# Self-attention

Hung-yi Lee
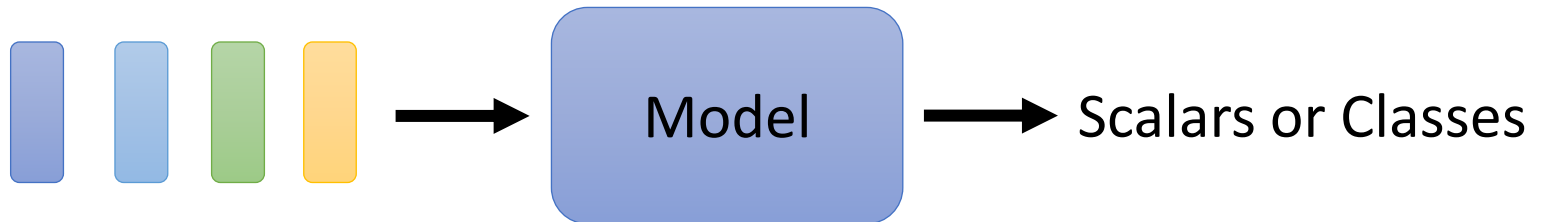
李宏毅

# Sophisticated Input

- Input is **a vector**

Model → Scalar or Class

- Input is **a set of vectors**

Model → Scalars or Classes

(may change length)

# Vector Set as Input

this    is    a    cat

向量长度与现有词汇数量相等
每个维度对应到一个词汇
没有联系(独立)

## *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*

run
jump

dog
rabbit
cat

tree
flower

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

# Vector Set as Input

10ms

1s → 100 frames

25ms

frame

400 sample points (16KHz)
39-dim MFCC
80-dim filter bank output

# Vector Set as Input

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



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 …… ]

C = [ 0   1   0   0   0 …… ]

O = [ 0   0   1   0   0 …… ]

⋮

One-hot vector

# *What is the output?*

- Each vector has a label. 输入与输出长度一样



N → Model → N

## *Example Applications*

I  saw  a  saw

↓   ↓   ↓   ↓

N   V  DET  N

**POS tagging**

词性标注



a   a   b   b

**HW2**



not

buy                    buy

# *What is the output?*
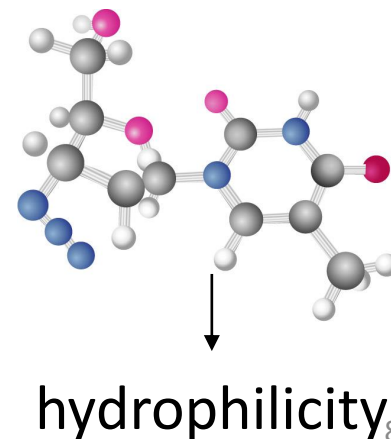
- Each vector has a label.



N → Model → N

- The whole sequence has a label.



Model →

## *Example Applications*

this is good

Sentiment analysis

一整个句子也一个 label

positive

**HW4**

speaker

hydrophilicity

8

# *What is the output?*

- Each vector has a label.    <span style="color:red">focus of this lecture</span>



- The whole sequence has a label.



- Model decides the number of labels itself.    seq2seq



**Translation (HW5)**

第一种. vector → label

# Sequence Labeling

Is it possible to consider the context?

FC can consider the neighbor

How to consider the whole sequence?

a window covers the whole sequence?

FC  Fully-connected

window
考虑 旁边

sequence
长度有长
有短

FC      FC      FC      FC

I       saw      a      saw

对 FC 来说, 两个 saw 完全一样

# *Self-attention*



with context

Self-attention

获得整个 sequence 的资讯

https://arxiv.org/abs/1706.03762

12

# *Self-attention*



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

X or a 都可以

# *Self-attention*



Find the relevant vectors in a sequence

找到每个a与a'的相关联程度

# Self-attention

## Dot-product

$$\alpha = \boldsymbol{q} \cdot \boldsymbol{k}$$

## Additive

$$\alpha$$

# *Self-attention*

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

$\alpha_{1,2}$ $\qquad$ $\alpha_{1,3}$ $\qquad$ $\alpha_{1,4}$

attention score

$\boldsymbol{q^1}$ query $\qquad$ $\boldsymbol{k^2}$ key $\qquad$ $\boldsymbol{k^3}$ $\qquad$ $\boldsymbol{k^4}$
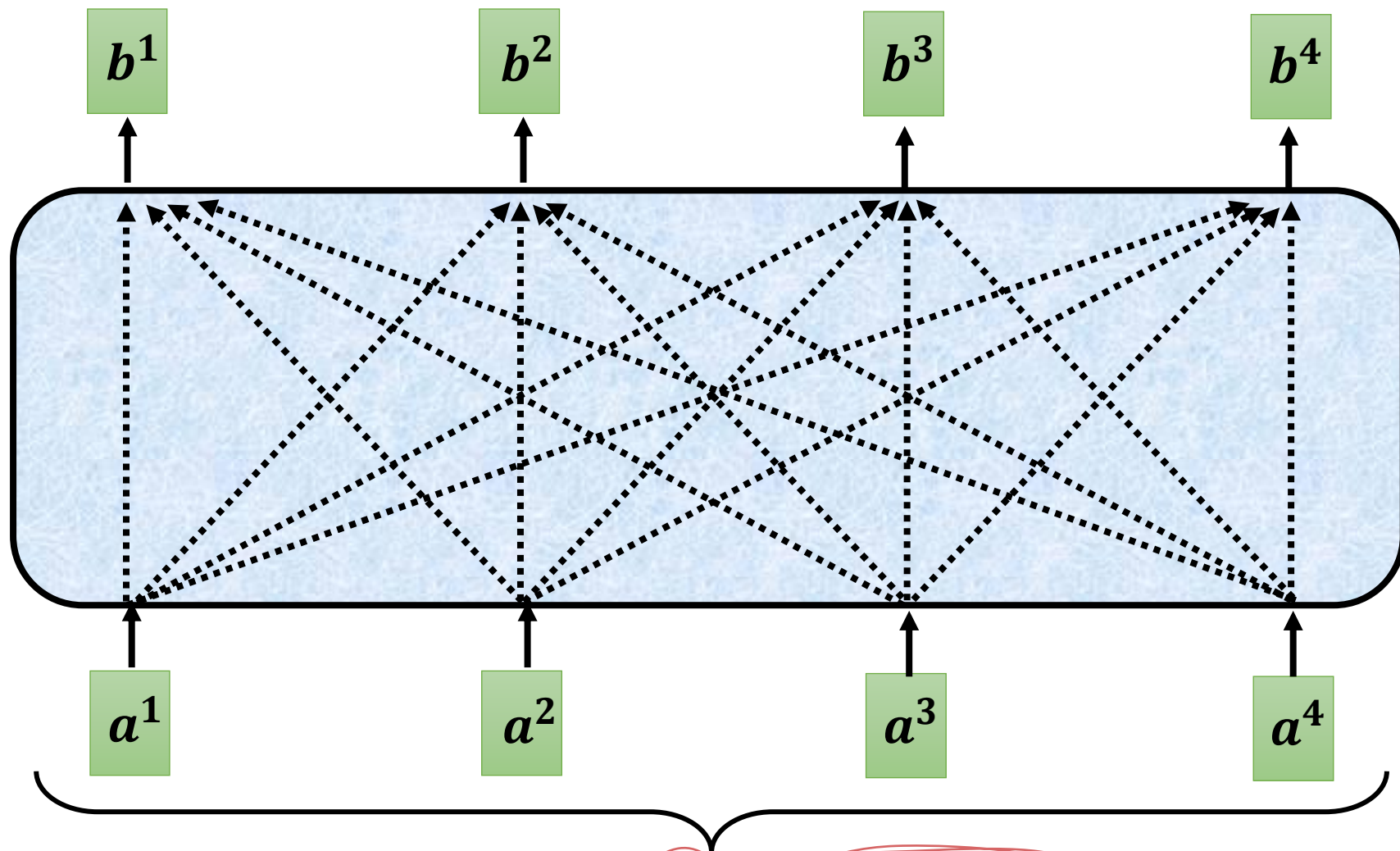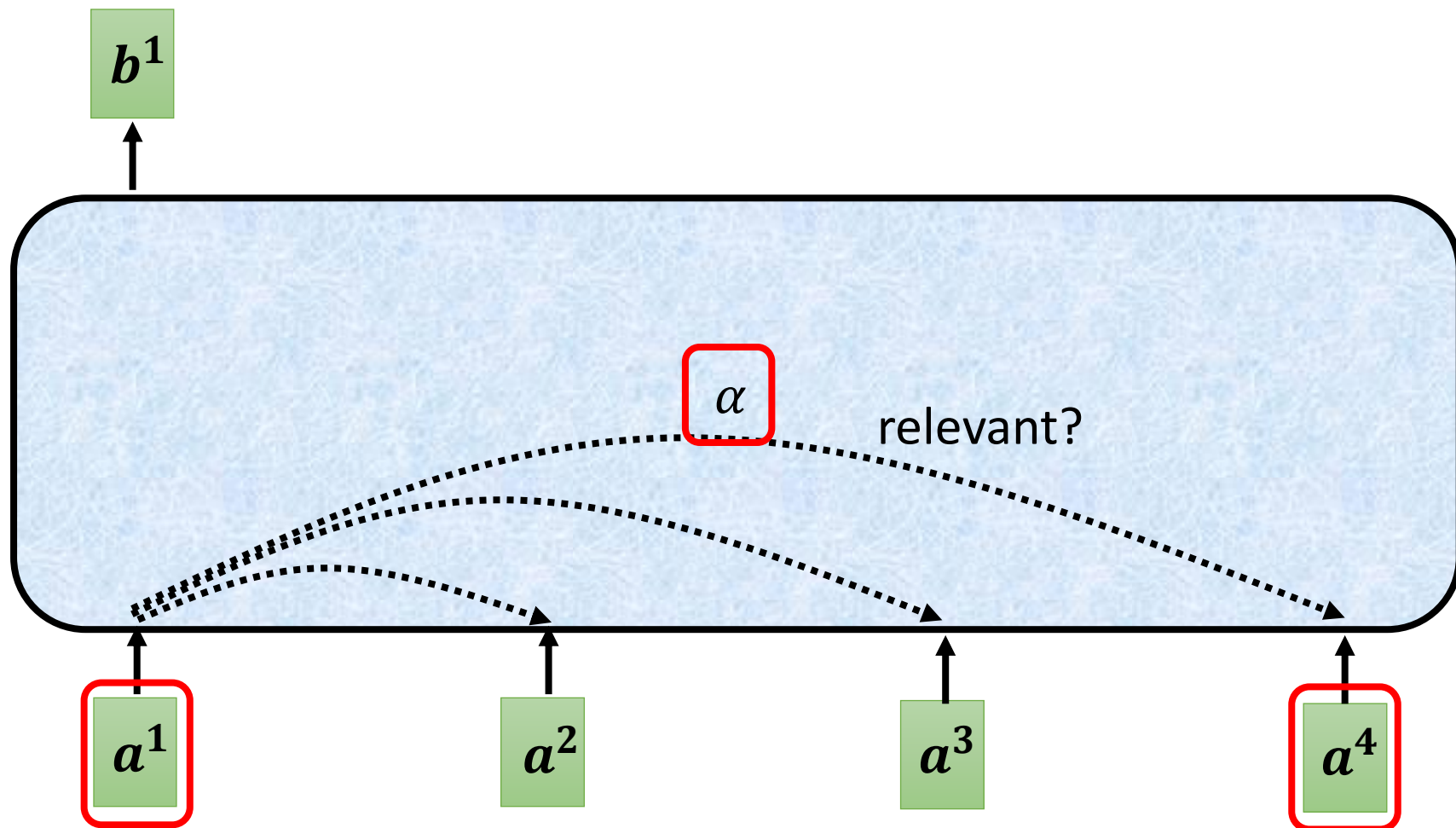
$\boldsymbol{a^1}$ $\qquad$ $\boldsymbol{a^2}$ $\qquad$ $\boldsymbol{a^3}$ $\qquad$ $\boldsymbol{a^4}$

$$\boldsymbol{q^1} = W^q \boldsymbol{a^1} \qquad \boldsymbol{k^2} = W^k \boldsymbol{a^2} \qquad \boldsymbol{k^3} = W^k \boldsymbol{a^3} \qquad \boldsymbol{k^4} = W^k \boldsymbol{a^4}$$

# Self-attention

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



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

$$k^2 = W^k a^2 \qquad k^3 = W^k a^3 \qquad k^4 = W^k a^4$$

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

17

# Self-attention

Extract information based on attention scores

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



$$v^1 = W^v a^1 \qquad v^2 = W^v a^2 \qquad v^3 = W^v a^3 \qquad v^4 = W^v a^4$$

# *Self-attention*

$b^1, b^2, b^3, b^4$ 一次同时被计算出，而非依次

parallel



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

# Self-attention

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



$k^1 = a^1 \cdot W^k$

$v^1 = a^1 \cdot W^v$

$q^2 = a_2 \cdot W^q$

$k^2 = a_2 \cdot W^k$

$\alpha'_{2,1} = k^1 \cdot q^2$

$b^2 = \sum \cdots \alpha'_{2,1} v^1 \cdots$

# Self-attention

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

$$\begin{array}{|c|c|c|c|}\hline q^1 & q^2 & q^3 & q^4 \\\hline\end{array} = \boxed{W^q} \begin{array}{|c|c|c|c|}\hline a^1 & a^2 & a^3 & a^4 \\\hline\end{array}$$

$Q$   I

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

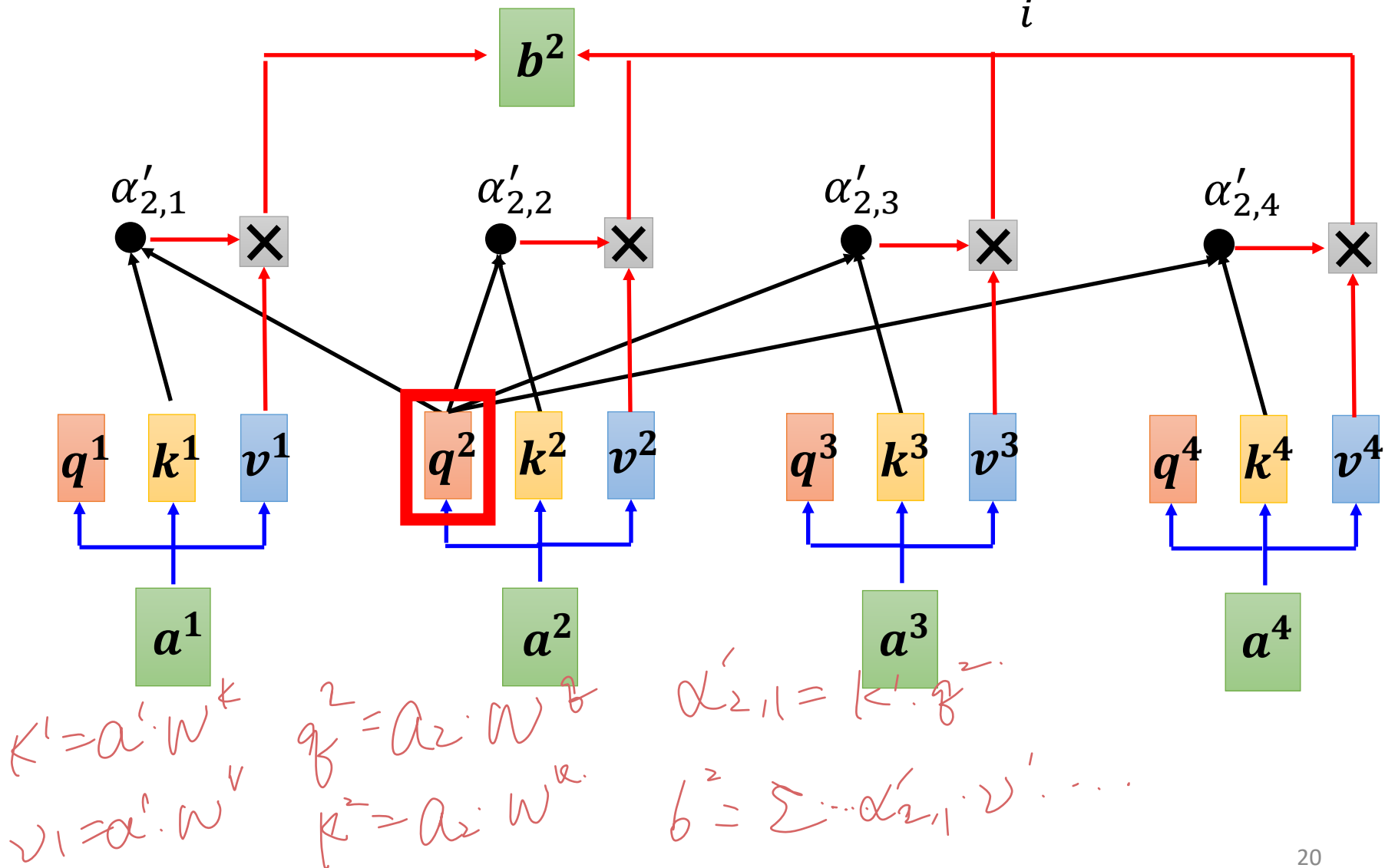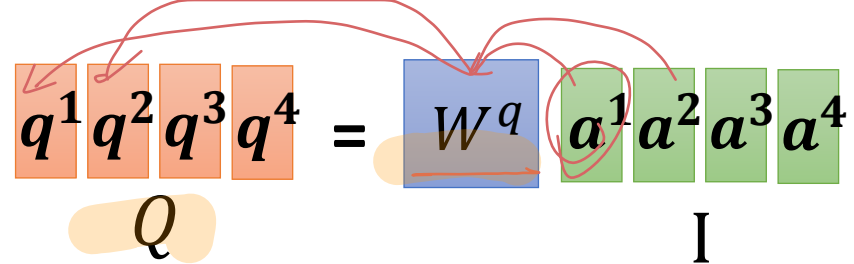$$\begin{array}{|c|c|c|c|}\hline k^1 & k^2 & k^3 & k^4 \\\hline\end{array} = \boxed{W^k} \begin{array}{|c|c|c|c|}\hline a^1 & a^2 & a^3 & a^4 \\\hline\end{array}$$

$K$   I

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

$$\begin{array}{|c|c|c|c|}\hline v^1 & v^2 & v^3 & v^4 \\\hline\end{array} = \boxed{W^v} \begin{array}{|c|c|c|c|}\hline a^1 & a^2 & a^3 & a^4 \\\hline\end{array}$$

$V$   I

$q^1$ $k^1$ $v^1$   $q^2$ $k^2$ $v^2$   $q^3$ $k^3$ $v^3$   $q^4$ $k^4$ $v^4$

$a^1$   $a^2$   $a^3$   $a^4$

# *Self-attention*

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

$$\alpha_{1,3} = \boxed{k^3}\,\boxed{q^1} \qquad \alpha_{1,4} = \boxed{k^4}\,\boxed{q^1}$$

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

# _Self-attention_

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

$$\alpha_{1,3} = \boxed{k^3}\ \boxed{q^1} \qquad \alpha_{1,4} = \boxed{k^4}\ \boxed{q^1}$$

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



$\alpha_{2,1}$  $\alpha_{2,2}$  $\alpha_{2,3}$  $\alpha_{2,4}$

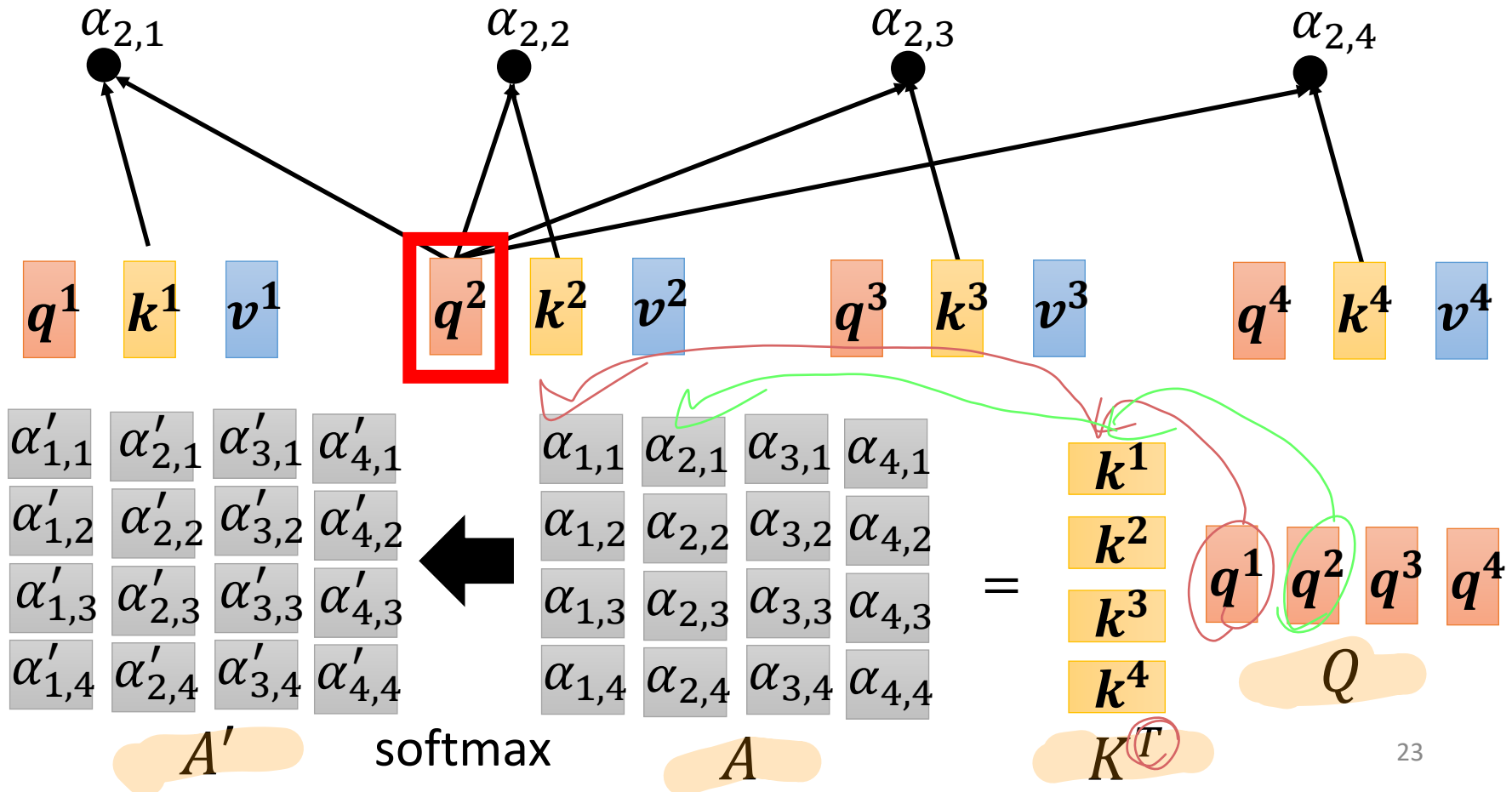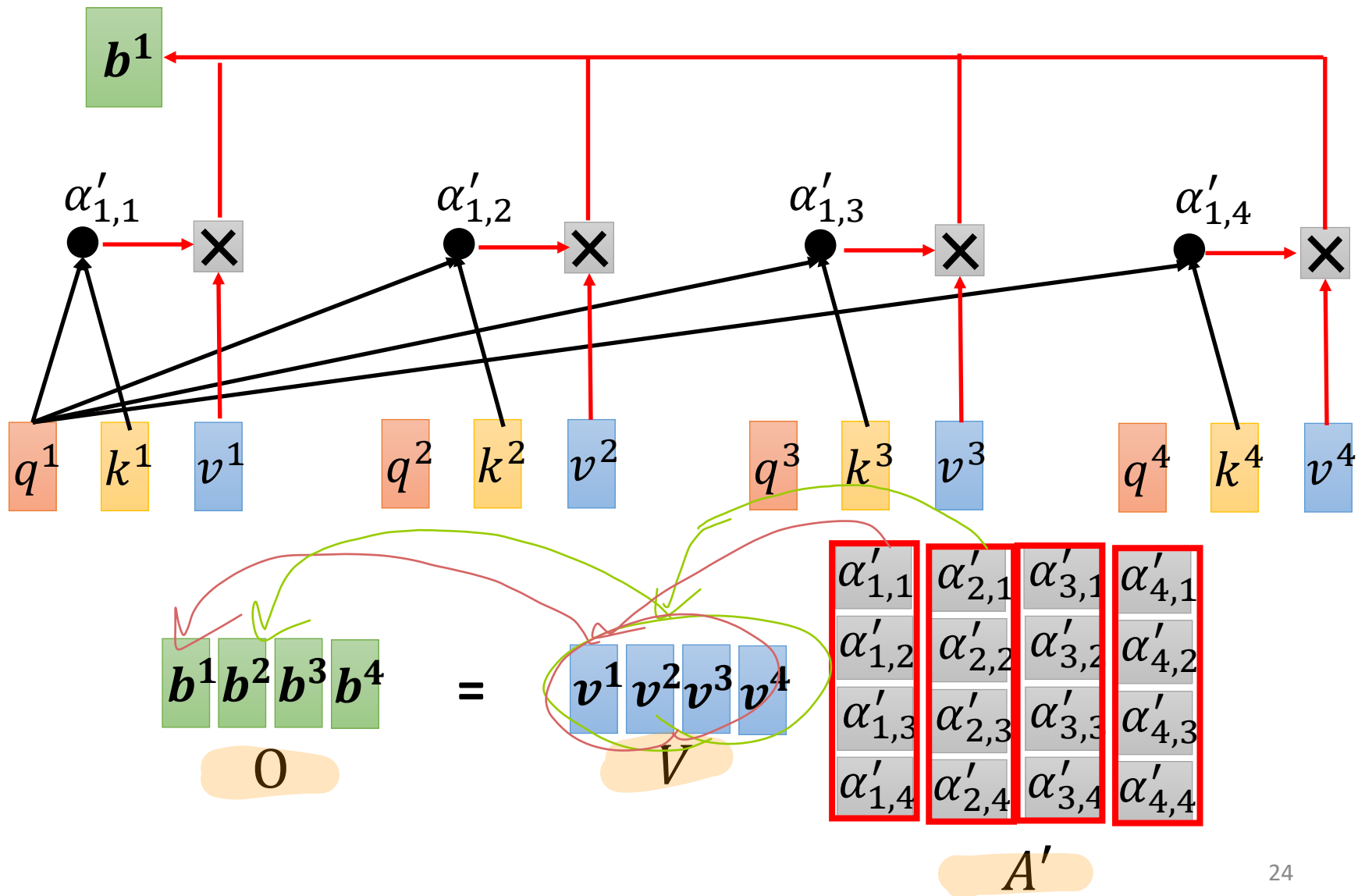$q^1\ k^1\ v^1\quad q^2\ k^2\ v^2\quad q^3\ k^3\ v^3\quad q^4\ k^4\ v^4$

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

$A'$ \qquad softmax \qquad $A$ \qquad $K^T$ \qquad $Q$

23

# *Self-attention*



$b^1$

$\alpha'_{1,1}$  $\times$  $\alpha'_{1,2}$  $\times$  $\alpha'_{1,3}$  $\times$  $\alpha'_{1,4}$  $\times$

$q^1$ $k^1$ $v^1$ $q^2$ $k^2$ $v^2$ $q^3$ $k^3$ $v^3$ $q^4$ $k^4$ $v^4$

$$\boldsymbol{b^1 b^2 b^3 b^4} \;=\; \boldsymbol{v^1 v^2 v^3 v^4} \begin{bmatrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{bmatrix}$$

$O$ $\qquad$ $V$ $\qquad\qquad\qquad$ $A'$

24

# *Self-attention*

$$Q = W^q I$$

$$K = W^k I$$

$$V = W^v I$$

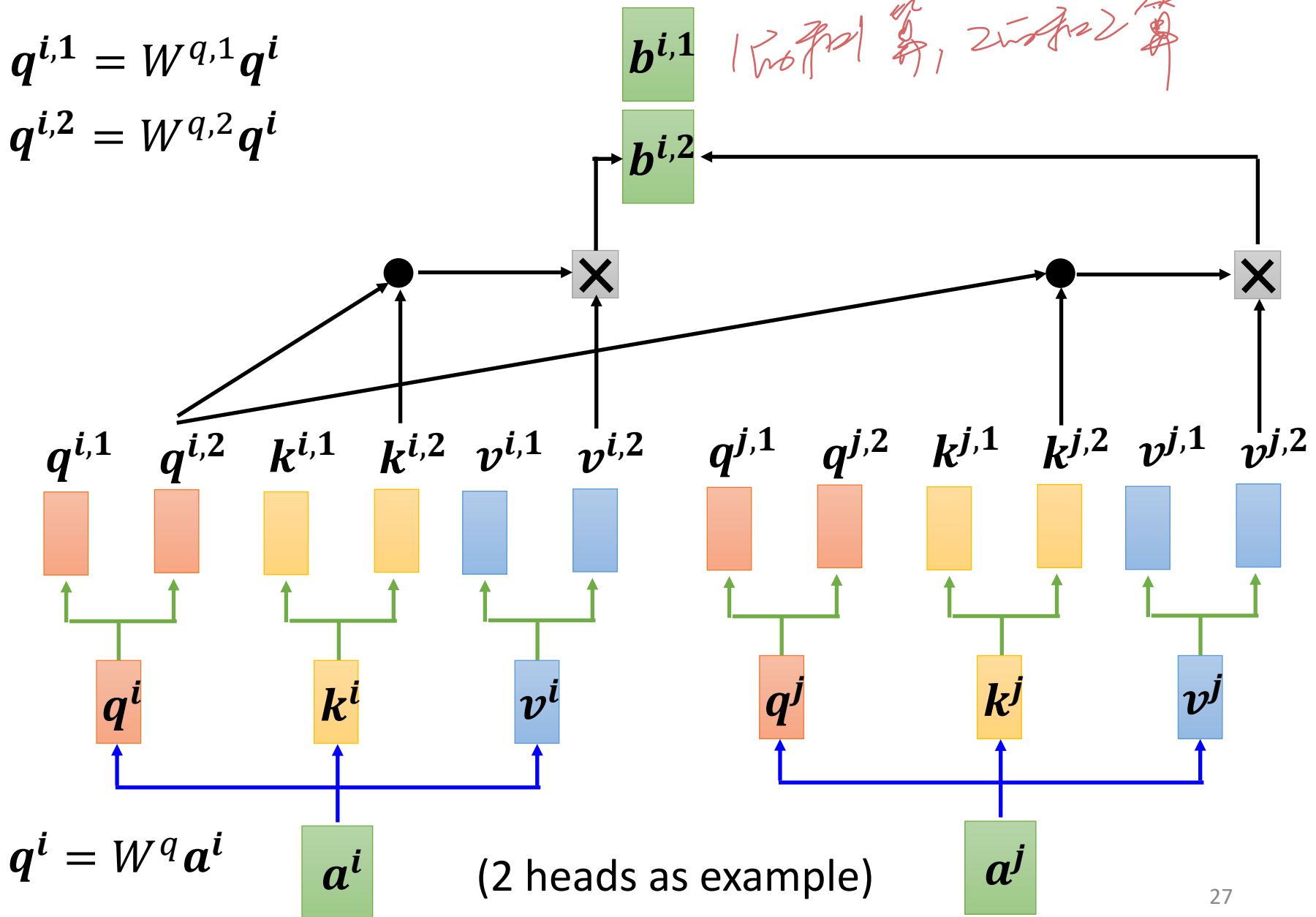Parameters to be learned

$$A = K^T Q$$

A' ← A

Attention Matrix

$$O = V A'$$

# *Multi-head Self-attention* Different types of relevance

$$q^{i,1} = W^{q,1} q^i$$
$$q^{i,2} = W^{q,2} q^i$$

$b^{i,1}$

两种不同的相关性

$q、k、v$ 分别乘上两个不同的矩阵，即可将 $q、k、v$ 分别拆为两个 head

$q^{i,1}$   $q^{i,2}$   $k^{i,1}$   $k^{i,2}$   $v^{i,1}$   $v^{i,2}$   $q^{j,1}$   $q^{j,2}$   $k^{j,1}$   $k^{j,2}$   $v^{j,1}$   $v^{j,2}$

$q^i$   $k^i$   $v^i$   $q^j$   $k^j$   $v^j$

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

$a^i$

用 $q$ 找相关的 $k$，把 相关不同，一种

(2 heads as example)

$a^j$

# Multi-head Self-attention   Different types of relevance

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

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

$b^{i,1}$

$b^{i,2}$

$q^{i,1}$  $q^{i,2}$  $k^{i,1}$  $k^{i,2}$  $v^{i,1}$  $v^{i,2}$    $q^{j,1}$  $q^{j,2}$  $k^{j,1}$  $k^{j,2}$  $v^{j,1}$  $v^{j,2}$

$q^i$          $k^i$          $v^i$              $q^j$          $k^j$          $v^j$

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

$a^i$          (2 heads as example)          $a^j$

27

# *Multi-head Self-attention* <span style="color:red">Different types of relevance</span>

$$b^i = \boxed{W^O \begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix}}$$

$q^{i,1}$  $q^{i,2}$  $k^{i,1}$  $k^{i,2}$  $v^{i,1}$  $v^{i,2}$    $q^{j,1}$  $q^{j,2}$  $k^{j,1}$  $k^{j,2}$  $v^{j,1}$  $v^{j,2}$

$q^i$  $k^i$  $v^i$    $q^j$  $k^j$  $v^j$

$q^i = W^q a^i$    $a^i$    (2 heads as example)    $a^j$

# Positional Encoding

Each column represents a positional vector $e^i$

- No position information in self-attention.

- Each position has a unique positional vector $e^i$

- **hand-crafted**

- **learned from data**



$q^i$  $k^i$  $v^i$

$e^i + a^i$

-1          29 1

**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

30

# Many applications …

**_Transformer_**

https://arxiv.org/abs/1706.03762

**_BERT_**

https://arxiv.org/abs/1810.04805

Widely used in Natural Langue Processing (NLP)!

# Self-attention for Speech

取一部分做 self-ate.

Attention in a range

Speech is a very long vector sequence.

10ms

If input sequence is length L

$L \times L$

L    A′

Attention Matrix

L

$b^1$   $b^2$   $b^3$   $b^4$

$a^1$   $a^2$   $a^3$   $a^4$

**Truncated Self-attention**

# Self-attention for Image

輸入是一排向量

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

This is a vector.



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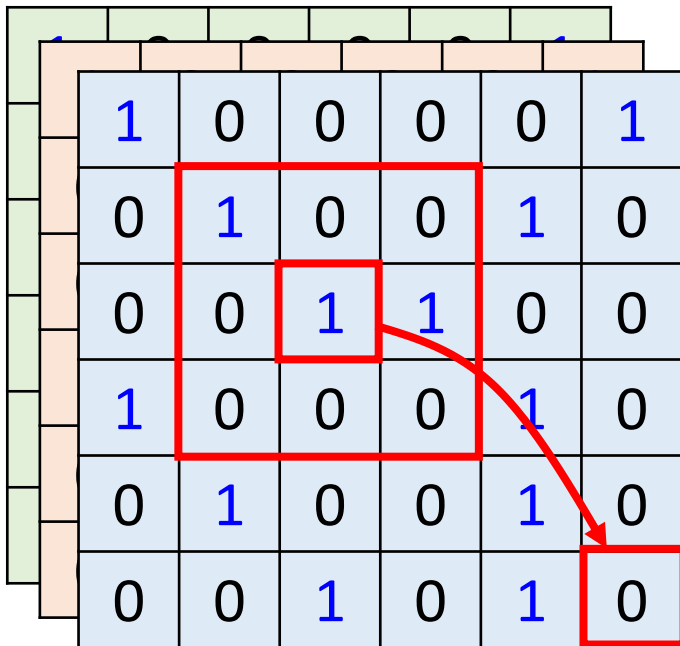
## Self-Attention GAN



convolution
feature maps (x)

1x1conv — f(x)

transpose

1x1conv — g(x)

1x1conv — h(x)

⊗ softmax → attention map

⊗ → self-attention feature maps (o)

https://arxiv.org/abs/1805.08318

## DEtection Transformer (DETR)



backbone — set of image features — CNN — positional encoding

encoder — transformer encoder

decoder — transformer decoder — object queries

prediction heads — FFN → class, box — FFN → no object — FFN → class, box — FFN → no object

https://arxiv.org/abs/2005.12872

# Self-attention v.s. CNN

CNN是简化版 self-attn.----



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

➢ CNN is simplified self-attention.

Self-attention: CNN with learnable receptive field

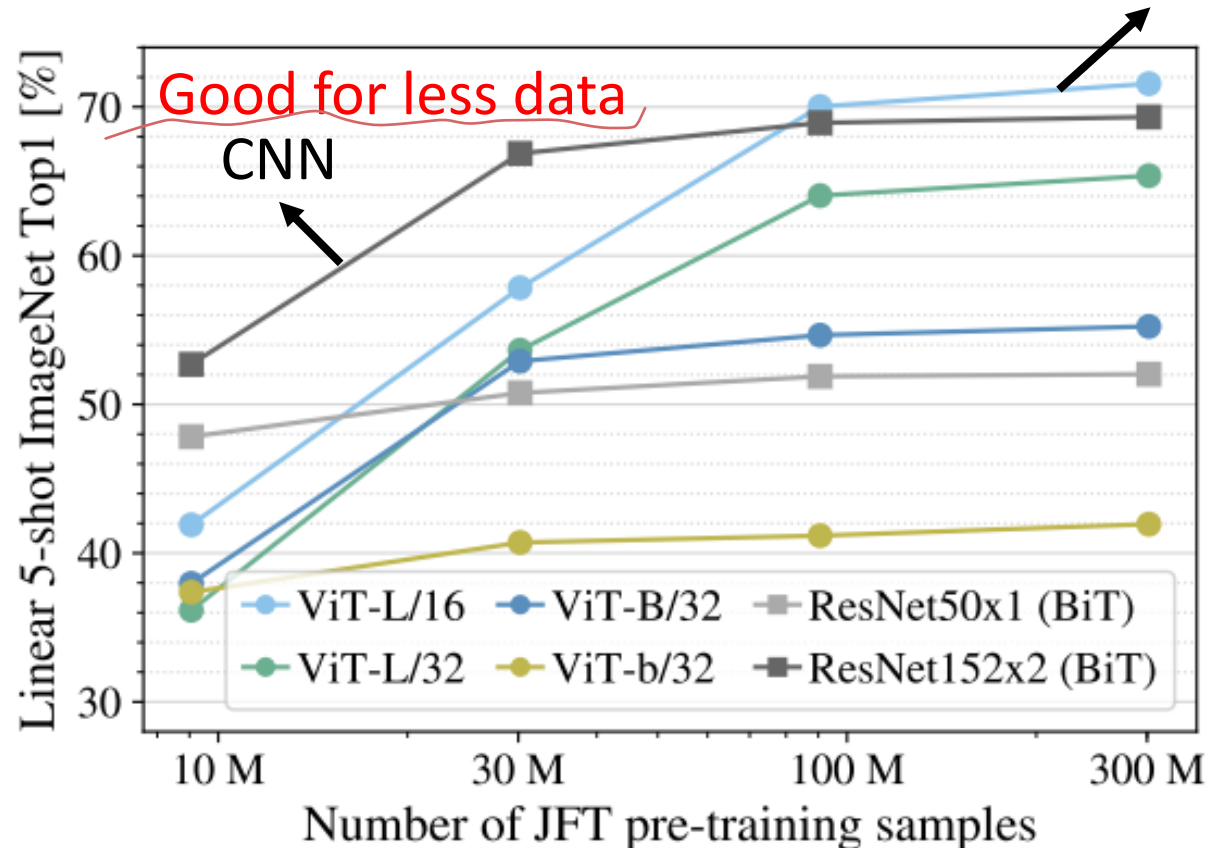➢ Self-attention is the complex version of CNN.

# Self-attention v.s. CNN



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

Good for less data

CNN



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

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

# *Self-attention v.s. RNN*

Recurrent Neural Network (RNN)



memory

hard to consider

nonparallel

parallel

Self-attention

**WIN**

easy to consider

Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention

# To learn more about RNN ......



https://youtu.be/xCGidAeyS4M

(in Mandarin)

https://youtu.be/Jjy6ER0bHv8

(in English)

# Self-attention for Graph
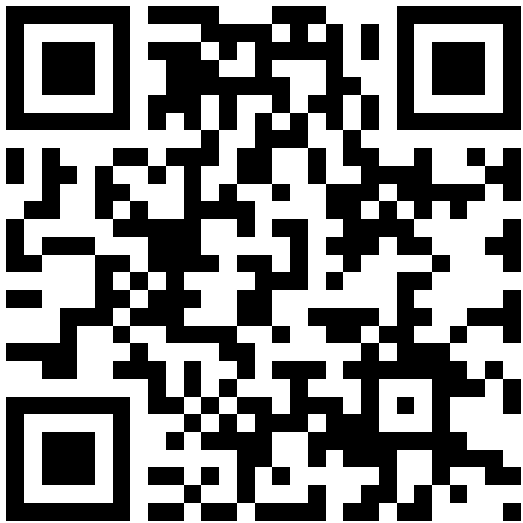
Consider **edge**: only attention
to connected nodes

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

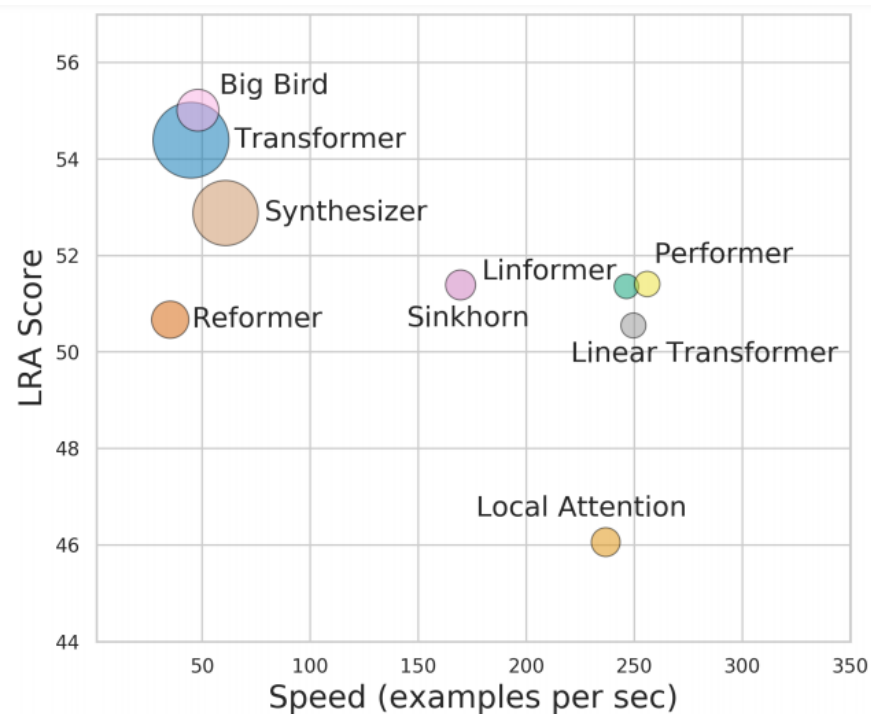# Self-attention for Graph

- To learn more about GNN ...



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



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

# *To Learn More …*



Long Range Arena: A
Benchmark for Efficient
Transformers

https://arxiv.org/abs/2011.04006



Efficient Transformers: A Survey

https://arxiv.org/abs/2009.06732