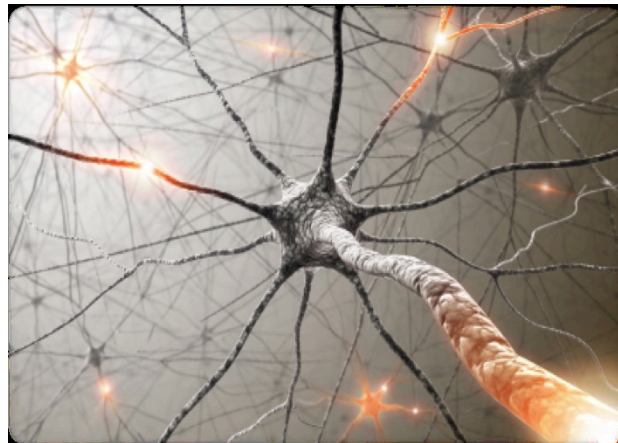


# 05 Neural Networks

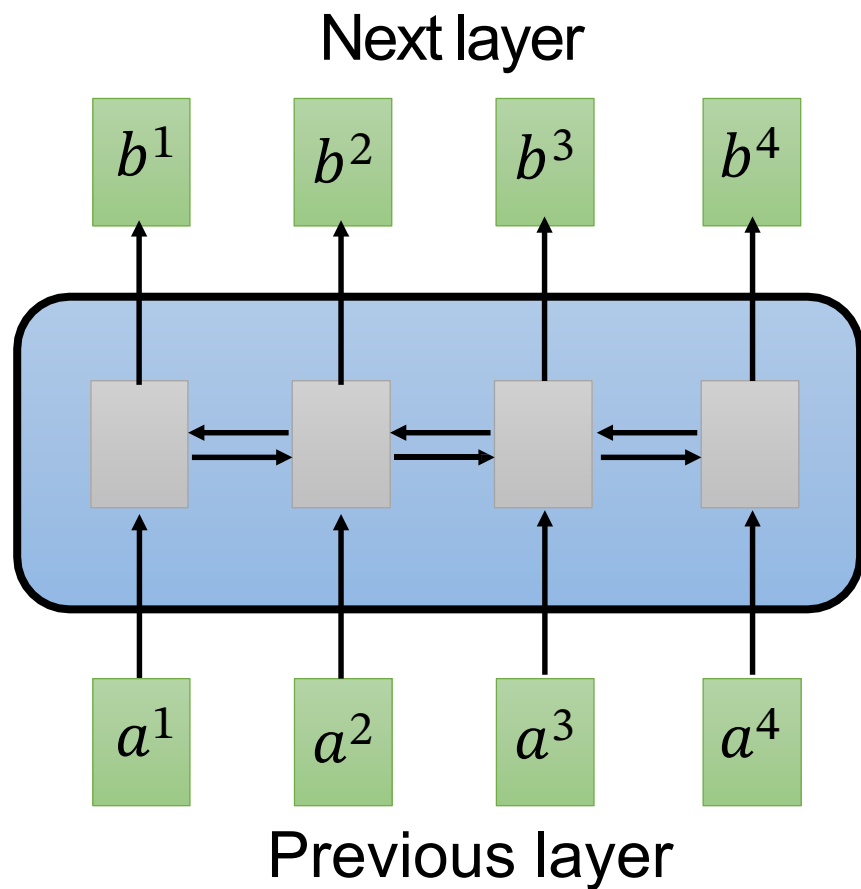
Transformer: Seq2seq model with “Self-attention”



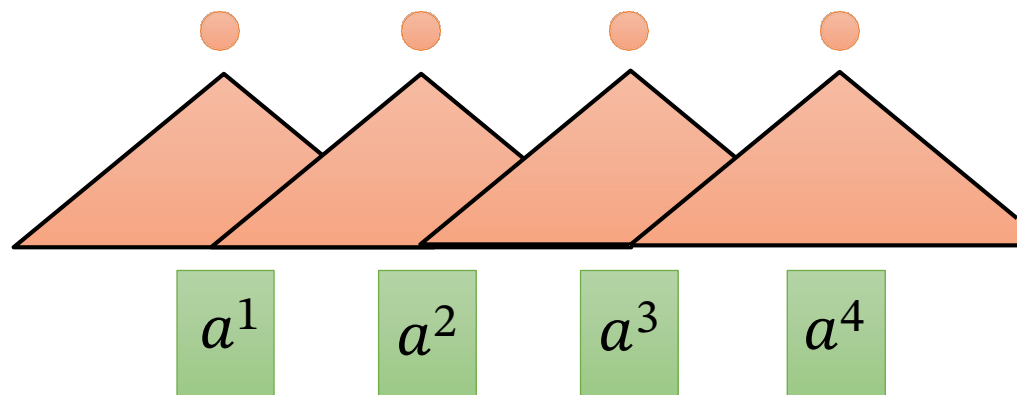


Transformer

# Sequence

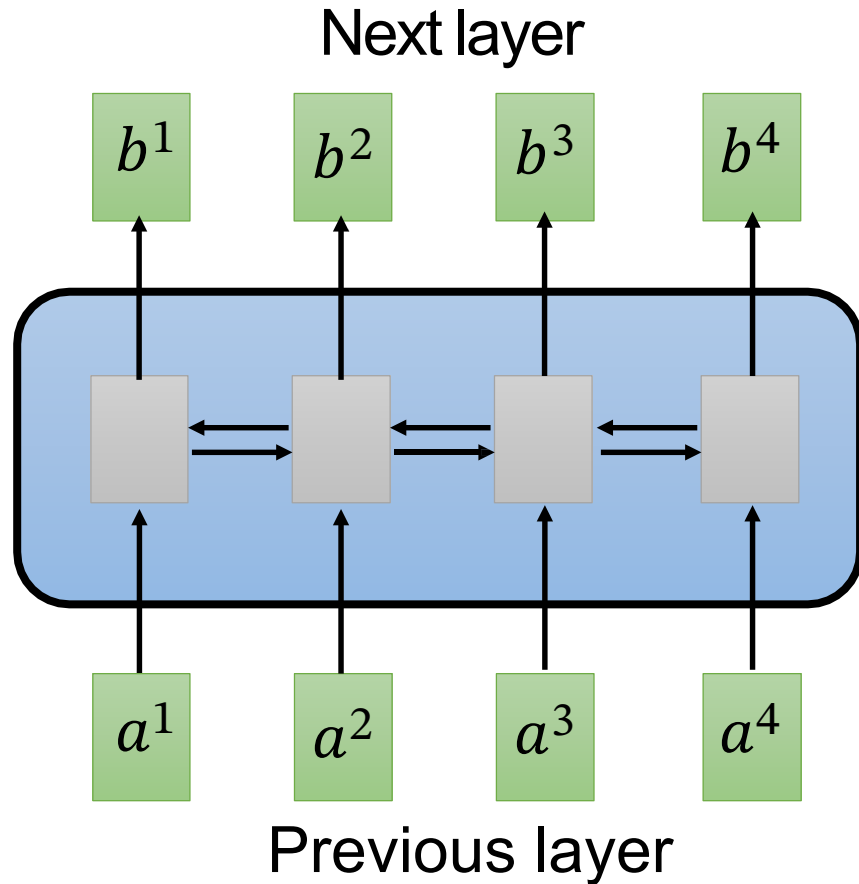


Hard to parallel !



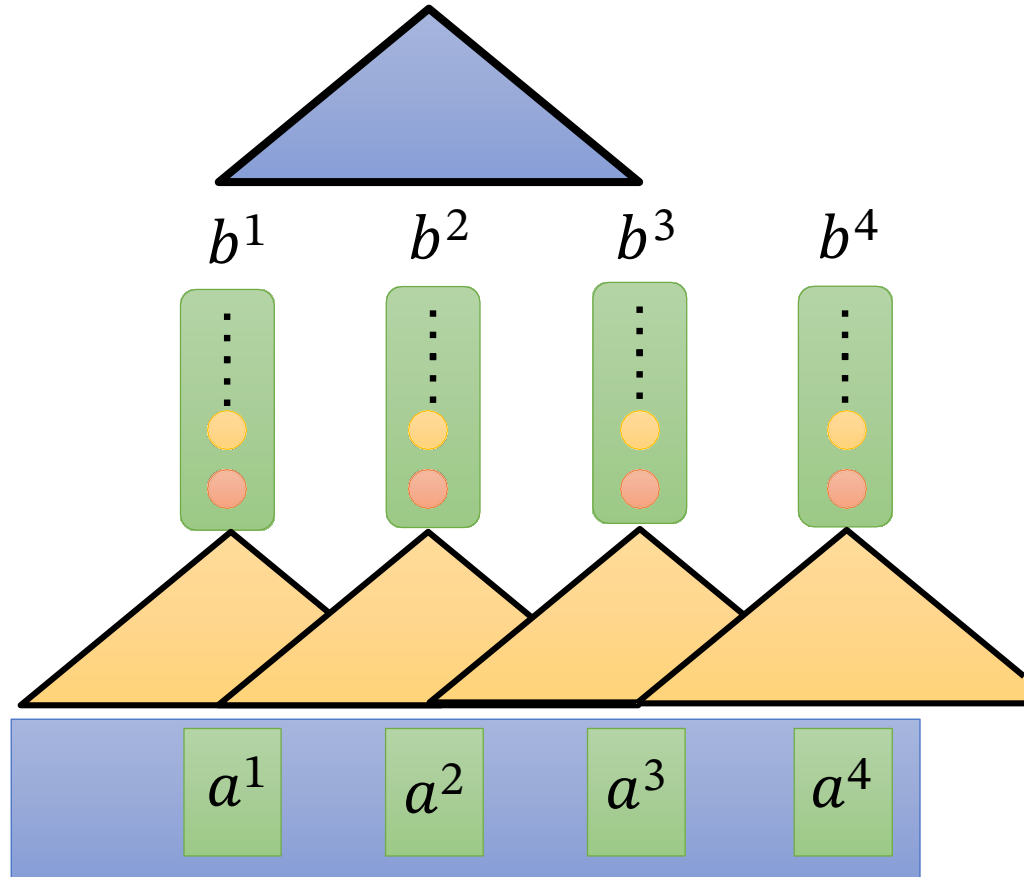
Using CNN to replace RNN

# Sequence



Hard to parallel

Filters in higher layer can consider longer sequence

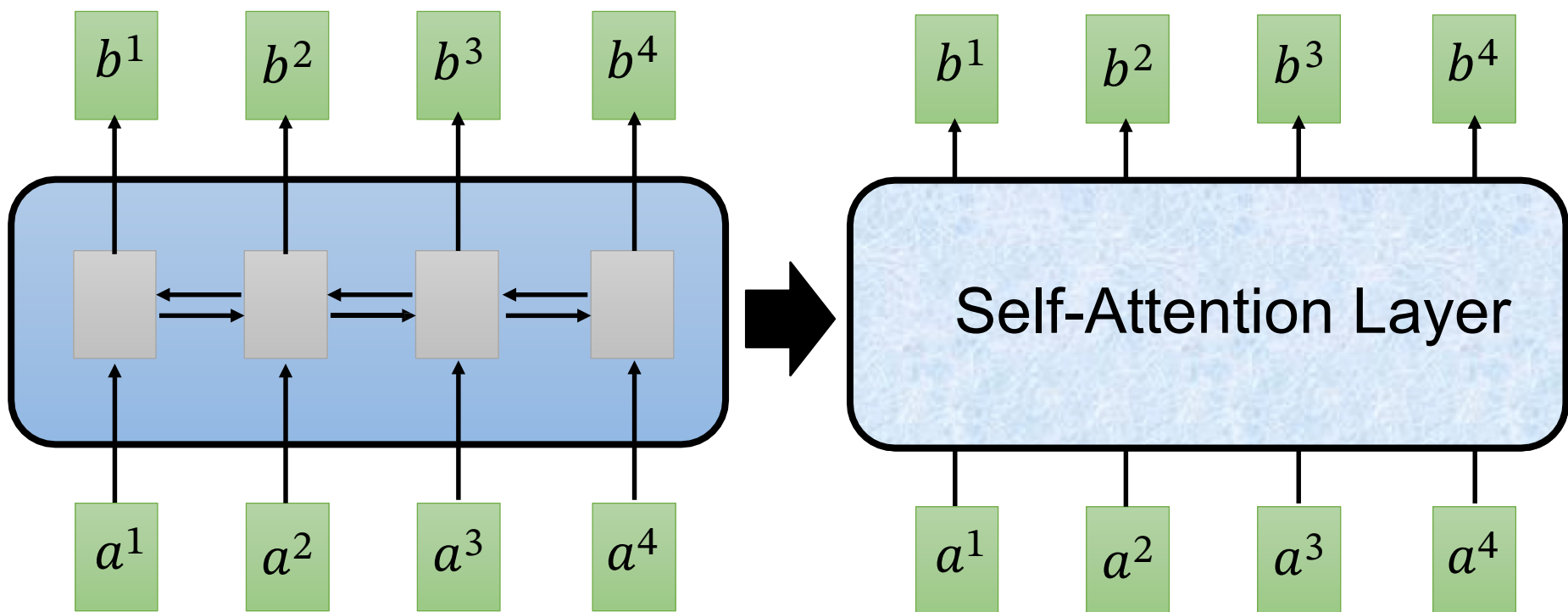


Using CNN to replace RNN  
(CNN can parallel)

# Self-Attention

$b^i$  is obtained based on the whole input sequence.

$b^1, b^2, b^3, b^4$  can be parallelly computed.



You can try to replace any thing that has been done by RNN with self-attention.

# Self-attention

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



$q$ : query (to match others)

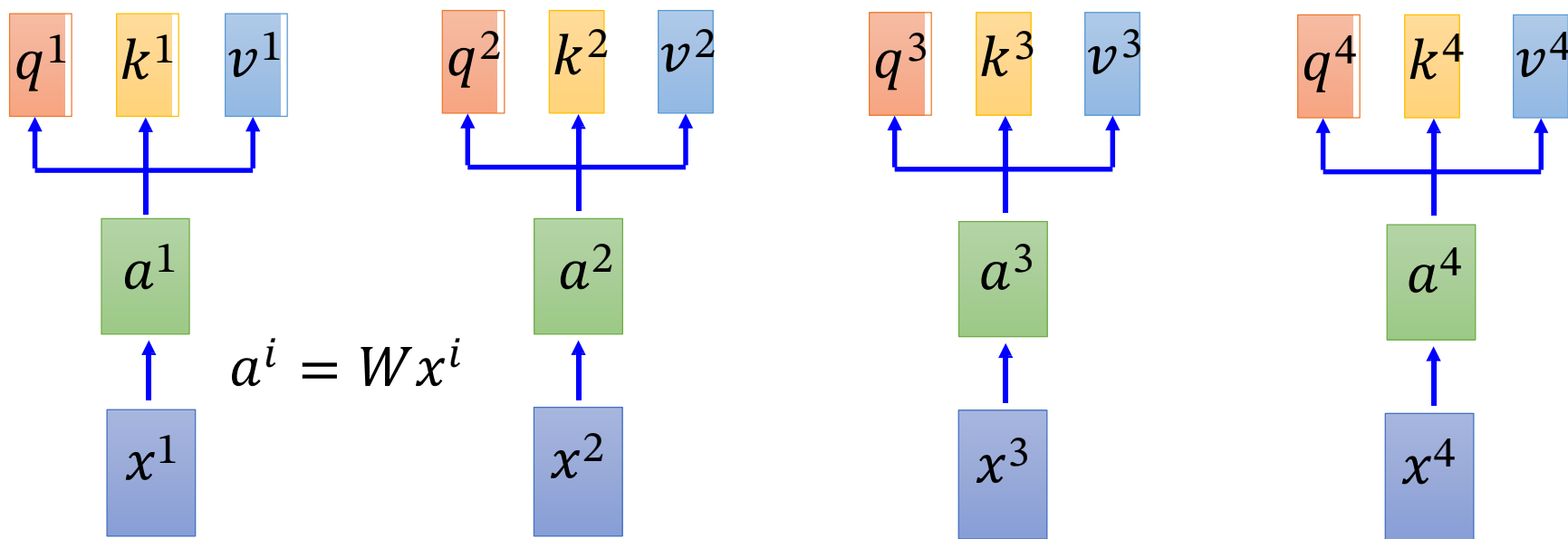
$$q^i = W^q a^i$$

$k$ : key (to be matched)

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

$v$ : information to be extracted

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

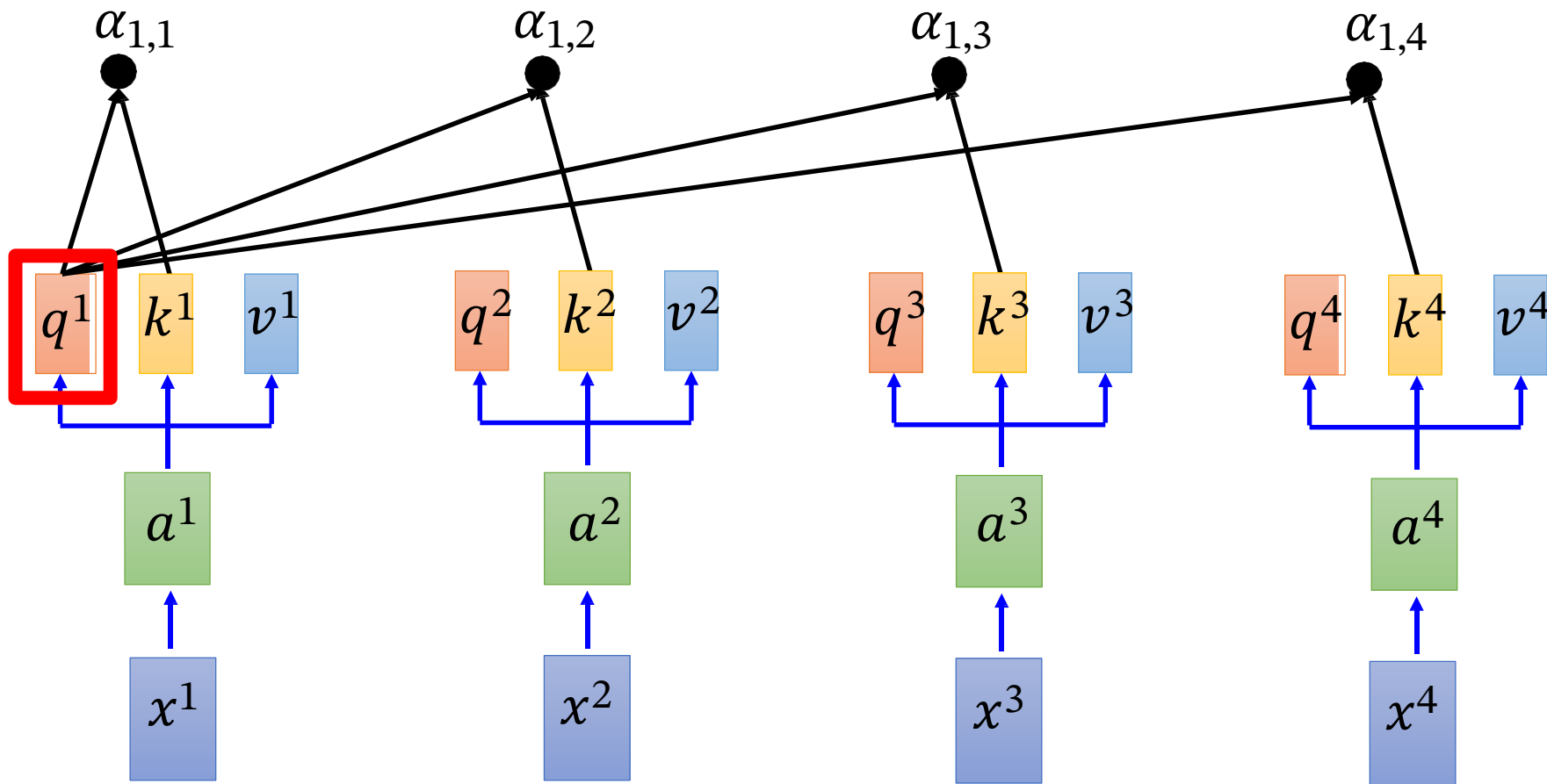


# Self-attention

拿每个 query  $q$  去对每个key  $k$  做 attention

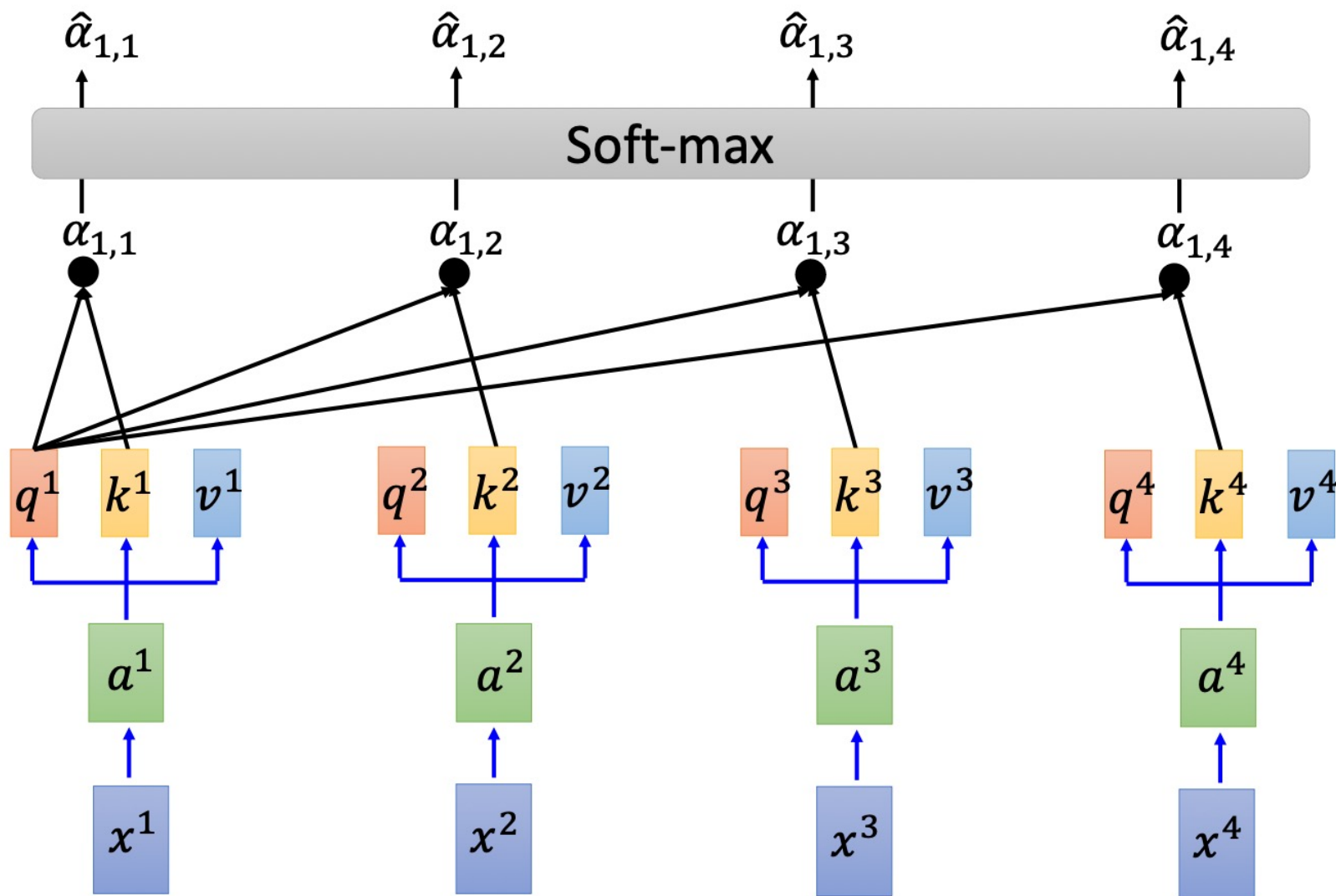
$d$  is the dim of  $q$  and  $k$

Scaled Dot-Product Attention:  $\alpha_{1,i} = \underbrace{q^1 \cdot k^i}_{\text{dot product}} / \sqrt{d}$



# Self-attention

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

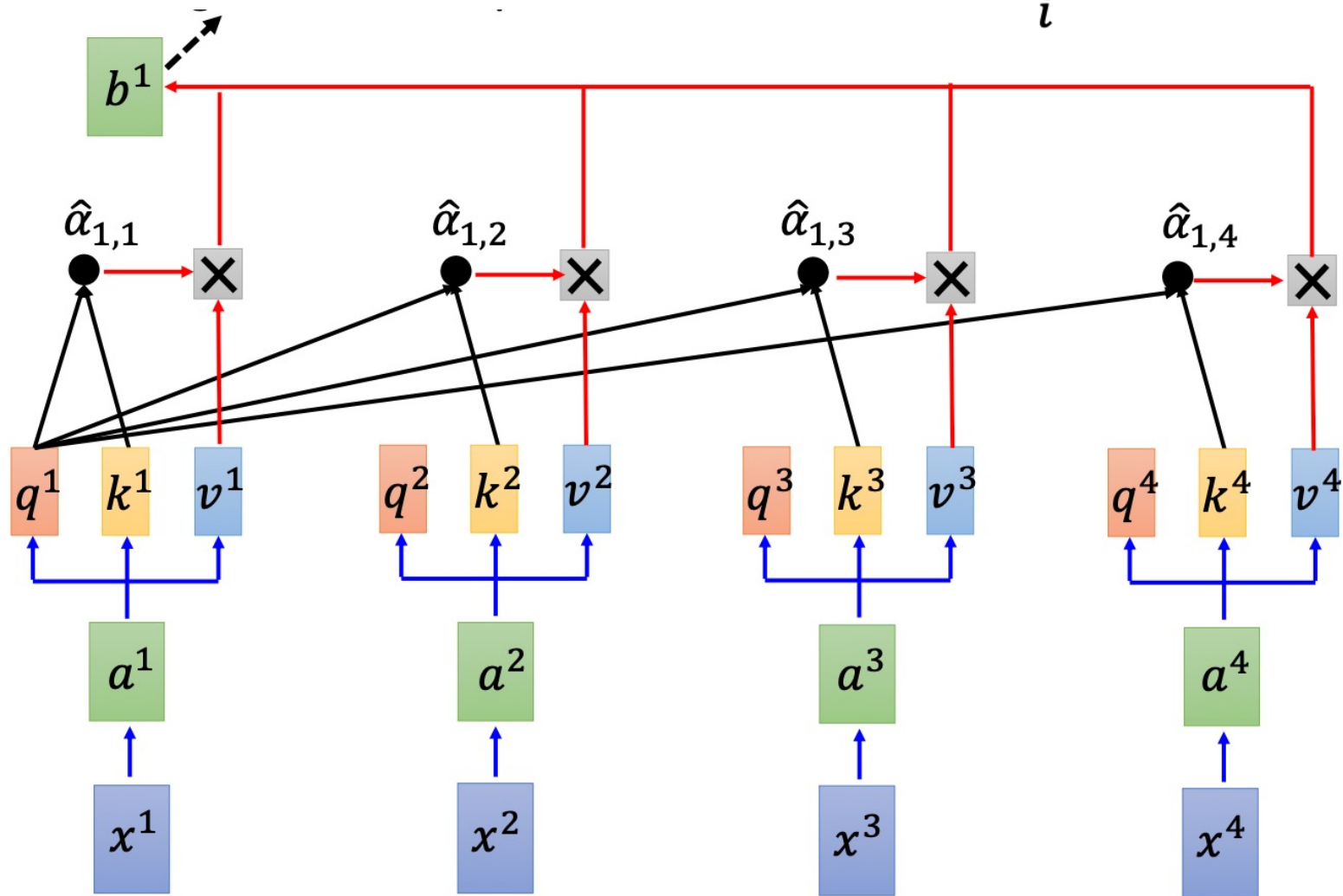




# Self-attention

Considering the whole sequence

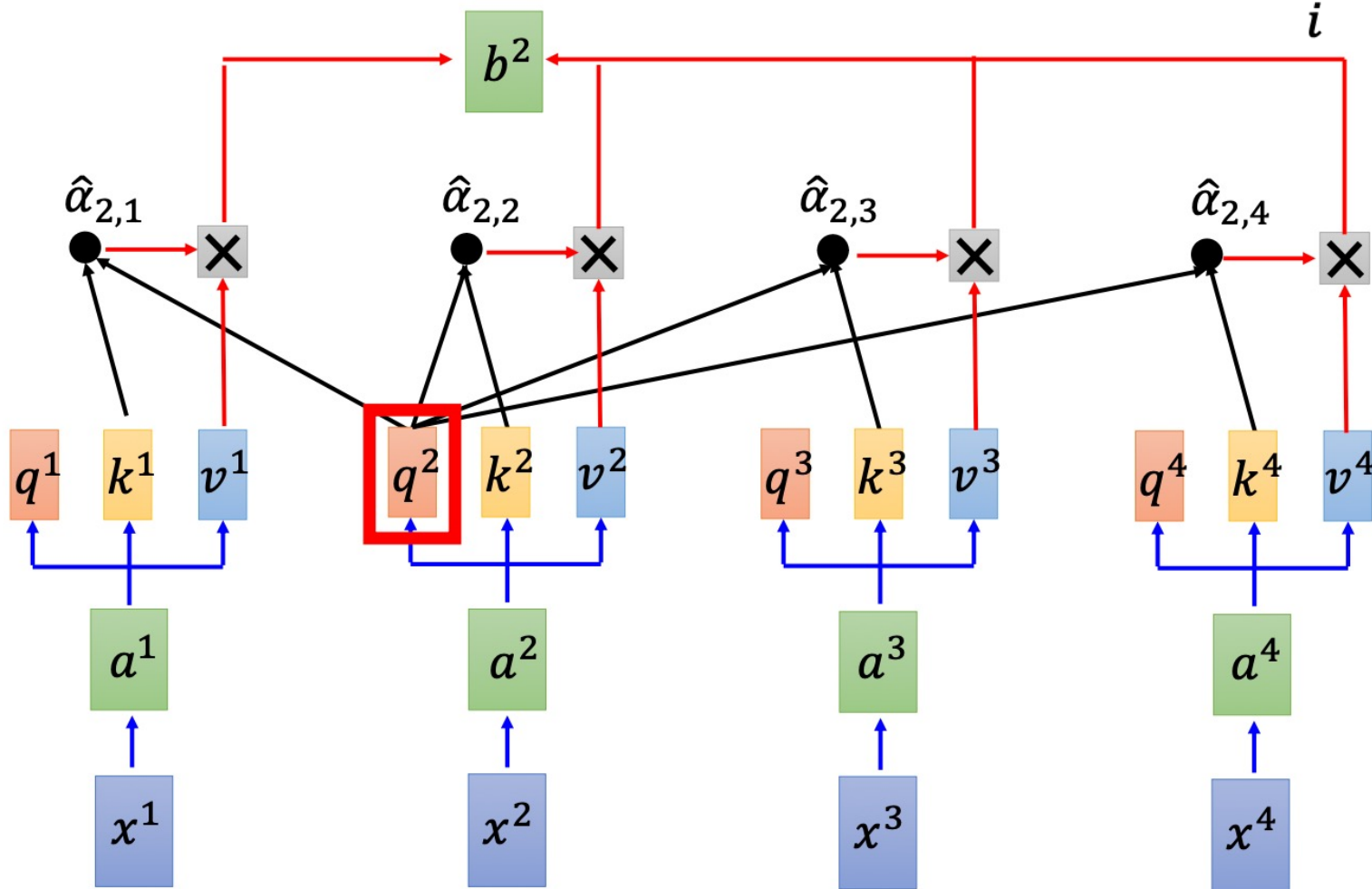
$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$



# Self-attention

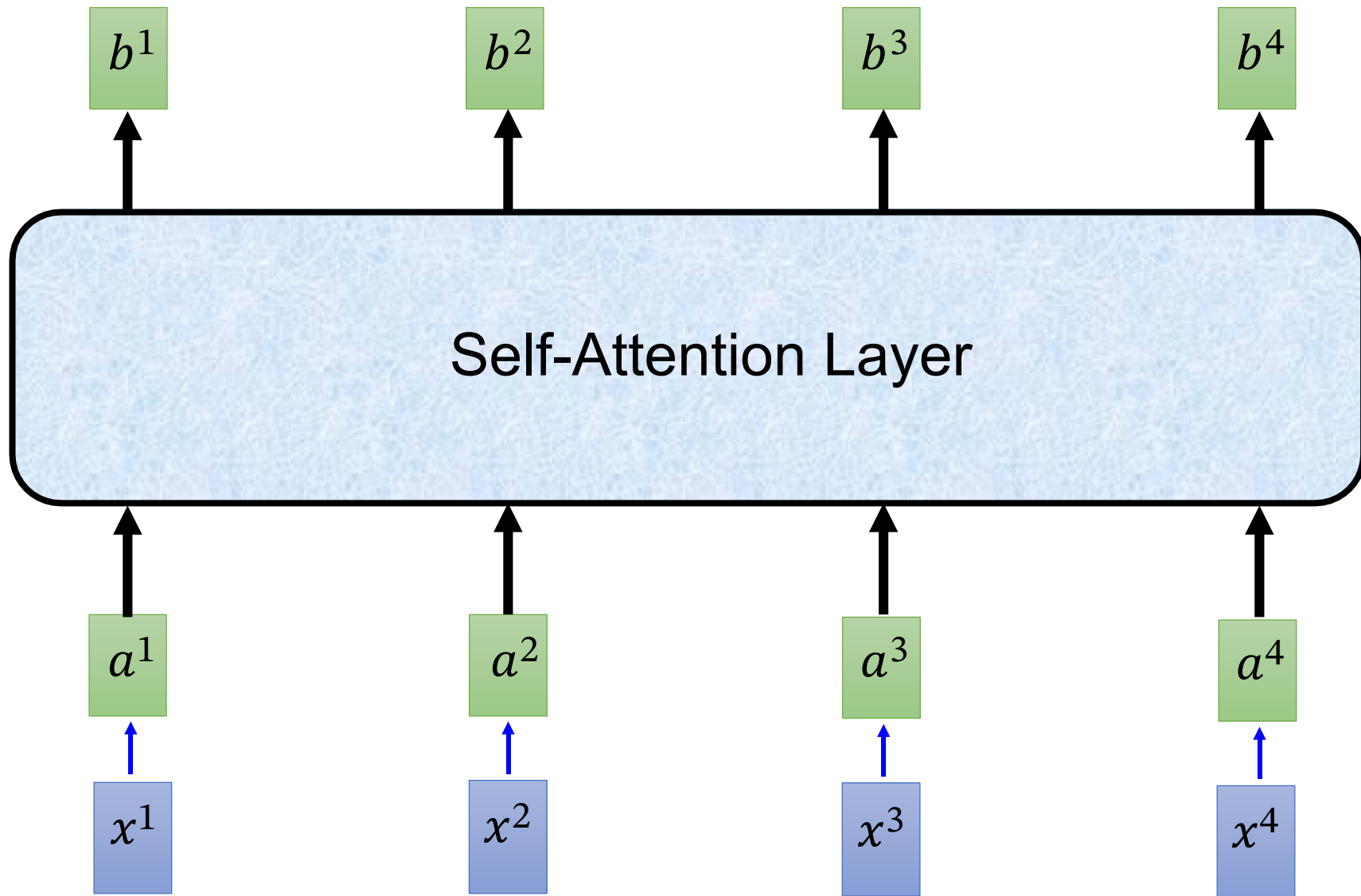
拿每个 query  $q$  去对每个 key  $k$  做 attention

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



# **Self-attention**

$b^1, b^2, b^3, b^4$  can be parallelly computed



# Self-attention

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

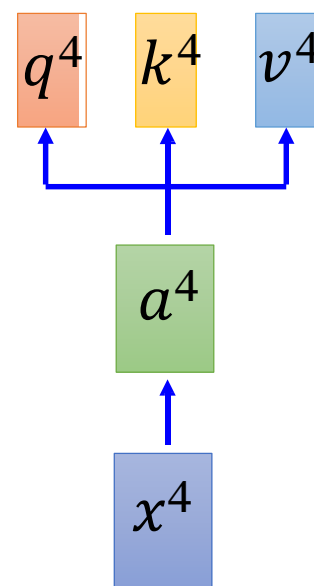
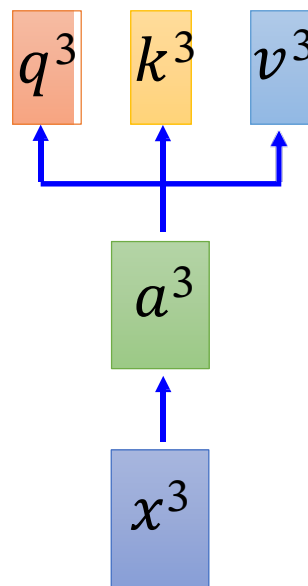
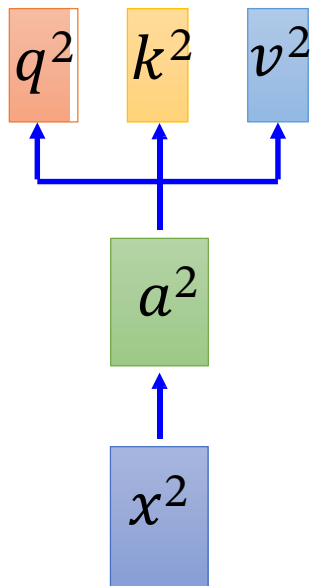
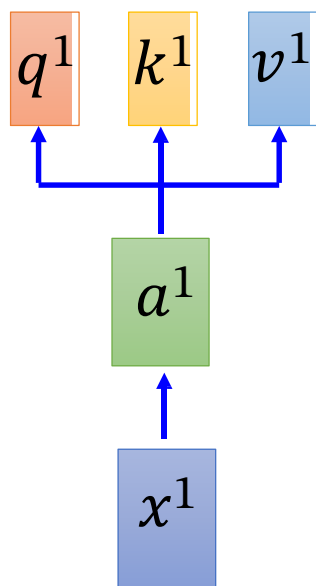
$$\begin{matrix} k^1 & k^2 & k^3 & k^4 \\ K \end{matrix} = \begin{matrix} W^k & \begin{matrix} a^1 & a^2 & a^3 & a^4 \\ I \end{matrix} \end{matrix}$$

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

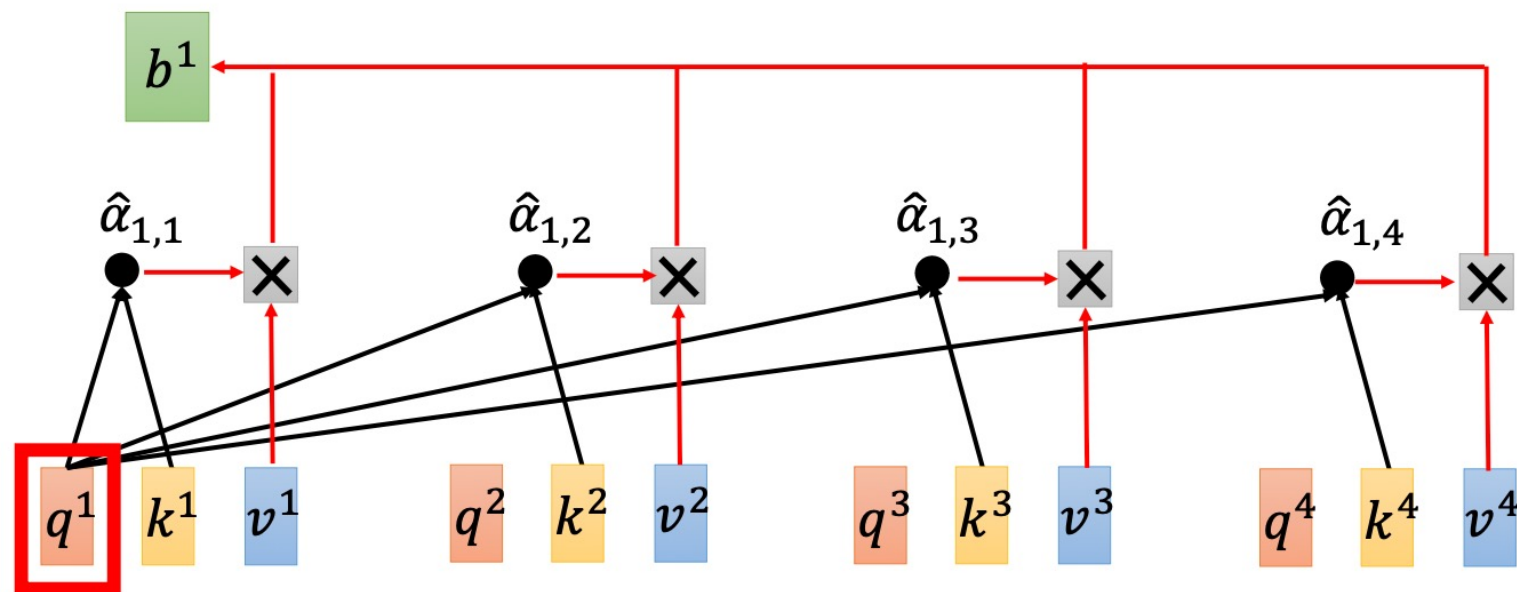
$$q^i = W^q a^i$$

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

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



# Self-attention



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

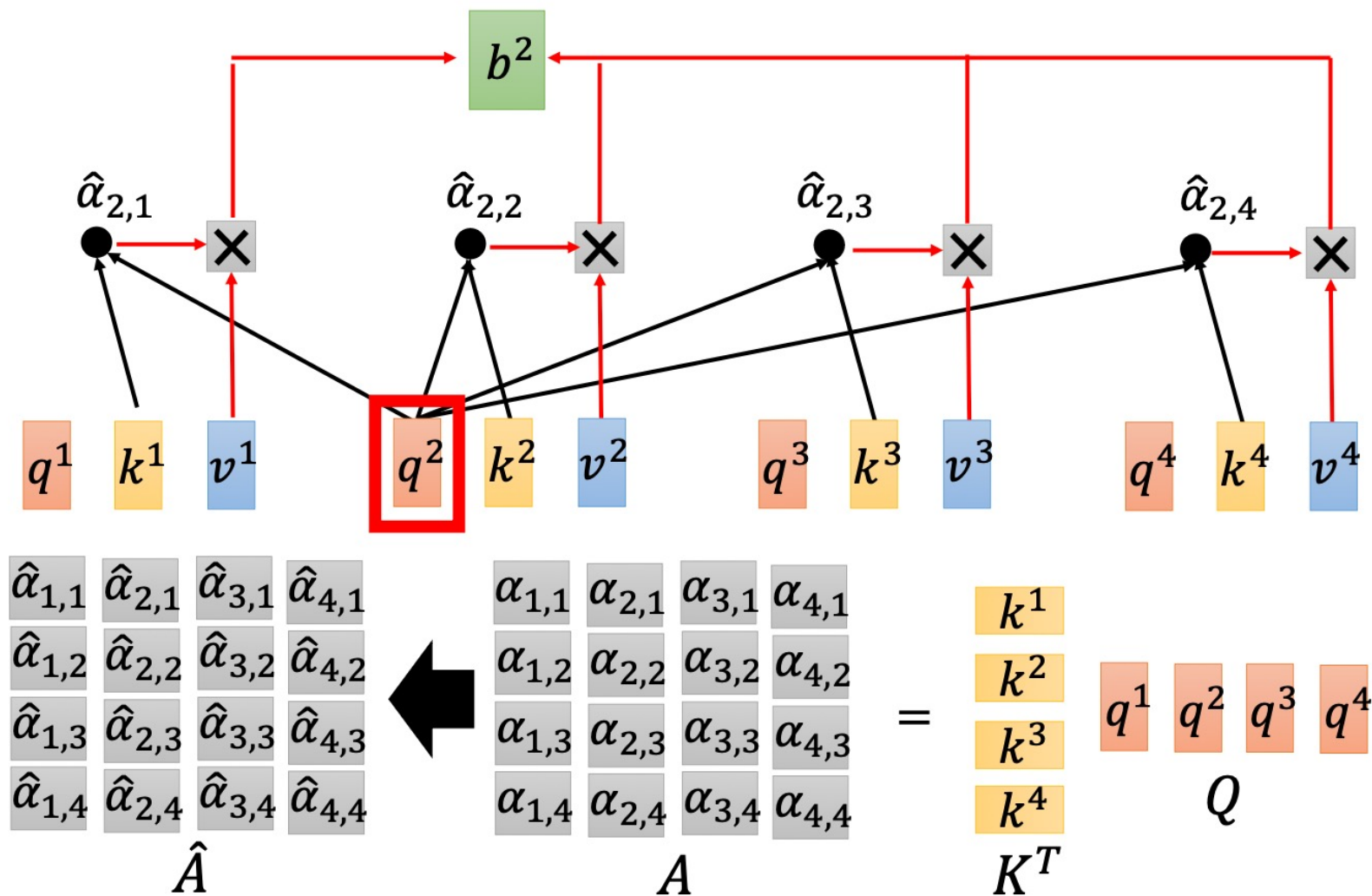
$$\alpha_{1,3} = k^3 q^1 \quad \alpha_{1,4} = k^4 q^1$$

(ignore  $\sqrt{d}$  for simplicity)

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

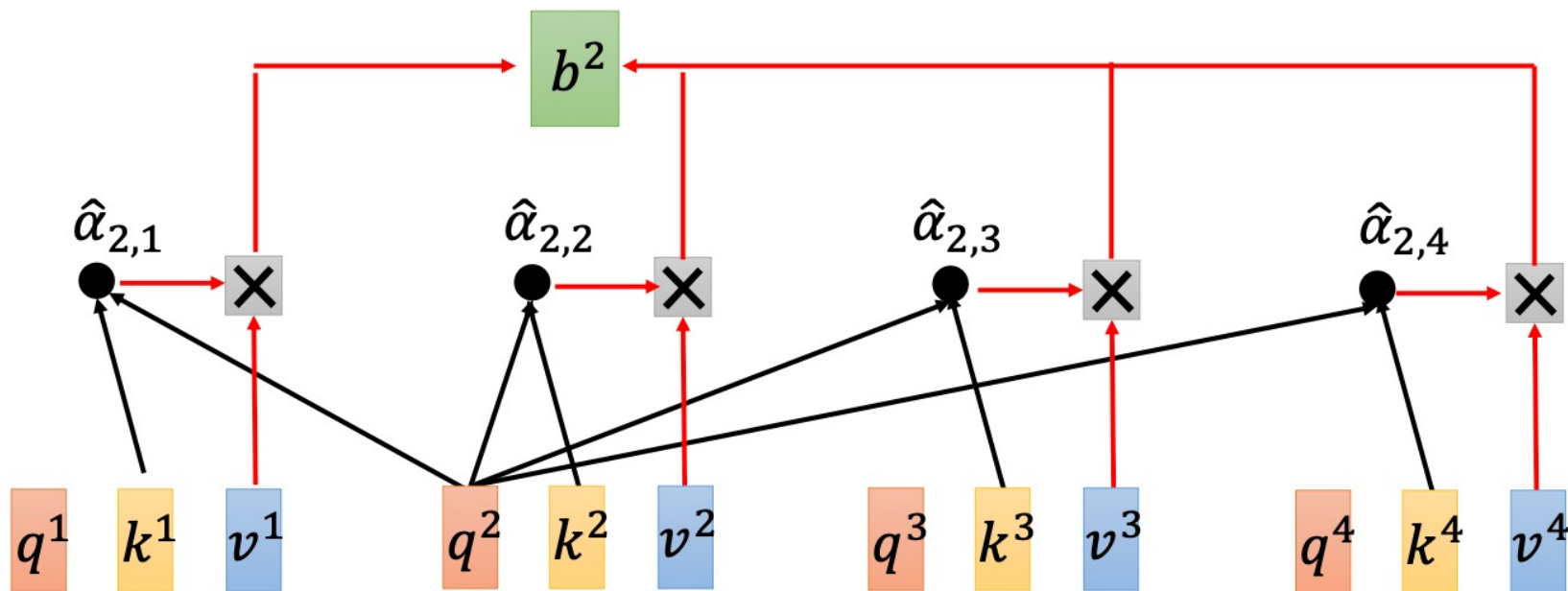
# Self-attention

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



# Self-attention

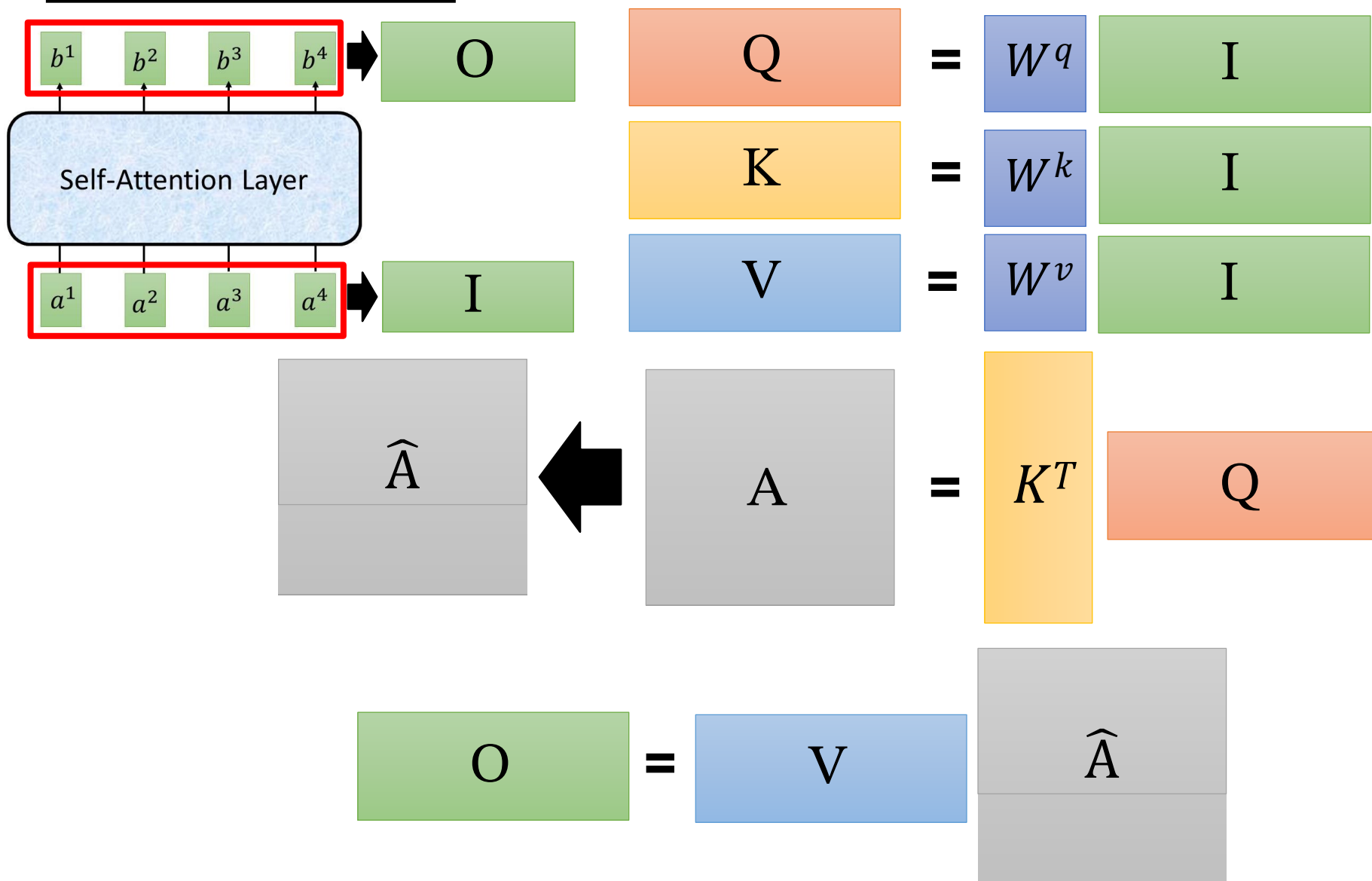
$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$



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

$\hat{A}$

# Self-attention



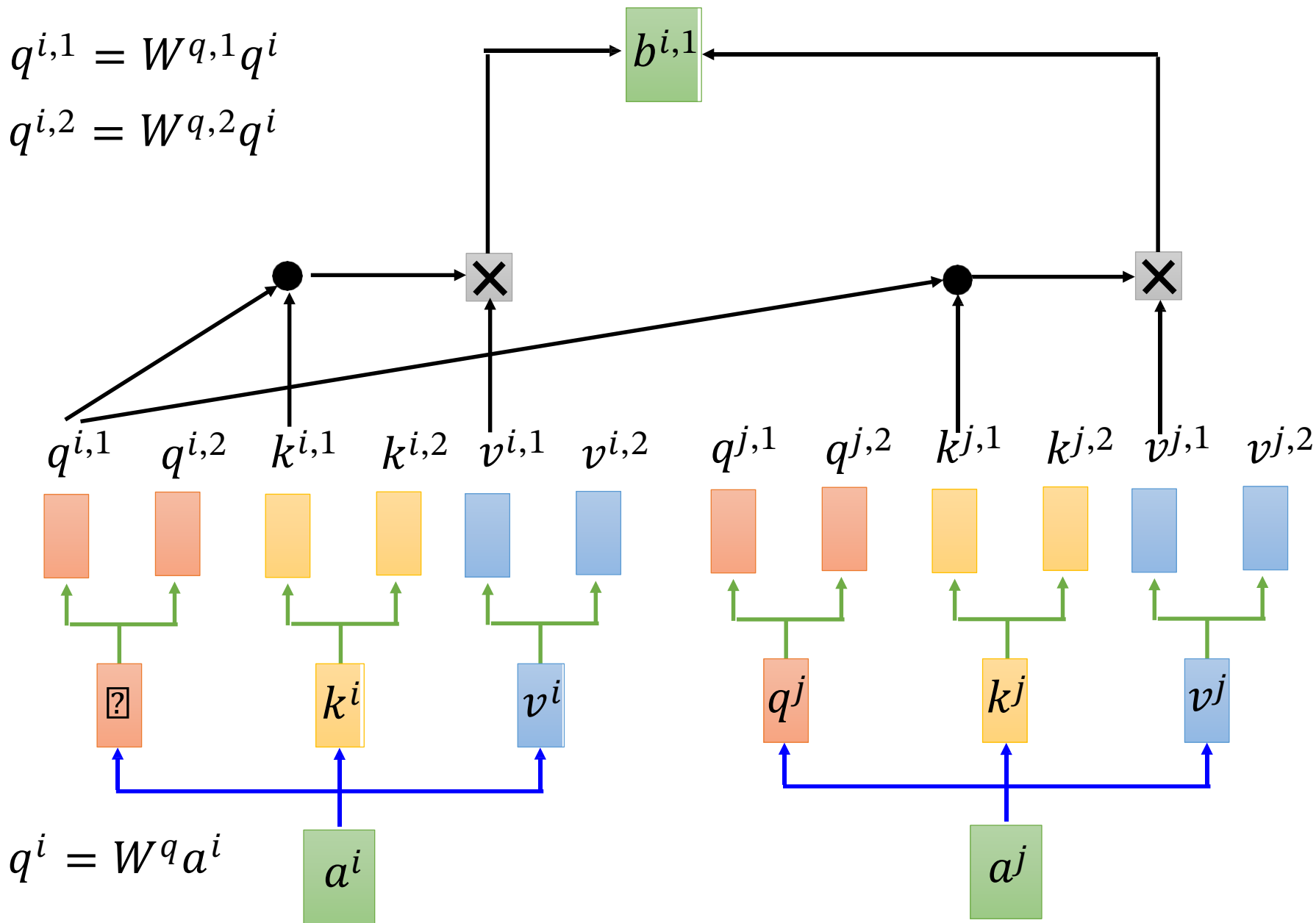
反正就是一堆矩阵乘法，用 GPU 可以加速



# Multi-head Self-attention (2 heads as example)

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

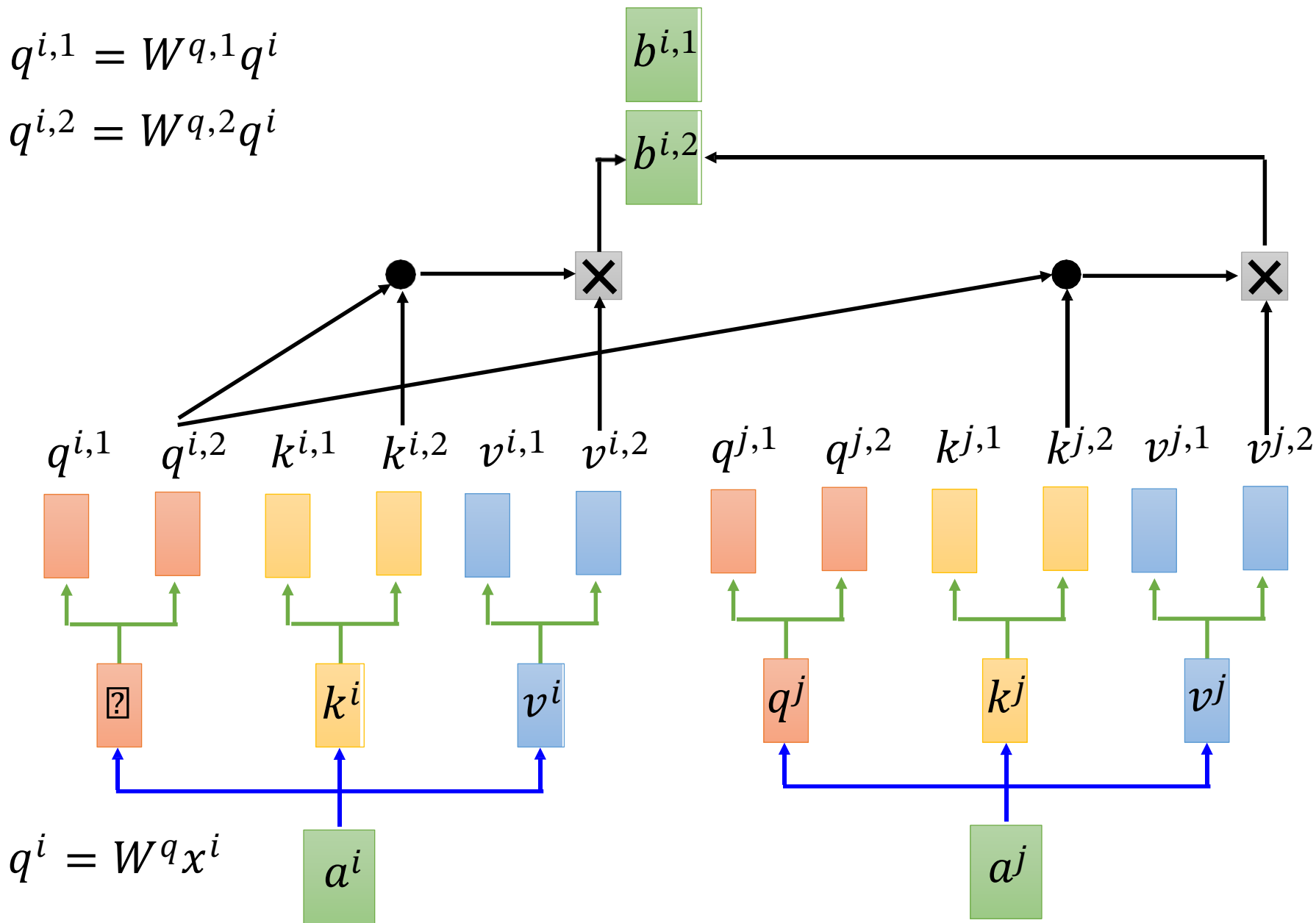
$$q^{i,2} = W^{q,2} q^i$$



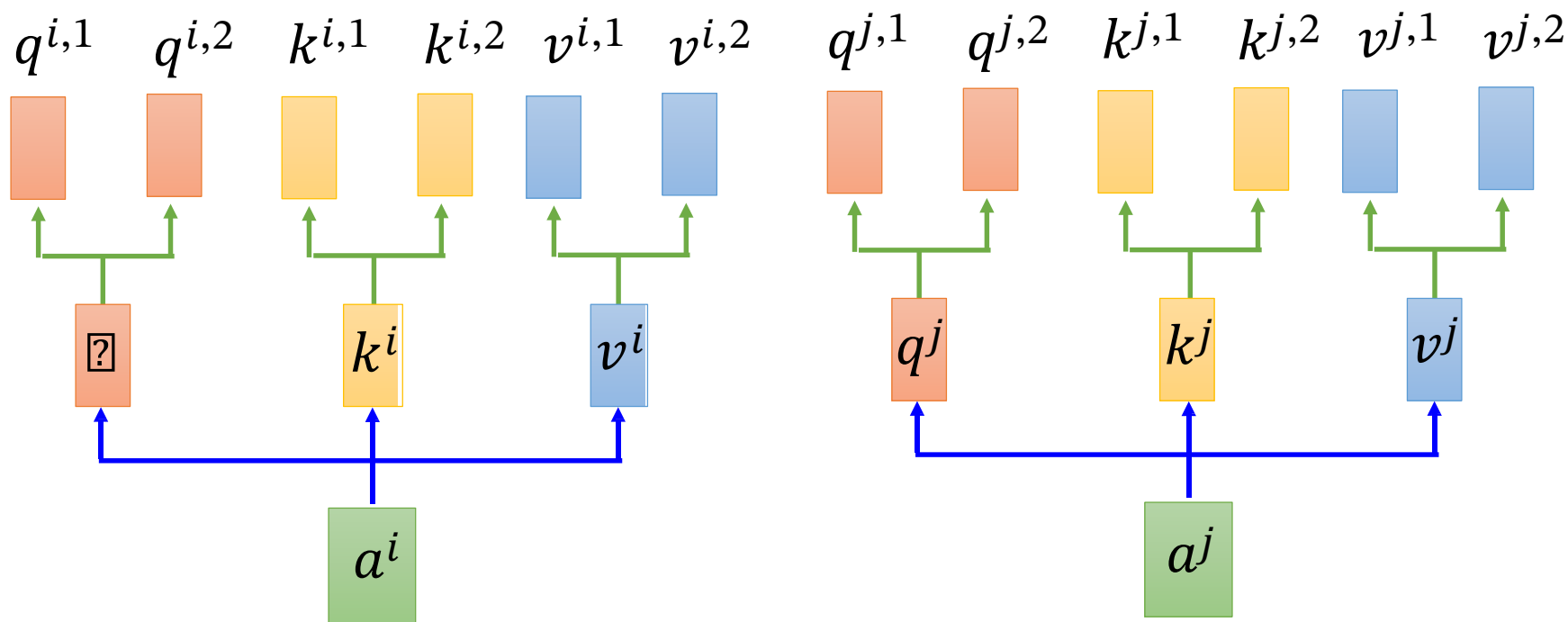
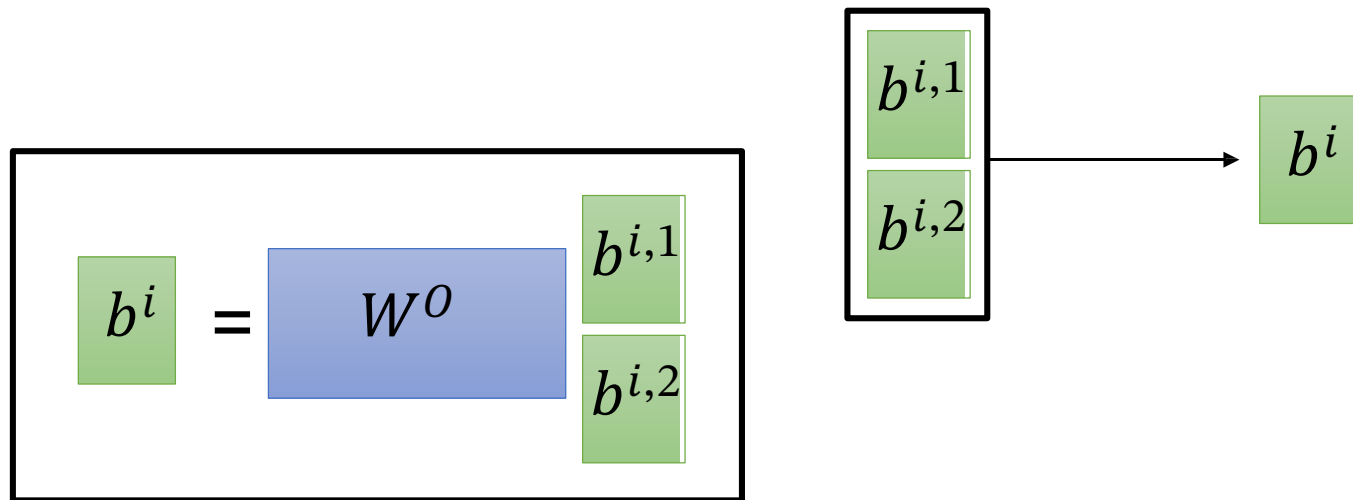
# Multi-head Self-attention (2 heads as example)

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

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

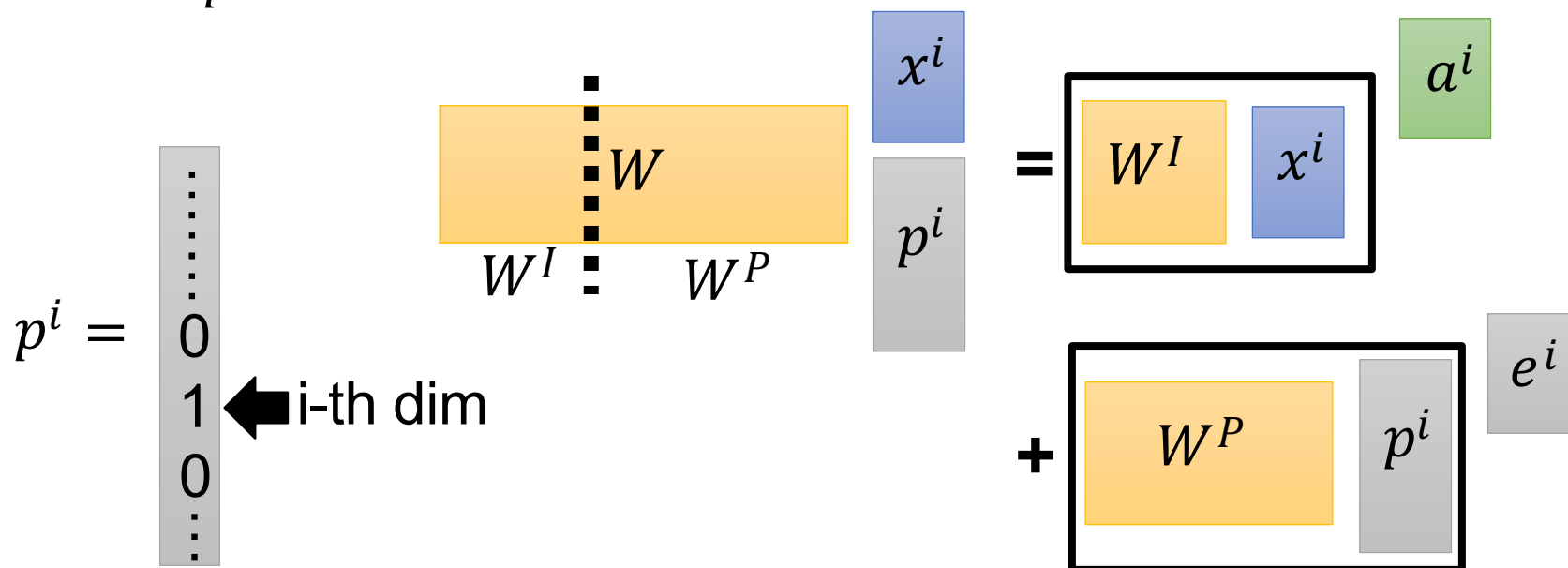
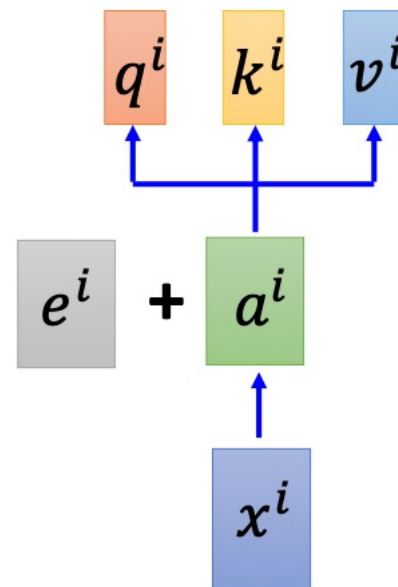


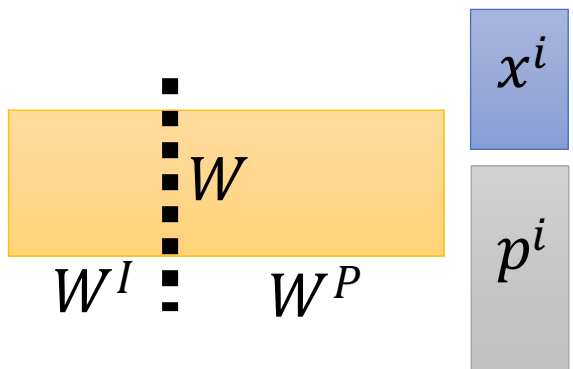
# Multi-head Self-attention (2 heads as example)



# Positional Encoding

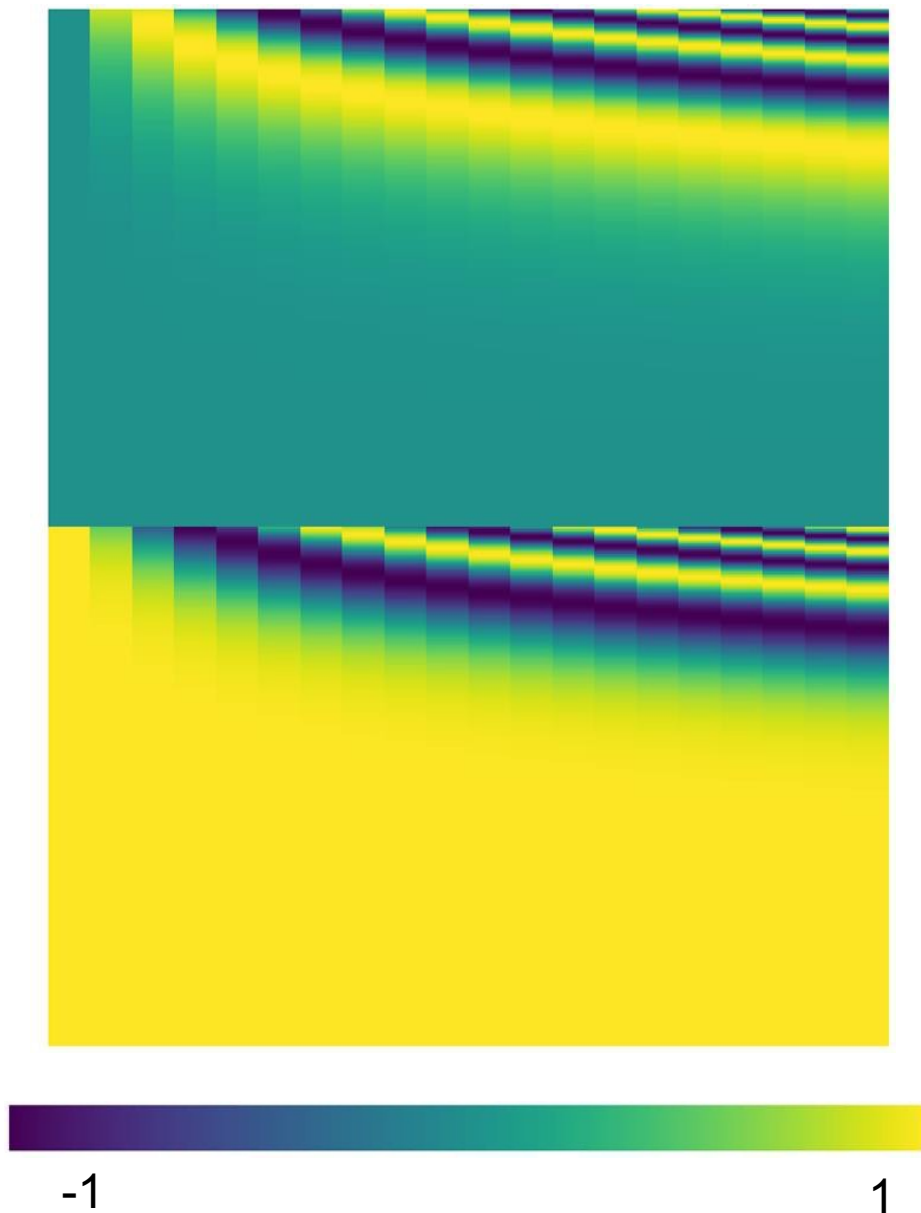
- No position information in self-attention.
- Original paper: each position has a unique positional vector  $e^i$  (not learned from data)
- In other words: each  $x^i$  appends a one-hot vector  $p^i$





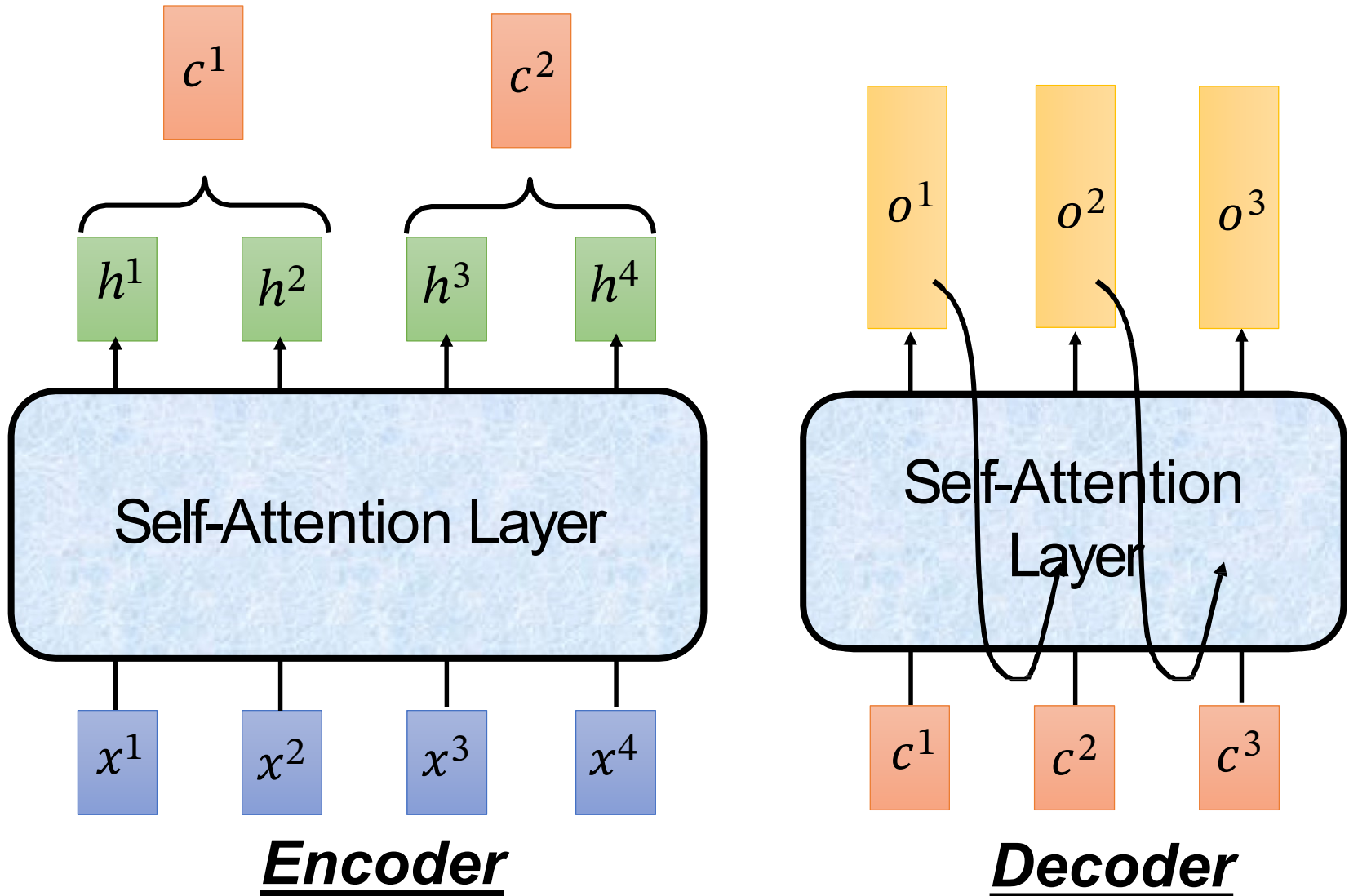
$$= \begin{bmatrix} W^I & x^i \end{bmatrix} a^i + \begin{bmatrix} W^P & p^i \end{bmatrix} e^i$$

Diagram illustrating the layer structure. The input is split into two parts,  $W^I$  and  $W^P$ , separated by a dashed line. The input is represented by a blue box labeled  $x^i$  and a gray box labeled  $p^i$ . The output is represented by a green box labeled  $a^i$  and a gray box labeled  $e^i$ .



source of image: <http://jalammar.github.io/illustrated-transformer/>

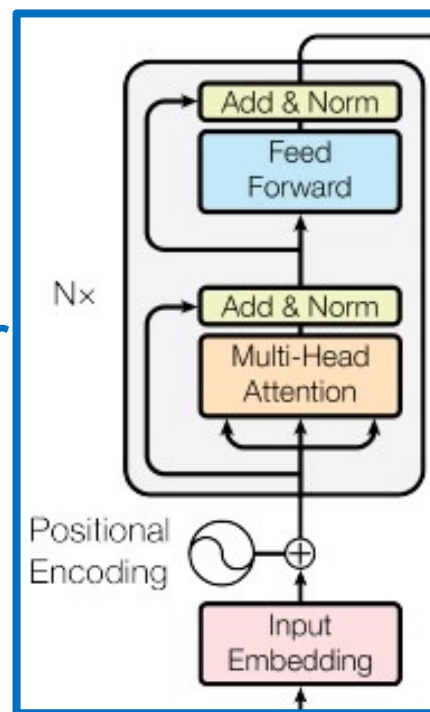
# Seq2seq with Attention



# Transformer

Using Chinese to English translation as example

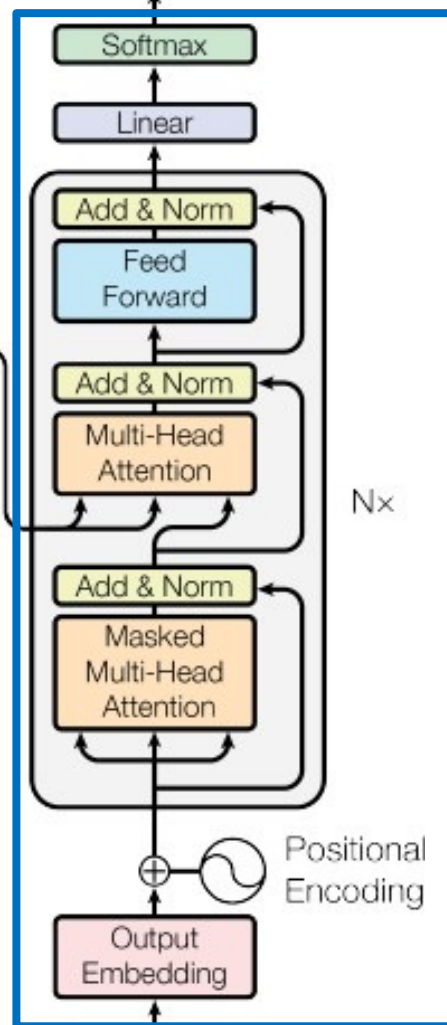
Encoder



Inputs

机器学习

machine learning



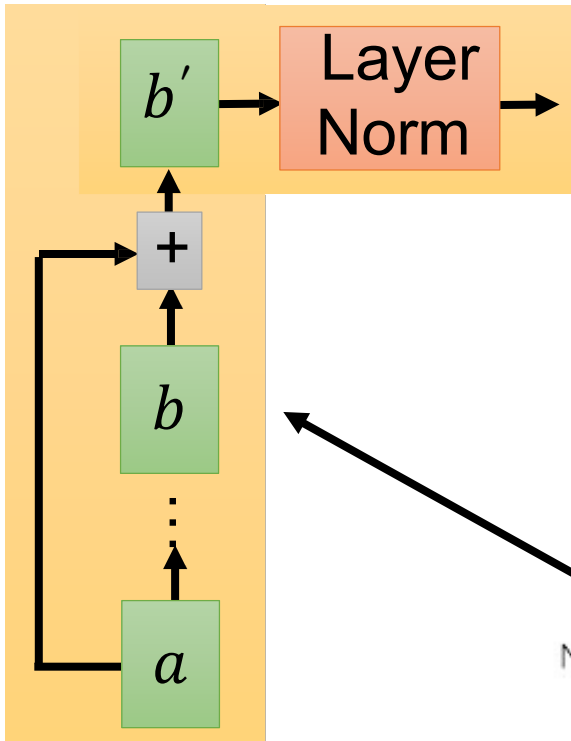
Outputs  
(shifted right)

<BOS>

machine

Decoder

# Transformer

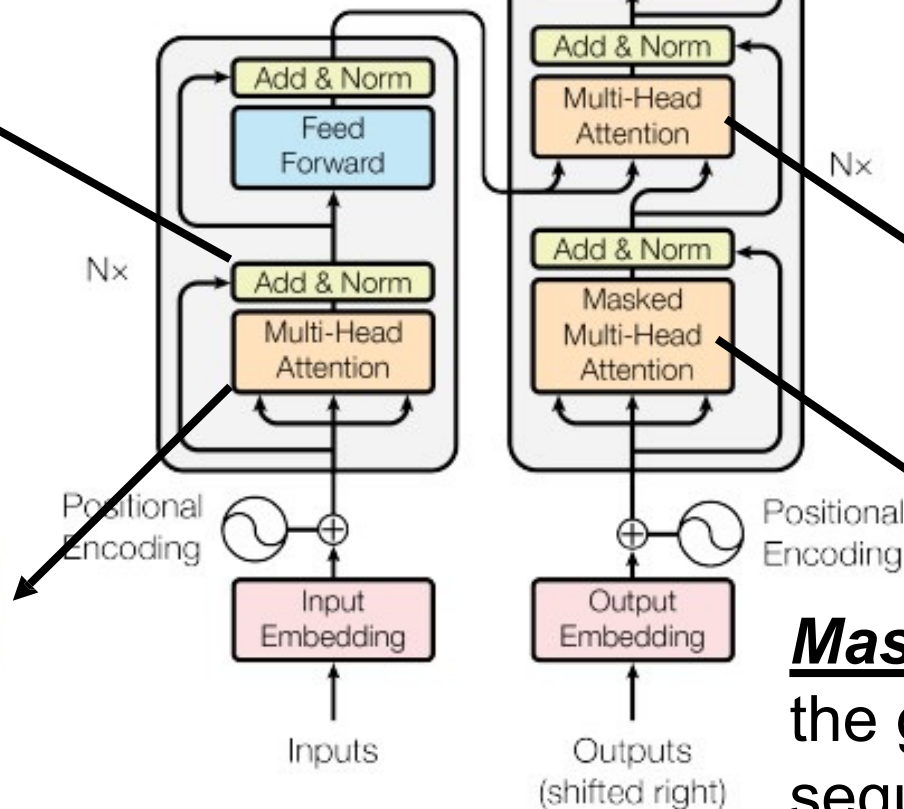
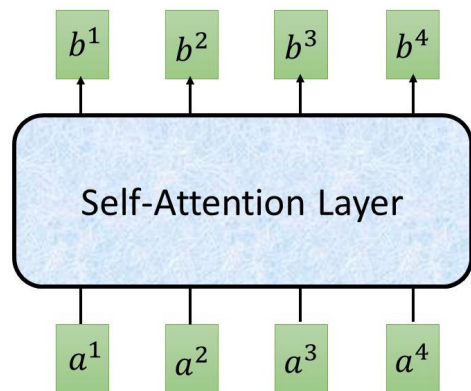


Layer Norm:

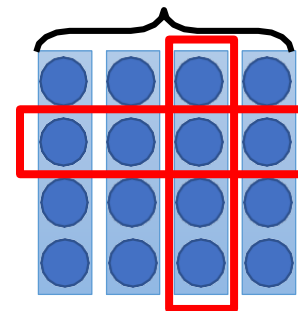
<https://arxiv.org/abs/1607.06450>

Batch Norm:

<https://www.youtube.com/watch?v=BZh1ltr5Rkg>



Batch Size



$\mu = 0,$   
 $\sigma = 1$   
Batch

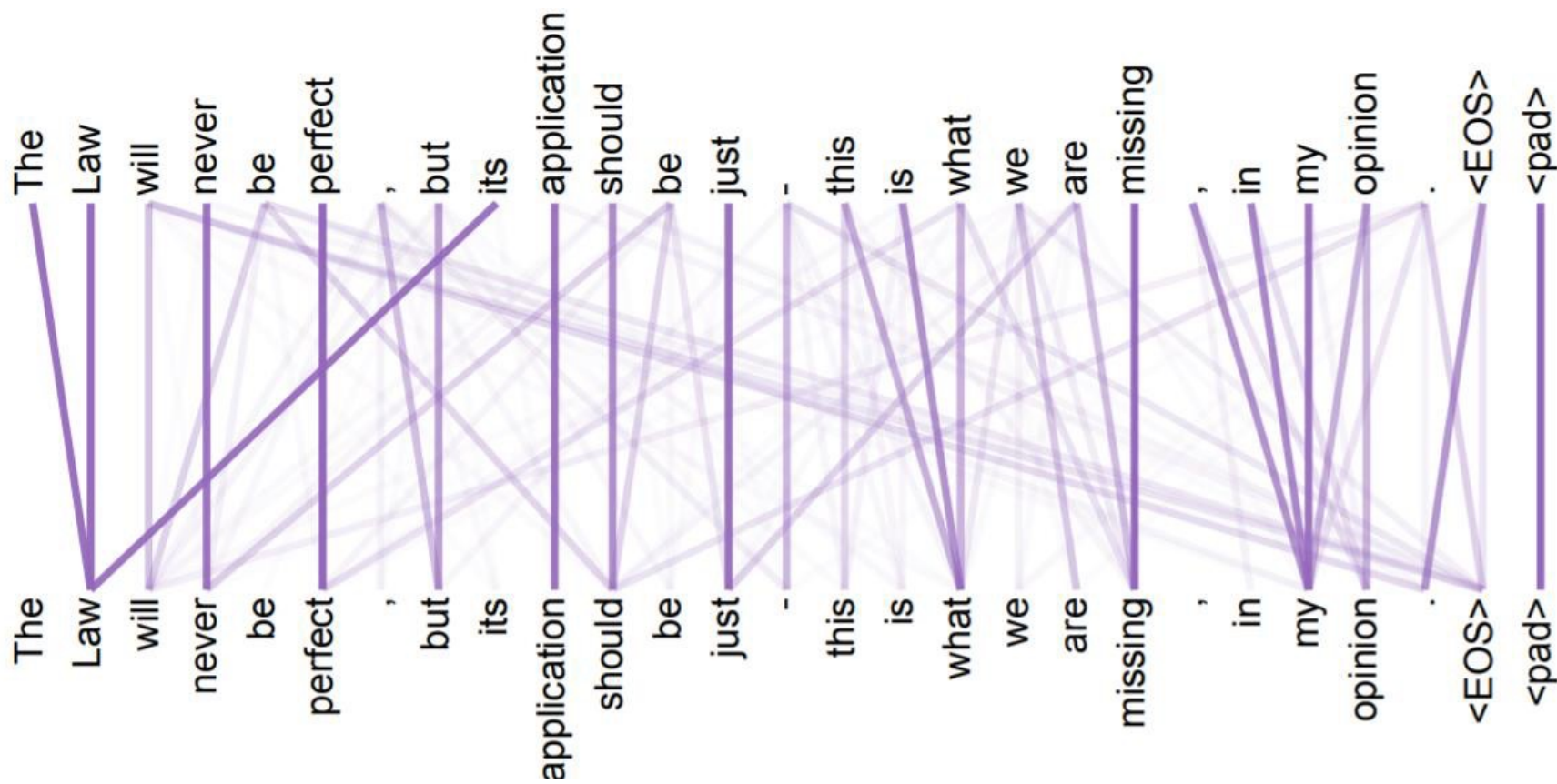
$\mu = 0, \sigma = 1$   
Layer

attend on the  
input sequence

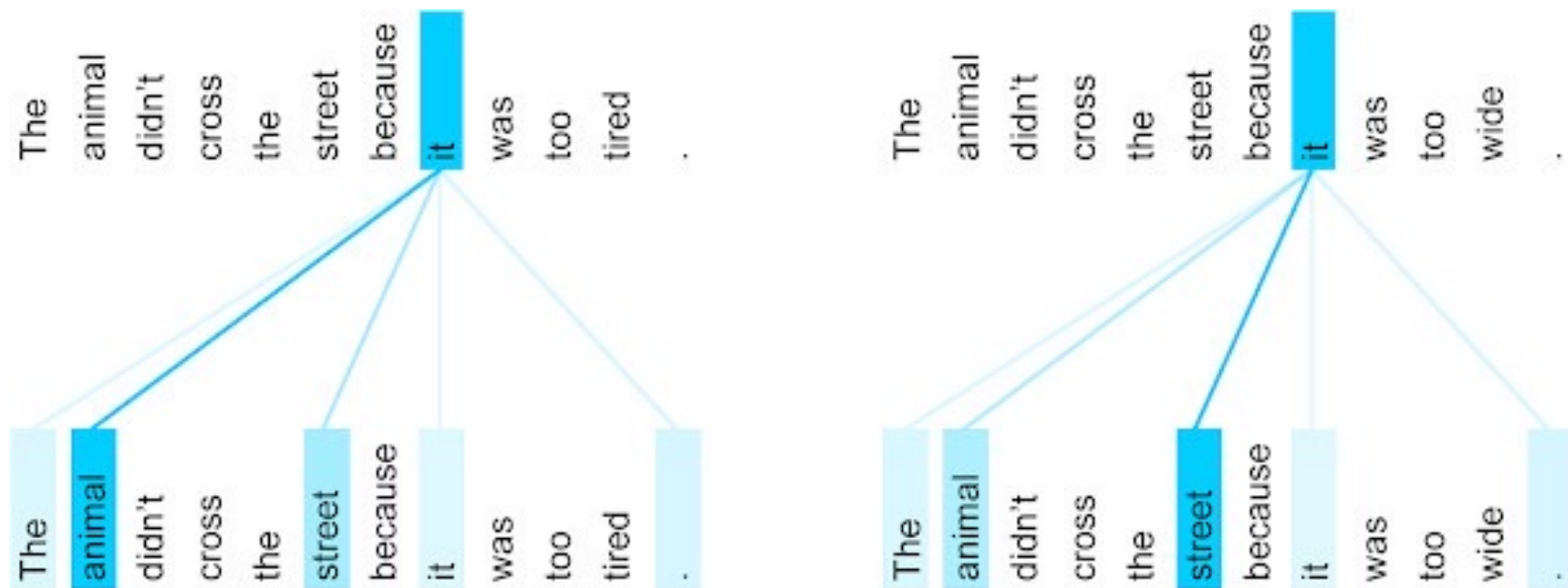
**Masked**: attend on  
the generated  
sequence



# Attention Visualization



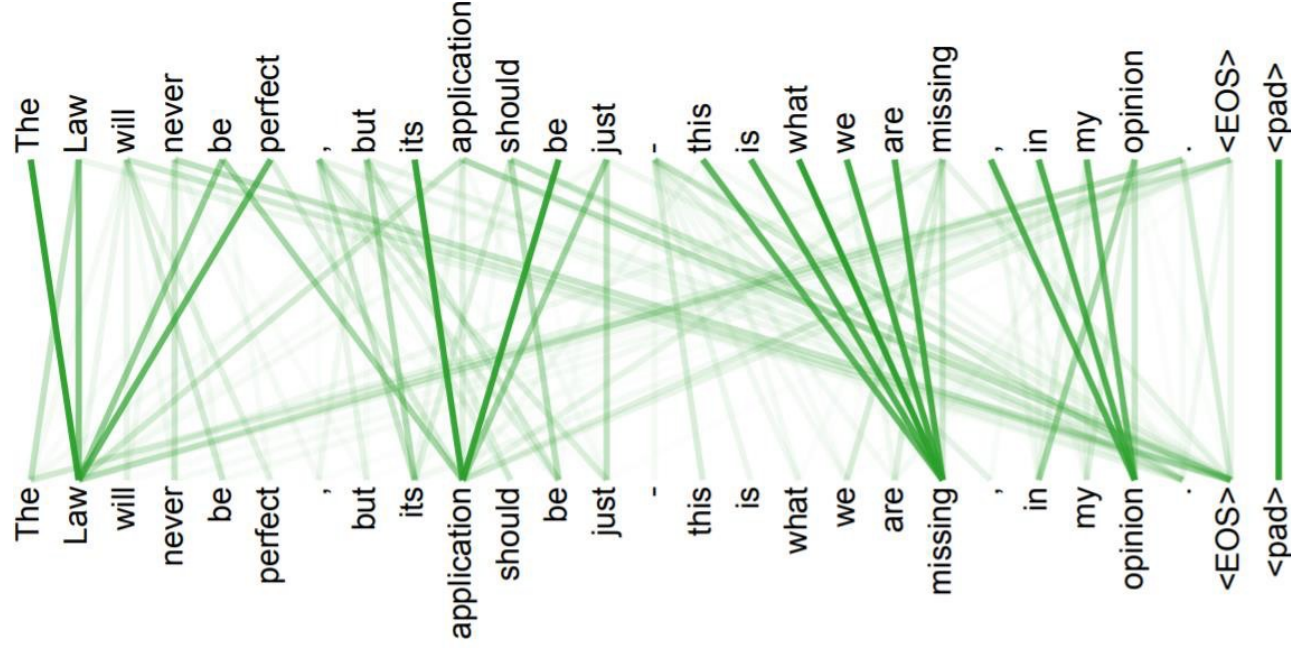
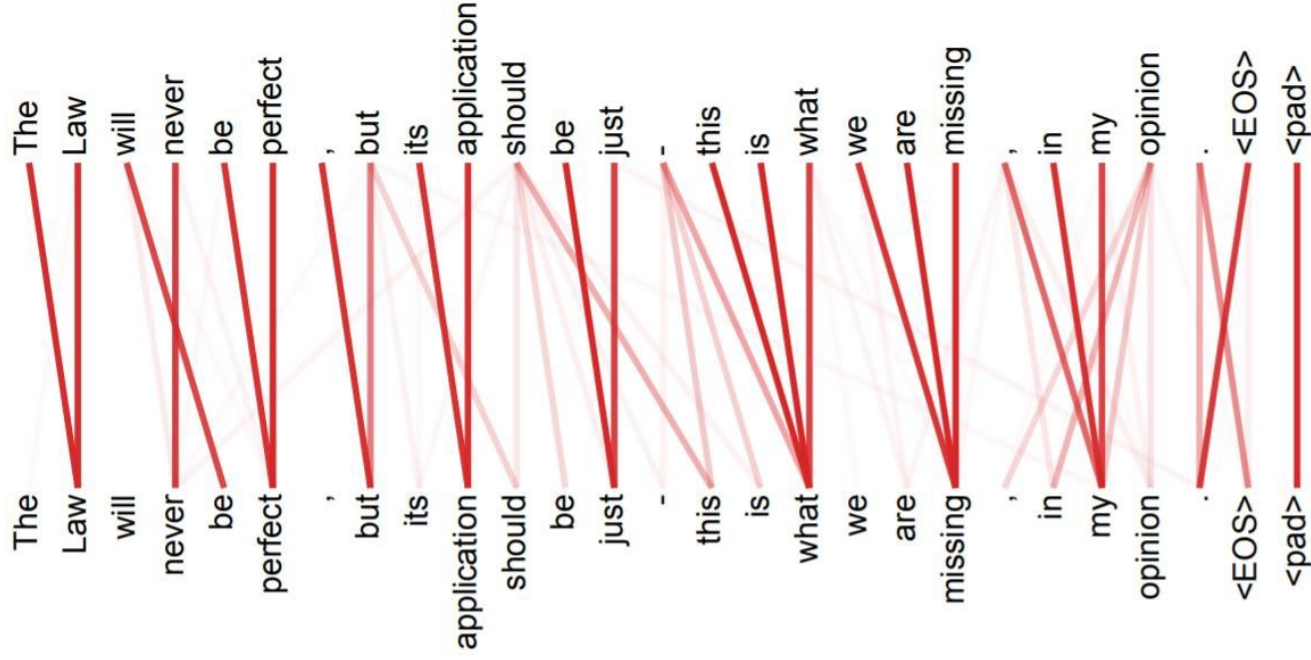
# Attention Visualization



The encoder self-attention distribution for the word "it" from the 5th to the 6th layer of a Transformer trained on English to French translation (one of eight attention heads).

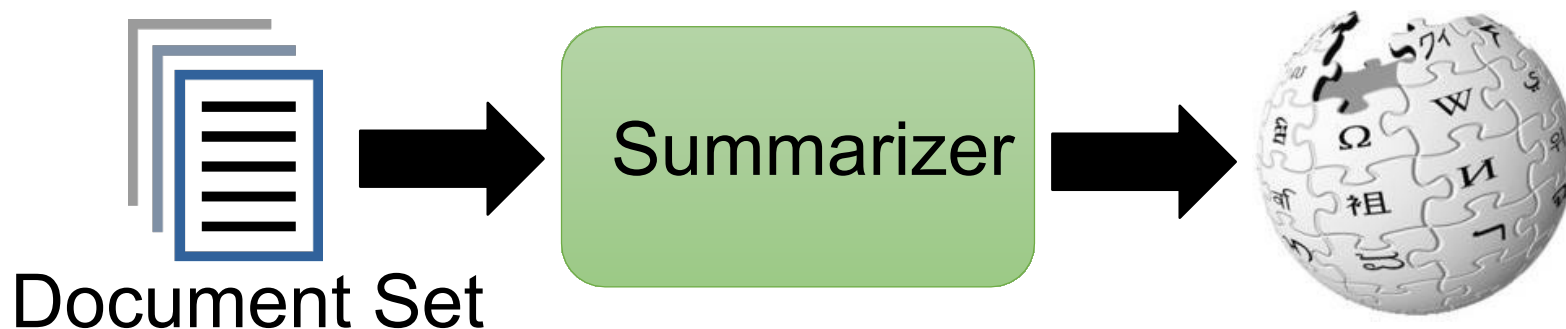
<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

# Multi-head Attention



# Example Application

- If you can use seq2seq, you can use transformer.



Dataset	Input	Output	# examples
Gigaword (Graff & Cieri, 2003)	$10^1$	$10^1$	$10^6$
CNN/DailyMail (Nallapati et al., 2016)	$10^2$ – $10^3$	$10^1$	$10^5$
WikiSum (ours)	$10^2$ – $10^6$	$10^1$ – $10^3$	$10^6$

# **Acknowledgement**

Reference and thanks to:

- **National Taiwan University ML2020 Course:**  
Machine Learning

<https://speech.ee.ntu.edu.tw/~hylee/ml/2020-spring.php>