/14	JFT-300M
<b>Ous-I21k</b>	ViT-L/16	ImageNet-21k
<b>BiT-L (Big Transfer)</b>	ResNet	JFT-300M
<b>Noisy Student</b>	EfficientNet	semi-supervised learning on ImageNet and JFT-300M with the labels removed

	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

# Comparison

- Why better than CNN?
  - Without inductive bias
  - But need sufficient data
- Inductive bias of CNN
  - translation equivariance
    - This means that if an object in an image is shifted, the same feature will be detected regardless of its position in the image
  - Locality bias
    - local receptive fields
    - weight sharing



# Outlines

- Attention Mechanism
- Transformer
- Vision Transformer
- Swin Transformer

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

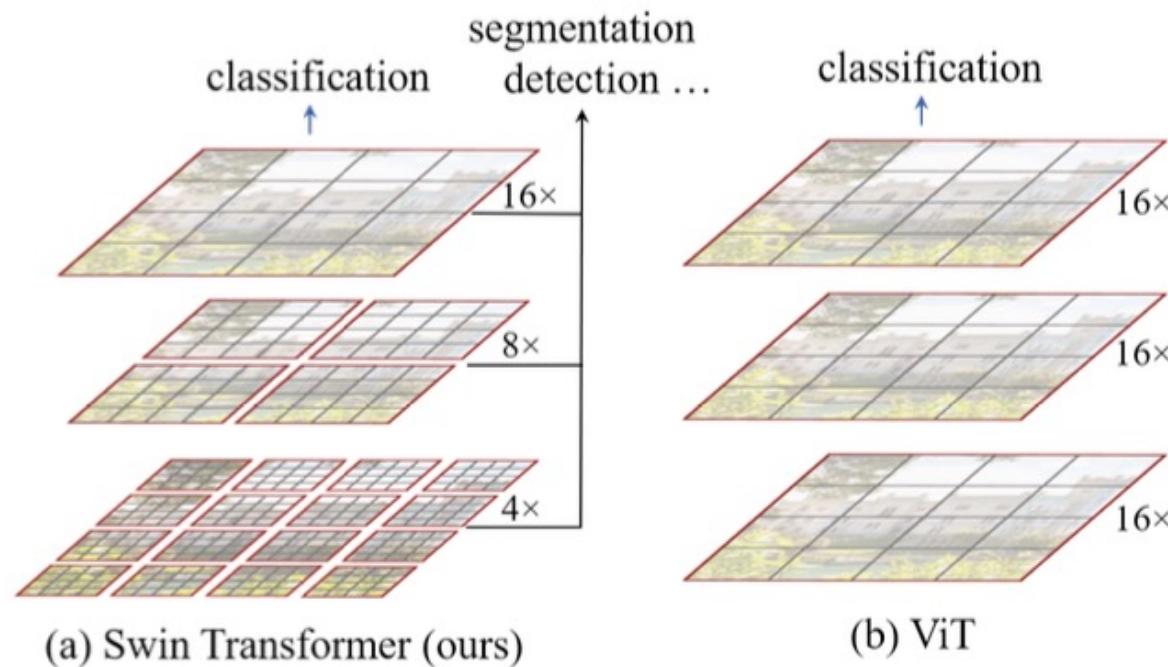
Ze Liu<sup>†\*</sup>   Yutong Lin<sup>†\*</sup>   Yue Cao<sup>\*</sup>   Han Hu<sup>\*‡</sup>   Yixuan Wei<sup>†</sup>  
Zheng Zhang   Stephen Lin   Baining Guo  
Microsoft Research Asia

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

ICCV 2021

# Swin Transformer

- Vision Transformers produce feature maps of a single low resolution cannot be used for advanced task.
- Swin transformer complexity becomes linear to image size.



# Swin Transformer

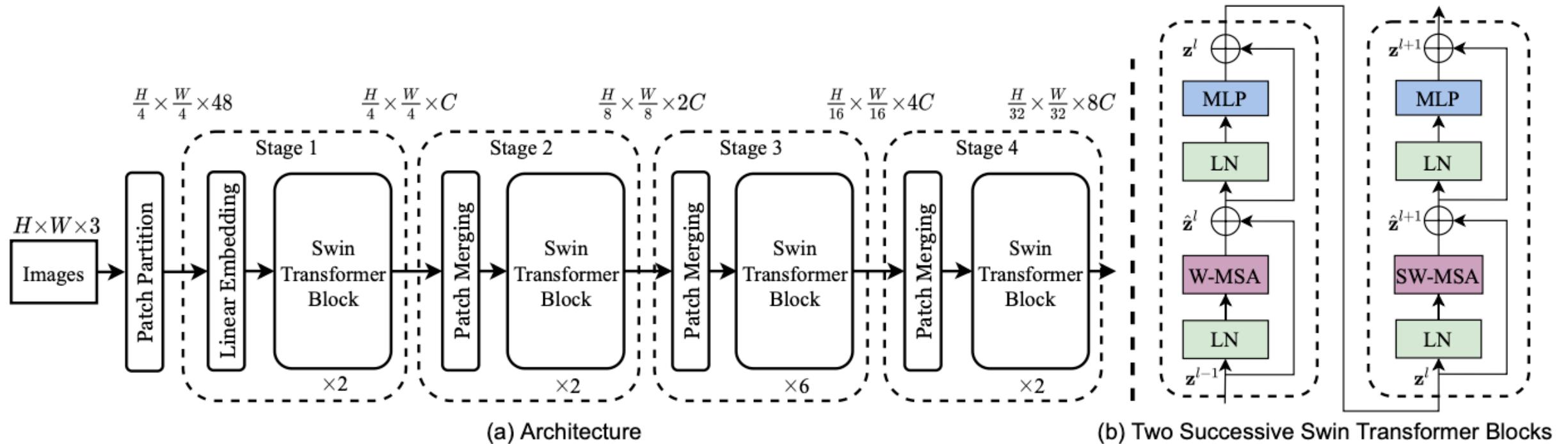
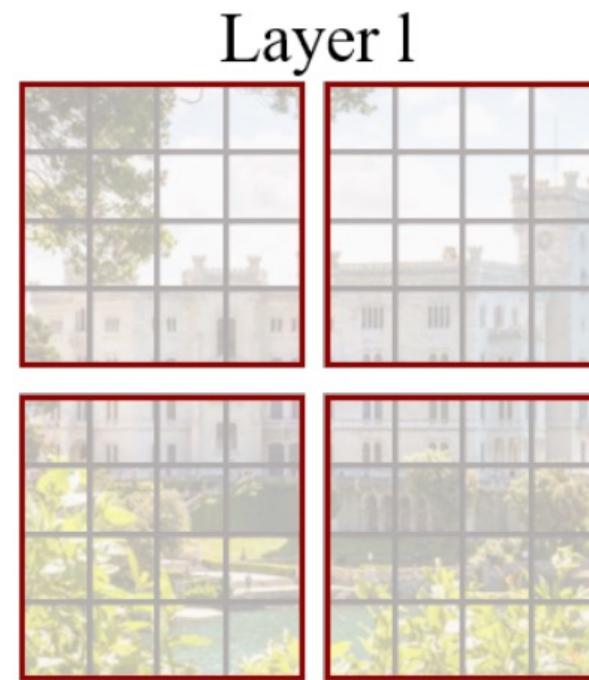


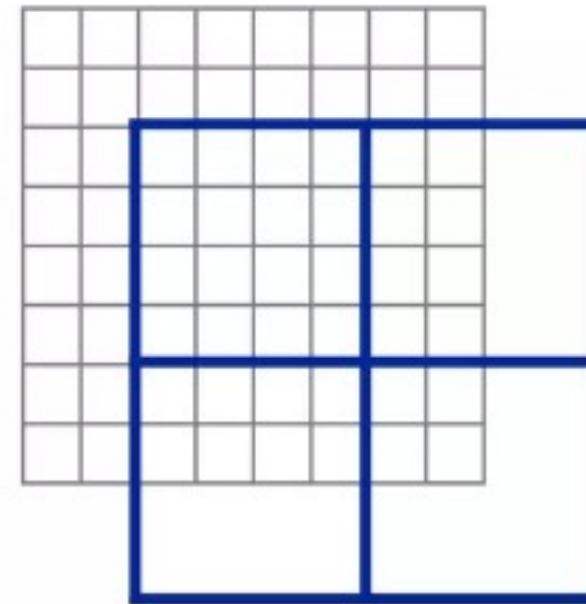
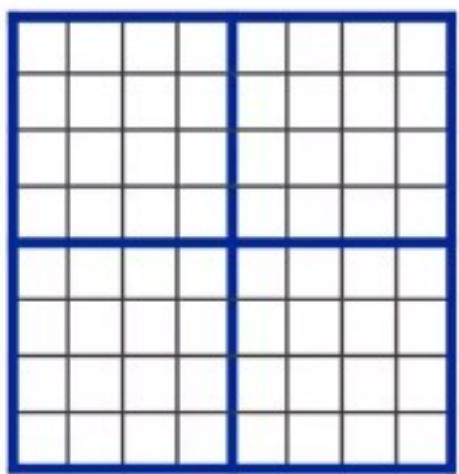
Figure 3. (a) The architecture of a Swin Transformer (Swin-T); (b) two successive Swin Transformer Blocks (notation presented with Eq. (3)). W-MSA and SW-MSA are multi-head self attention modules with regular and shifted windowing configurations, respectively.

# W-MSA

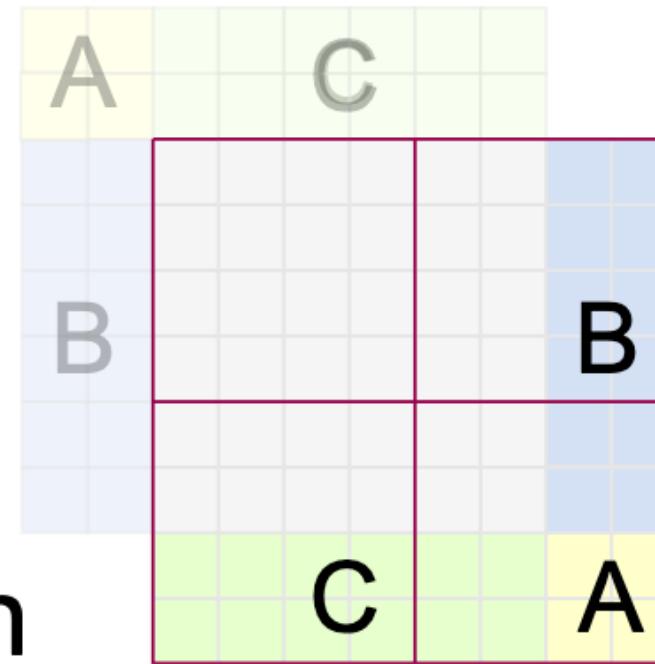
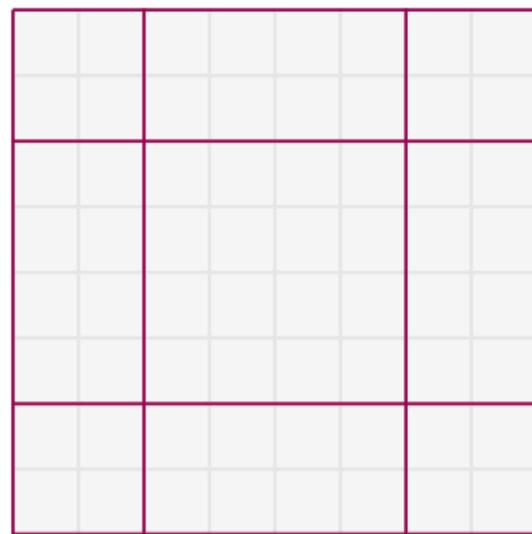
- Layer  $l$  has 4 windows, which means 4 attention mechanism



# SW-MSA



# SW-MSA

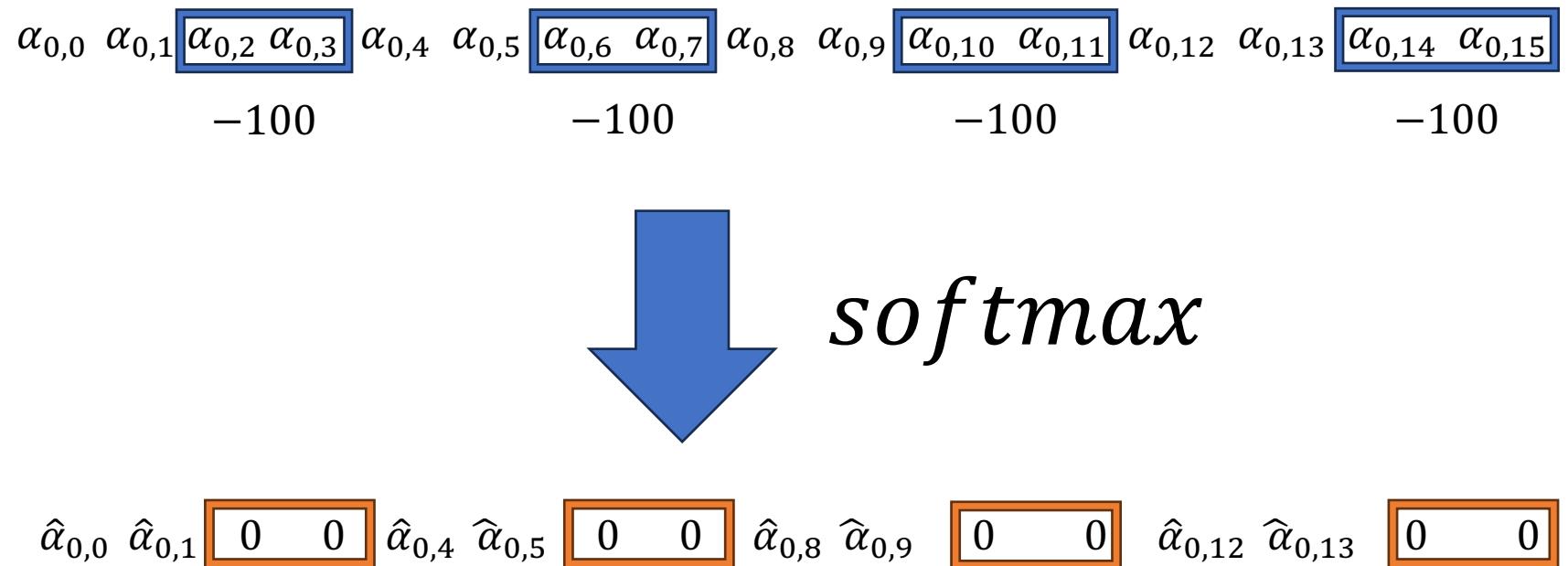


window partition

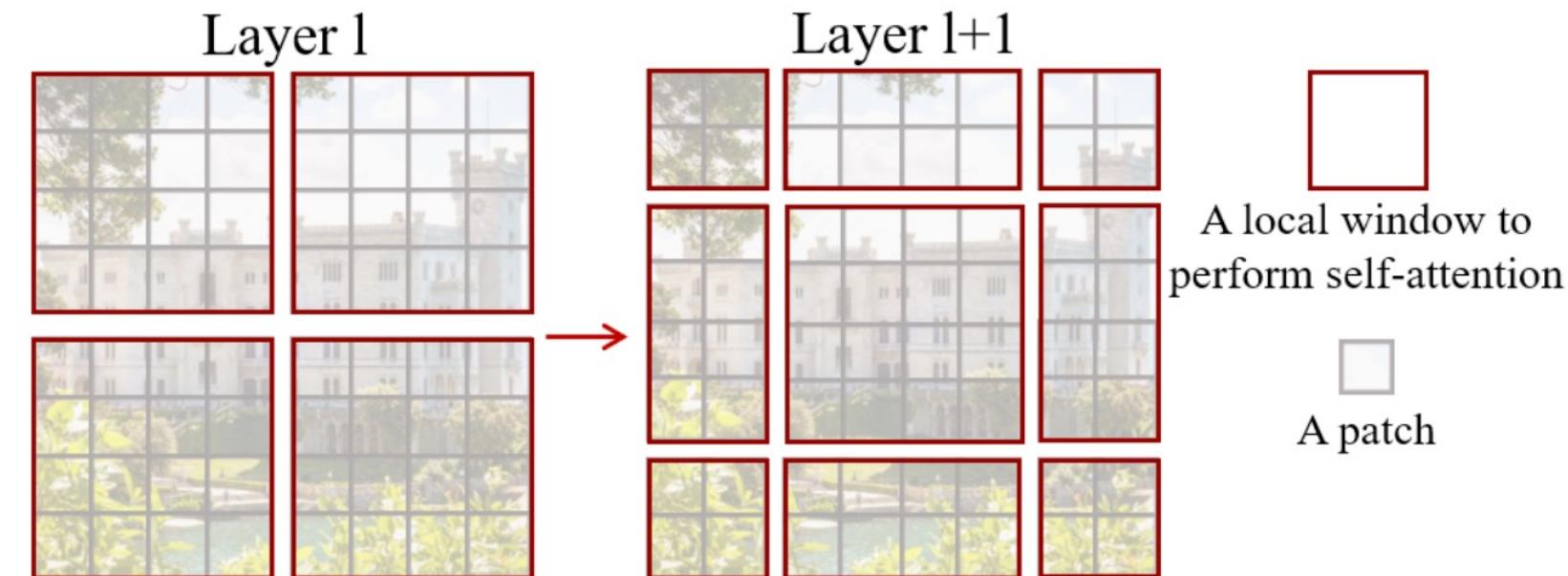
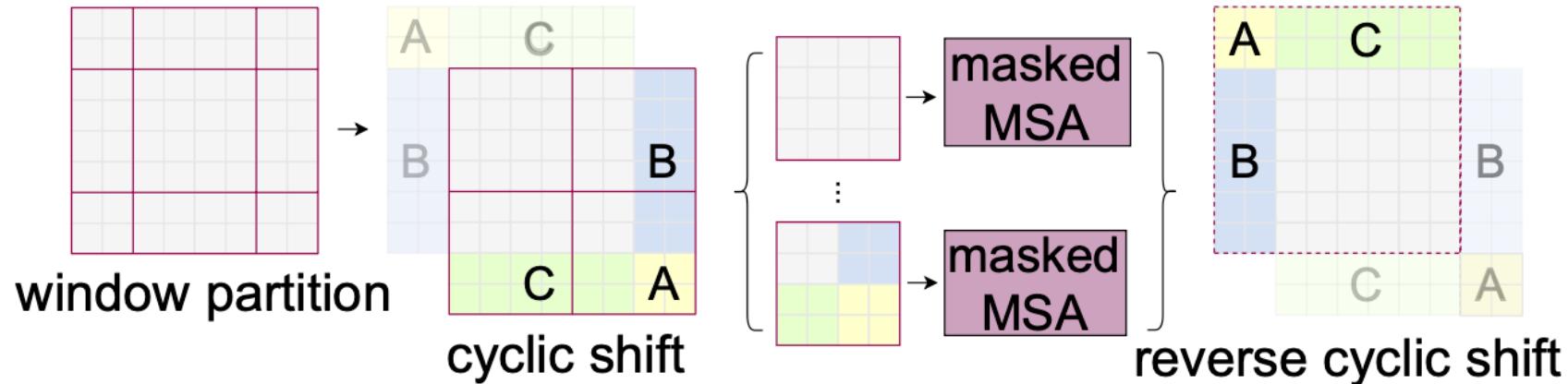
cyclic shift

# Masked MSA

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15



# SW-MSA

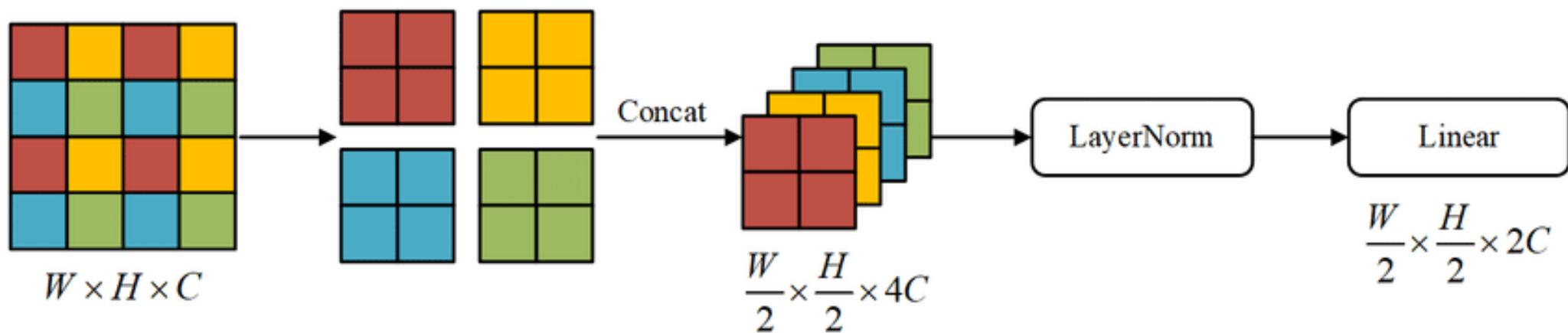


# Relative positive bias

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

名稱	依據	學習	添加位置
Transformer	sin/cos 函數	不可以學習	concat 在開頭
Visual Transformer	絕對位置	可學習	加在開頭
Swin Transformer	相對位置	可學習	加在 attention matrix

# Patch Merging



# Detailed architecture specifications

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

# Scores

(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

# Ablation Study

- Window
- Shifted window
- Patch merging
- Relatively Position Embedding

# Reference

- Attention Is All You Need(<https://arxiv.org/abs/1706.03762>)
- Applied Deep Learning 2022(<https://www.csie.ntu.edu.tw/~miulab/f111-adl/>)
- An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale(<https://arxiv.org/abs/2010.11929>)
- Swin Transformer: Hierarchical Vision Transformer using Shifted Windows(<https://arxiv.org/abs/2103.14030>)
- 210523 swin transformer v1.5(<https://www.slideshare.net/taeseonryu/210523-swin-transformer-v15>)