

The background of the slide is a dark, textured image featuring a large crowd of stylized human figures. Most figures are dark brown or black, but one figure in the center is a lighter, greyish-brown color and has its arms raised in a 'V' shape, making it stand out from the rest of the crowd.

Self-attention

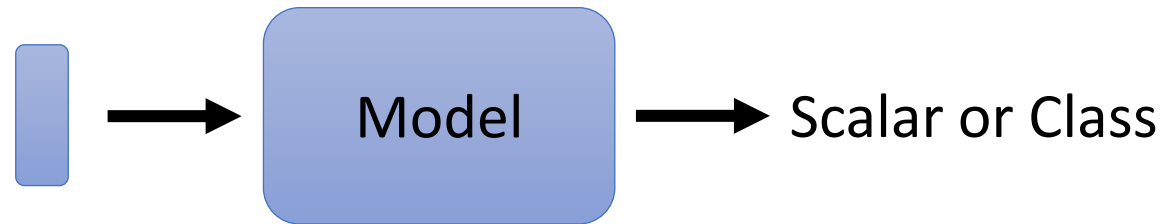
自注意力机制

Hung-yi Lee

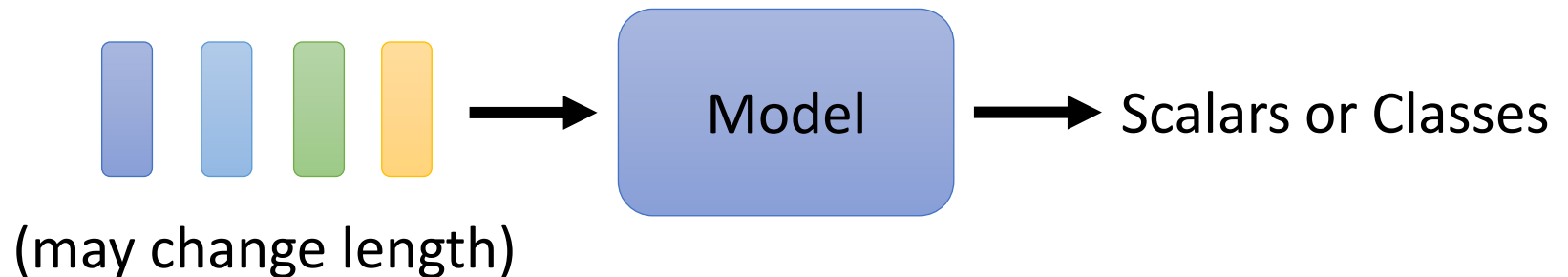
李宏毅

Sophisticated Input

- Input is a **vector** 之前的输入：一个向量




- Input is a **set of vectors** 现在我们开始考虑的输入：一排向量



Vector Set as Input

this is a cat



1
One-hot Encoding

apple = [1 0 0 0 0]

bag = [0 1 0 0 0]

cat = [0 0 1 0 0]

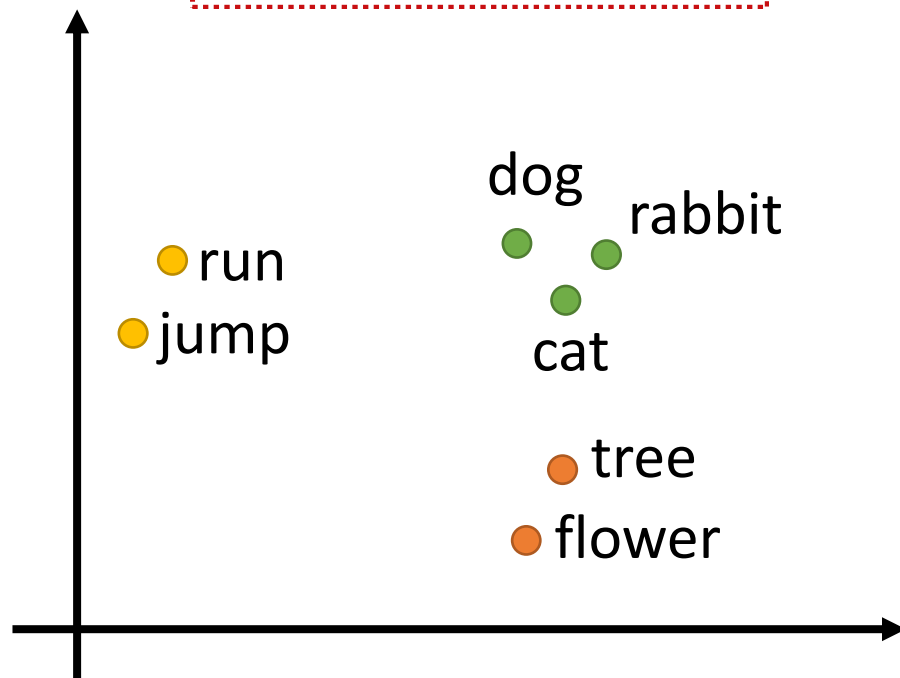
dog = [0 0 0 1 0]

elephant = [0 0 0 0 1]

1. length of vectors = # of words

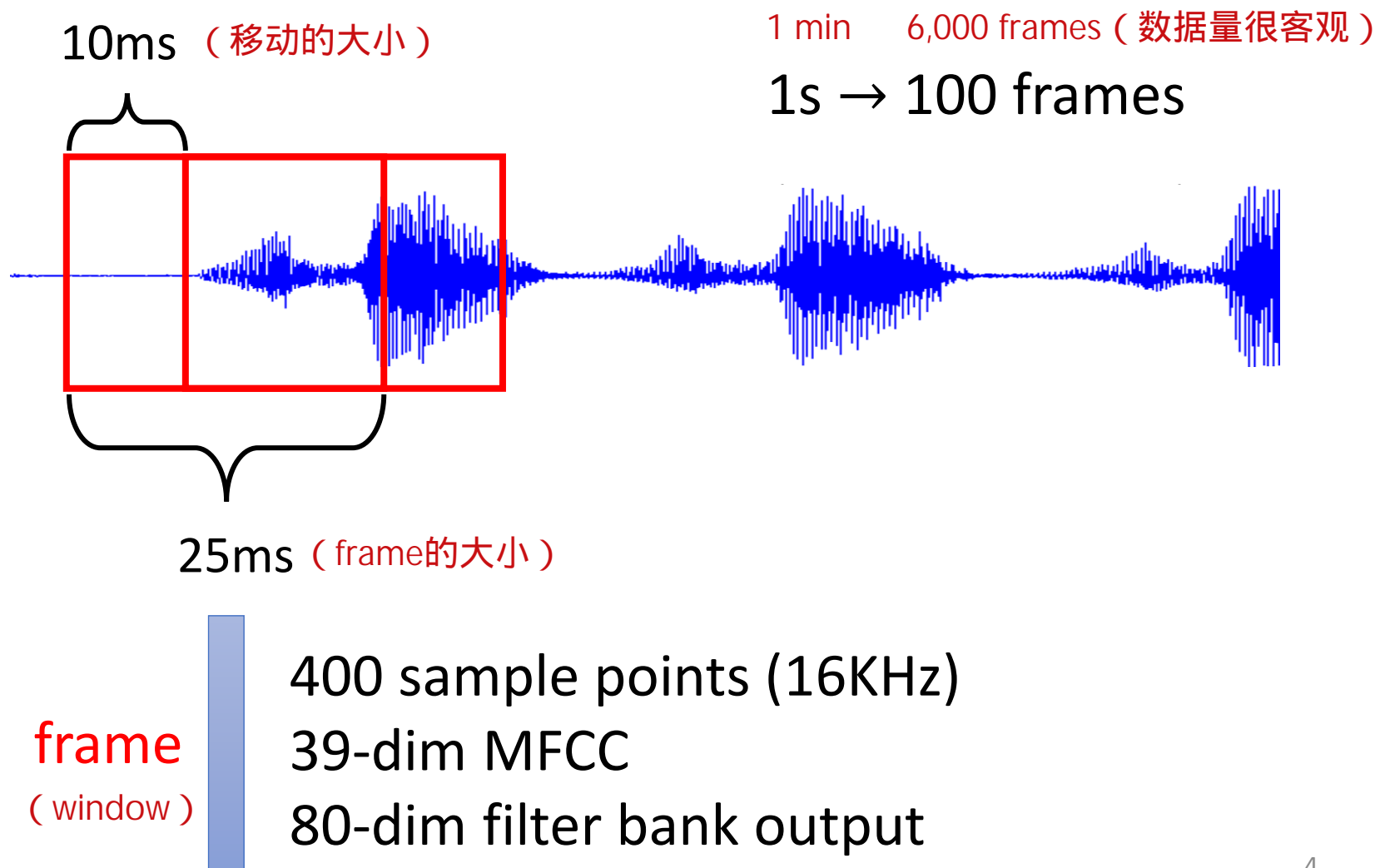
2. 问题：假设word之间没有关系

2
Word Embedding



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

Vector Set as Input



Vector Set as Input

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

node : a vector
edge : 有edge的两个
nodes会计算关联程度 α



Each profile
is a vector

Vector Set as Input

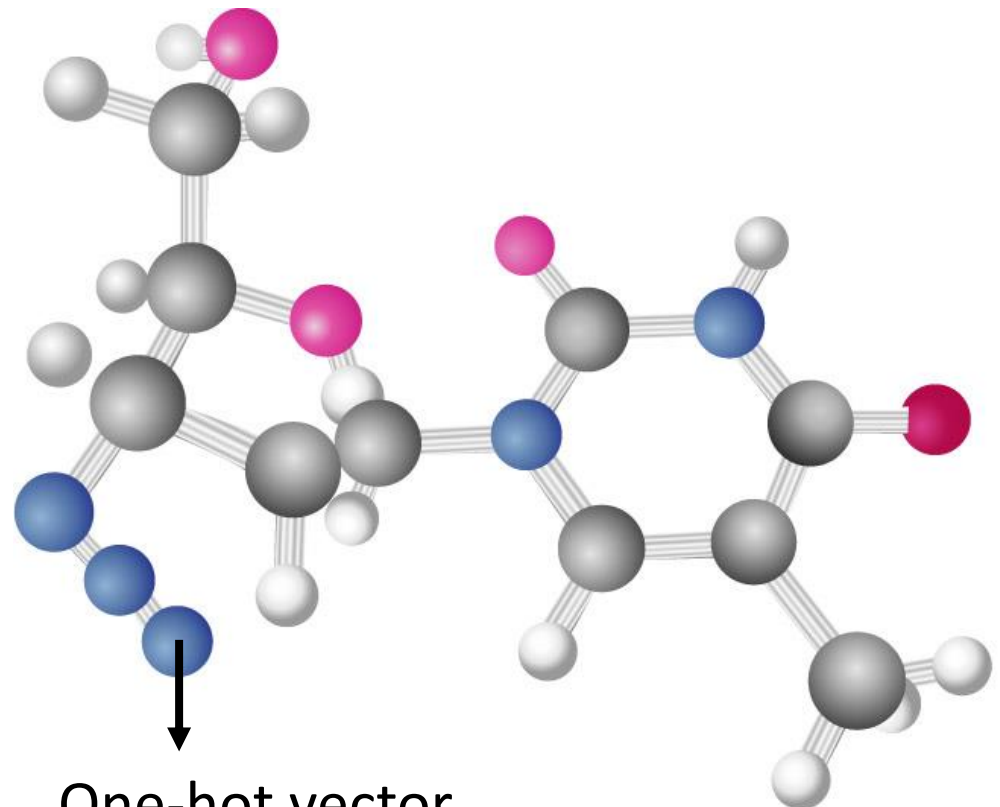
- Graph is also a set of vectors (consider each **node** as a **vector**)

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

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

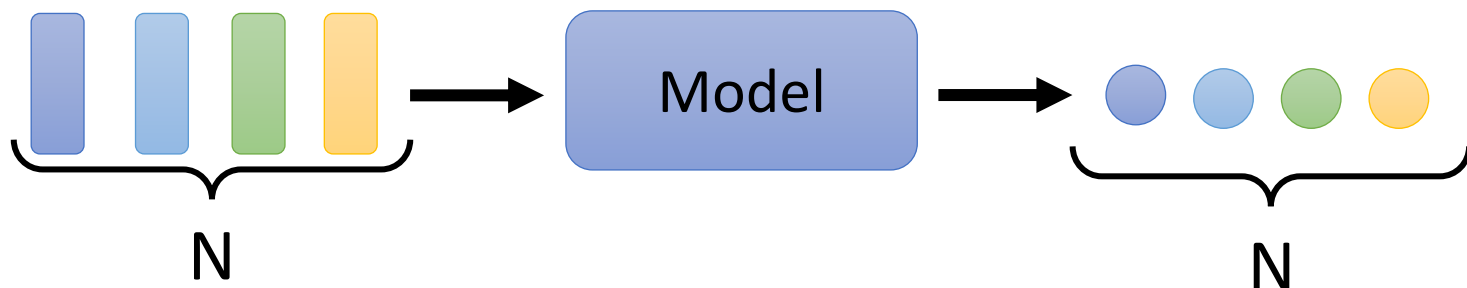
$$O = [0 \ 0 \ 1 \ 0 \ 0 \ \dots]$$

⋮

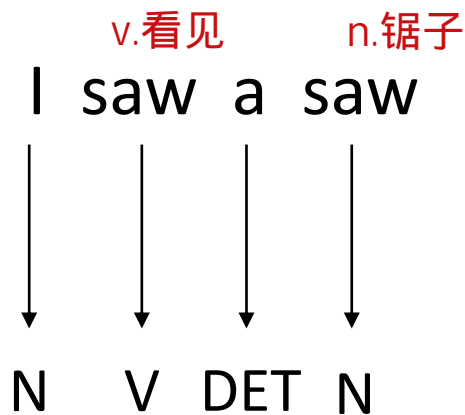


What is the output? 有三种可能性。

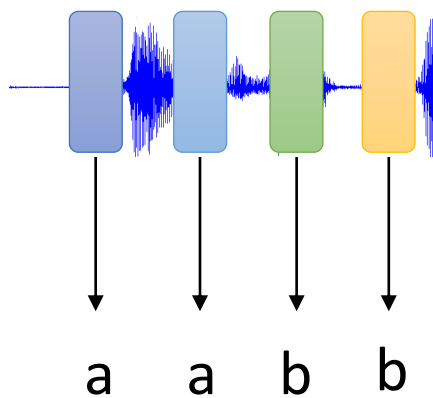
- Each vector has a label.



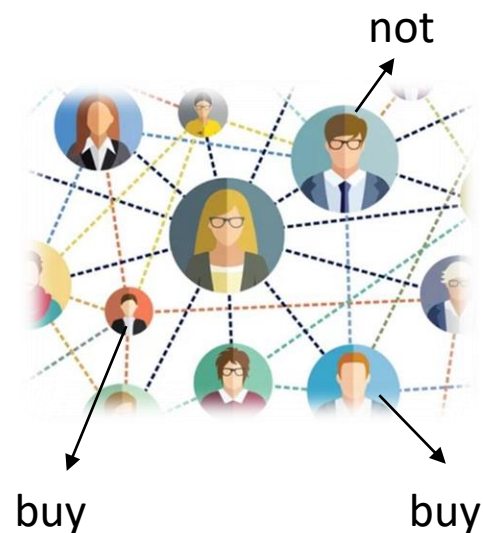
Example Applications



POS tagging



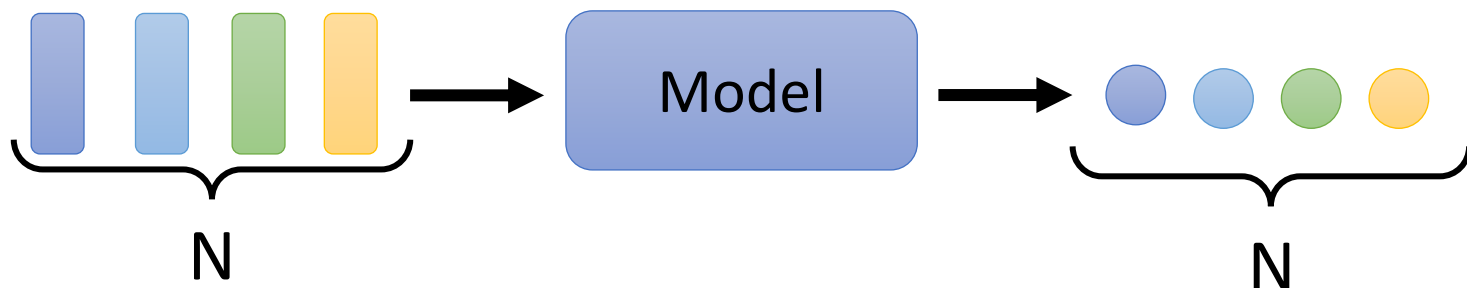
HW2



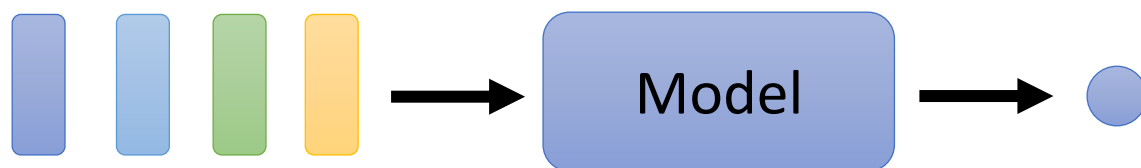
语音识别的简化版

What is the output?

1. • Each vector has a label.



2. • The whole sequence has a label. Sequence-to-vector model



Example Applications

this is good
Sentiment analysis
情感分析
positive

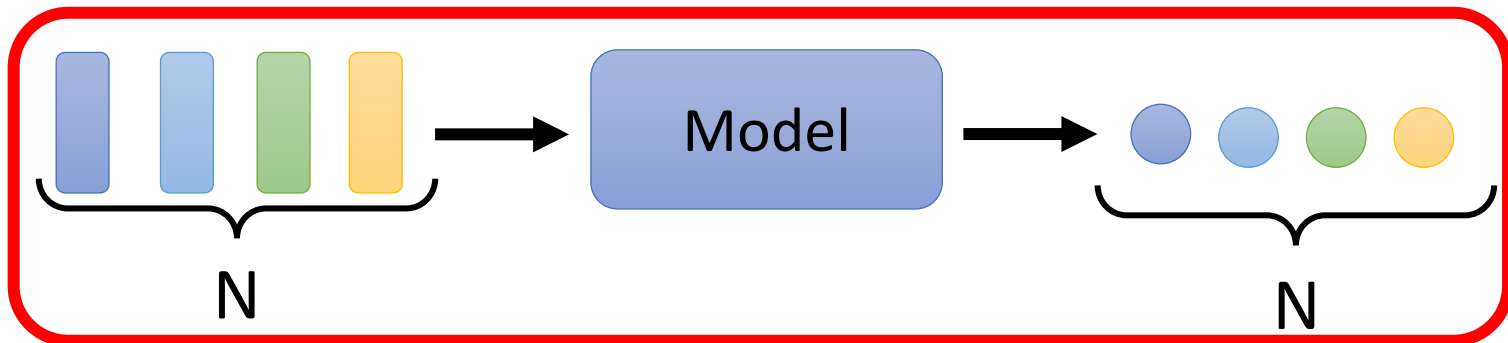
speaker 语者分析

hydrophilicity₈ 亲水性

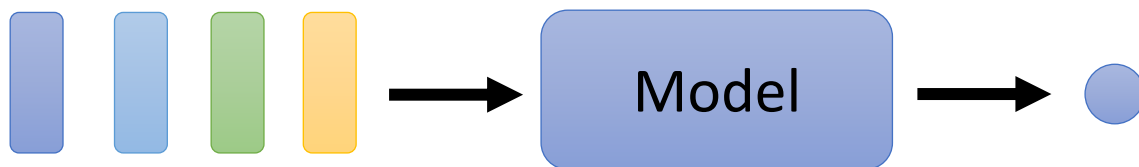
What is the output?

- 1 • Each vector has a label.

focus of this lecture

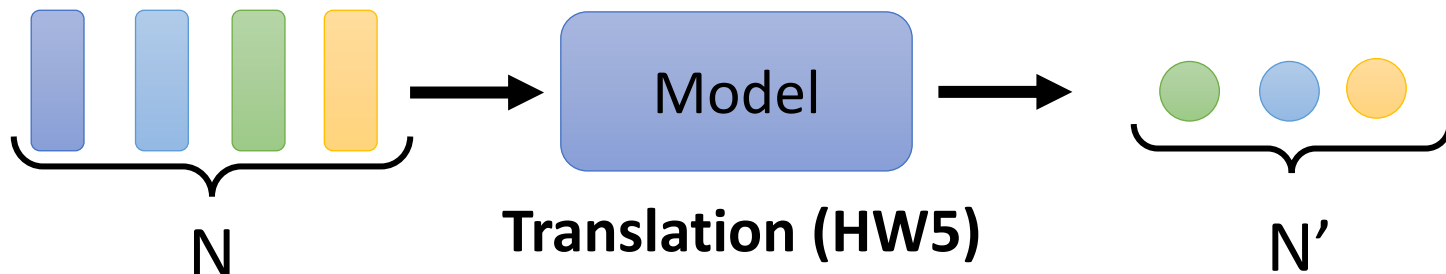


- 2 • The whole sequence has a label.



- 3 • Model decides the number of labels itself.

seq2seq



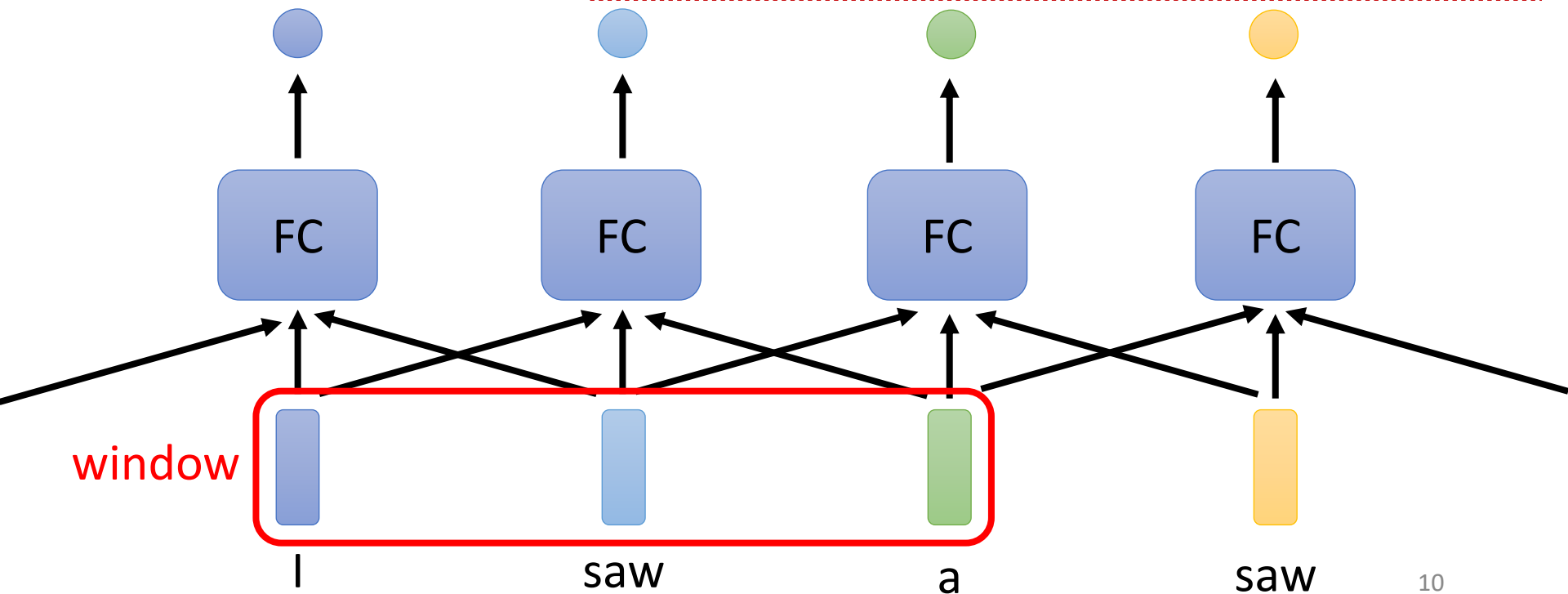
Sequence Labeling

FC Fully-connected

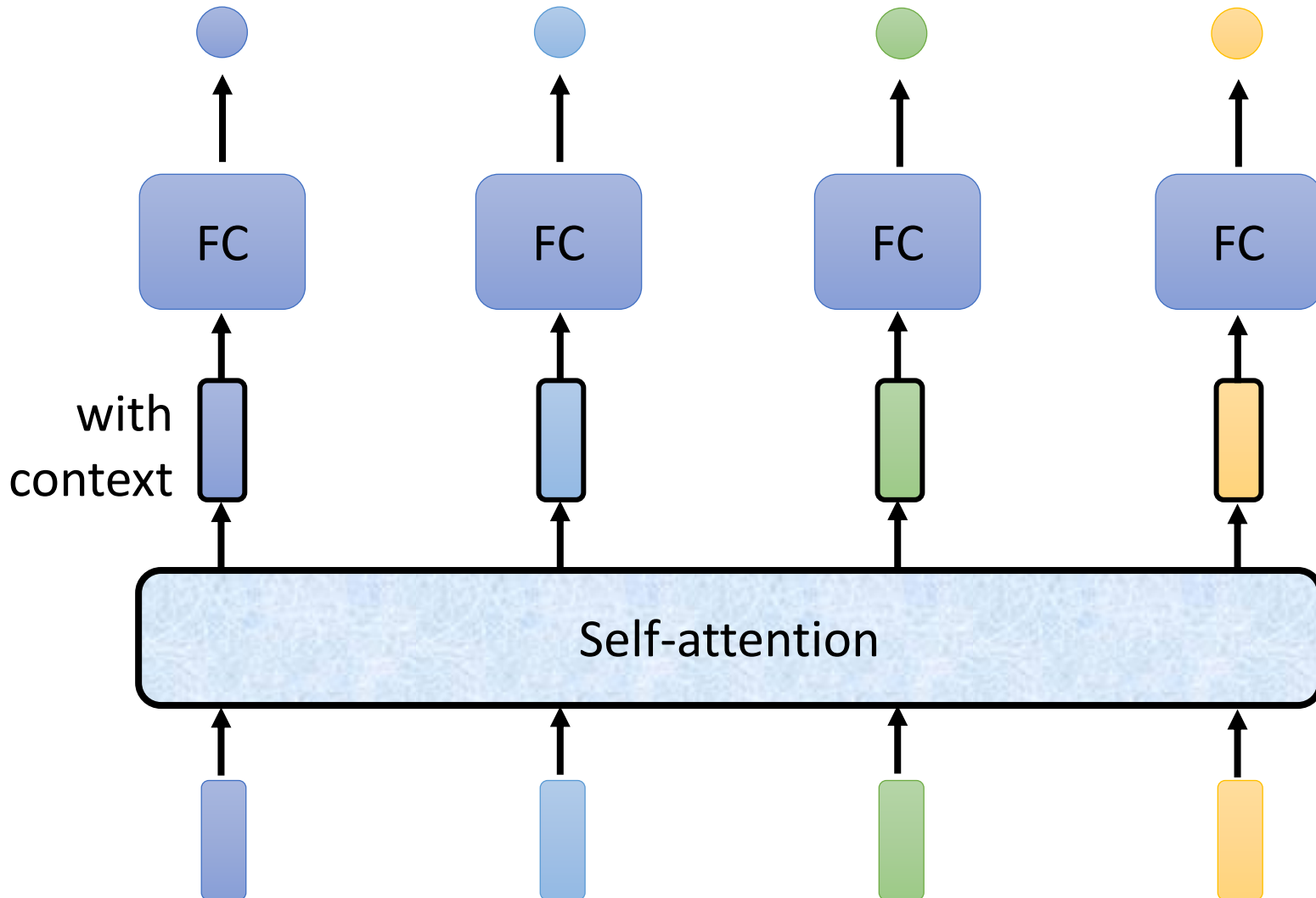
Is it possible to consider the context?

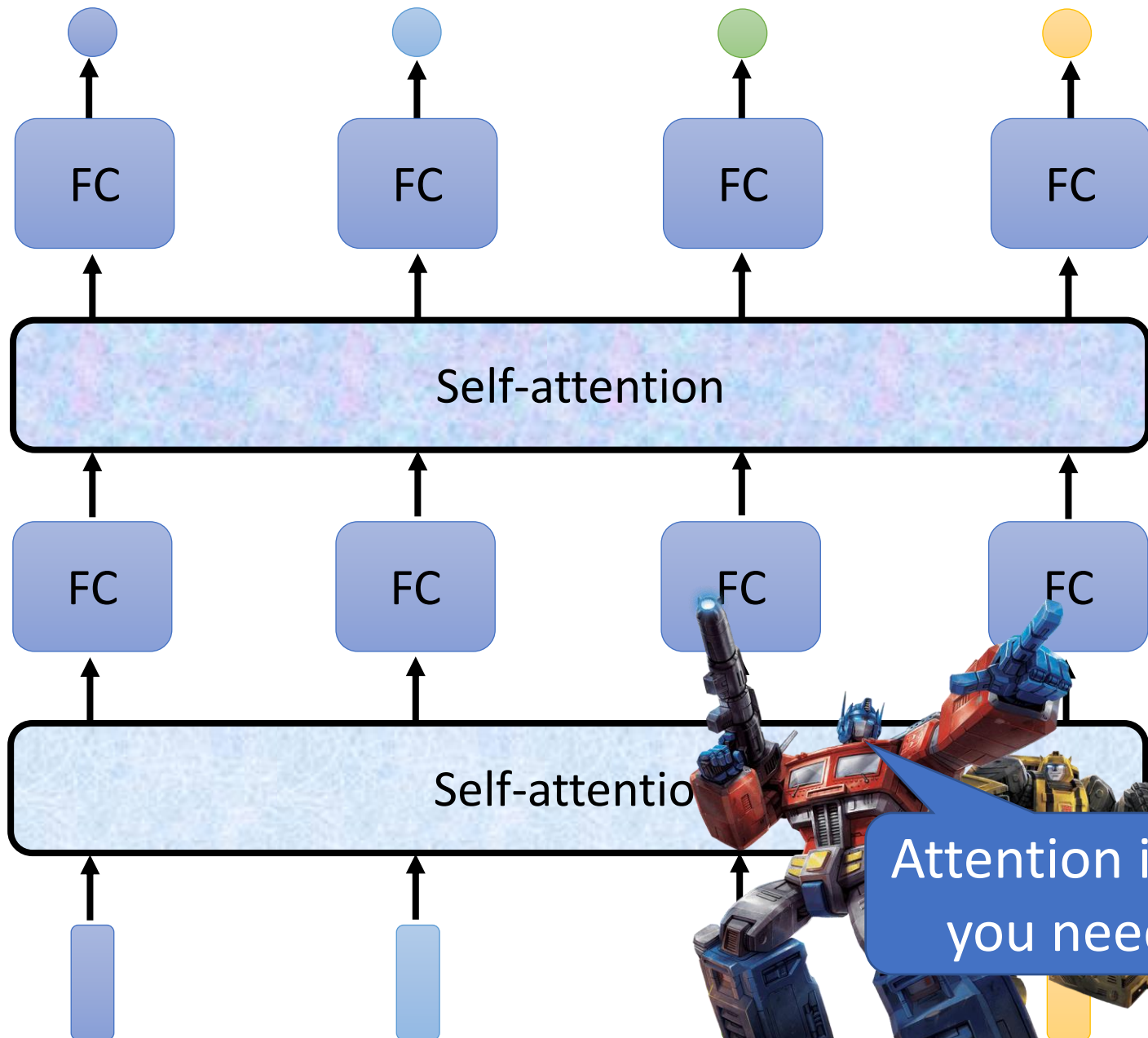
FC can consider the neighbor

How to consider the whole sequence?
a window covers the whole sequence?



Self-attention



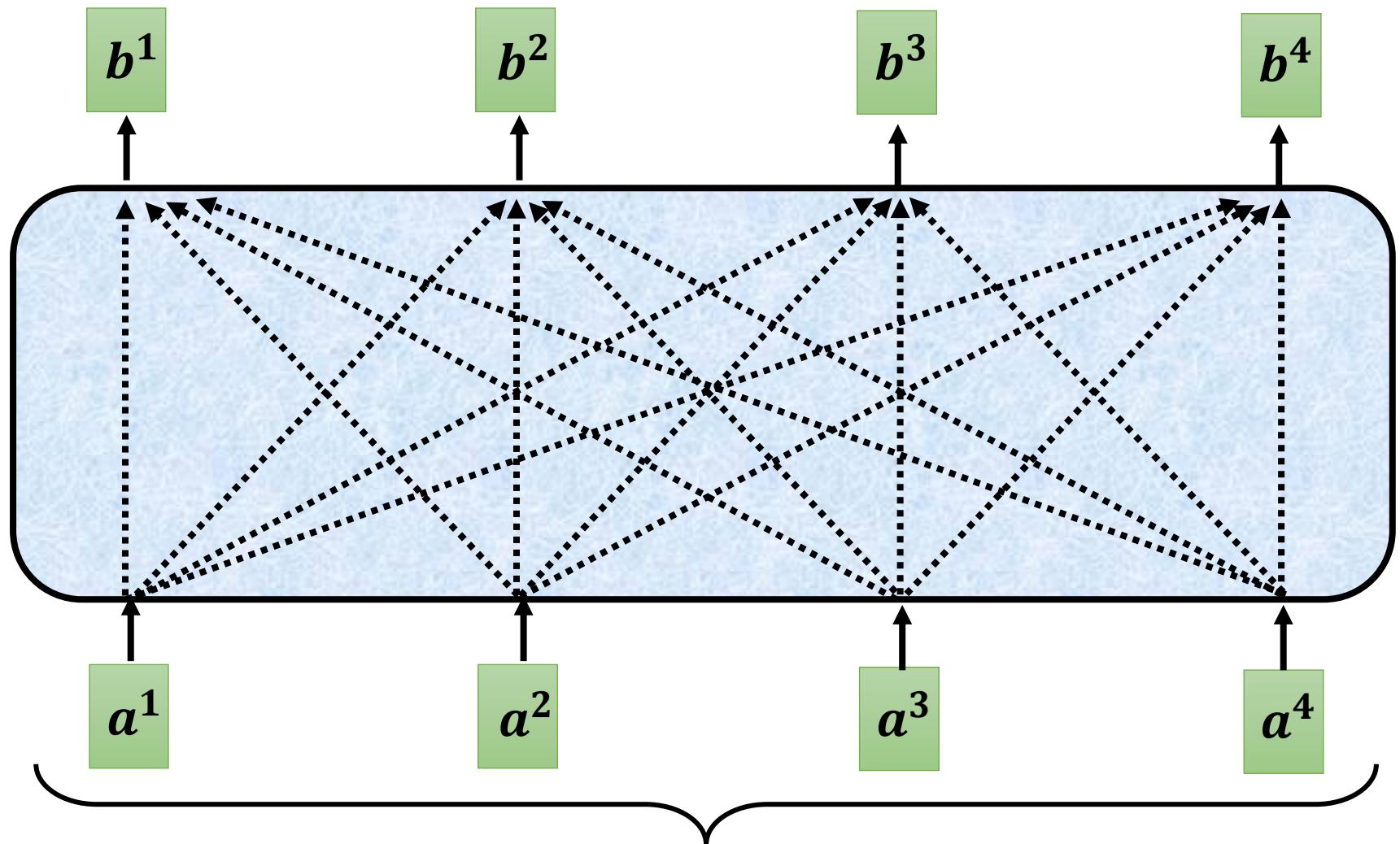


Attention is all
you need.

self-attention层可以叠加

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

Self-attention

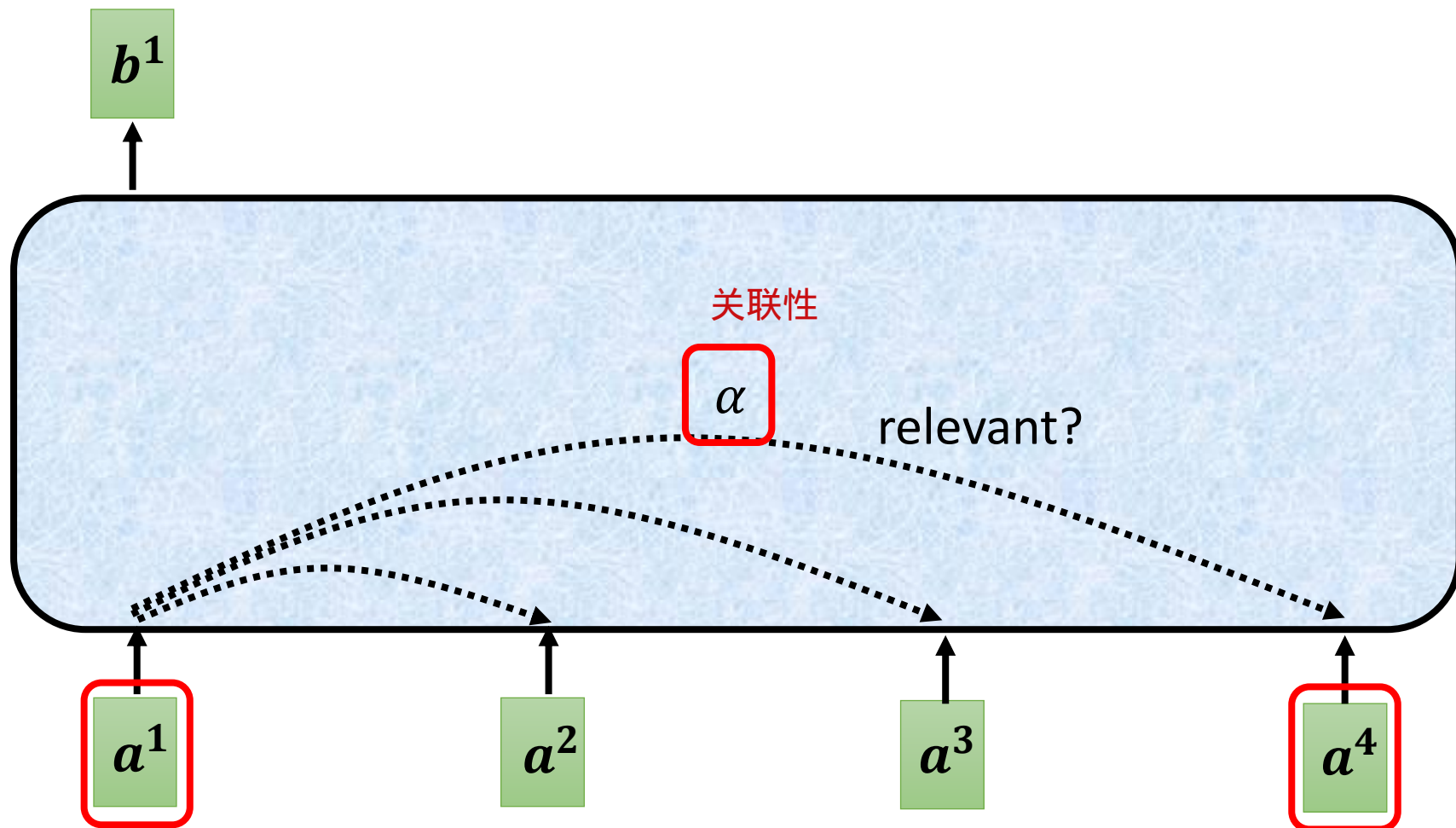


Can be either **input** or a **hidden layer**

Self-attention

怎么产生 b^1 呢？

第一步：根据找出 a^1 和其他向量的关联程度 α



Find the relevant vectors in a sequence

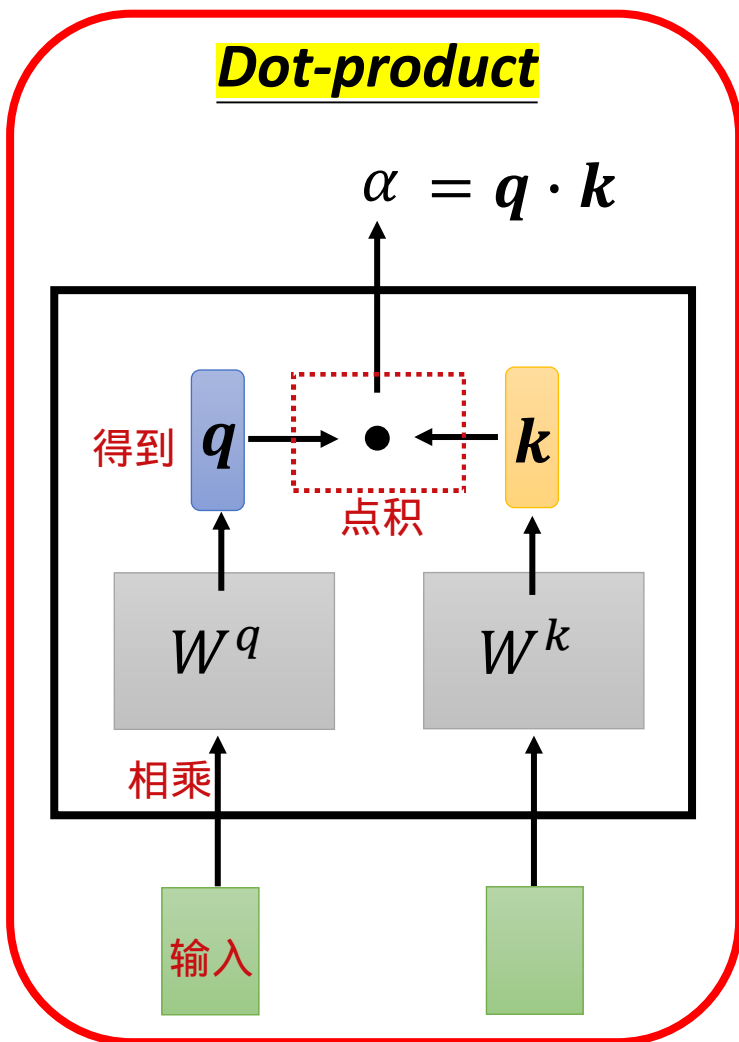
如何计算 α 呢？两种常见的方法。

1. dot-product

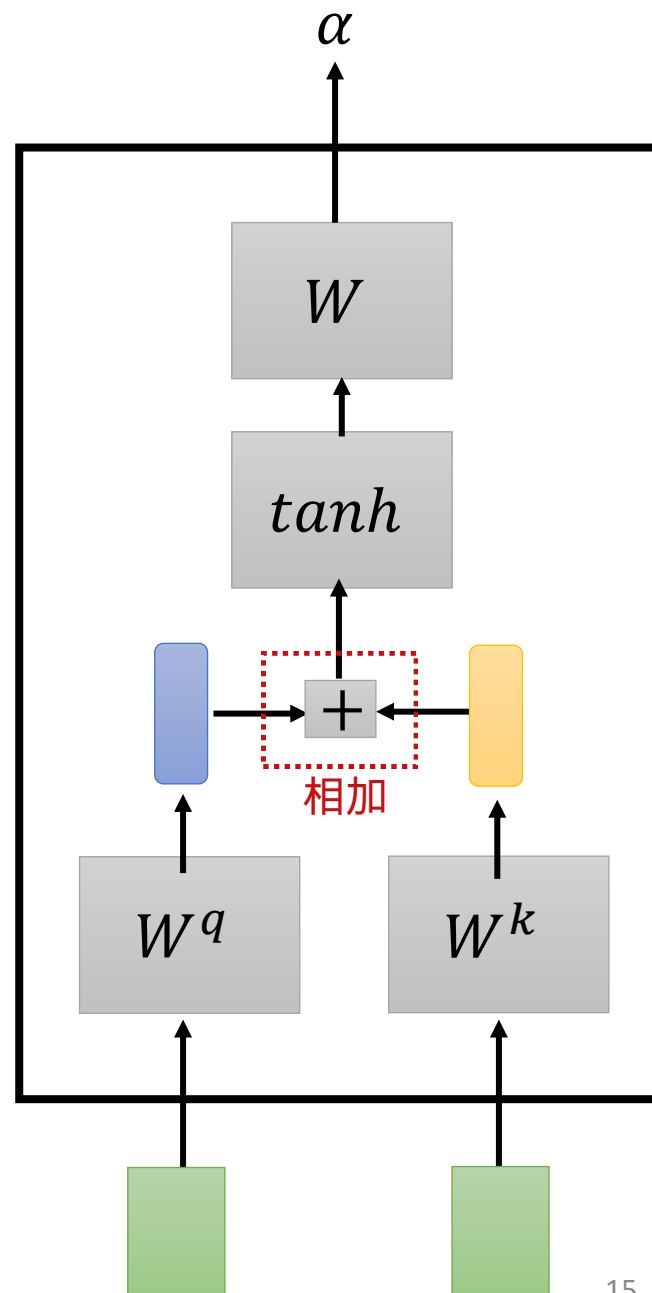
2. additive

Self-attention

Dot-product

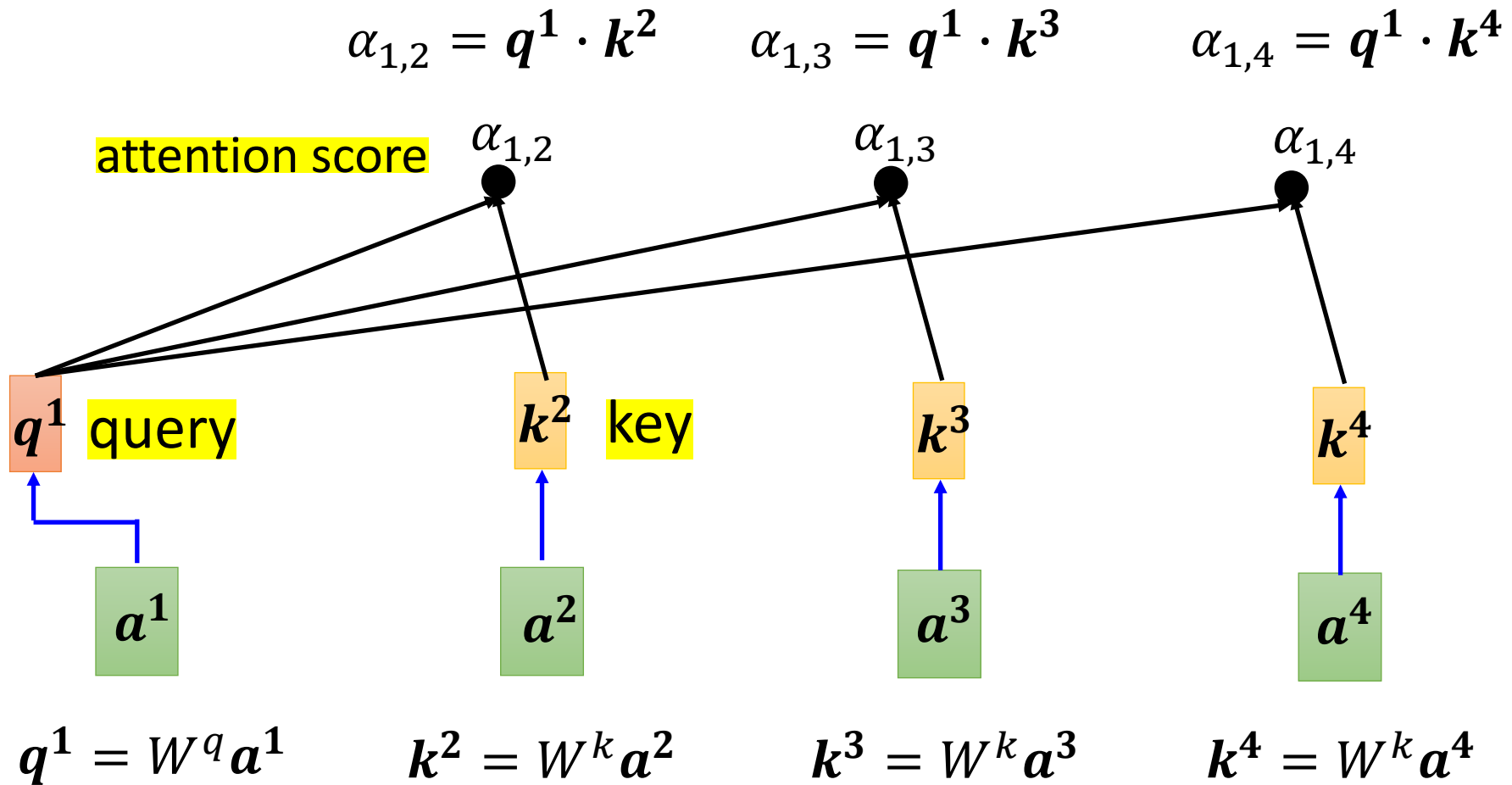


Additive



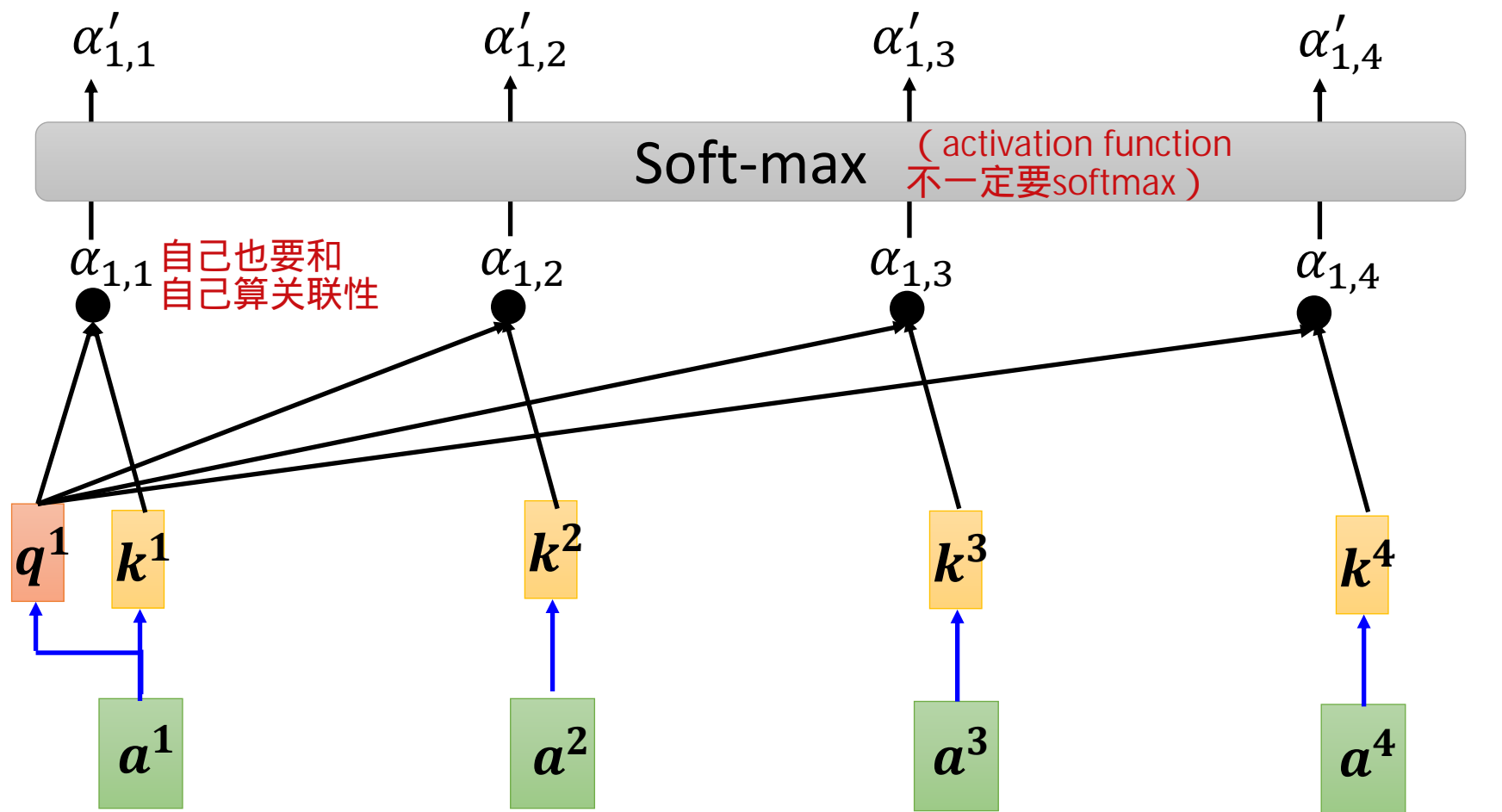
Self-attention

1. q : query
2. k : key
3. attention score



Self-attention

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



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

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

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

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

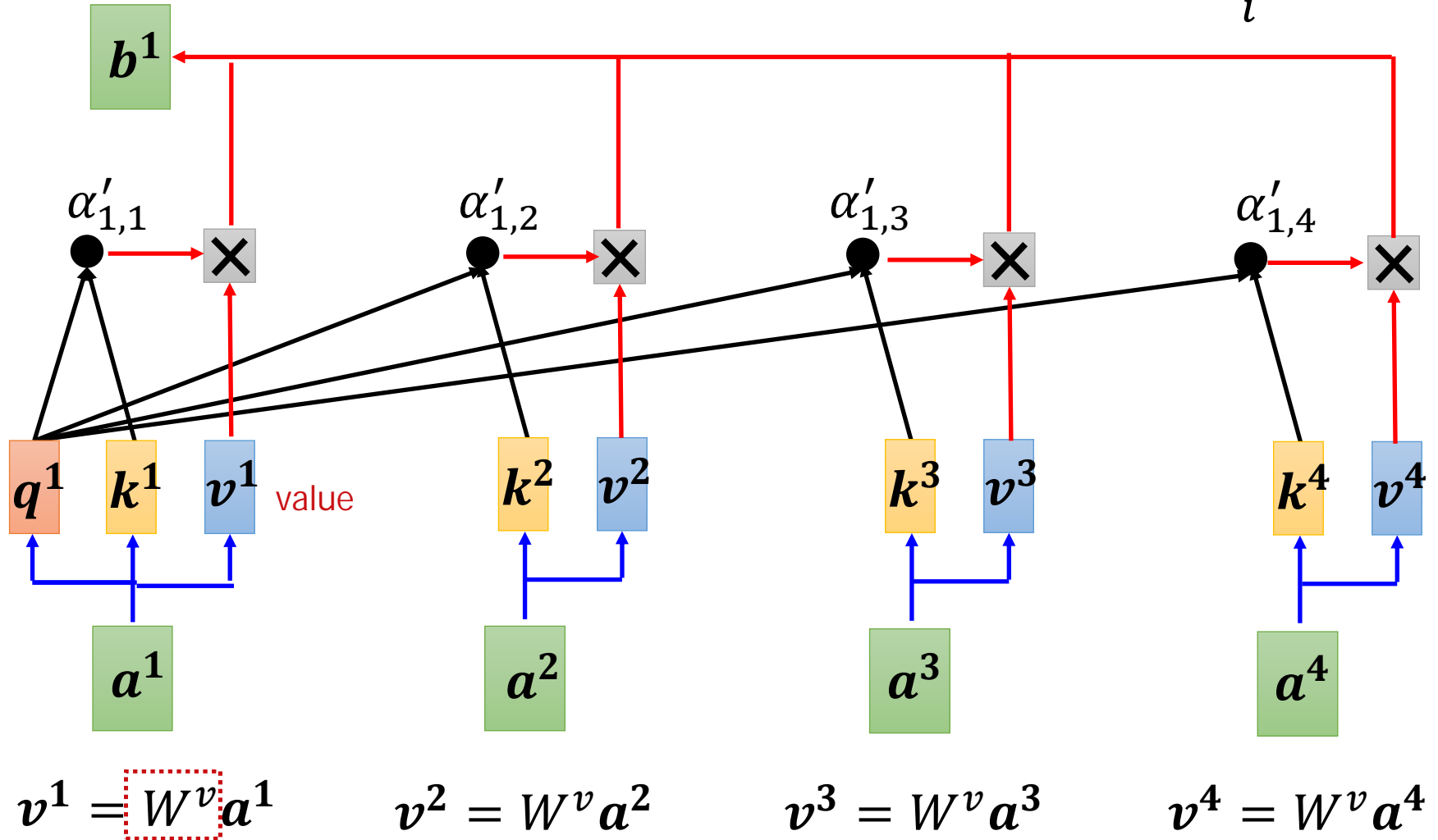
$$k^4 = W^k a^4$$

Self-attention

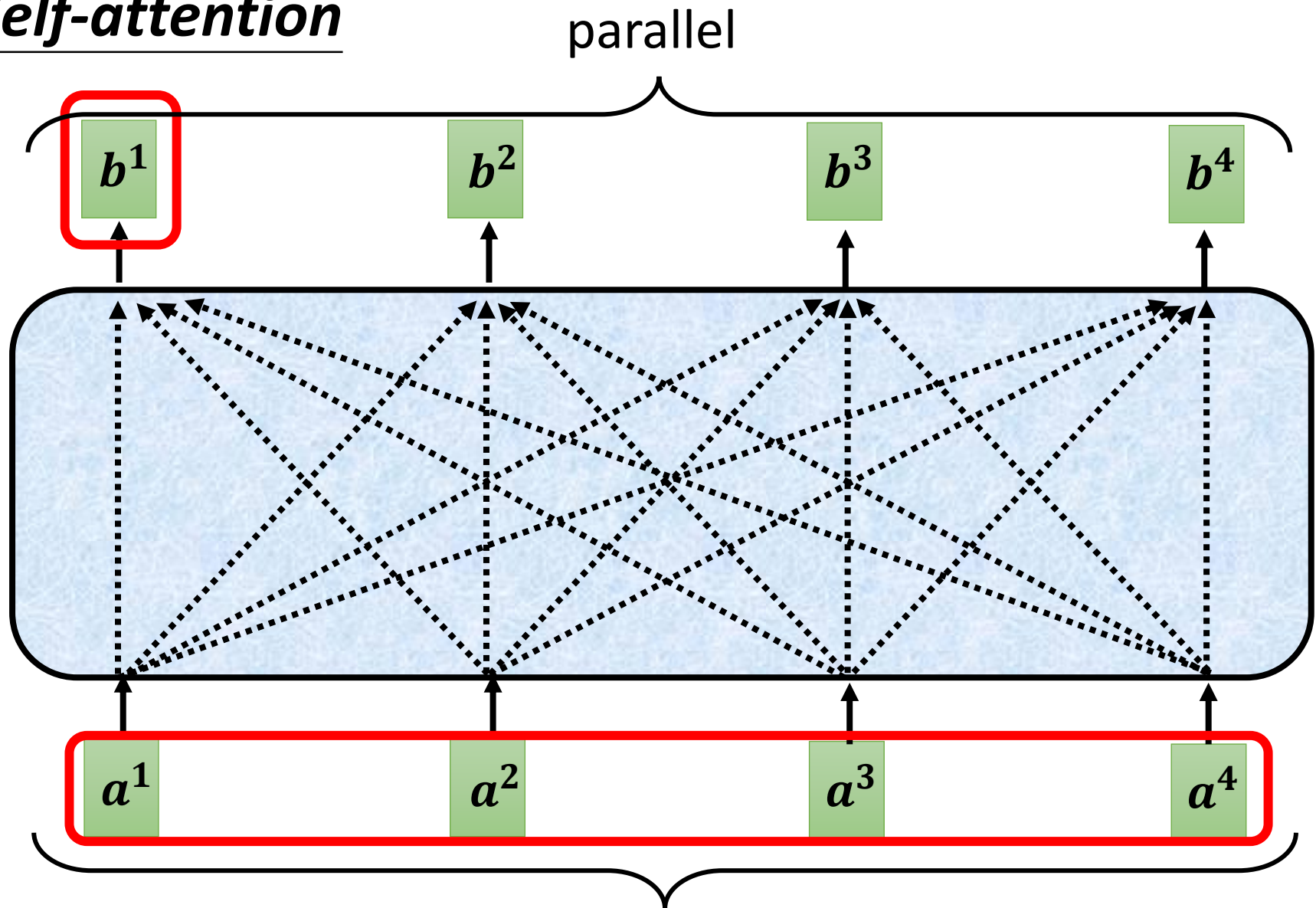
value

Extract information based
on attention scores

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



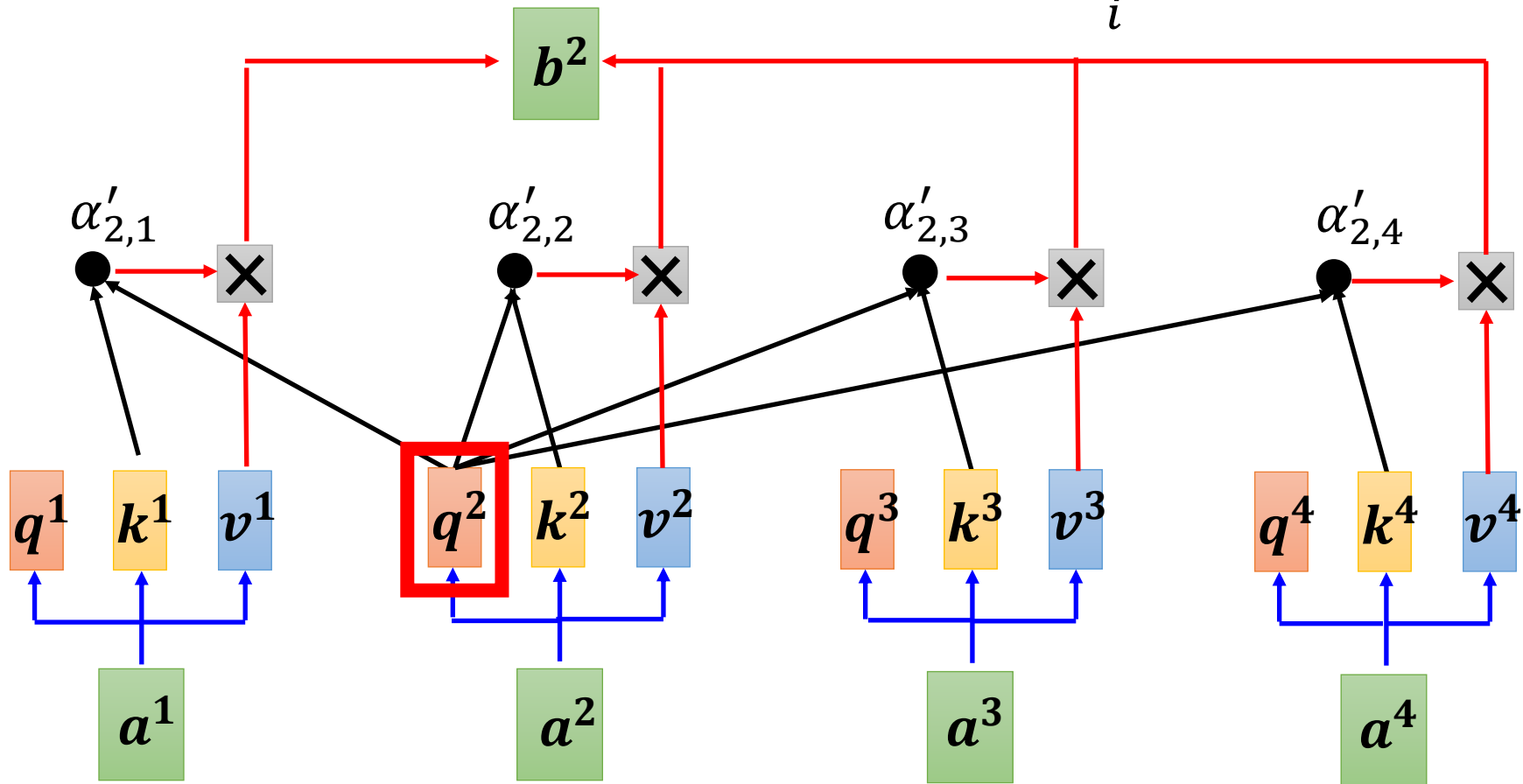
Self-attention



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

Self-attention

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



Self-attention

矩阵运算的角度：

a Q/K/V

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

$$\begin{matrix} q^1 & q^2 & q^3 & q^4 \\ \hline Q \end{matrix} = \begin{matrix} W^q & a^1 & a^2 & a^3 & a^4 \\ \hline I \end{matrix}$$

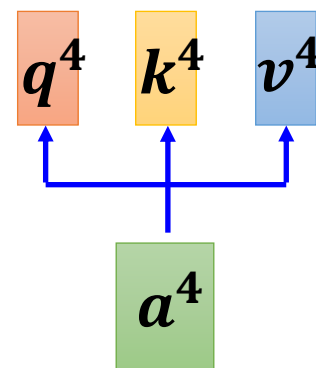
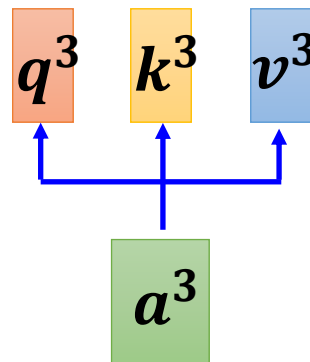
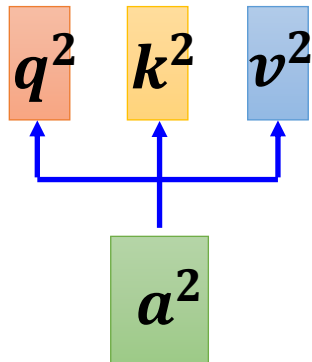
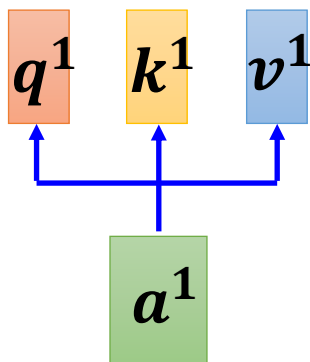
network的参数，
会被learn出来

$$k^i = W^k a^i$$

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

$$v^i = W^v a^i$$

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

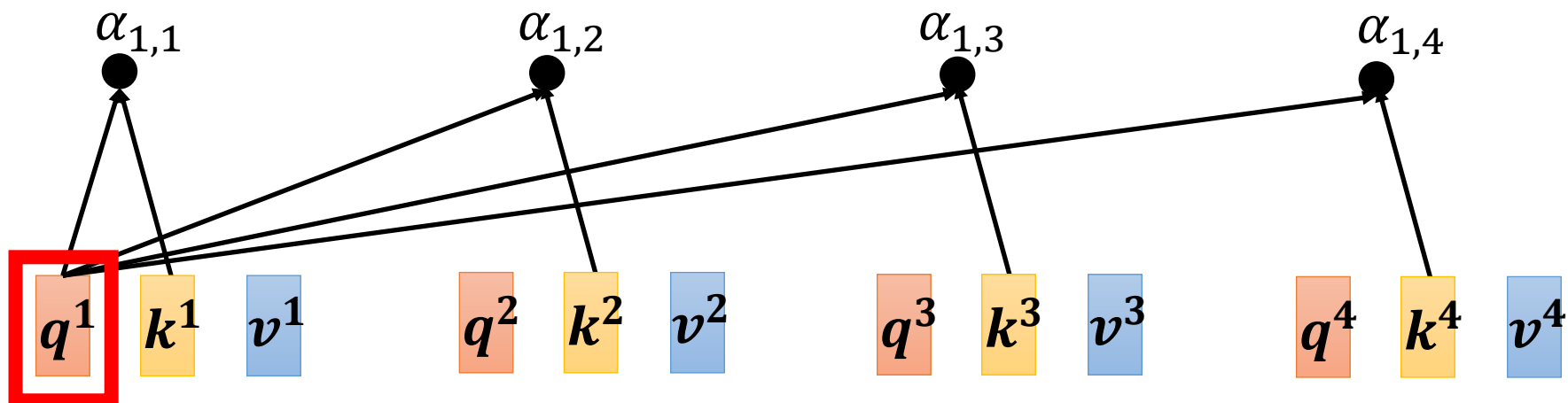


Self-attention

计算: Q/K α

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

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

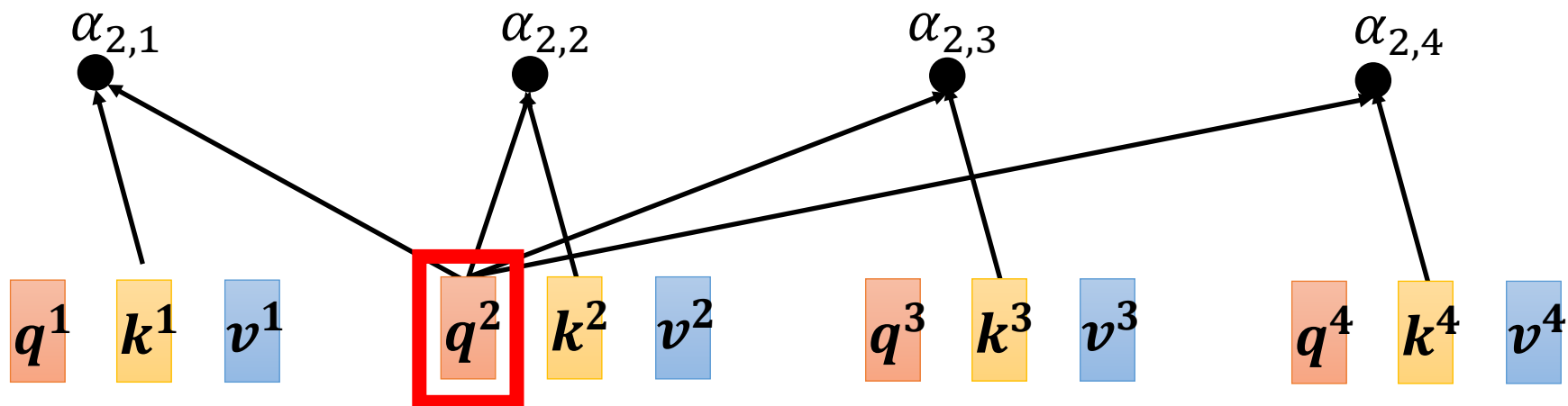


Self-attention

Q/K α α'

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

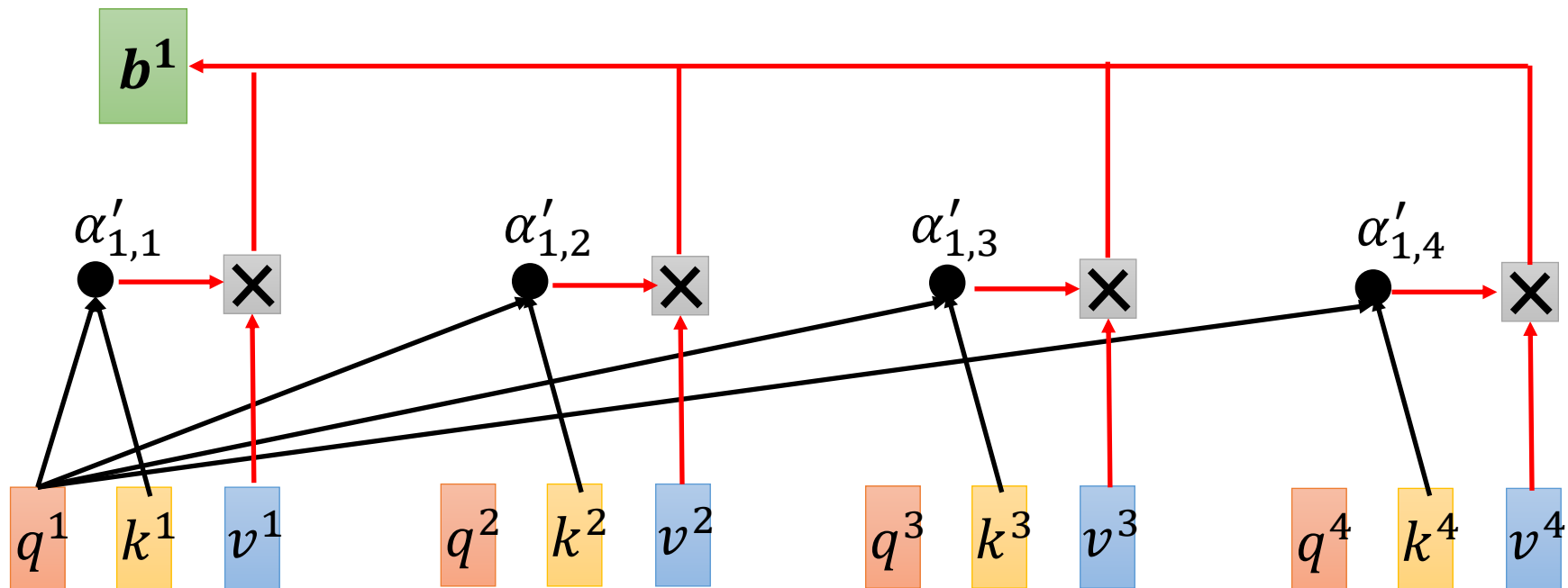
$$\begin{bmatrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{bmatrix} = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{bmatrix} q^1$$



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

A' A K^T Q

Self-attention V/α' b



$$\begin{matrix} b^1 & b^2 & b^3 & b^4 \\ \hline 0 \end{matrix} = \begin{matrix} v^1 & v^2 & v^3 & v^4 \\ \hline V \end{matrix} \begin{matrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{matrix}$$

A'

Self-attention

$$\begin{aligned} Q &= W^q I \\ K &= W^k I \\ V &= W^v I \end{aligned}$$

Parameters
to be learned
只有这三个参数需要学习获得

$$A'$$



$$A$$

$$=$$

$$K^T$$

$$Q$$

Attention Matrix

$$O$$

$$=$$

$$V$$

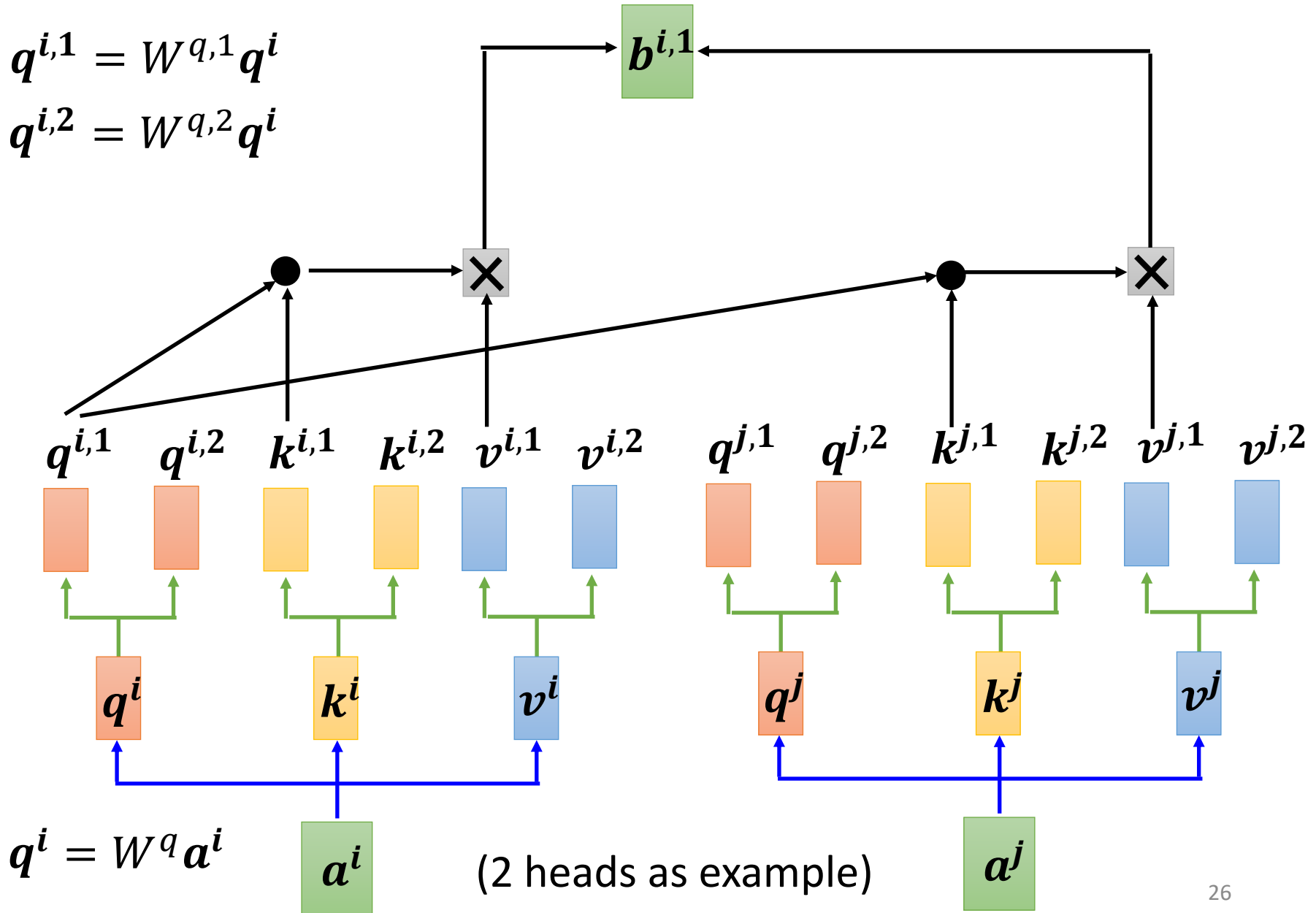
$$A'$$

Multi-head Self-attention

Different types of relevance

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

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

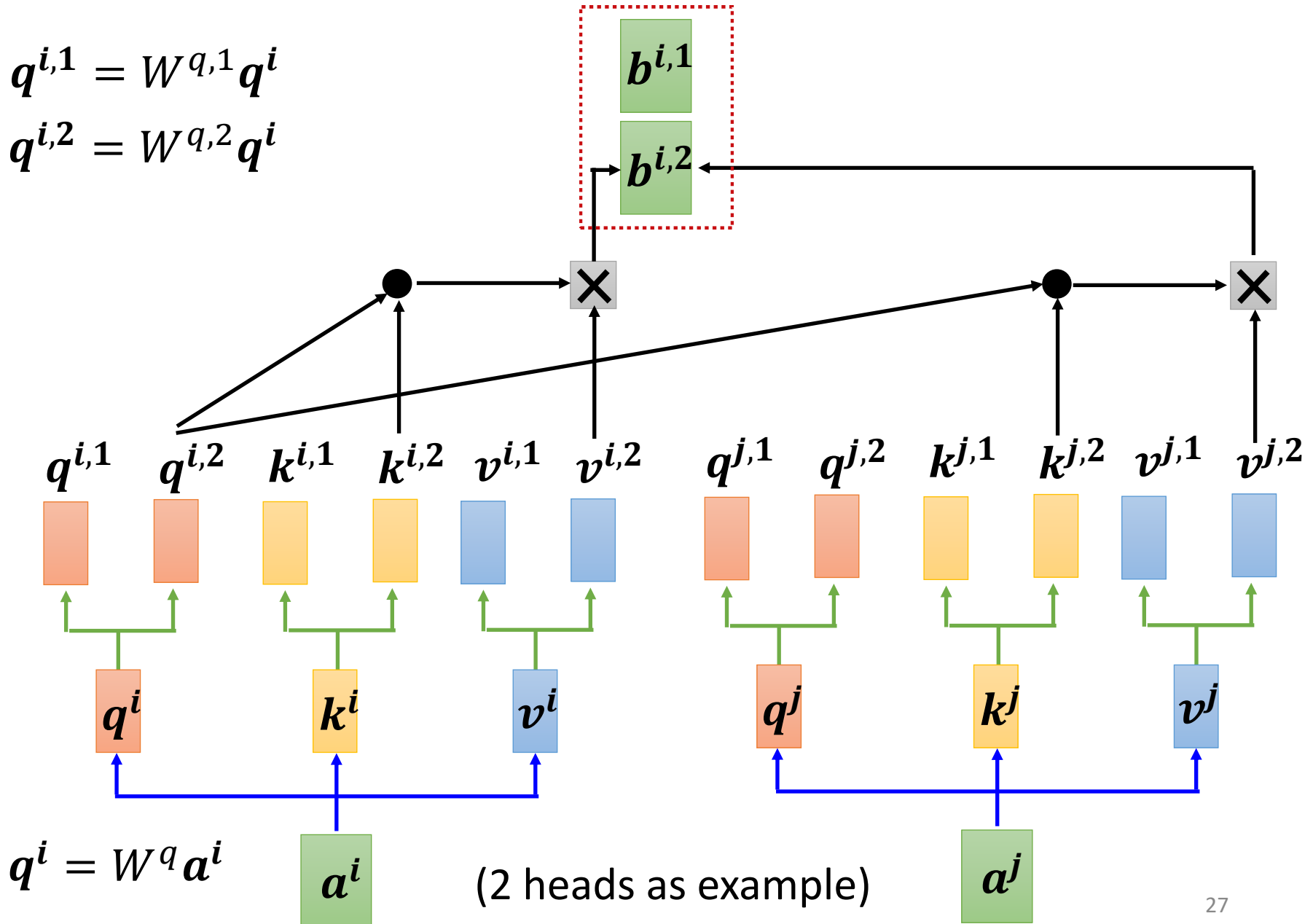


Multi-head Self-attention

Different types of relevance

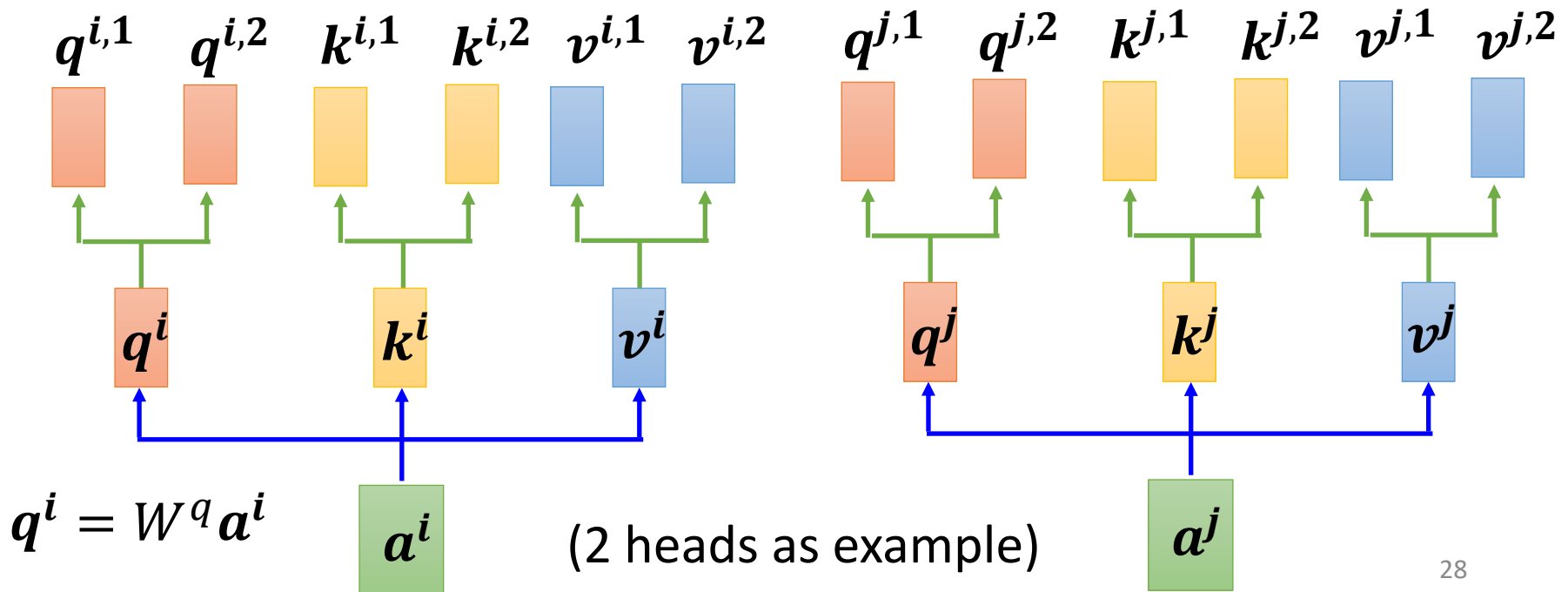
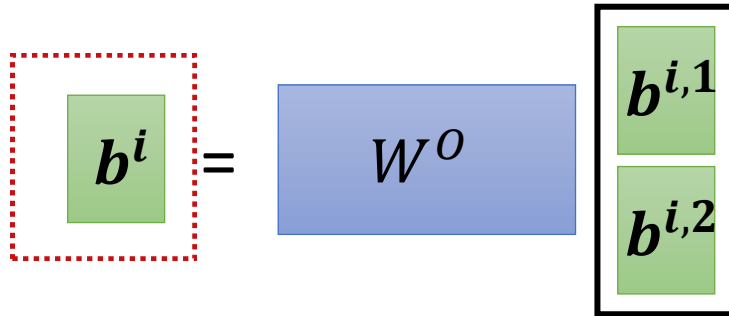
$$q^{i,1} = W^{q,1} q^i$$

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



Multi-head Self-attention

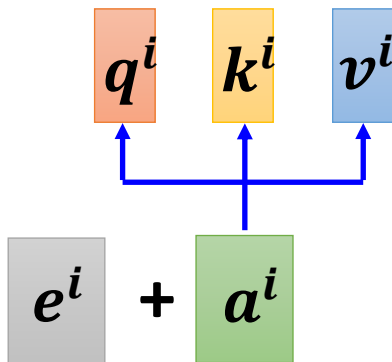
Different types of relevance



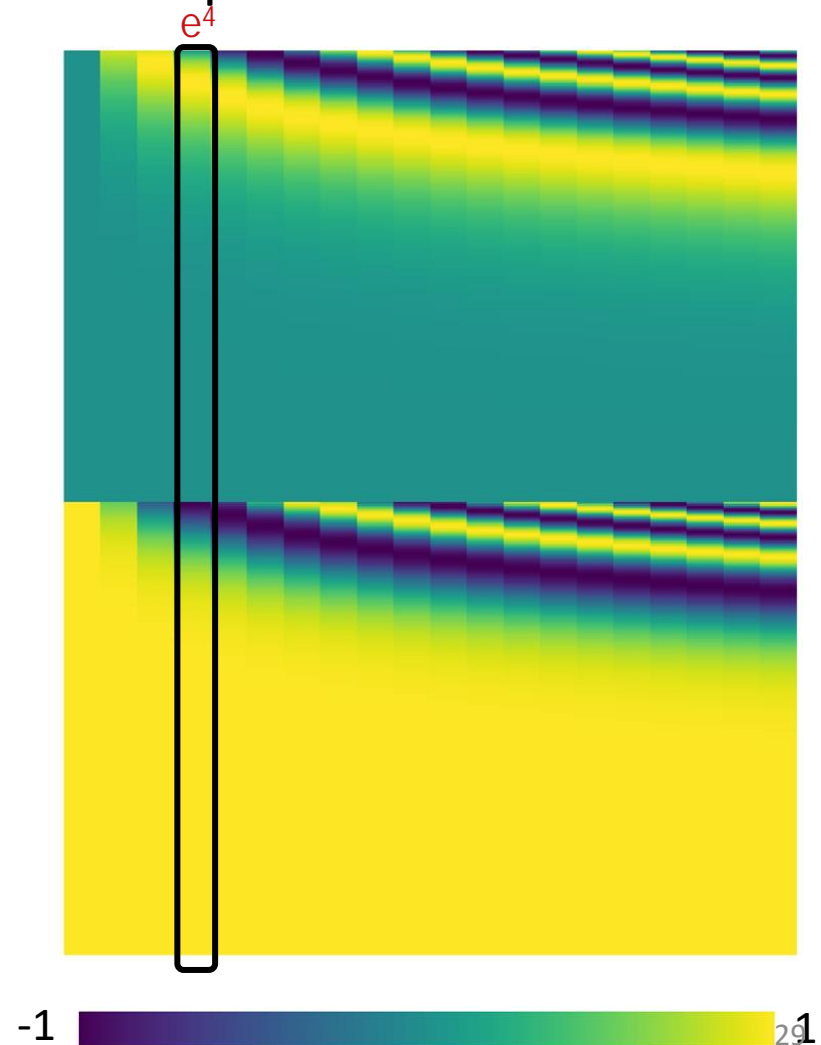
Attention is all you need里
positional encoding的设定

Positional Encoding

- No position information in self-attention.
- Each position has a unique positional vector e^i
- **hand-crafted**
- **learned from data**



Each column represents a positional vector e^i

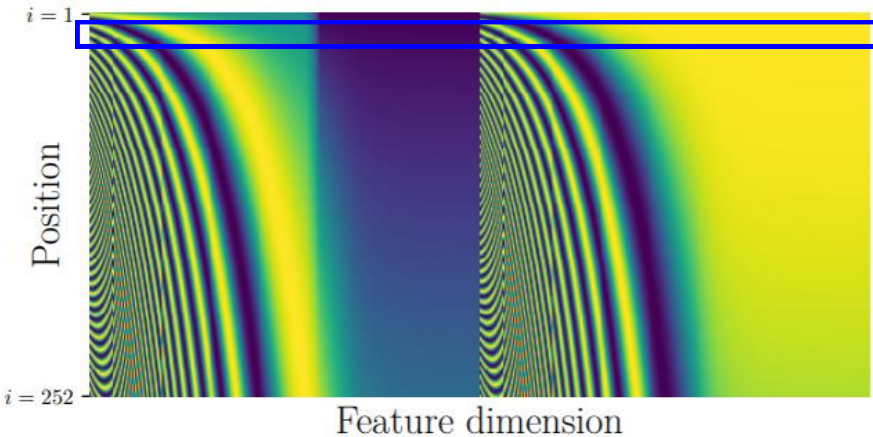


<https://arxiv.org/abs/2003.09229>

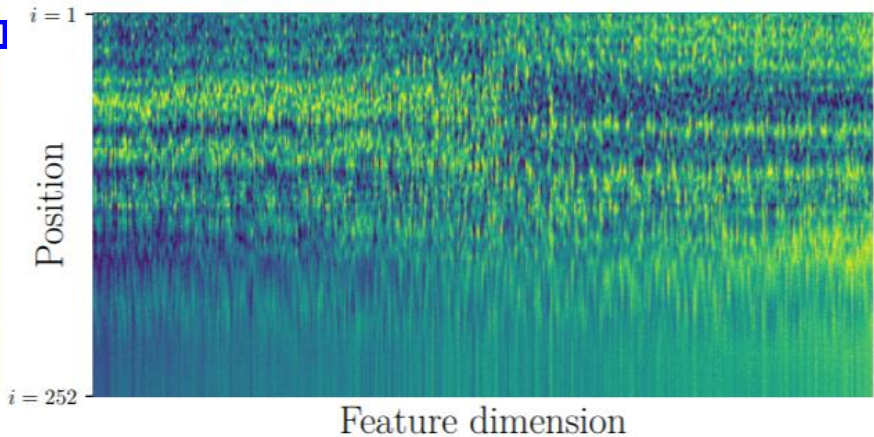
Table 1. Comparing position representation methods

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

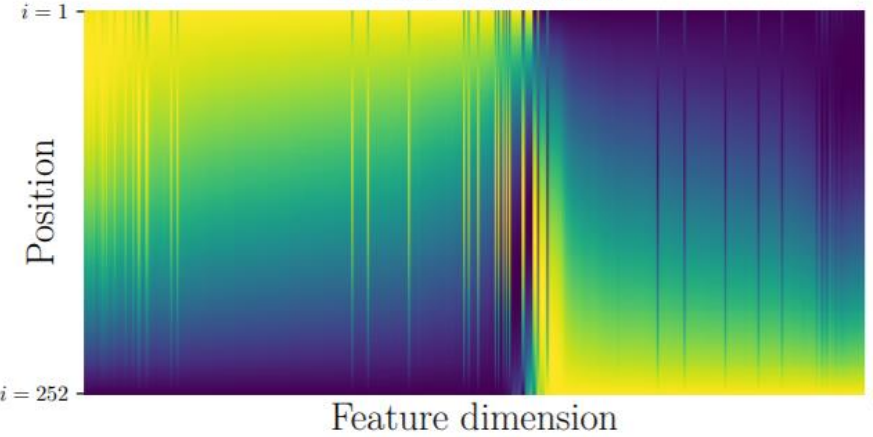
(a) Sinusoidal



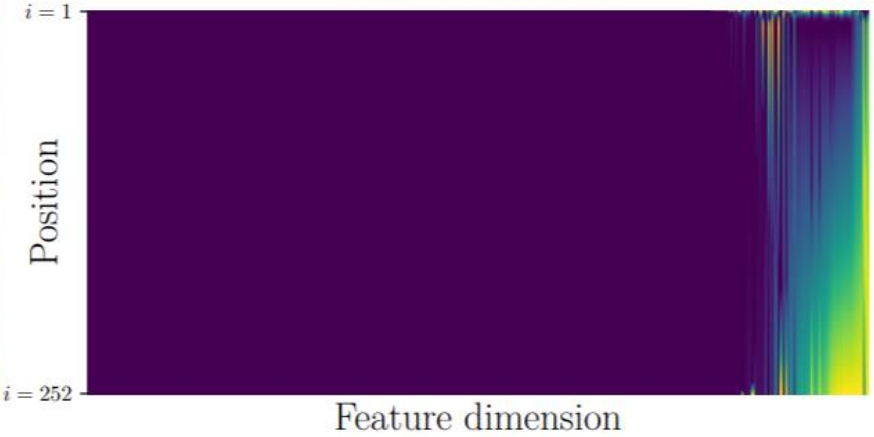
(b) Position embedding



(c) FLOATER



(d) RNN



Many applications ...



Transformer

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



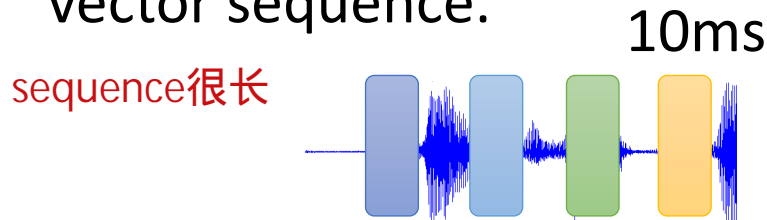
BERT

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

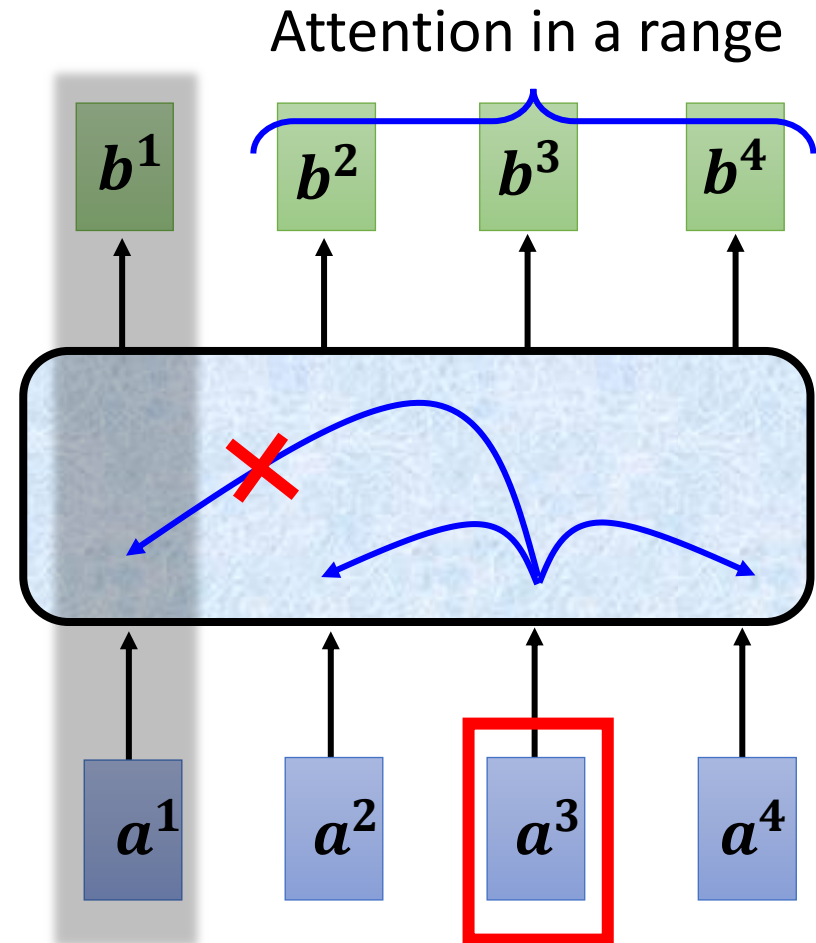
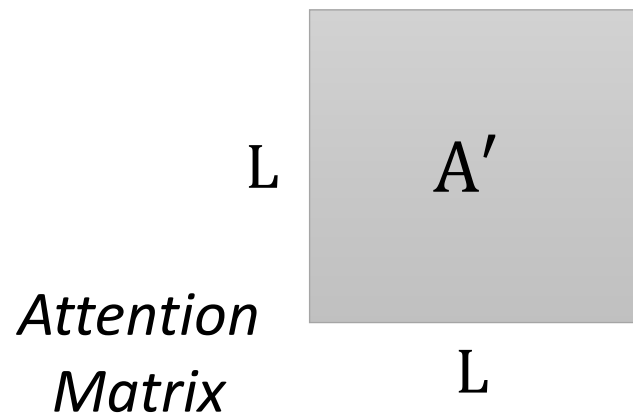
Widely used in Natural Language Processing (NLP)!

Self-attention for Speech 语音上的应用

Speech is a very long vector sequence.



If input sequence is length L

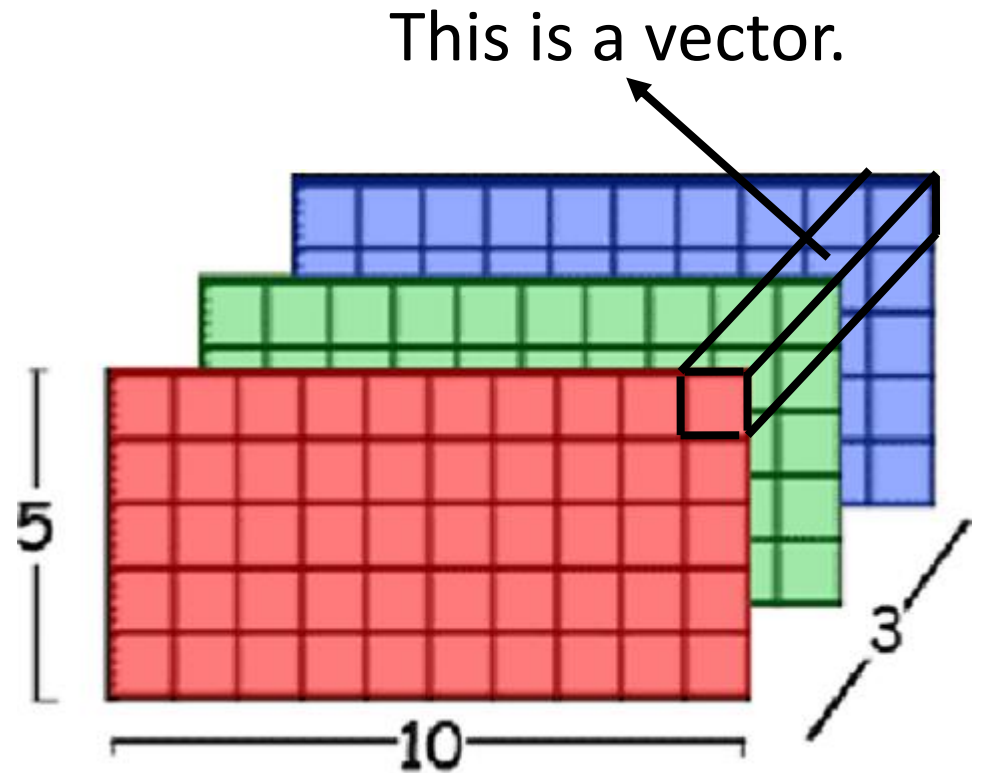
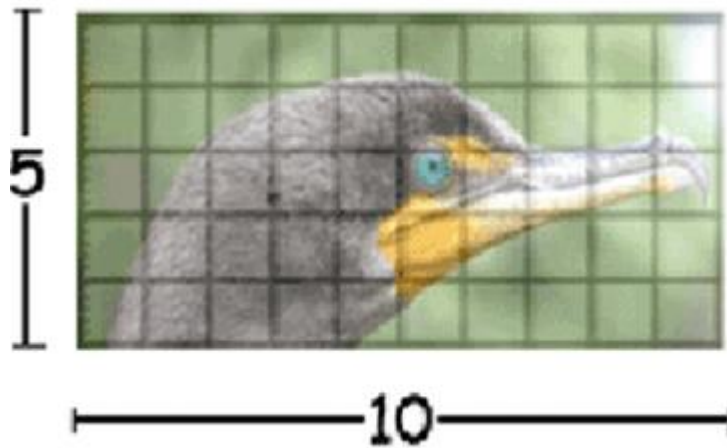


Truncated Self-attention

小范围内做self-attention

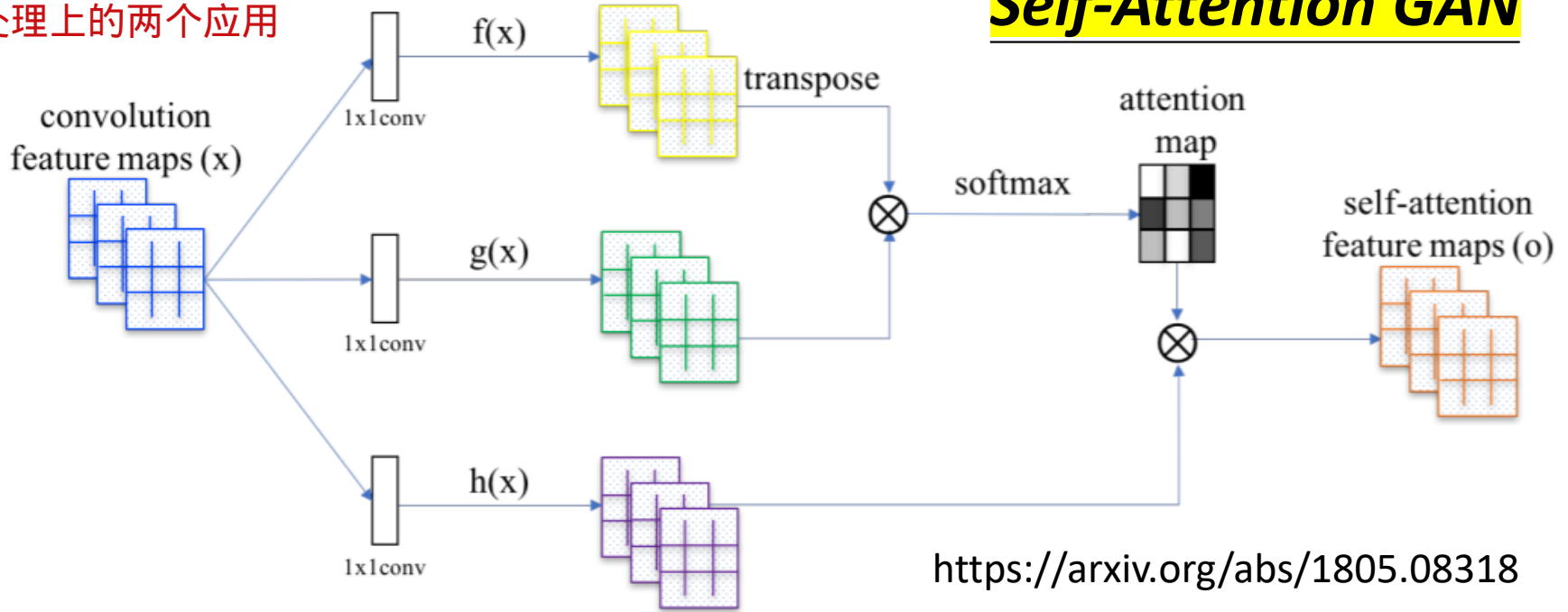
Self-attention for Image 影像处理上的应用

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

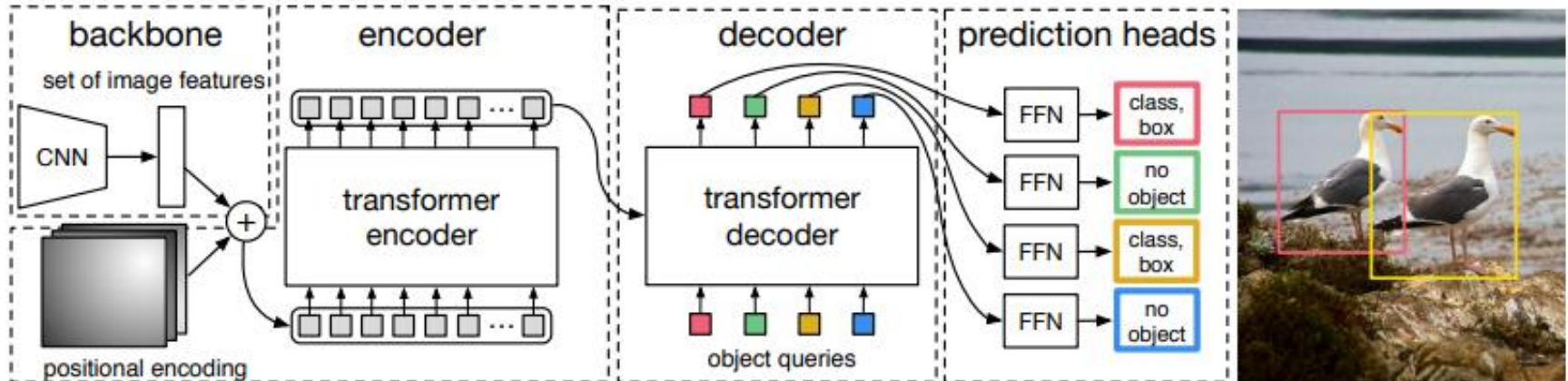


Source of image: https://www.researchgate.net/figure/Color-image-representation-and-RGB-matrix_fig15_282798184

self-attention在影像
处理上的两个应用

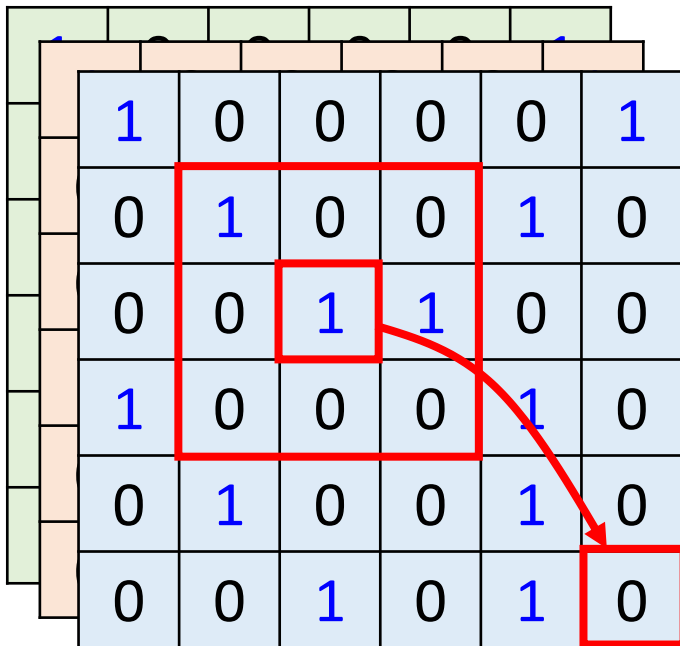


DEtection Transformer (DETR)



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

Self-attention v.s. CNN



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

- CNN is simplified self-attention.

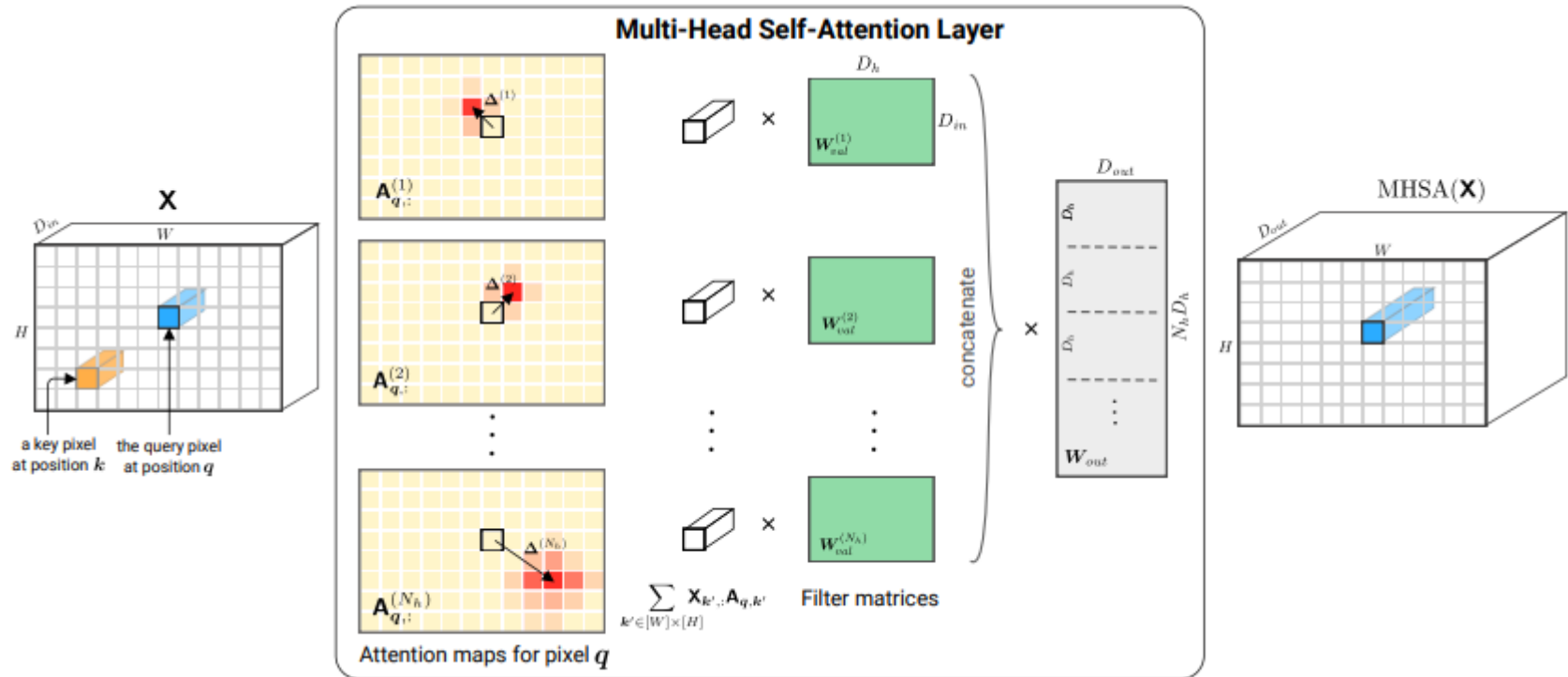
Self-attention: CNN with learnable receptive field

- Self-attention is the complex version of CNN.

Self-attention v.s. CNN

Self-attention

CNN



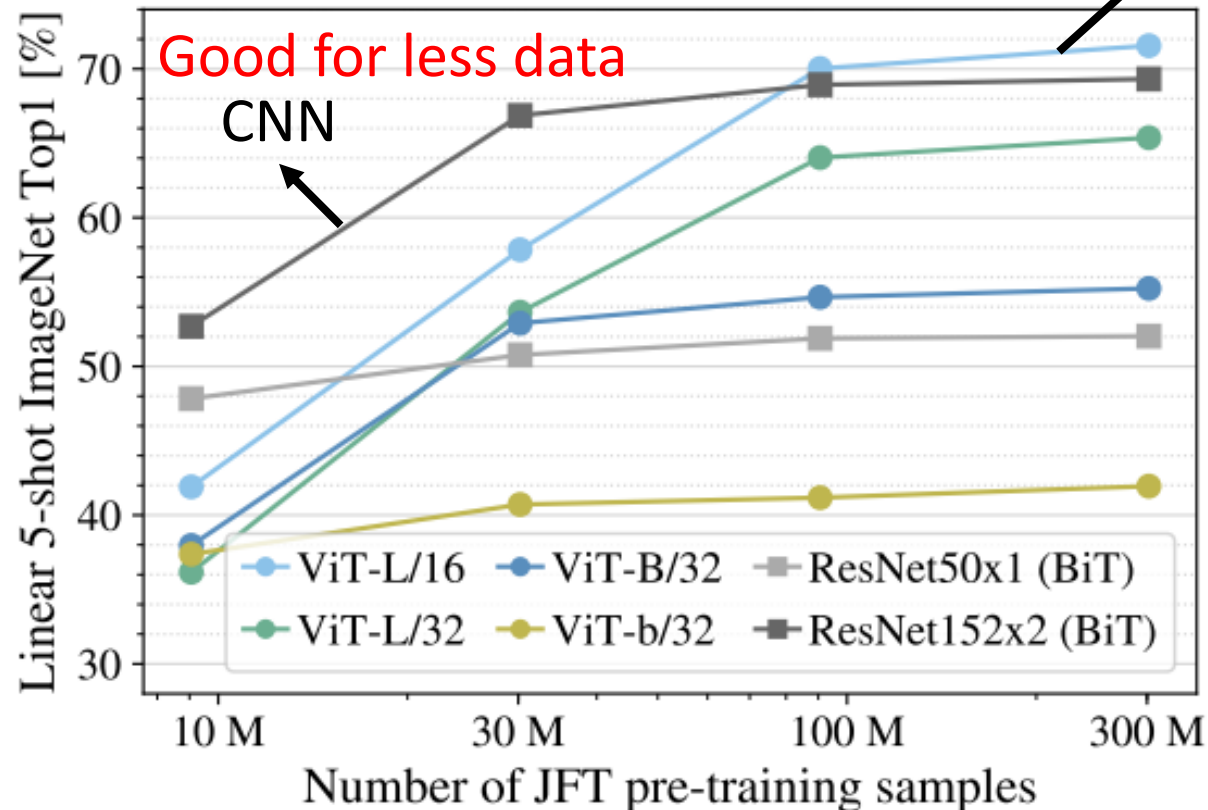
On the Relationship between Self-Attention and Convolutional Layers

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

Self-attention v.s. CNN

Good for more data

Self-attention

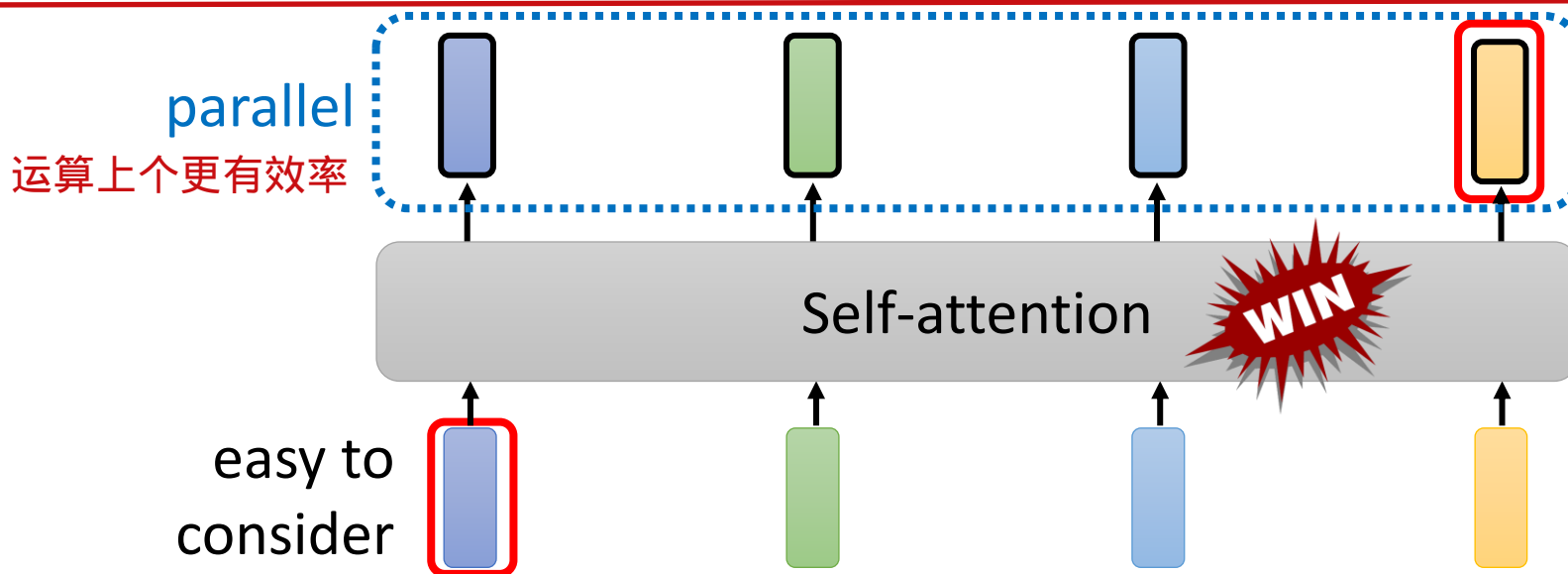
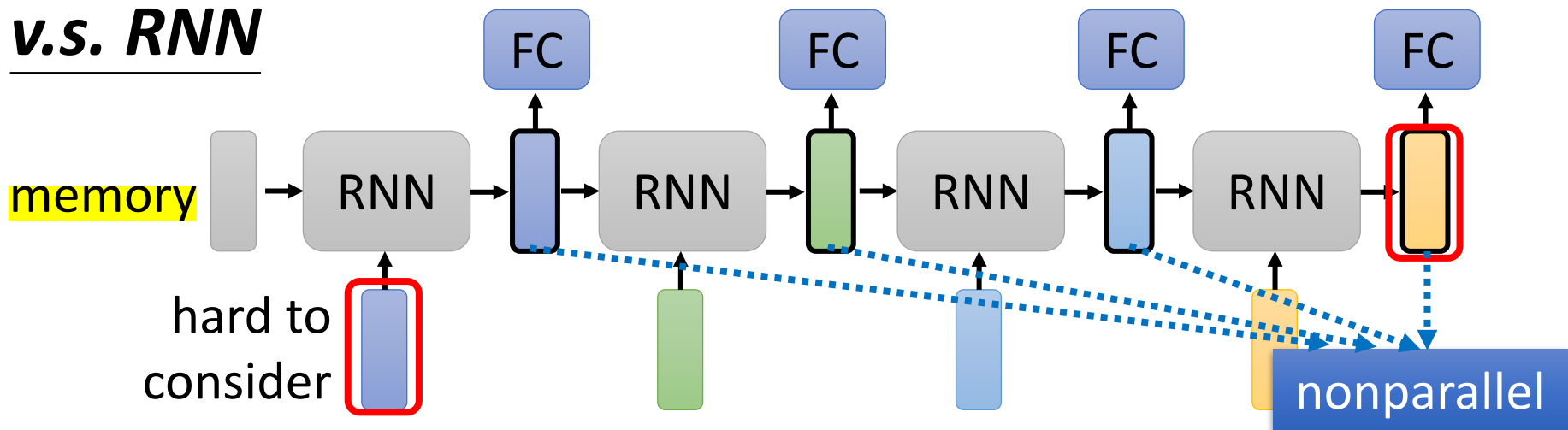


An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale

<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/2006.16236>

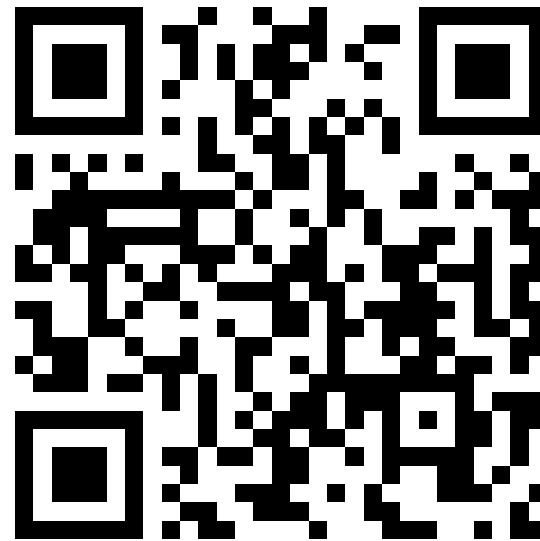
To learn more about RNN



<https://youtu.be/xCGidAeyS4M>

(in Mandarin)

2017年ML课RNN

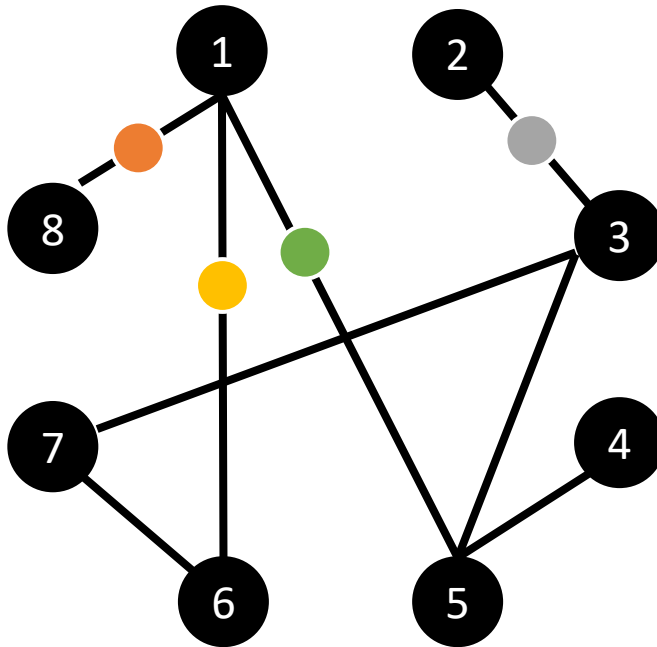


<https://youtu.be/Jjy6ER0bHv8>

(in English)

Graph可以被看作是set of vectors

Self-attention for Graph



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

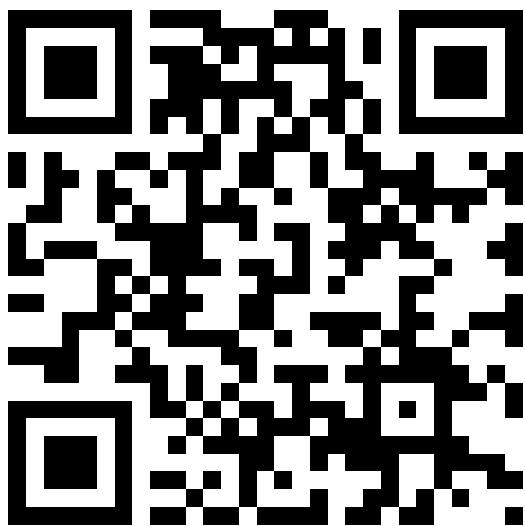
Attention Matrix

	1	2	3	4	5	6	7	8
1								
2								
3								
4								
5								
6								
7								
8							0	

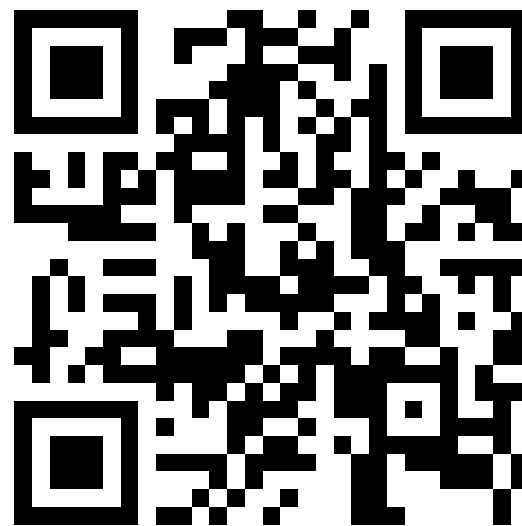
This is one type of **Graph Neural Network (GNN)**.

Self-attention for Graph

- To learn more about GNN ...



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



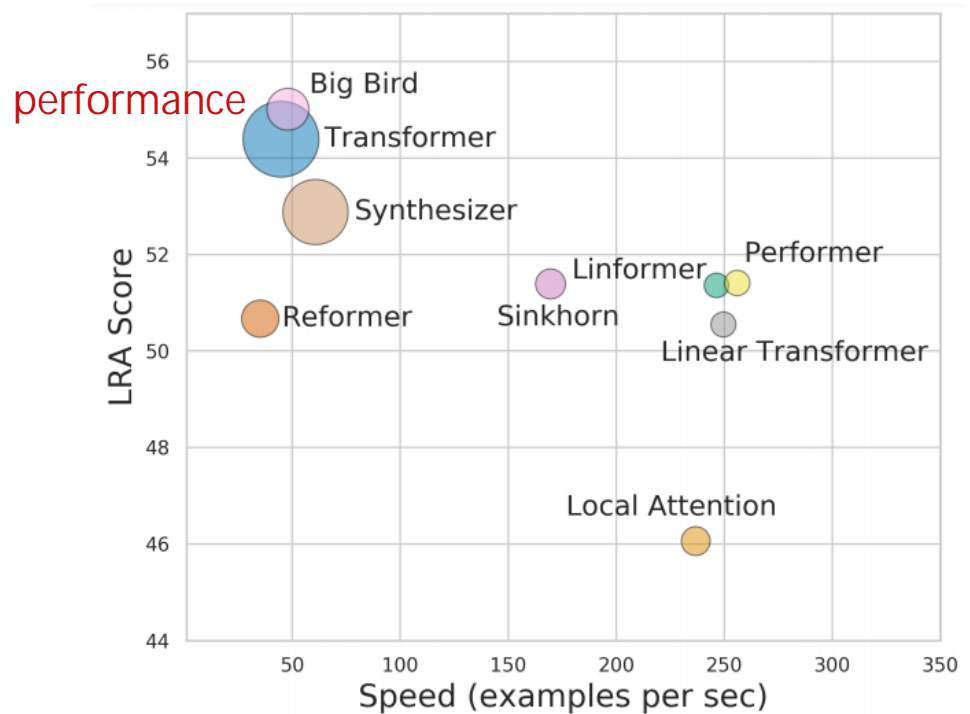
<https://youtu.be/M9ht8vsVEw8>
(in Mandarin)

To Learn More ...

self-attention的变形
self-attention的运算量很大，
怎么较少运算量是未来的一个课题。

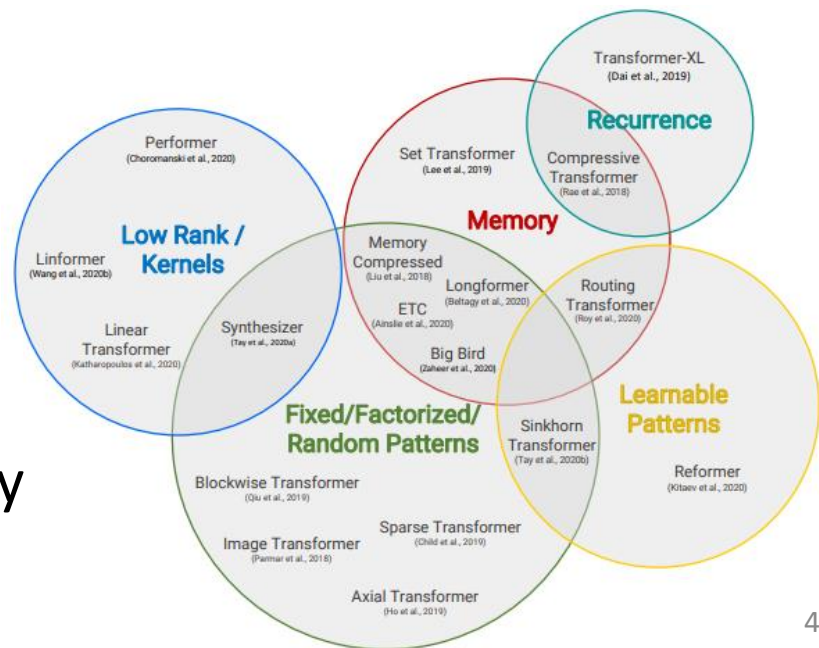
Long Range Arena: A Benchmark for Efficient Transformers

<https://arxiv.org/abs/2011.04006>



Efficient Transformers: A Survey

<https://arxiv.org/abs/2009.06732>



Q&A