



大模型网络结构

魏明强、宫丽娜

计算机科学与技术学院

智周万物·道济天下

目录



- Transformer
 - 注意力机制
 - 编码器-解码器结构
 - 大模型中的编码器-解码器结构
- 编码器结构----BERT家族
 - BERT结构
 - 预训练策略
 - BERT的变体
- 解码器结构----GPT家族
 - GPT结构
 - 自回归预训练
 - 后续改进
- 思考

目录



□ Transformer

- 注意力机制
- 编码器-解码器结构
- 大模型中的编码器-解码器结构

□ 编码器结构----BERT家族

- BERT结构
- 预训练策略
- BERT的变体

□ 解码器结构----GPT家族

- GPT结构
- 自回归预训练
- 后续改进

□ 思考

Transformer



- 面对问题: 记录输入序列中的长期依赖关系
- Transformer 利用注意力机制完成对源语言序列和目标语言序列全局依赖的建模

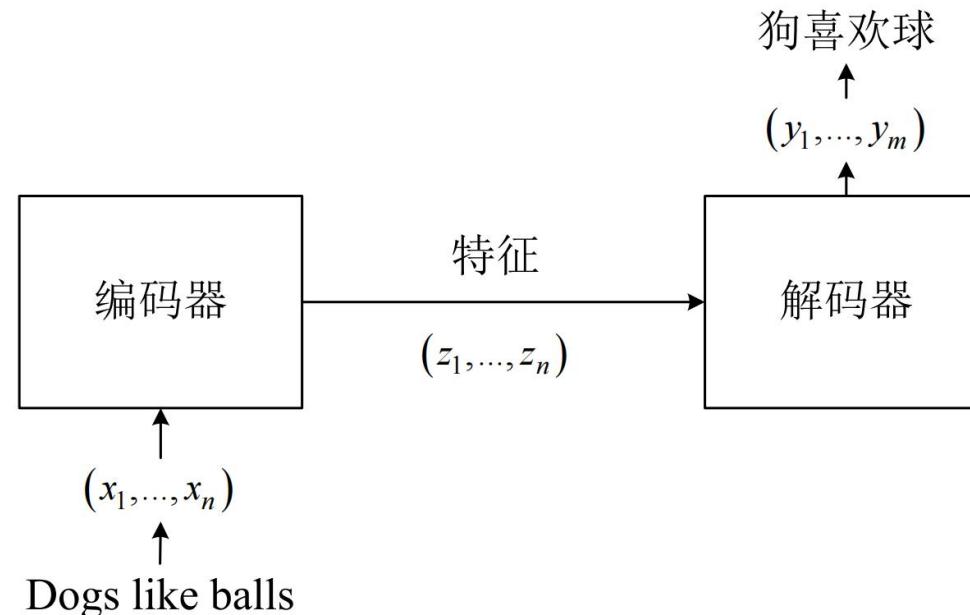


Transformer 是自然语言处理领域的颠覆者，为后续大模型网络结构（BERT、GPT）的发展奠定了基础

Transformer



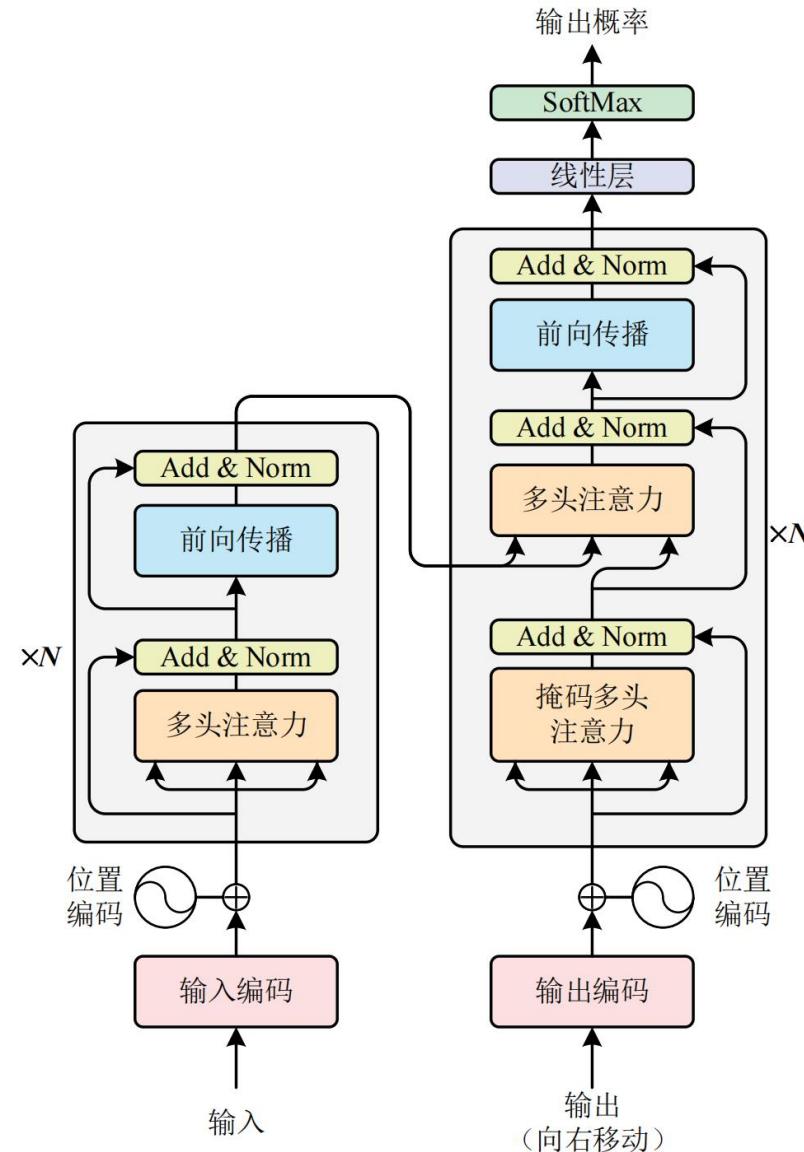
- 基本结构: 编码器-解码器结构
- 编码器输入, 解码器输出



Transformer



- 嵌入表示层
- 注意力层
- 位置前馈感知层
- 残差连接
- 层归一化



目录



□ Transformer

- 注意力机制
- 编码器-解码器结构
- 大模型中的编码器-解码器结构

□ 编码器结构----BERT家族

- BERT结构
- 预训练策略
- BERT的变体

□ 解码器结构----GPT家族

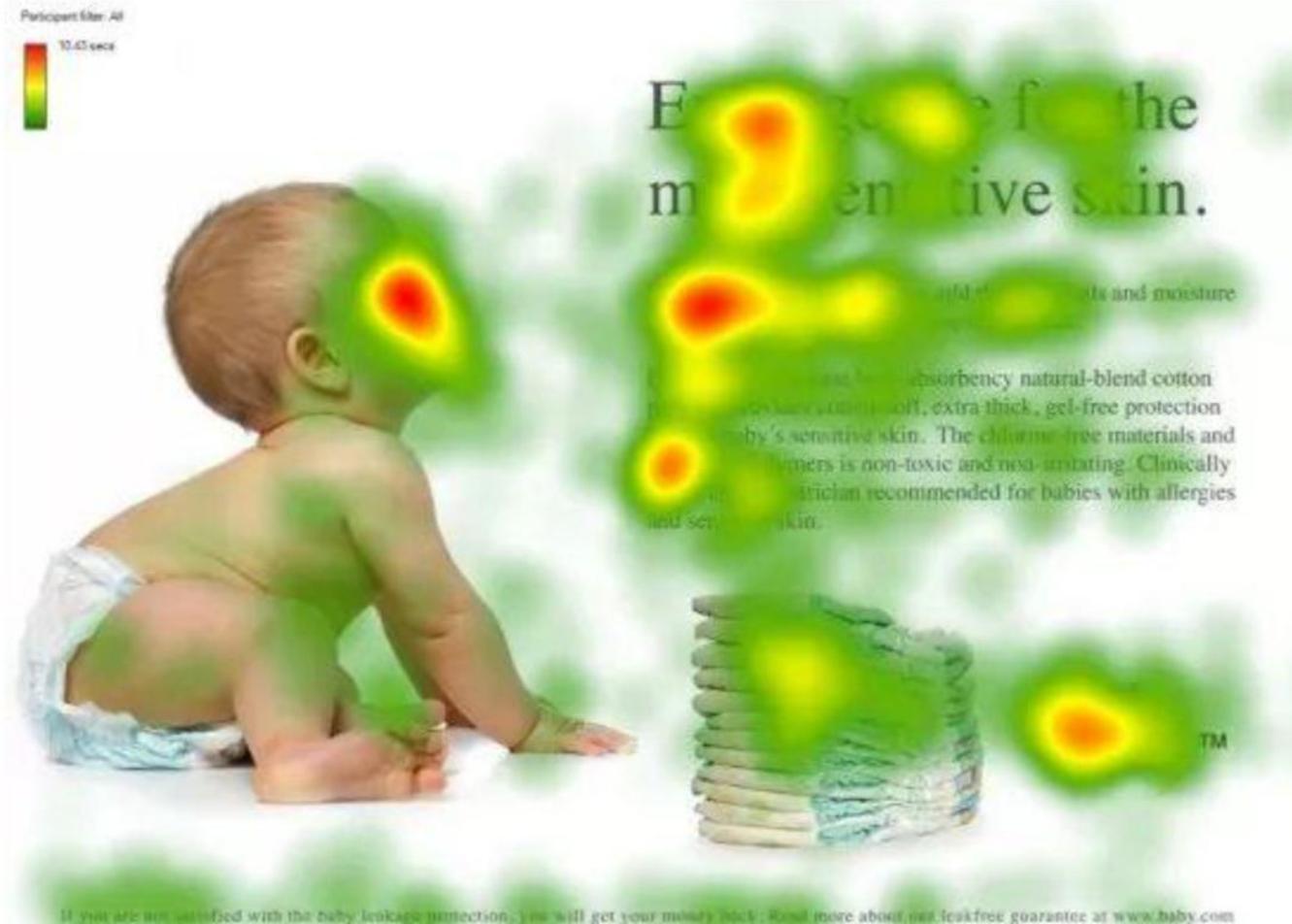
- GPT结构
- 自回归预训练
- 后续改进

□ 思考

注意力机制

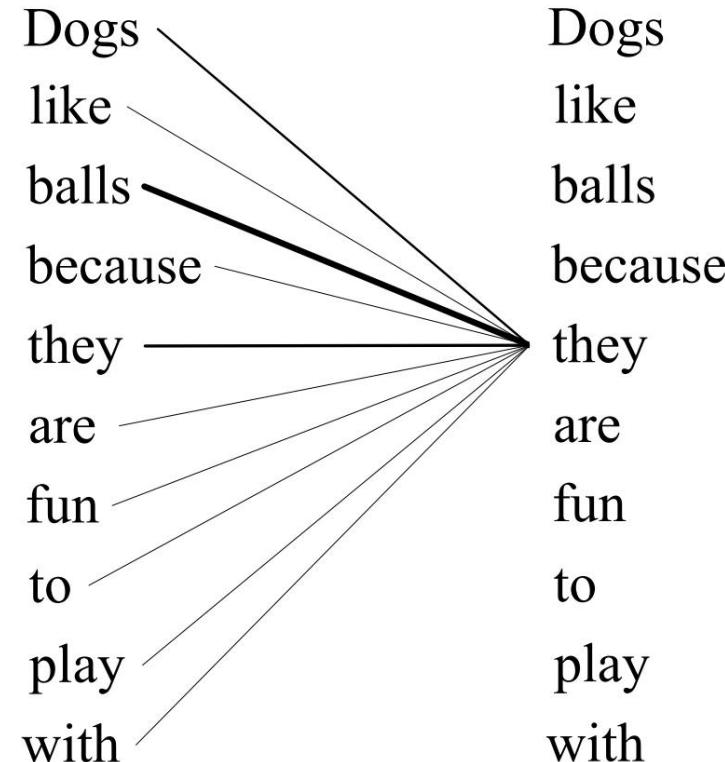


南京航空航天大學
NANJING UNIVERSITY OF AERONAUTICS AND ASTRONAUTICS



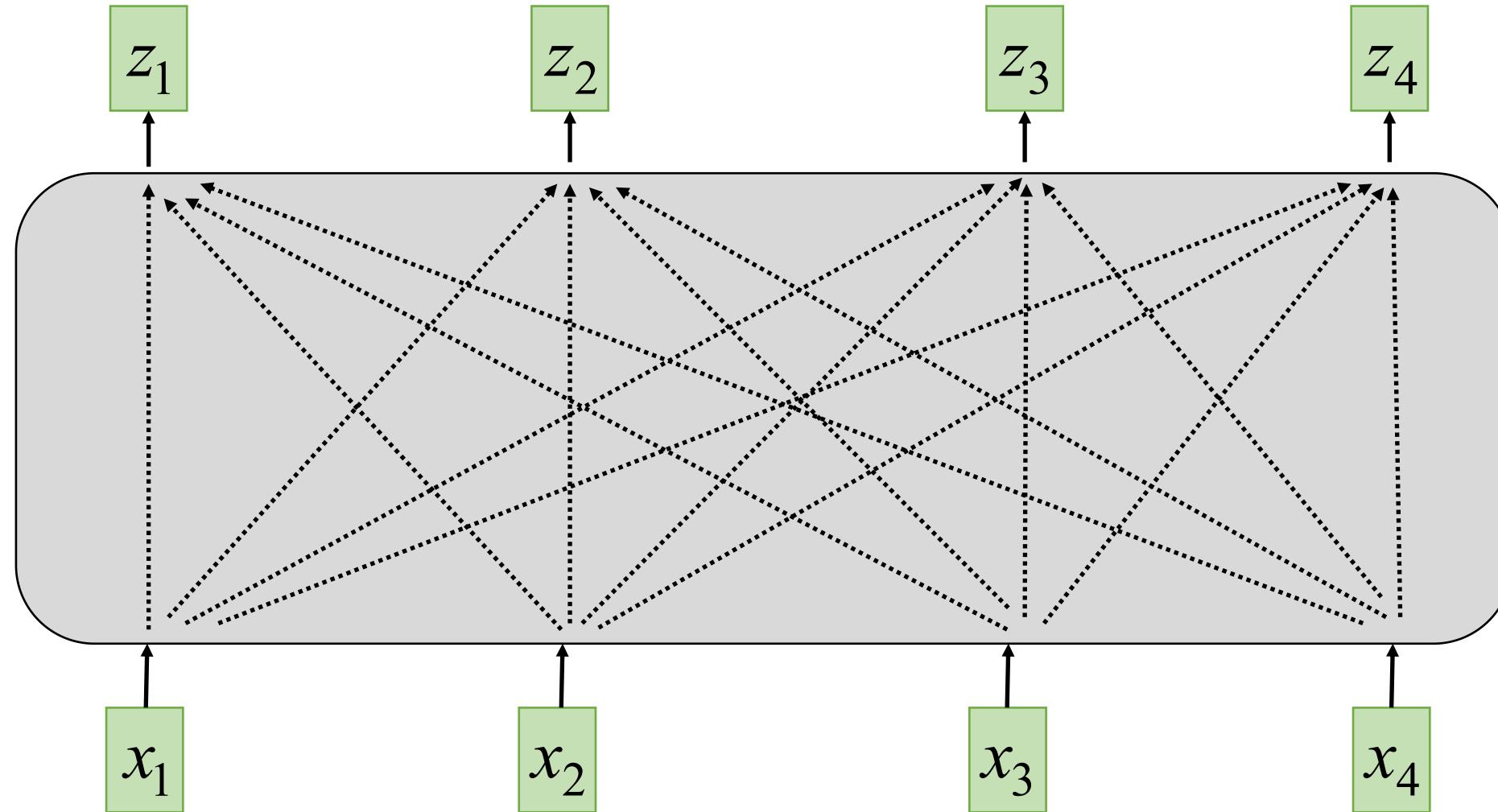
注意力机制是对人类行为的一种仿生，起源于对人类视觉注意机制的研究

1. 自注意力模块

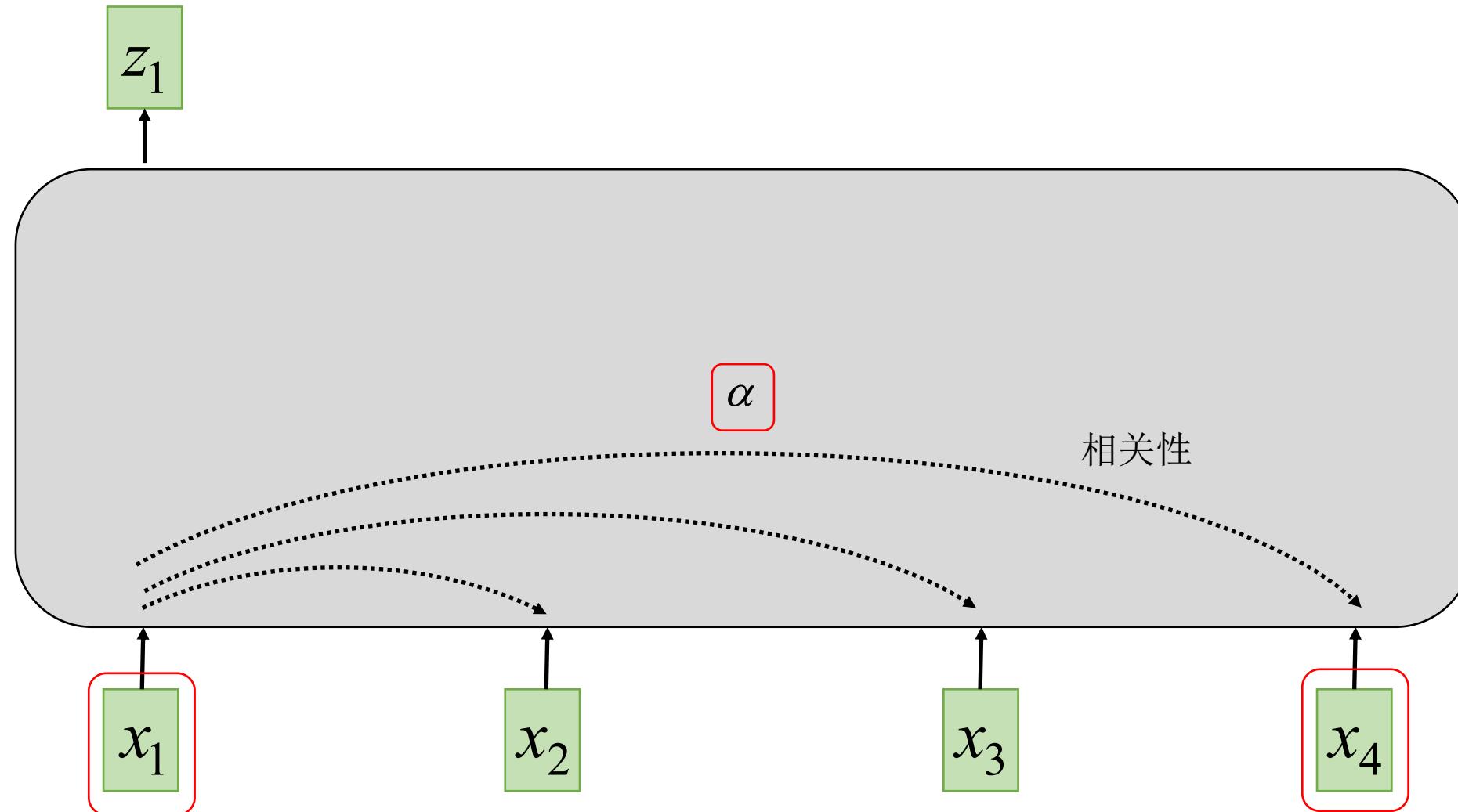


不同的单词通过不同权重计算影响

注意力机制



注意力机制

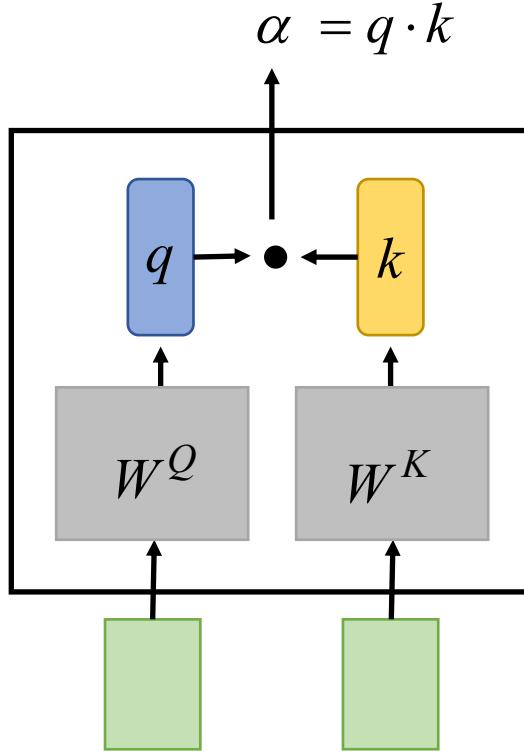


通过词与词之间的关系来更好地理解当前词的意思

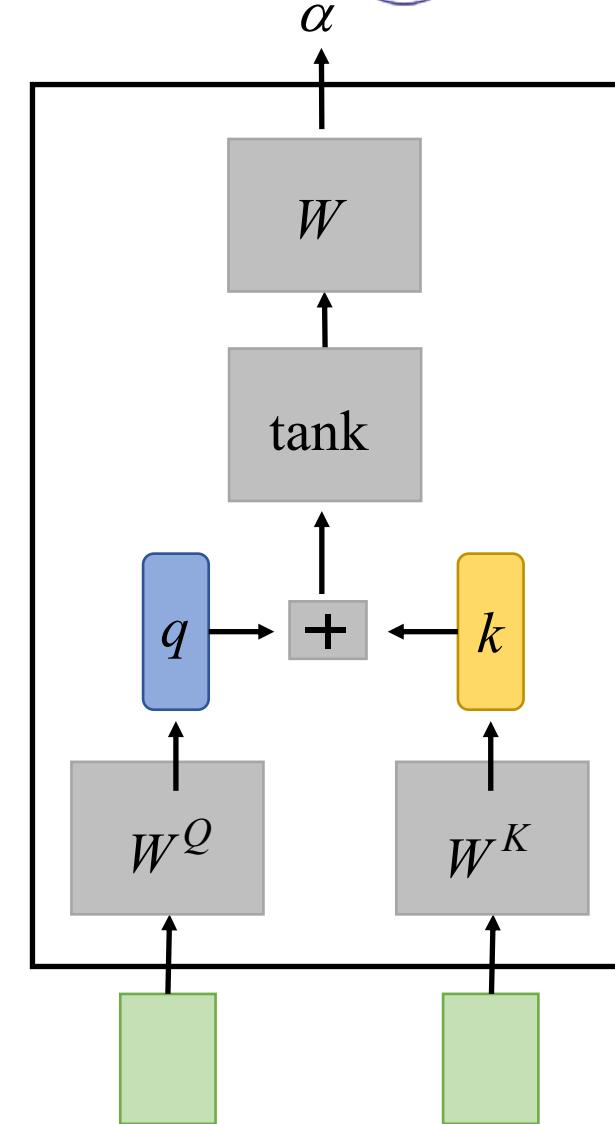
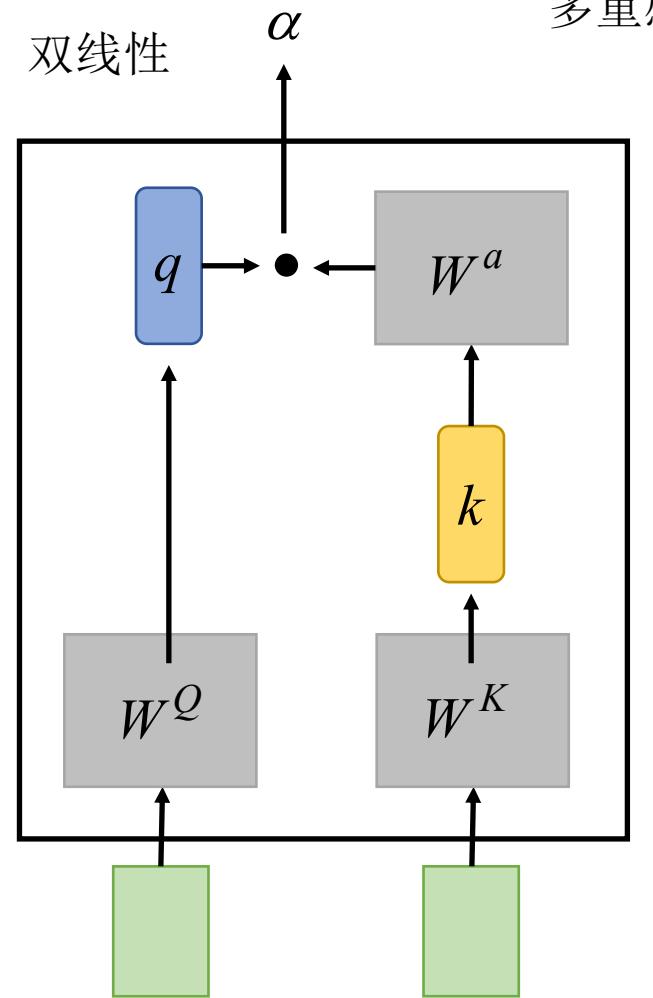
注意力机制



点积



双线性
多重感知机



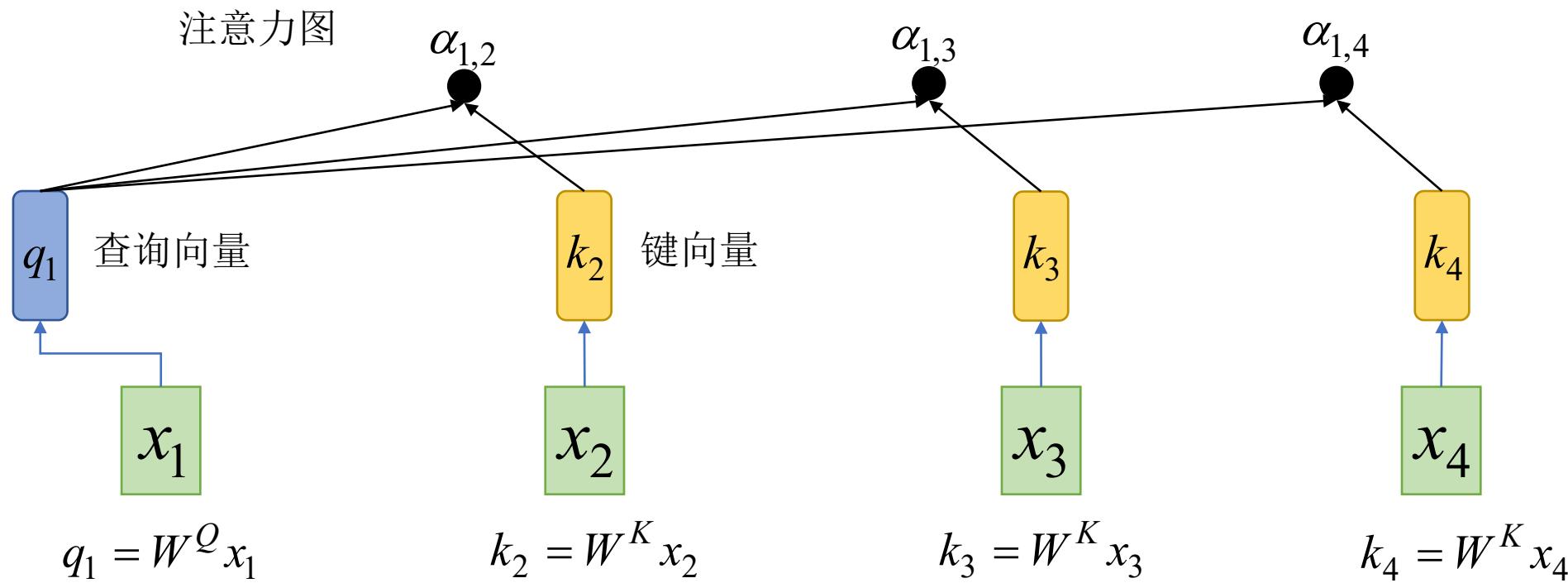
注意力机制



$$\alpha_{1,2} = q_1 \cdot k_2$$

$$\alpha_{1,3} = q_1 \cdot k_3$$

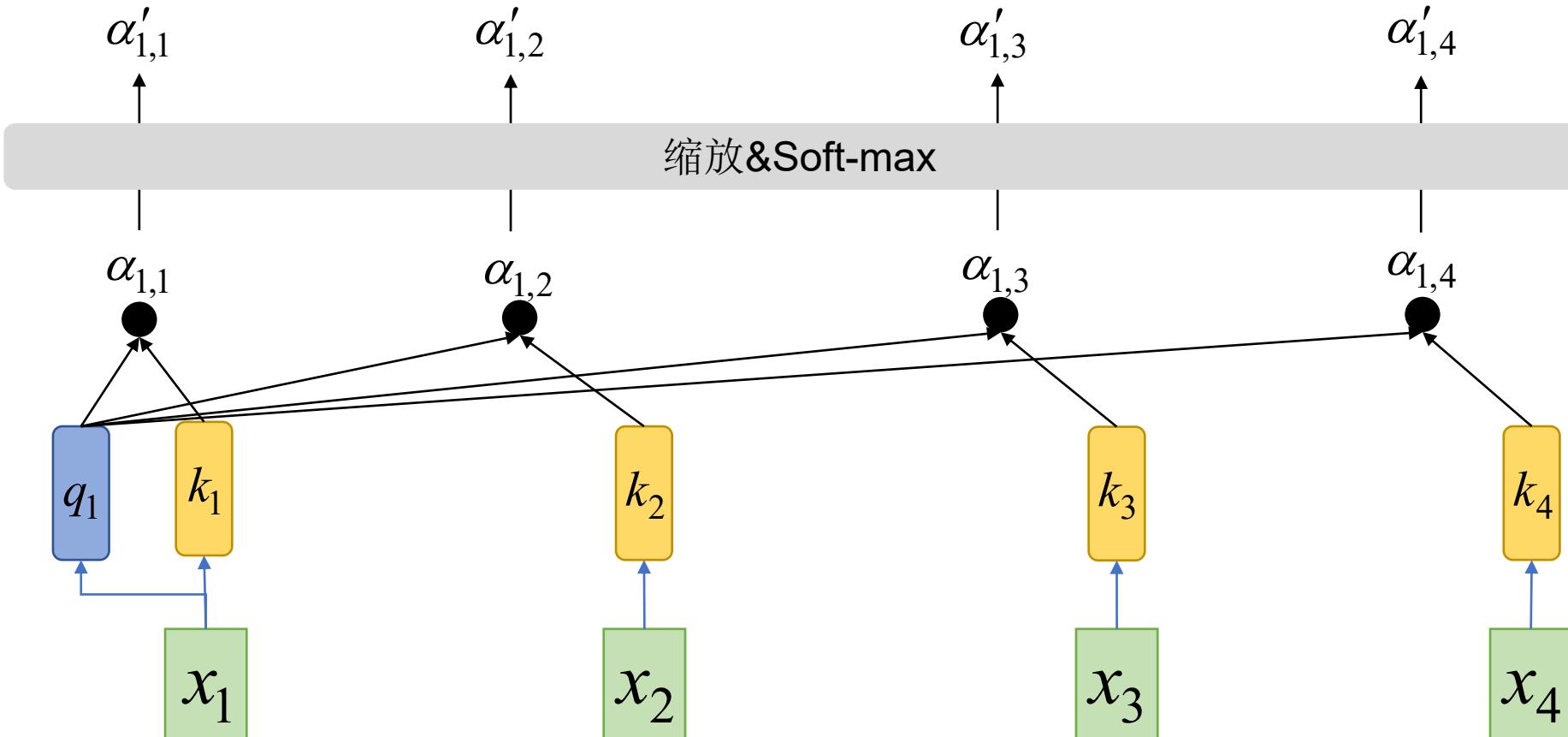
$$\alpha_{1,4} = q_1 \cdot k_4$$



注意力机制



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



$$q_1 = W^Q x_1$$

$$k_1 = W^K x_1$$

$$k_2 = W^K x_2$$

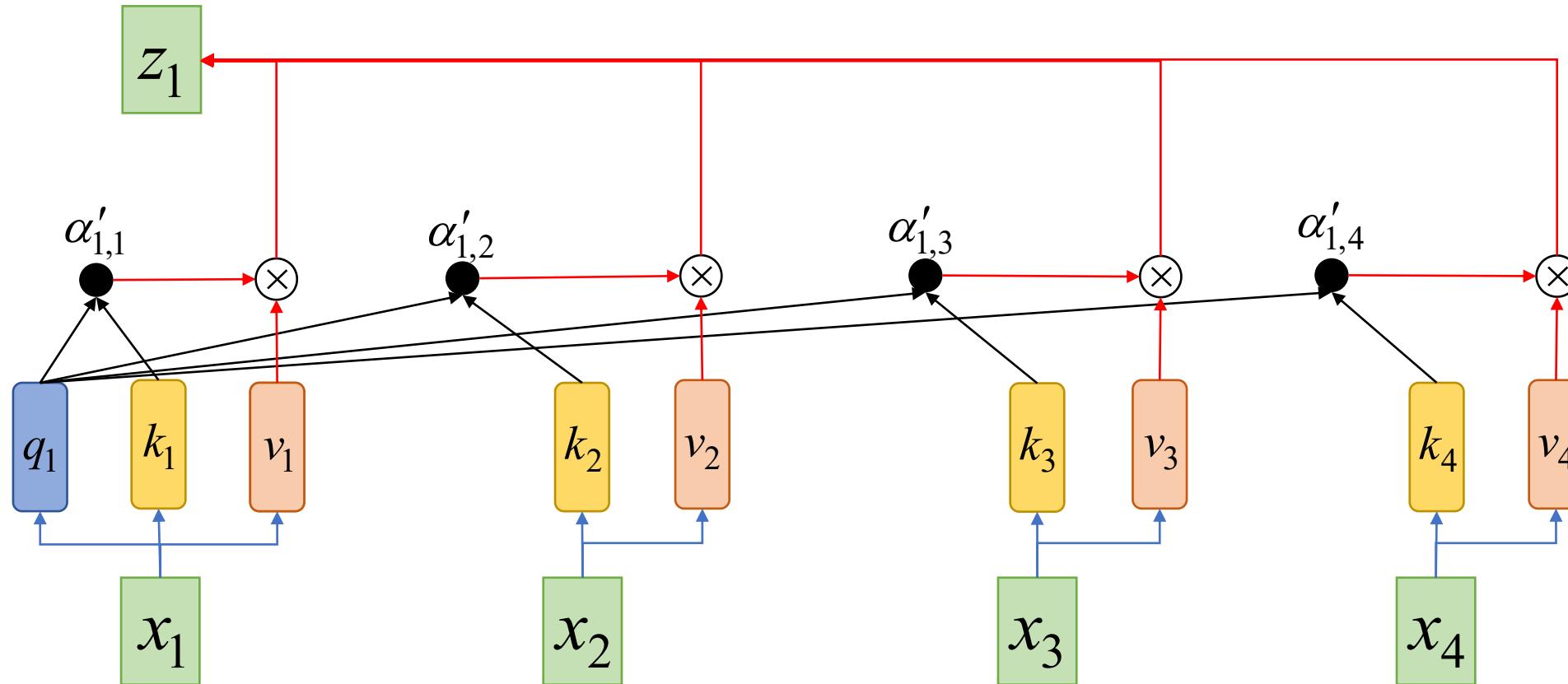
$$k_3 = W^K x_3$$

$$k_4 = W^K x_4$$

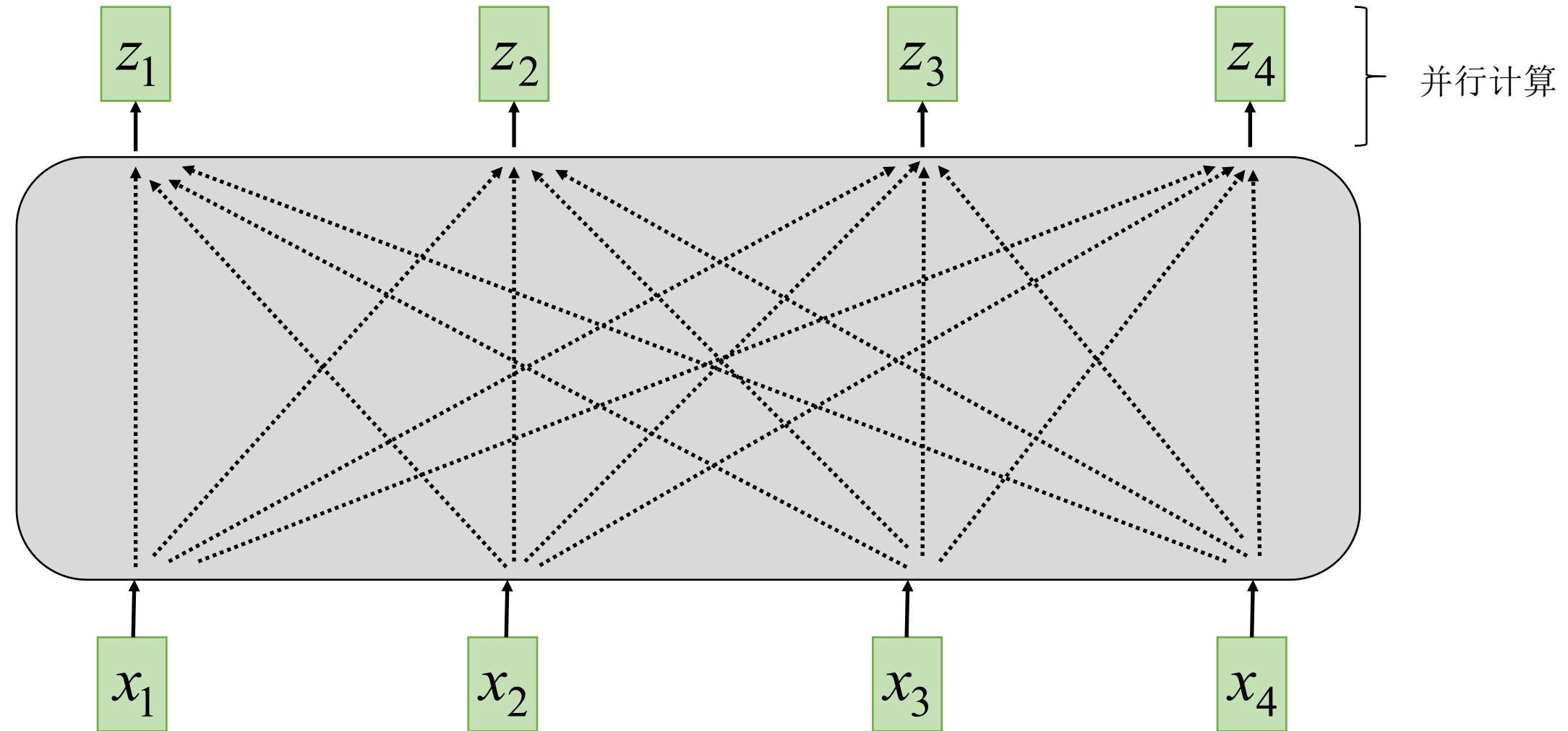
注意力机制



$$z_1 = \sum_i \alpha'_{1,i} v_i$$



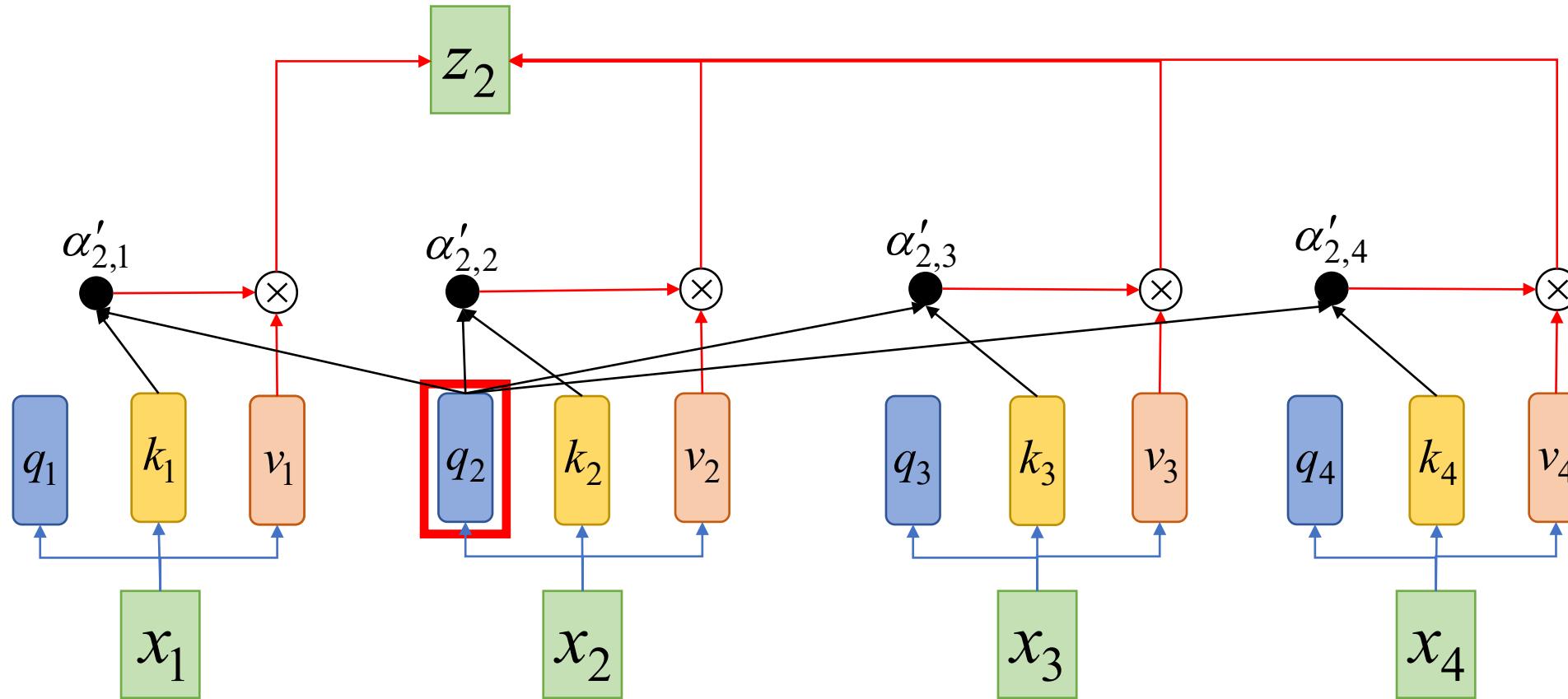
注意力机制



注意力机制

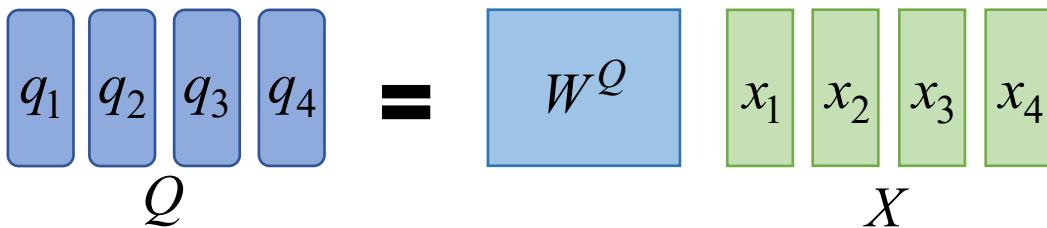


$$z_2 = \sum_i \alpha'_{2,i} v_i$$

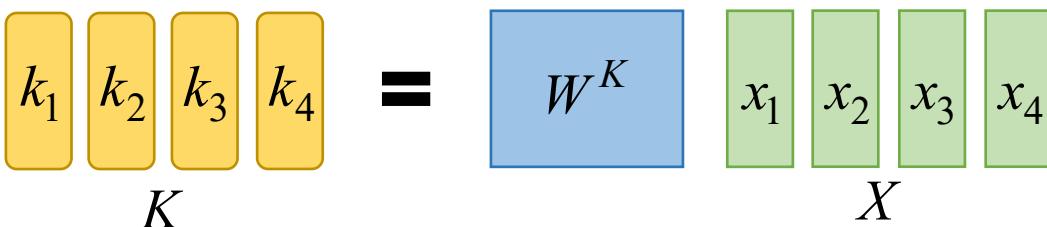


注意力机制

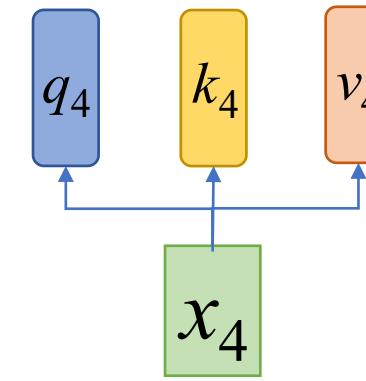
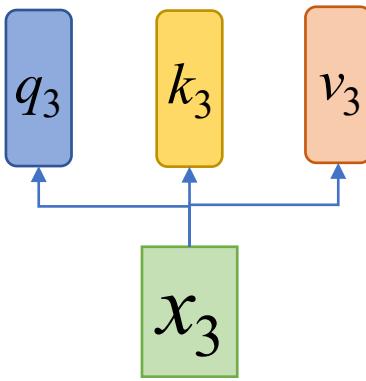
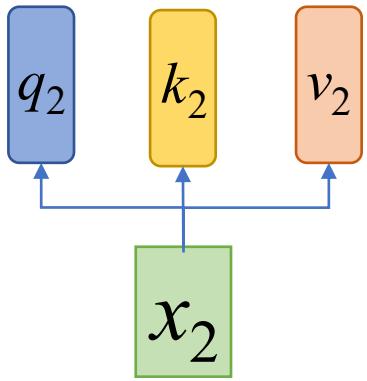
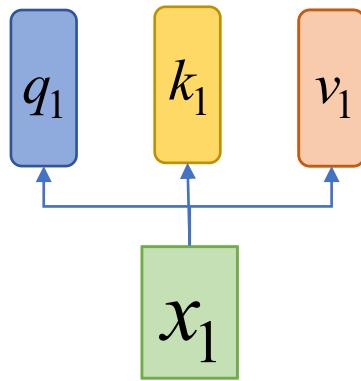
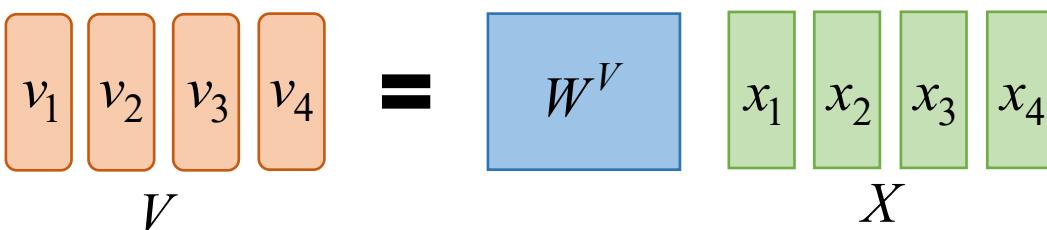
$$q_i = W^Q x_i$$



$$k_i = W^K x_i$$



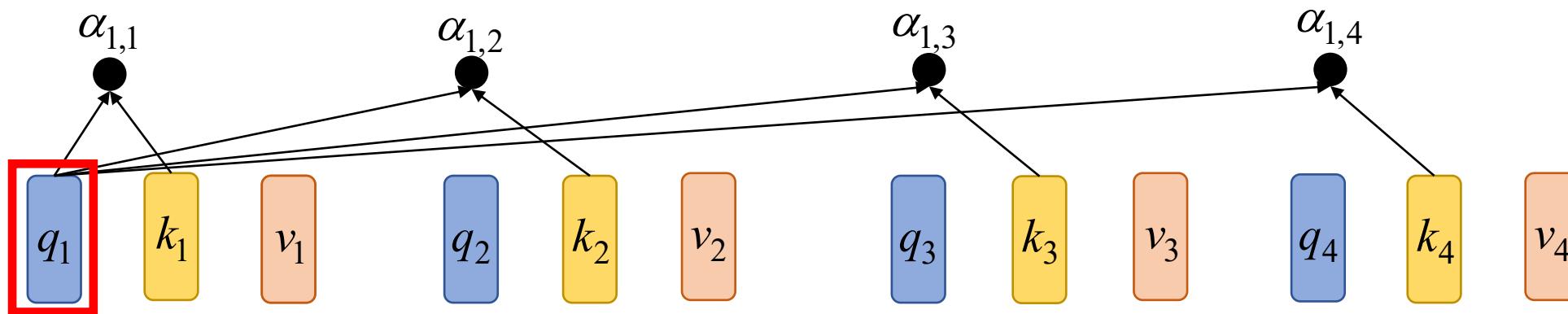
$$v_i = W^V x_i$$



注意力机制



$$\begin{array}{ll} \alpha_{1,1} = \begin{matrix} k_1 \\ q_1 \end{matrix} & \alpha_{1,2} = \begin{matrix} k_2 \\ q_1 \end{matrix} \\ \alpha_{1,3} = \begin{matrix} k_3 \\ q_1 \end{matrix} & \alpha_{1,4} = \begin{matrix} k_4 \\ q_1 \end{matrix} \end{array} \quad \begin{matrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{matrix} = \begin{matrix} k_1 \\ k_2 \\ k_3 \\ k_4 \end{matrix} \quad q_1$$



注意力机制



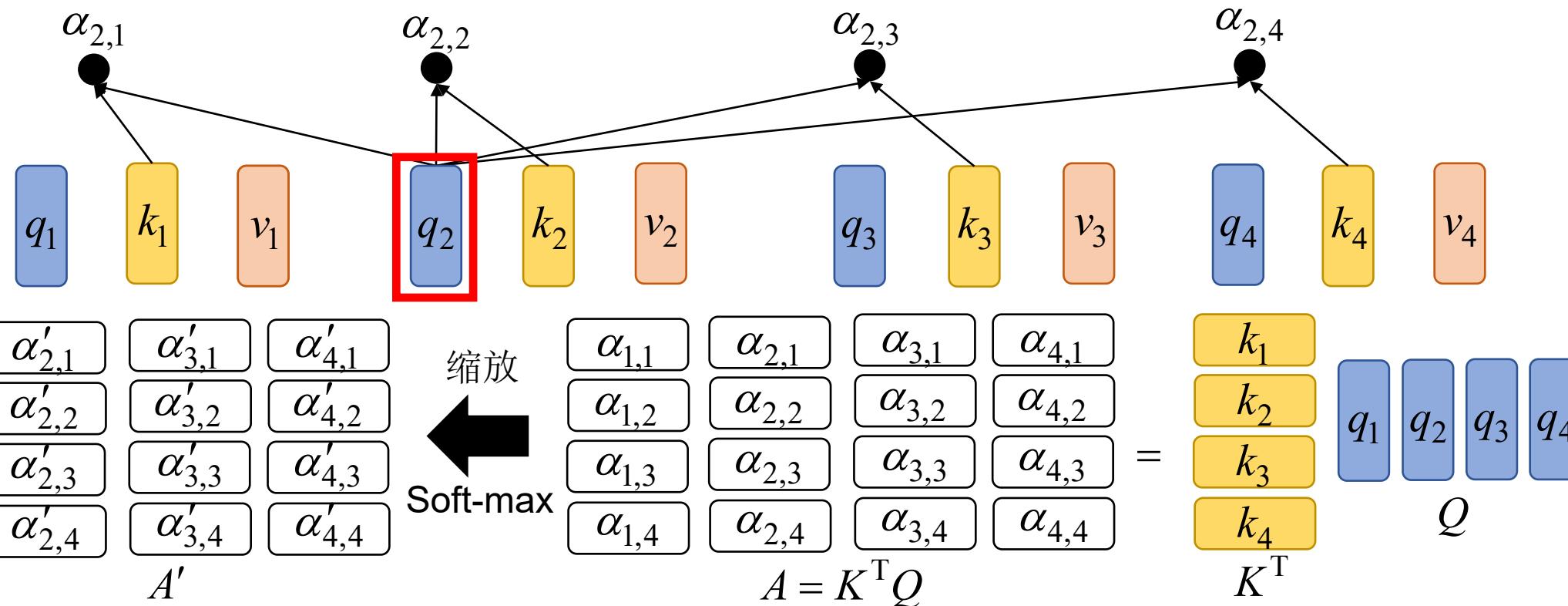
$$\alpha_{1,1} = \begin{matrix} k_1 \\ q_1 \end{matrix}$$

$$\alpha_{1,2} = \begin{matrix} k_2 \\ q_1 \end{matrix}$$

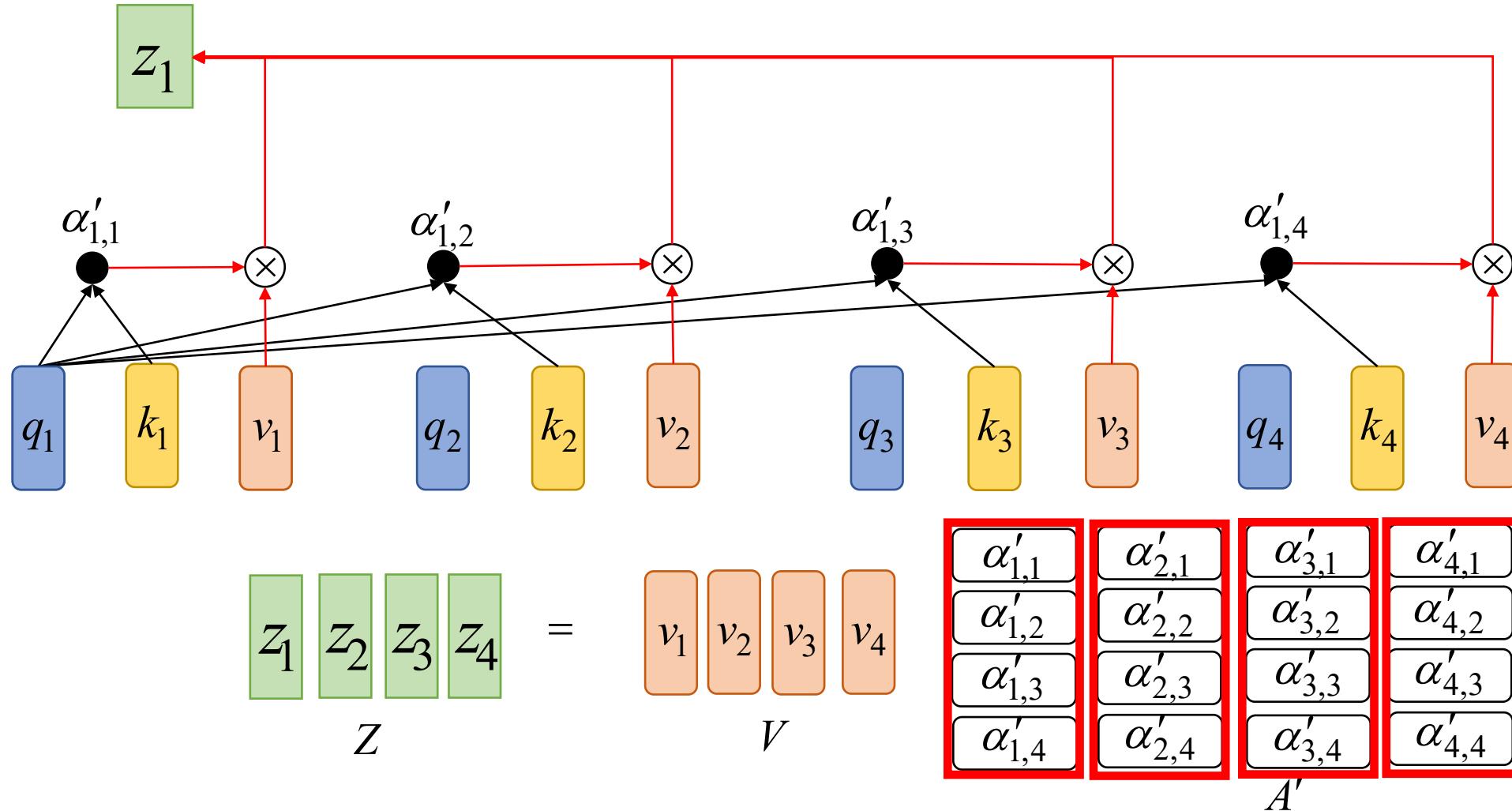
$$\alpha_{1,3} = \begin{matrix} k_3 \\ q_1 \end{matrix}$$

$$\alpha_{1,4} = \begin{matrix} k_4 \\ q_1 \end{matrix}$$

$$\begin{matrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{matrix} = \begin{matrix} k_1 \\ k_2 \\ k_3 \\ k_4 \end{matrix} \begin{matrix} q_1 \end{matrix}$$



注意力机制



注意力机制

$$Z = \begin{matrix} Dogs & like & balls \\ Dogs & [0.90 & 0.03 & 0.07] \\ like & [0.02 & 0.95 & 0.03] \\ balls & [0.02 & 0.01 & 0.97] \end{matrix} \times \begin{bmatrix} v_{11} & v_{12} & \cdots & v_{1d'} \\ v_{21} & v_{22} & \cdots & v_{2d'} \\ v_{31} & v_{32} & \cdots & v_{3d'} \end{bmatrix} v_1 \\ softmax\left(\frac{QK^T}{\sqrt{d_k}}\right) \quad V \end{math>$$

$$QK^T = \begin{matrix} Dogs & like & balls \\ Dogs & [q_{11} & q_{12} & \cdots & q_{1d'}] \\ like & [q_{21} & q_{22} & \cdots & q_{2d'}] \\ balls & [q_{31} & q_{32} & \cdots & q_{3d'}] \end{matrix} q_1 \times \begin{bmatrix} k_{11} & k_{21} & k_{31} \\ k_{12} & k_{22} & k_{32} \\ \vdots & \vdots & \vdots \\ k_{1d'} & k_{2d'} & k_{3d'} \end{bmatrix} k_1 \\ Q \quad K^T \end{math>$$

需要学习的参数

$$Q = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ x_{31} & x_{32} & \cdots & x_{3d} \end{bmatrix} x_1 \times \begin{bmatrix} w_{11}^Q & w_{12}^Q & \cdots & w_{1d'}^Q \\ w_{21}^Q & w_{22}^Q & \cdots & w_{2d'}^Q \\ \vdots & \vdots & \ddots & \vdots \\ w_{d1}^Q & w_{d2}^Q & \cdots & w_{dd'}^Q \end{bmatrix} W^Q$$

$$K = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ x_{31} & x_{32} & \cdots & x_{3d} \end{bmatrix} x_2 \times \begin{bmatrix} w_{11}^K & w_{12}^K & \cdots & w_{1d'}^K \\ w_{21}^K & w_{22}^K & \cdots & w_{2d'}^K \\ \vdots & \vdots & \ddots & \vdots \\ w_{d1}^K & w_{d2}^K & \cdots & w_{dd'}^K \end{bmatrix} W^K$$

$$V = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ x_{31} & x_{32} & \cdots & x_{3d} \end{bmatrix} x_3 \times \begin{bmatrix} w_{11}^V & w_{12}^V & \cdots & w_{1d'}^V \\ w_{21}^V & w_{22}^V & \cdots & w_{2d'}^V \\ \vdots & \vdots & \ddots & \vdots \\ w_{d1}^V & w_{d2}^V & \cdots & w_{dd'}^V \end{bmatrix} W^V$$

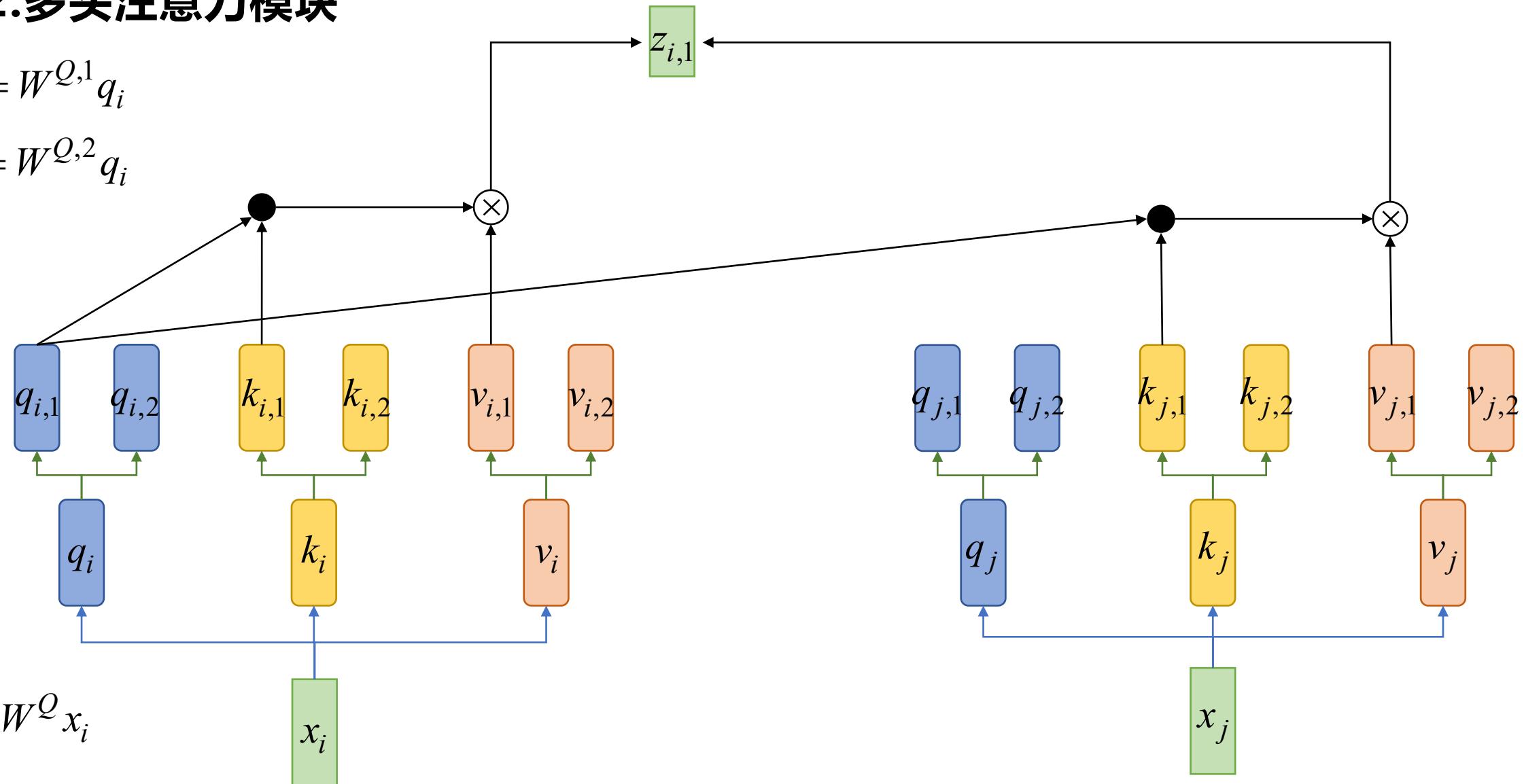
注意力机制



2. 多头注意力模块

$$q_{i,1} = W^{Q,1} q_i$$

$$q_{i,2} = W^{Q,2} q_i$$



$$q_i = W^Q x_i$$

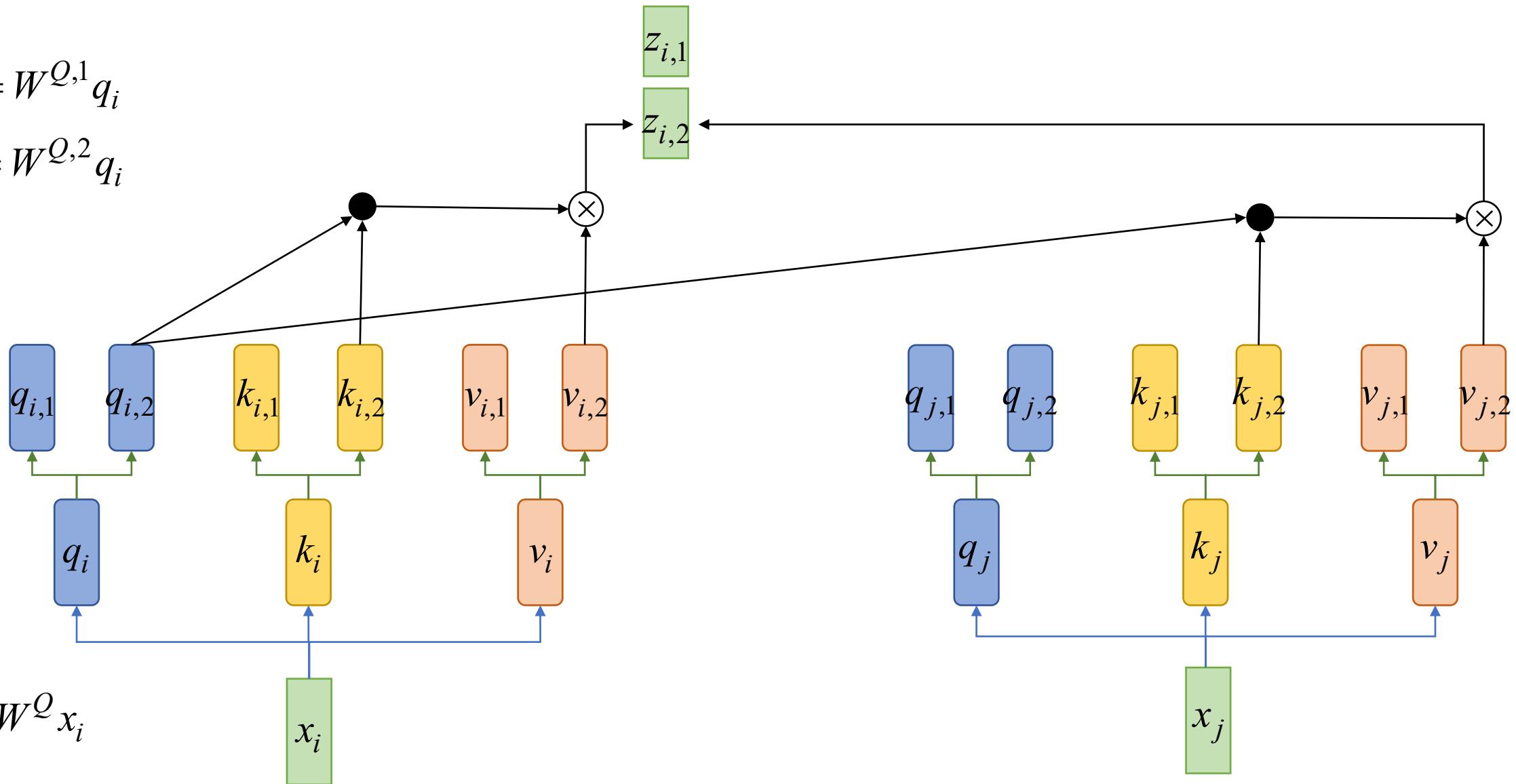
x_i

注意力机制



$$q_{i,1} = W^{Q,1} q_i$$

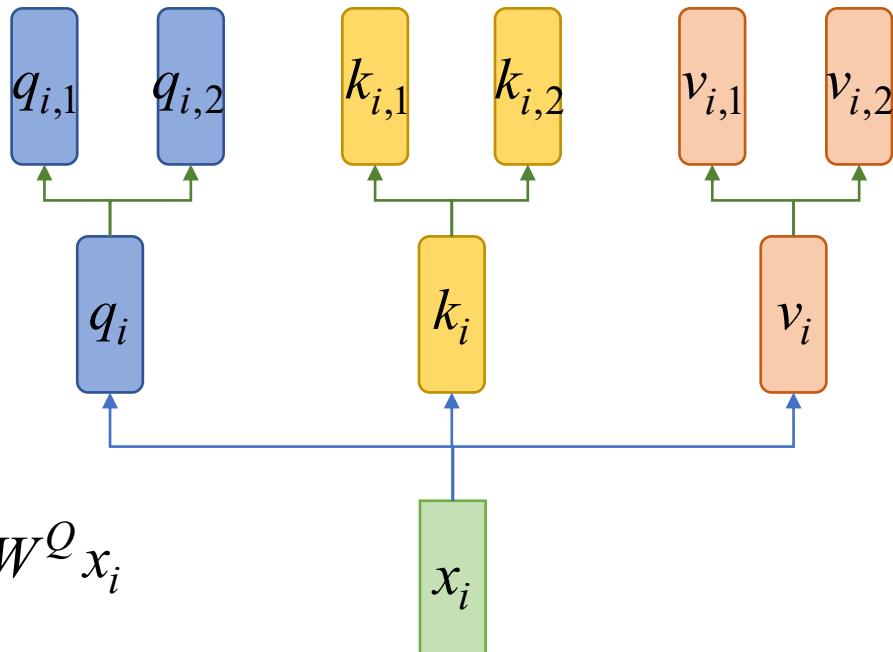
$$q_{i,2} = W^{Q,2} q_i$$



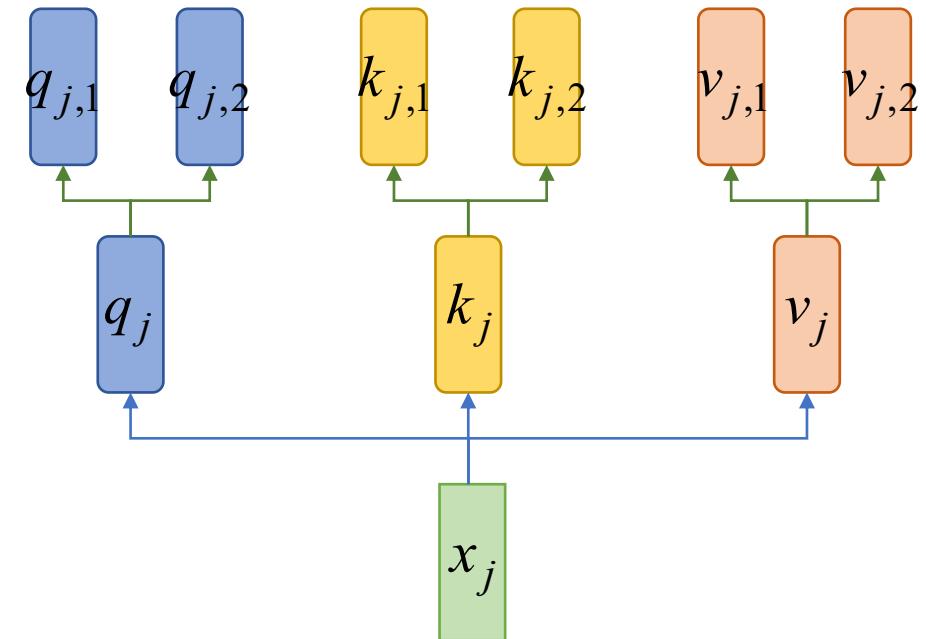
$$q_i = W^Q x_i$$

注意力机制

$$z_i = W^O \begin{bmatrix} z_{i,1} \\ z_{i,2} \end{bmatrix}$$



$$q_i = W^Q x_i$$



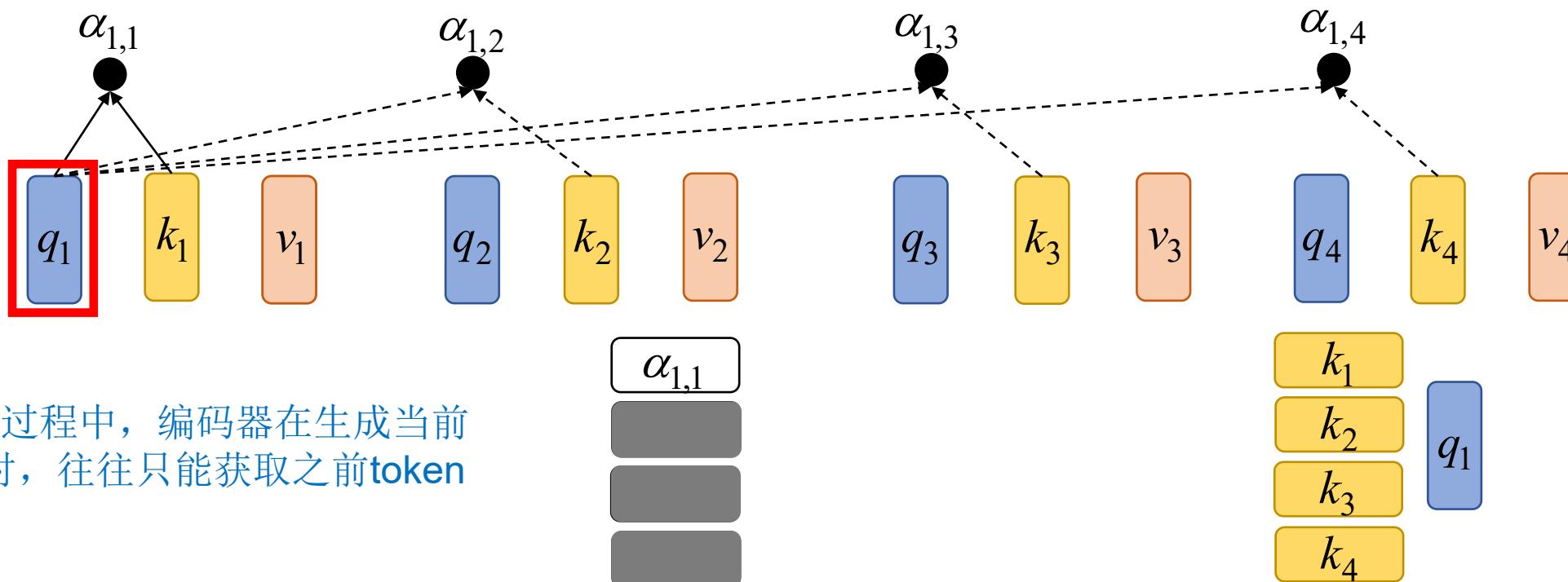
3. 掩码多头注意力模块



$$\alpha_{1,1} = \begin{matrix} k_1 \\ q_1 \end{matrix} \quad \alpha_{1,2} = \begin{matrix} k_2 \\ q_1 \end{matrix} = 0$$

$$\alpha_{1,3} = \begin{matrix} k_3 \\ q_1 \end{matrix} = 0 \quad \alpha_{1,4} = \begin{matrix} k_4 \\ q_1 \end{matrix} = 0$$

$$\begin{matrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{matrix} = \begin{matrix} k_1 \\ k_2 \\ k_3 \\ k_4 \end{matrix} \begin{matrix} q_1 \end{matrix}$$



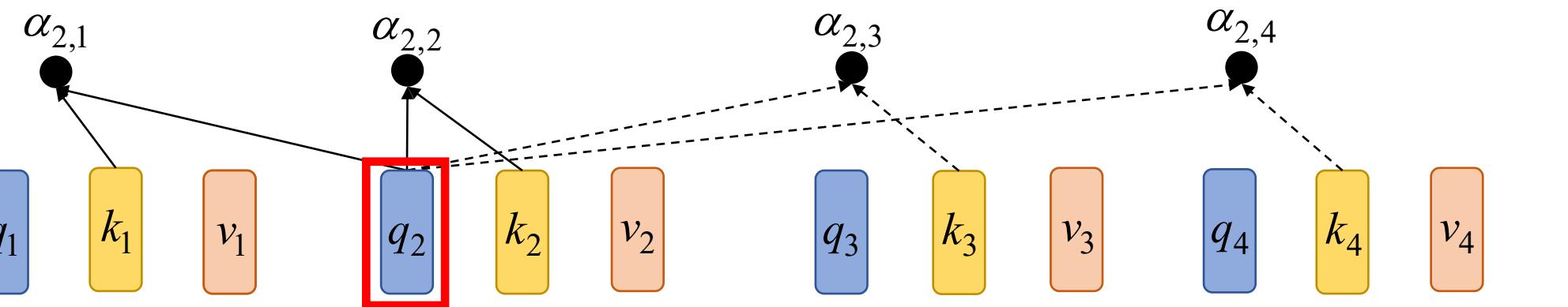
在推理过程中，编码器在生成当前 token 时，往往只能获取之前 token 的信息

注意力机制



$$\begin{aligned}\alpha_{1,1} &= \begin{matrix} k_1 \\ q_1 \end{matrix} & \alpha_{1,2} &= \begin{matrix} k_2 \\ q_1 \end{matrix} & \alpha_{1,3} &= \begin{matrix} k_3 \\ q_1 \end{matrix} & \alpha_{1,4} &= \begin{matrix} k_4 \\ q_1 \end{matrix} \\ &&&&&&& \\ \alpha_{1,1} &= \begin{matrix} k_1 \\ q_1 \end{matrix} & \alpha_{1,2} &= \begin{matrix} k_2 \\ q_1 \end{matrix} & \alpha_{1,3} &= \begin{matrix} k_3 \\ q_1 \end{matrix} & \alpha_{1,4} &= \begin{matrix} k_4 \\ q_1 \end{matrix}\end{aligned}$$

$$\begin{matrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{matrix} = \begin{matrix} k_1 \\ k_2 \\ k_3 \\ k_4 \end{matrix} \begin{matrix} q_1 \end{matrix}$$



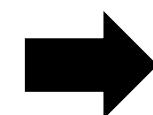
$$\begin{matrix} \alpha'_{1,1} \\ \alpha'_{2,1} \\ \alpha'_{3,1} \\ \alpha'_{4,1} \\ \hline \alpha'_{2,2} \\ \alpha'_{3,2} \\ \alpha'_{4,2} \\ \hline \alpha'_{3,3} \\ \alpha'_{4,3} \\ \hline \alpha'_{4,4} \end{matrix} \xleftarrow[\text{Soft-max}]{\text{缩放}} \begin{matrix} \alpha_{1,1} \\ \alpha_{2,1} \\ \alpha_{3,1} \\ \alpha_{4,1} \\ \hline \alpha_{2,2} \\ \alpha_{3,2} \\ \alpha_{4,2} \\ \hline \alpha_{3,3} \\ \alpha_{4,3} \\ \hline \alpha_{4,4} \end{matrix} = \begin{matrix} k_1 \\ k_2 \\ k_3 \\ k_4 \end{matrix} \begin{matrix} q_1 \\ q_2 \\ q_3 \\ q_4 \end{matrix} Q$$

注意力机制



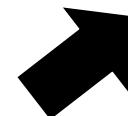
$\alpha'_{1,1}$	$\alpha'_{2,1}$	$\alpha'_{3,1}$	$\alpha'_{4,1}$
	$\alpha'_{2,2}$	$\alpha'_{3,2}$	$\alpha'_{4,2}$
		$\alpha'_{3,3}$	$\alpha'_{4,3}$
			$\alpha'_{4,4}$

A'



$\alpha'_{1,1}$	$\alpha'_{2,1}$	$\alpha'_{3,1}$	$\alpha'_{4,1}$
0	$\alpha'_{2,2}$	$\alpha'_{3,2}$	$\alpha'_{4,2}$
0	0	$\alpha'_{3,3}$	$\alpha'_{4,3}$
0	0	0	$\alpha'_{4,4}$

A'



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

A'



1	1	1	1
0	1	1	1
0	0	1	1
0	0	0	1

$mask$

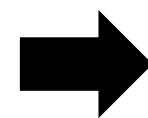
构造掩码矩阵

注意力机制



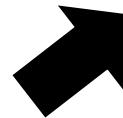
$\alpha'_{1,1}$	$\alpha'_{2,1}$	$\alpha'_{3,1}$	$\alpha'_{4,1}$
	$\alpha'_{2,2}$	$\alpha'_{3,2}$	$\alpha'_{4,2}$
		$\alpha'_{3,3}$	$\alpha'_{4,3}$
			$\alpha'_{4,4}$

A'



$\alpha'_{1,1}$	$\alpha'_{2,1}$	$\alpha'_{3,1}$	$\alpha'_{4,1}$
0			
0	0		
0	0	0	

A'



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

$A = K^T Q$



$-\infty$	$-\infty$	$-\infty$	$-\infty$
0	$-\infty$	$-\infty$	$-\infty$
0	0	$-\infty$	$-\infty$
0	0	0	$-\infty$

$mask$

4.位置编码

在自注意力模块中缺乏句子词序

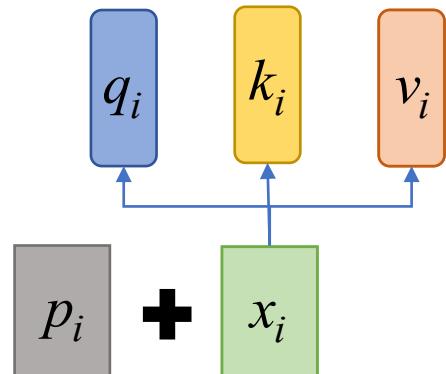
位置编码：每个位置有一个独一无二的位置向量

$$P_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

手工设计的

$$P_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

也可以是学习得到的



$$P = \begin{bmatrix} \sin\left(\frac{pos}{10000^0}\right) & \cos\left(\frac{pos}{10000^0}\right) & \sin\left(\frac{pos}{10000^{2/3}}\right) \\ \sin\left(\frac{pos}{10000^0}\right) & \cos\left(\frac{pos}{10000^0}\right) & \sin\left(\frac{pos}{10000^{2/3}}\right) \\ \sin\left(\frac{pos}{10000^0}\right) & \cos\left(\frac{pos}{10000^0}\right) & \sin\left(\frac{pos}{10000^{2/3}}\right) \end{bmatrix}$$

$$P = \begin{bmatrix} \sin(0) & \cos(0) & \sin\left(\frac{0}{10000^{2/3}}\right) \\ \sin(1) & \cos(1) & \sin\left(\frac{1}{10000^{2/3}}\right) \\ \sin(2) & \cos(2) & \sin\left(\frac{2}{10000^{2/3}}\right) \end{bmatrix}$$

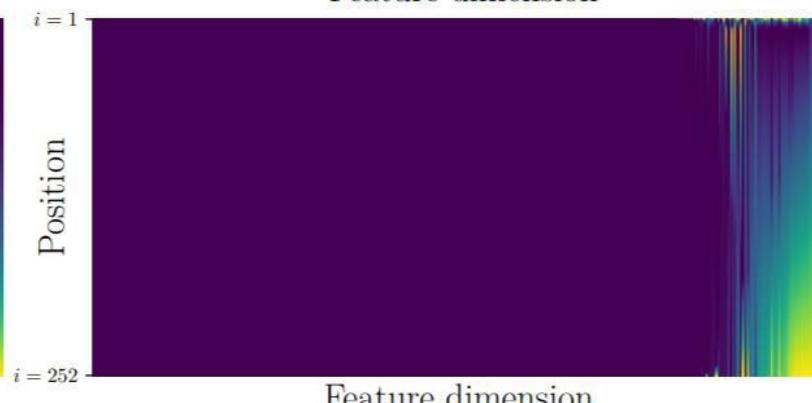
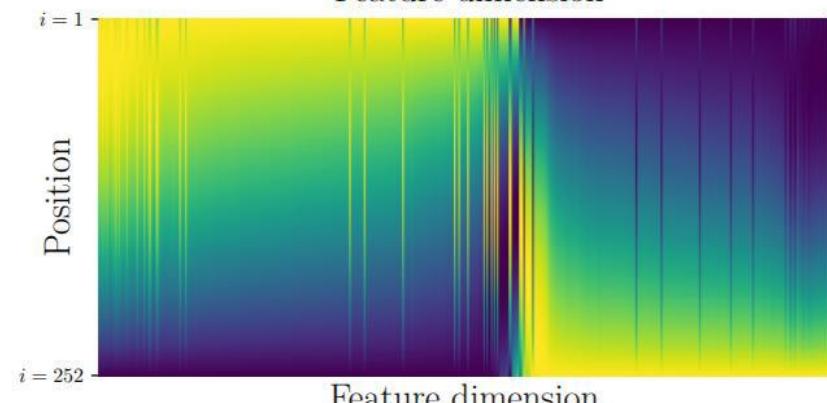
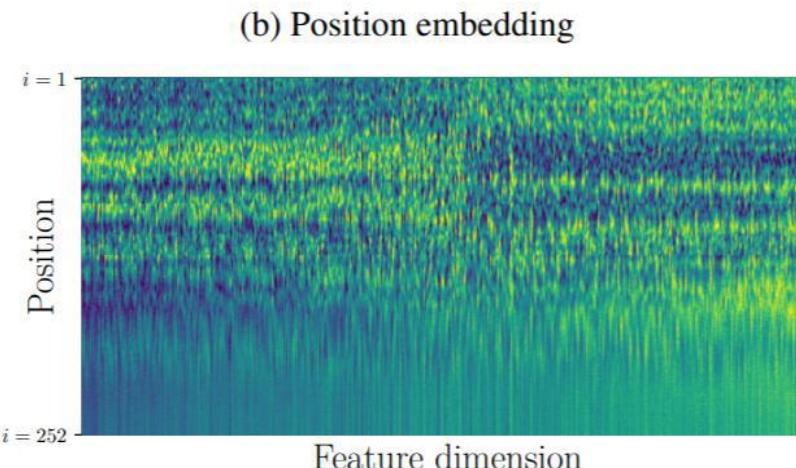
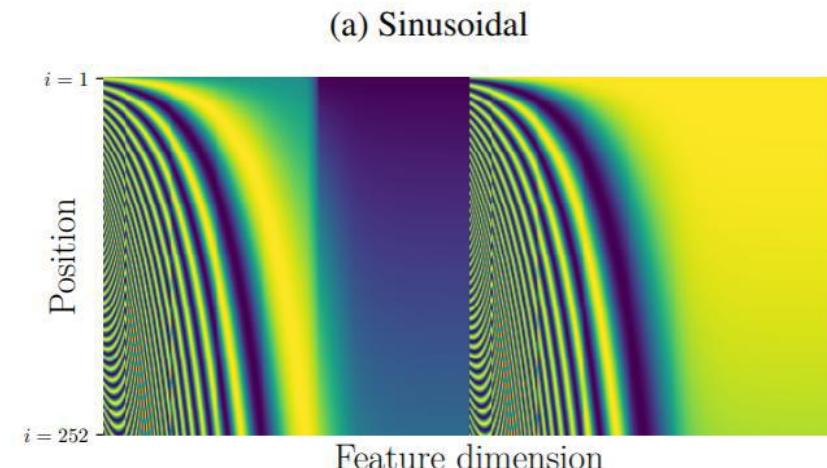
注意力机制



<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	✓	✓	✓



目录



□ Transformer

- 注意力机制
- 编码器-解码器结构
- 大模型中的编码器-解码器结构

□ 编码器结构----BERT家族

- BERT结构
- 预训练策略
- BERT的变体

□ 解码器结构----GPT家族

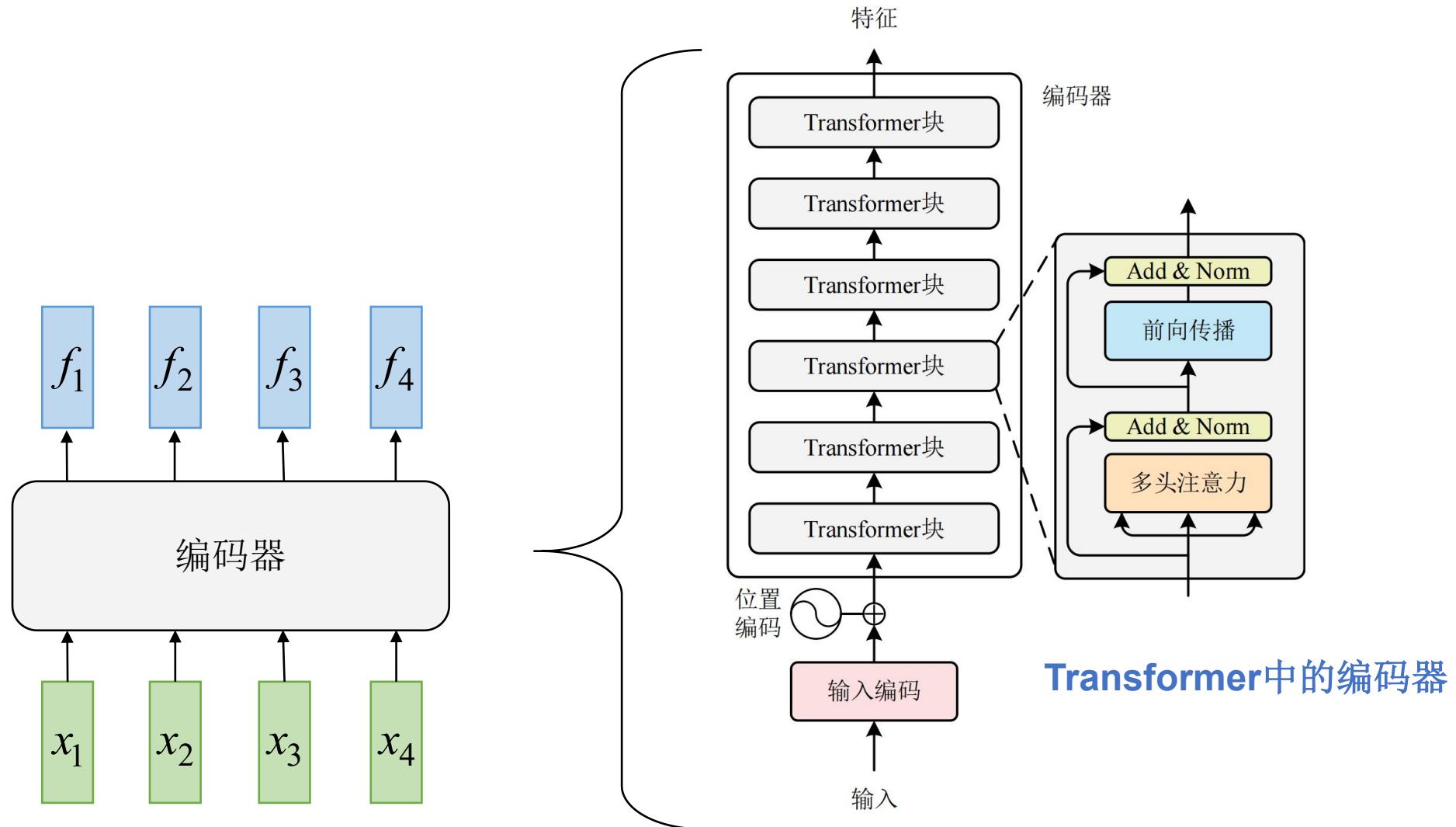
- GPT结构
- 自回归预训练
- 后续改进

□ 思考

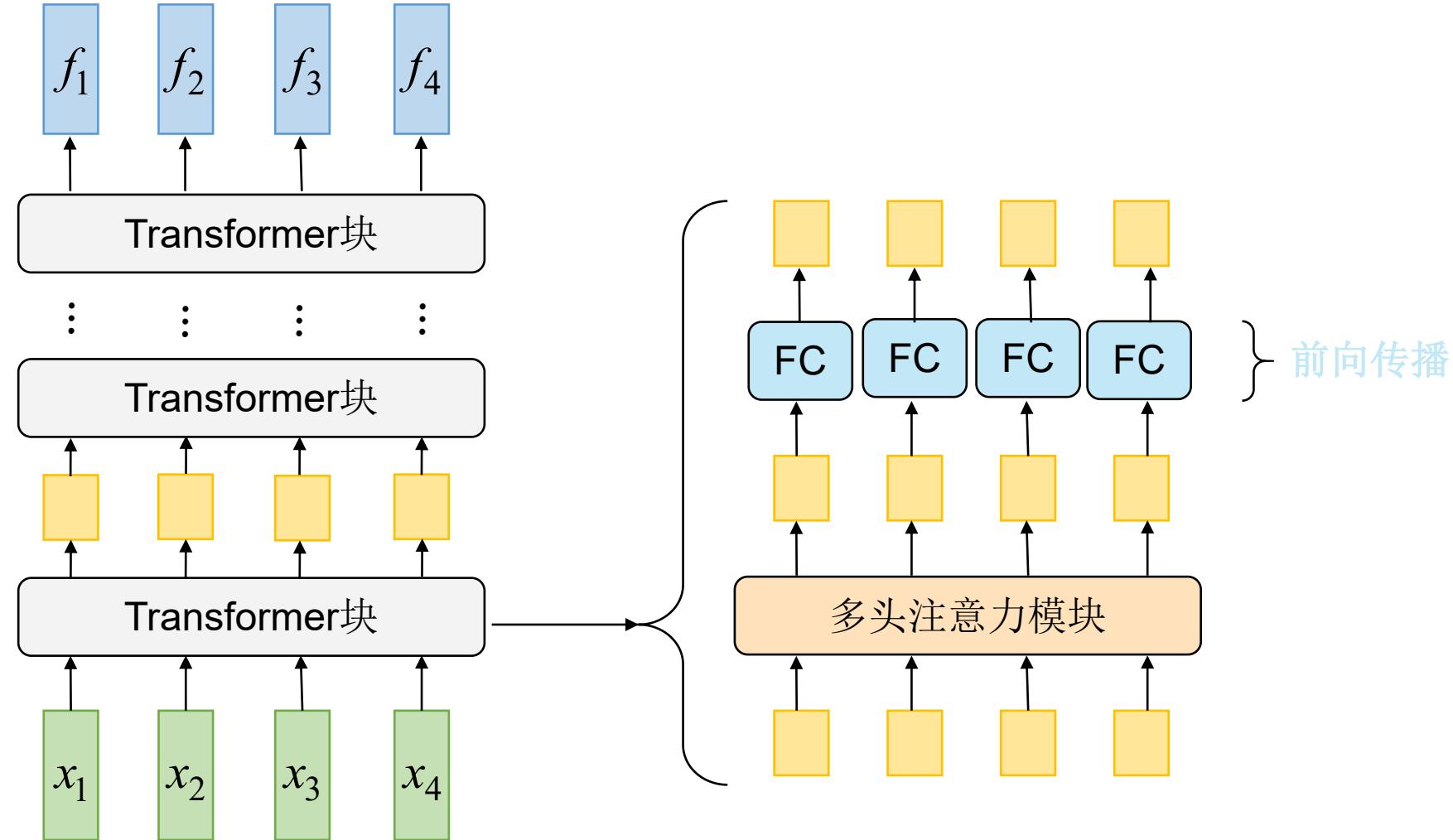
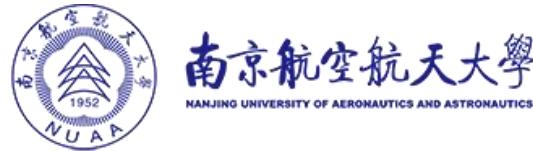
编码器-解码器结构



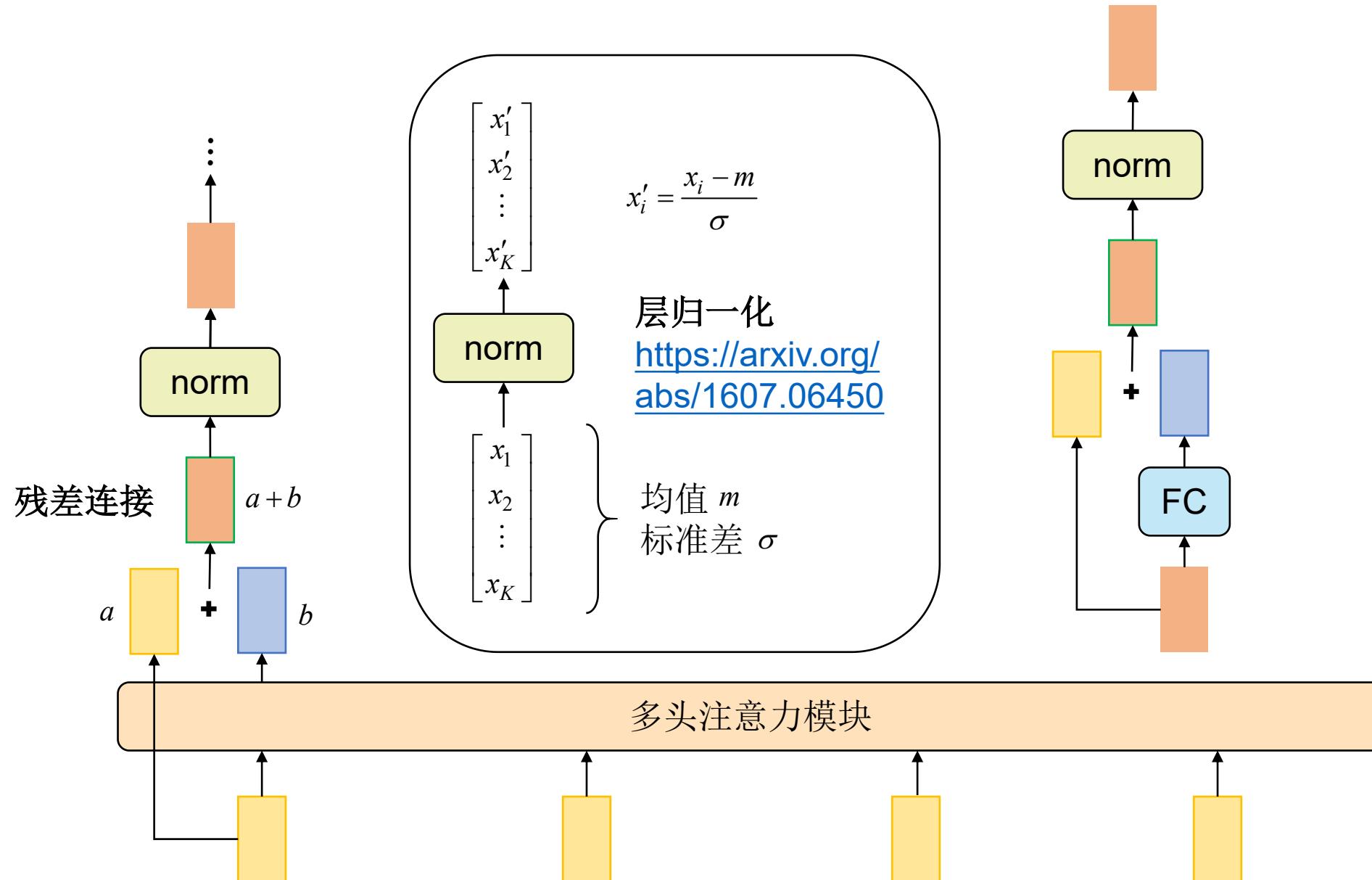
1. 编码器



编码器-解码器结构



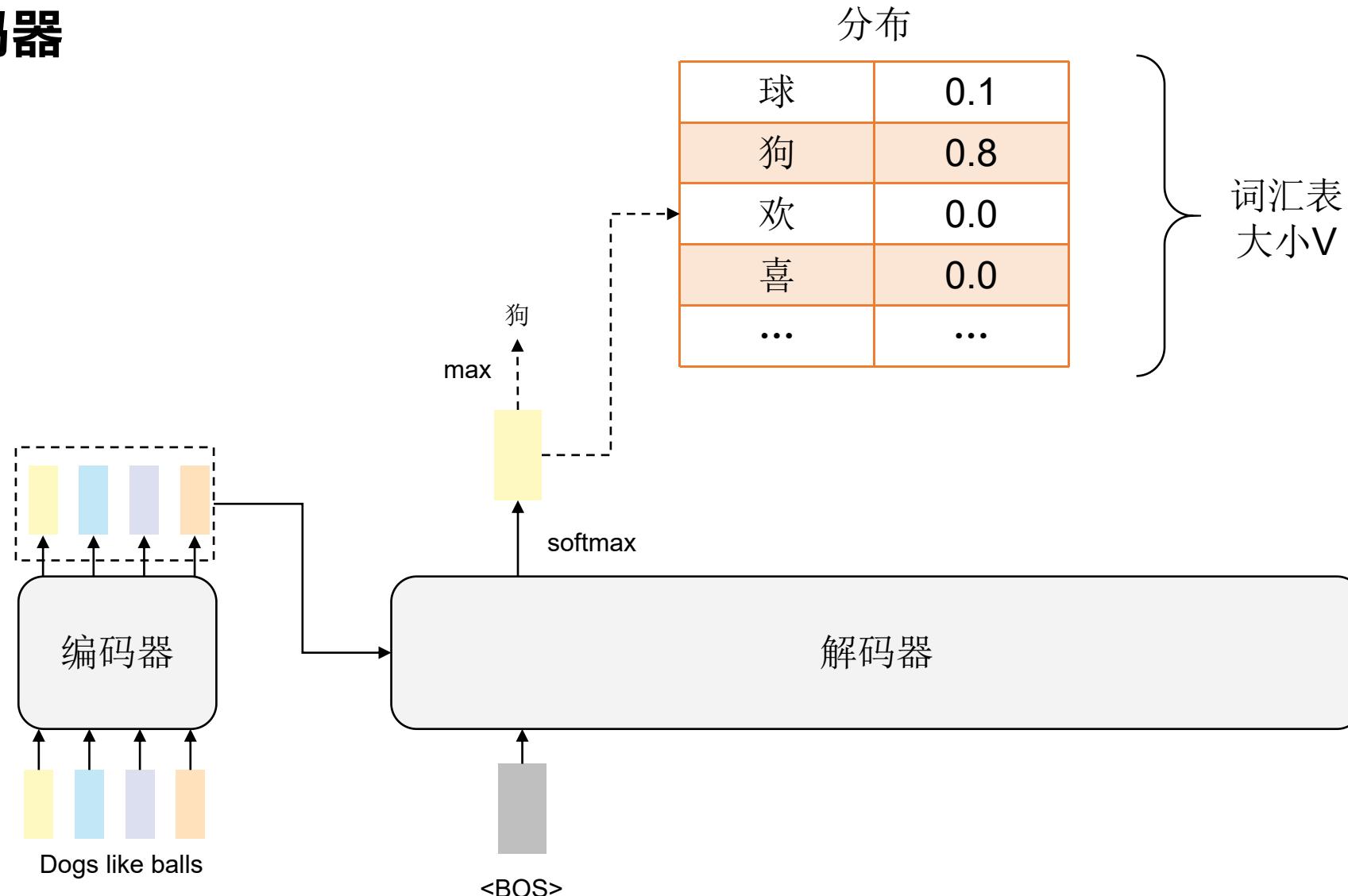
编码器-解码器结构



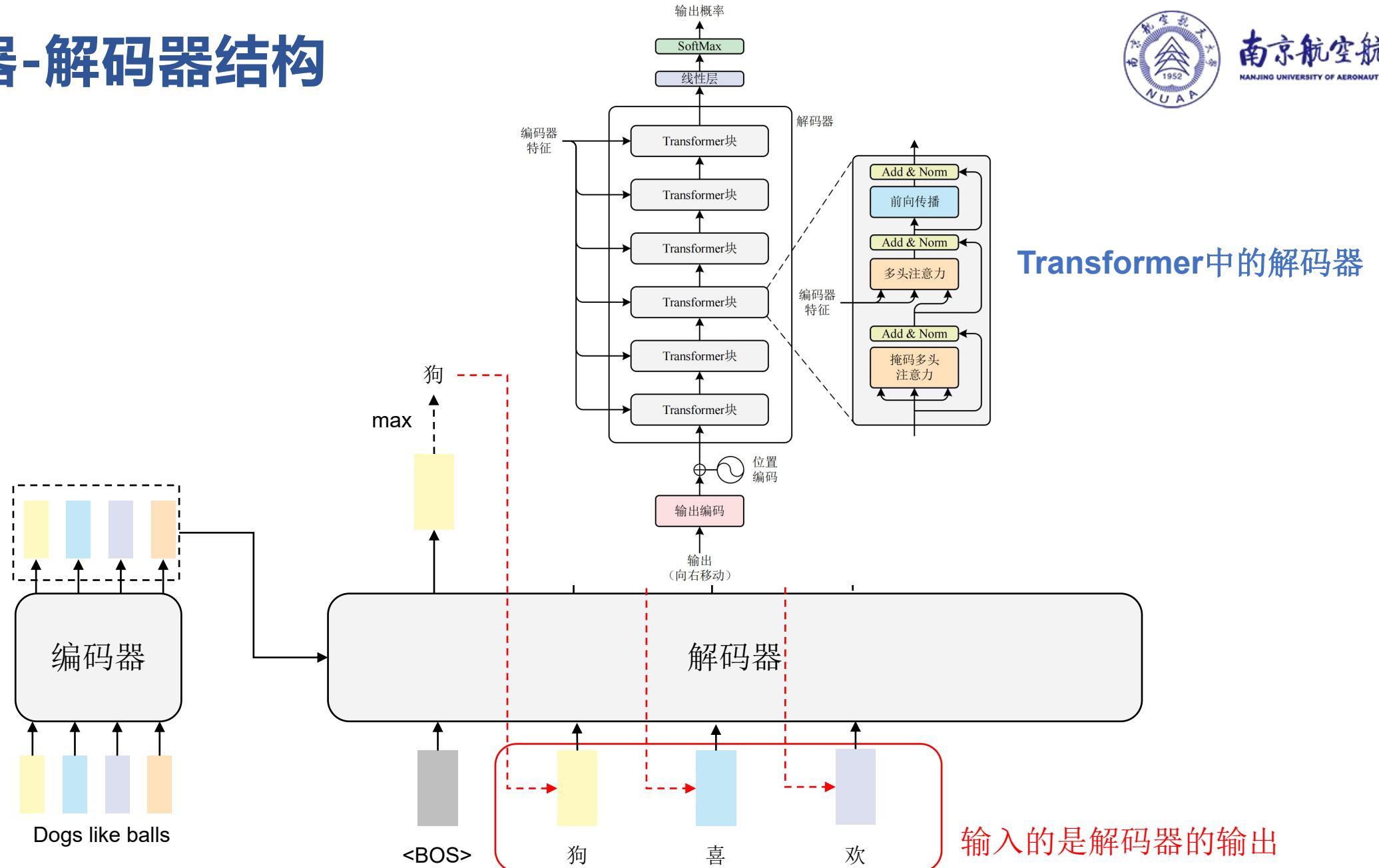
编码器-解码器结构



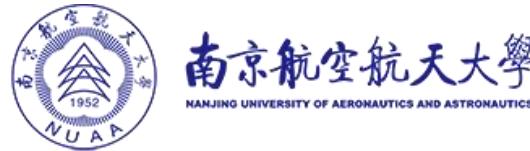
2. 解码器



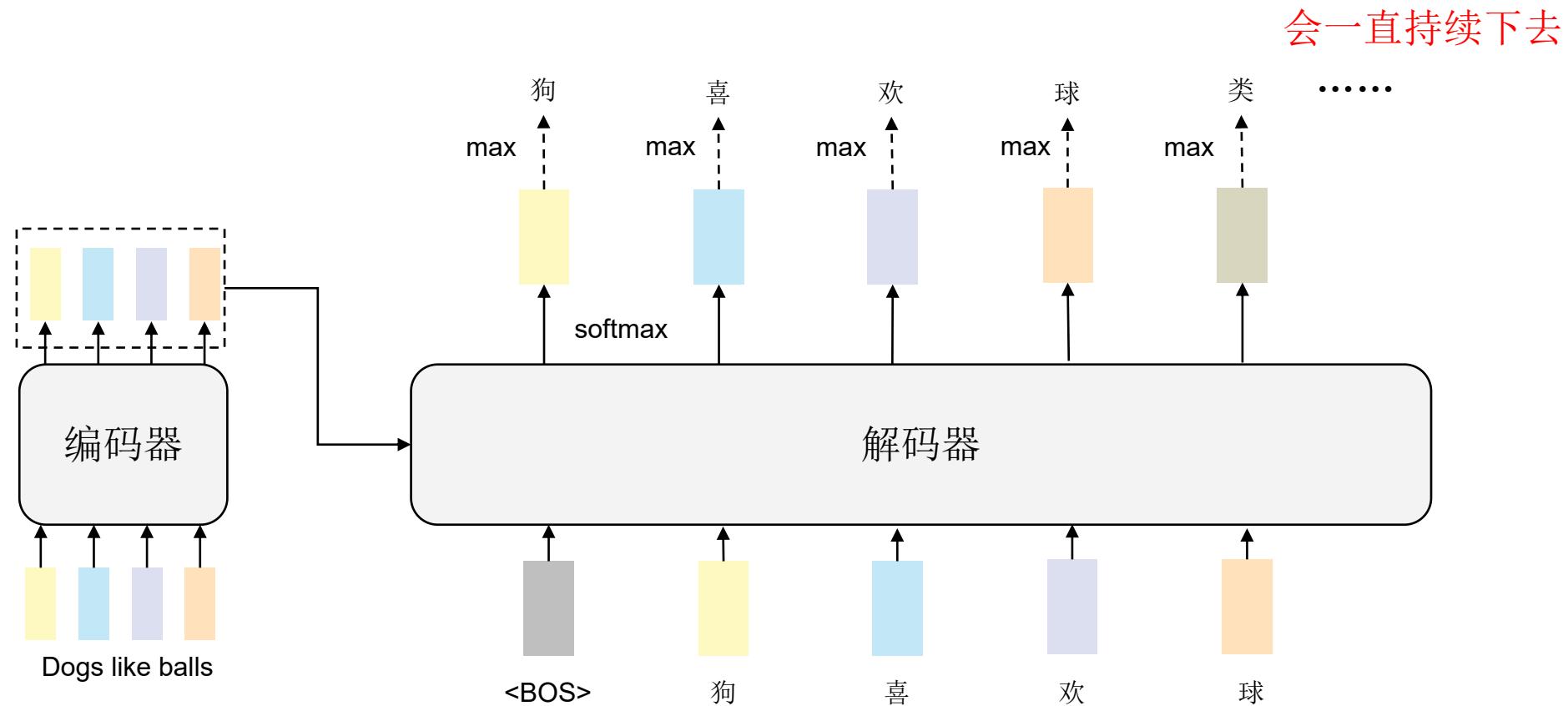
编码器-解码器结构



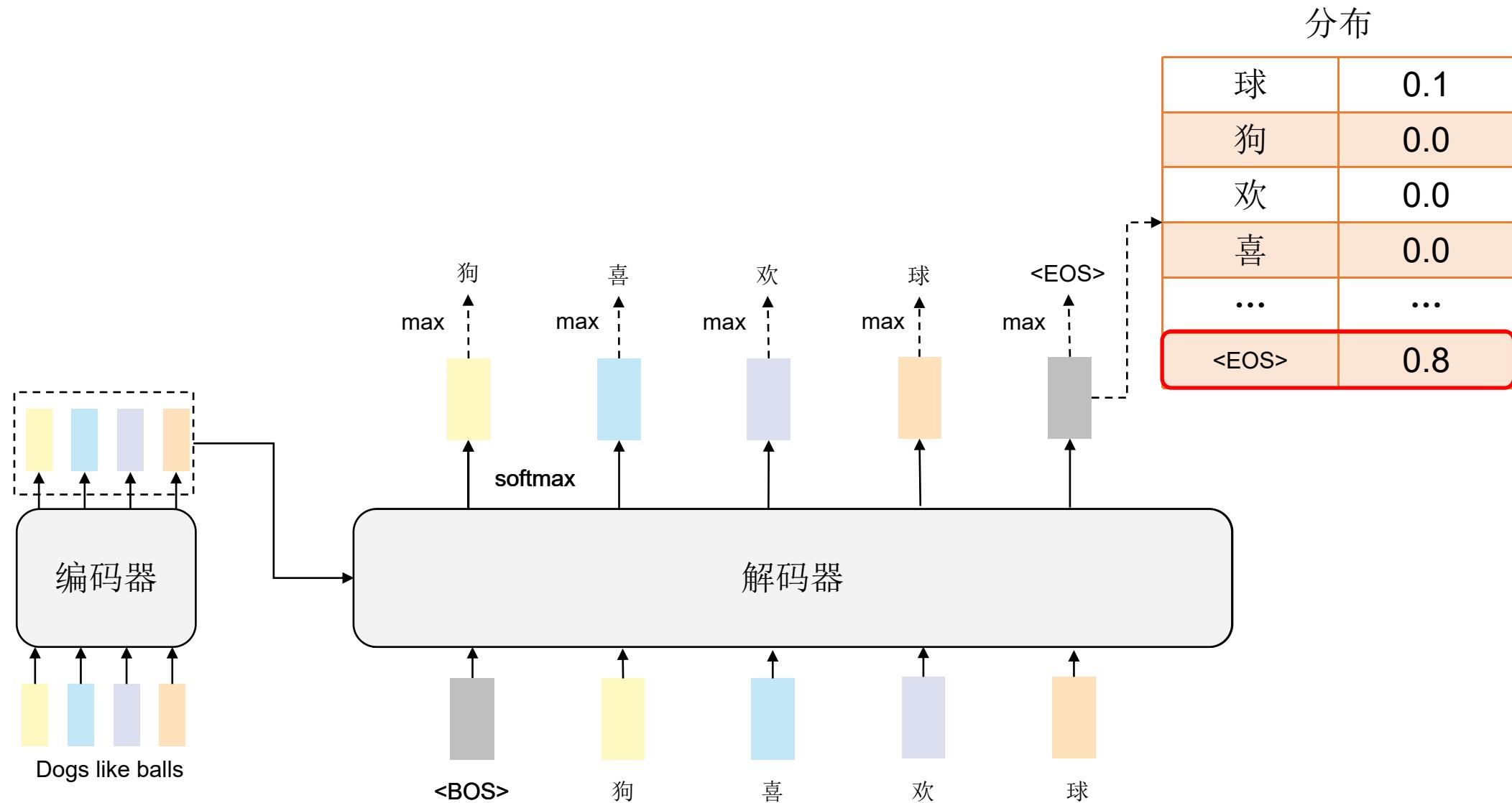
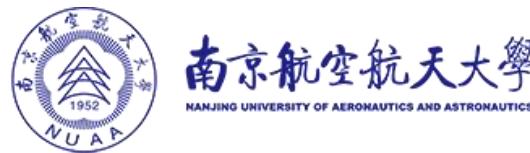
编码器-解码器结构



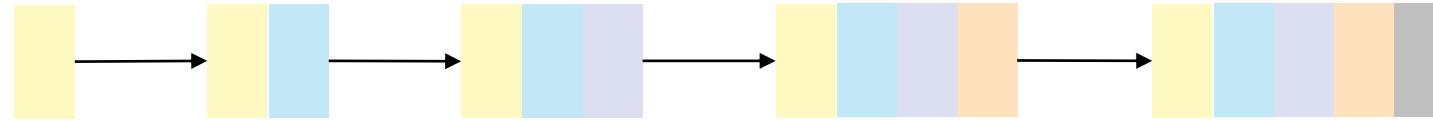
- 解码器需要自己决定生成的目标句的长度
- 实际上：机器并不能确定正确的输出长度



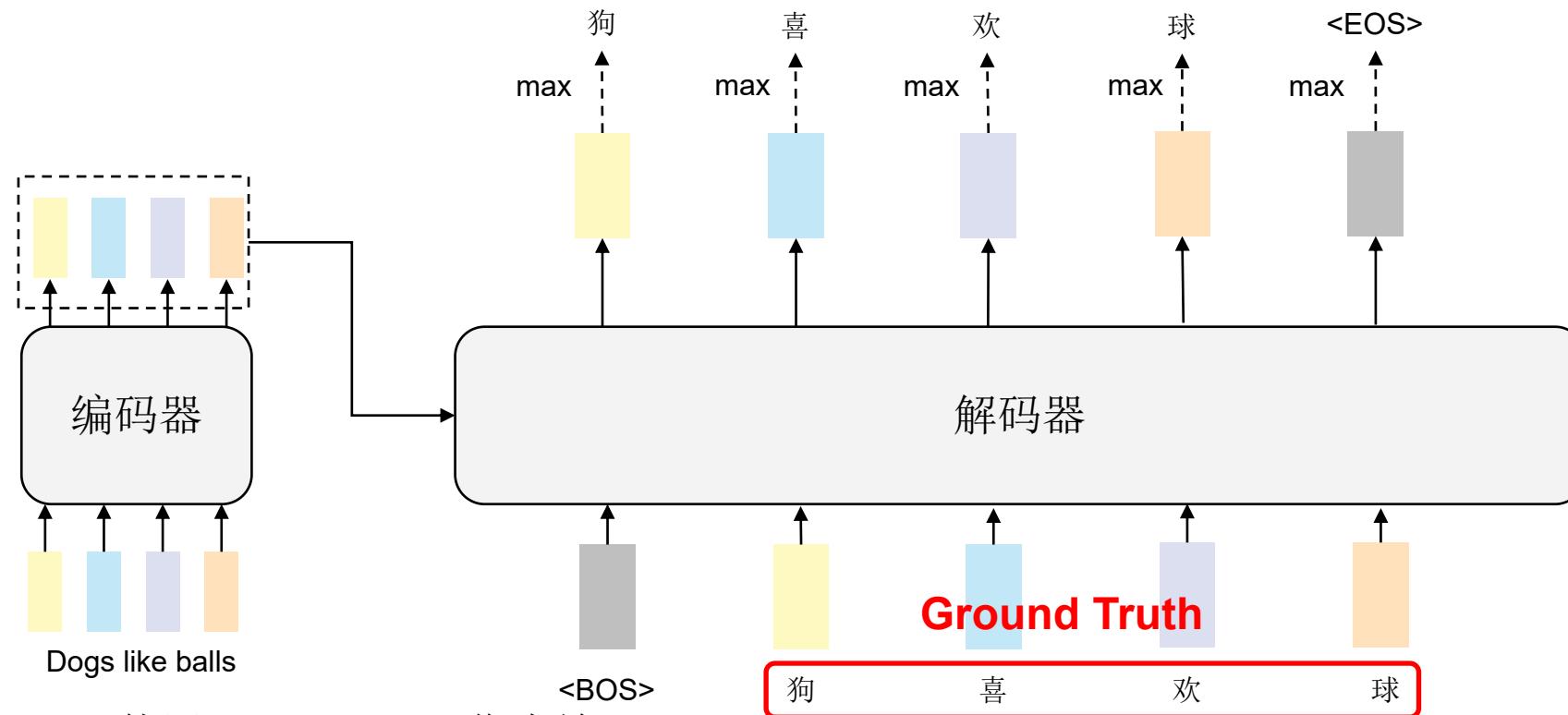
编码器-解码器结构



编码器-解码器结构

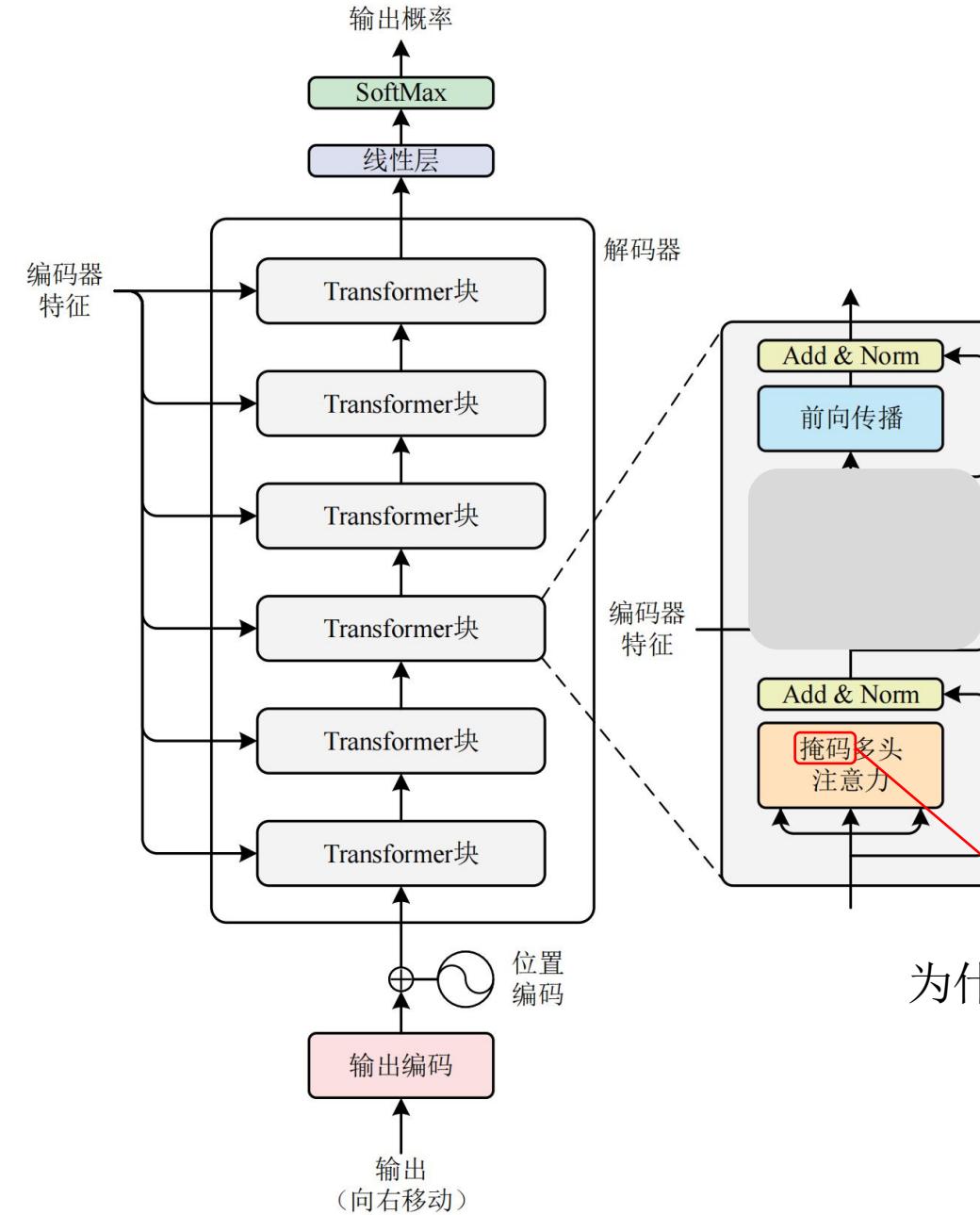
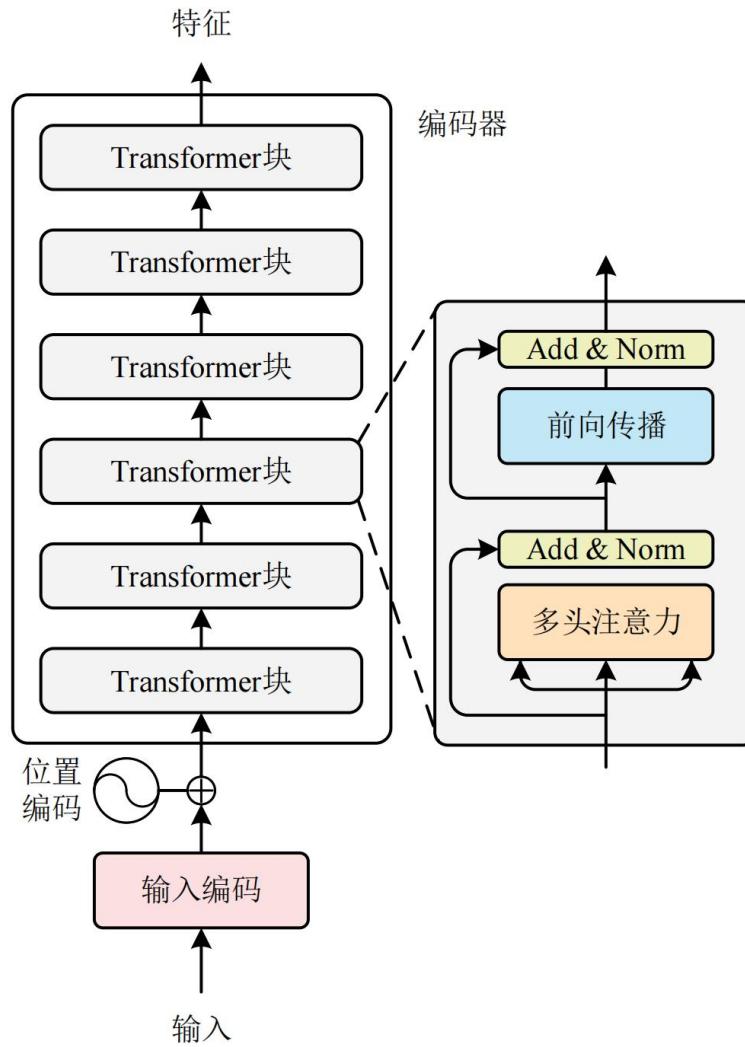


在训练的过程中，不仅降低模型并行程度，使得训练时间变长，而且训练更加困难



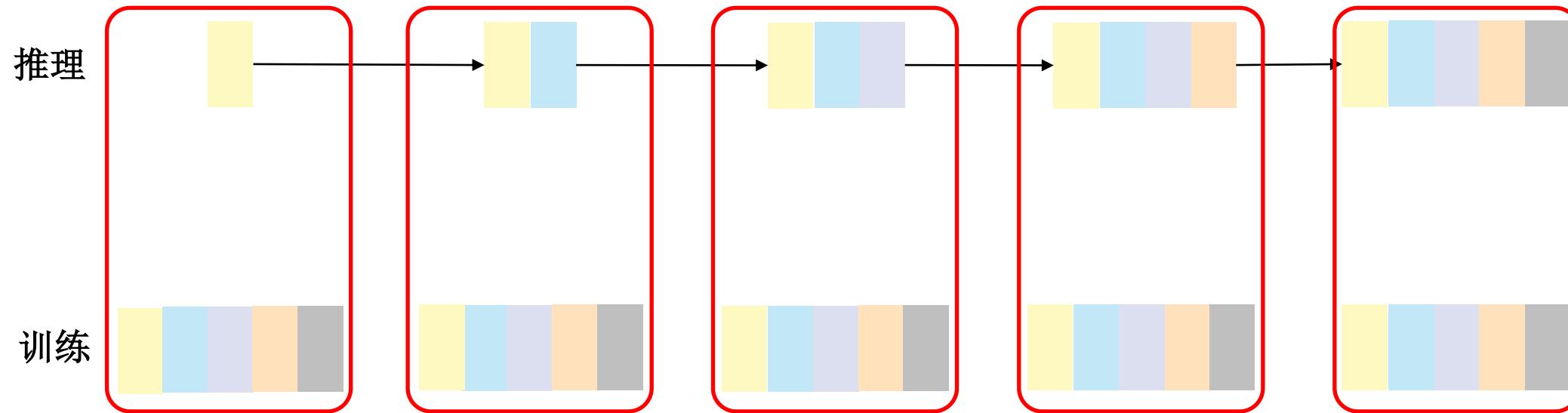
Teacher Forcing: 使用Ground Truth作为输入

编码器-解码器结构



为什么需要掩码?

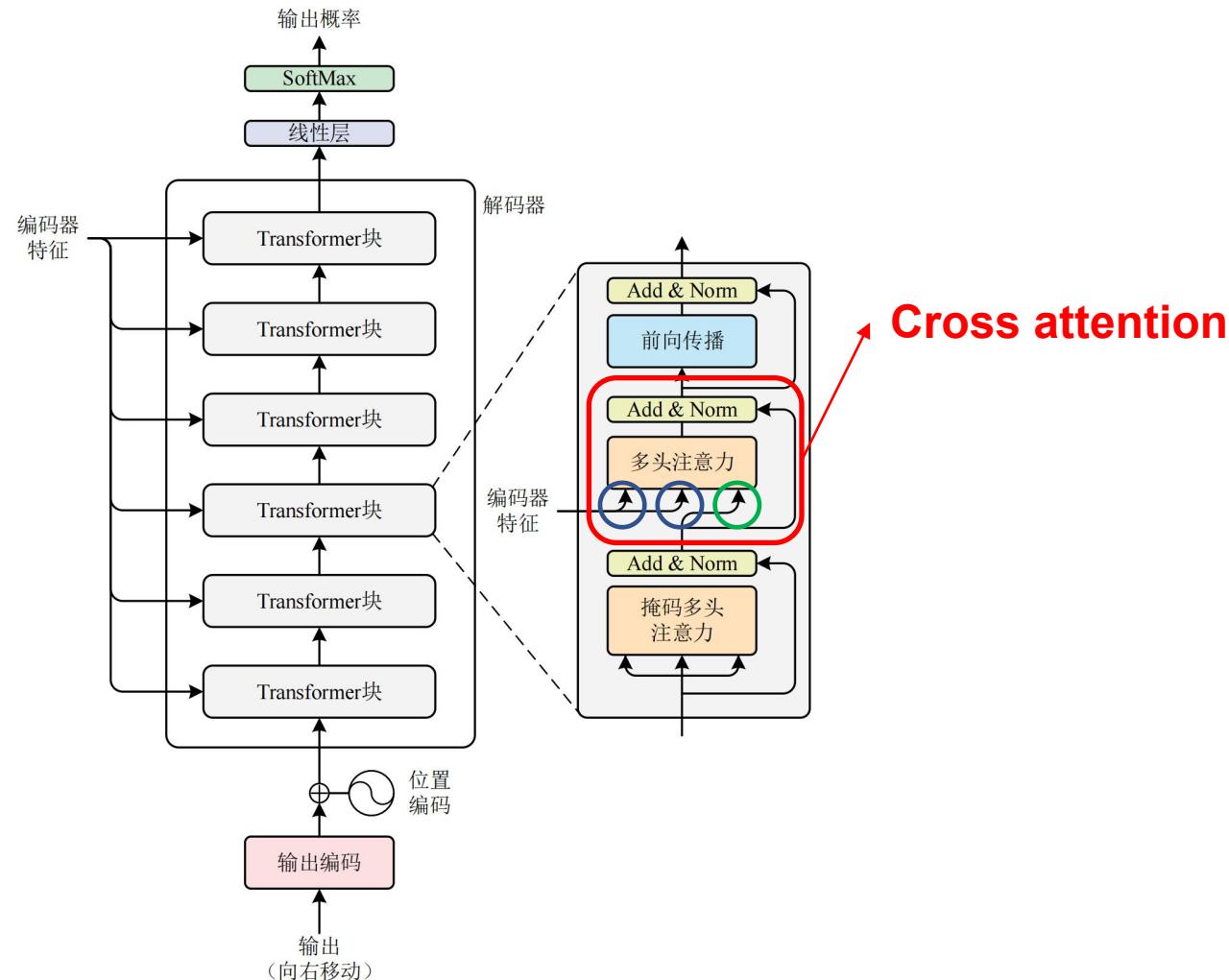
编码器-解码器结构



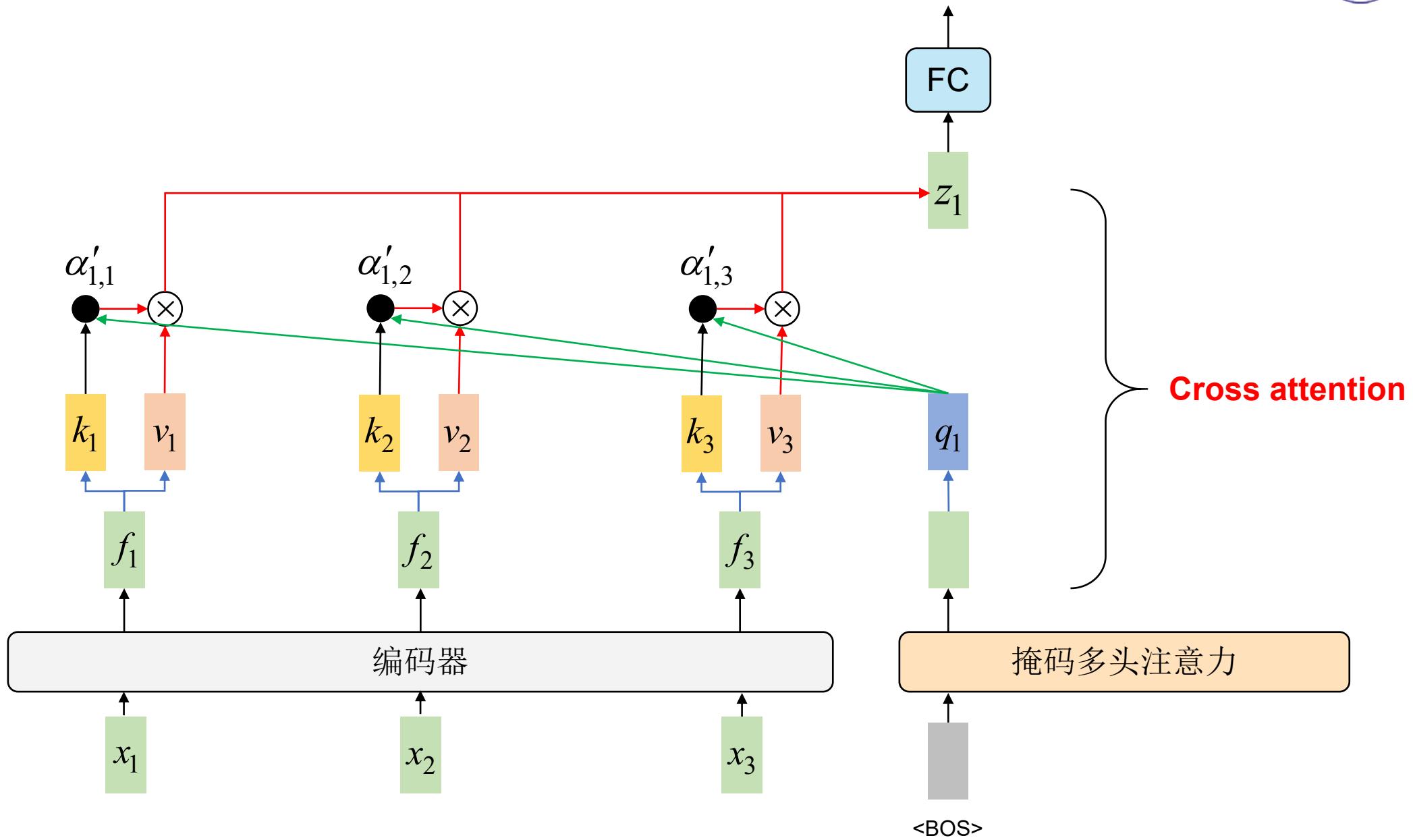
在训练的过程中，使用掩码多头注意力模块

编码器-解码器结构

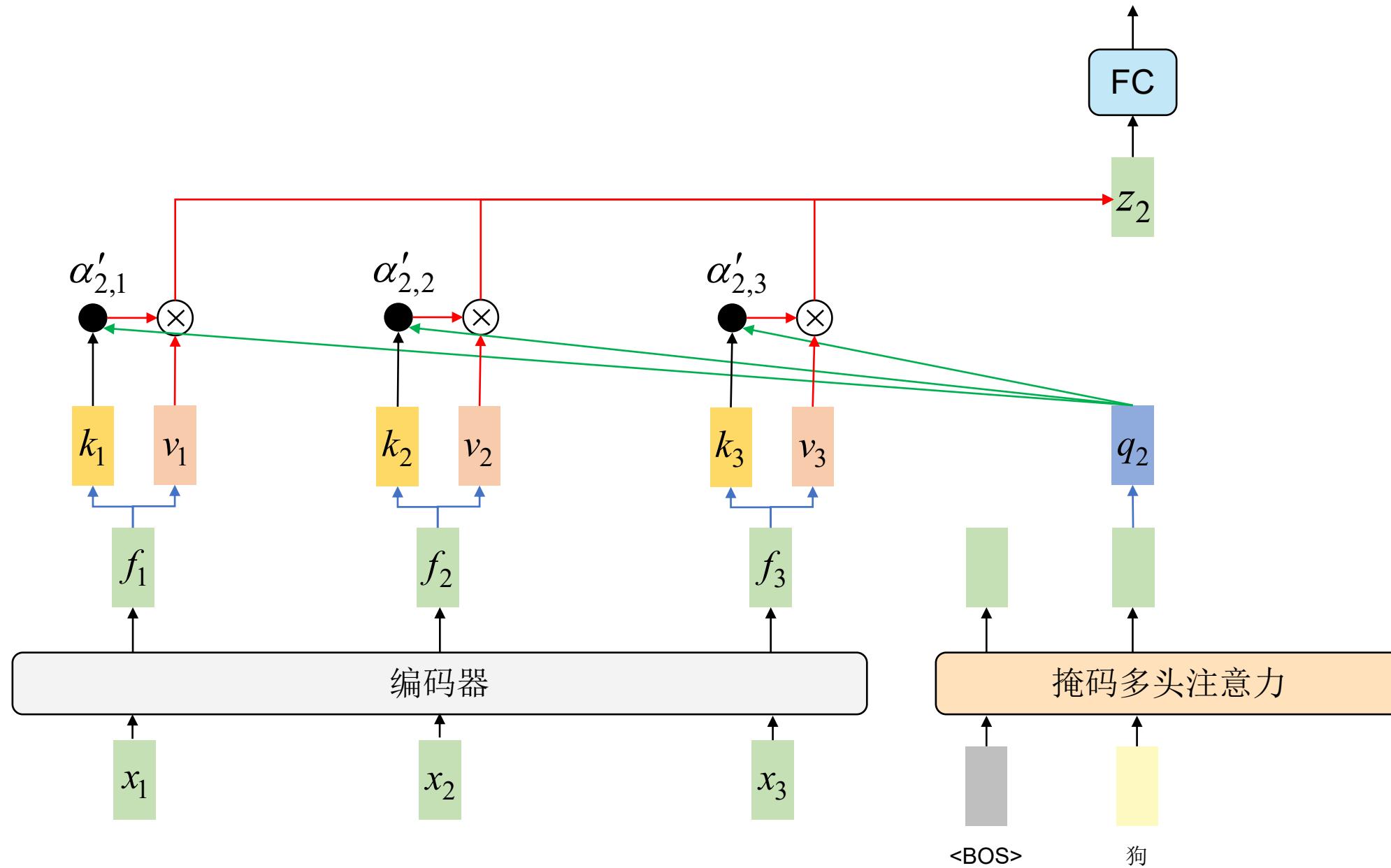
3. 编码器与解码器之间的信息传递



编码器-解码器结构



编码器-解码器结构



目录



□ Transformer

- 注意力机制
- 编码器-解码器结构
- 大模型中的编码器-解码器结构

□ 编码器结构----BERT家族

- BERT结构
- 预训练策略
- BERT的变体

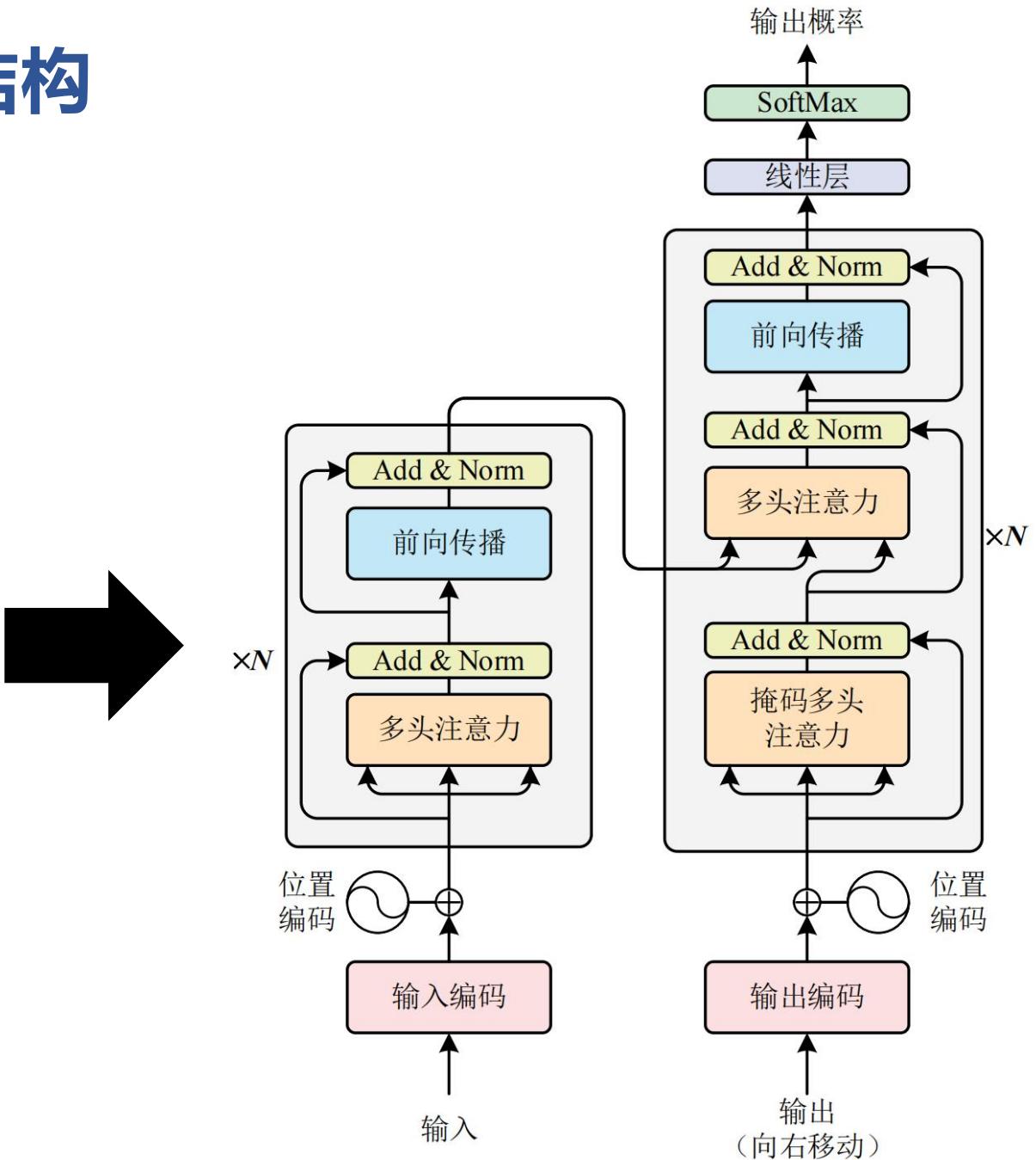
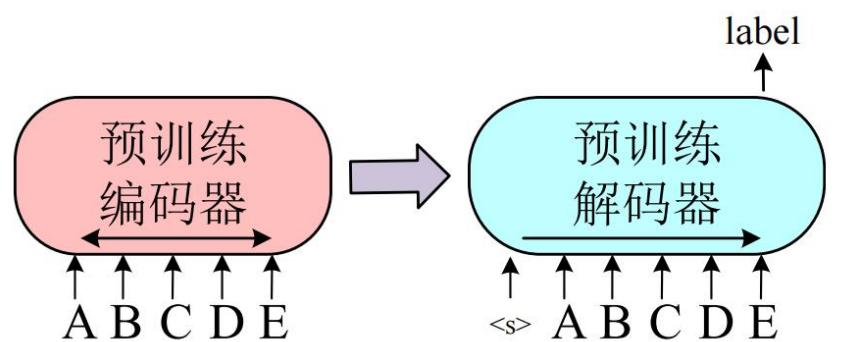
□ 解码器结构----GPT家族

- GPT结构
- 自回归预训练
- 后续改进

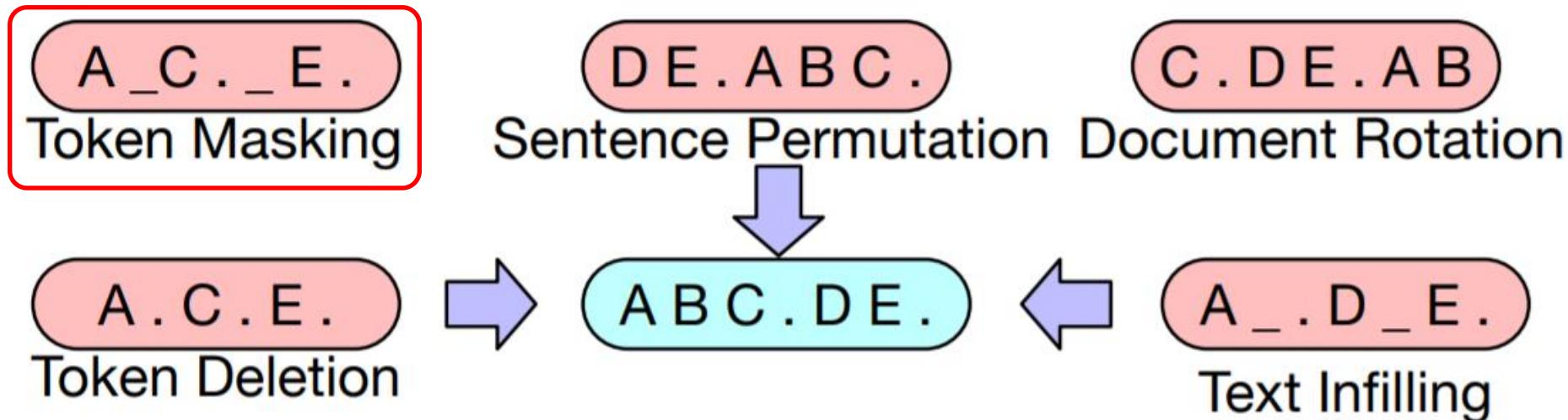
□ 思考

大模型中的编码器-解码器结构

1. BART



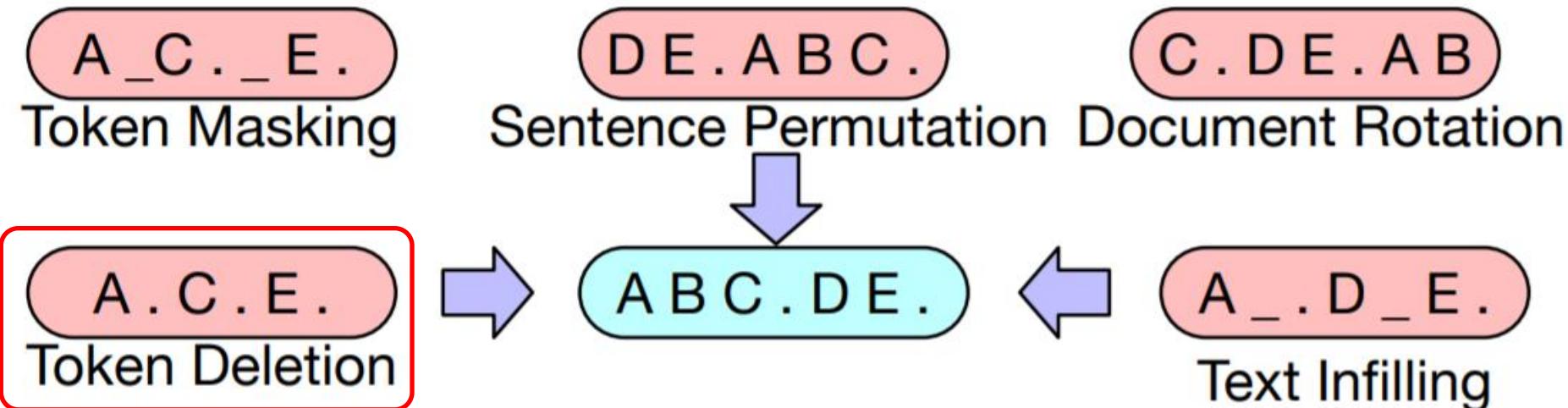
大模型中的编码器-解码器结构



Token Masking

A [MASK] C . [MASK] E .
↓ ↓
A B C . D E .

大模型中的编码器-解码器结构

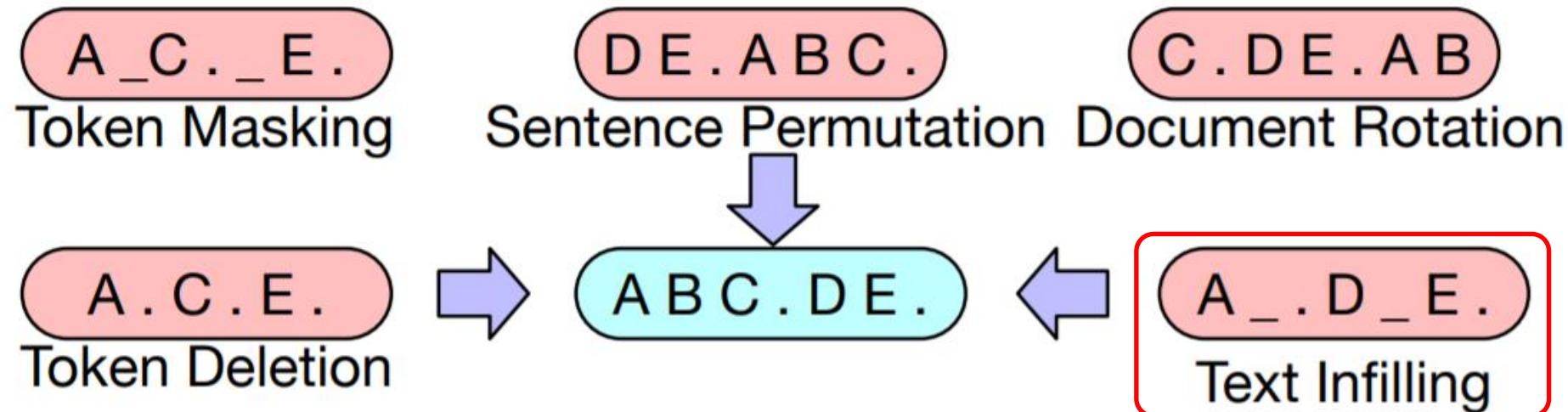
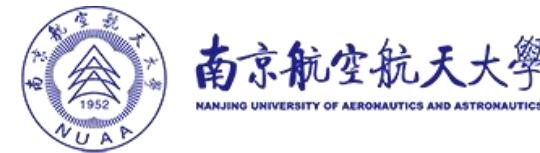


Token Masking

Token Deletion

A C . E .
↓ ↓
A $\not\sim$ C . $\not\sim$ E .

大模型中的编码器-解码器结构



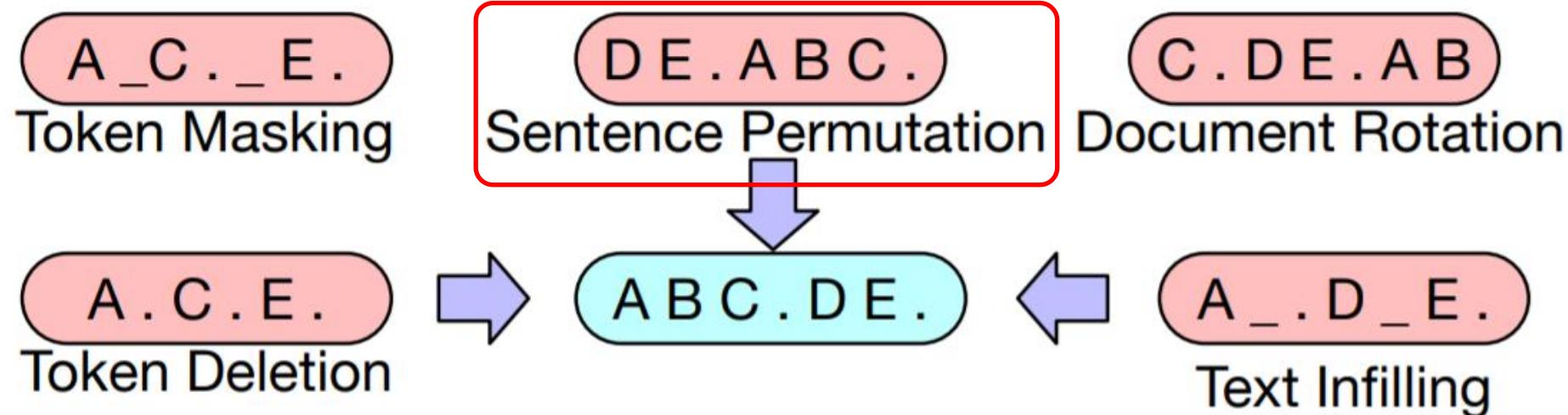
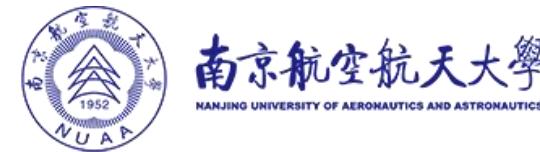
Token Masking

Token Deletion

Token Infilling

A [MASK] . D [MASK] E .
A 2 . D 0 E .

大模型中的编码器-解码器结构



Token Masking

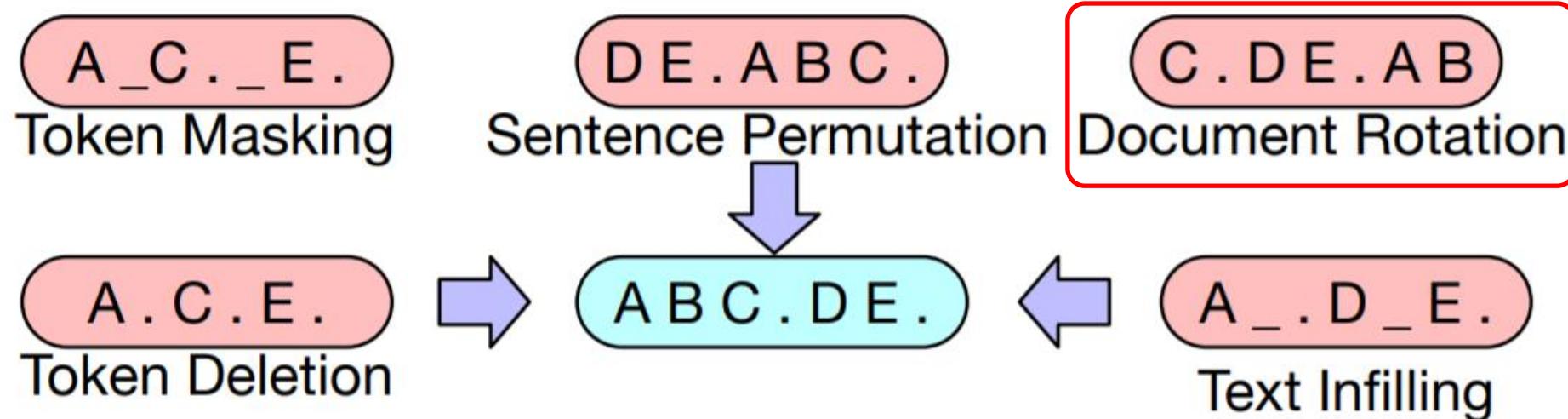
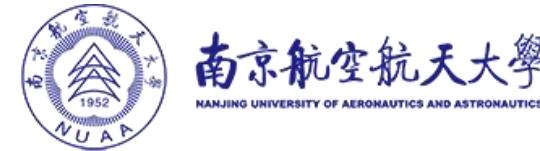
Token Deletion

Token Infilling

Sentence Permutation

D E . A B C .
↓ ↓
2 1

大模型中的编码器-解码器结构



Token Masking

Token Deletion

Token Infilling

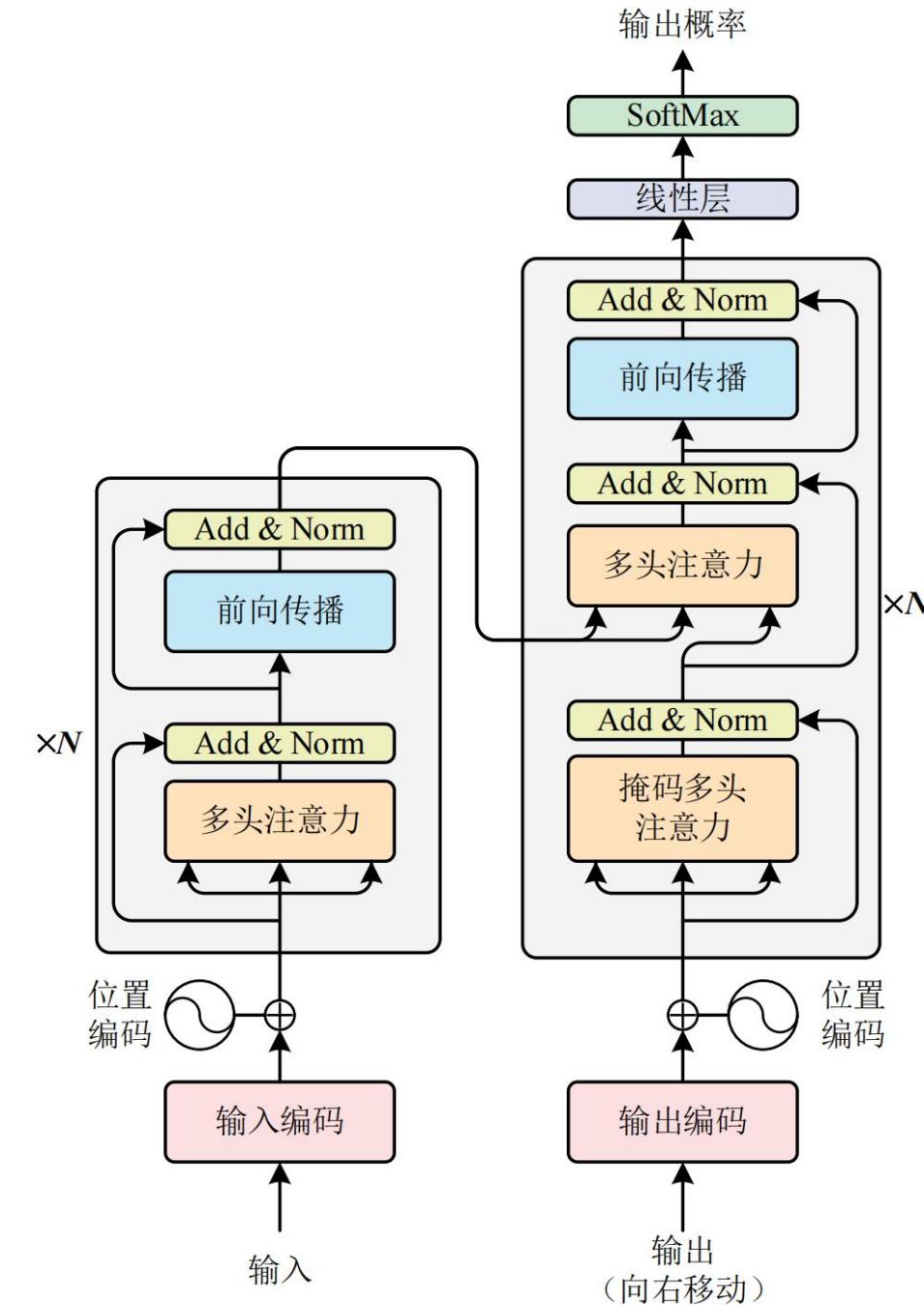
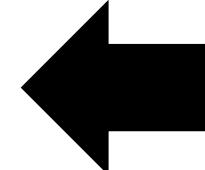
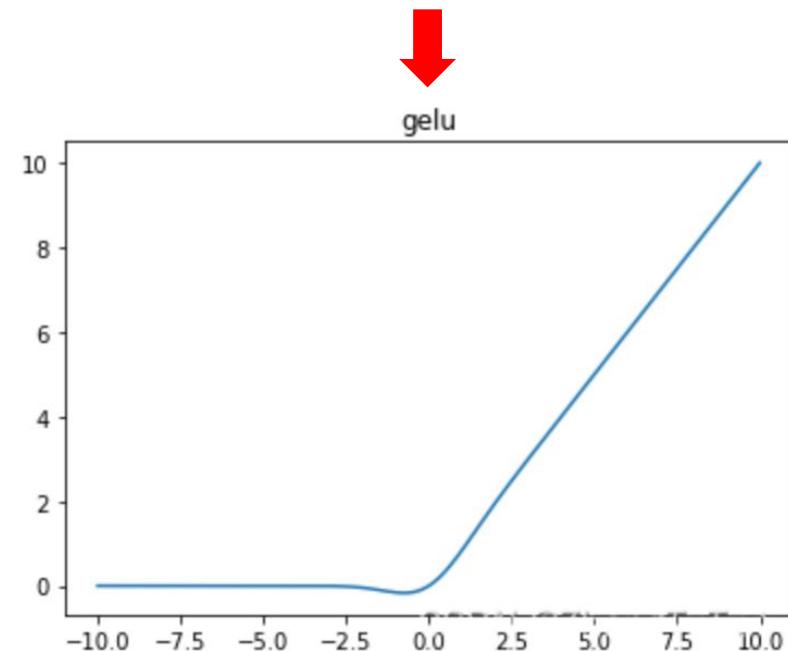
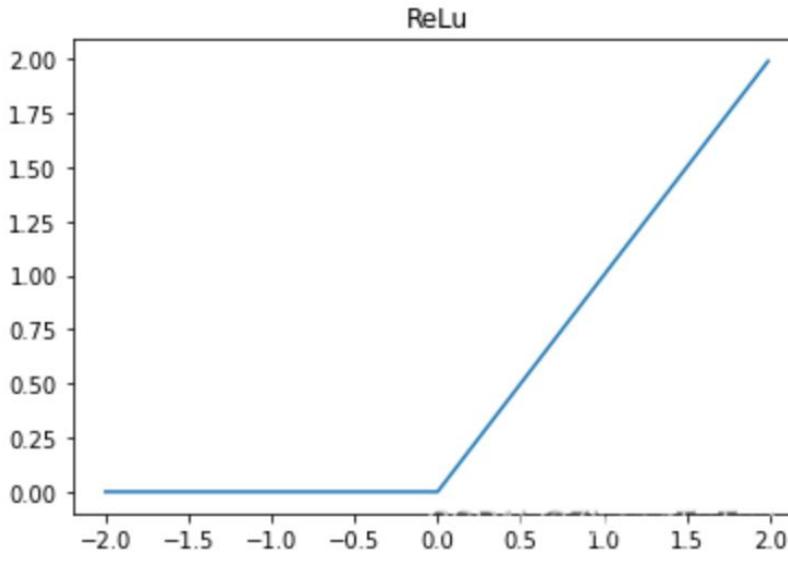
Sentence Permutation

Document Rotation

C . D E . A B

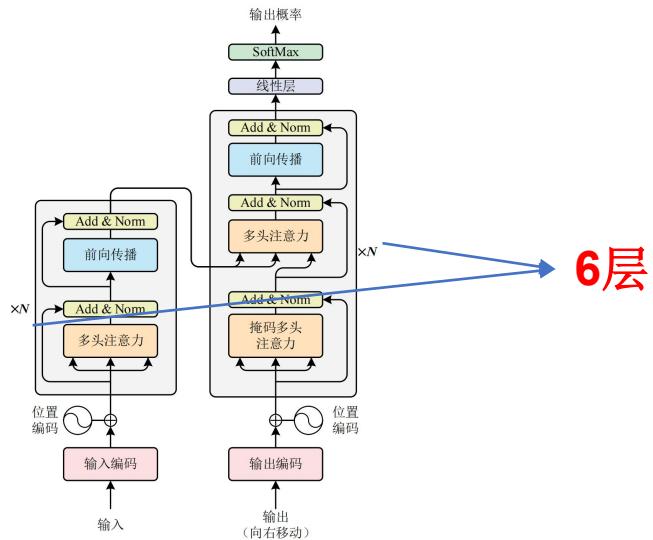
↓
Begin

大模型中的编码器-解码器结构



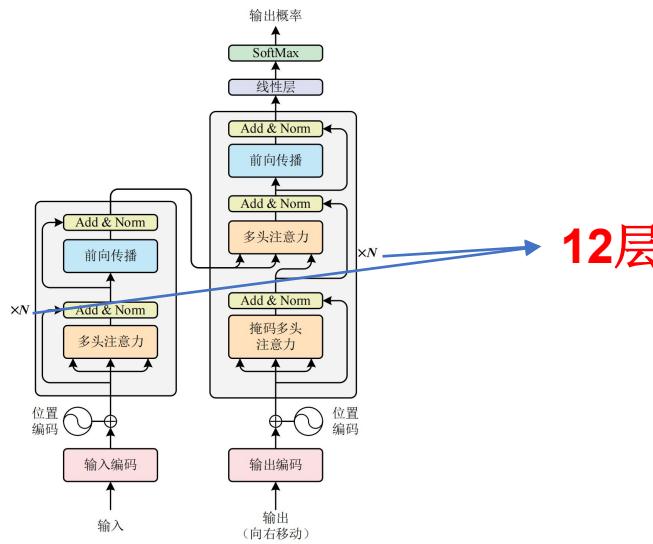
大模型中的编码器-解码器结构

Base

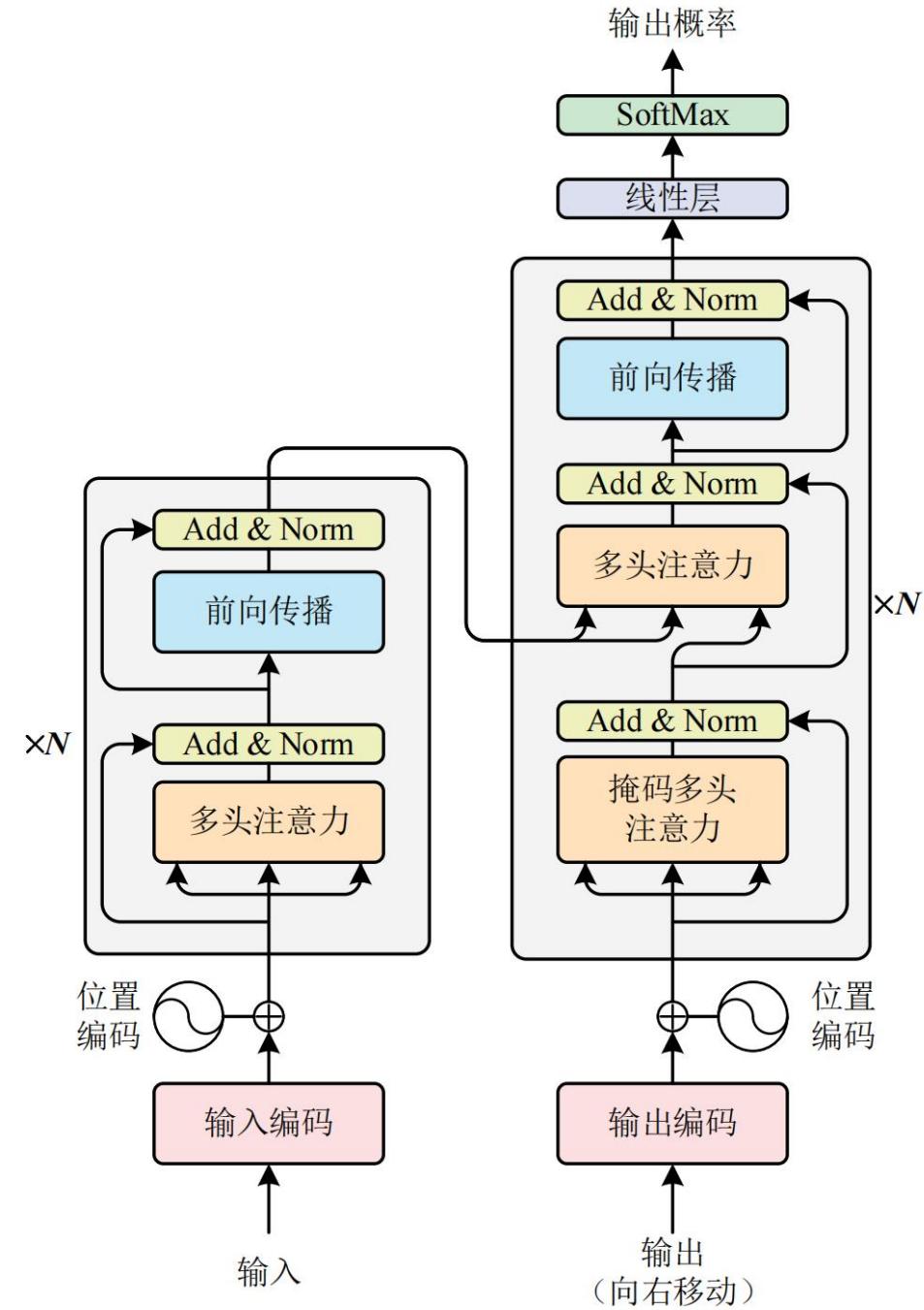
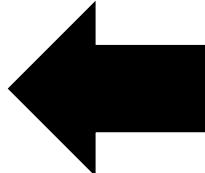


6层

Large

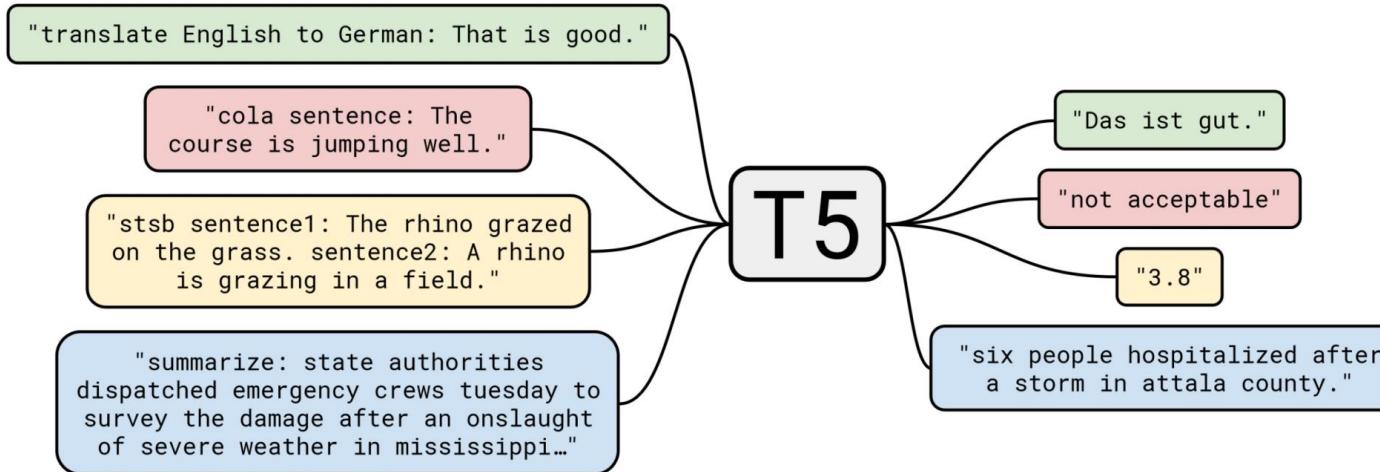


12层

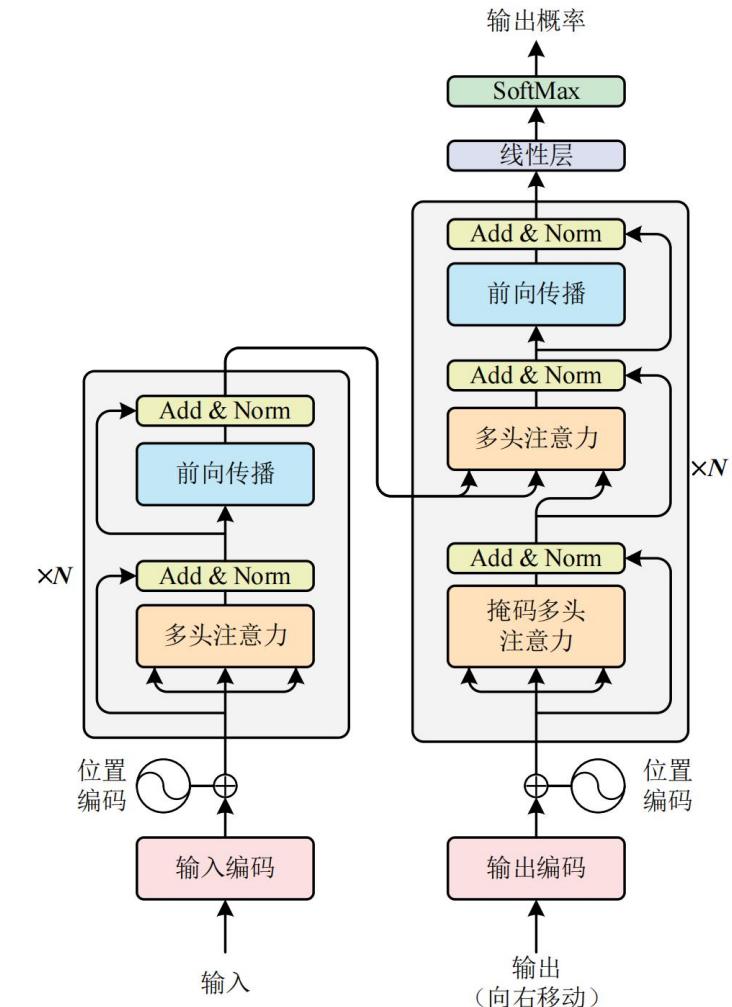


大模型中的编码器-解码器结构

1.T5



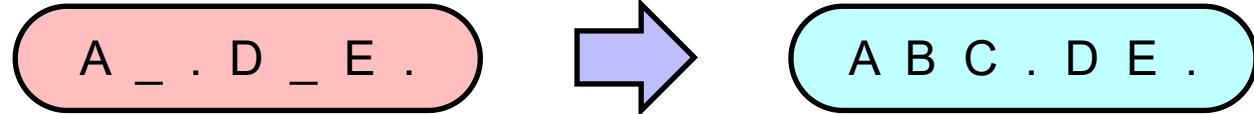
将所有文本处理问题转化为“文本到文本”的问题



大模型中的编码器-解码器结构



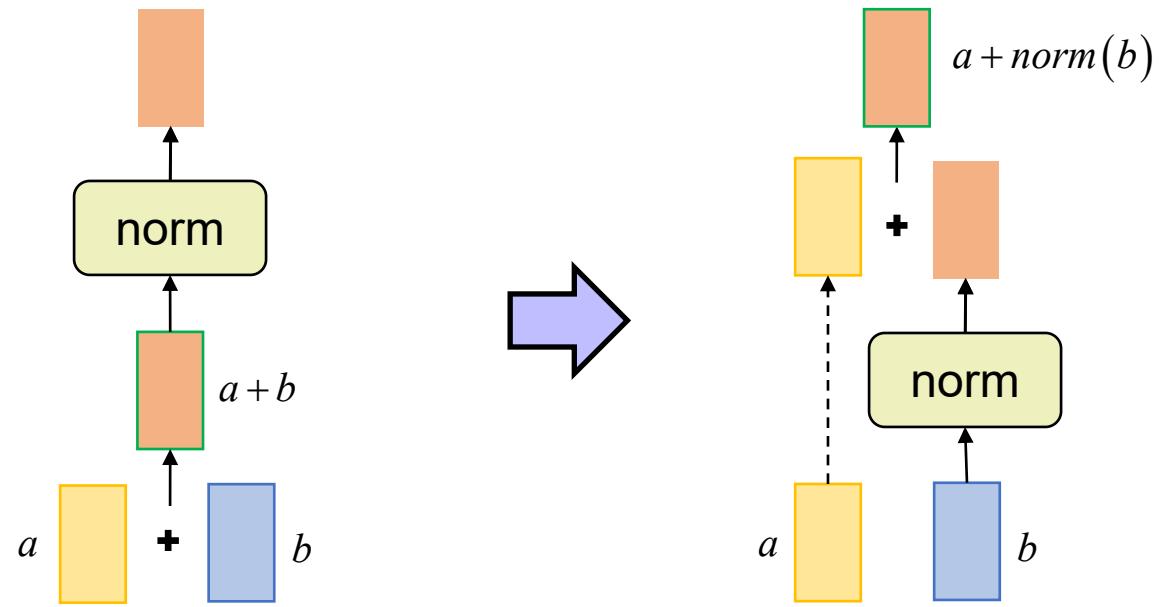
Replace Span



层归一化

$$\text{LayerNorm}(x) = \frac{x - \mu}{\sqrt{\sigma^2 + \varepsilon}} * \gamma + \beta$$

残差连接

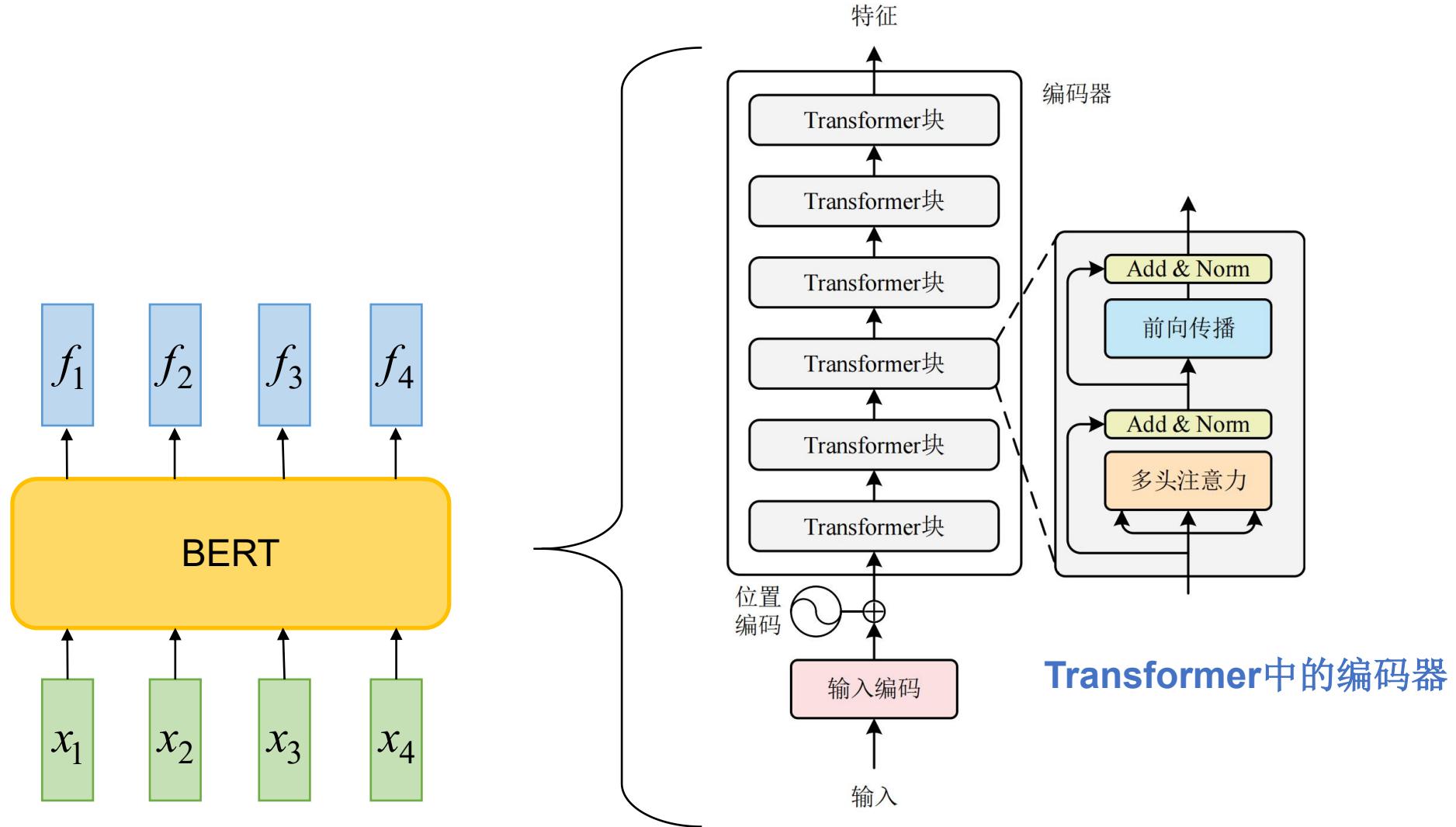


目录

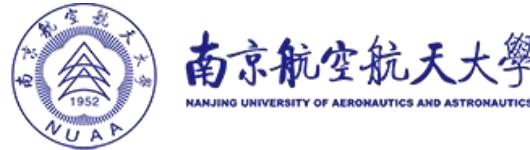


- Transformer
 - 注意力机制
 - 编码器-解码器结构
 - 大模型中的编码器-解码器结构
- 编码器结构----BERT家族
 - BERT结构
 - 预训练策略
 - BERT的变体
- 解码器结构----GPT家族
 - GPT结构
 - 自回归预训练
 - 后续改进
- 思考

编码器结构-BERT家族



编码器结构-BERT家族



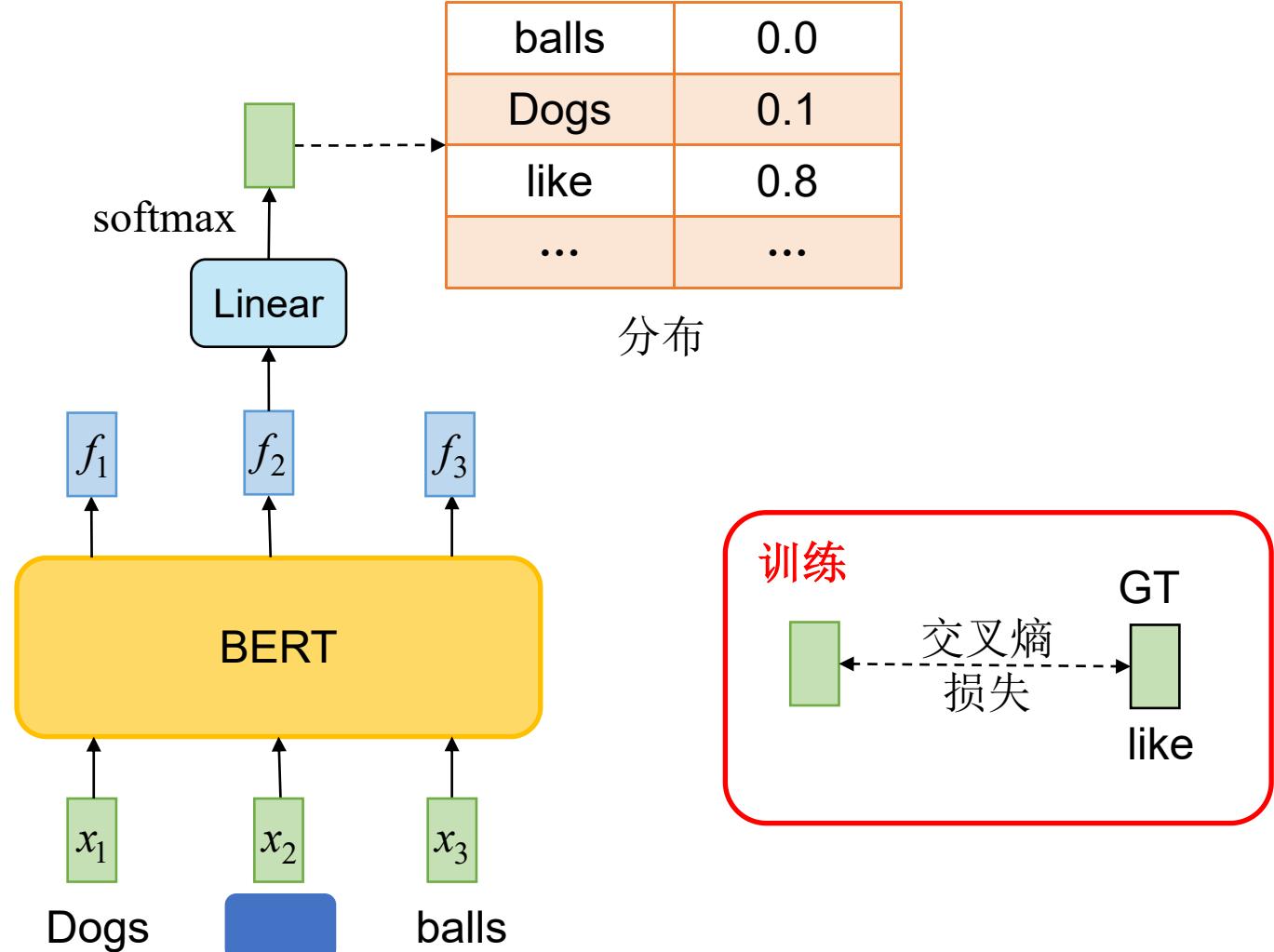
2. 预训练策略

自监督学习策略的提出使得BERT的训练成为可能

“掩码语言建模”

= “[MASK]”

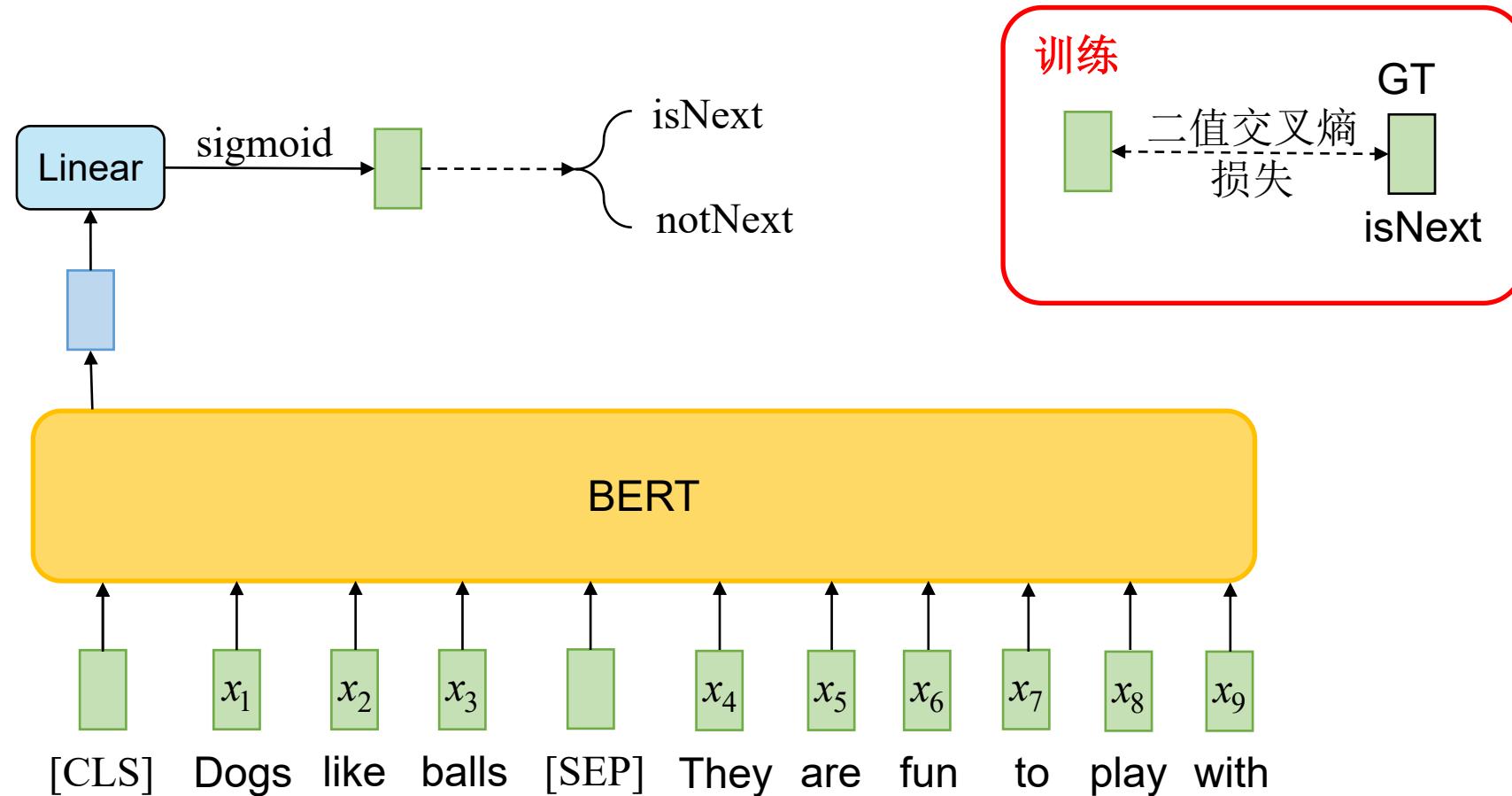
随机掩盖其中一些token



编码器结构-BERT家族



“下句预测”

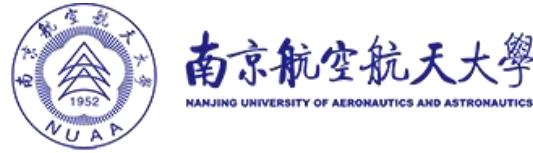


目录

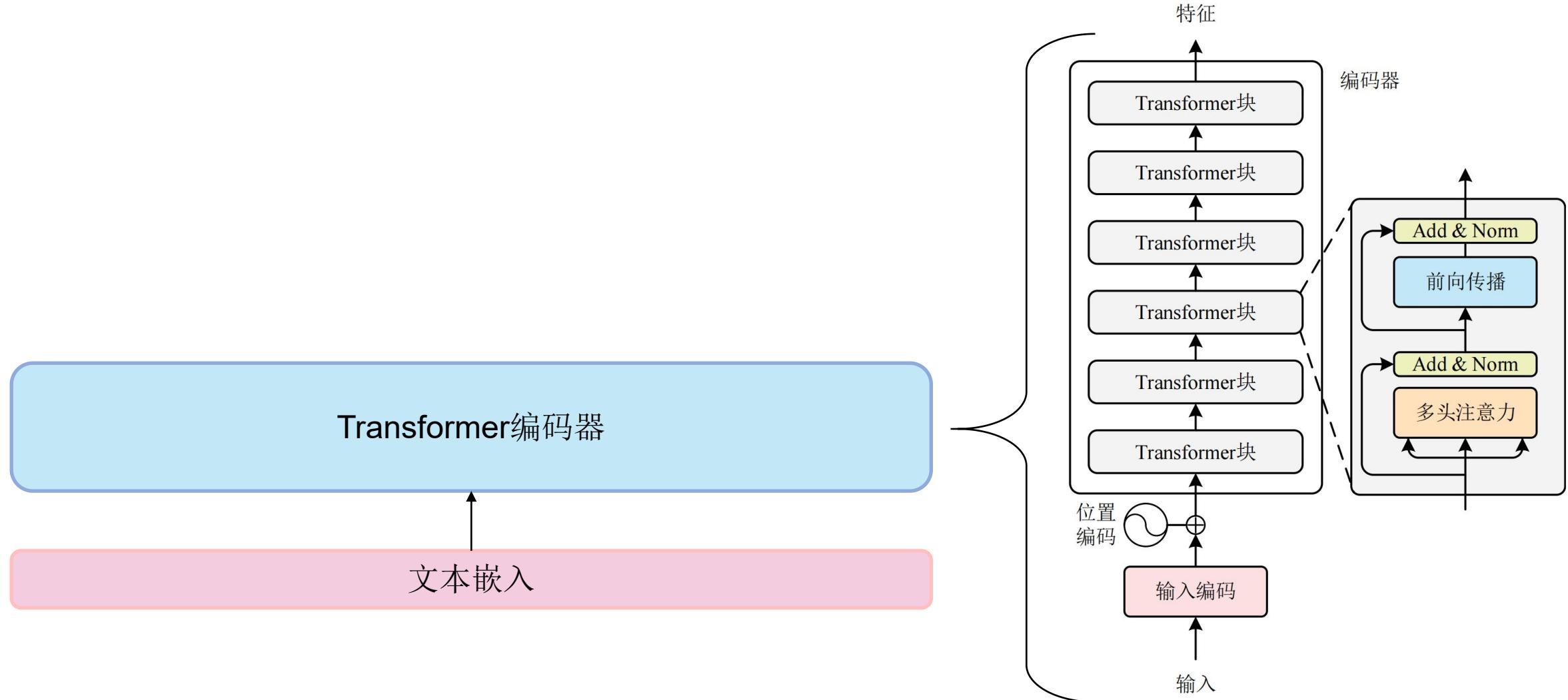


- Transformer
 - 注意力机制
 - 编码器-解码器结构
 - 大模型中的编码器-解码器结构
- 编码器结构----BERT家族
 - BERT结构
 - 预训练策略
 - BERT的变体
- 解码器结构----GPT家族
 - GPT结构
 - 自回归预训练
 - 后续改进
- 思考

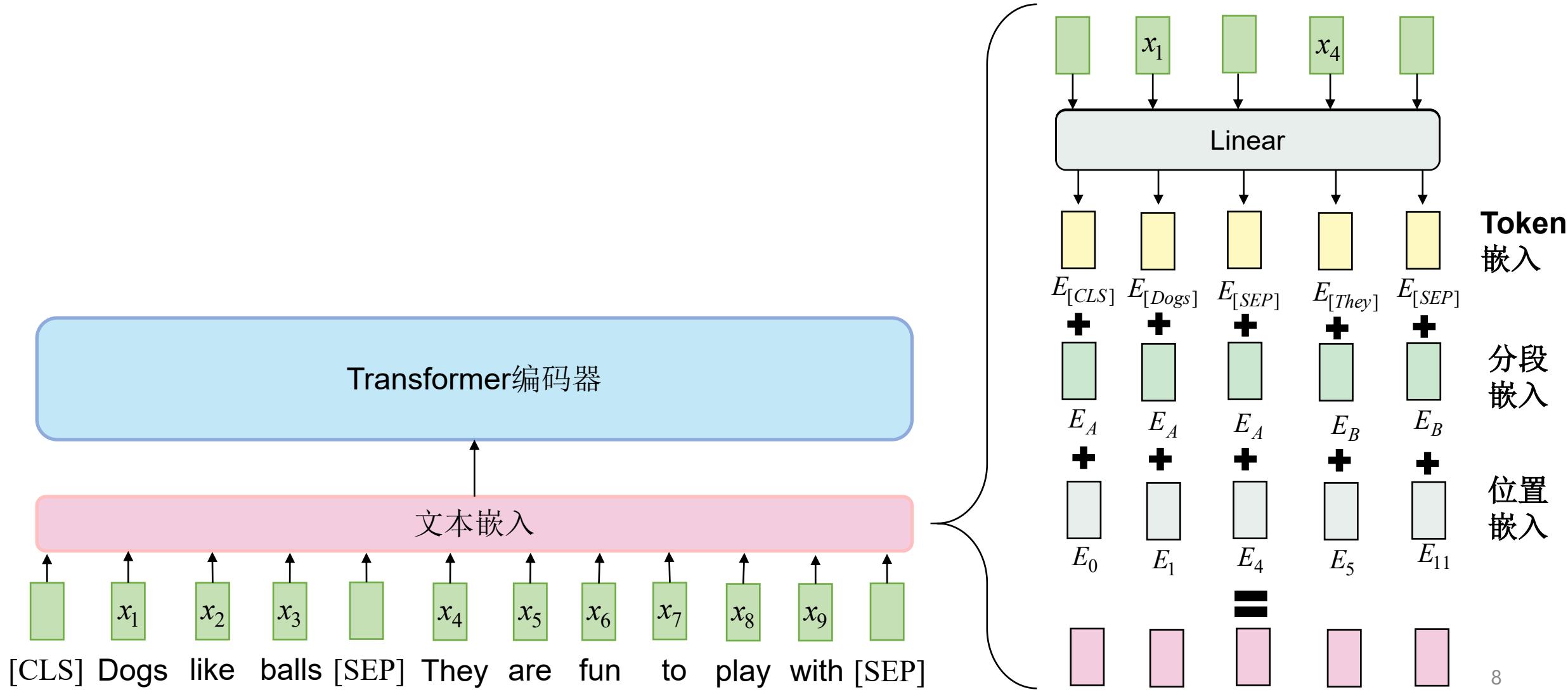
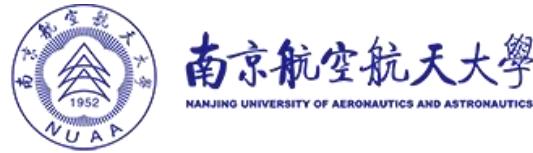
编码器结构-BERT家族



1. BERT结构



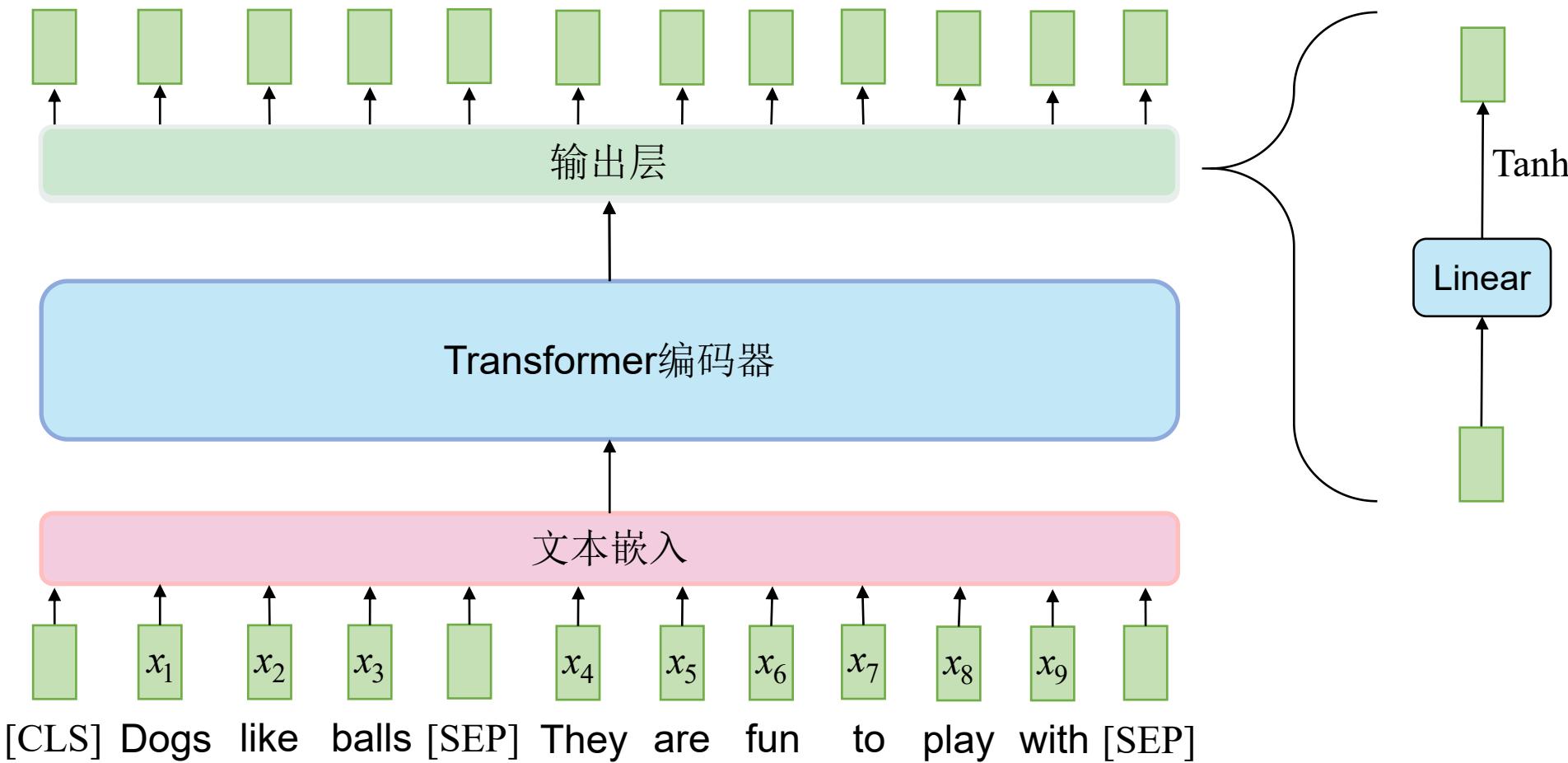
编码器结构-BERT家族



编码器结构-BERT家族



1. BERT结构



目录

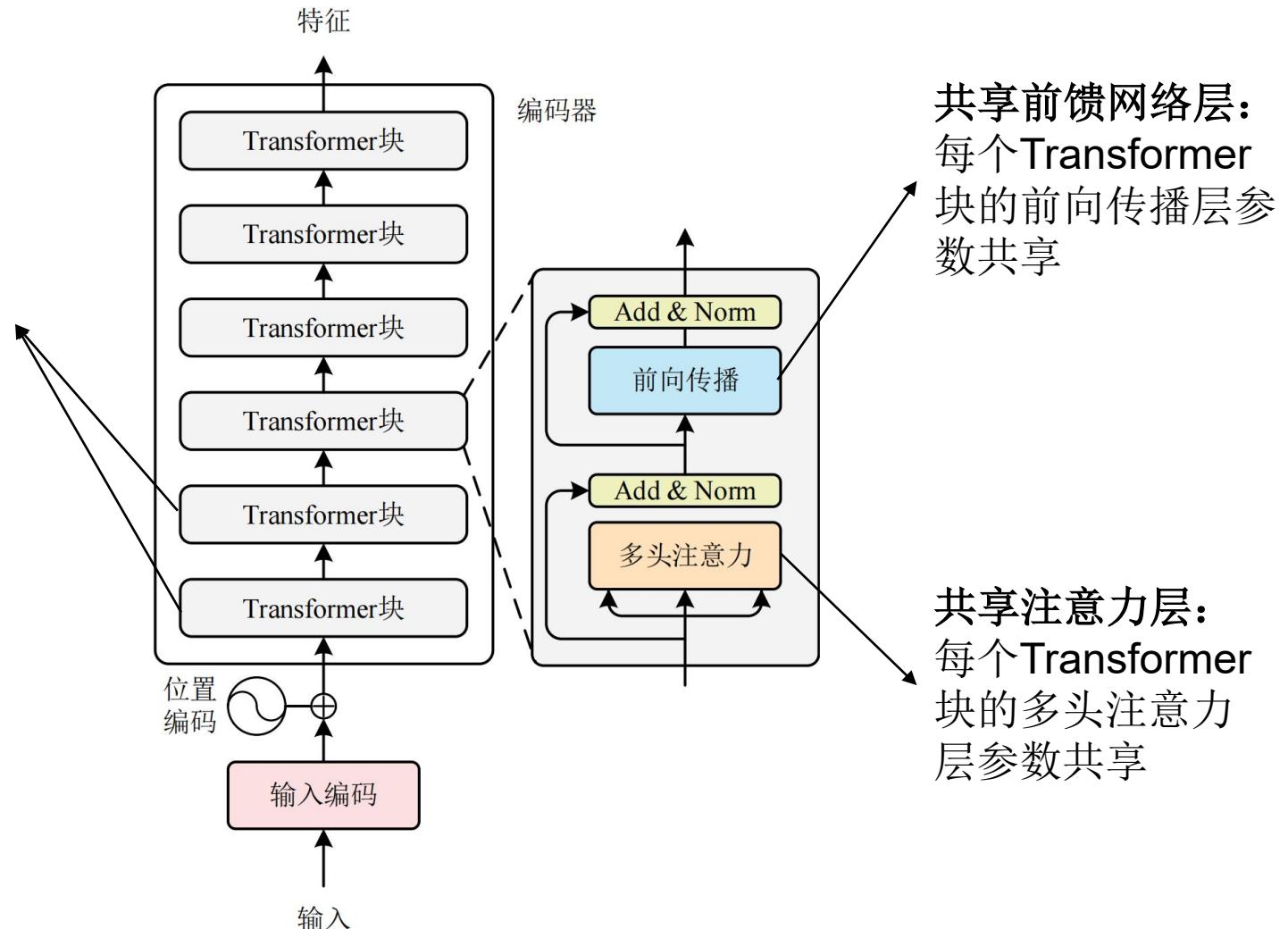


- Transformer
 - 注意力机制
 - 编码器-解码器结构
 - 大模型中的编码器-解码器结构
- 编码器结构----BERT家族
 - BERT结构
 - 预训练策略
 - BERT的变体
- 解码器结构----GPT家族
 - GPT结构
 - 自回归预训练
 - 后续改进
- 思考

3.BERT的变体

BERT的参数量太大--ALBERT

全共享：
Transformer
块的参数共享
“跨层参数共享”



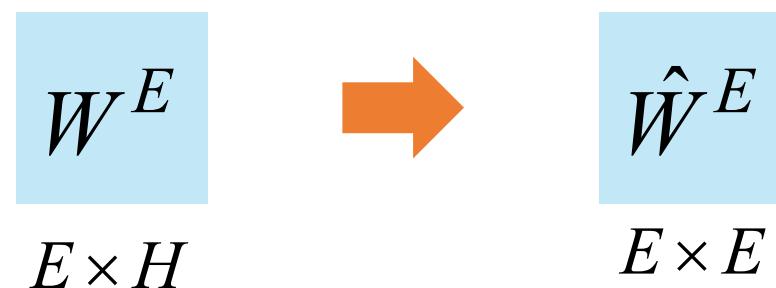
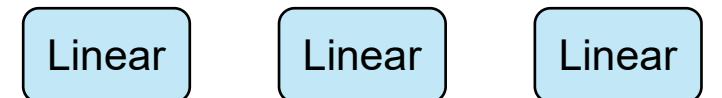


3.BERT的变体

BERT的参数量太大--ALBERT

文本嵌入

“嵌入层参数因子分解”



编码器结构-BERT家族



模型	参数量	层数 N	隐藏神经元数量	嵌入层
BERT-base	1.1 亿	12	768	768
BERT-large	3.4 亿	24	1024	1024
ALBERT-base	1200 万	12	768	128
ALBERT-large	1800 万	24	1024	128
ALBERT-xlarge	6000 万	24	2048	128
ALBERT-xxlarge	2.35 亿	12	4096	128

RoBERTa: 改进BERT预训练

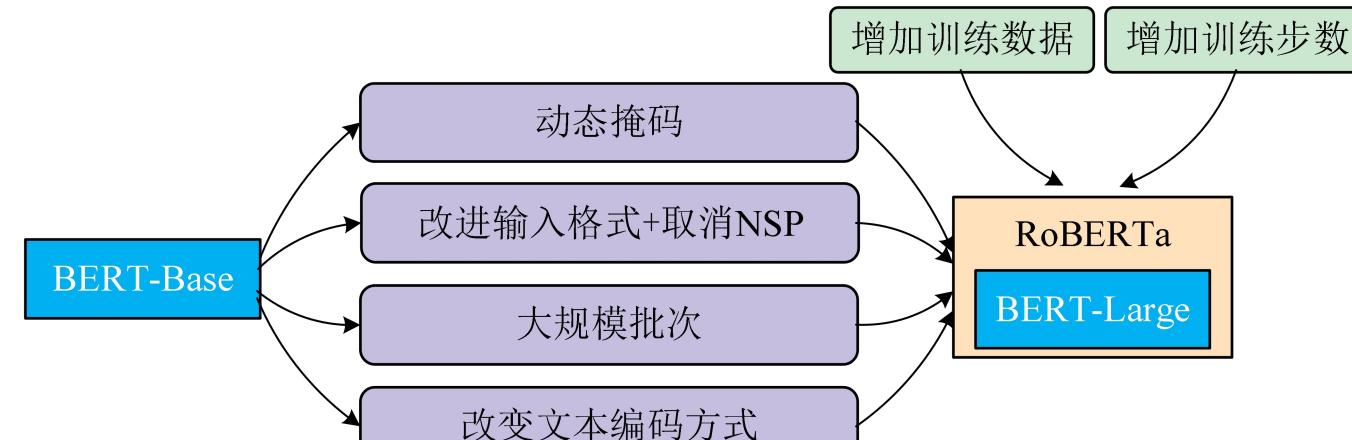
动态掩码

移除下句预测任务

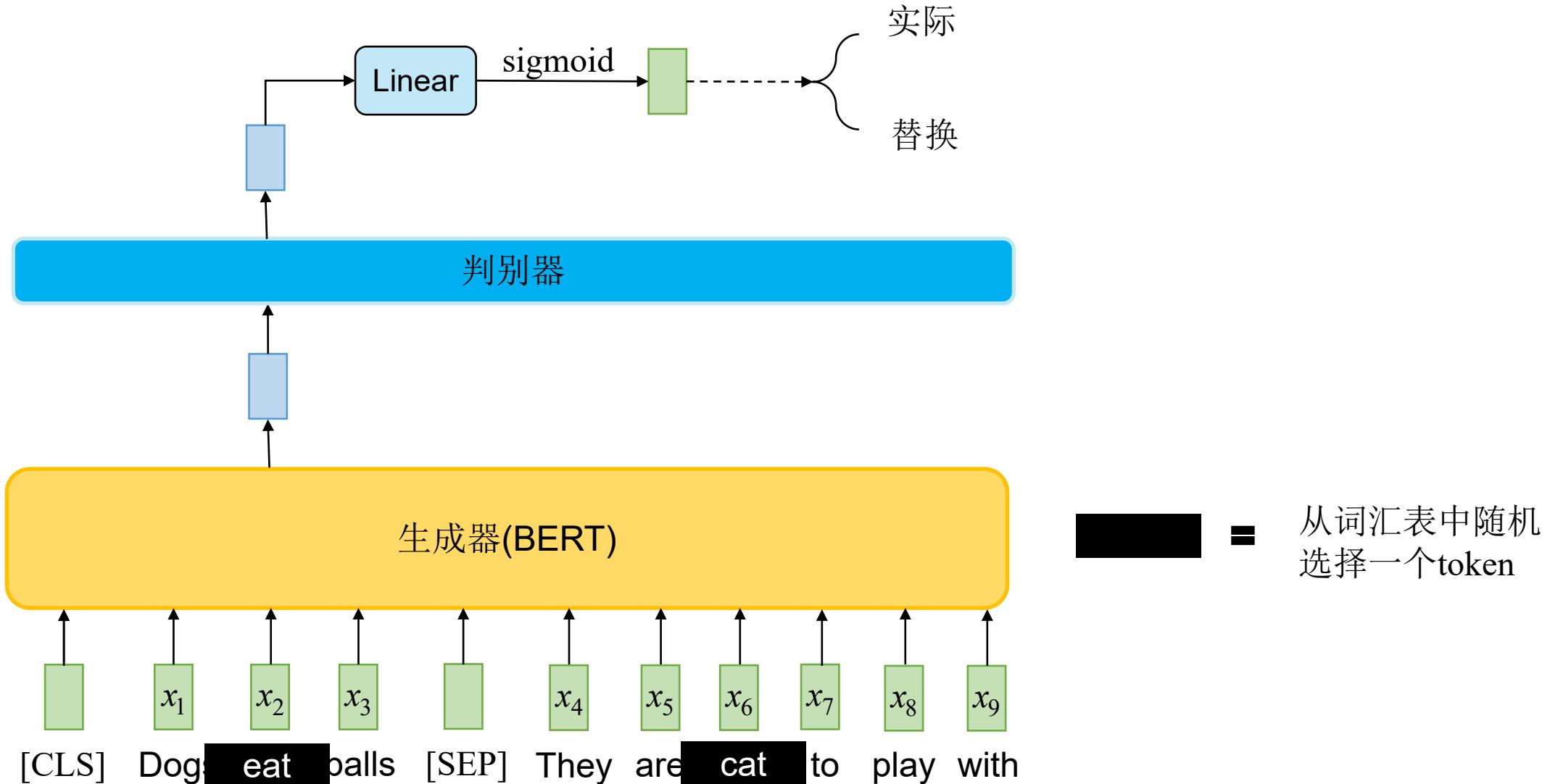
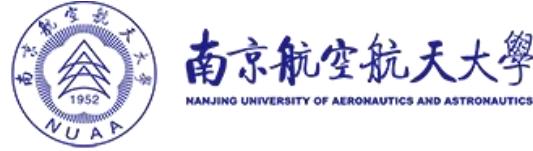
增加数据和训练步数

增大batch size

使用字节级字节对编码



编码器结构-BERT家族



目录



□ Transformer

- 注意力机制
- 编码器-解码器结构
- 大模型中的编码器-解码器结构

□ 编码器结构----BERT家族

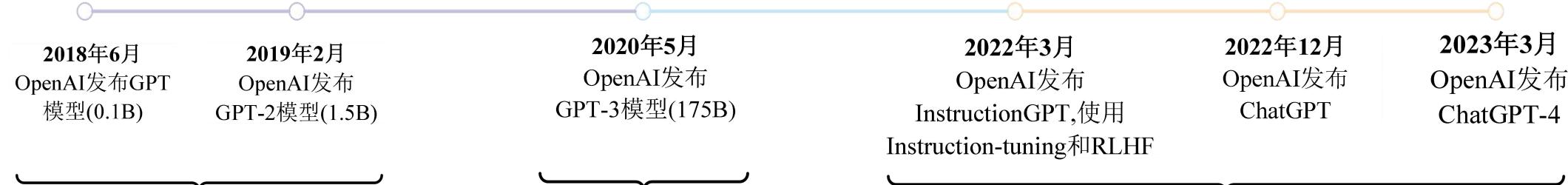
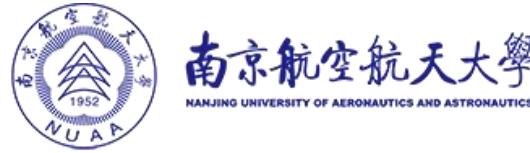
- BERT结构
- 预训练策略
- BERT的变体

□ 解码器结构----GPT家族

- GPT结构
- 自回归预训练
- 后续改进

□ 思考

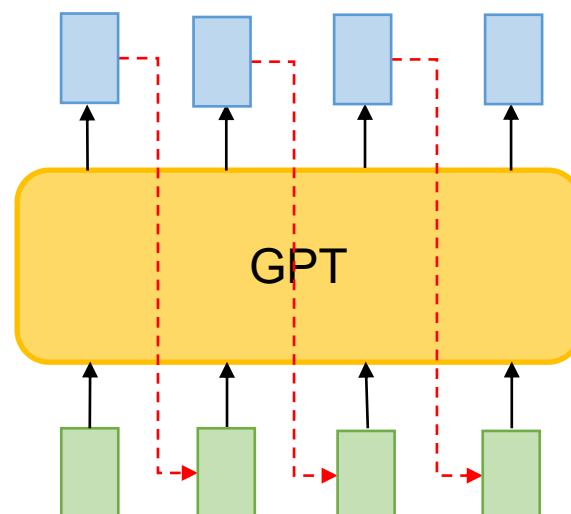
解码器结构-GPT家族



第一阶段:
提出新颖的自监督训练目标
设计模型架构
针对某部分领域数据单独Fine-tuning

第二阶段:
逐步扩大模型参数和训练语料规模
模型架构偏向生成式任务
Prompt技术在Fine-tuning阶段得以展现

第三阶段:
进一步增大模型参数量和训练语料规模
完全采用单向因果模型
注重人类反馈信息、对话交互
关注于模型可靠性、事实性、有偏性、安全性
更多丰富的Prompt技术用于Fine-tuning或Inference



Transformer解码器样式

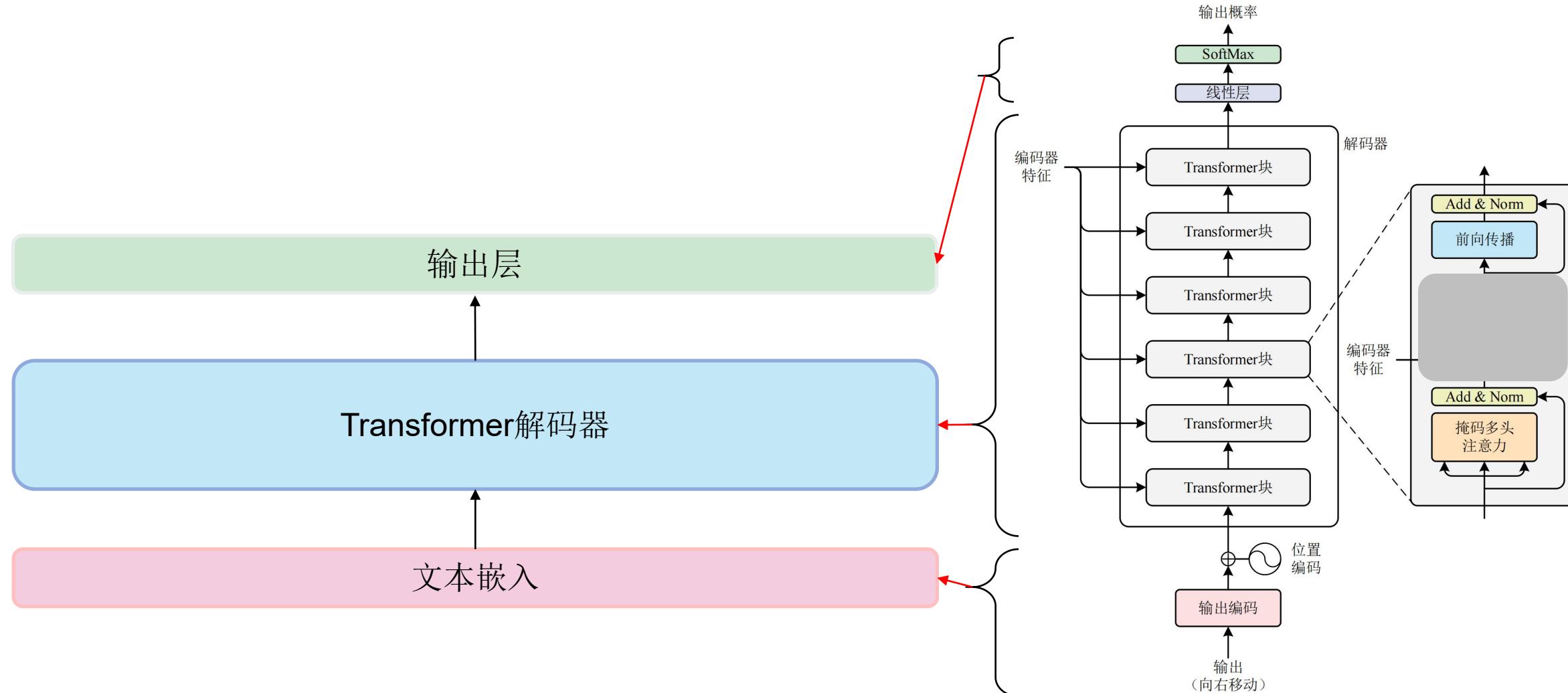


OpenAI

解码器结构-GPT家族

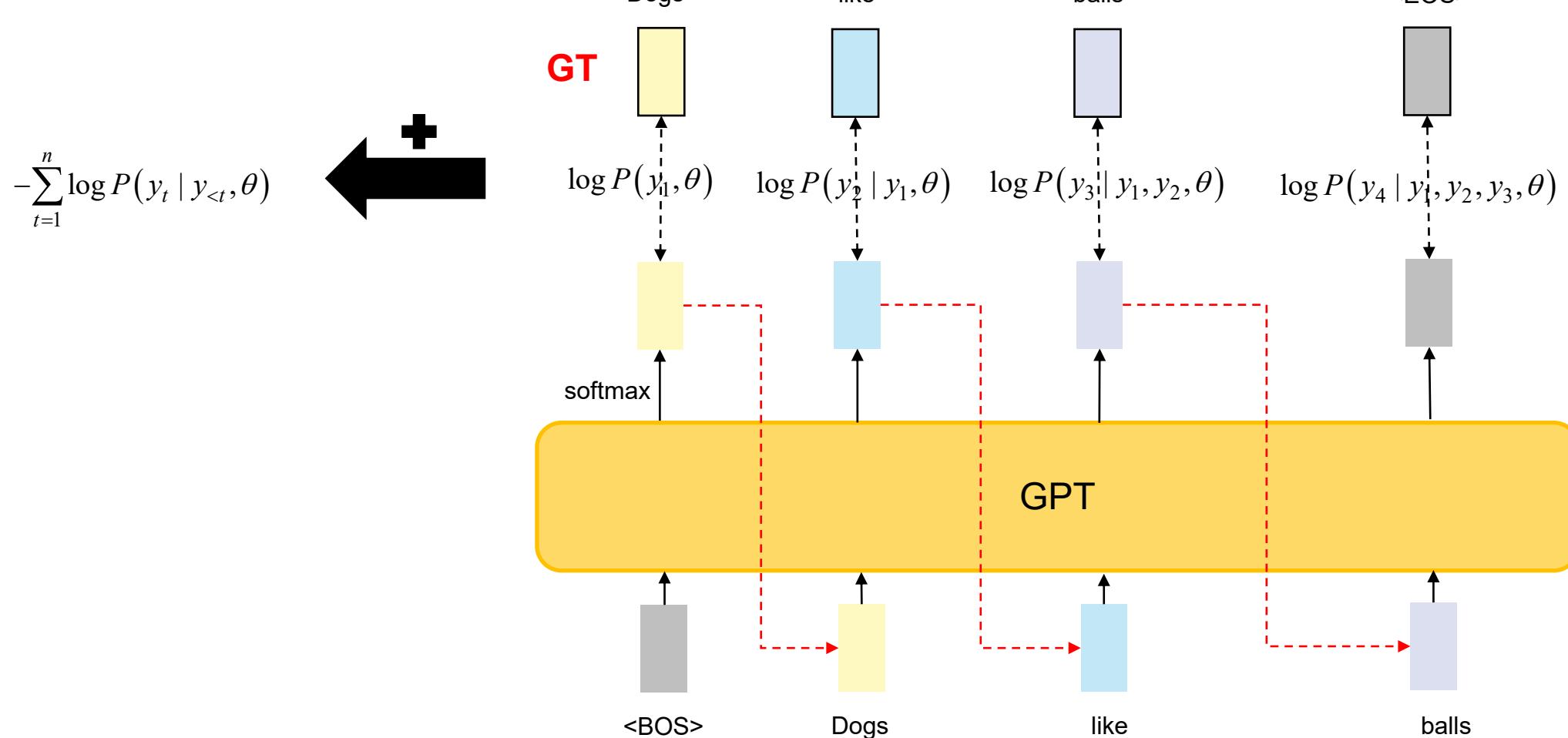


1.GPT结构

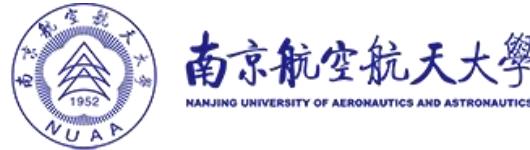


2. 自回归预训练

自监督学习

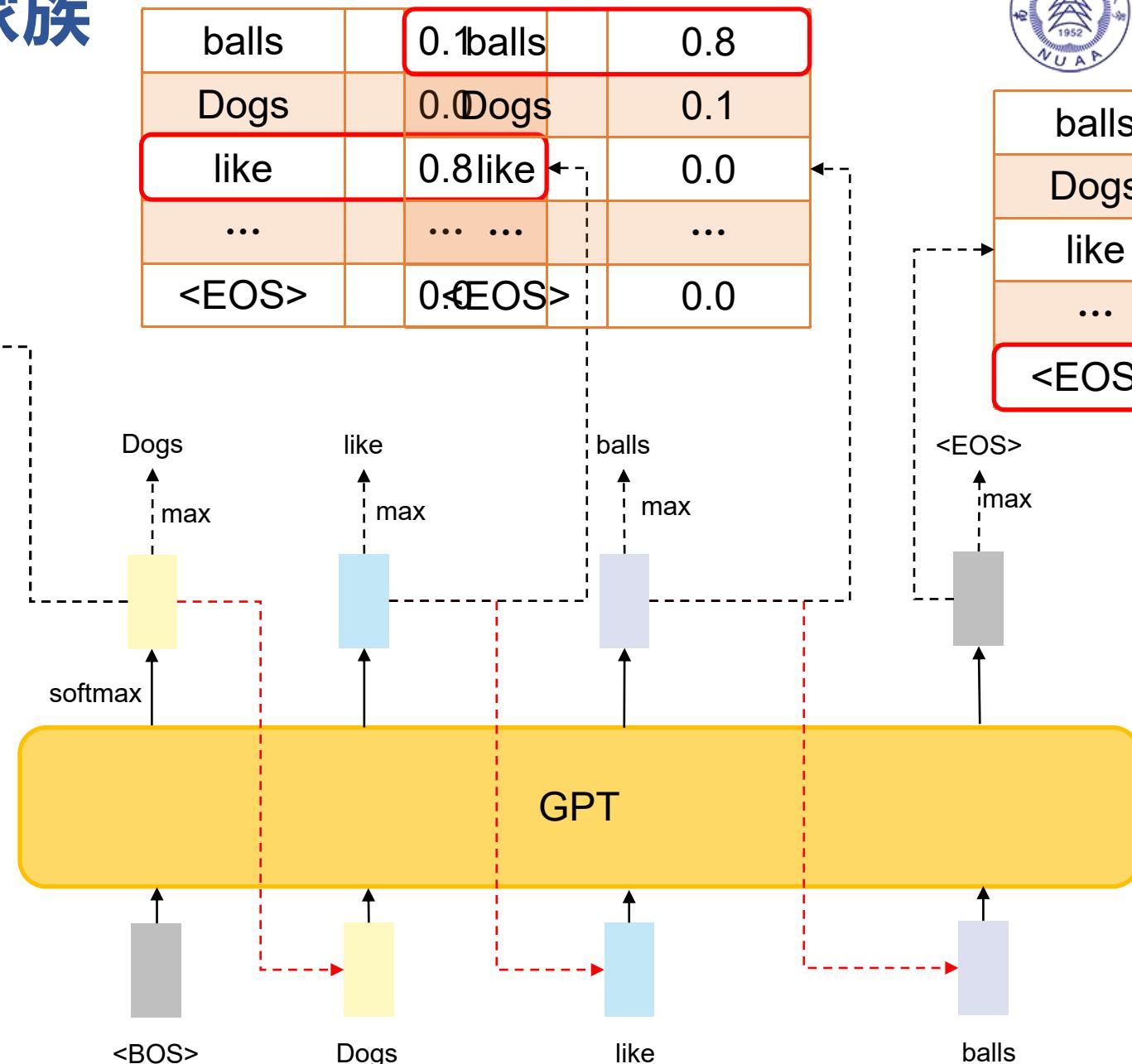


解码器结构-GPT家族

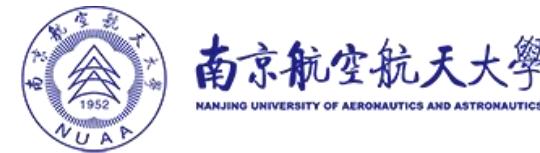


南京航空航天大學
NANJING UNIVERSITY OF AERONAUTICS AND ASTRONAUTICS

balls	0.1
Dogs	0.8
like	0.0
...	...
<EOS>	0.0

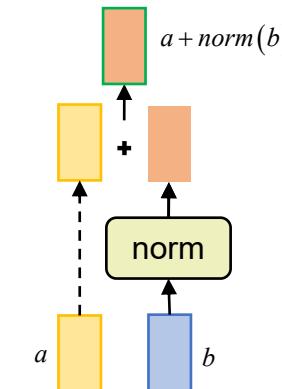
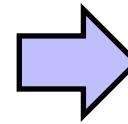
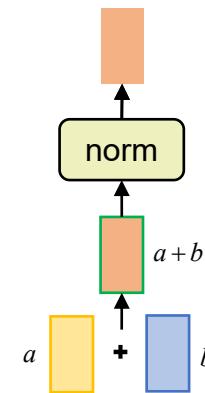


解码器结构-GPT家族

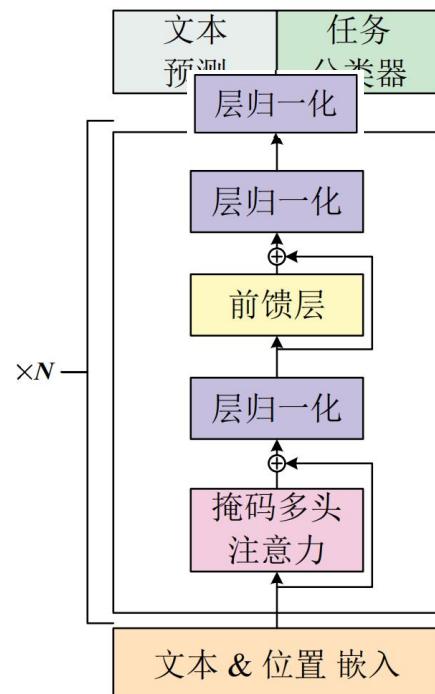


3. 后续改进

GPT-2



1)



缩放参数

增大输入序列长度

解码器结构-GPT家族



GPT-2移除了微调模型，完全只做无监督预训练

zero-shot 零样本学习

2)

$$-\sum_{t=1}^n \log P(y_t|y_{<t}, \theta) \quad \rightarrow \quad P(x) = \prod_{i=1}^n P(x_i|x_1, x_2, \dots, x_{i-1})$$

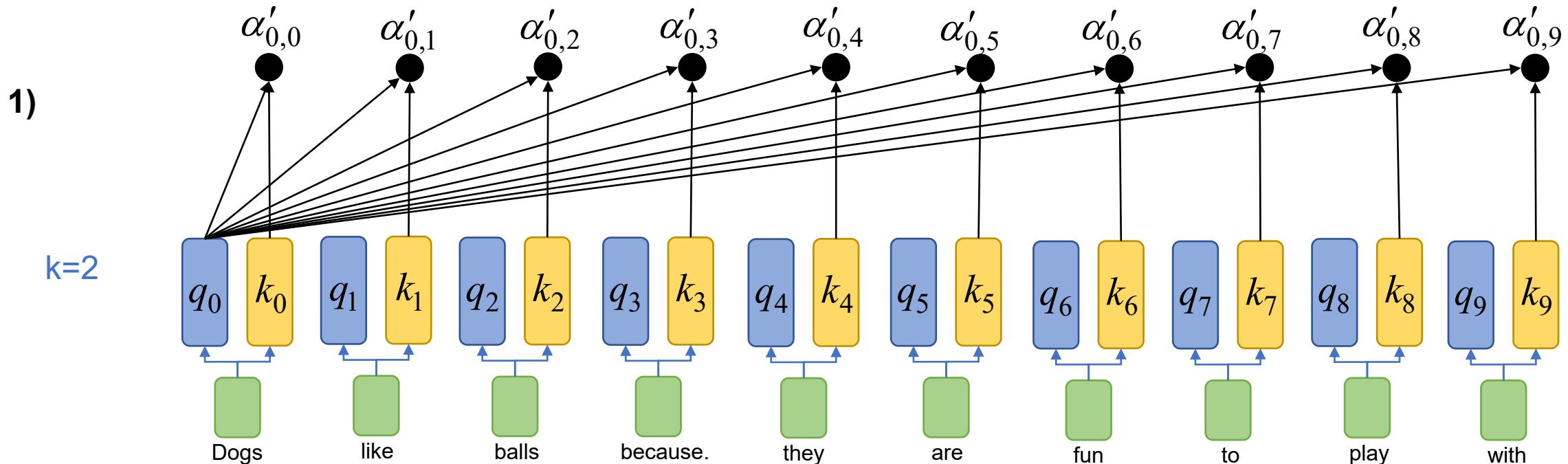
3)

更大规模的预训练数据集

解码器结构-GPT家族



GPT-3



解码器结构-GPT家族



2)

是否进行微调和使用
多少数据进行微调

In-Context Learning的三种设置

零样本: 该模型仅根据任务的自然语言描述来预测答案, 不进行梯度更新。

- 1 把英文翻译成中文: → 任务描述
- 2 Large Model → 提示

一样本: 除了任务描述外, 模型还能看到任务的单个示例, 不进行梯度更新。

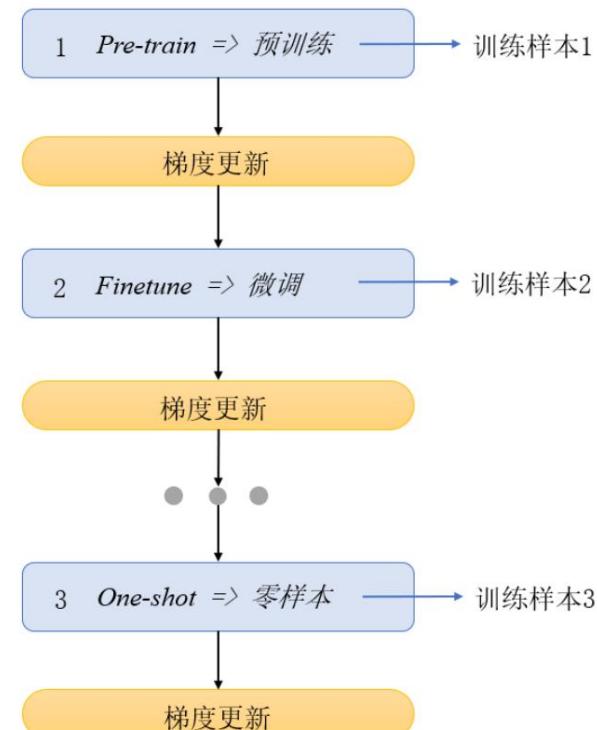
- 1 把英文翻译成中文: → 任务描述
- 2 Pre-train => 预训练 → 例子
- 3 Large Model => 提示

少样本: 除了任务描述外, 模型还能看到一些任务示例(两个及以上), 不进行梯度更新。

- 1 把英文翻译成中文: → 任务描述
- 2 Pre-train => 预训练 → 例子1
- 3 Finetune => 微调 → 例子2
- 4 Large Model => 提示

传统方法中的微调设置(不在GPT-3里使用)

微调: 该模型通过使用大量的示例任务库进行反复梯度更新来训练。



谢谢!
Thanks!