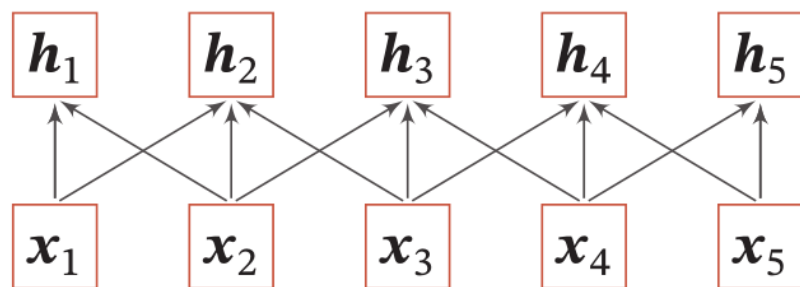


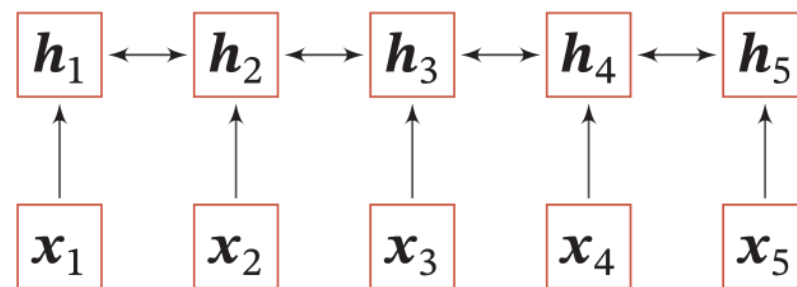
## 自注意力模型

# 自注意力模型

- 当使用神经网络来处理一个变长的向量序列时，我们通常可以使用卷积网络或循环网络进行编码来得到一个相同长度的输出向量序列。



(a) 卷积网络



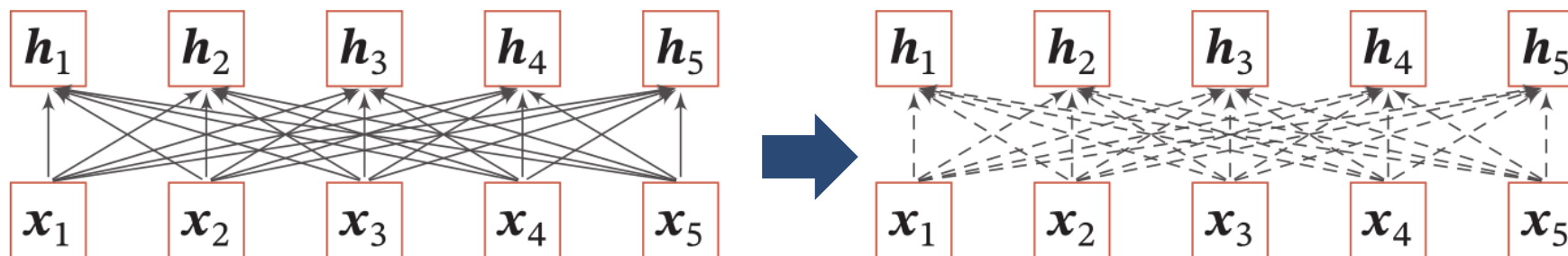
(b) 双向循环网络

只建模了输入信息的局部依赖关系

# 自注意力模型

► 如何建立非局部 (Non-local) 的依赖关系

► 全连接?



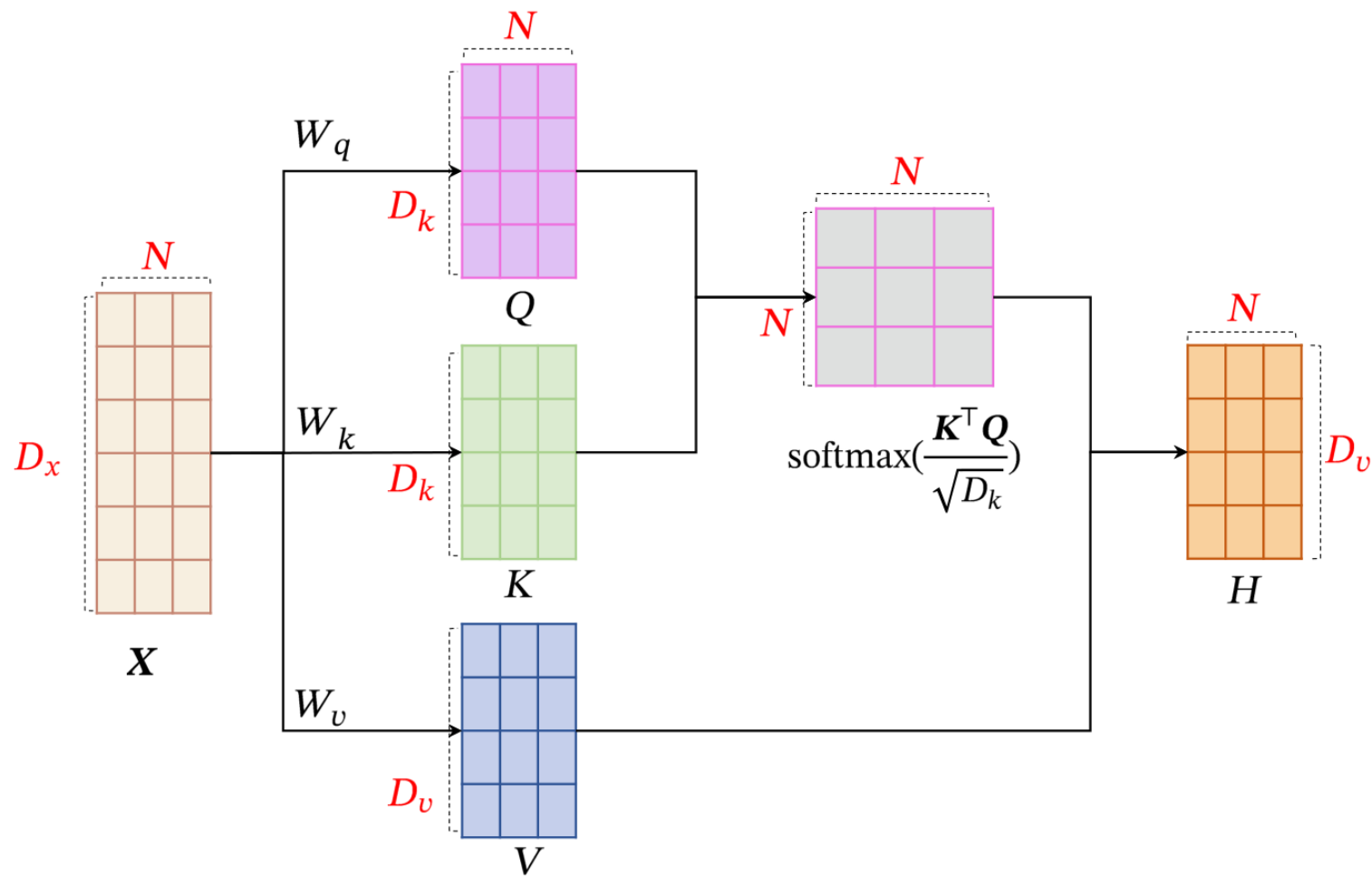
(a) 全连接模型

(b) 自注意力模型

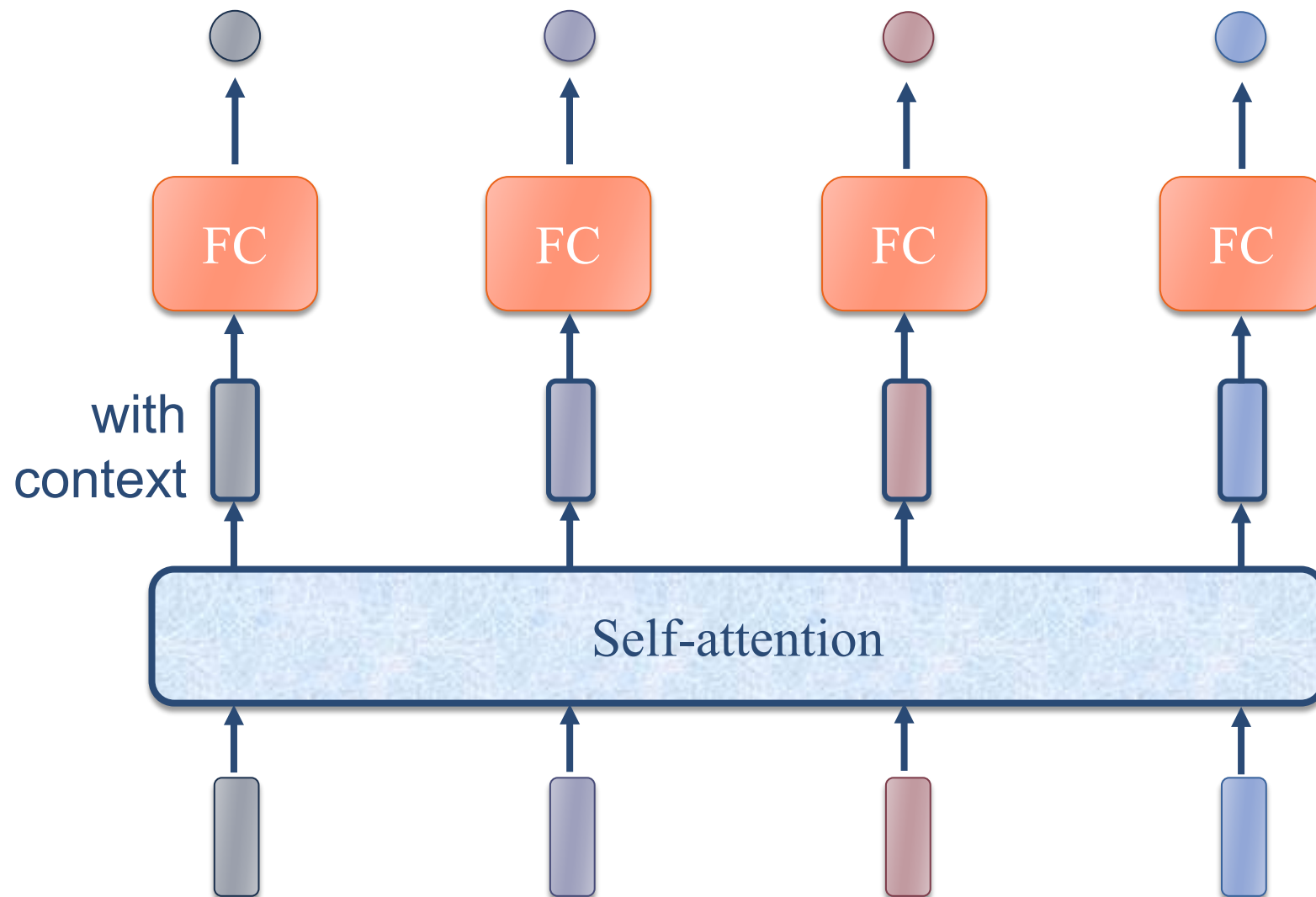
无法处理变长问题

连接权重 $\alpha_{ij}$  由注意力机制动态生成

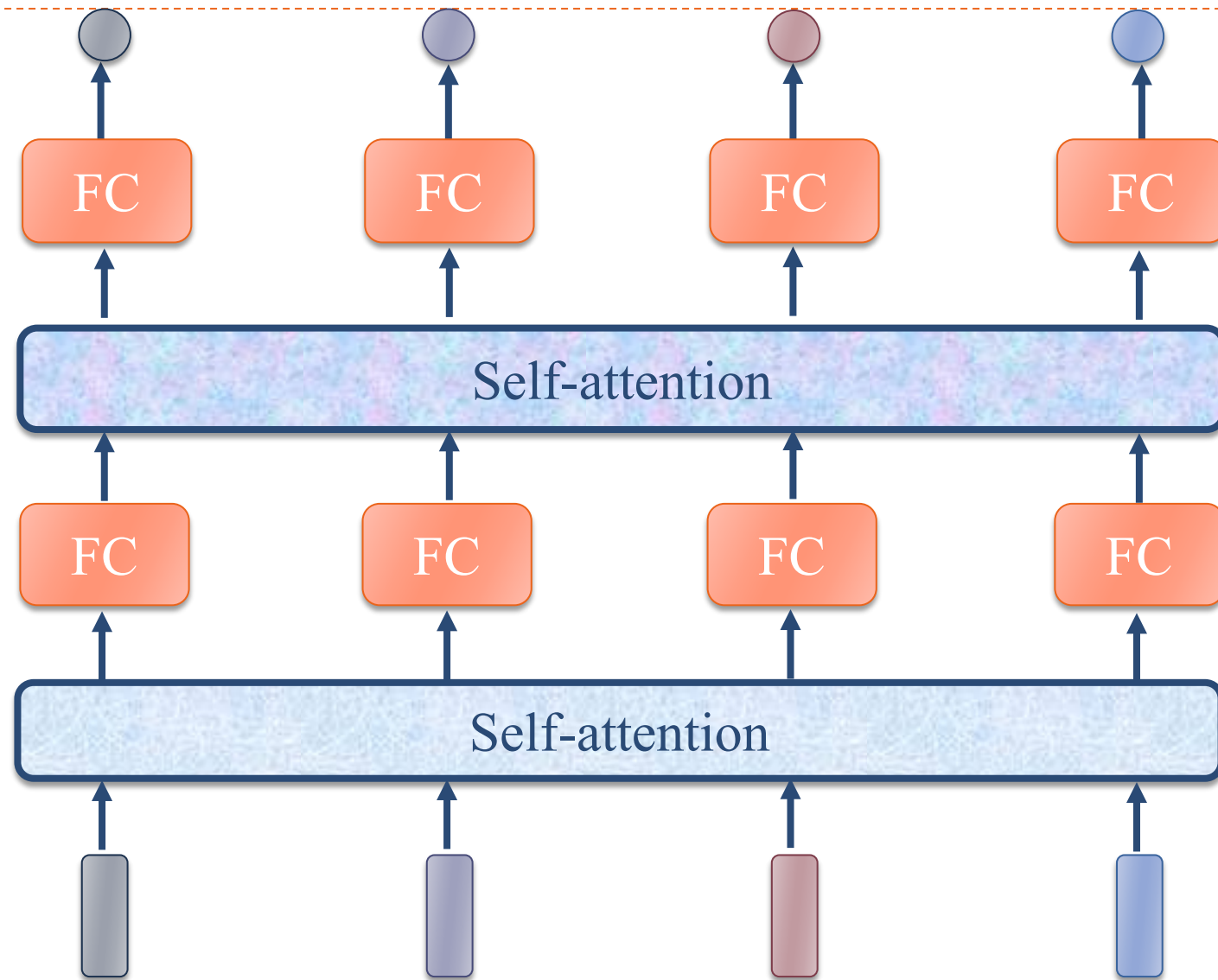
# QKV模式 (Query-Key-Value)



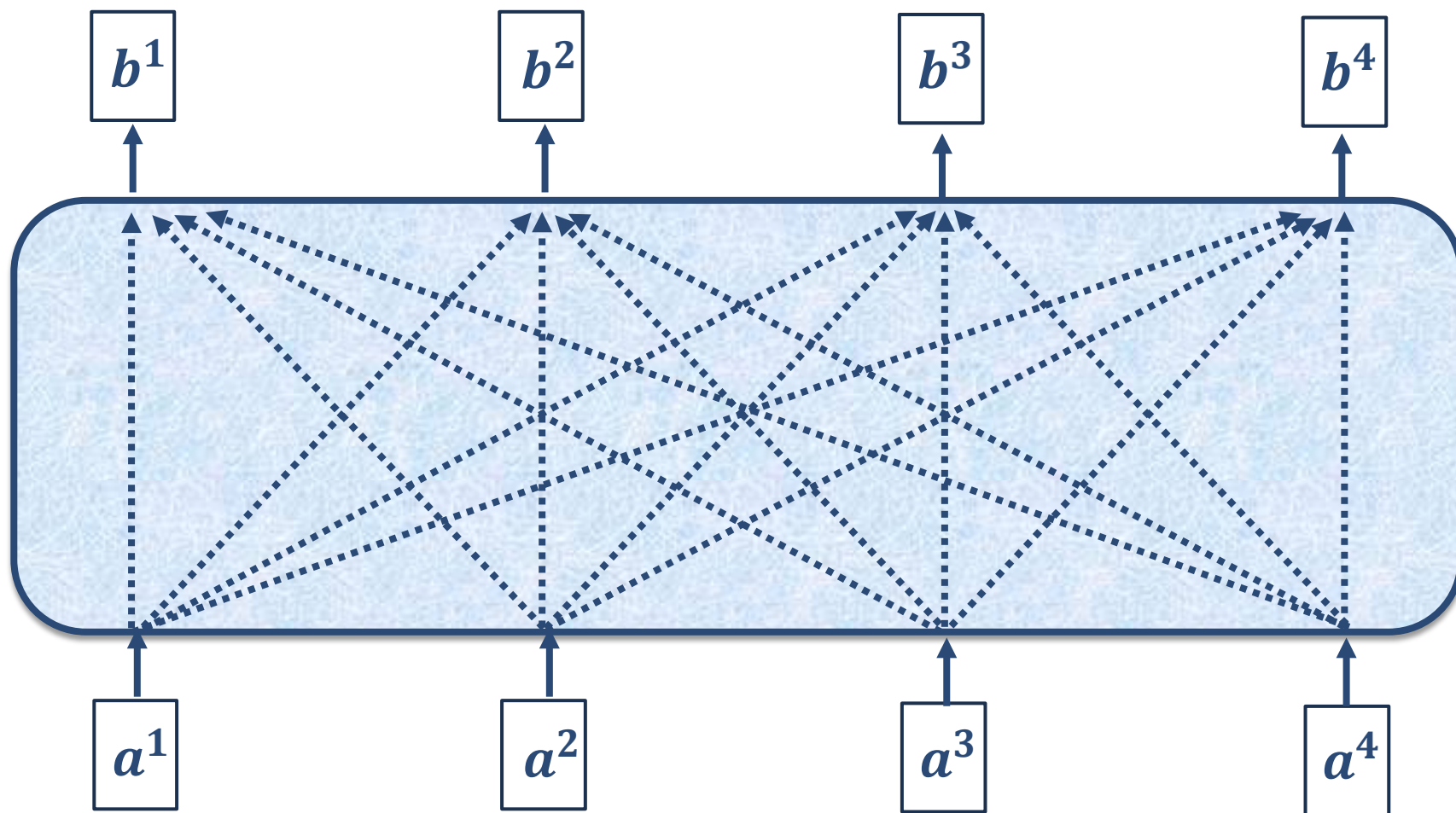
# 自注意力模型



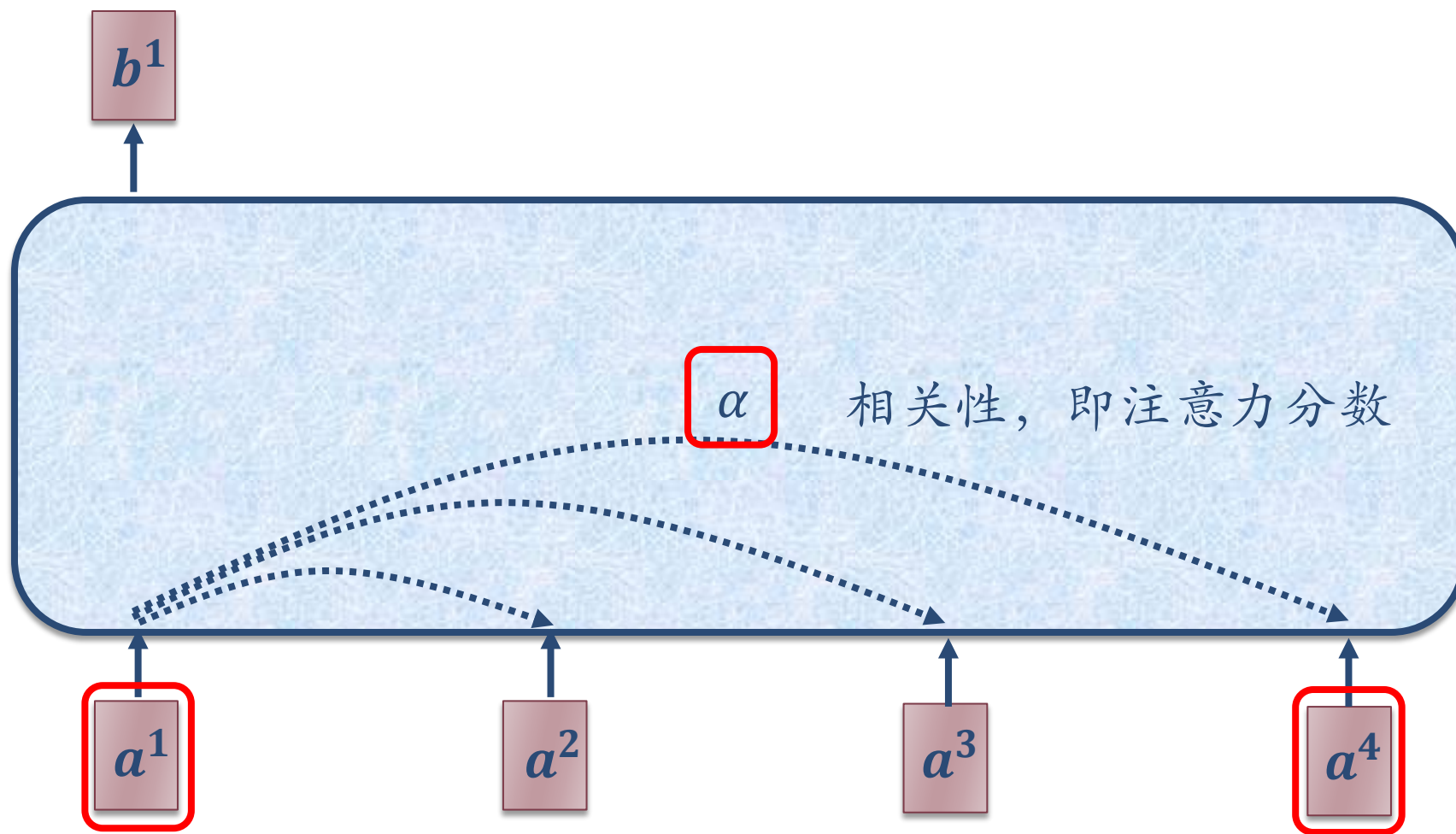
# 自注意力模型



# 自注意力模型



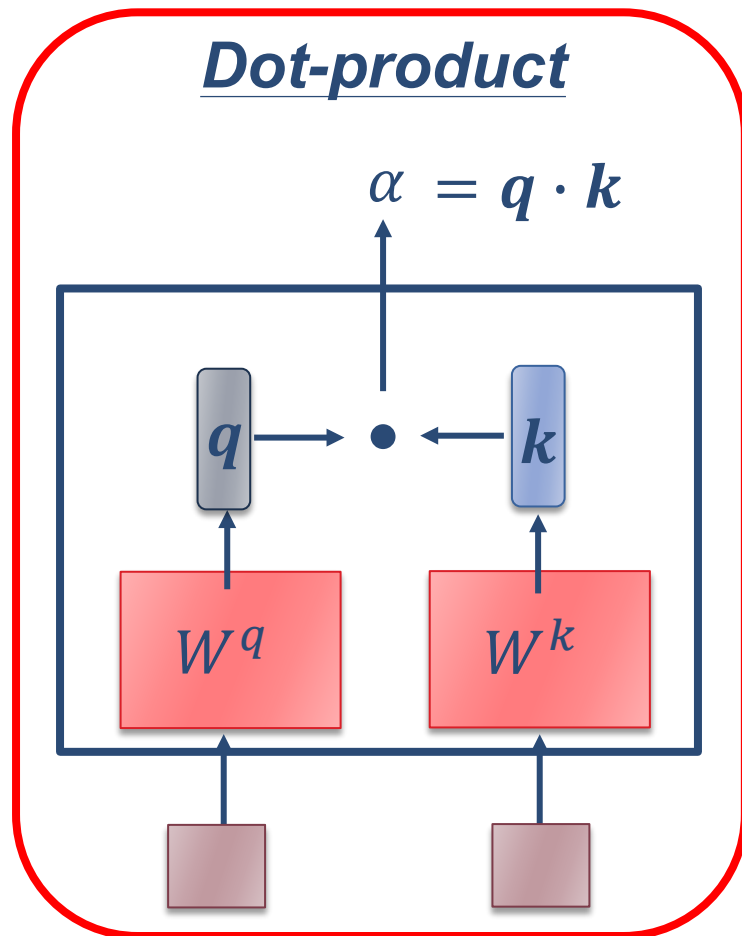
# 自注意力模型



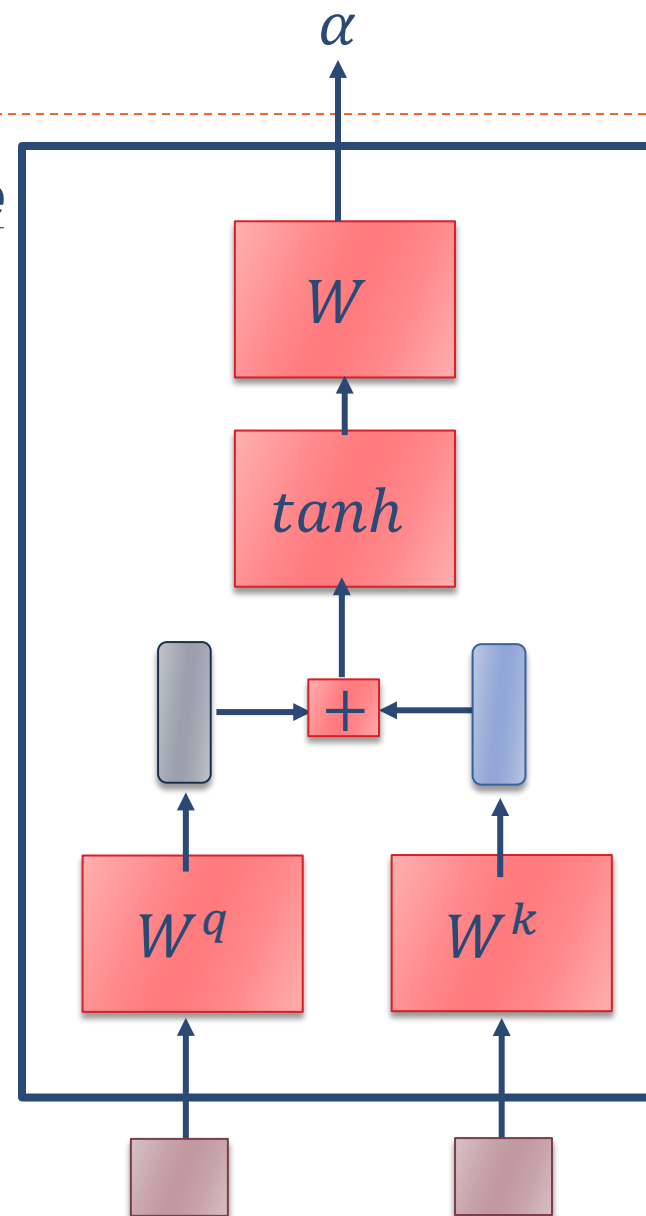
找出序列中相关的向量



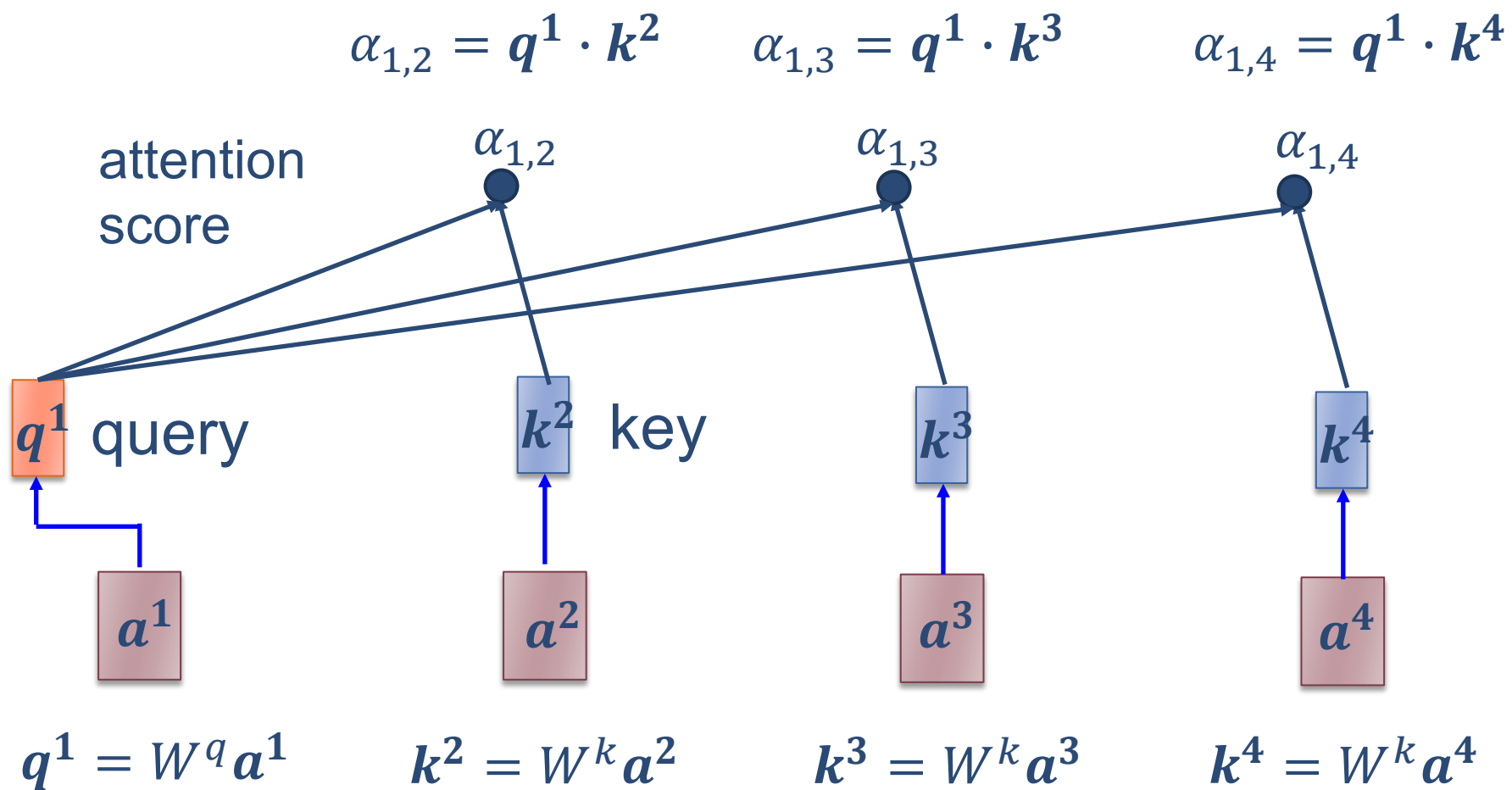
# 自注意力模型



Additive

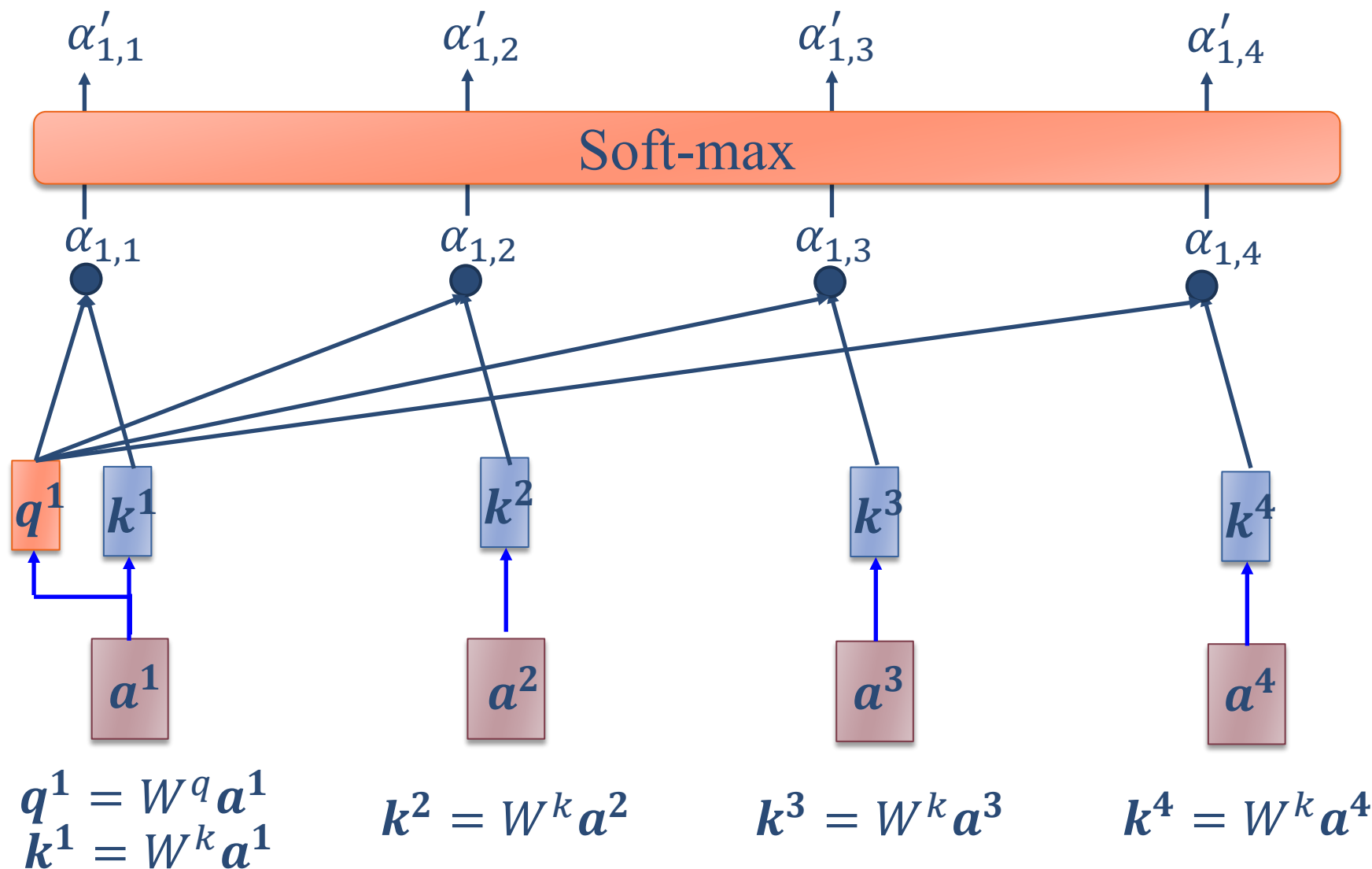


# 自注意力模型



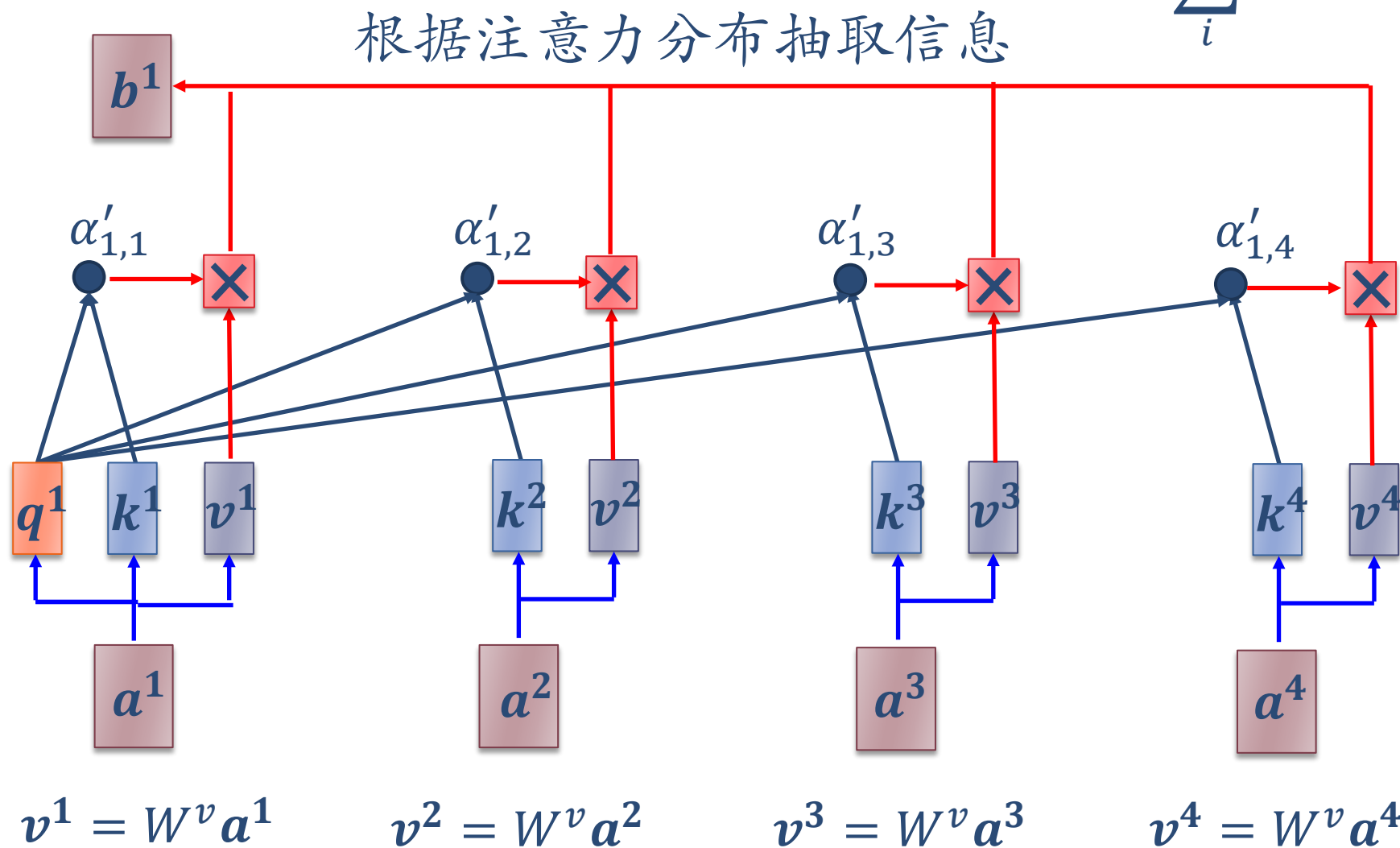
# 自注意力模型

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

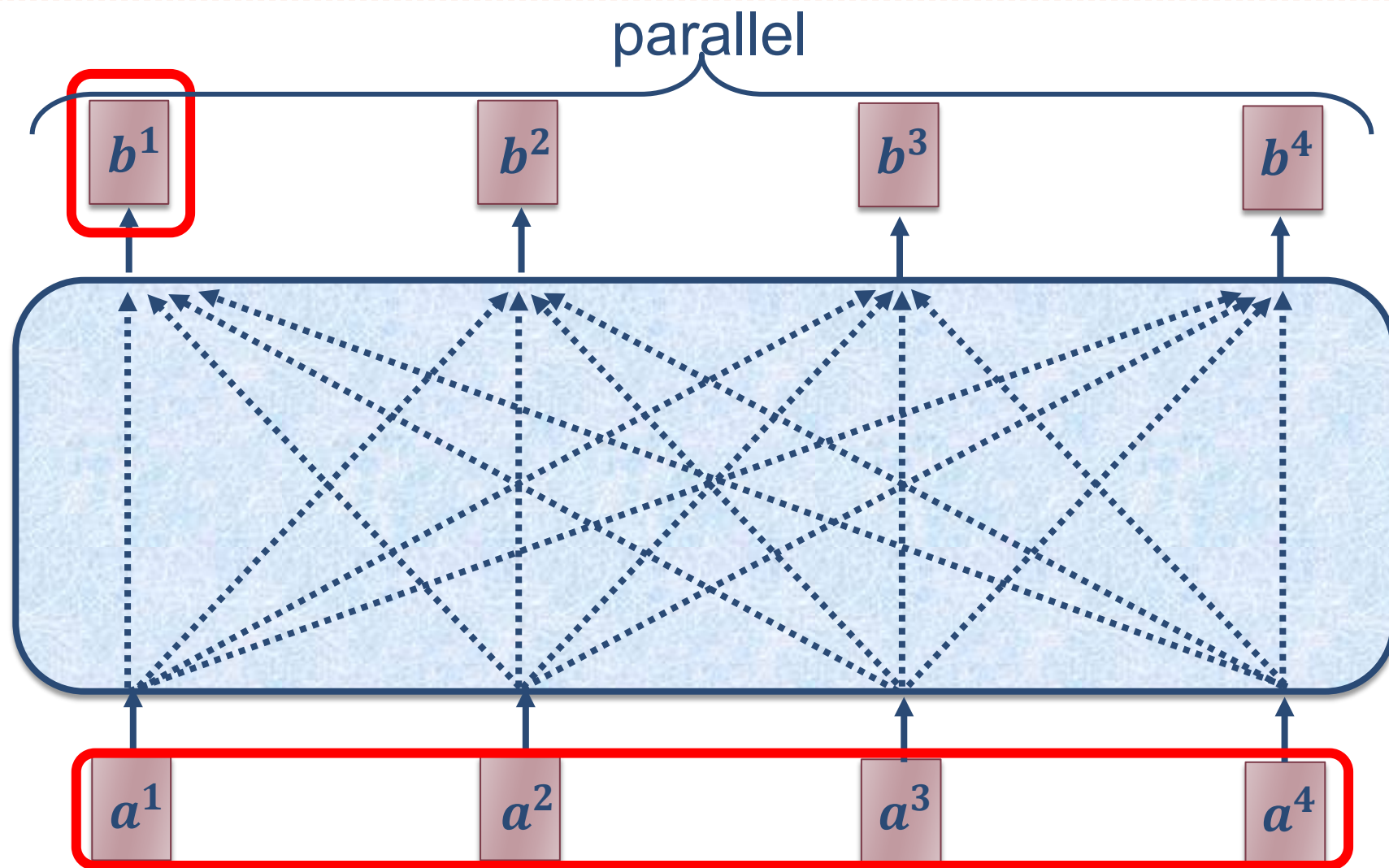


# 自注意力模型

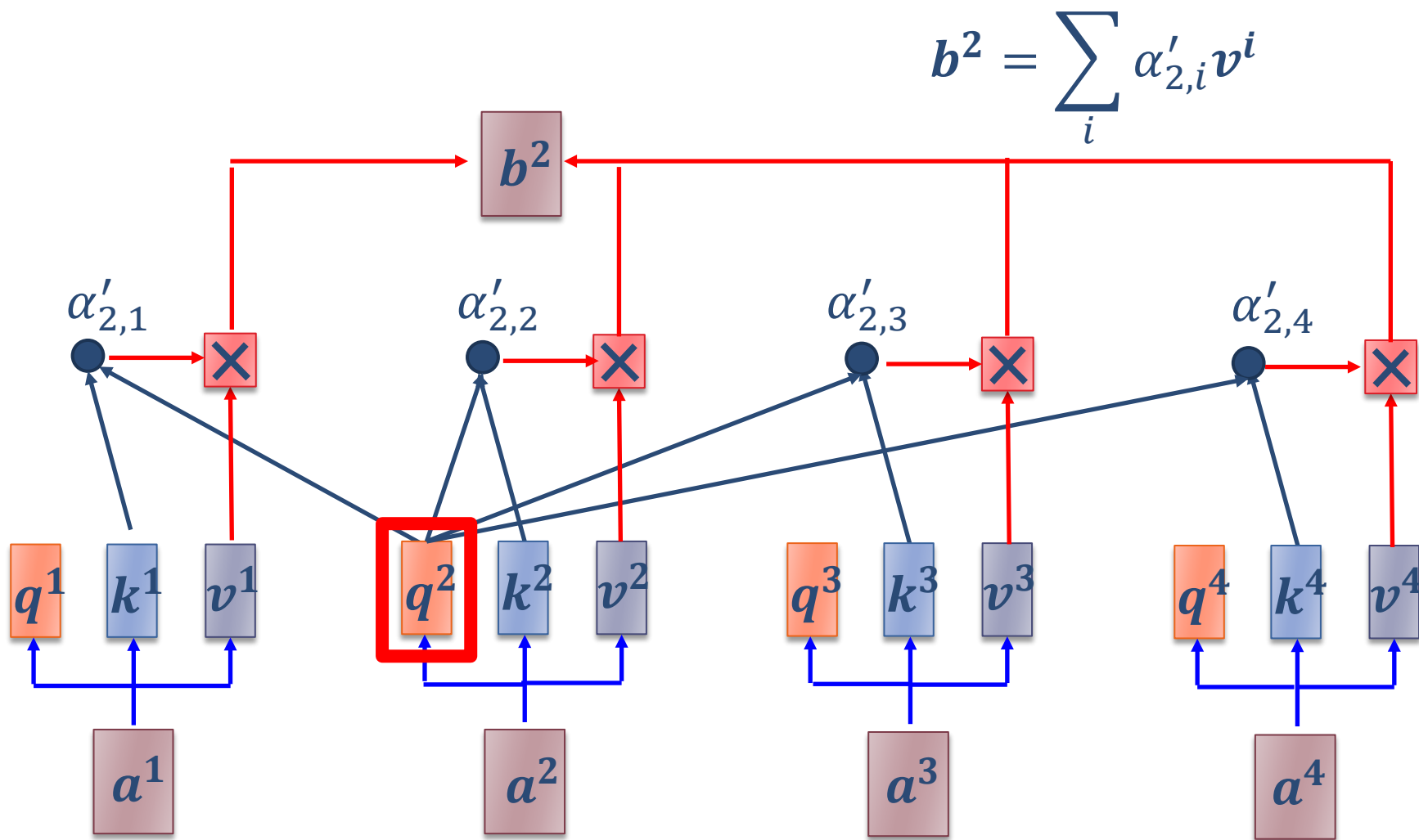
$$b^1 = \sum_i \alpha'_{1,i} v^i$$



# 自注意力模型



# 自注意力模型

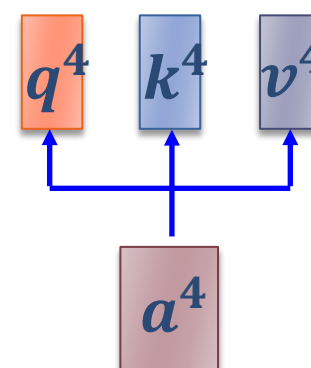
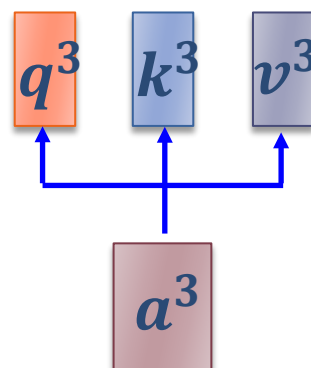
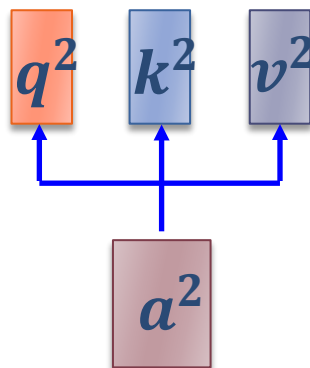
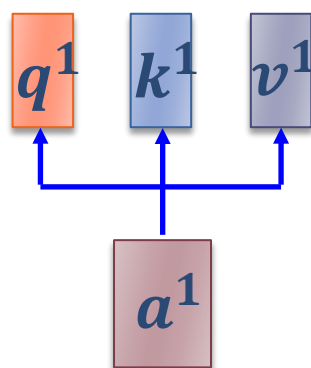


# 自注意力模型

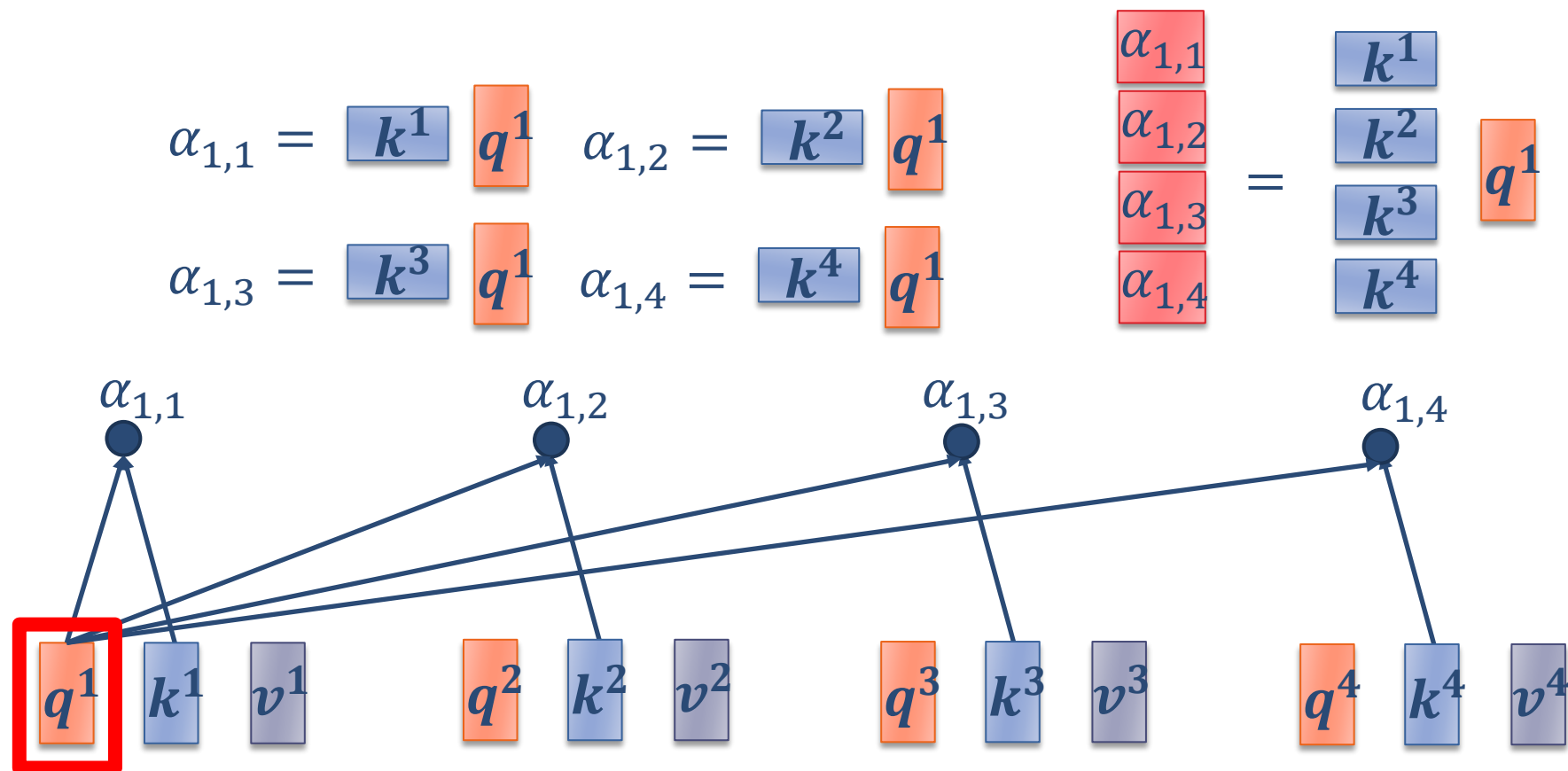
$$q^i = W^q a^i \quad \underbrace{q^1 q^2 q^3 q^4}_Q = \underbrace{W^q}_{\text{matrix}} \underbrace{a^1 a^2 a^3 a^4}_I$$

$$k^i = W^k a^i \quad \underbrace{k^1 k^2 k^3 k^4}_K = \underbrace{W^k}_{\text{matrix}} \underbrace{a^1 a^2 a^3 a^4}_I$$

$$v^i = W^v a^i \quad \underbrace{v^1 v^2 v^3 v^4}_V = \underbrace{W^v}_{\text{matrix}} \underbrace{a^1 a^2 a^3 a^4}_I$$

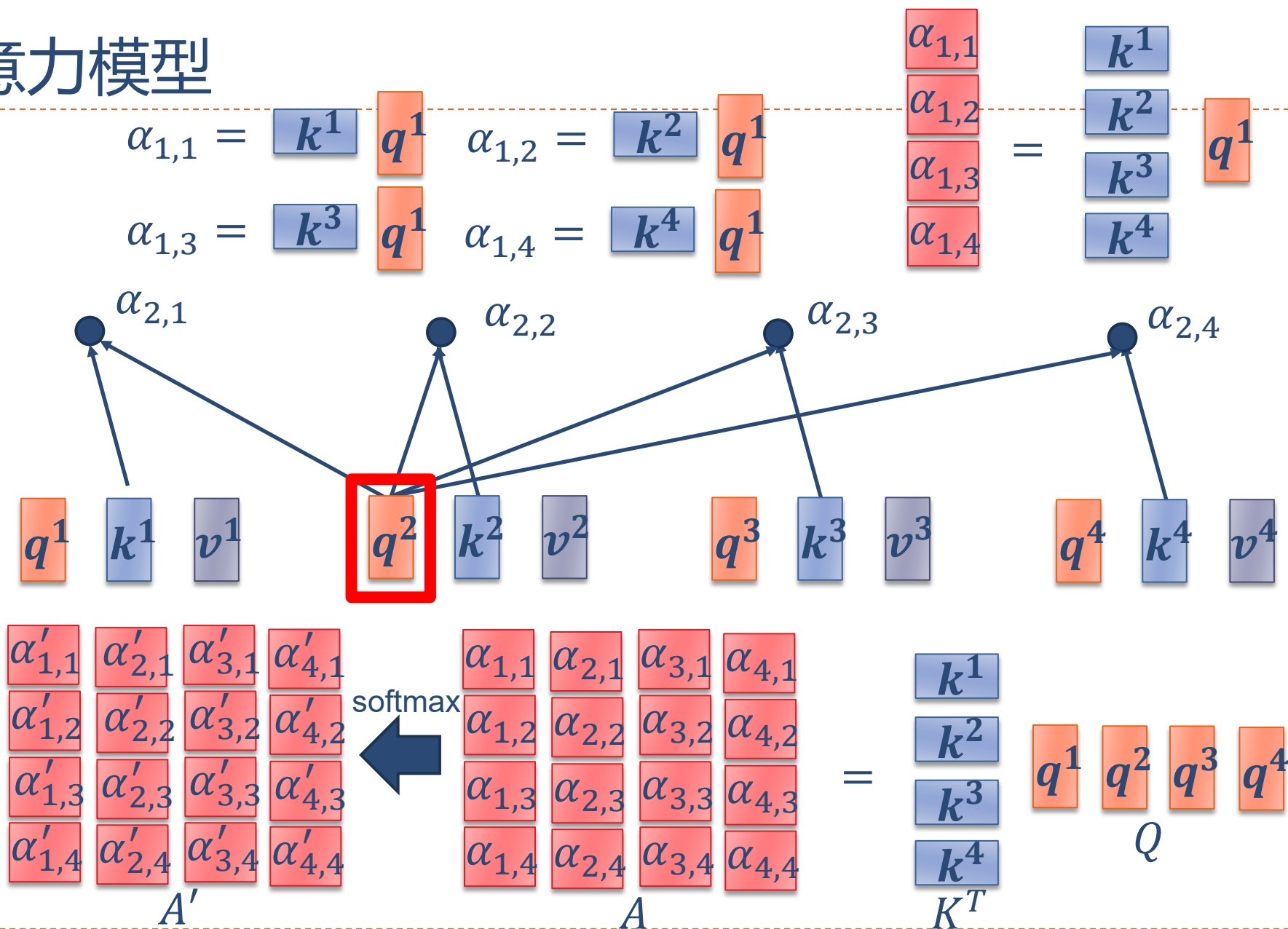


# 自注意力模型

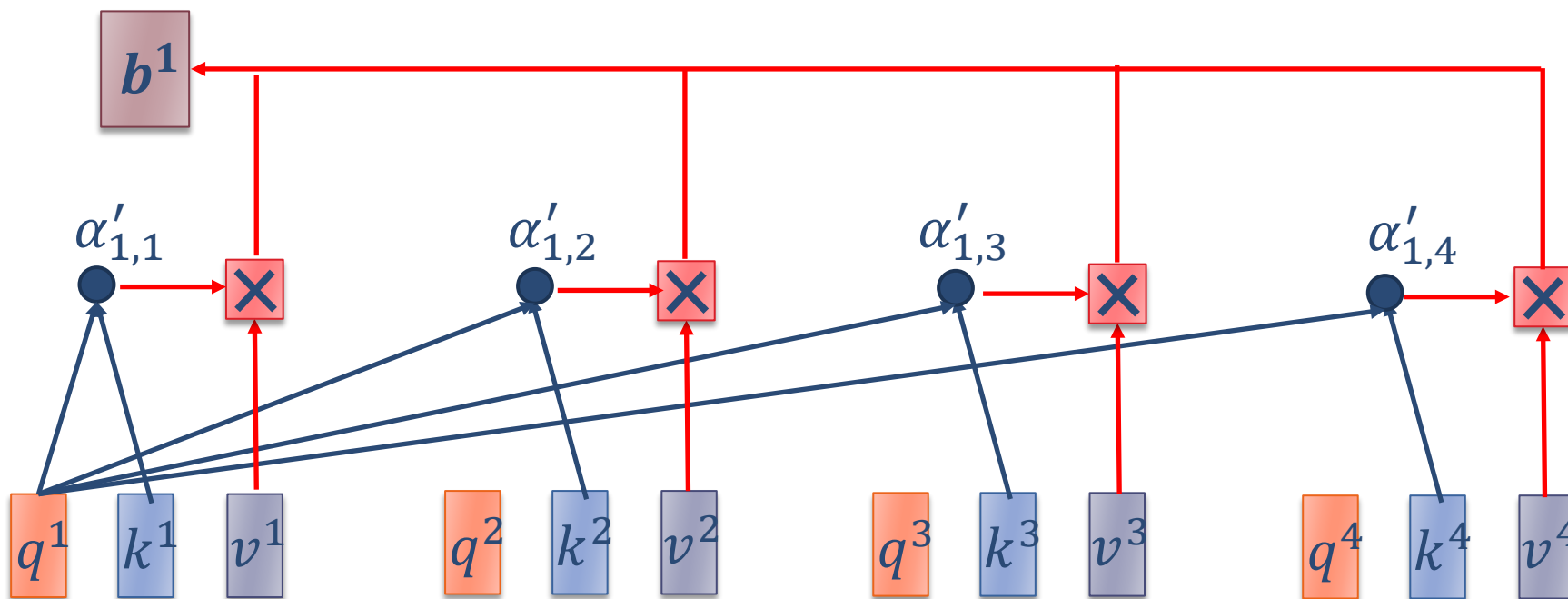




# 自注意力模型

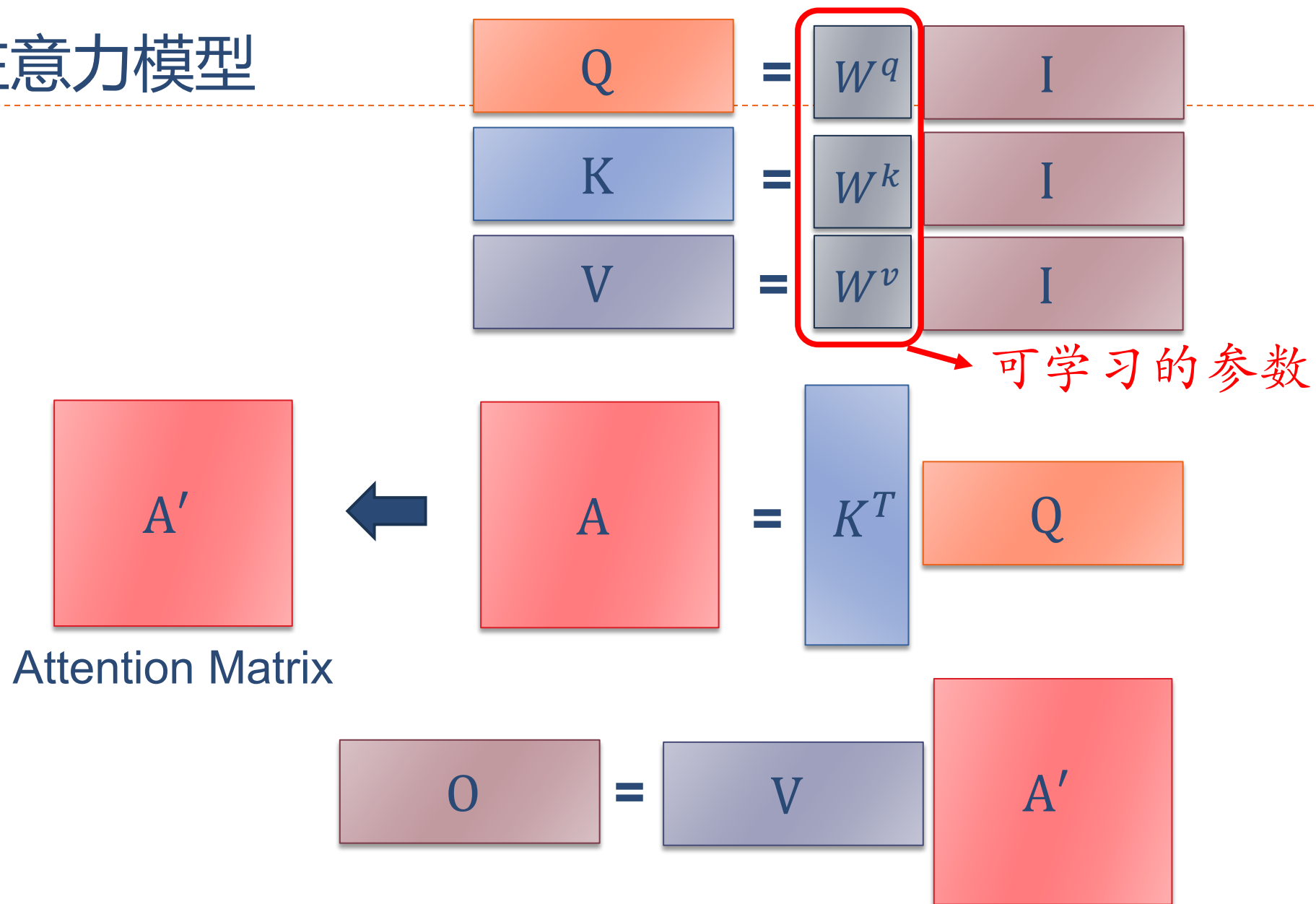


# 自注意力模型

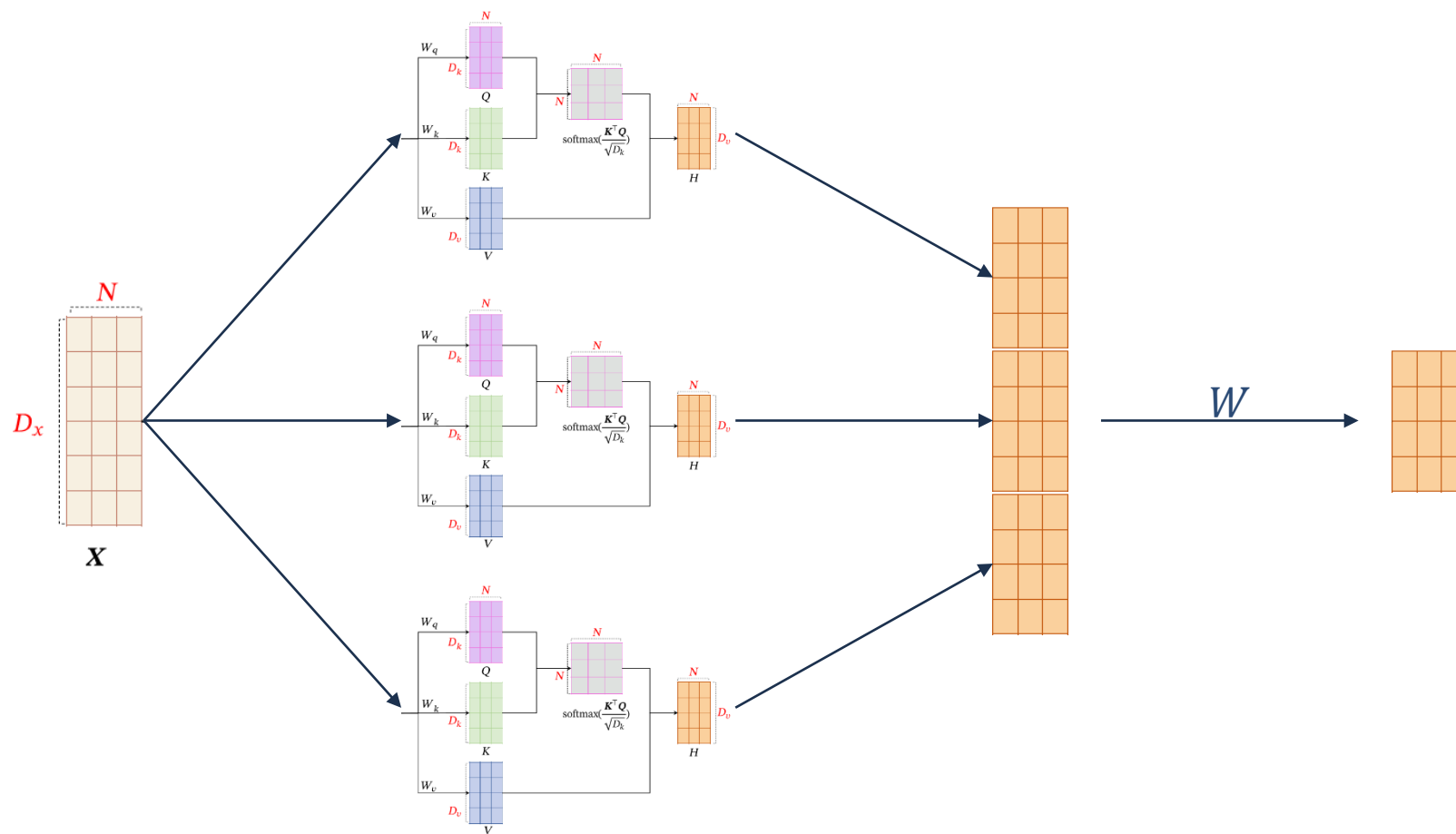


$$\begin{matrix} b^1 & b^2 & b^3 & b^4 \\ \hline \end{matrix} = \begin{matrix} v^1 & v^2 & v^3 & v^4 \\ \hline \end{matrix} \begin{matrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} &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

# 自注意力模型



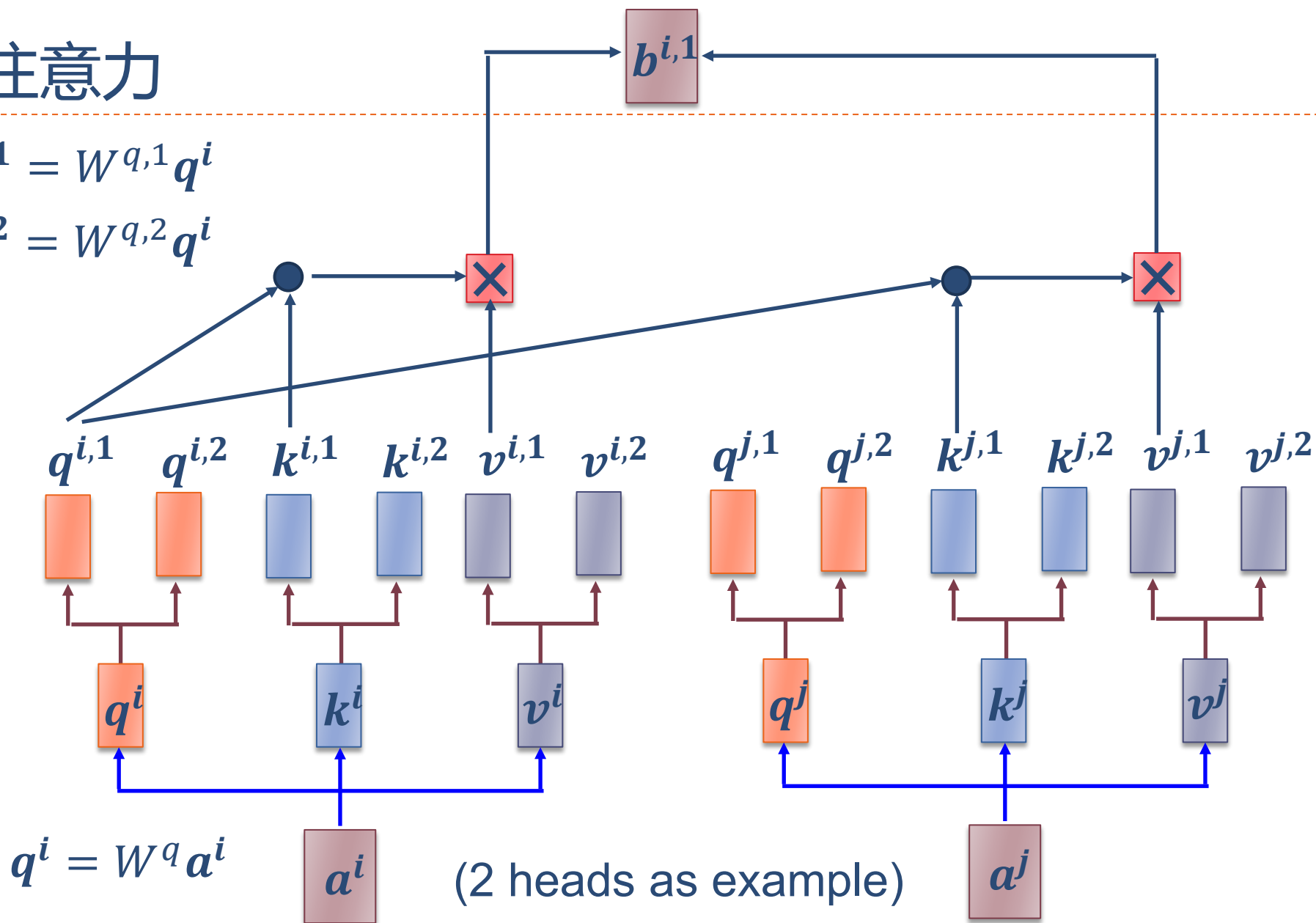
# 多头 (multi-head) 自注意力模型



# 多头自注意力

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

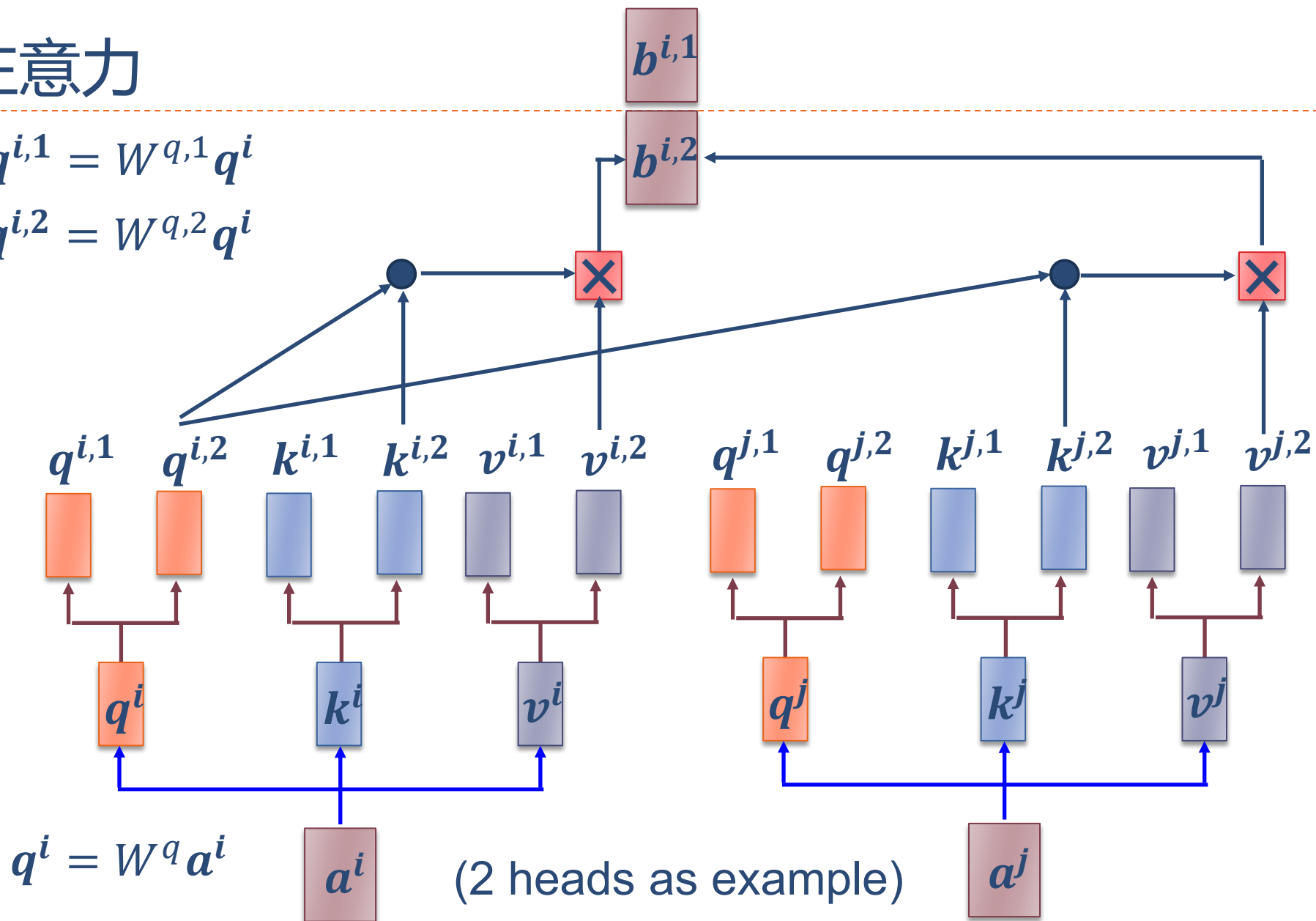
$$q^{i,2} = W^{q,2} q^i$$



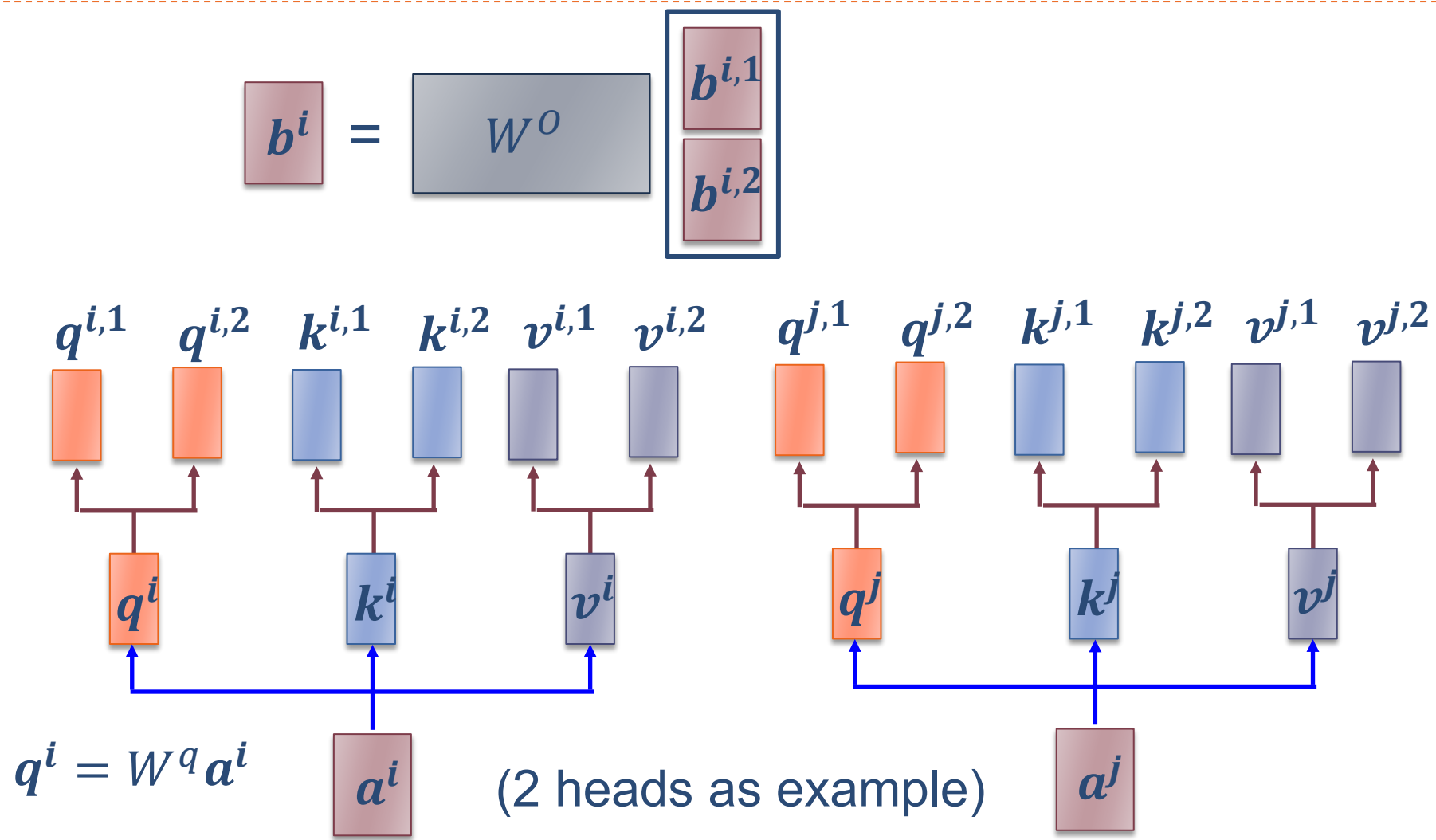
# 多头自注意力

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

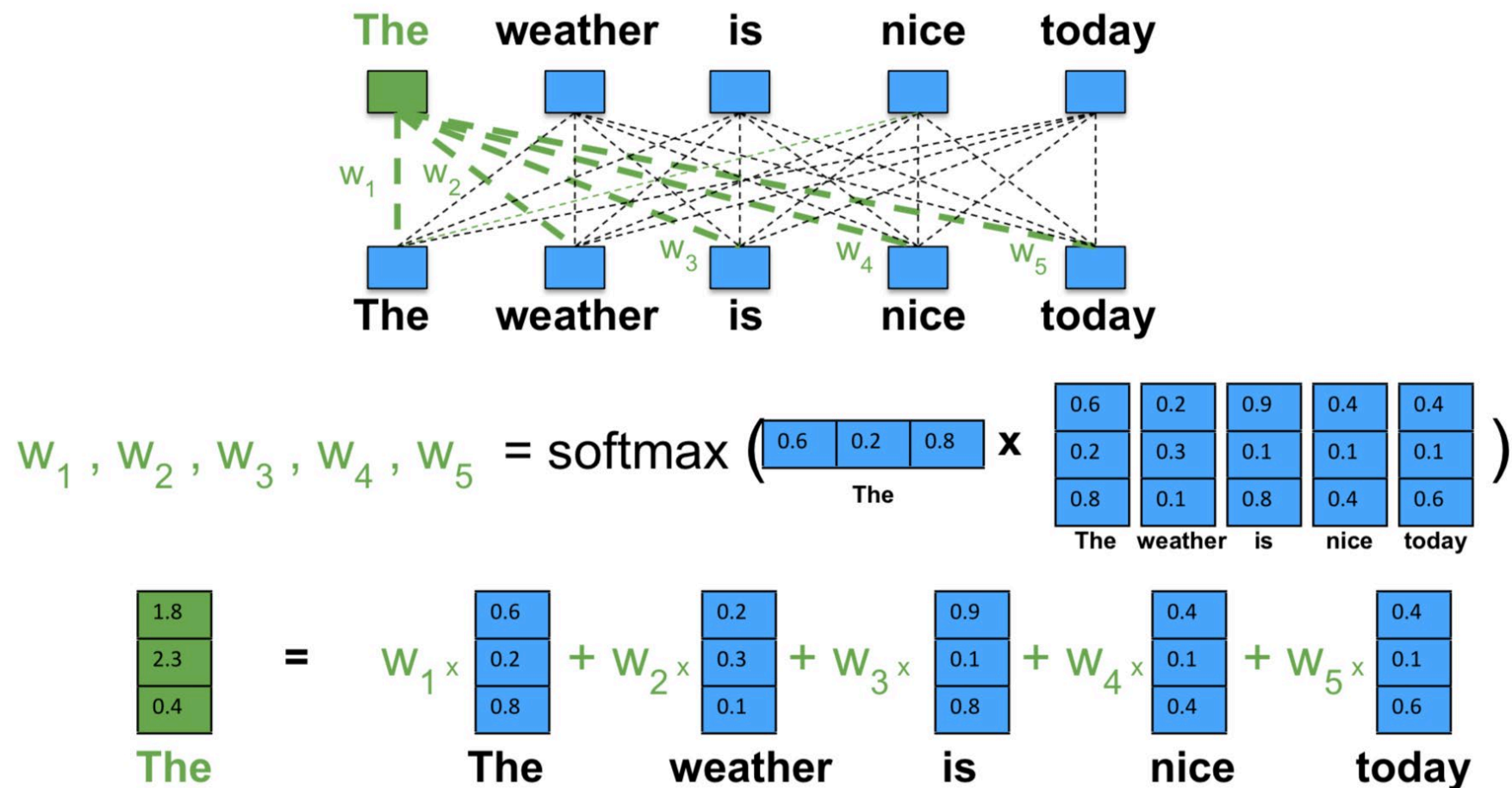
$$q^{i,2} = W^{q,2} q^i$$



# 多头自注意力



# 自注意力示例



图片来源: [http://fuyw.top/NLP\\_02\\_QANet/](http://fuyw.top/NLP_02_QANet/)



# 自注意力模型的形式化表示

▶ 输入序列为  $X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D_x \times N}$

▶ 首先生成三个向量序列

$$\mathbf{Q} = \mathbf{W}_q \mathbf{X} \in \mathbb{R}^{D_k \times N},$$

$$\mathbf{K} = \mathbf{W}_k \mathbf{X} \in \mathbb{R}^{D_k \times N},$$

$$\mathbf{V} = \mathbf{W}_v \mathbf{X} \in \mathbb{R}^{D_v \times N},$$

▶ 计算  $\mathbf{h}_n$

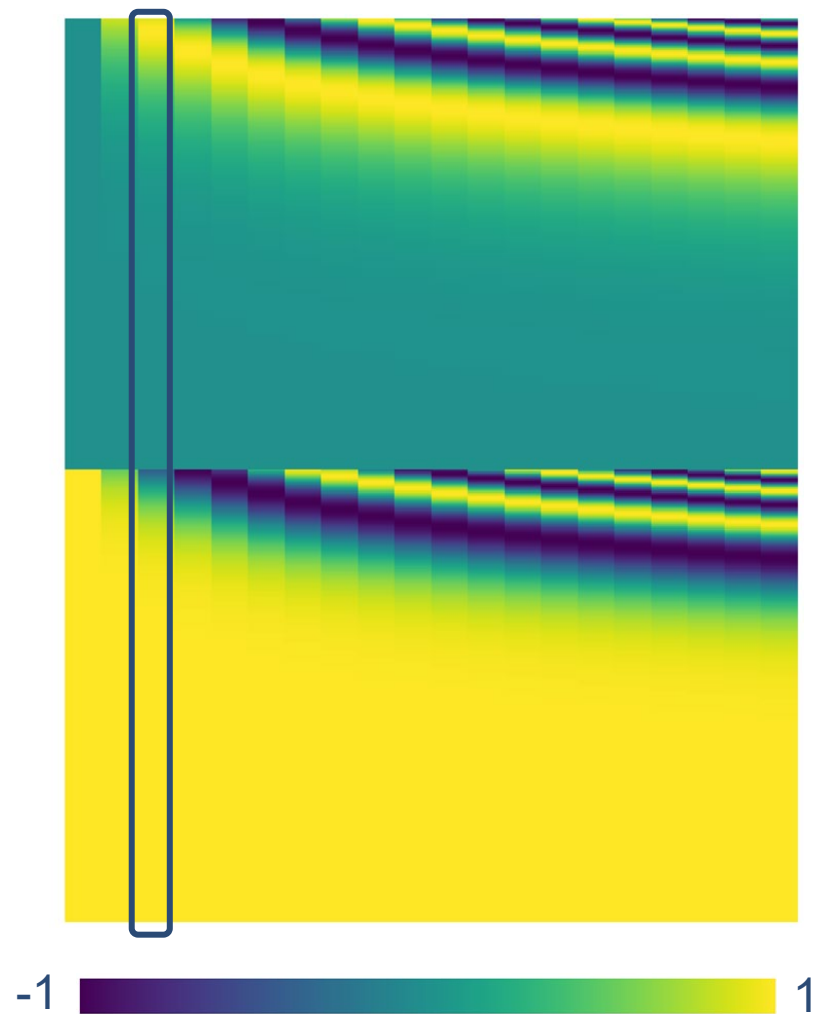
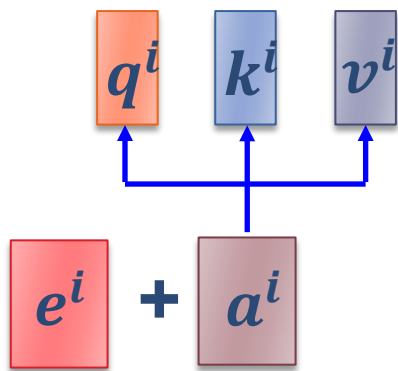
$$\mathbf{h}_n = \text{att}((\mathbf{K}, \mathbf{V}), \mathbf{q}_n)$$

▶ 如果使用缩放点积来作为注意力打分函数，输出向量序列可以简写为

$$\mathbf{H} = \mathbf{V} \text{softmax}\left(\frac{\mathbf{K}^\top \mathbf{Q}}{\sqrt{D_k}}\right),$$

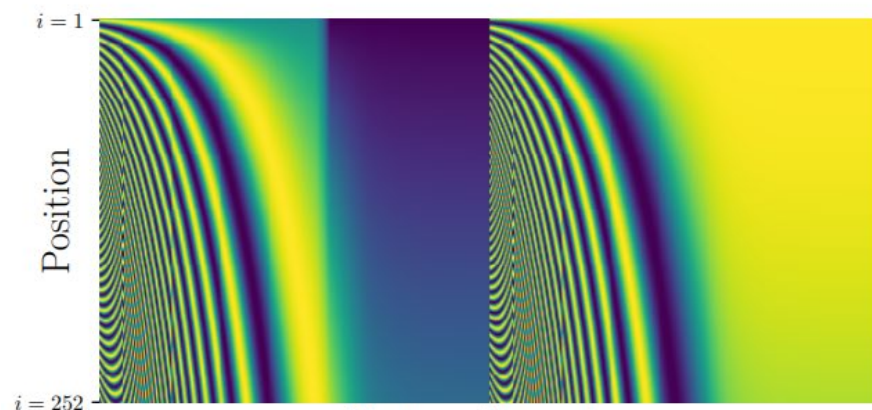
# 位置编码

- ▶ 自注意力机制没有位置信息
- ▶ 解决办法：为每个位置设置一个位置向量  $e^i$
- ▶ 把位置向量跟原始向量相加
- ▶ 如何设计位置编码？
  - ▶ 人工设计
  - ▶ 从数据中学习



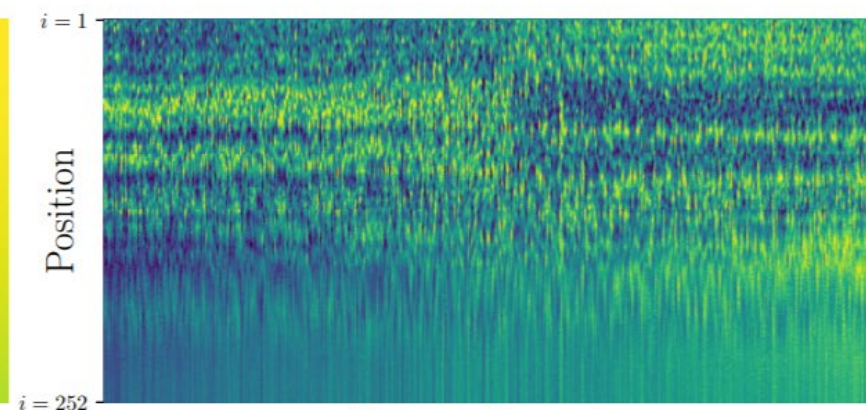
# 位置编码

(a) Sinusoidal

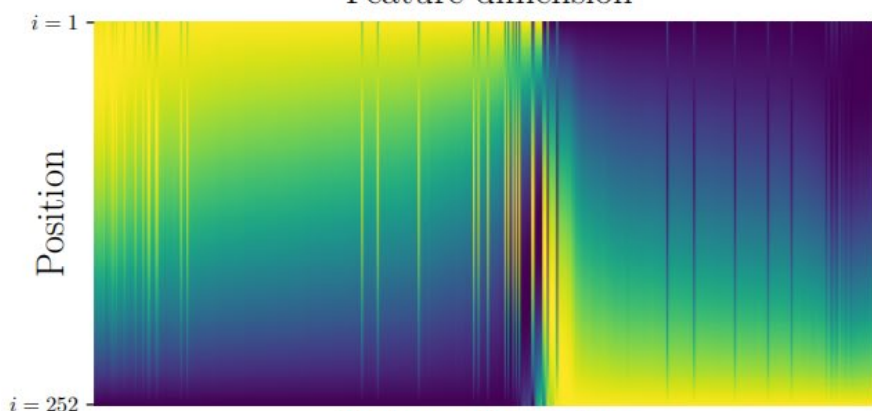


Feature dimension

(b) Position embedding



Feature dimension



Feature dimension

(c) FLOATER



Feature dimension

(d) RNN

# 自注意力的应用

---



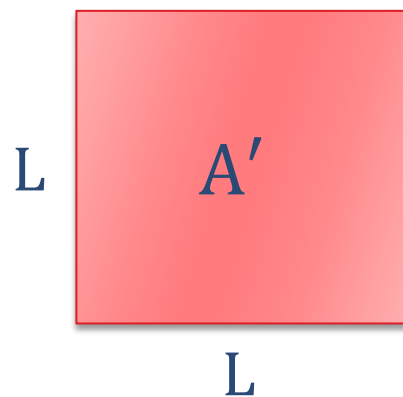
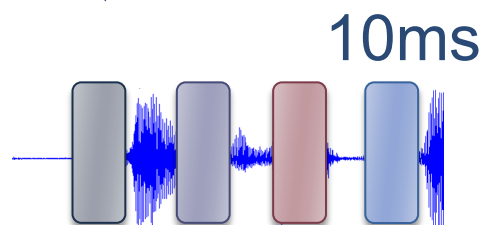
Transformer



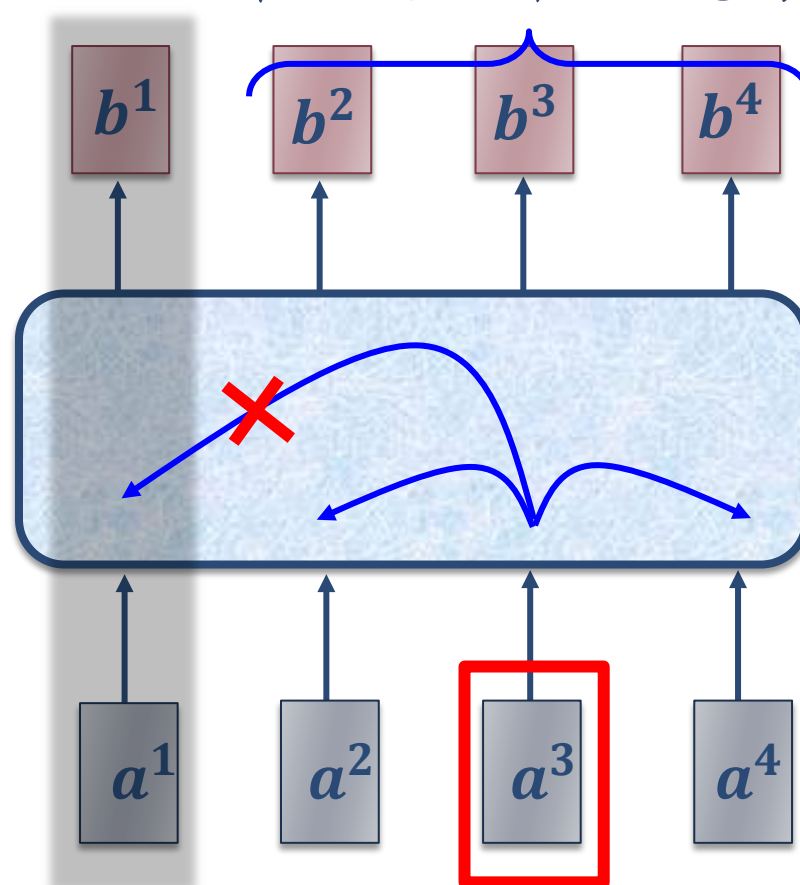
BERT

# 自注意力的应用

语音序列非常长



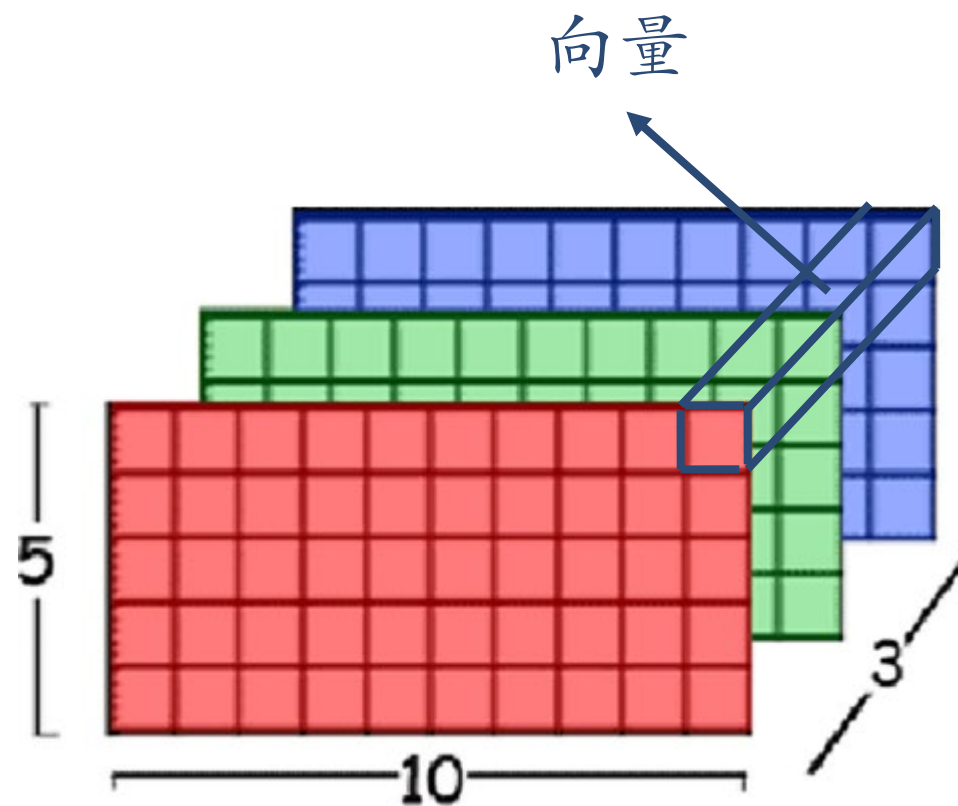
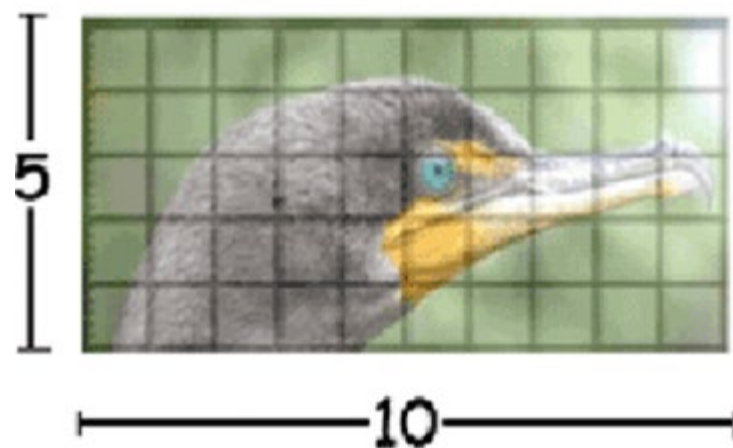
取一定范围的注意力



Truncated Self-attention

# 自注意力的应用

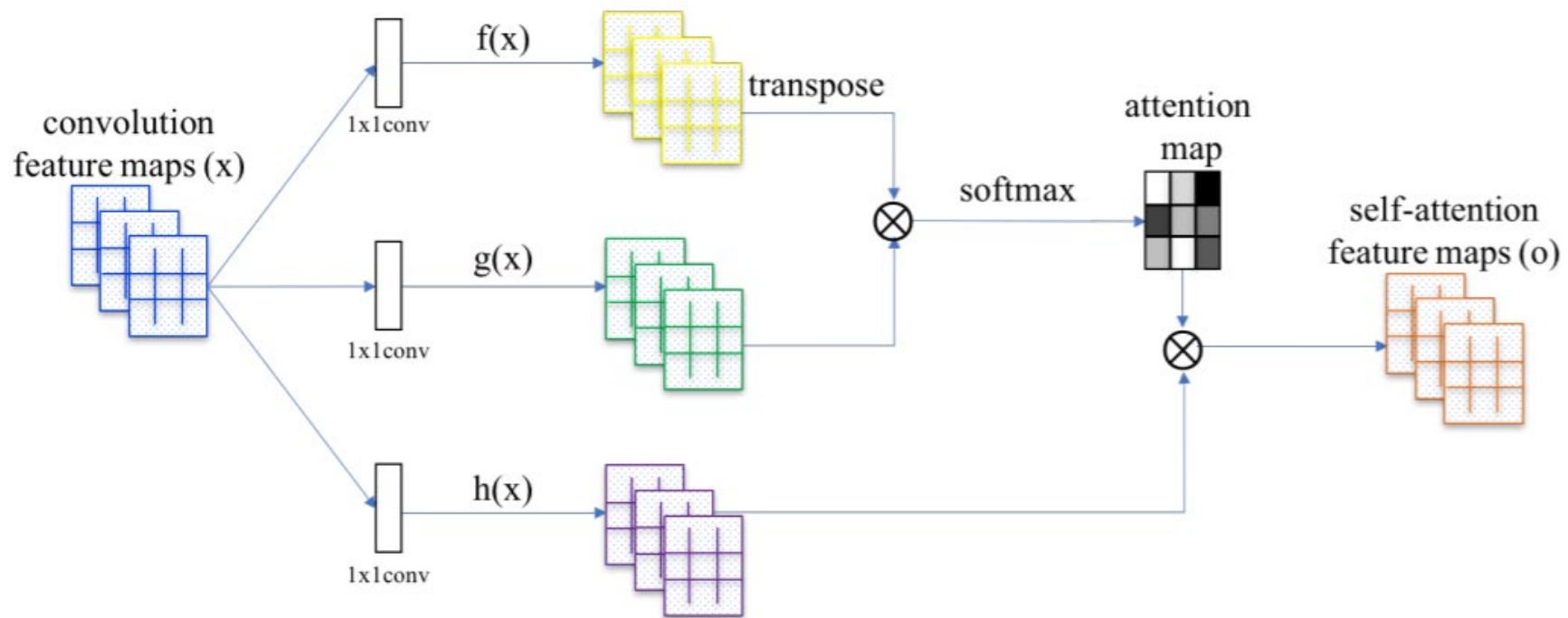
把图片看成向量的集合





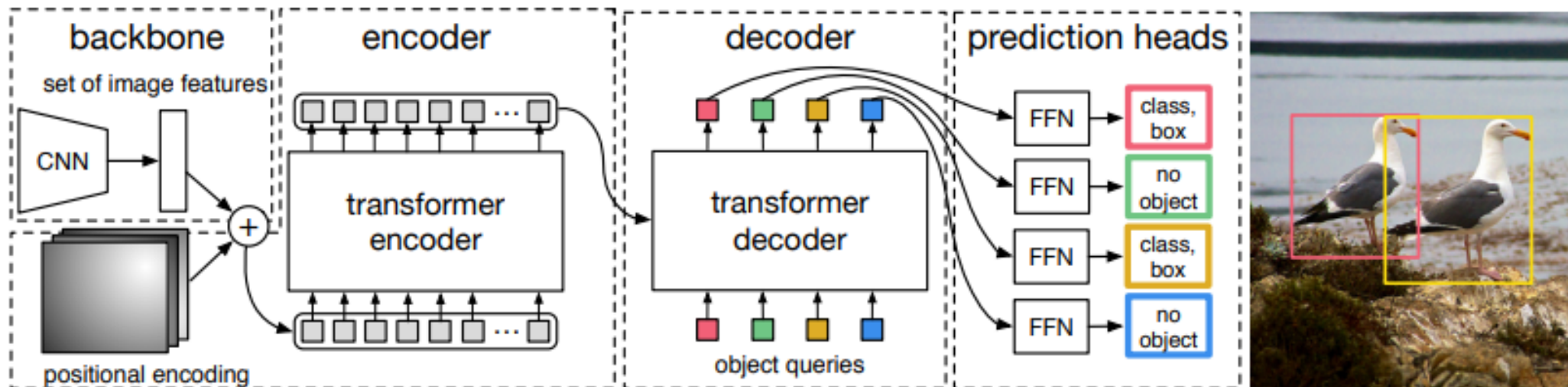
# 自注意力的应用

## Self-Attention GAN



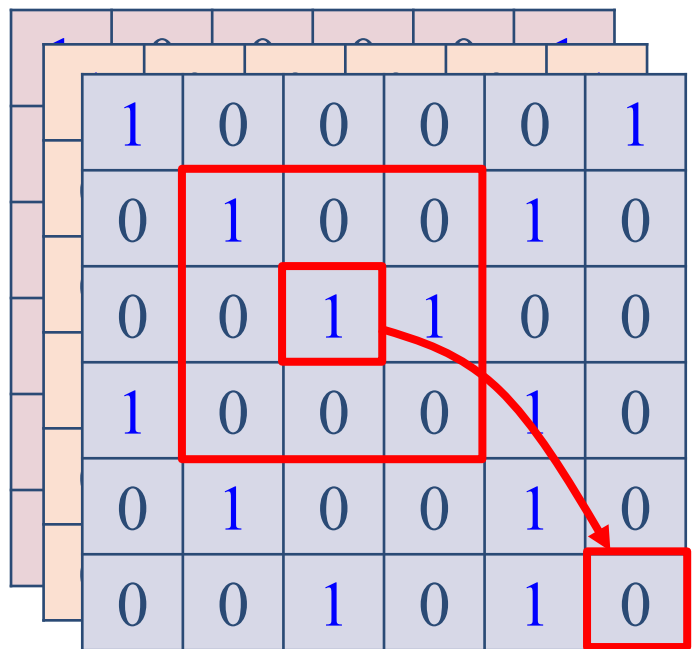
# 自注意力的应用

## DEtection Transformer (DETR)





# 比较：自注意力与卷积神经网络



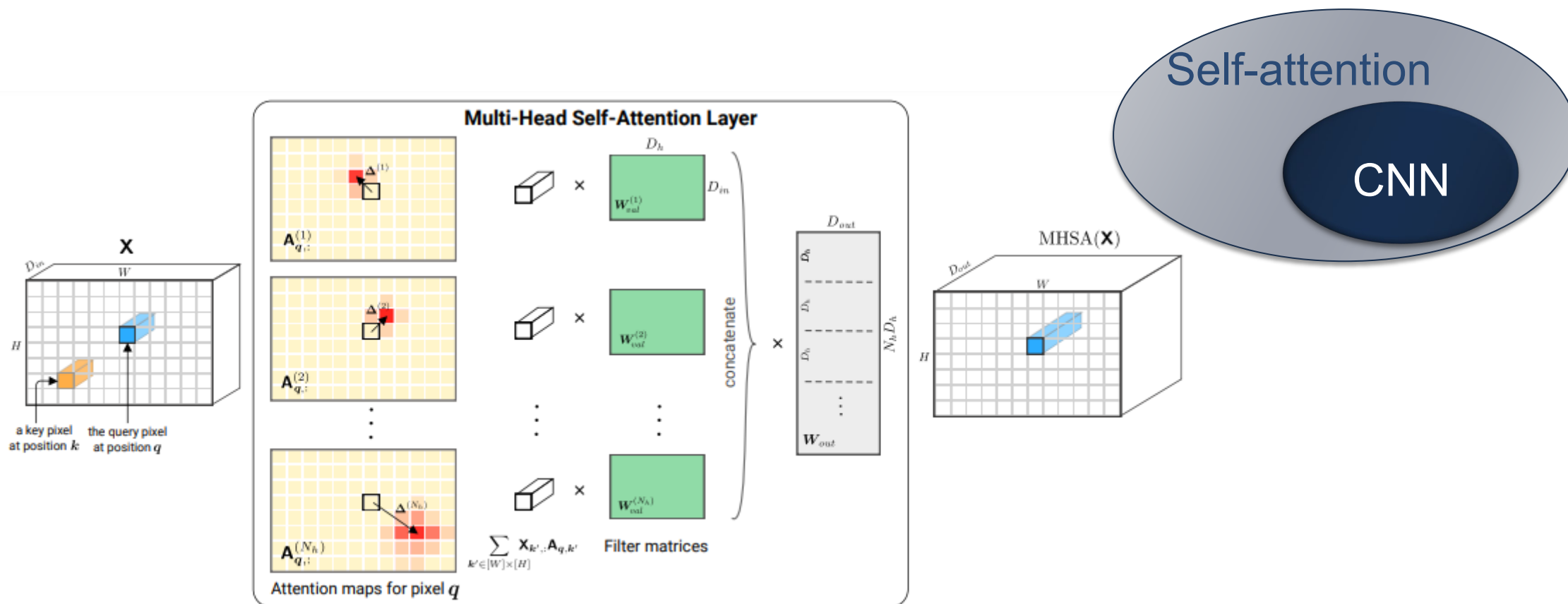
CNN: 可以看成一种限制在局部感受野的自注意力

➤ CNN是一种简化版的自注意力

自注意力: 有可学习感受野的CNN

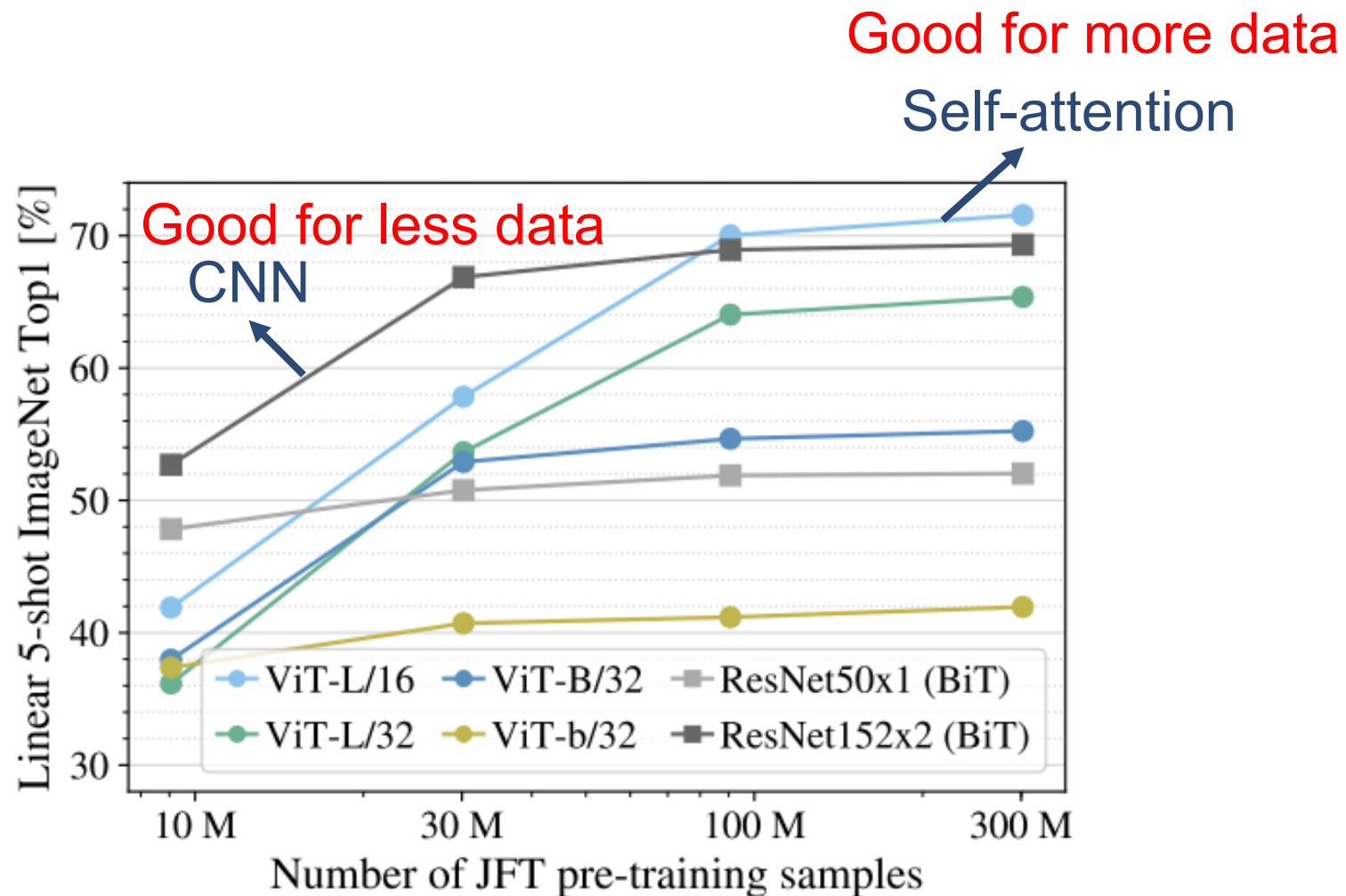
➤ 自注意力是一种复杂版本的CNN

# 比较：自注意力与卷积神经网络

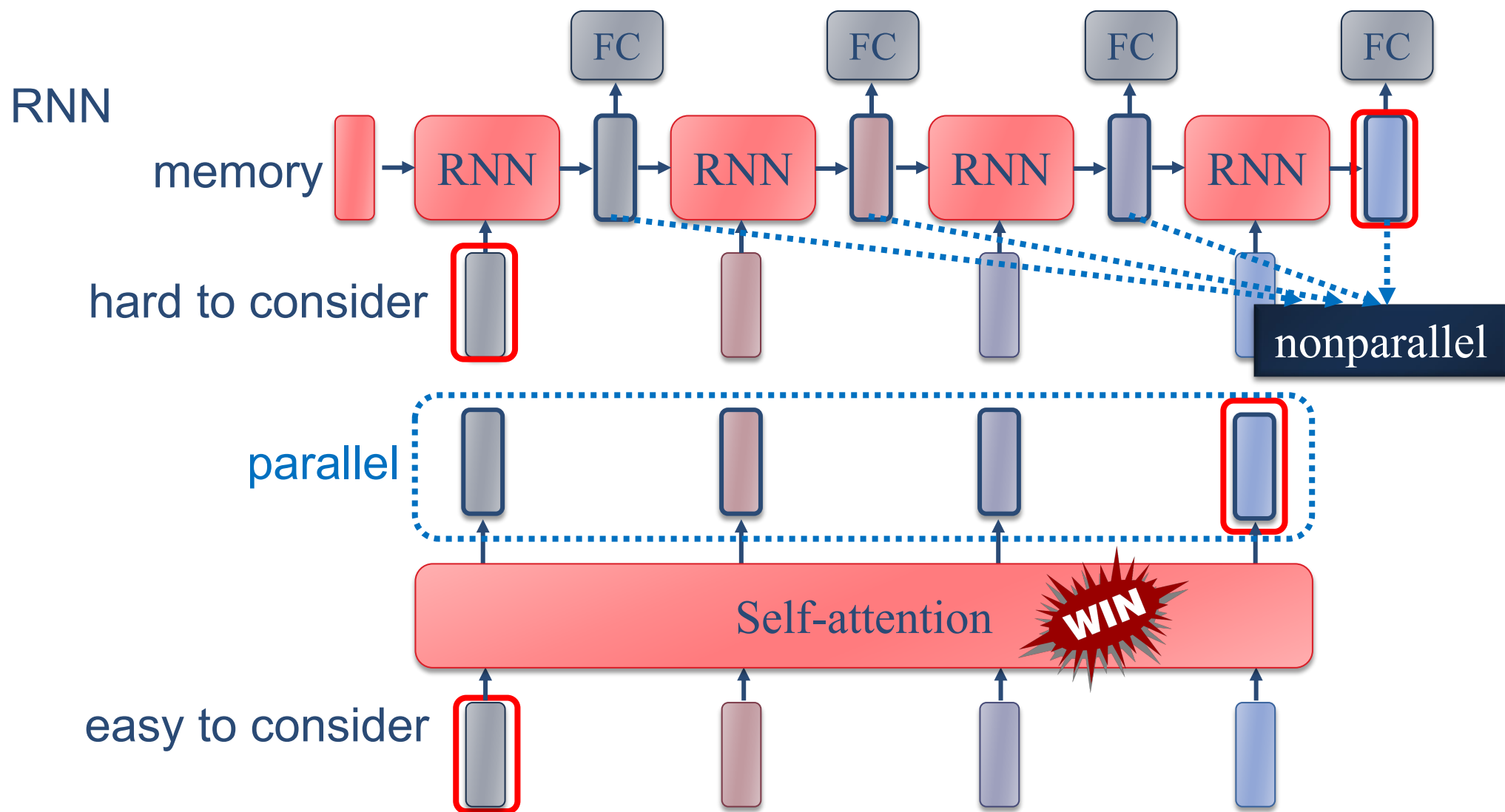


On the Relationship between Self-Attention and Convolutional Layers

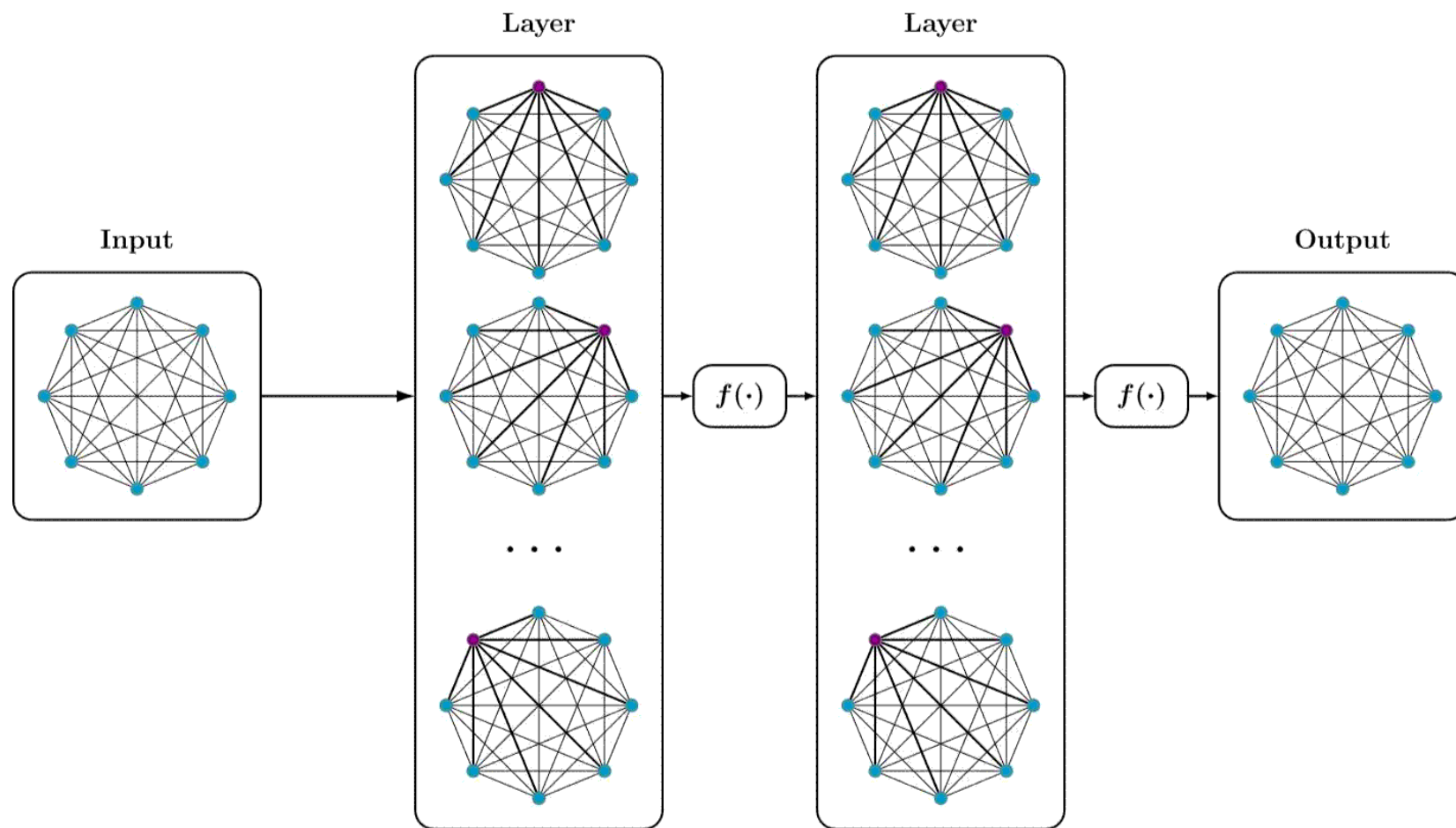
# 比较：自注意力与卷积神经网络



# 比较：自注意力与循环神经网络



# 比较：自注意力与循环神经网络



# 复杂度分析

模型	每层复杂度	序列操作数	最大路径长度
CNN	$O(kLd^2)$	$O(1)$	$O(\log_k(L))$
RNN	$O(Ld^2)$	$O(L)$	$O(L)$
自注意力	$O(L^2d)$	$O(1)$	$O(1)$

$k$ 卷积核大小     $L$ 序列长度     $d$ 维度



Transformer

---

## Attention Is All You Need

---

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\*<sup>†</sup>**  
University of Toronto  
aidan@cs.toronto.edu

**Łukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

**Illia Polosukhin\*<sup>‡</sup>**  
illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions



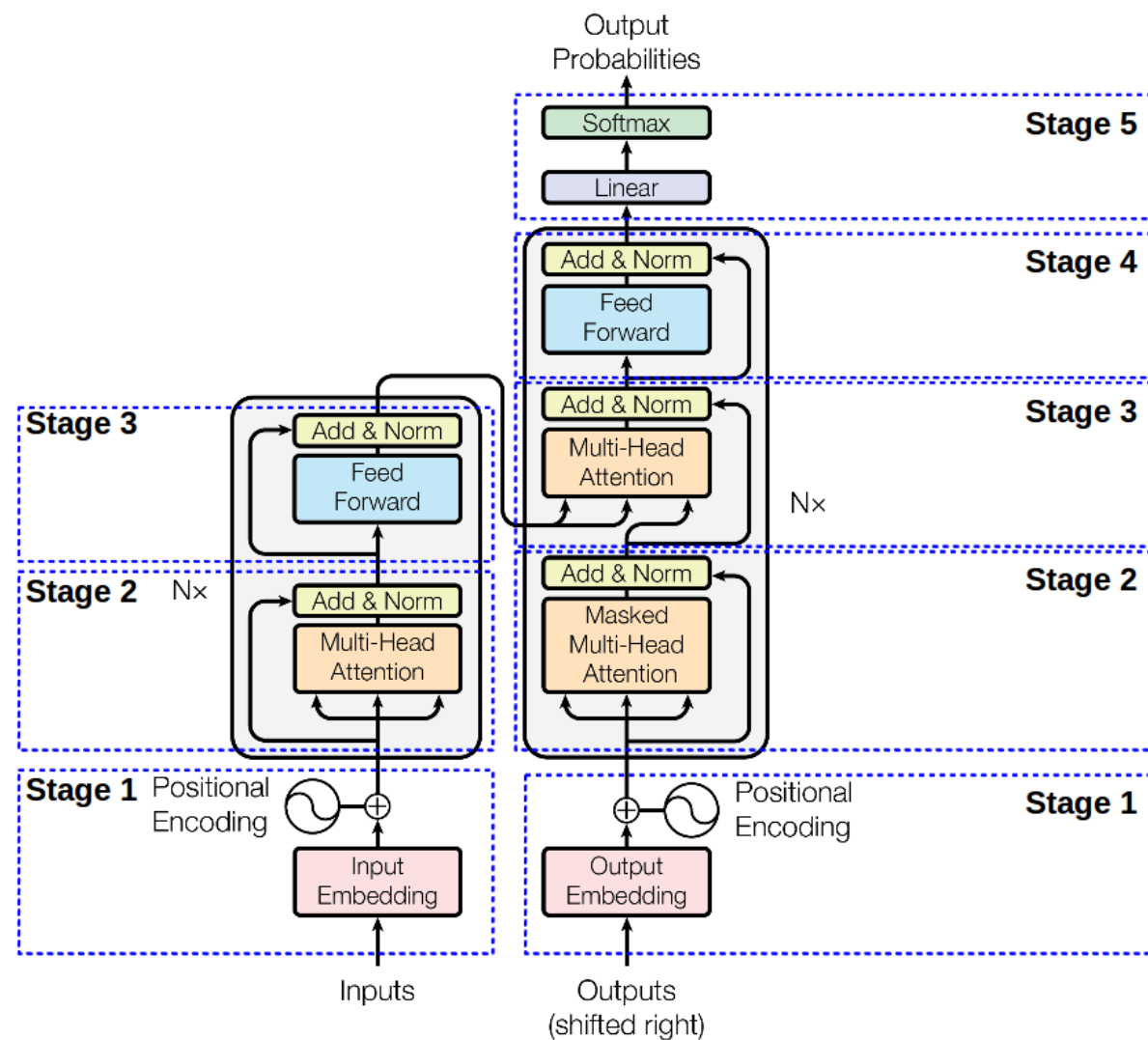


**Transformer**

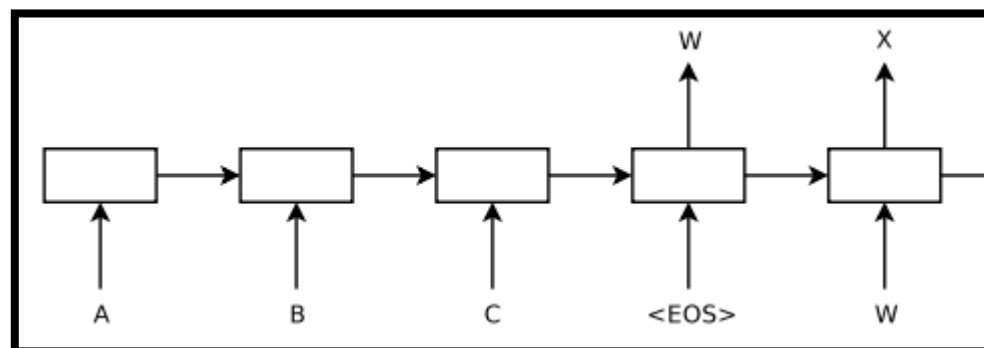
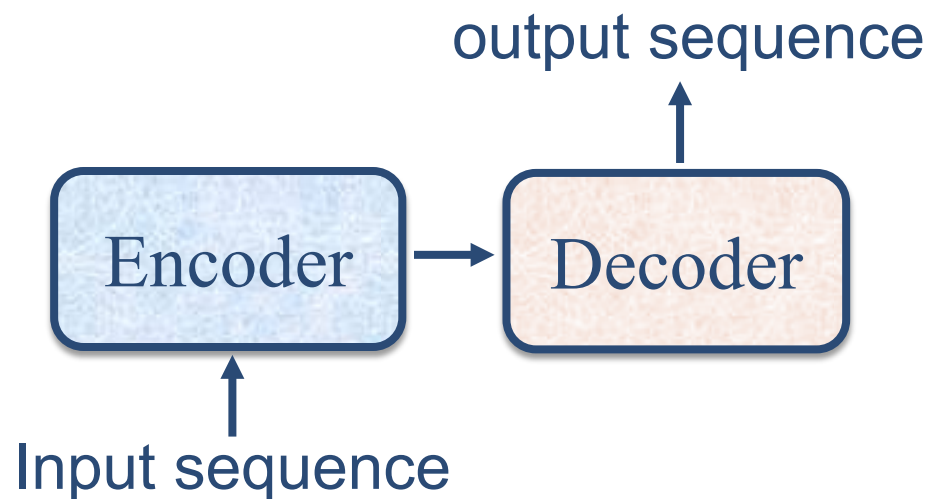


# Transformer

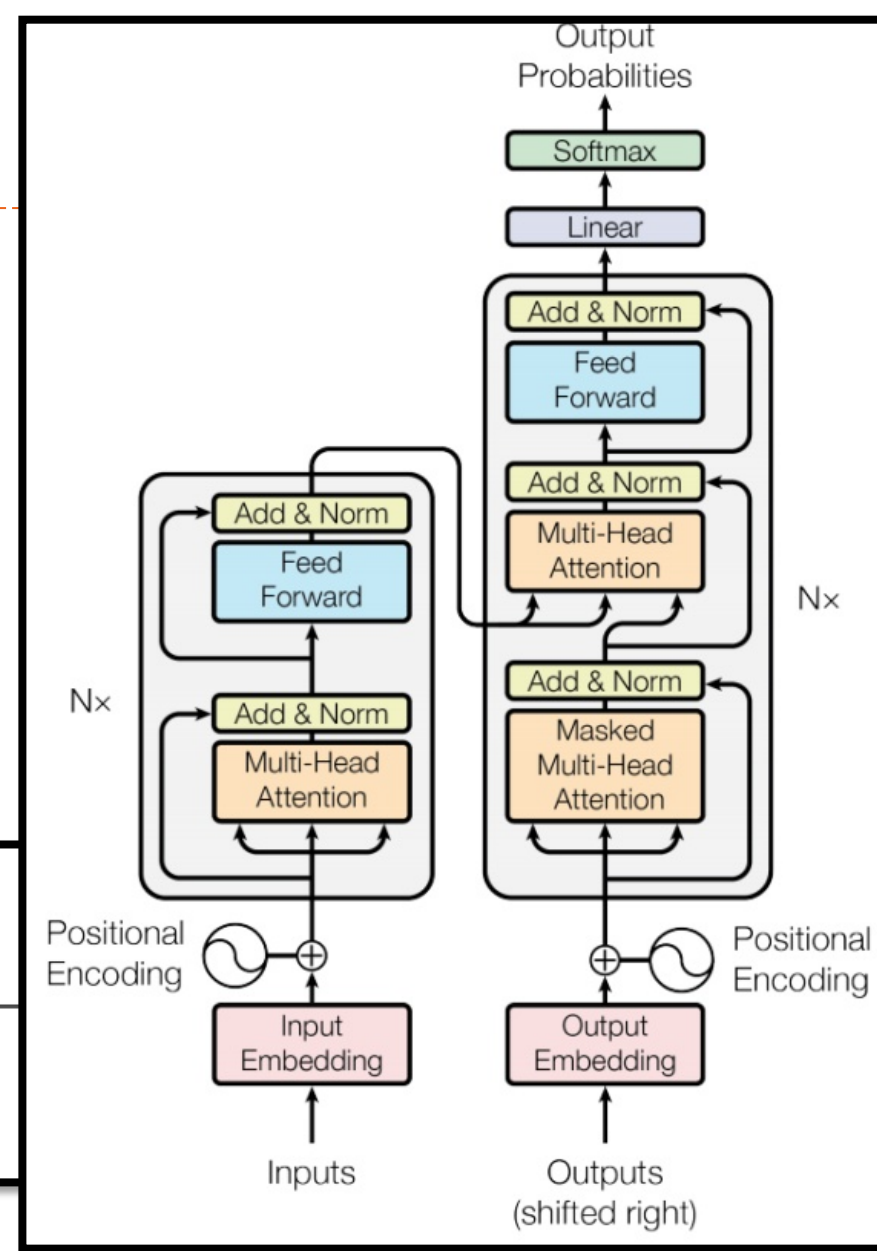
<https://mchromiak.github.io/articles/2017/Sep/12/Transformer-Attention-is-all-you-need/#.W90QB5Mzabg>



# Seq2Seq



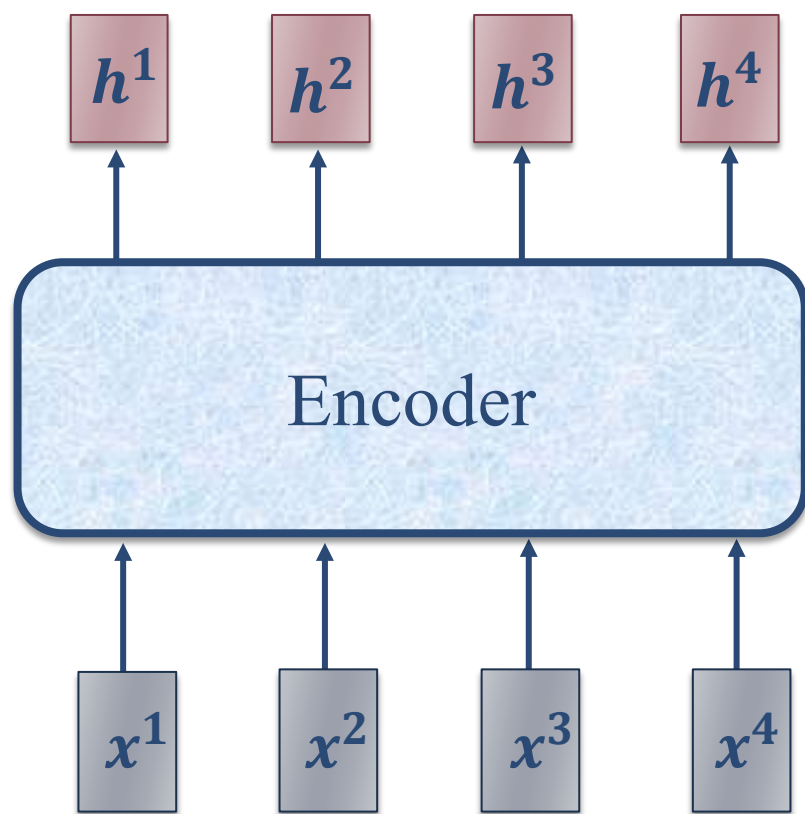
Sequence to Sequence Learning with Neural Networks



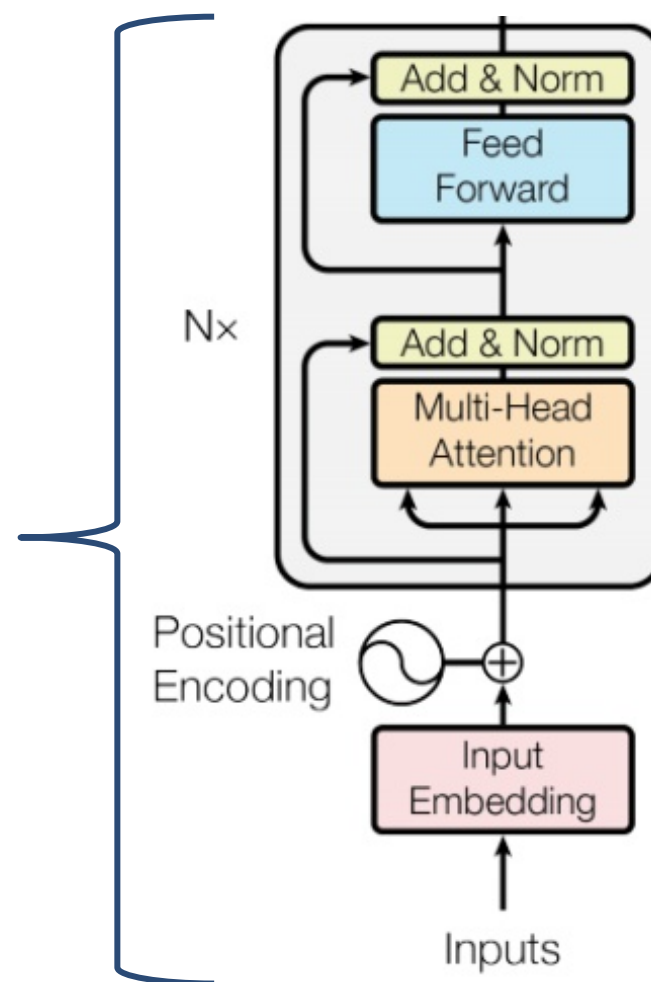
Transformer

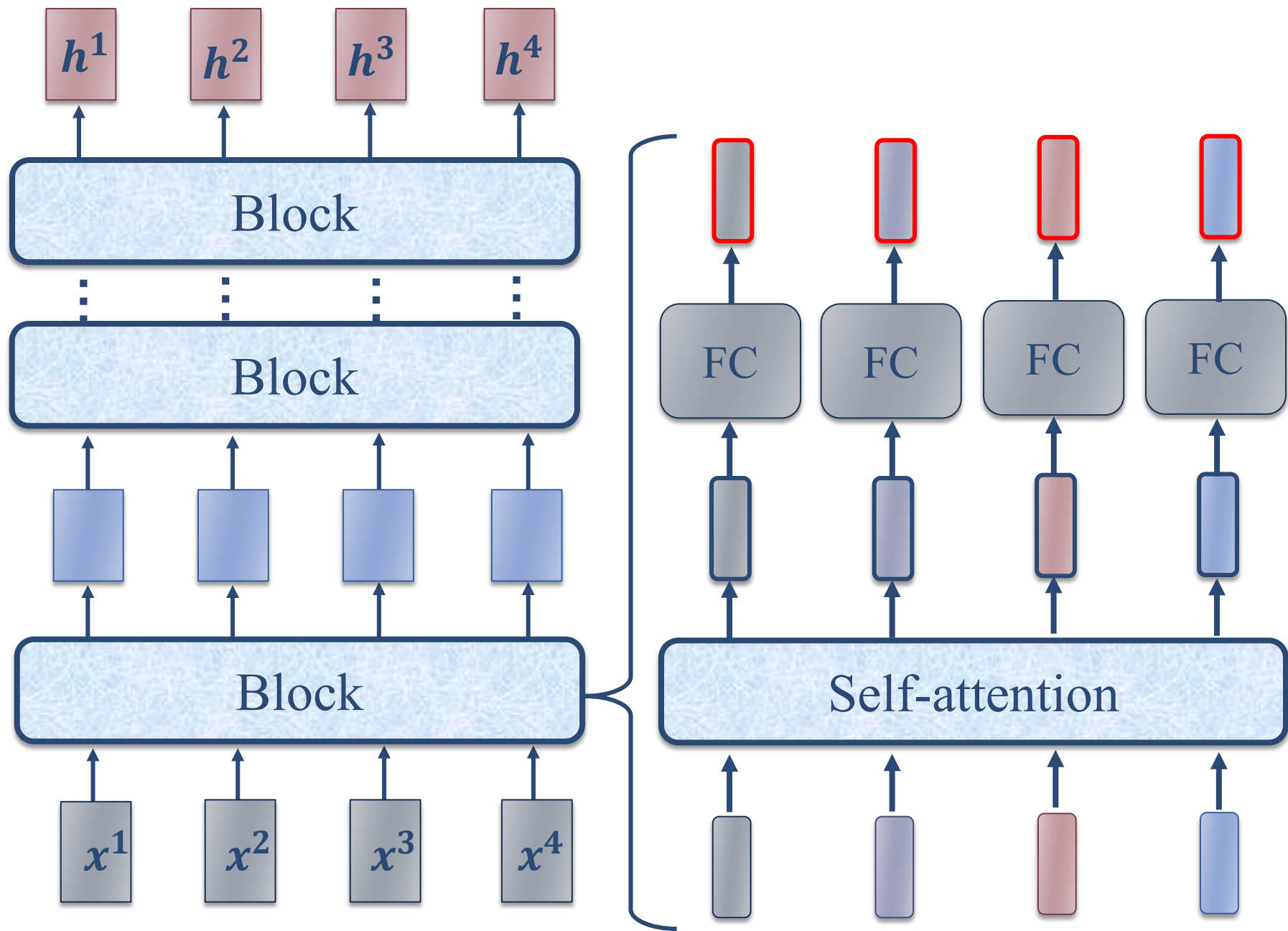
# Transformer's Encoder

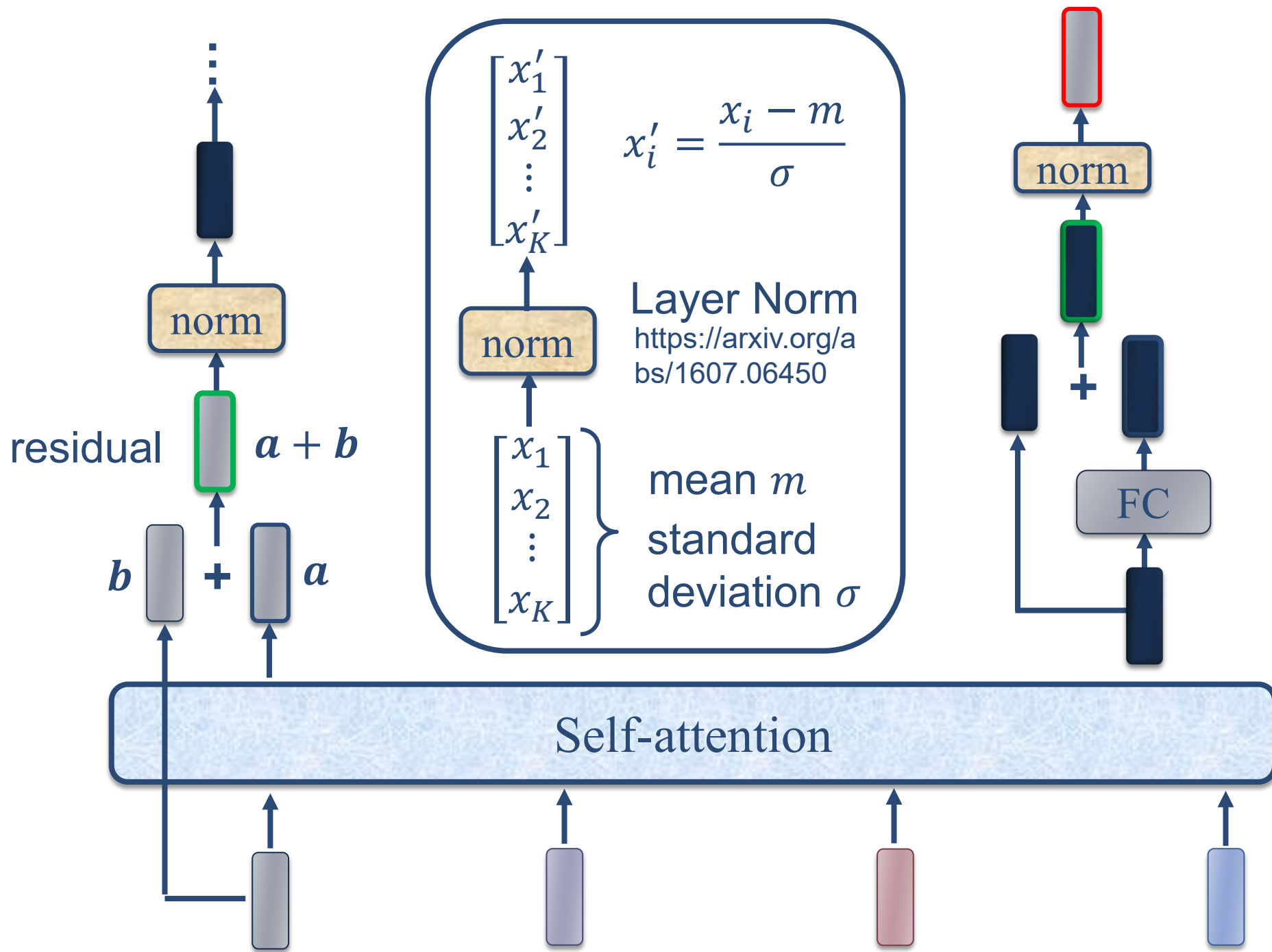
RNN, CNN, or ... ?



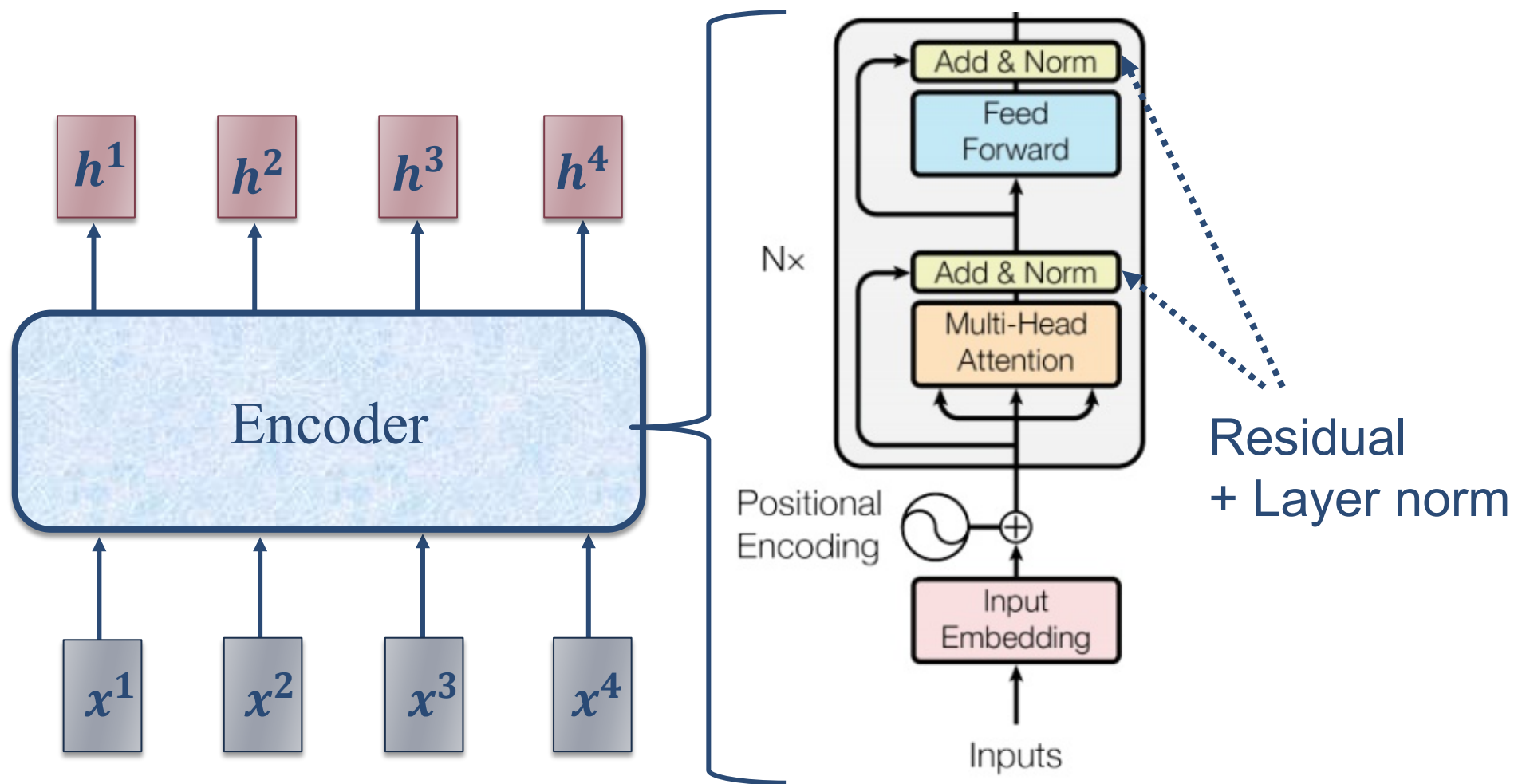
## Transformer's Encoder



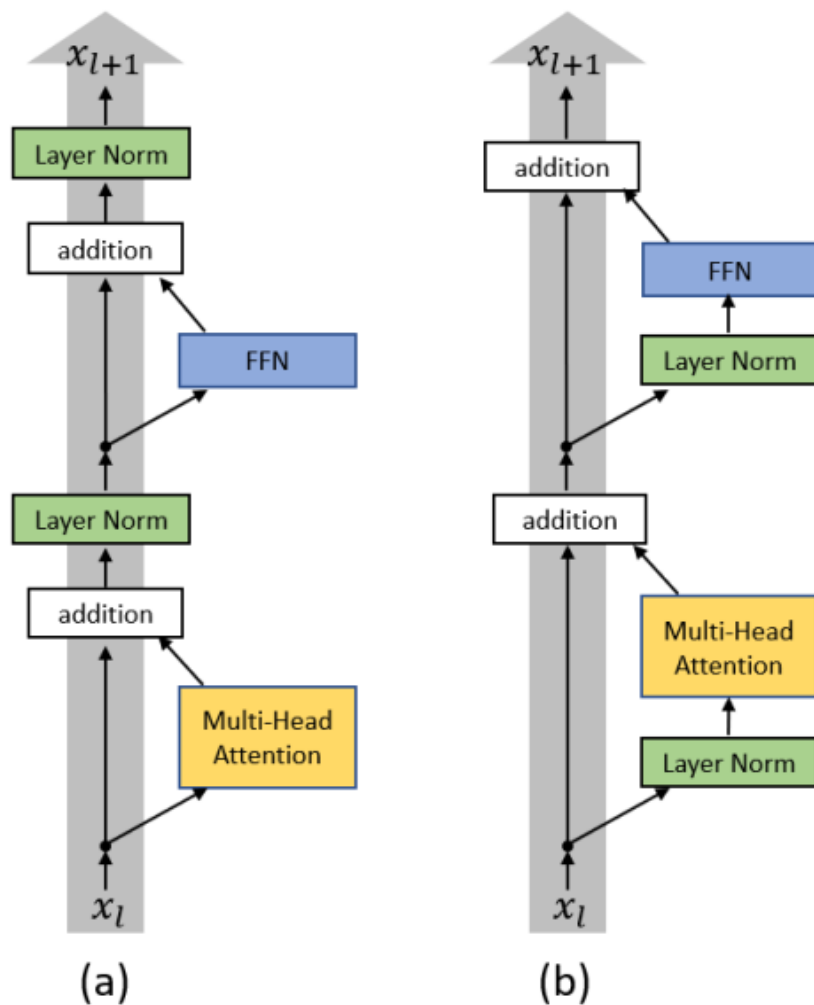




# Transformer's Encoder



# Transformer's Encoder





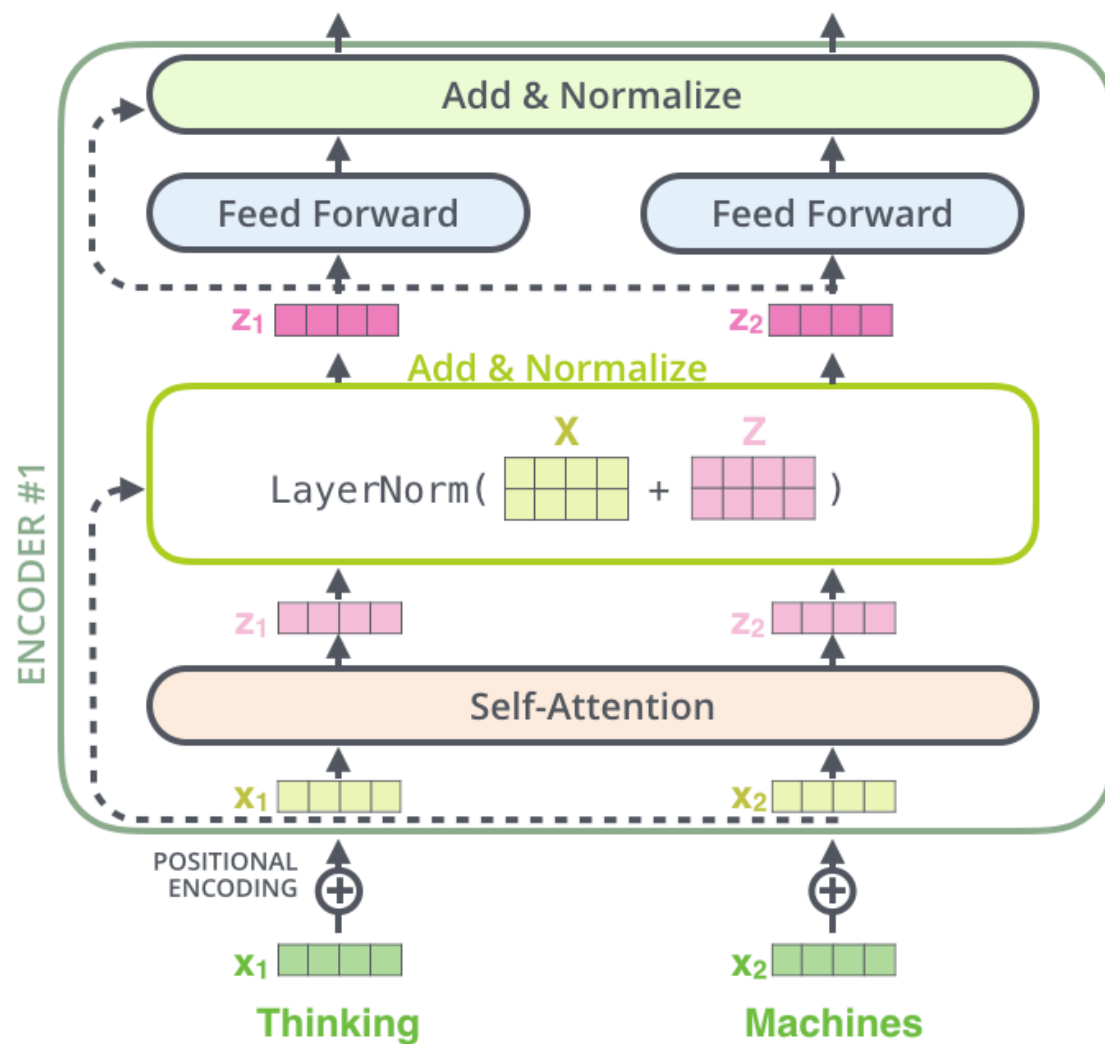
# Transformer's Encoder

图片来源: <http://jalammr.github.io/illustrated-transformer/>

## ► 基于自注意力的模型

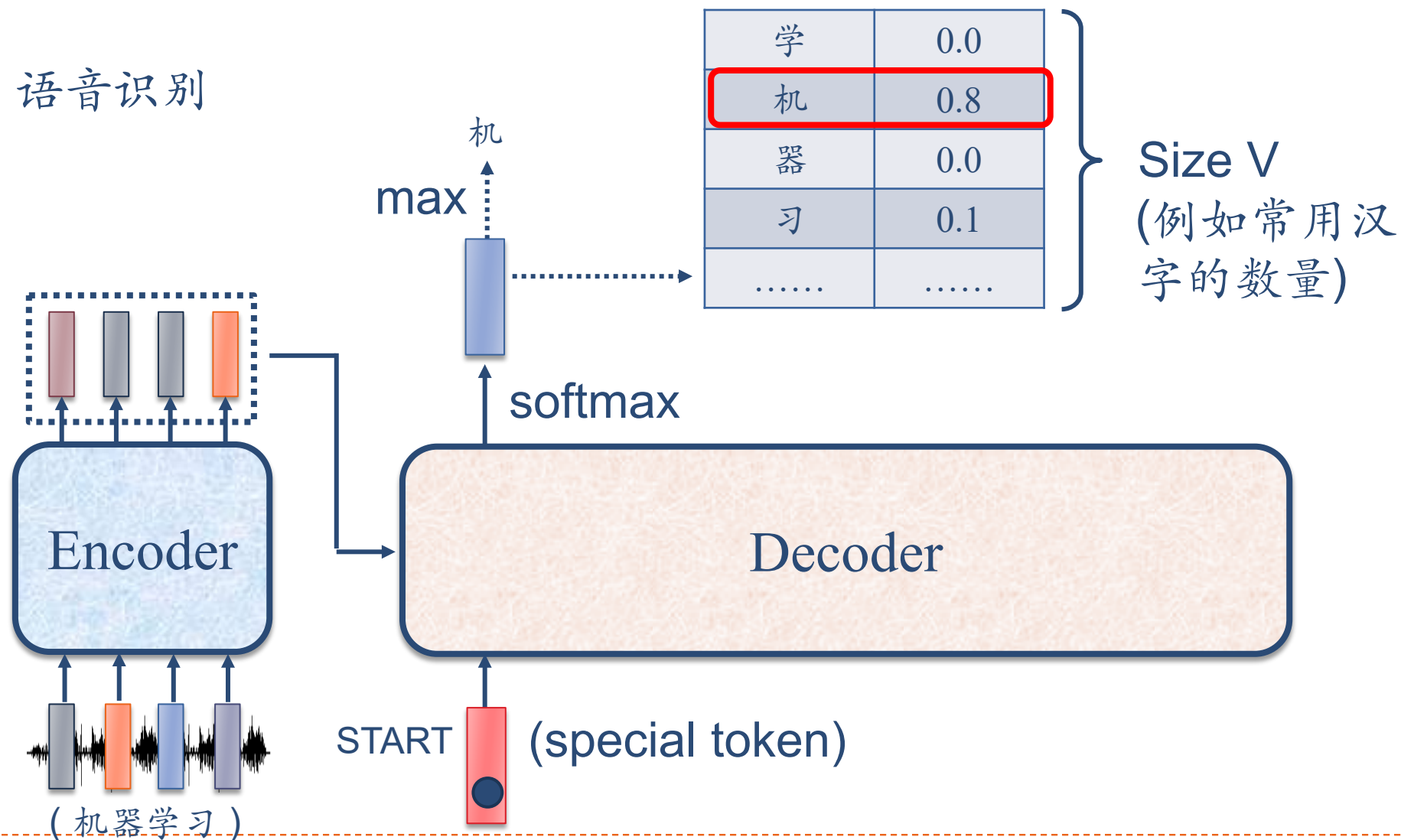
## ► 引入其它操作

- 位置编码
- 层归一化
- 直连边
- 逐位的FNN

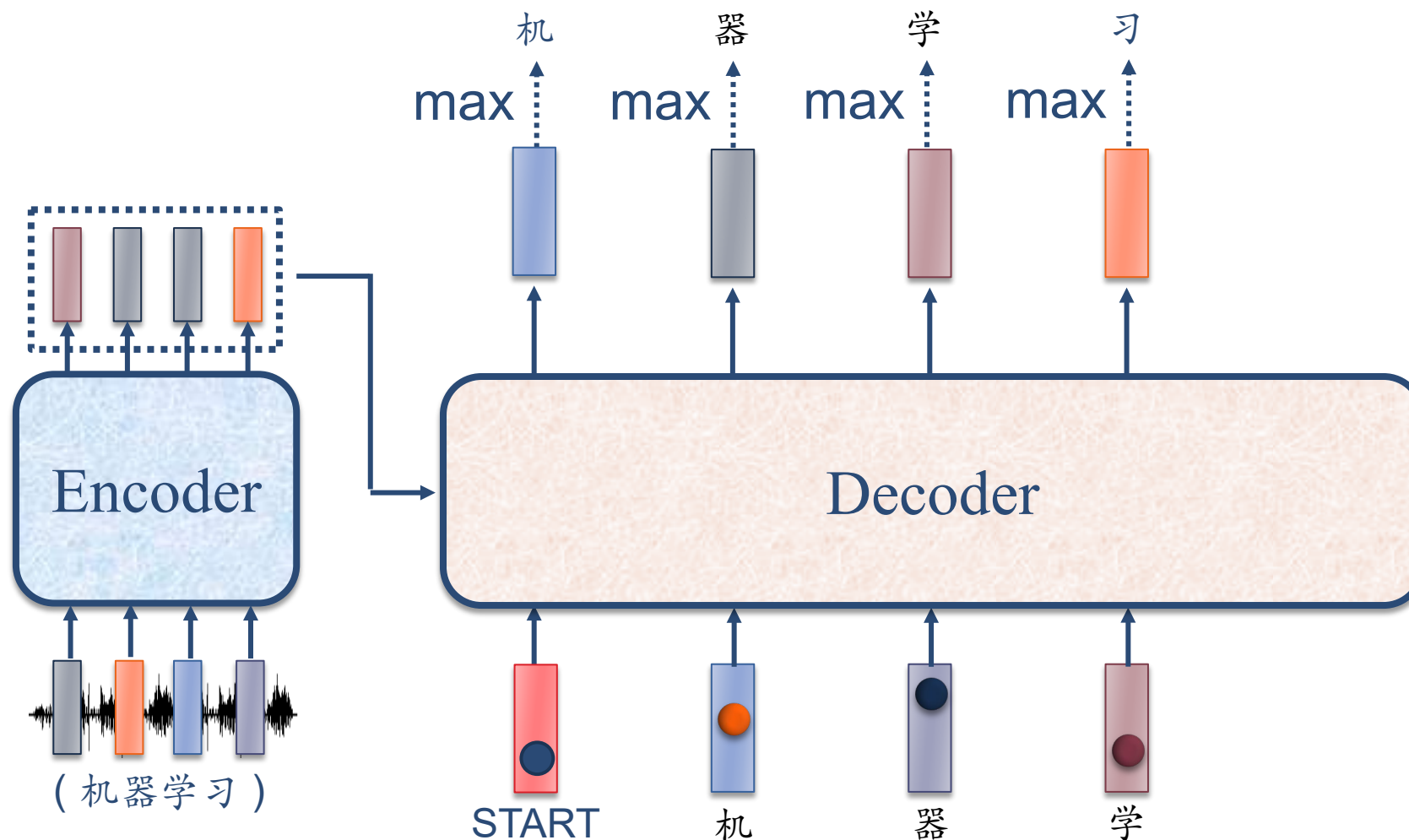


# Transformer's Decoder——Autoregressive (AT)

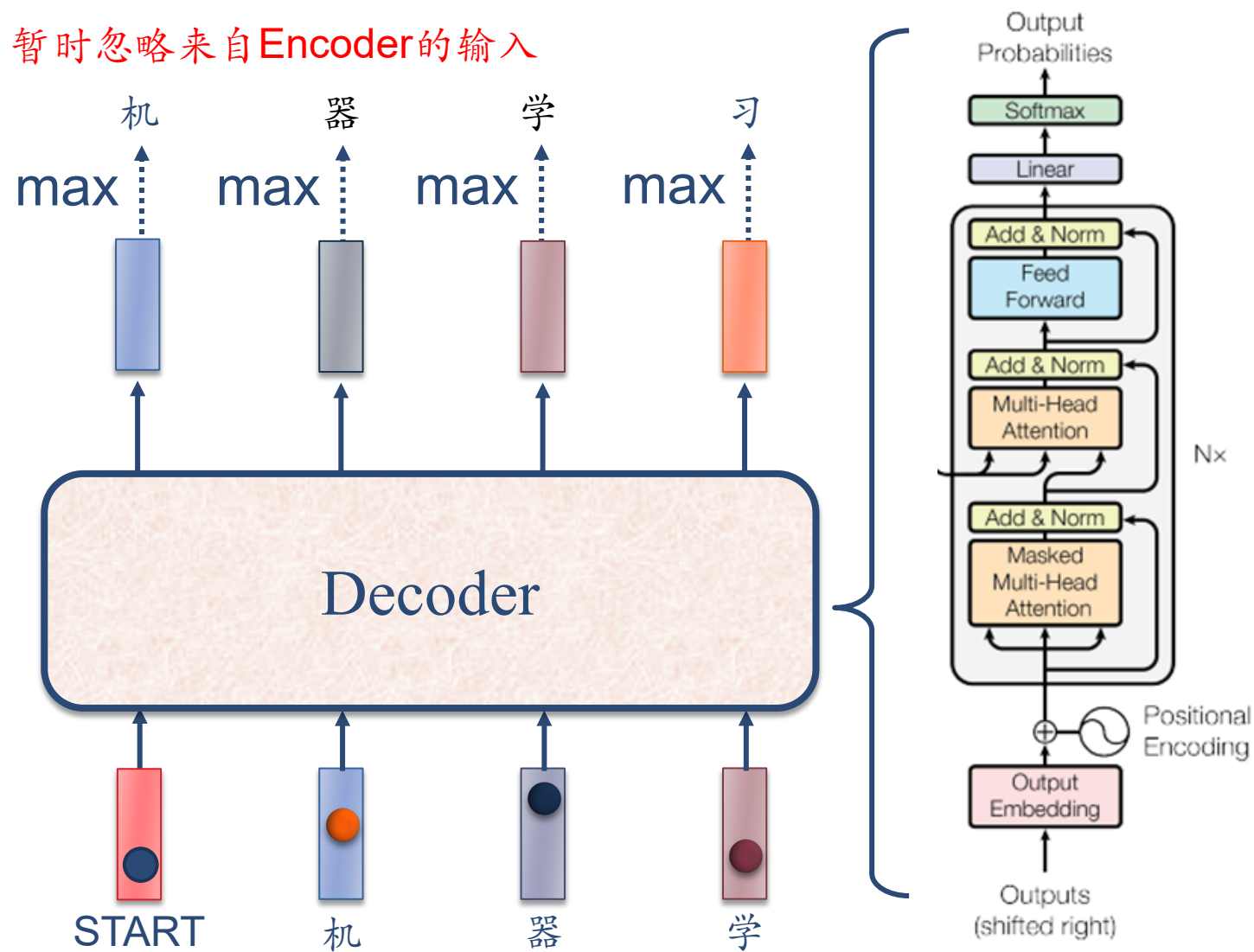
案例：语音识别



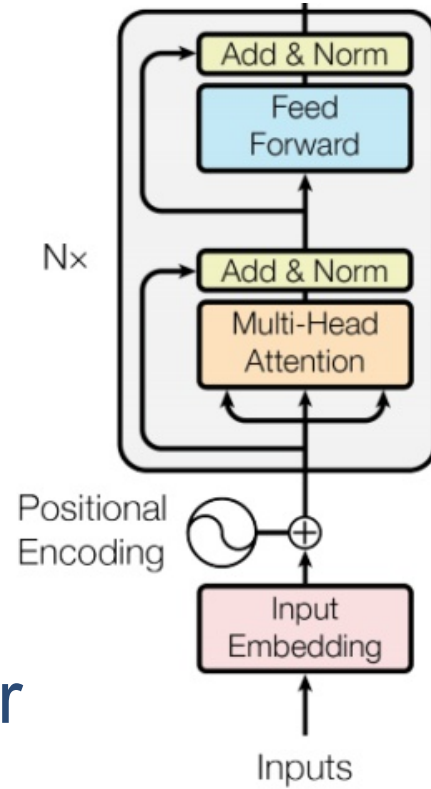
# Transformer's Decoder——Autoregressive (AT)



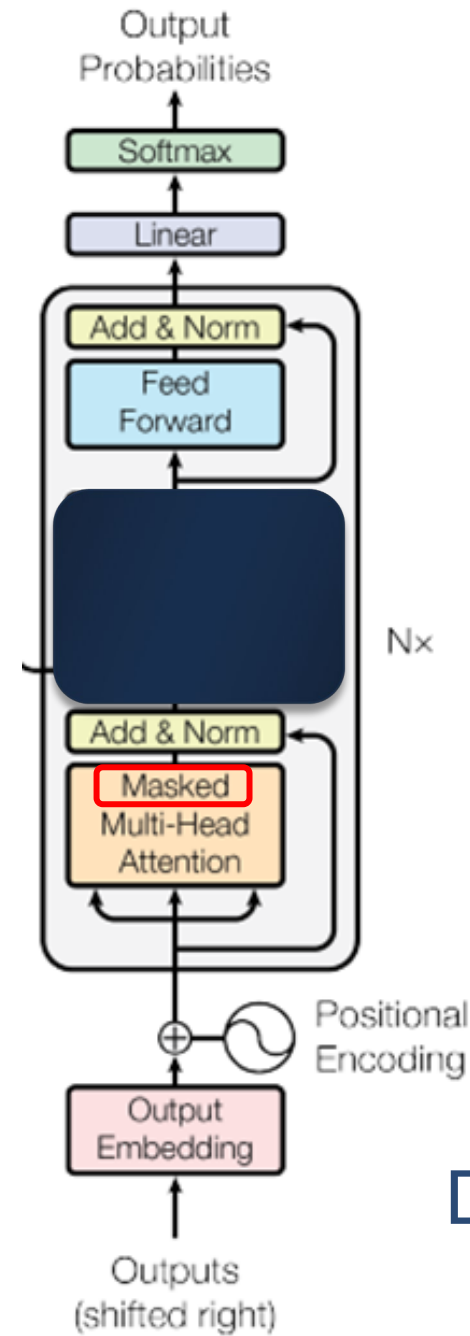
# Transformer's Decoder—Autoregressive (AT)



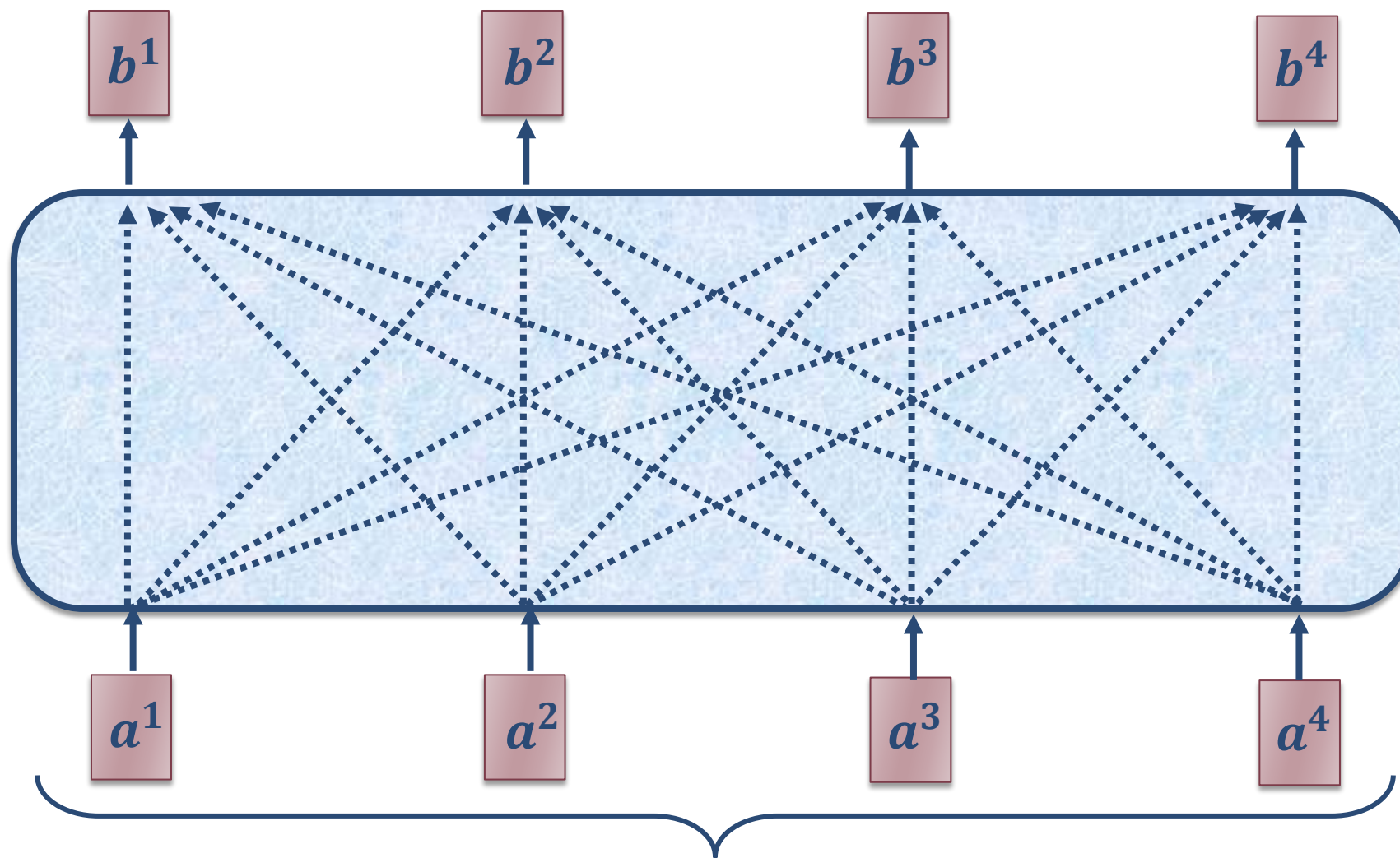
Encoder



Decoder

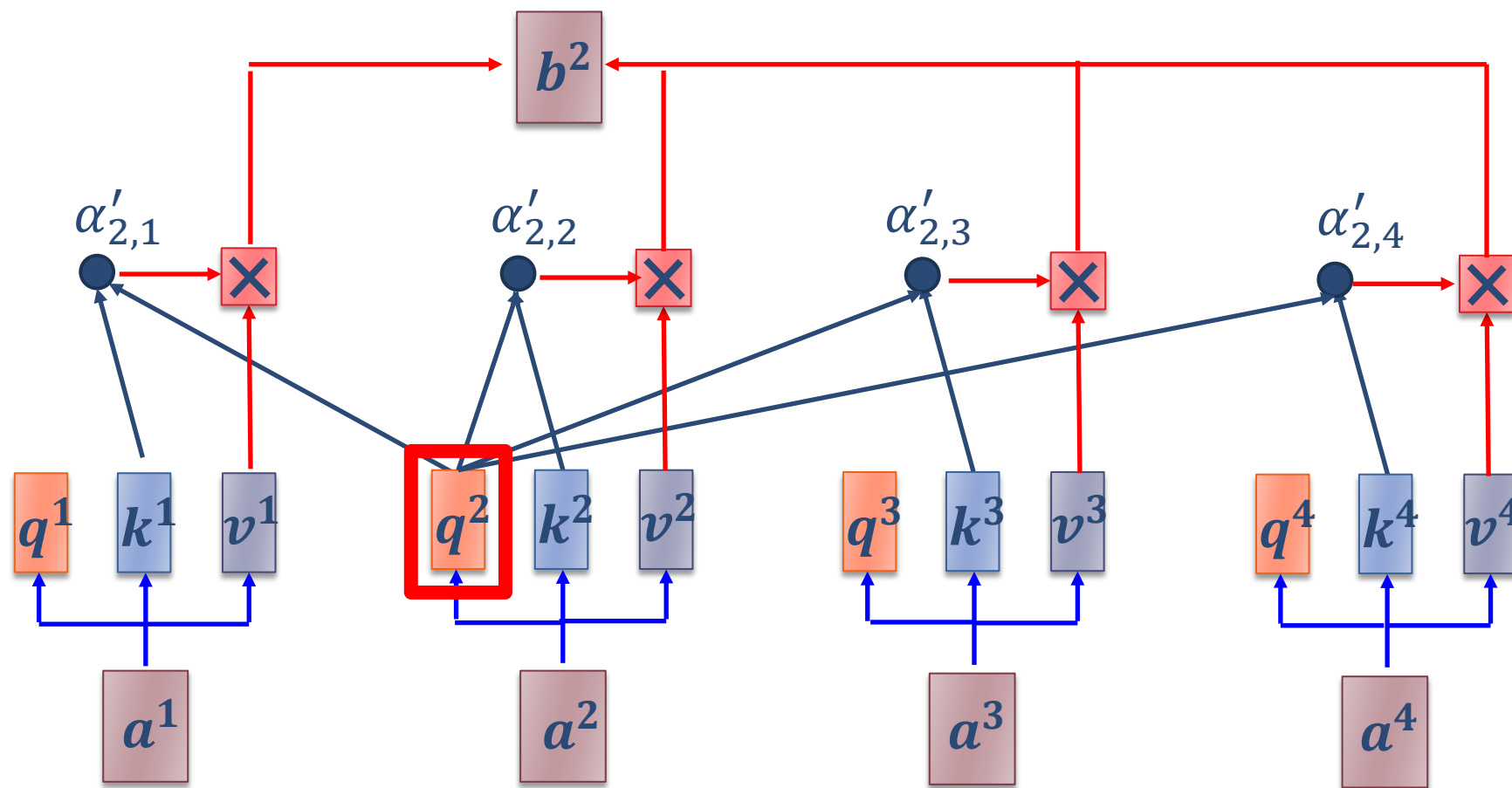


# Self-attention → Masked Self-attention

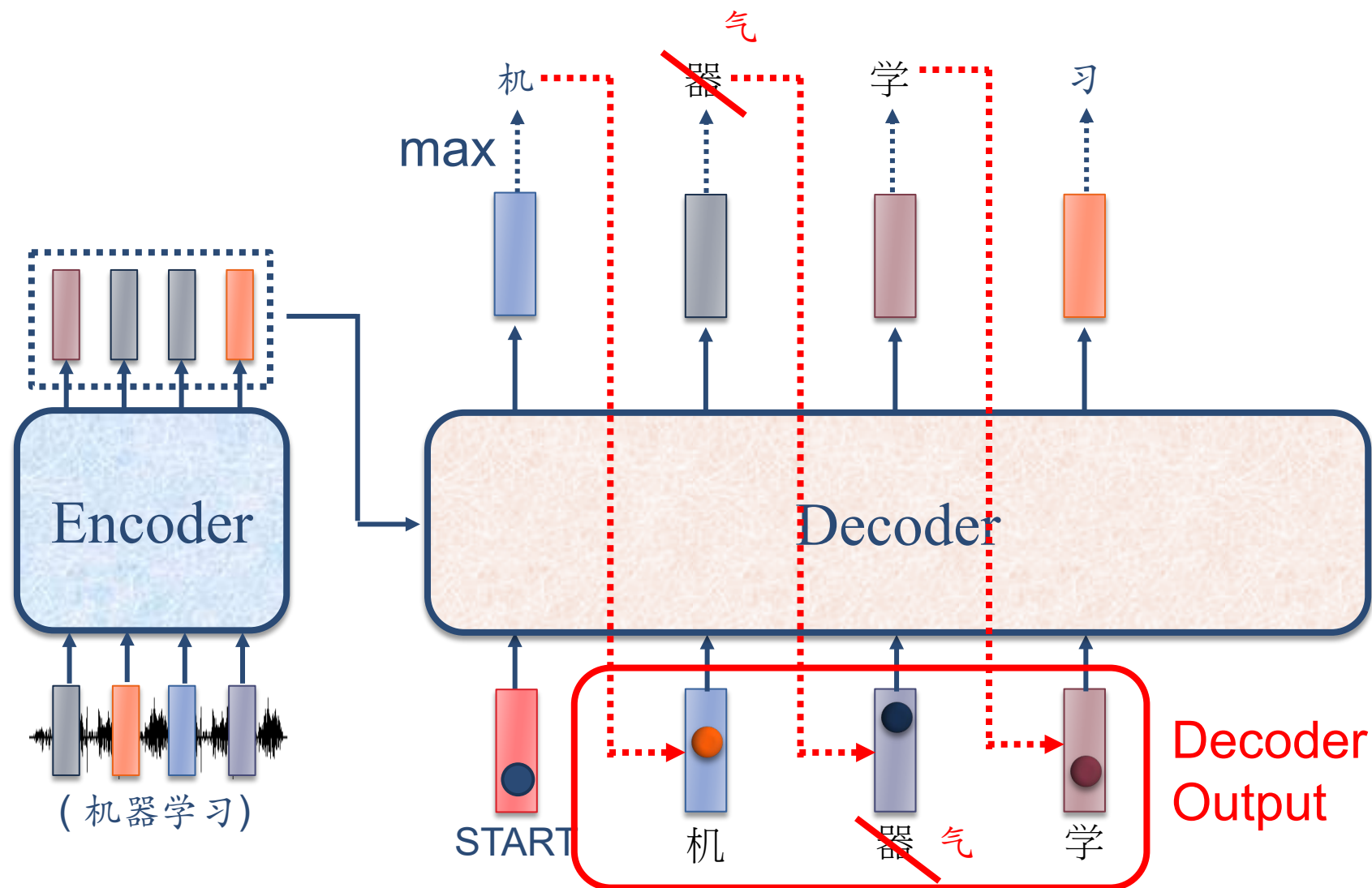


原始输入或来自隐层的输入

# Self-attention → Masked Self-attention



# Transformer's Decoder——Autoregressive (AT)

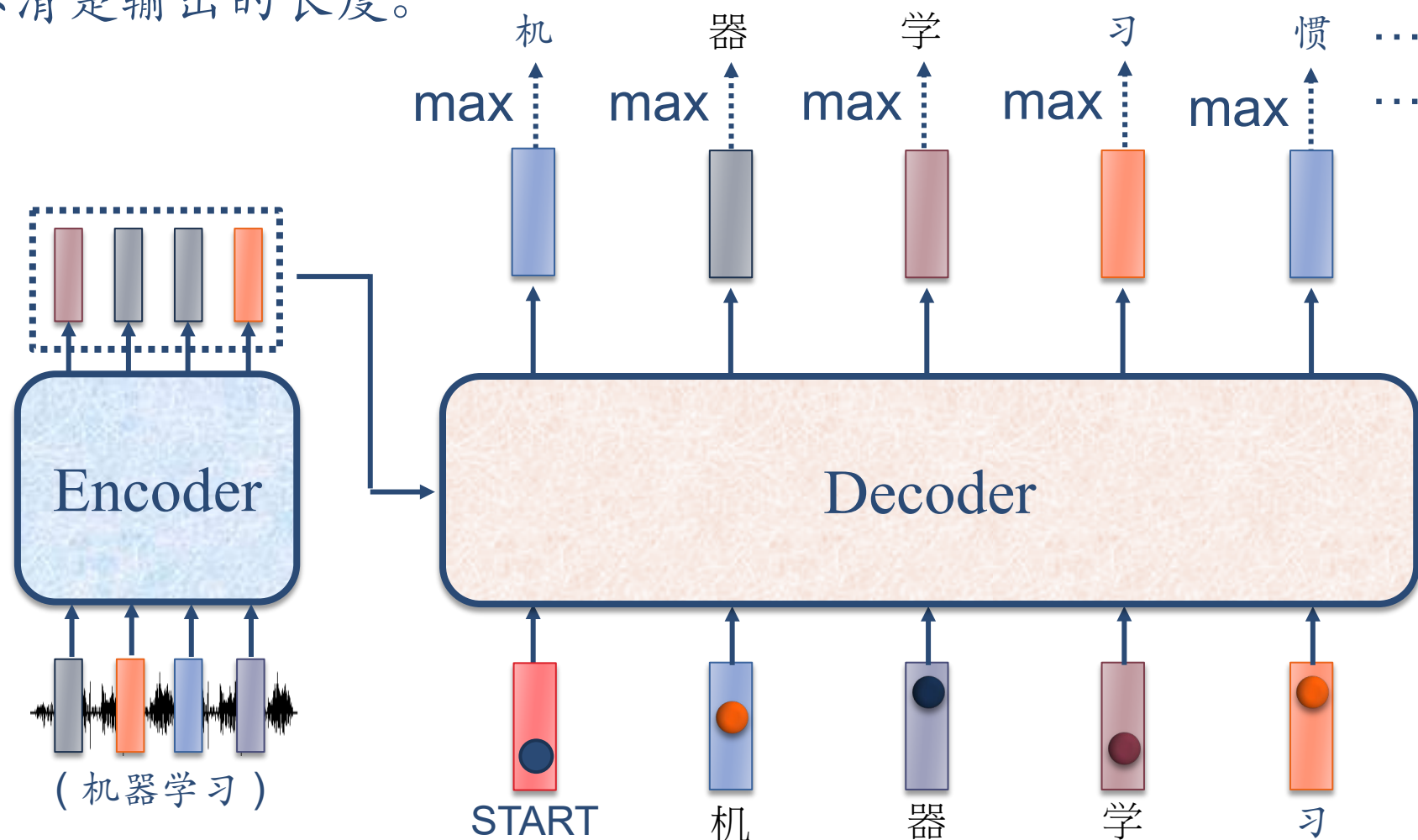




# Transformer's Decoder——Autoregressive (AT)

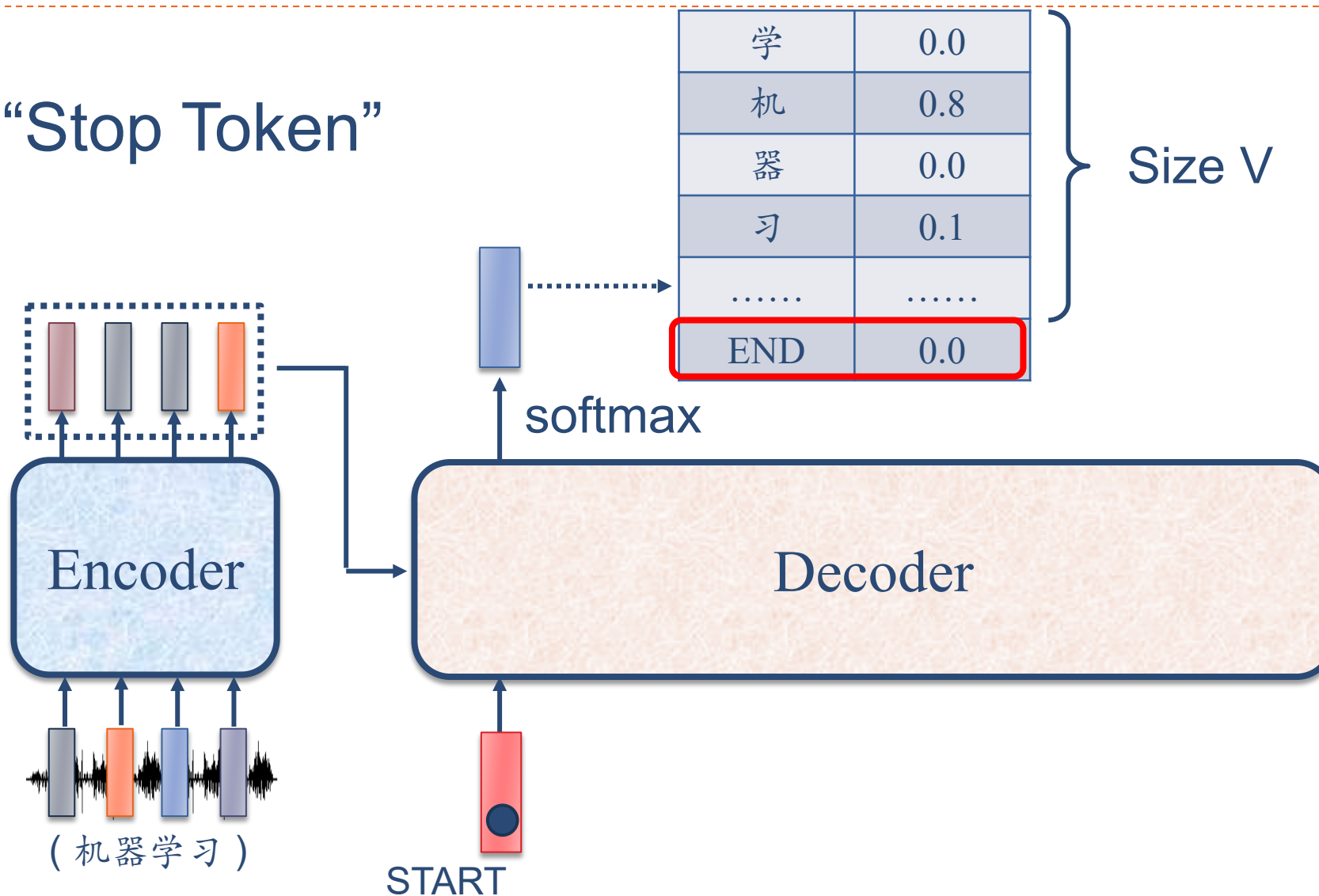
我们并不清楚输出的长度。

Never stop!

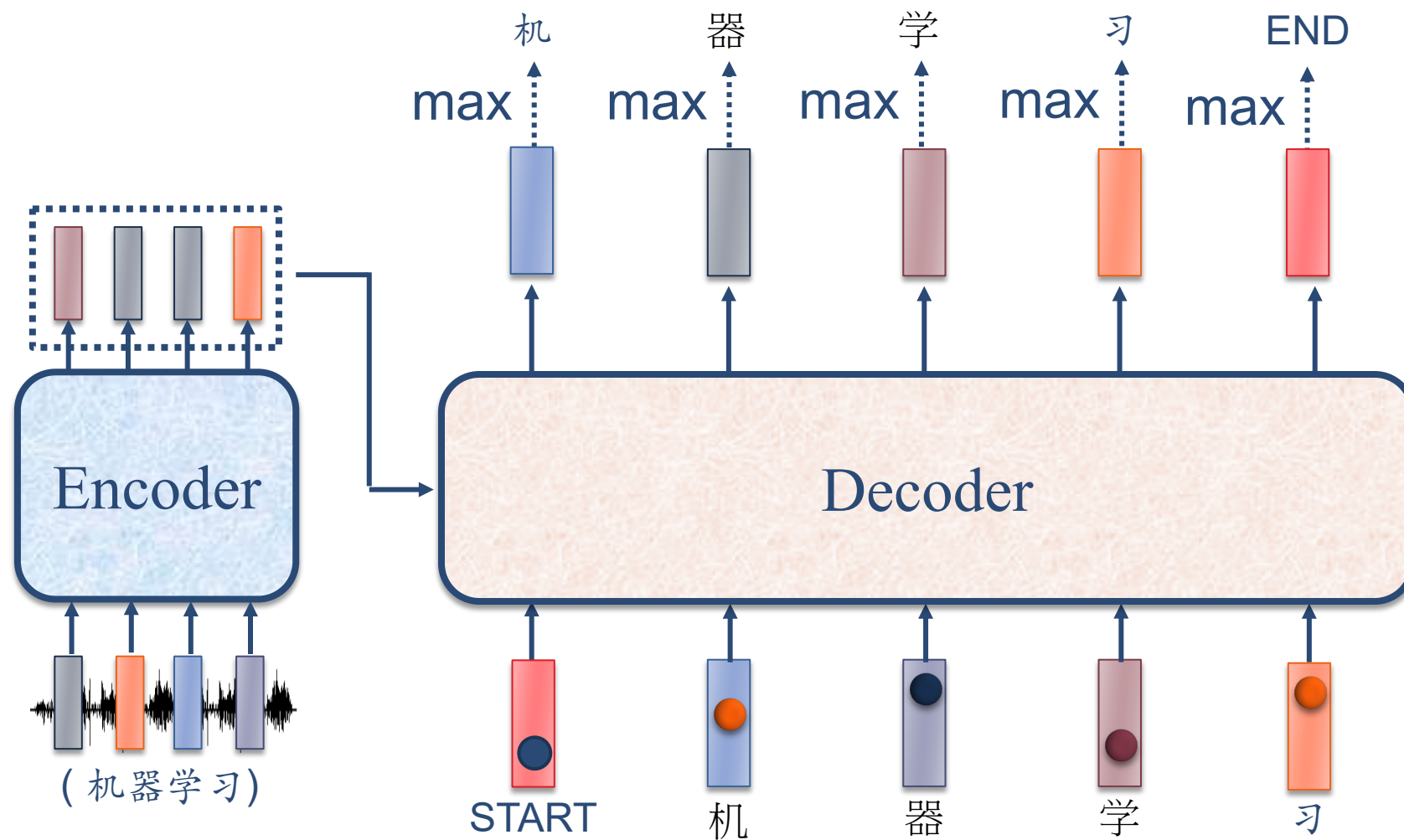


# Transformer's Decoder——Autoregressive (AT)

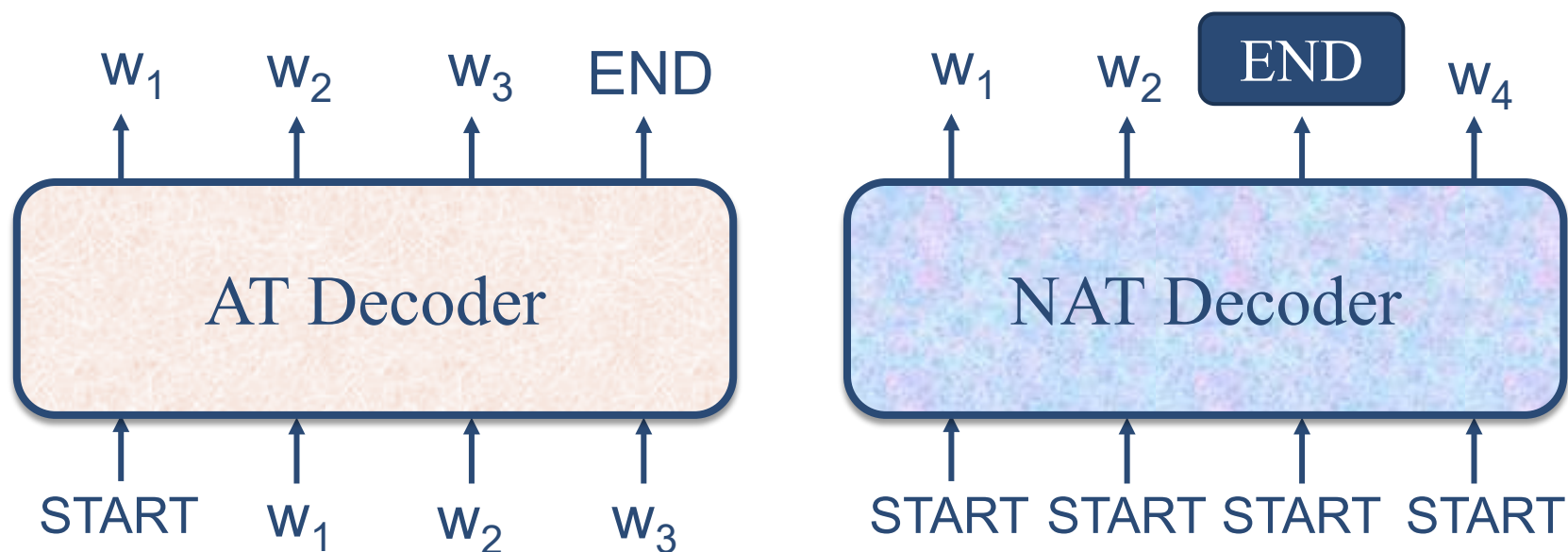
添加“Stop Token”



# Transformer's Decoder——Autoregressive (AT)

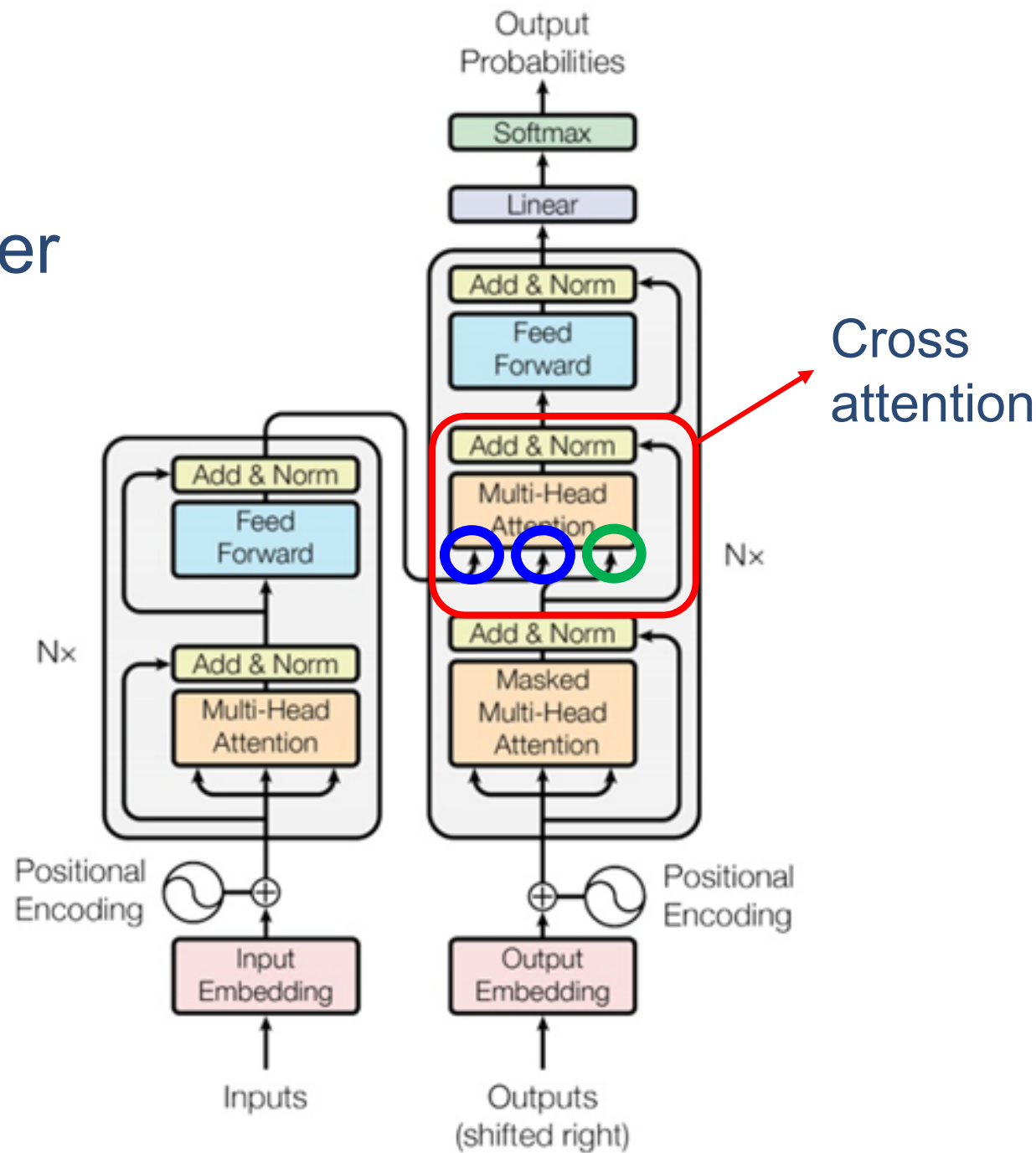


# Transformer's Decoder——Non-Autoregressive (NAT)

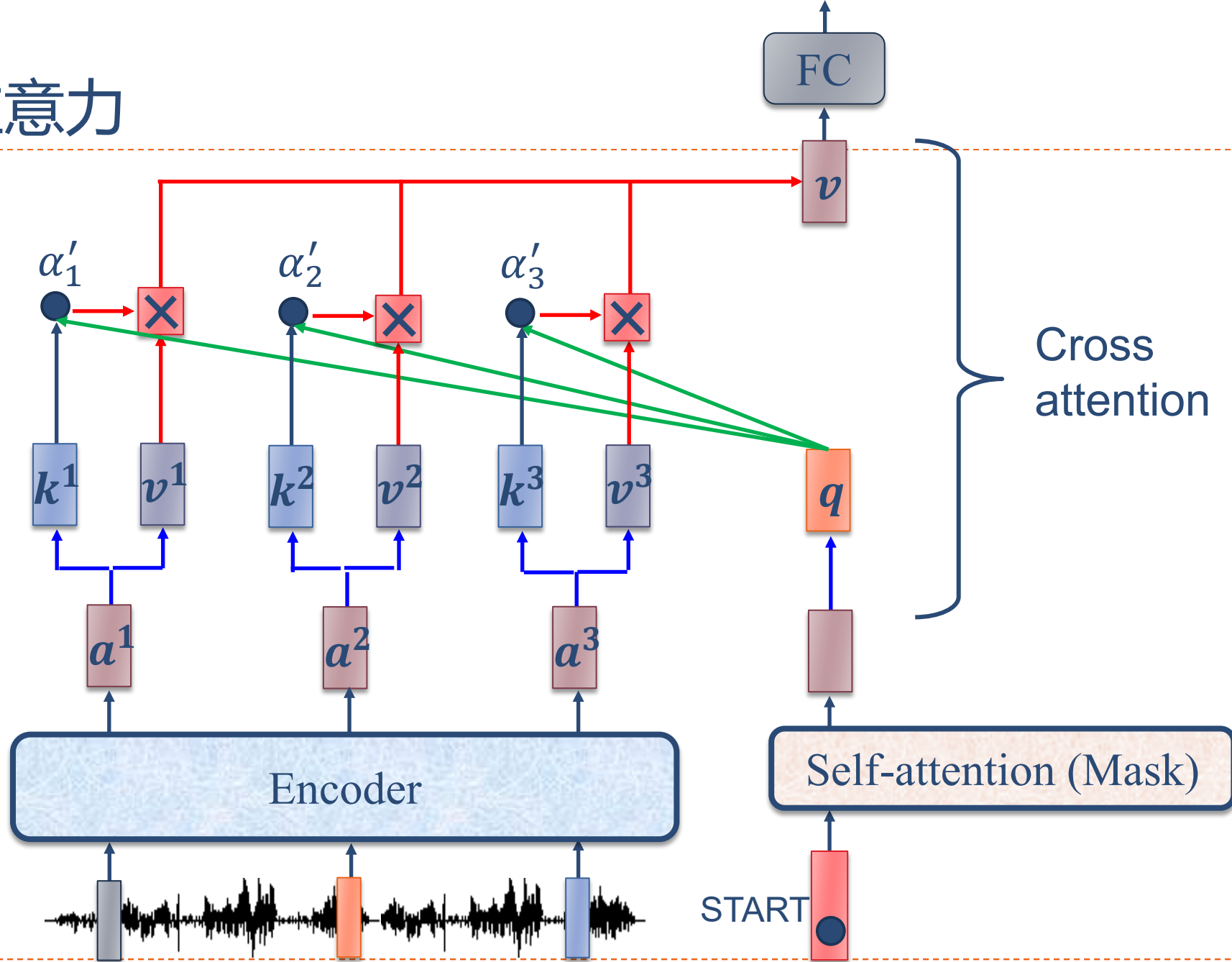


- NAT解码器一次生成整个序列；
- NAT怎样决定输出的长度？用另一个模型预测输出长度；或截取END之前的内容。
- NAT的优势：可并行（而AT只能逐个Token生成）；能够很好地控制输出长度；
- NAT的劣势：通常性能比AT差。

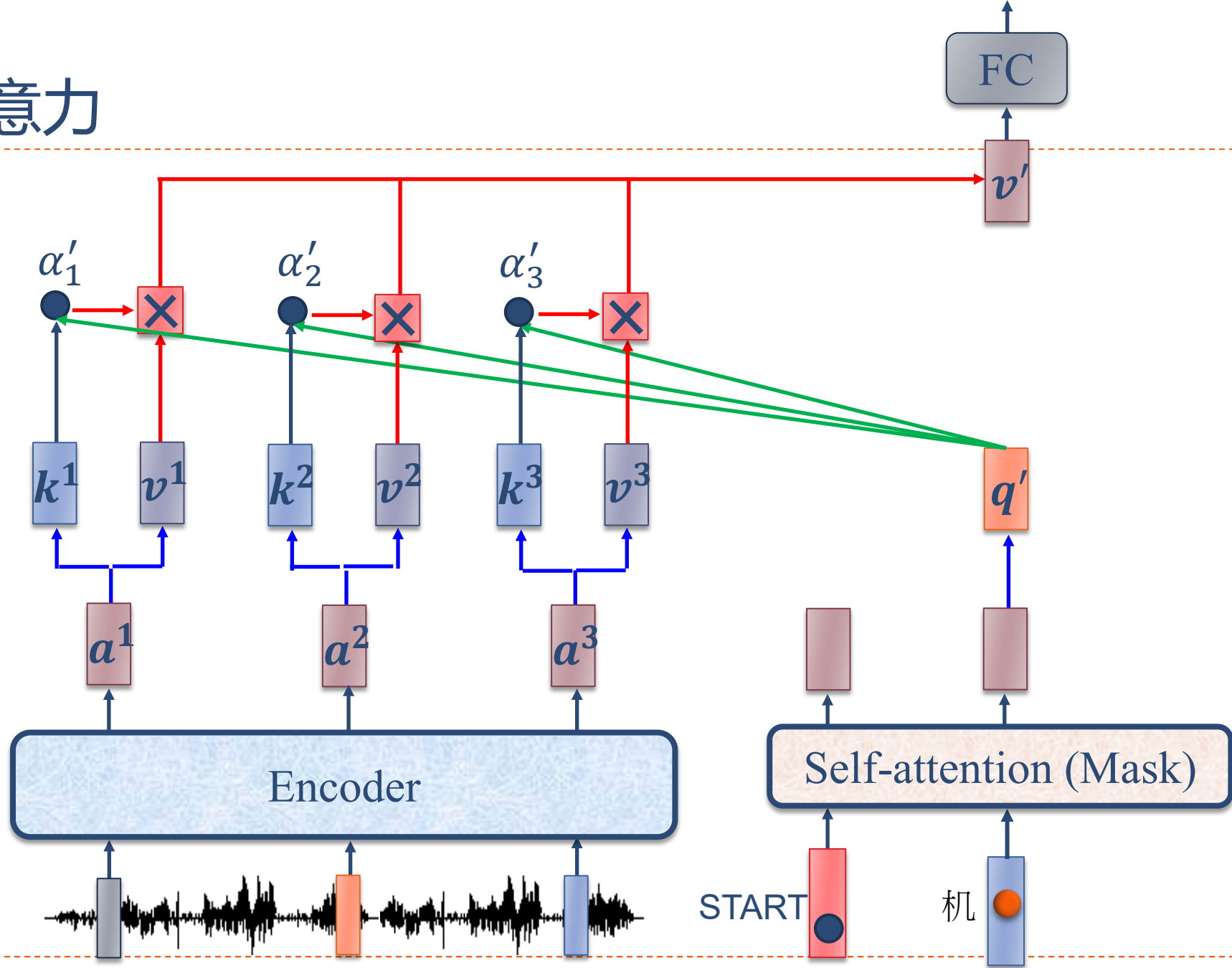
# Encoder-Decoder



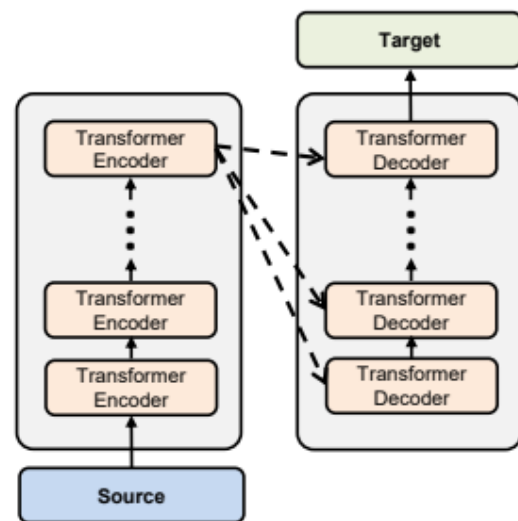
# 交叉注意力



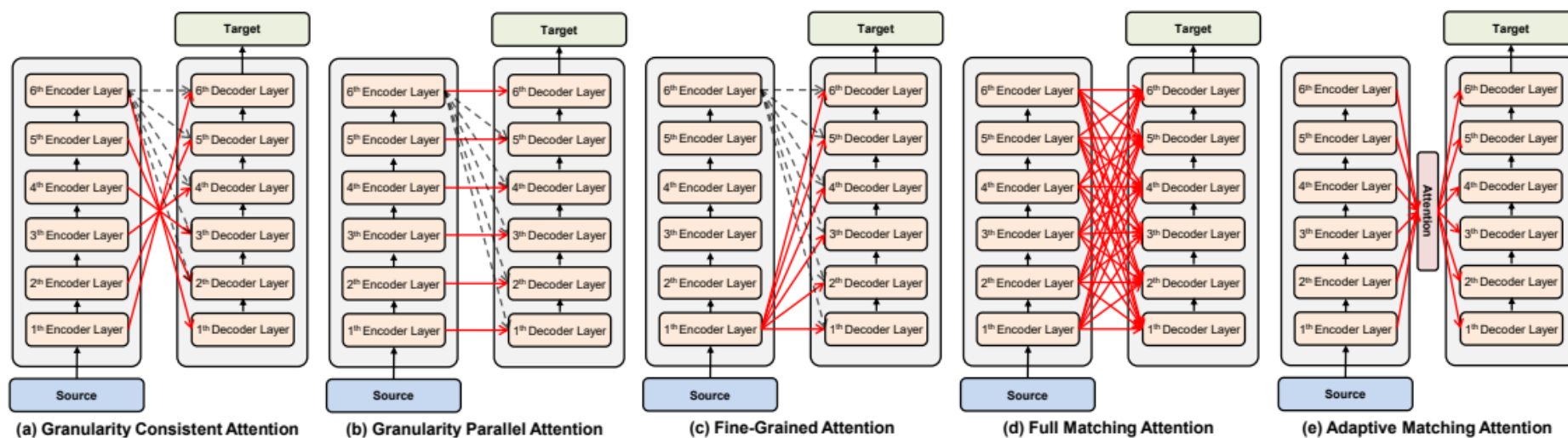
# 交叉注意力



# 交叉注意力



(a) Conventional Transformer



(a) Granularity Consistent Attention

(b) Granularity Parallel Attention

(c) Fine-Grained Attention

(d) Full Matching Attention

(e) Adaptive Matching Attention



# Transformer的训练

**Teacher Forcing:** 用真实标签作为Decoder输入。

