

Self-attention

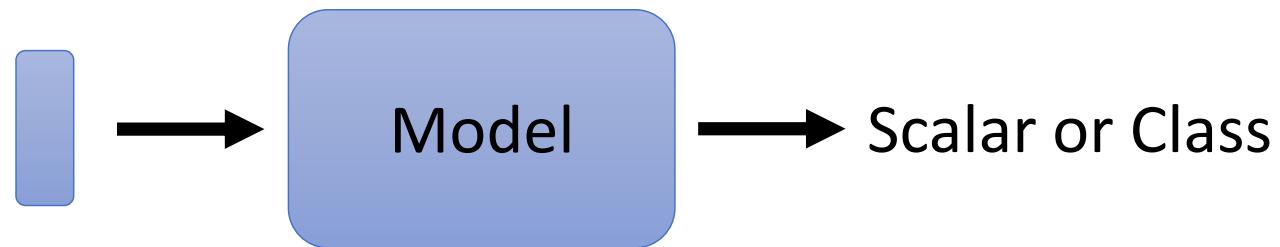
Hung-yi Lee

李宏毅

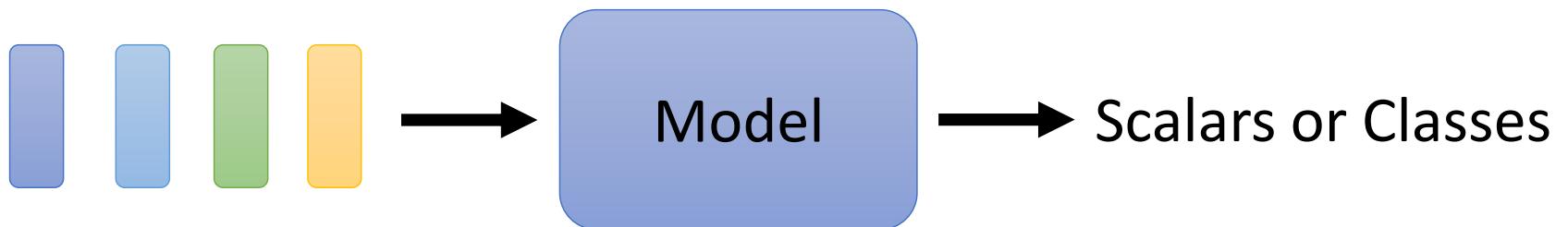
Sophisticated Input

→ Input は -J 向き

- Input is a **vector**



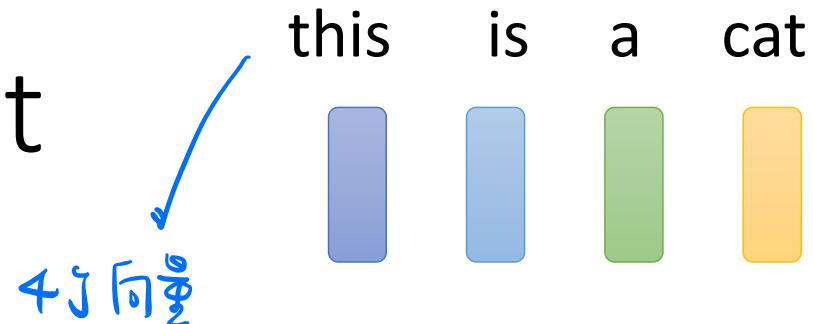
- Input is a **set of vectors** ⇒ Input は "-J 集合" 向き



(may change length)

而且、首次食糧追手の向き取、却不一様

Vector Set as Input



One-hot Encoding

apple = [1 0 0 0 0]

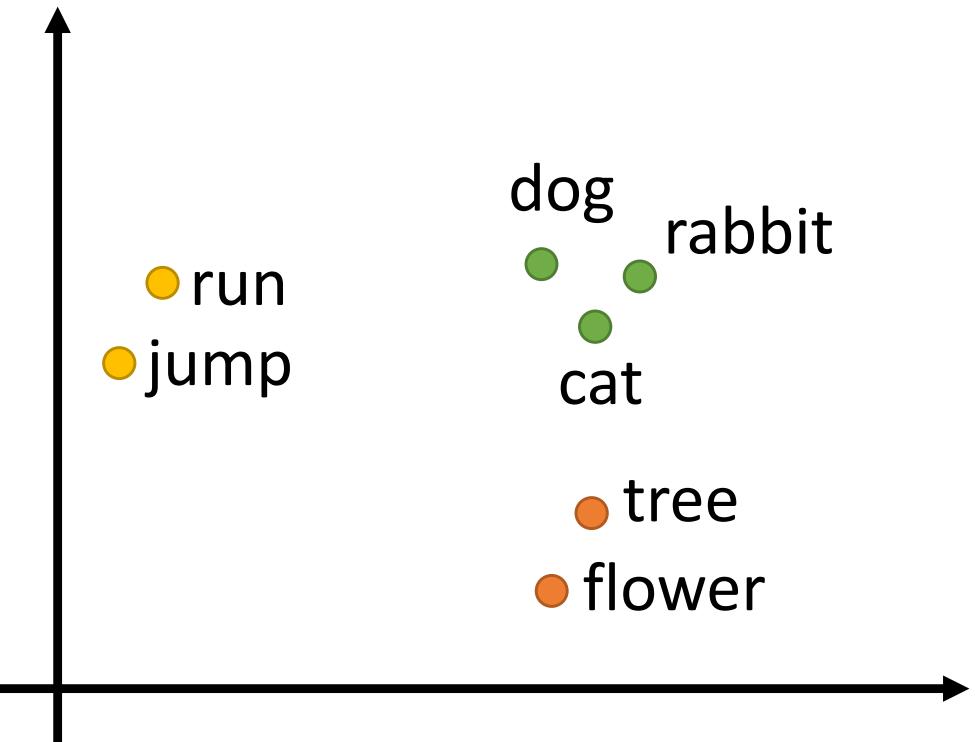
bag = [0 1 0 0 0]

cat = [0 0 1 0 0]

dog = [0 0 0 1 0]

elephant = [0 0 0 0 1]

Word Embedding



To learn more: <https://youtu.be/X7PH3NuYW0Q> (in Mandarin)

Vector Set as Input

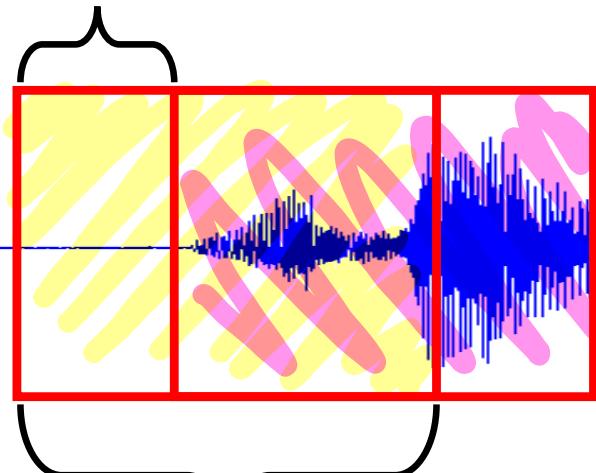
③ 向右移 10 毫升后，再做一次，又得到一个向量

10ms

③ 一路行下去，就得到心的向量

1s → 100 frames

Q: 今日は
10ms
⇒ 十聖半臘
音試験後
後の結果



① 間隔 25ms by window, 把這段

frame

~~400 sample points (16KHz)~~

39-dim MFCC ✓ 的向量

80-dim filter bank output

Vector Set as Input

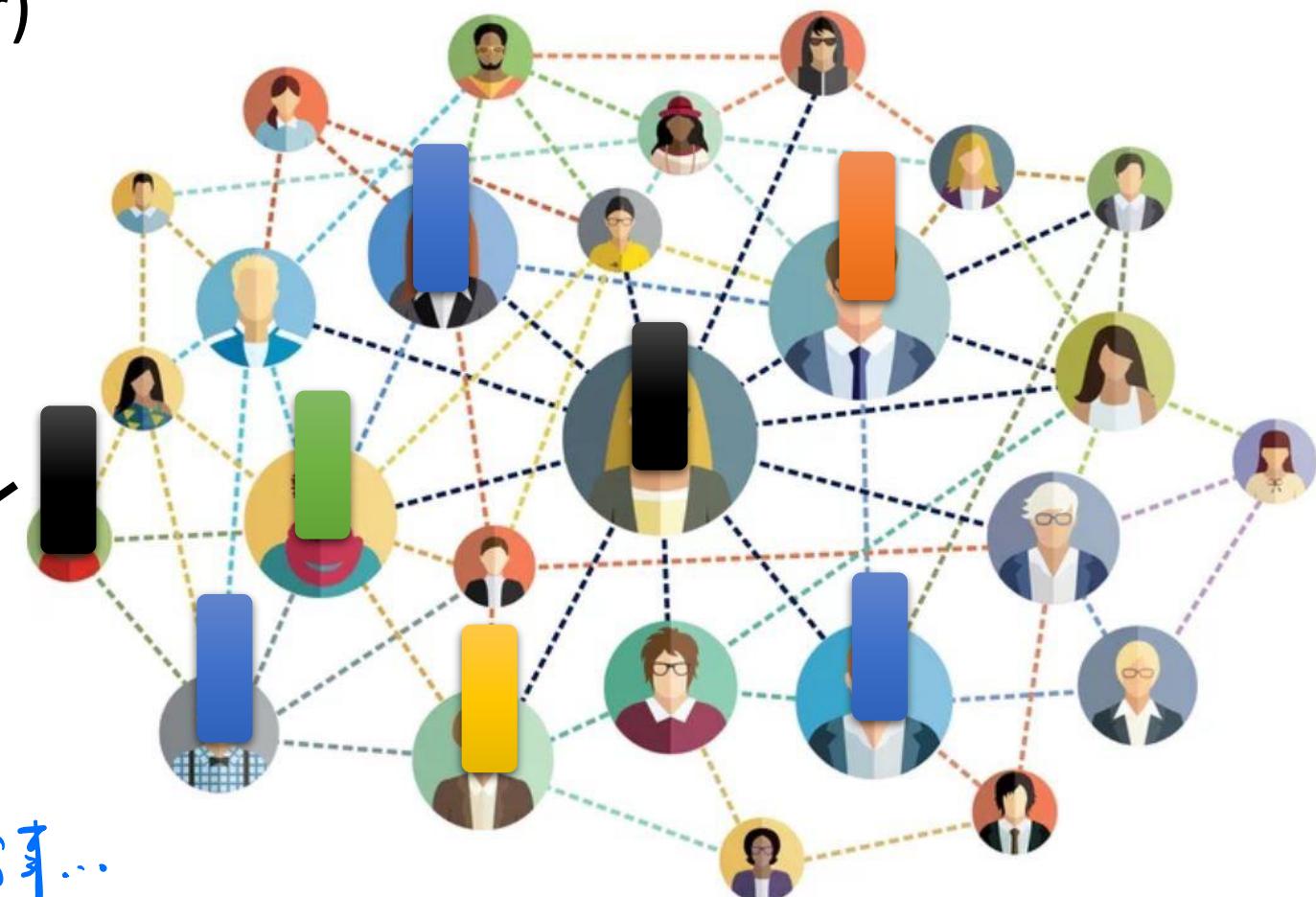
①

- Graph is also a set of vectors (consider each **node** as a **vector**)

②

- Each profile is a vector

例如，小明这个人
他的学历，做过啥...
可转成一个向量



Vector Set as Input

* 分子，可看成一张图(graph)，表格 -> 序列，是一个向量

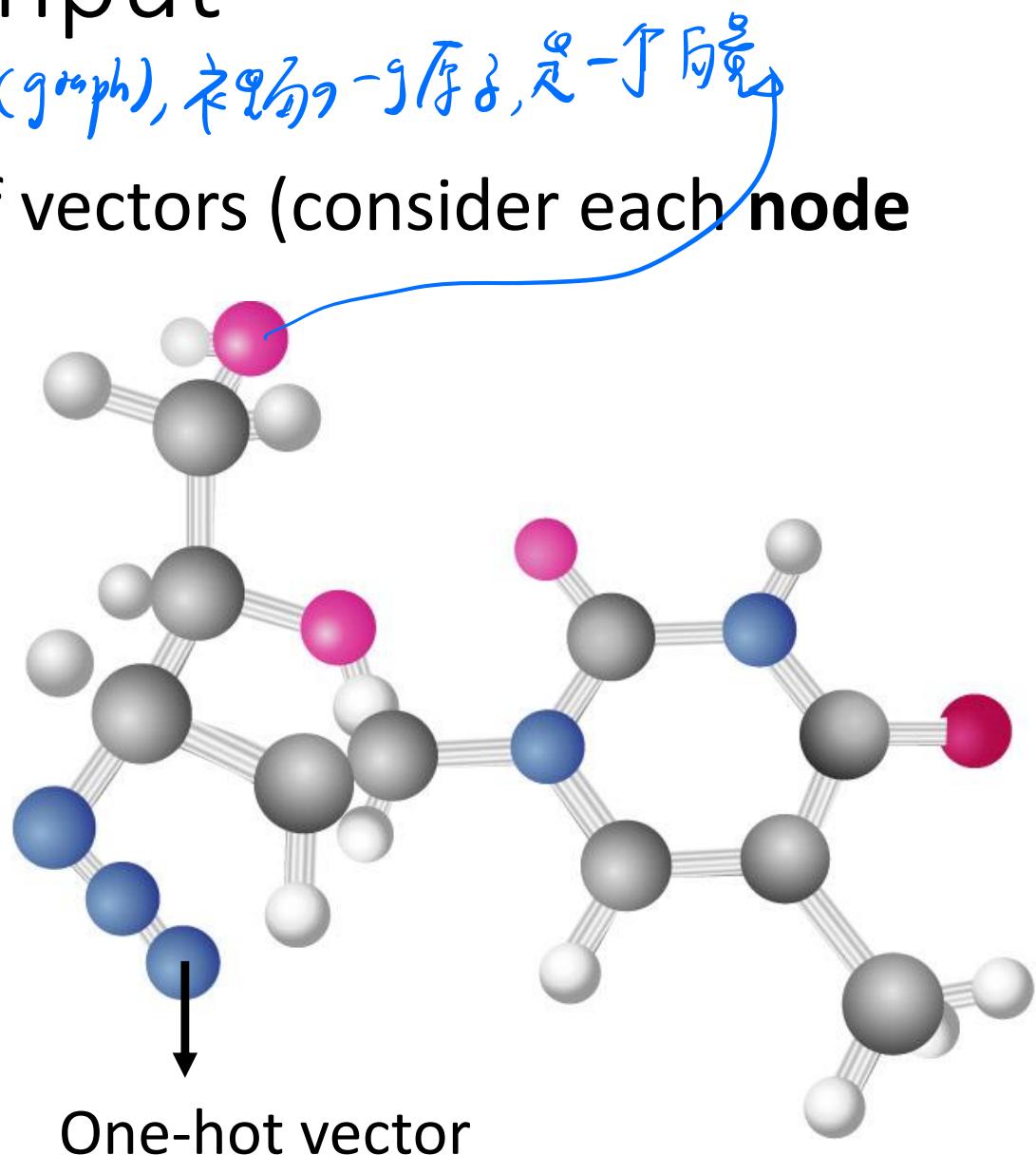
- Graph is also a set of vectors (consider each node as a vector)

$$H = [1 \ 0 \ 0 \ 0 \ 0 \dots]$$

$$C = [0 \ 1 \ 0 \ 0 \ 0 \dots]$$

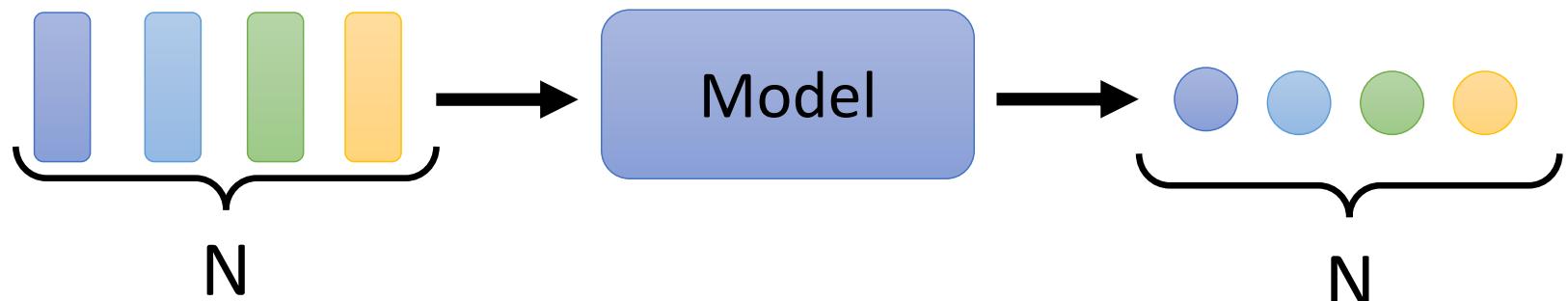
$$O = [0 \ 0 \ 1 \ 0 \ 0 \dots]$$

⋮



What is the output? → 3种可能

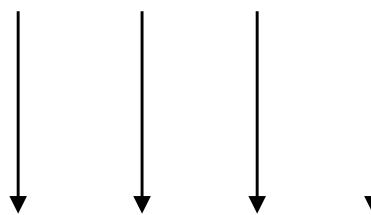
- Each vector has a label.



② 例如，词性标注

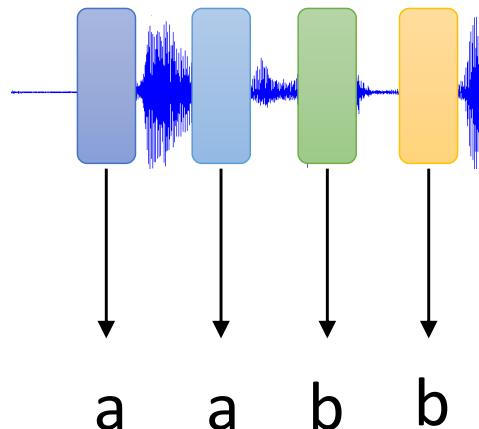
Example Applications

input ⇒ I saw a saw



output ⇒

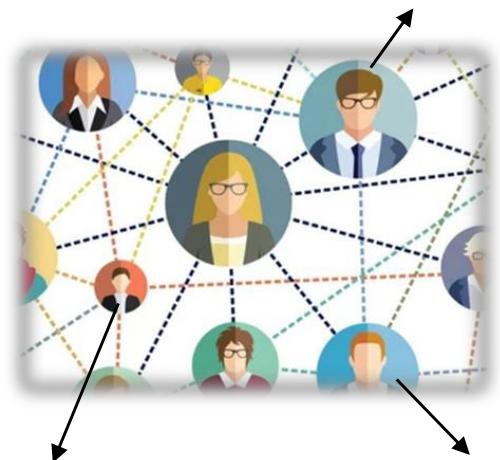
POS tagging



HW2

input ⇒ - 哪管是

not



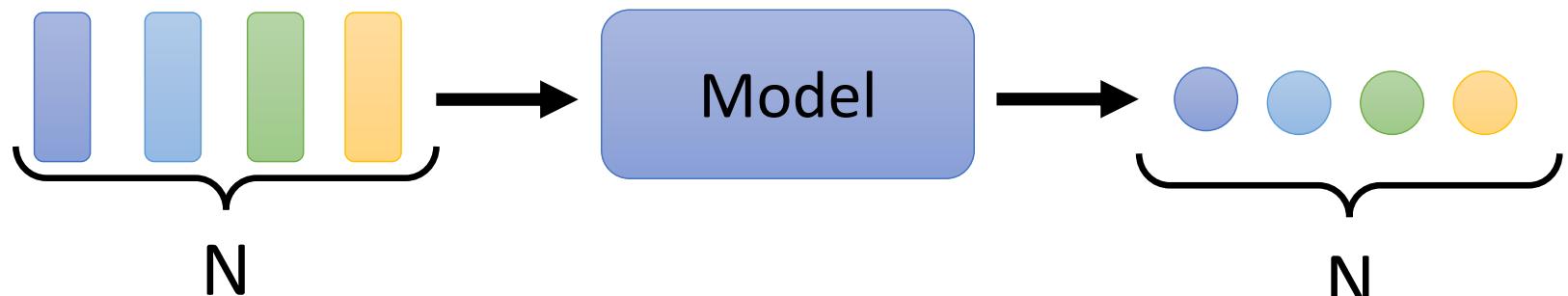
buy

buy

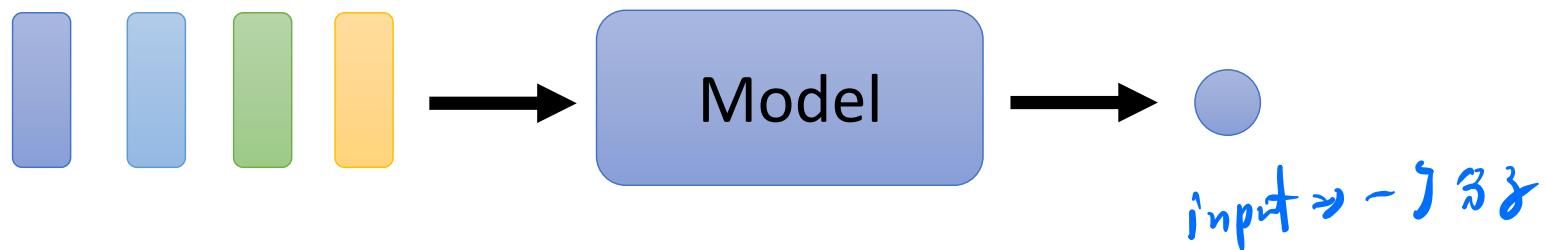
输出 ⇒ 的节点
今不会卖商品

What is the output?

- Each vector has a label.



- The whole sequence has a label. \rightarrow 第二种是整个序列有 label



Example Applications

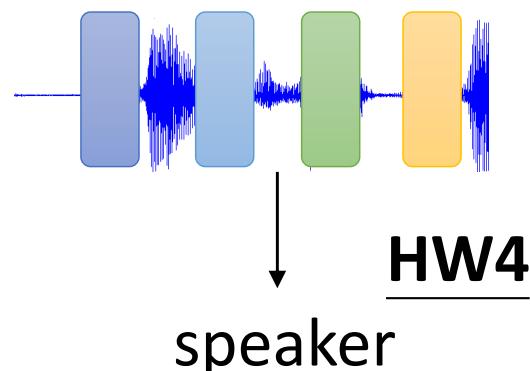
例子

this is good

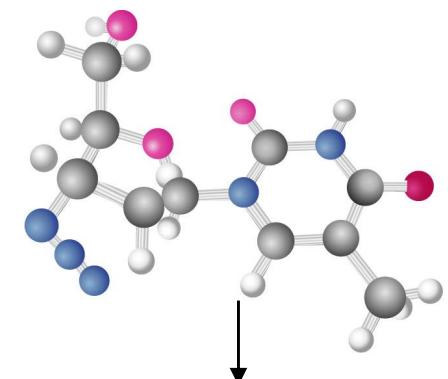
Sentiment analysis

positive

像imdb影评 \Rightarrow 喜欢 or 厌恶



speaker



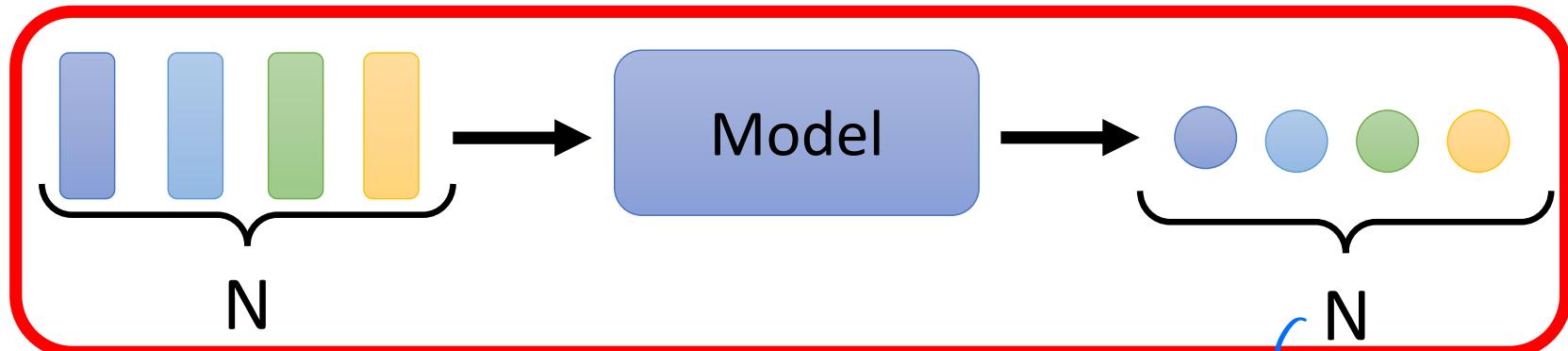
hydrophilicity
output \Rightarrow 有亲水性

What is the output?

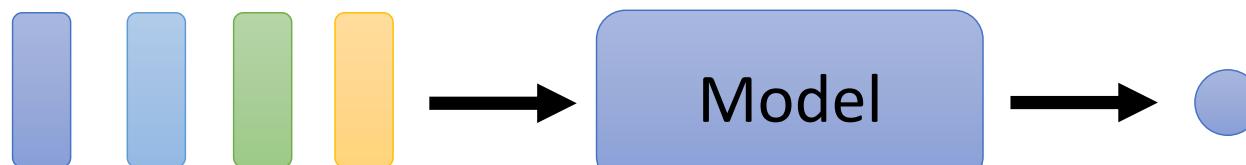
- Each vector has a label.

④

focus of this lecture



- The whole sequence has a label.

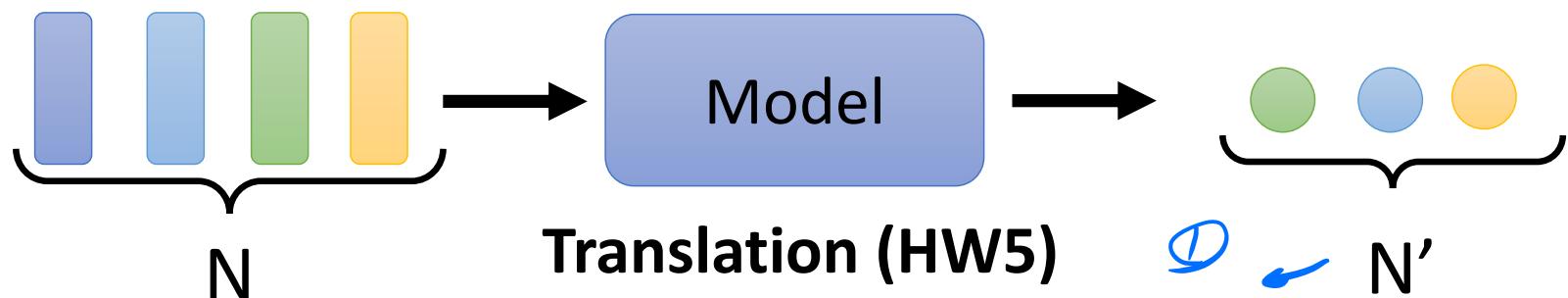


③ 逐種叫 seq2seq

③ 逐種叫 seq2seq

- Model decides the number of labels itself.

seq2seq



Translation (HW5)

① $\leftarrow N'$

第三步输出 N'

Sequence Labeling

Is it possible to consider the context?

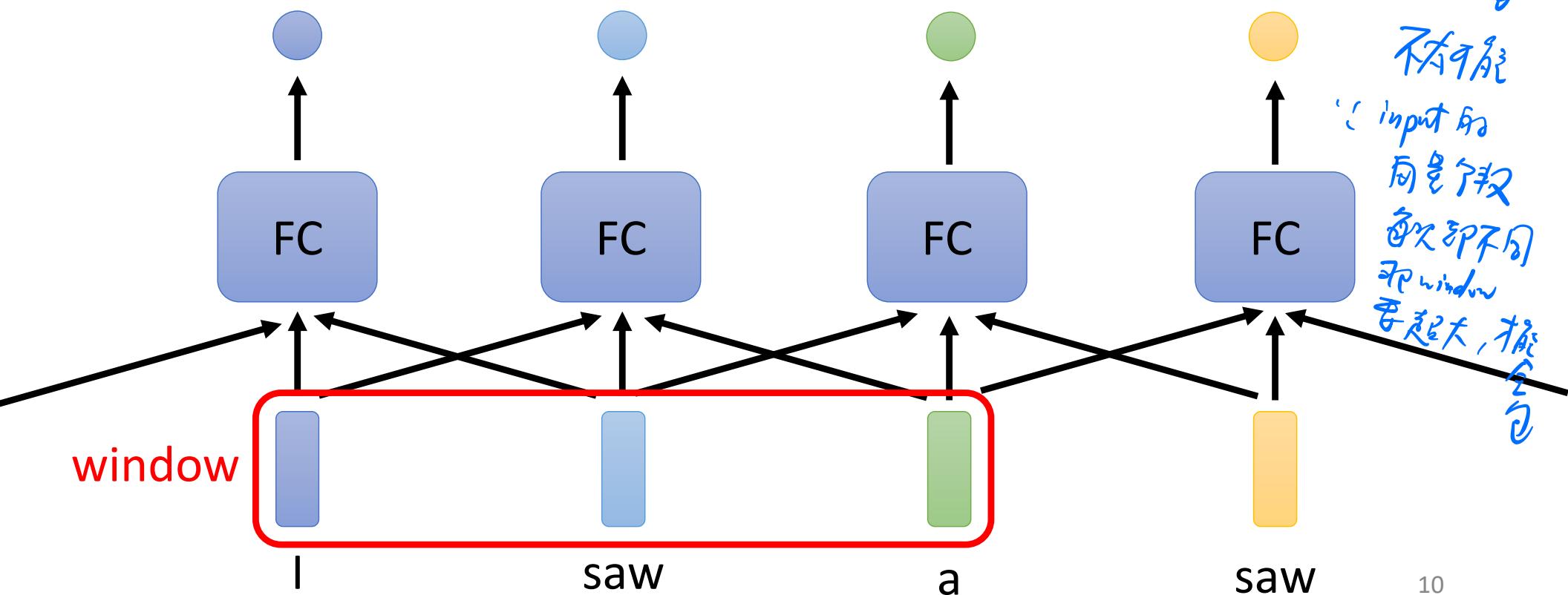


Fully-connected

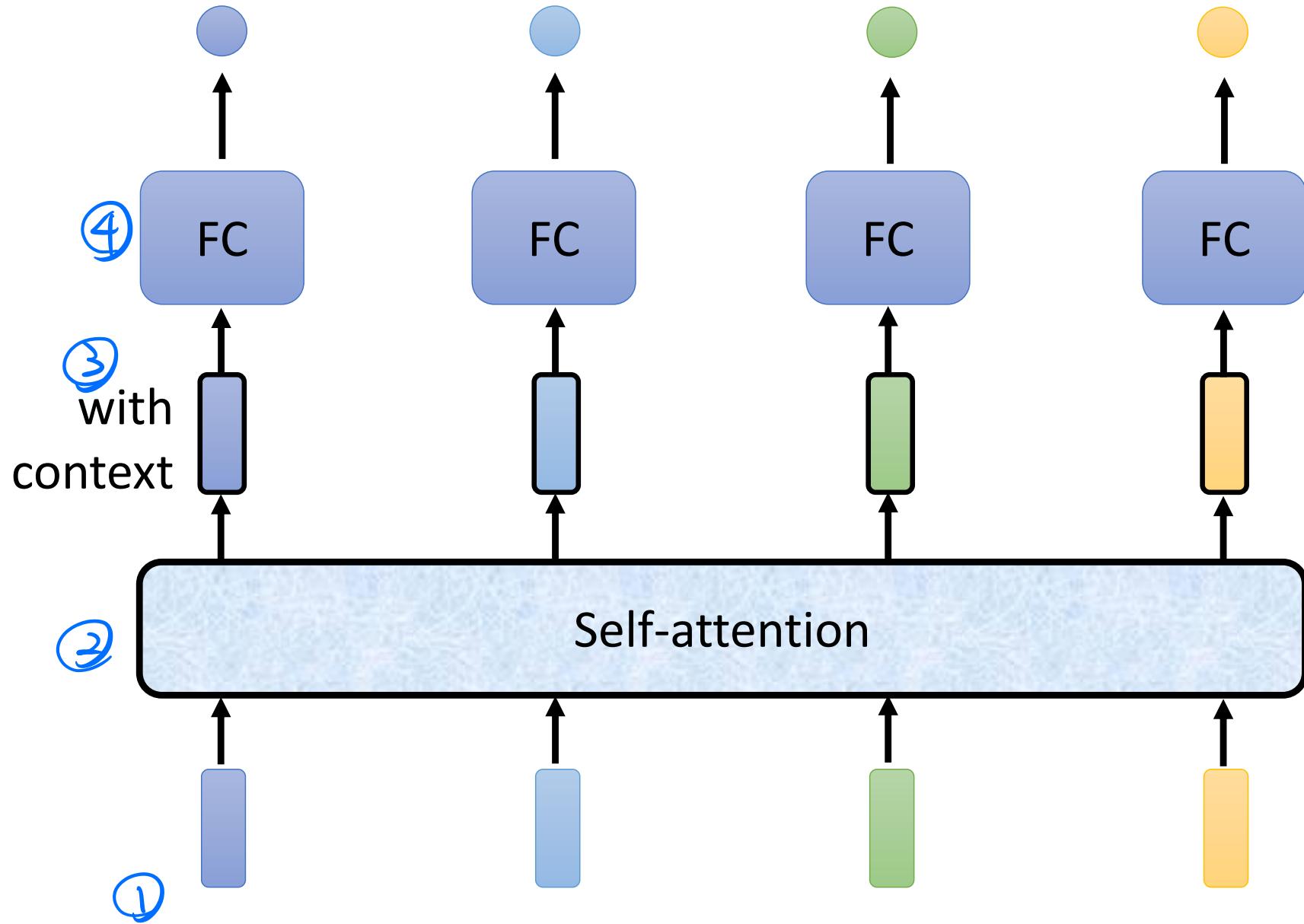
FC can consider the neighbor

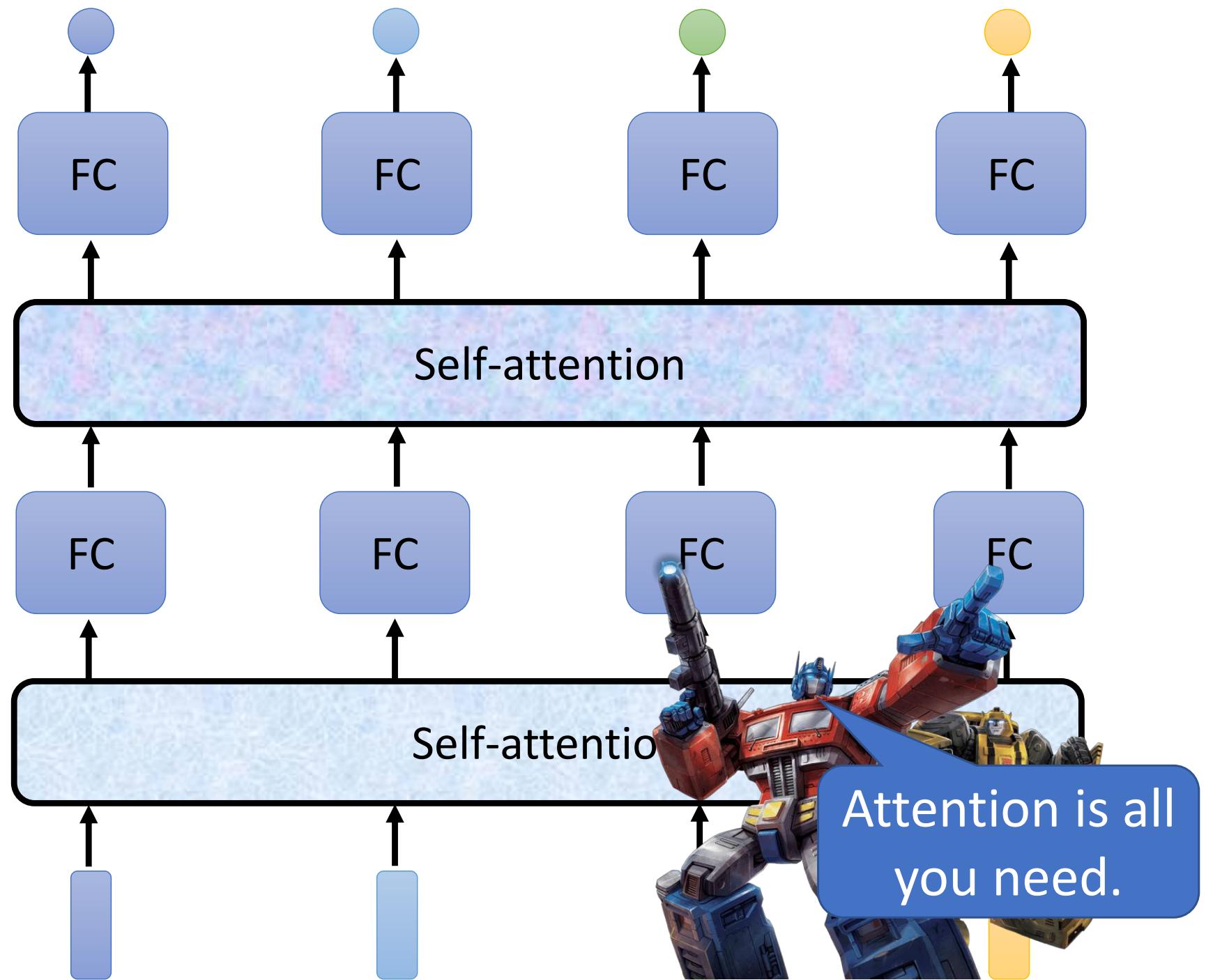
How to consider the whole sequence?

a window covers the whole sequence?

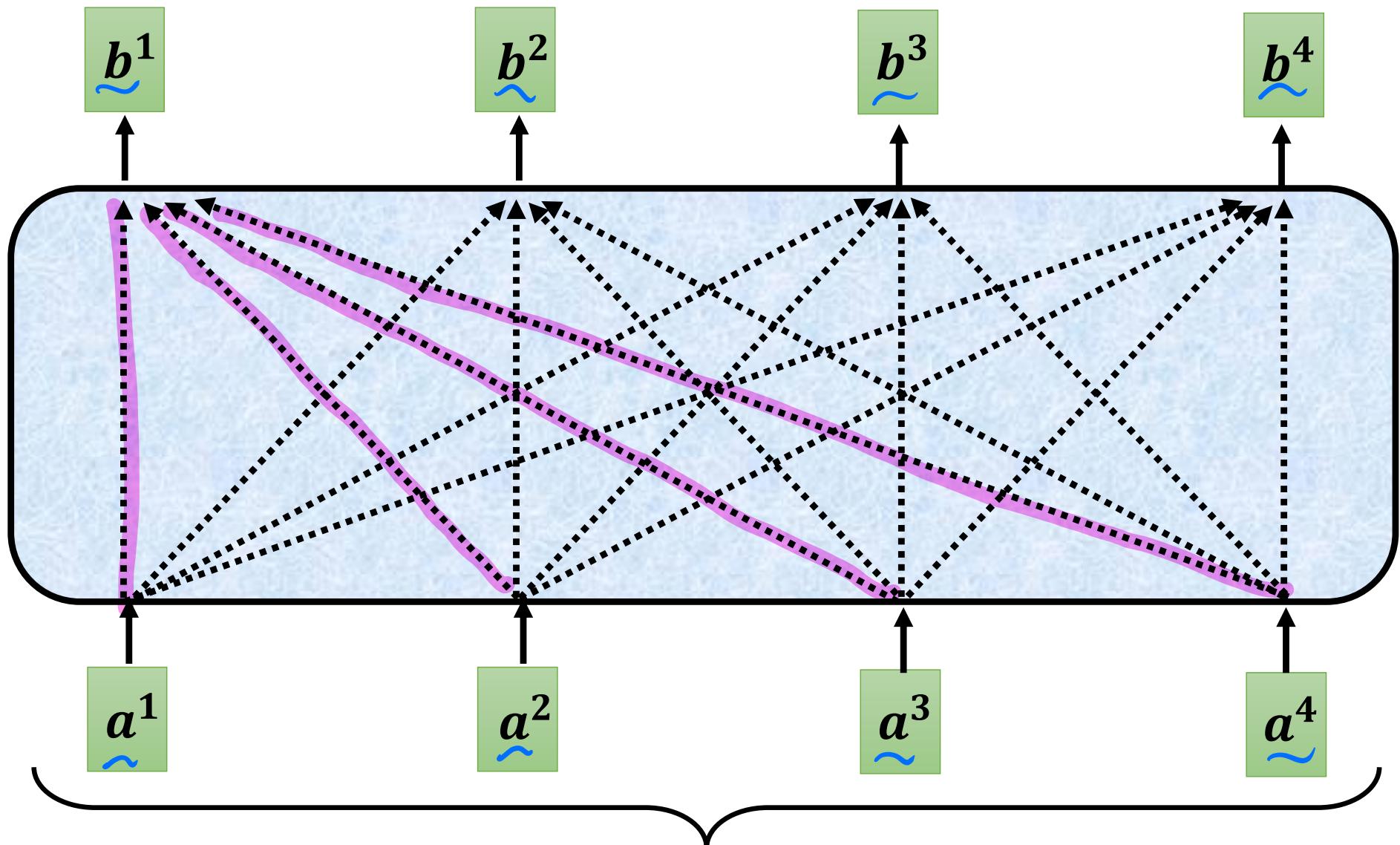


Self-attention



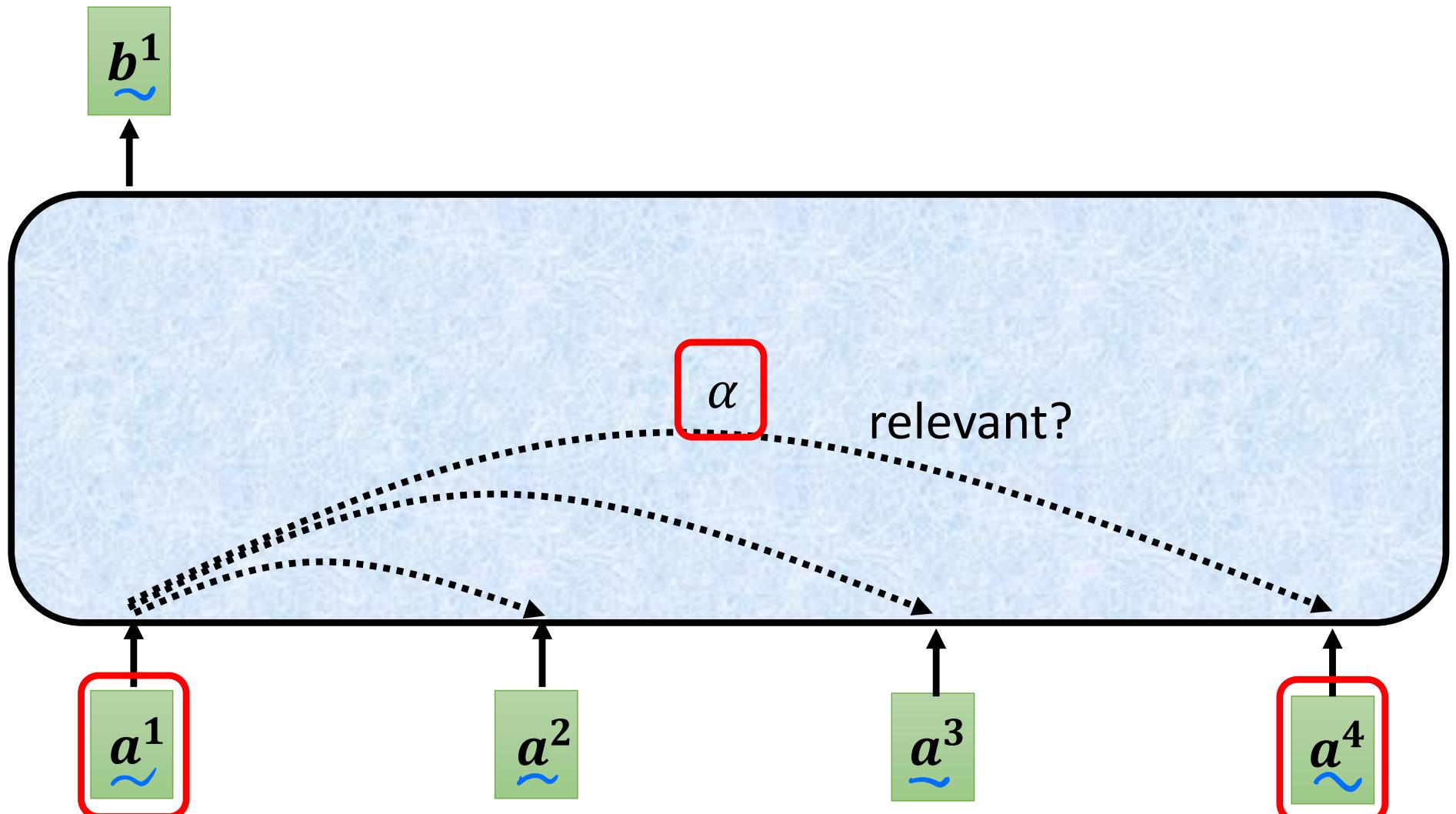


Self-attention



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

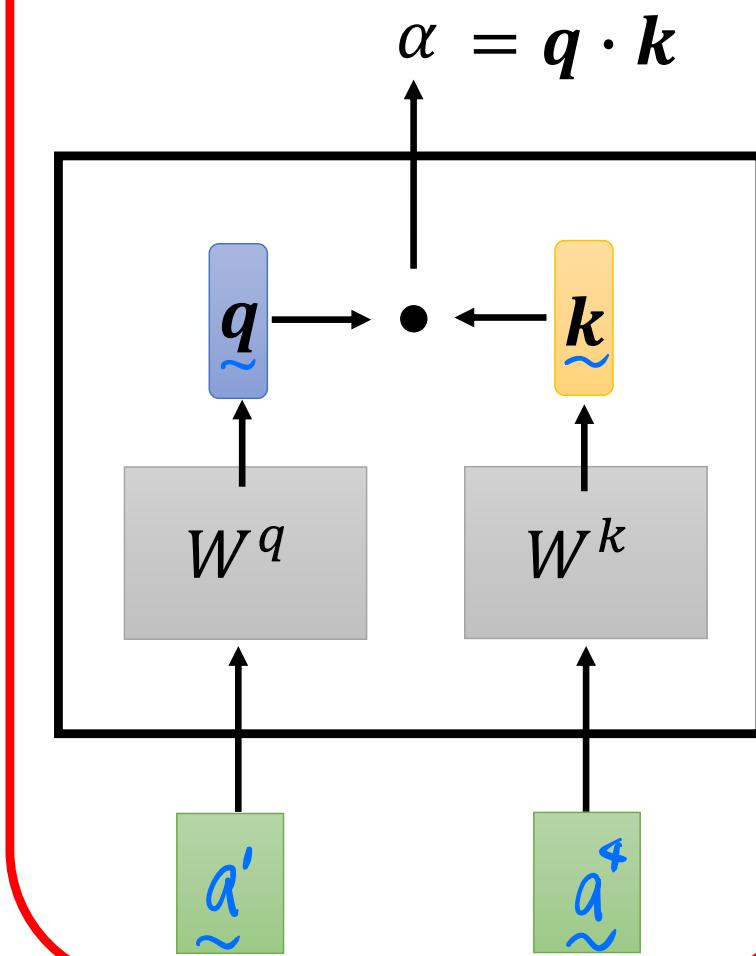
Self-attention



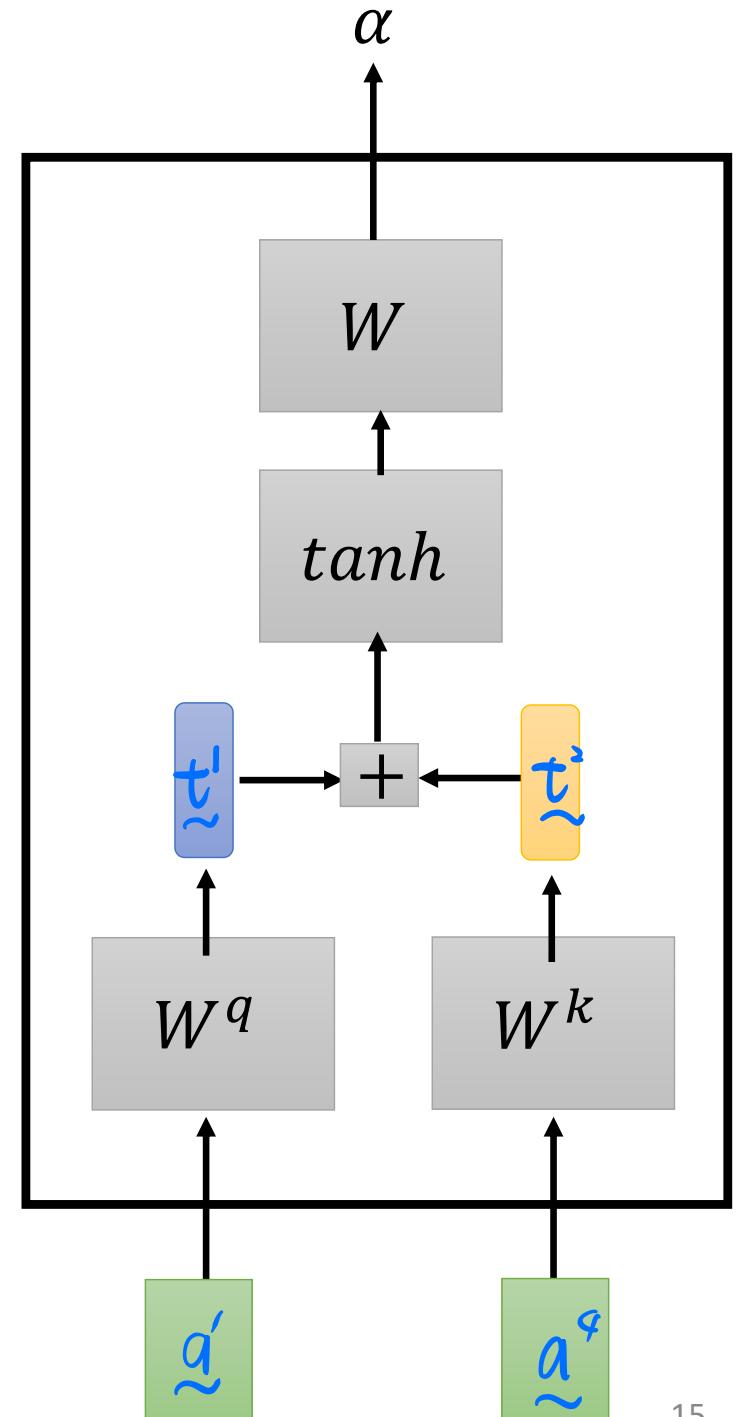
Find the relevant vectors in a sequence

Self-attention

Dot-product
内積



$\alpha =$
Additive

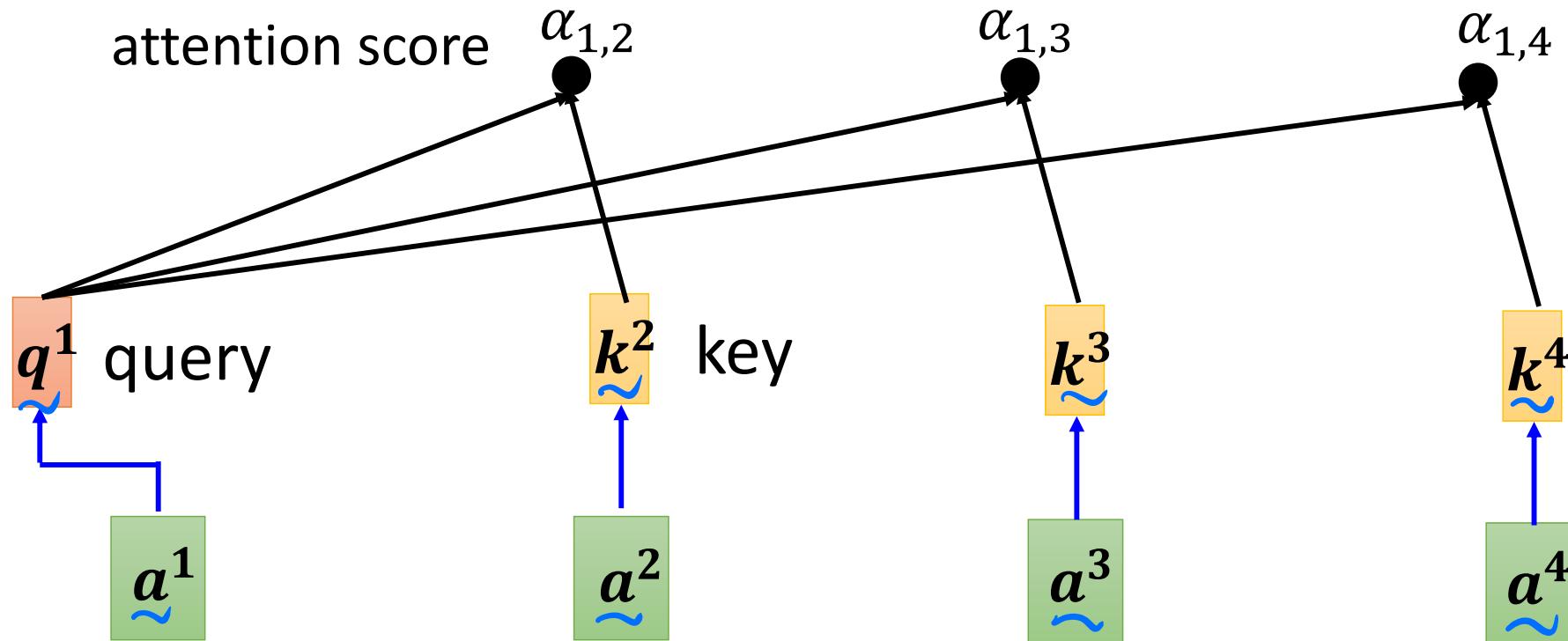


Self-attention

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

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

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



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

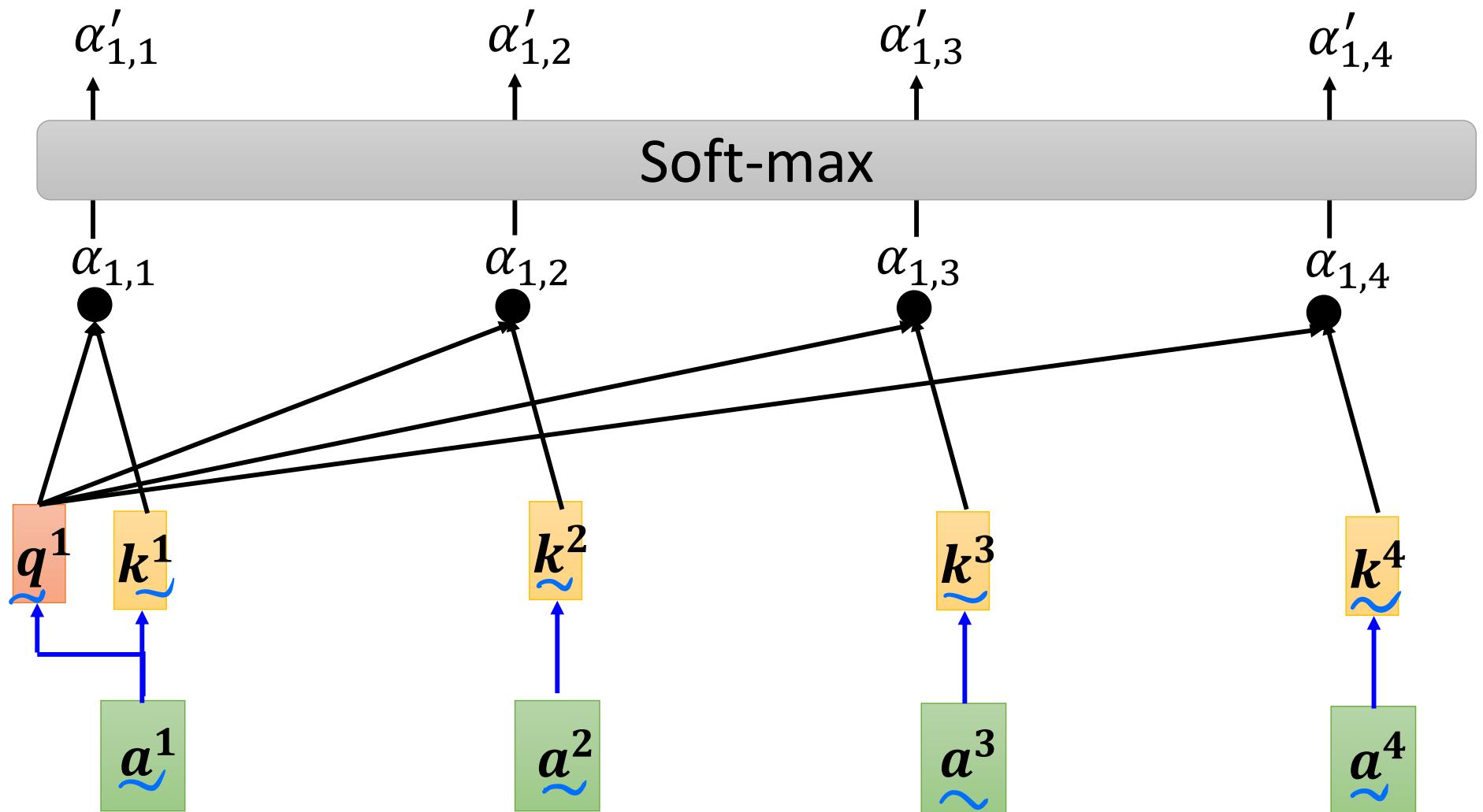
$$k^2 = W^k a^2$$

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

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

Self-attention

$$\alpha'_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



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

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

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

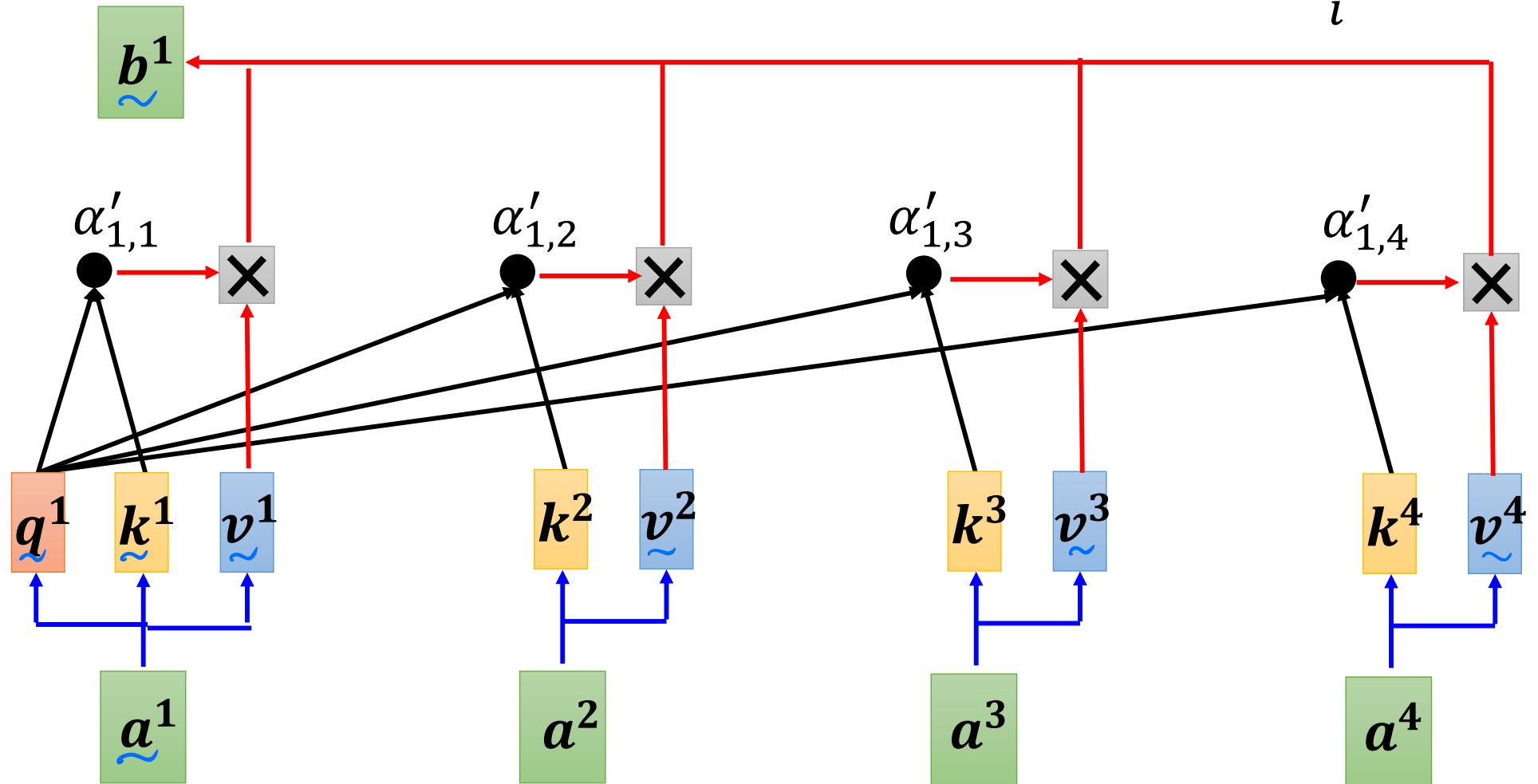
$$k^4 = W^k a^4$$

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

Self-attention

Extract information based
on attention scores

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



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

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

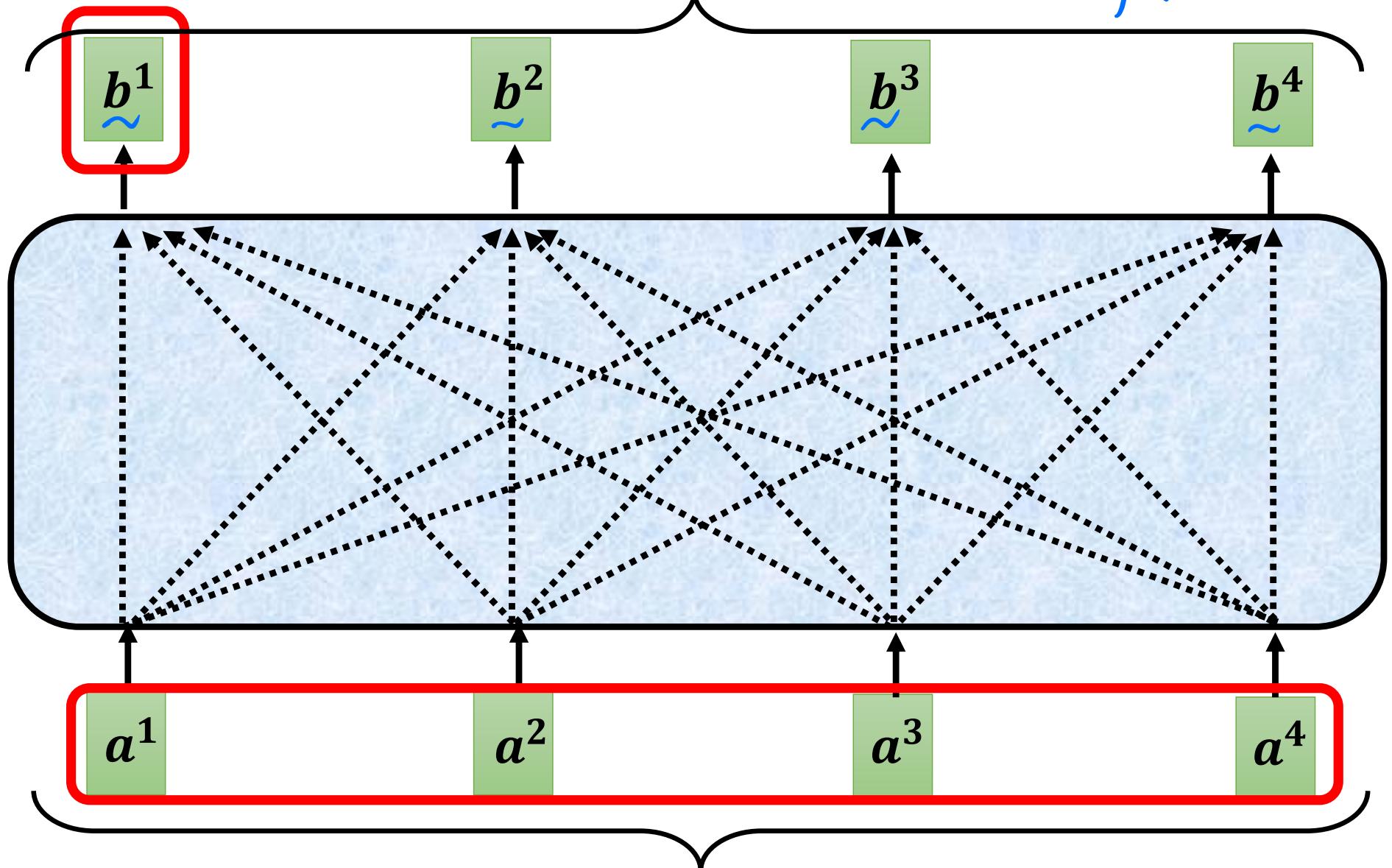
$$v^3 = W^v a^3$$

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

Self-attention

parallel

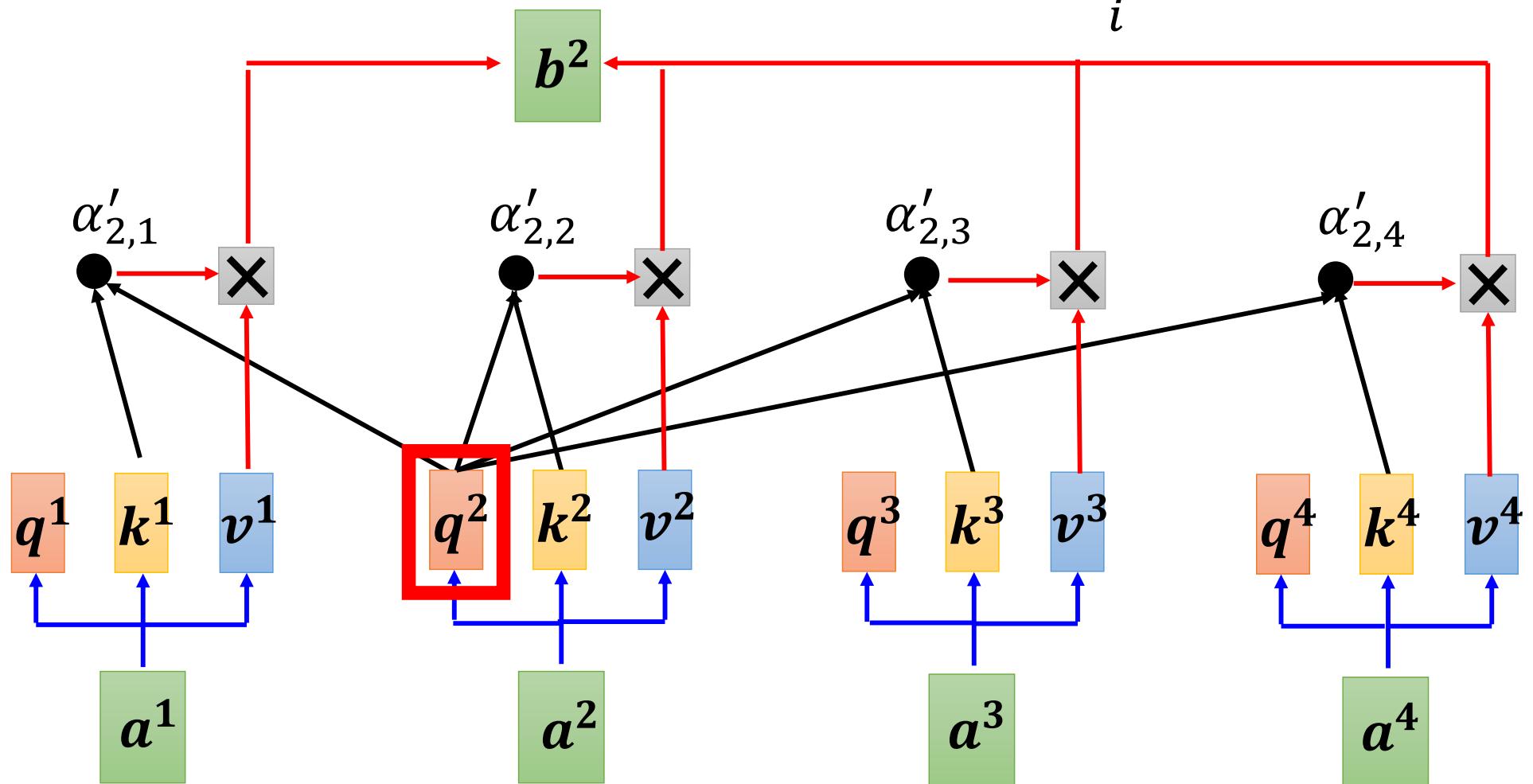
→ 算 b^1 , 不用等 b^2 等
他們的計算是獨立的



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

Self-attention

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



Self-attention

$$\tilde{q}^i = W^q \tilde{a}^i$$

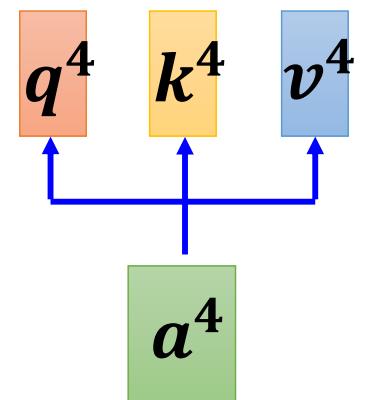
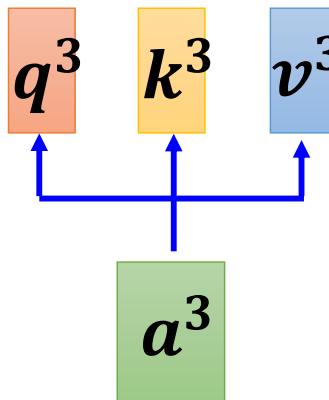
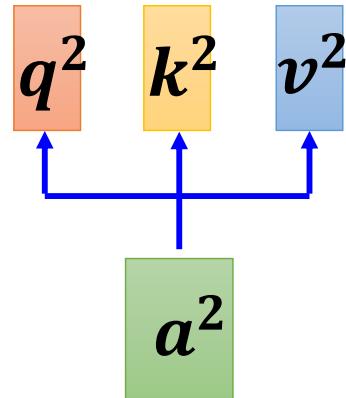
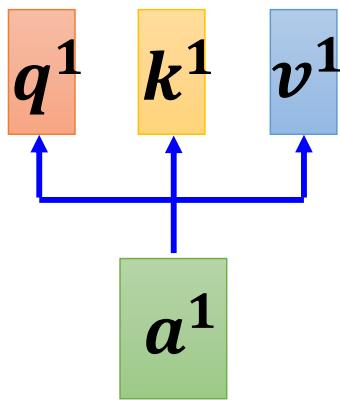
$$q^1 | q^2 | q^3 | q^4 = \begin{matrix} W^q \\ Q \end{matrix} \quad a^1 | a^2 | a^3 | a^4 = \begin{matrix} I \\ \text{input} \end{matrix}$$

$$\tilde{k}^i = W^k \tilde{a}^i$$

$$k^1 | k^2 | k^3 | k^4 = \begin{matrix} W^k \\ K \end{matrix} \quad a^1 | a^2 | a^3 | a^4 = \begin{matrix} I \\ I \end{matrix}$$

$$\tilde{v}^i = W^v \tilde{a}^i$$

$$v^1 | v^2 | v^3 | v^4 = \begin{matrix} W^v \\ V \end{matrix} \quad a^1 | a^2 | a^3 | a^4 = \begin{matrix} I \\ I \end{matrix}$$

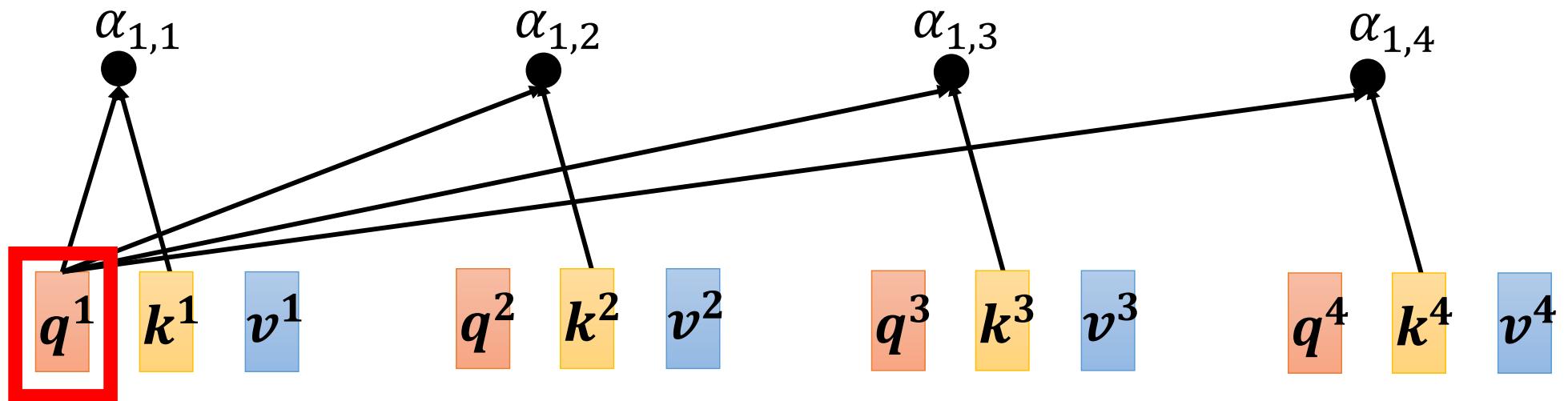


Self-attention

$$\alpha_{1,1} = \begin{matrix} k^1 \\ q^1 \end{matrix} \quad \alpha_{1,2} = \begin{matrix} k^2 \\ q^1 \end{matrix}$$

$$\alpha_{1,3} = \begin{matrix} k^3 \\ q^1 \end{matrix} \quad \alpha_{1,4} = \begin{matrix} k^4 \\ q^1 \end{matrix}$$

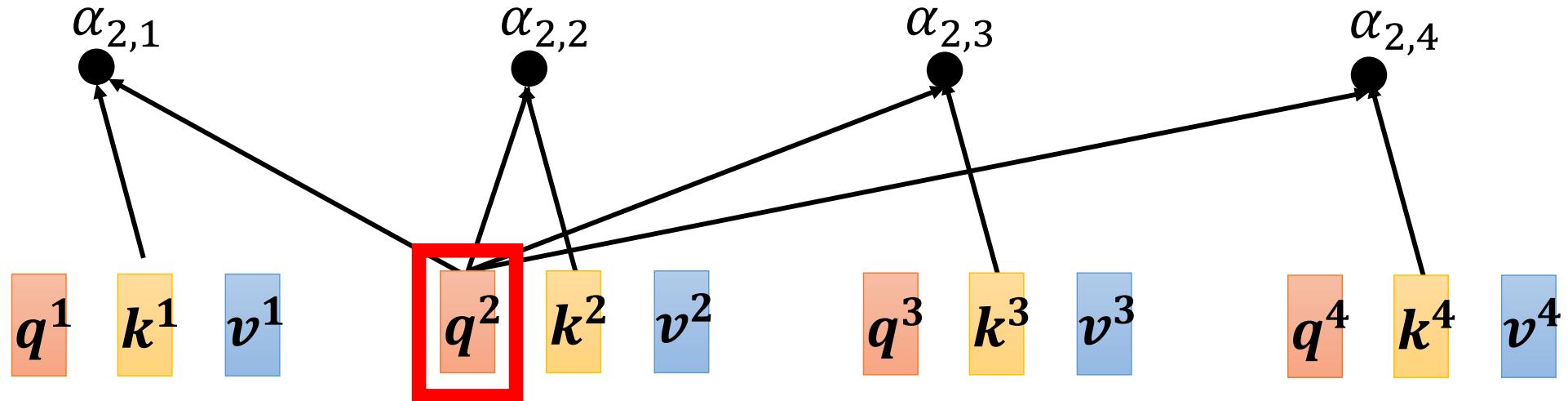
$$\begin{matrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{matrix} = \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} \quad \begin{matrix} q^1 \\ q^1 \\ q^1 \\ q^1 \end{matrix}$$



Self-attention

$$\alpha_{1,1} = \begin{matrix} k^1 \\ q^1 \end{matrix} \quad \alpha_{1,2} = \begin{matrix} k^2 \\ q^1 \end{matrix} \quad \alpha_{1,3} = \begin{matrix} k^3 \\ q^1 \end{matrix} \quad \alpha_{1,4} = \begin{matrix} k^4 \\ q^1 \end{matrix}$$

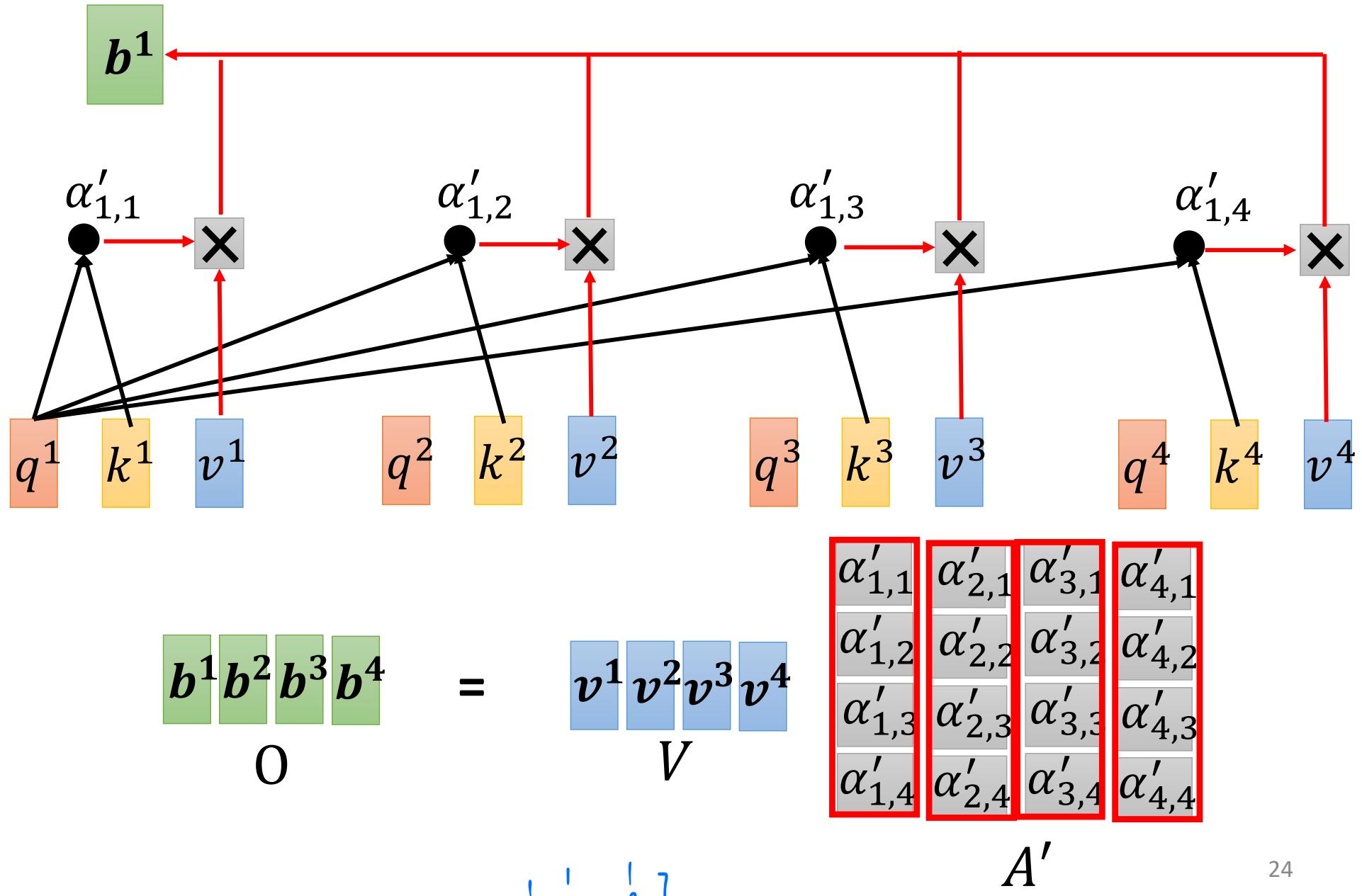
$$\alpha_{1,1} \quad \alpha_{1,2} \quad \alpha_{1,3} \quad \alpha_{1,4} = \begin{matrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{matrix} = \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} \quad \begin{matrix} q^1 \\ q^1 \\ q^1 \\ q^1 \end{matrix}$$



$$\begin{matrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{matrix} \xleftarrow{\text{softmax}} \begin{matrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{matrix} = \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} \quad Q = \begin{matrix} q^1 \\ q^2 \\ q^3 \\ q^4 \end{matrix}$$

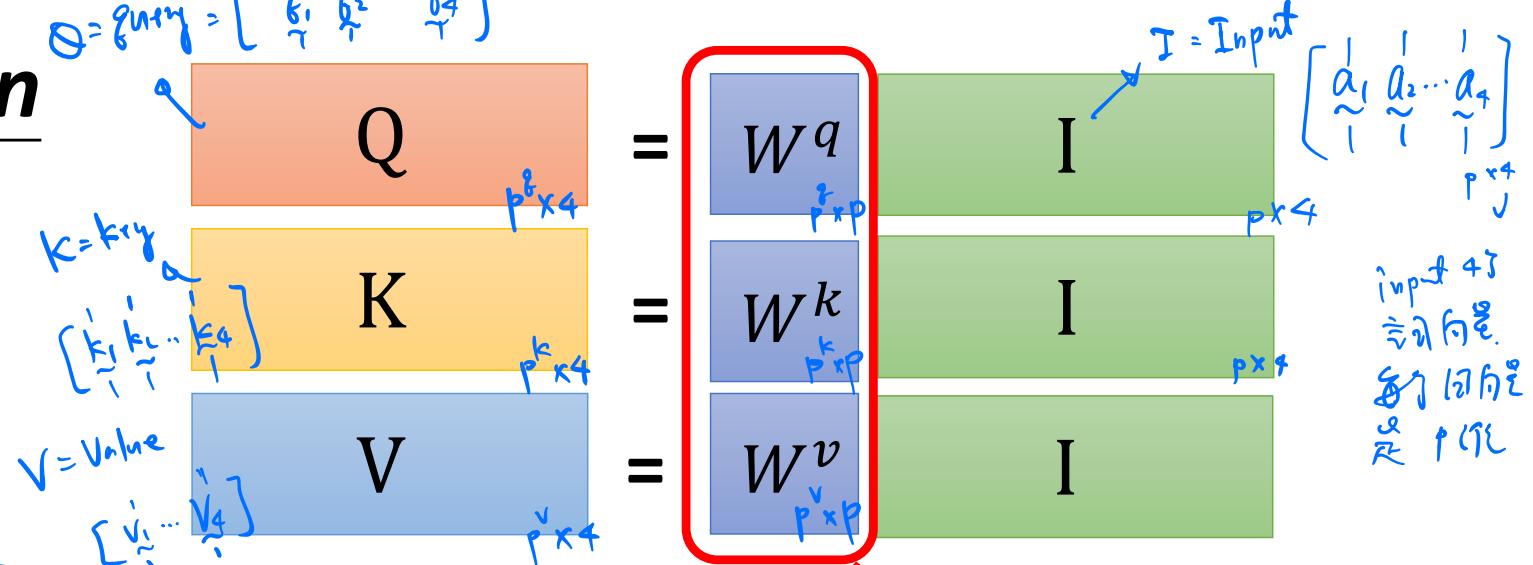
A' softmax A K^T Q

Self-attention

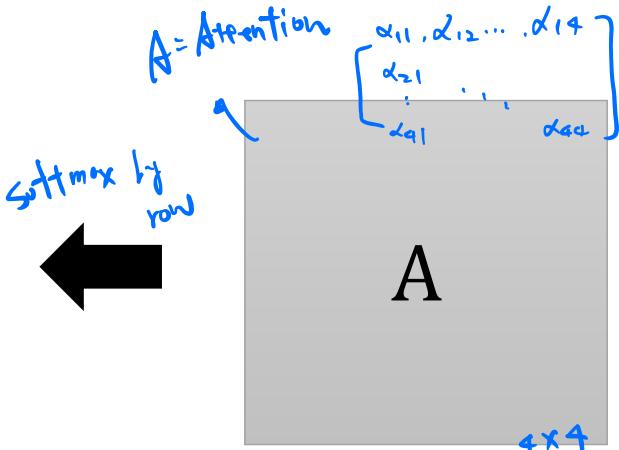
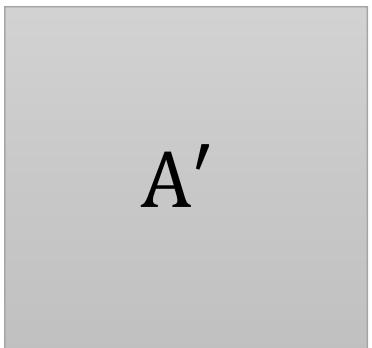


Self-attention

Step 1 \Rightarrow



Step 2



$$A' = Q K^T$$

$$Q = p \times 4$$

$$K^T = 4 \times p$$

$$A' = \begin{bmatrix} k_1^T & k_2^T & \dots & k_4^T \end{bmatrix} \begin{bmatrix} q_1 & q_2 & \dots & q_4 \end{bmatrix}^T = \begin{bmatrix} k_1^T q_1 & k_2^T q_2 & \dots & k_4^T q_4 \end{bmatrix}$$

Attention Matrix

Step 3:

$$O = V$$

$$O = p \times 4$$

$$V = p \times 4$$

$$A' = 4 \times 4$$

O is a $p \times 4$ matrix.

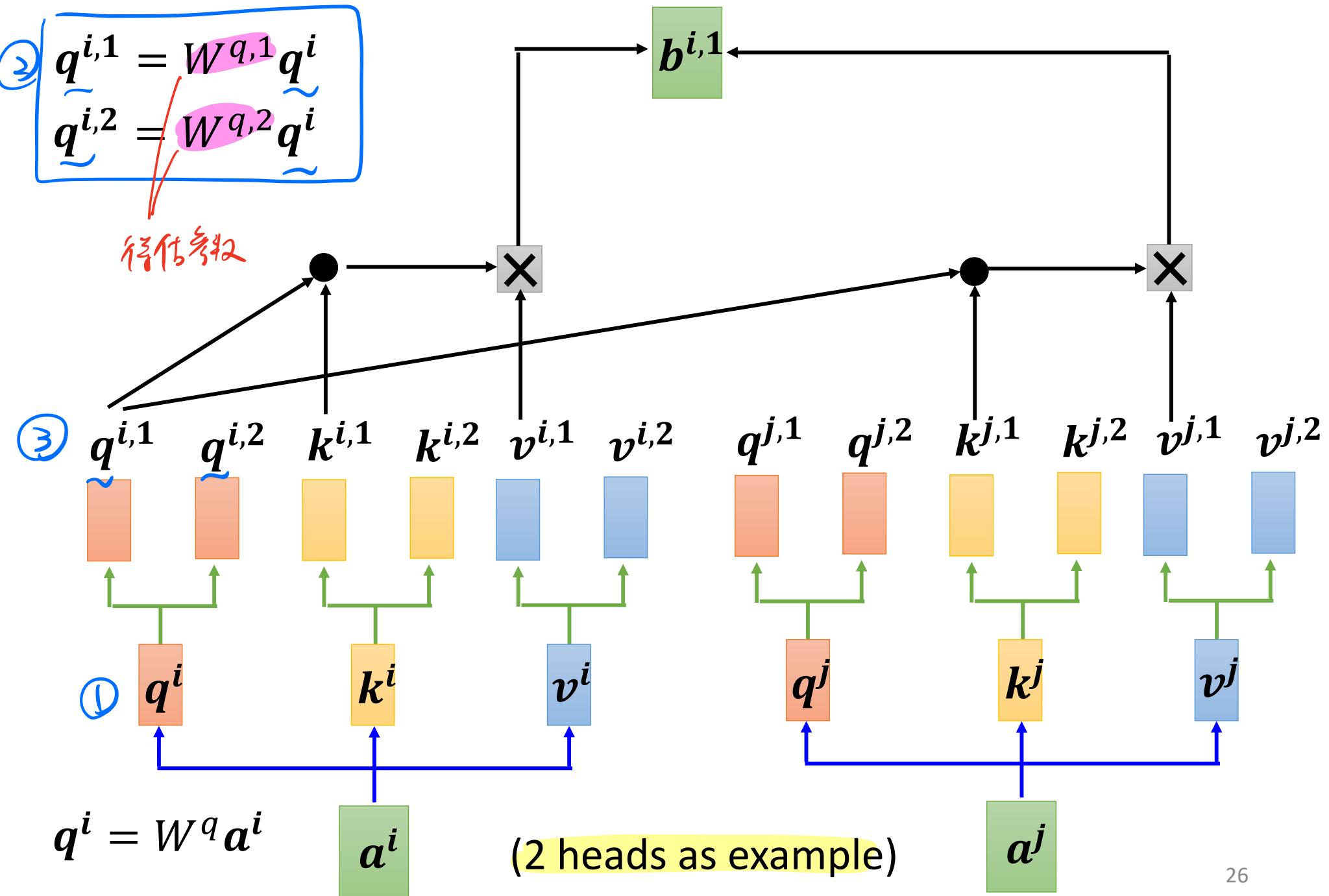
$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ v_1 & v_2 & \dots & v_4 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ \alpha_1 & \alpha_2 & \dots & \alpha_4 \end{bmatrix} = \begin{bmatrix} 0_1 & 0_2 & 0_3 & 0_4 \end{bmatrix}$$

Multi-head Self-attention

Different types of relevance

$$\alpha_{i,1}v_1 + \alpha_{i,2}v_2 + \dots + \alpha_{i,q}v_q$$

用 attention score 权重

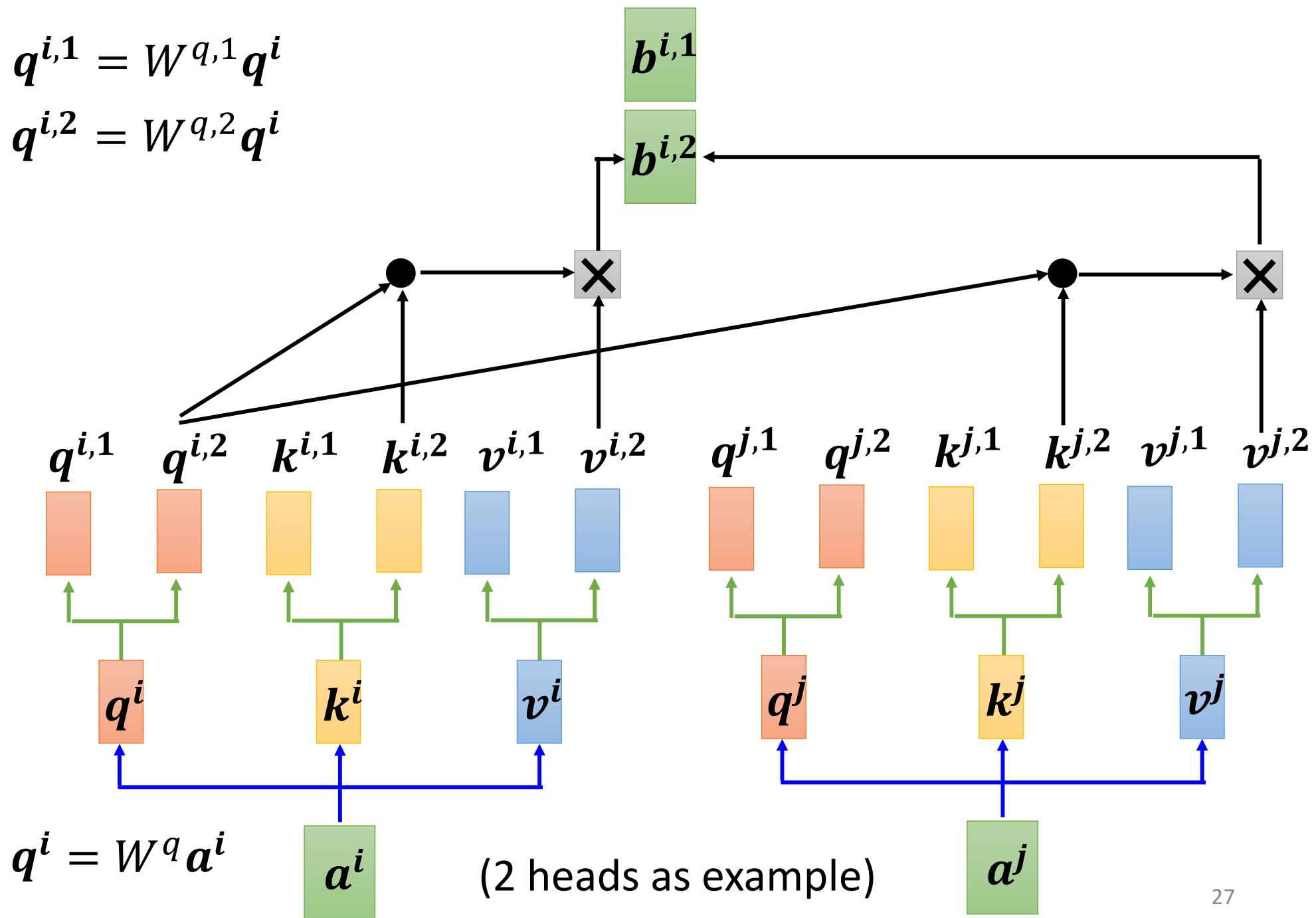


Multi-head Self-attention

Different types of relevance

$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$



$$q^i = W^q a^i$$

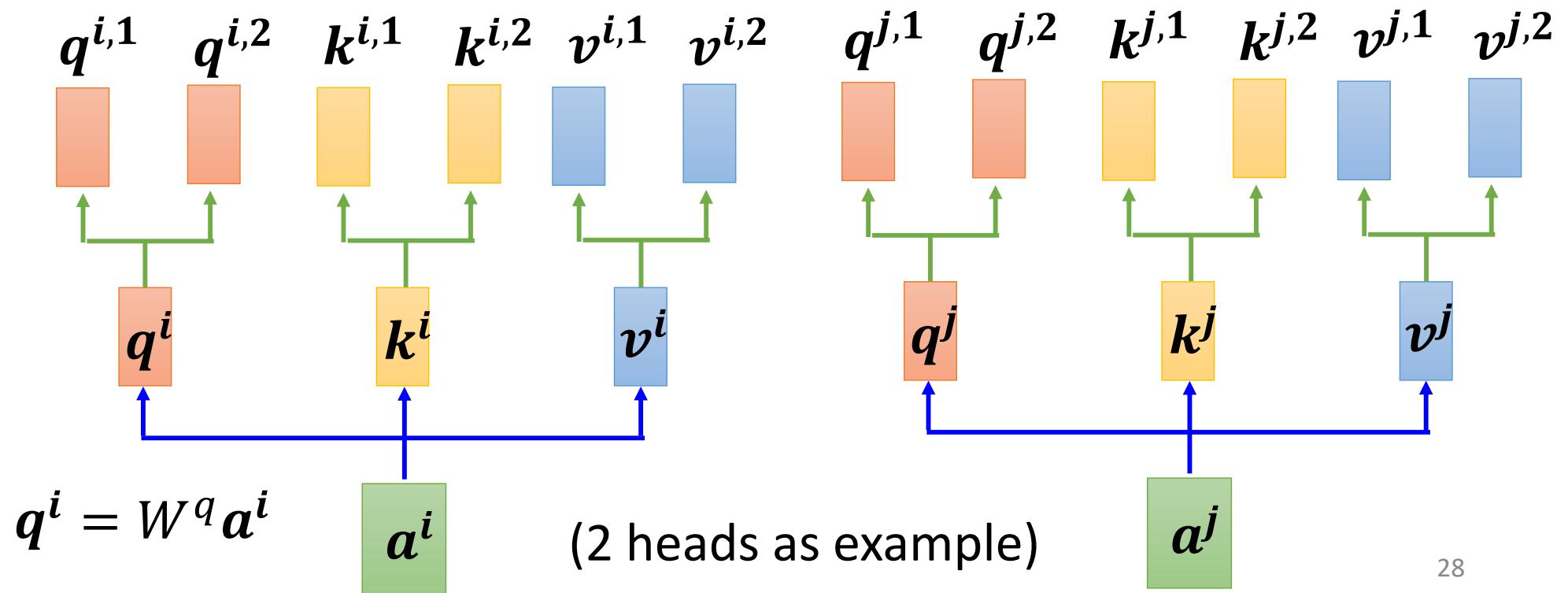
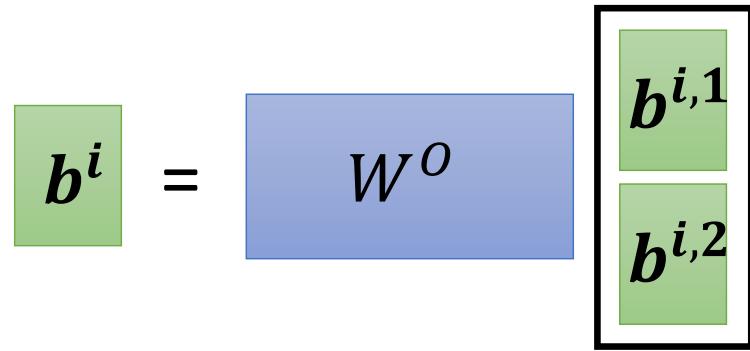
a^i

(2 heads as example)

a^j

Multi-head Self-attention

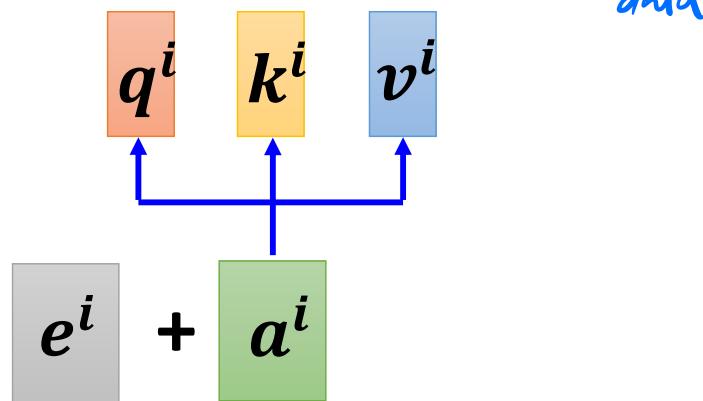
Different types of relevance



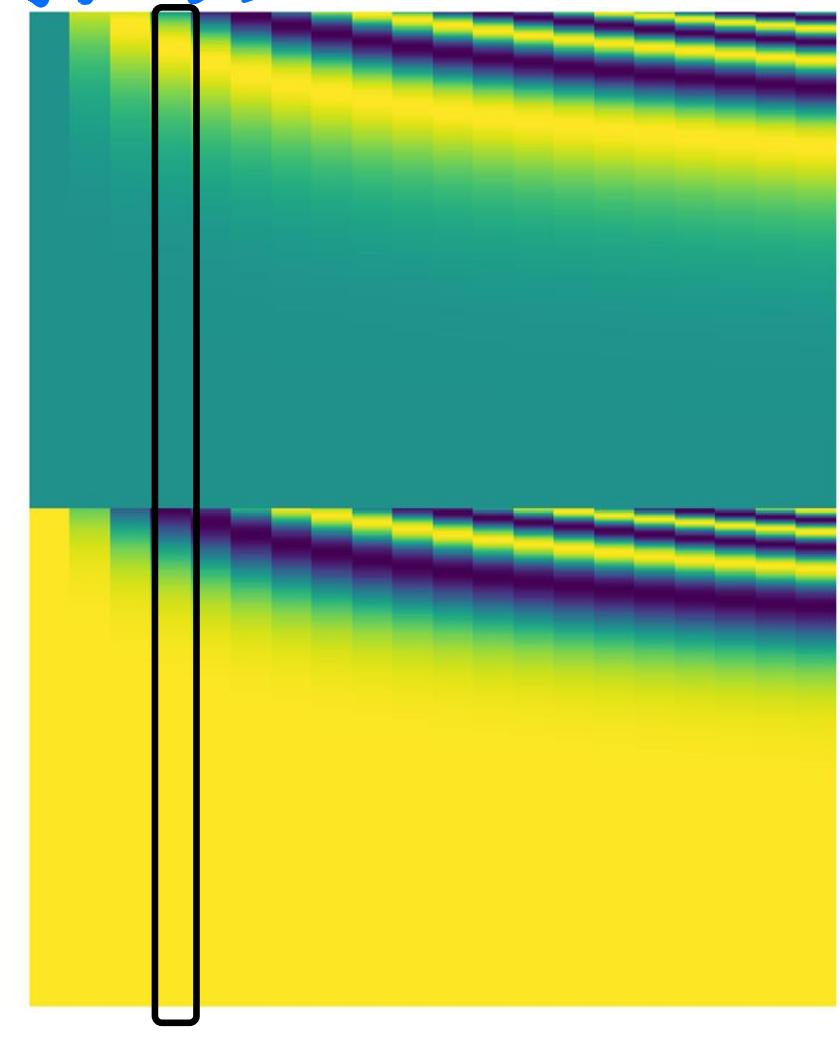
Positional Encoding

①

- No position information in self-attention.
- Each position has a unique positional vector e^i $\xrightarrow{i\text{ 位置}}$
- hand-crafted \rightarrow 人工設計
- learned from data \rightarrow $e^i \xrightarrow{\text{learned from data}}$



Each column represents a positional vector e^i

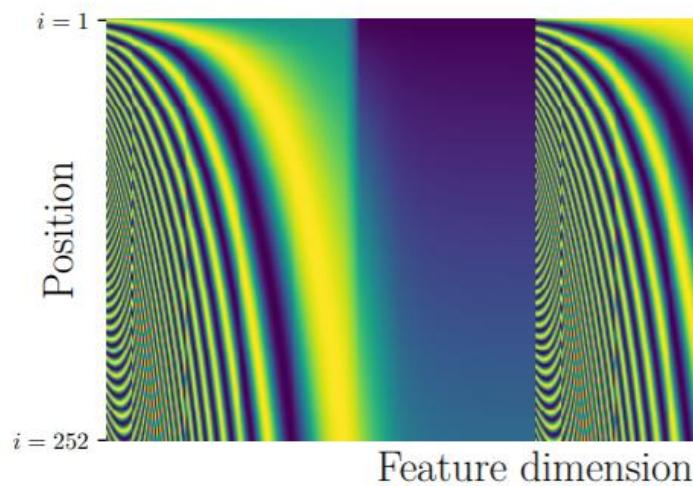


[https://arxiv.org/abs/
2003.09229](https://arxiv.org/abs/2003.09229)

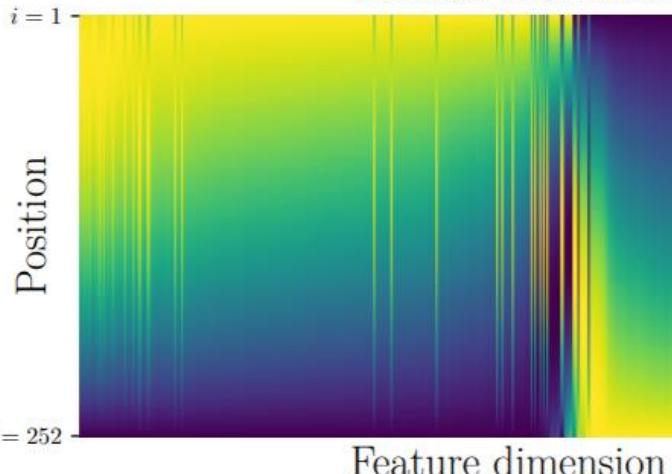
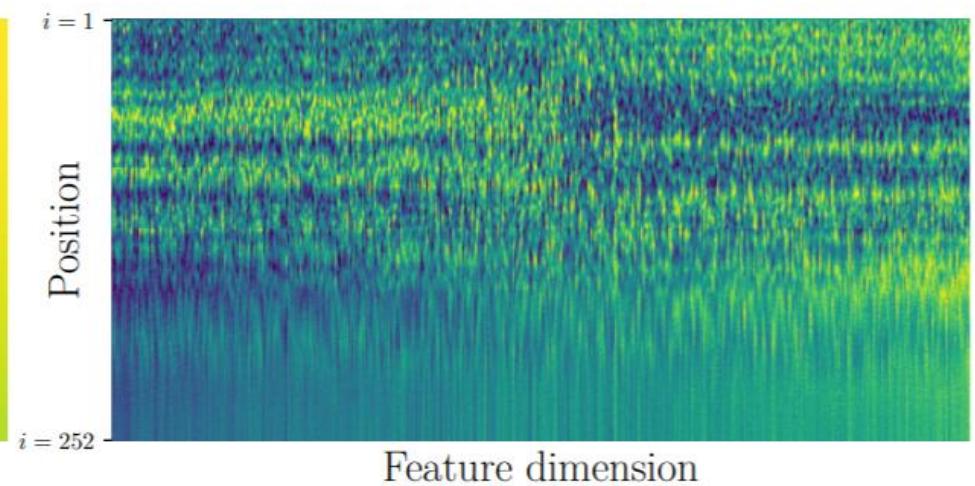
Table 1. Comparing position representation methods

Methods	Inductive	Data-Driven	Parameter Efficient
Sinusoidal (Vaswani et al., 2017)	✓	✗	✓
Embedding (Devlin et al., 2018)	✗	✓	✗
Relative (Shaw et al., 2018)	✗	✓	✓
This paper	✓	✓	✓

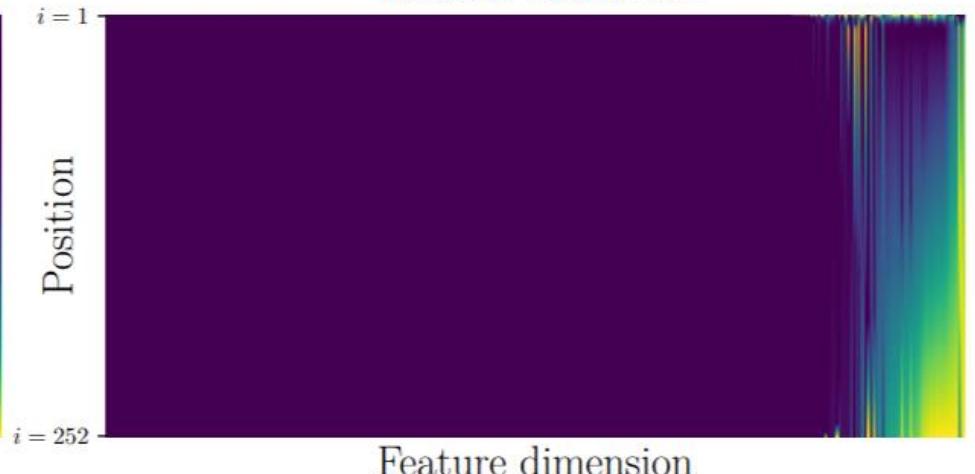
(a) Sinusoidal



(b) Position embedding



(c) FLOATER



(d) RNN

Many applications ...



Transformer

<https://arxiv.org/abs/1706.03762>

BERT

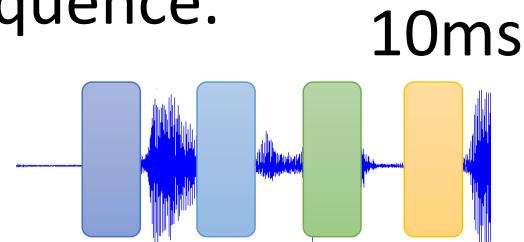
<https://arxiv.org/abs/1810.04805>

Widely used in Natural Language Processing (NLP)!

Self-attention for Speech

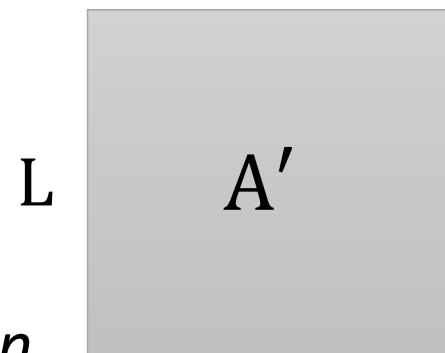
① Speech is a very long vector sequence.

→ 音長以後
可達未成 LJ 向量
而 L 通常很大



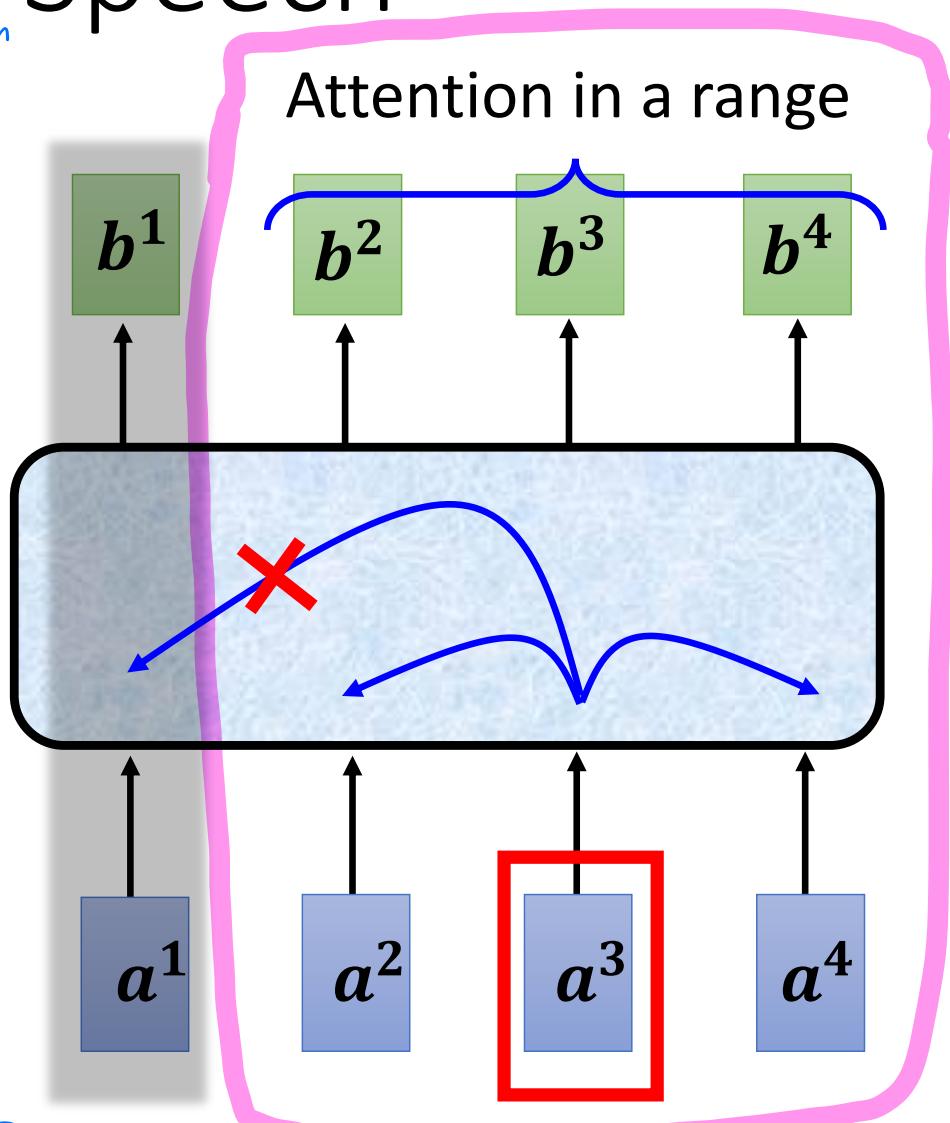
If input sequence is length L

② NP-J



就会很大

Attention in a range

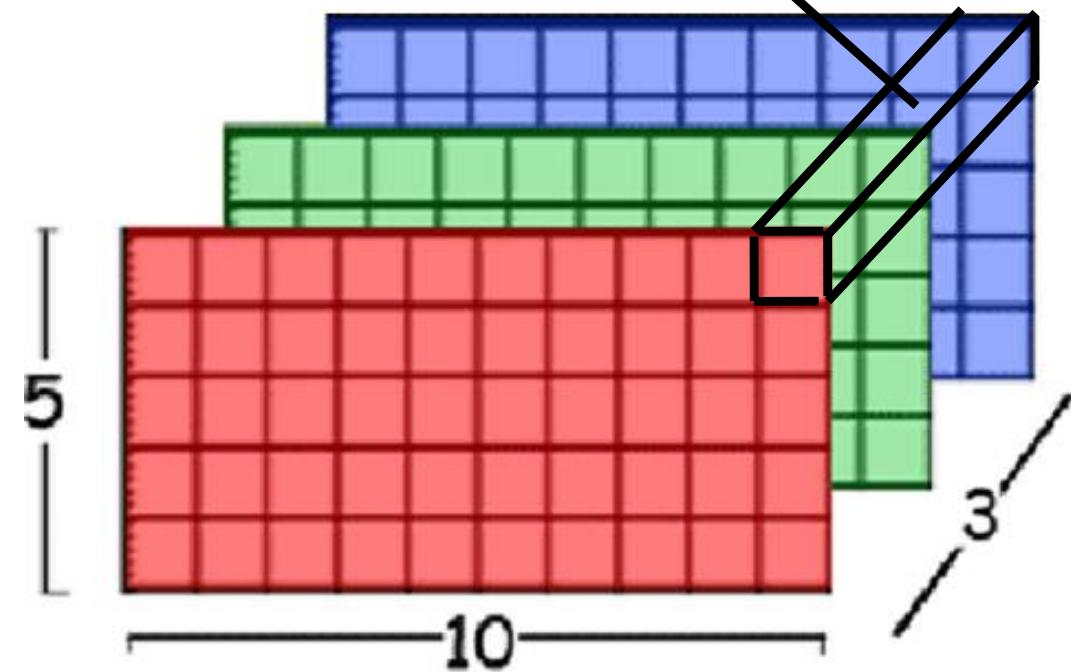
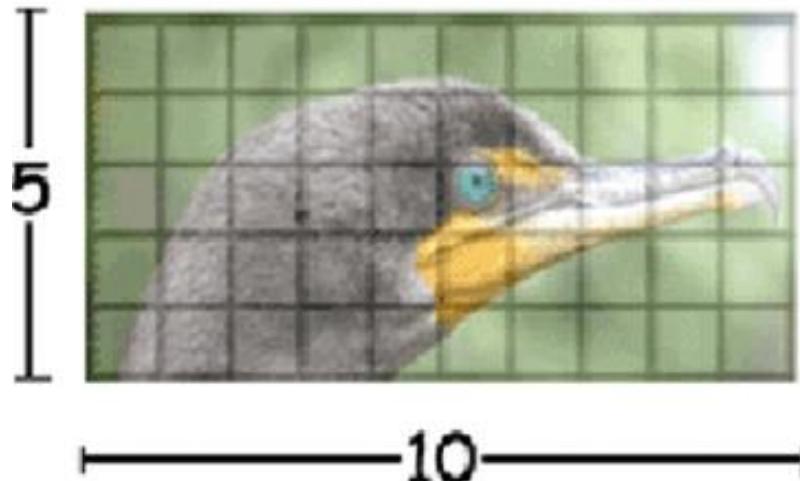


③ Truncated Self-attention

解快得多, 可以以 window, 窗口 9

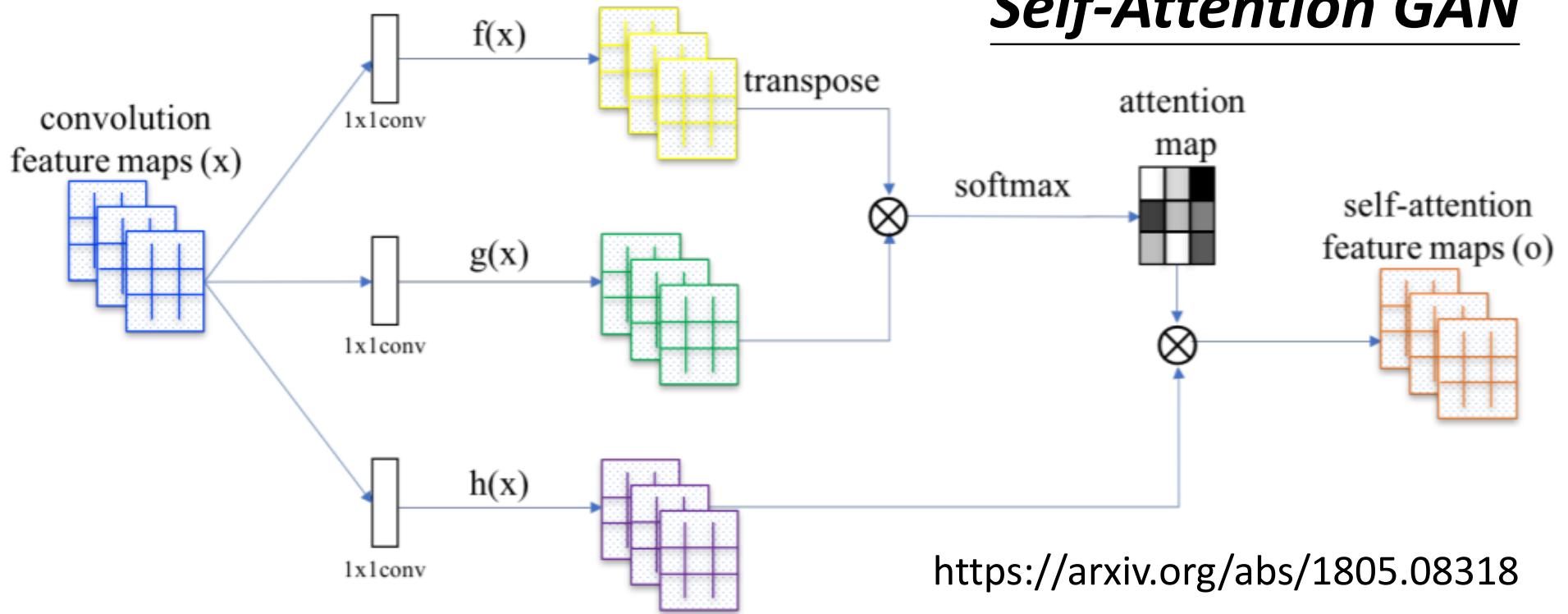
Self-attention for Image

An **image** can also be considered as a **vector set**.

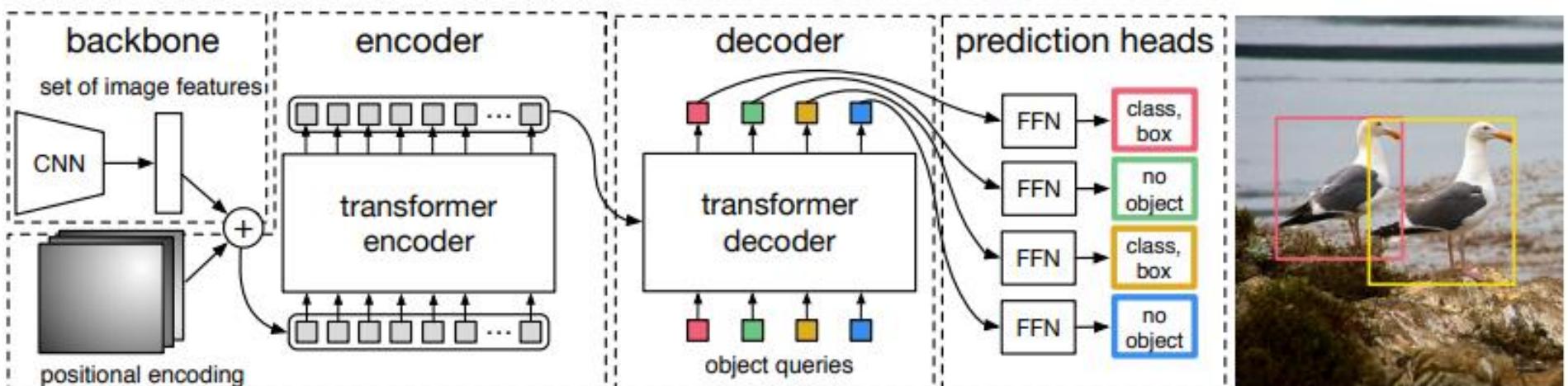


Source of image: https://www.researchgate.net/figure/Color-image-representation-and-RGB-matrix_fig15_282798184

Self-Attention GAN



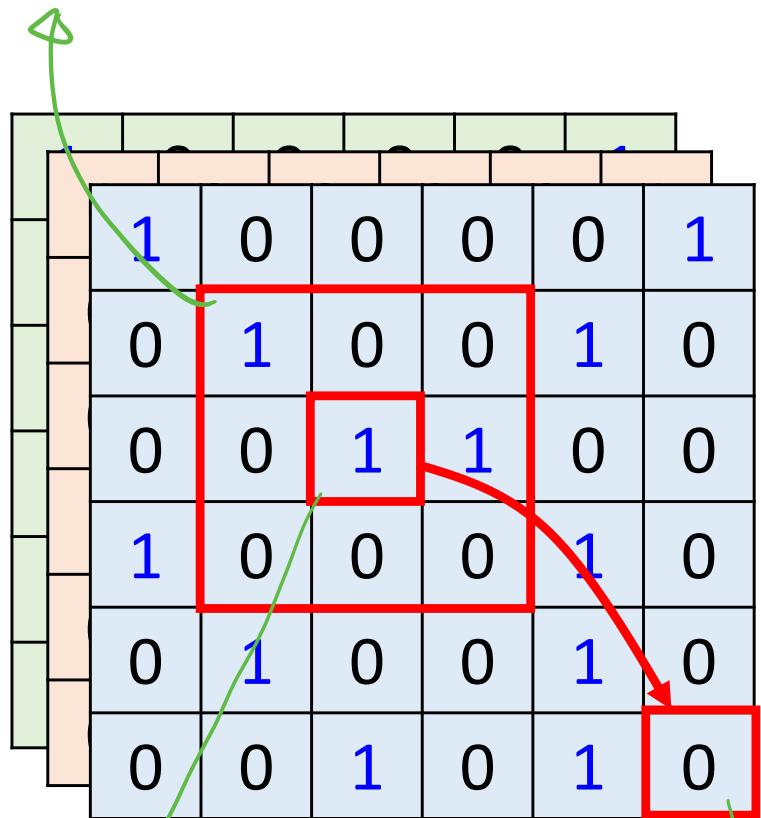
DEtection Transformer (DETR)



<https://arxiv.org/abs/2005.12872>

Self-attention v.s. CNN

③ 但如果是 CNN, 他只關注 query 和 window 內係數的相似性



① attention 的, 這是 query

② 他会找遍整個窗口，看多個係數和 query 的關聯性

④ 之謎

CNN: self-attention that can only attend in a receptive field

⑤ 答案

➤ CNN is simplified self-attention.

⑥ 答案

Self-attention: CNN with learnable receptive field

$x_{i,j}$ 取決於
window 大小
及形狀

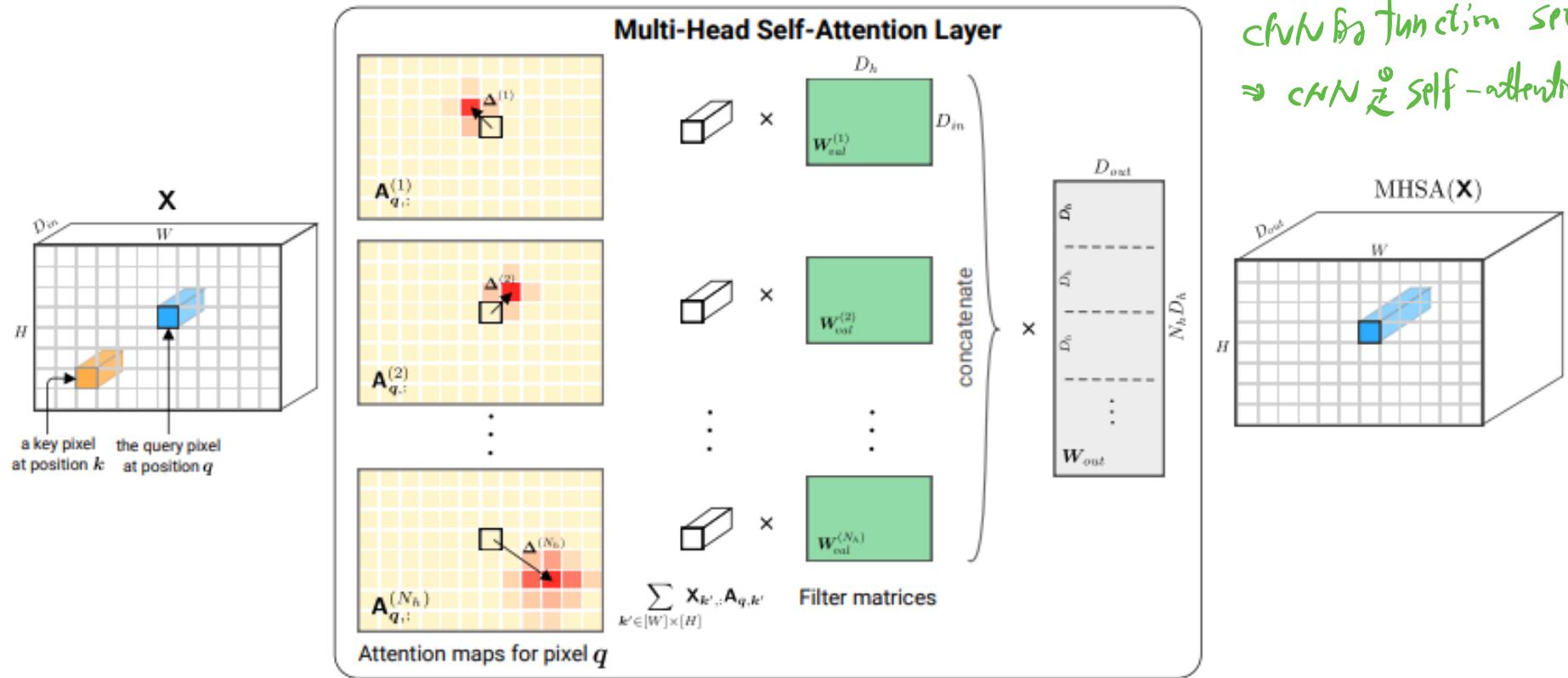
➤ Self-attention is the complex version of CNN.

Self-attention v.s. CNN

Self-attention

CNN

convolution function set
⇒ CNN ≠ self-attention
y3
fb.

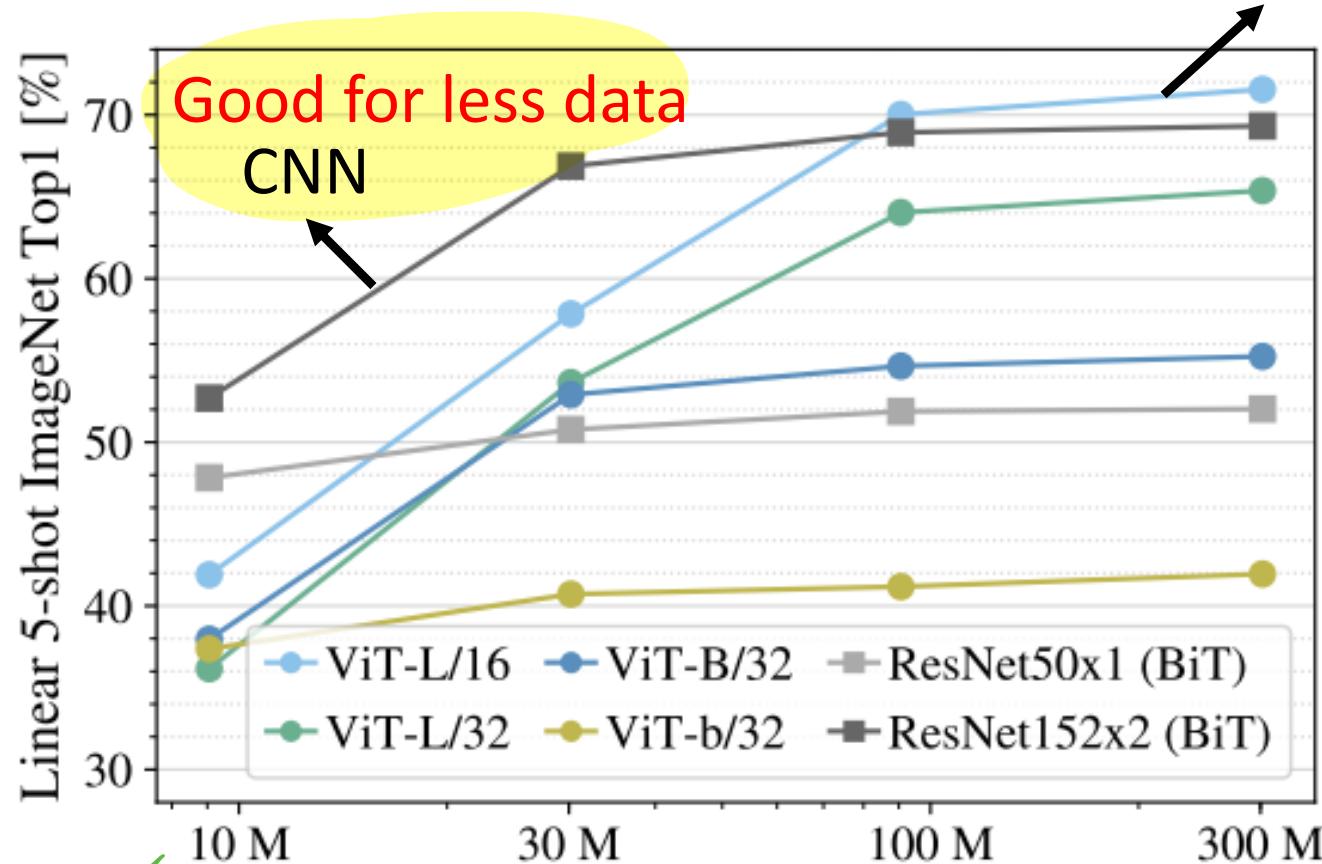


On the Relationship between Self-Attention and Convolutional Layers

<https://arxiv.org/abs/1911.03584>

Self-attention v.s. CNN

Good for more data
Self-attention



千萬張圖，第一次「只見文字」
data3

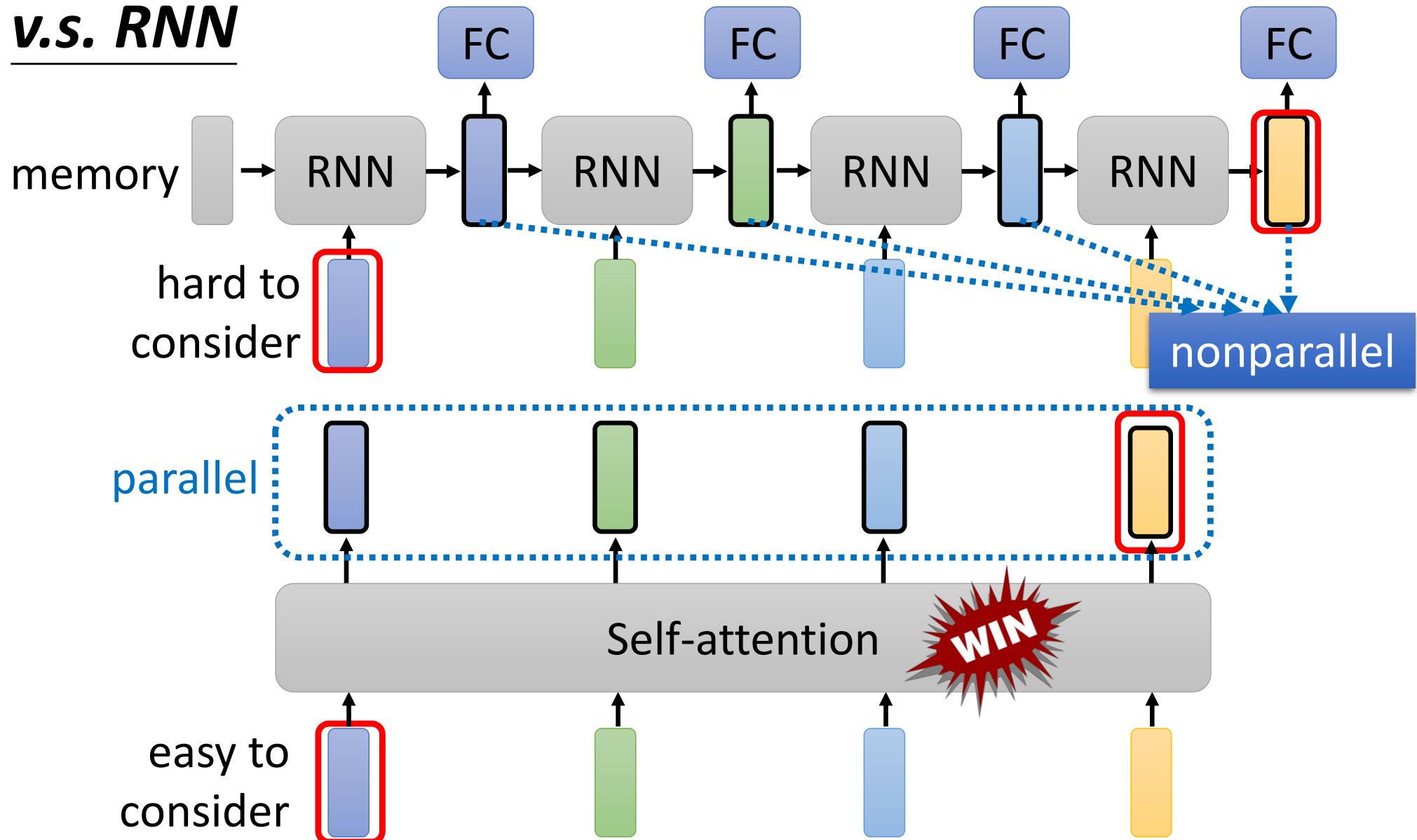
Number of JFT pre-training samples

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

<https://arxiv.org/pdf/2010.11929.pdf>

Self-attention

v.s. RNN



Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention
<https://arxiv.org/abs/2006.15236>

To learn more about RNN



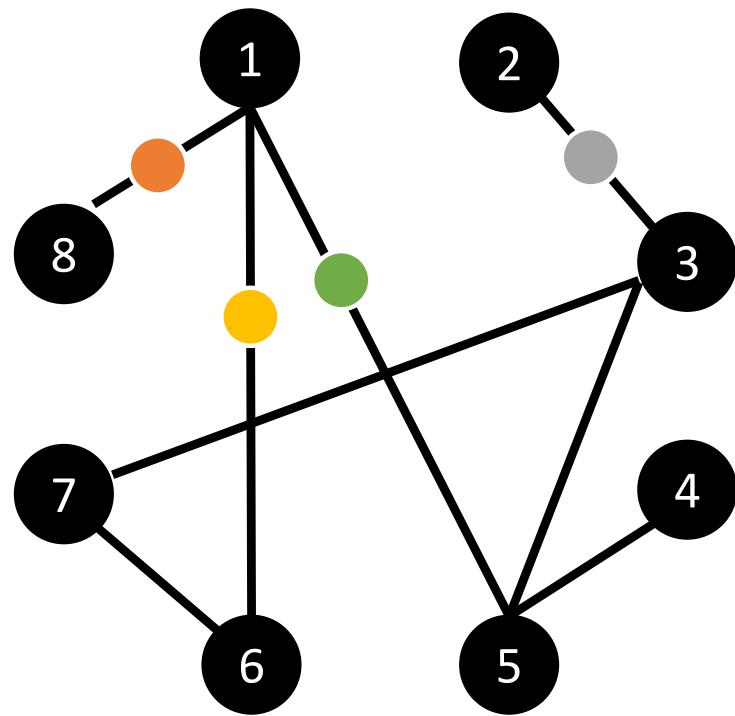
<https://youtu.be/xCGidAeyS4M>
(in Mandarin)



<https://youtu.be/Jjy6ER0bHv8>
(in English)

Self-attention for Graph

只對有 edge 相連的 node
→ i.e. 邊所接觸的



Consider **edge**: only attention to connected nodes

Attention Matrix									
		1	2	3	4	5	6	7	8
1	1					●	●		
	2			●					
3		●							
4									
5	●								
6		●							
7									
8	●							0	

This is one type of **Graph Neural Network (GNN)**.

Self-attention for Graph

- To learn more about GNN ...



<https://youtu.be/eybCCtNKwzA>
(in Mandarin)

约 3hrs, 水很深~



<https://youtu.be/M9ht8vsVEw8>
(in Mandarin)

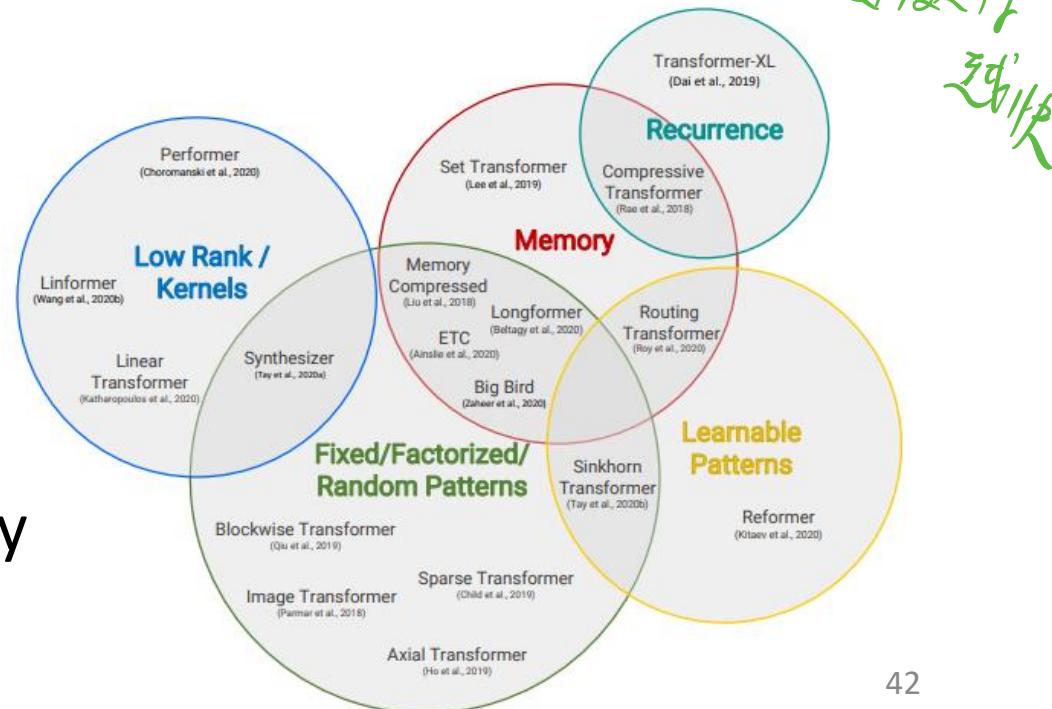
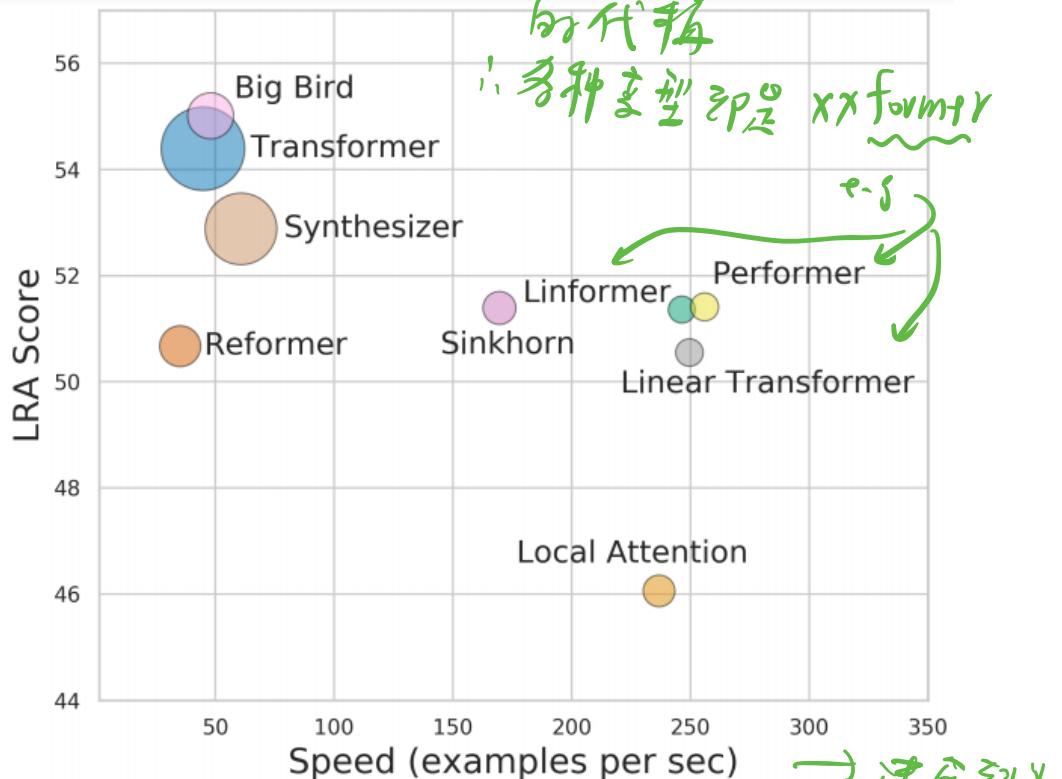
→ Transformer 已被集成 self-attention
41

To Learn More ...

這篇 paper, 用 self-attention 的變型
Long Range Arena: A Benchmark for Efficient Transformers
<https://arxiv.org/abs/2011.04006>

降低記憶量
(減少可怕的
計算量)

Efficient Transformers: A Survey
<https://arxiv.org/abs/2009.06732>



Q&A