

Transformer

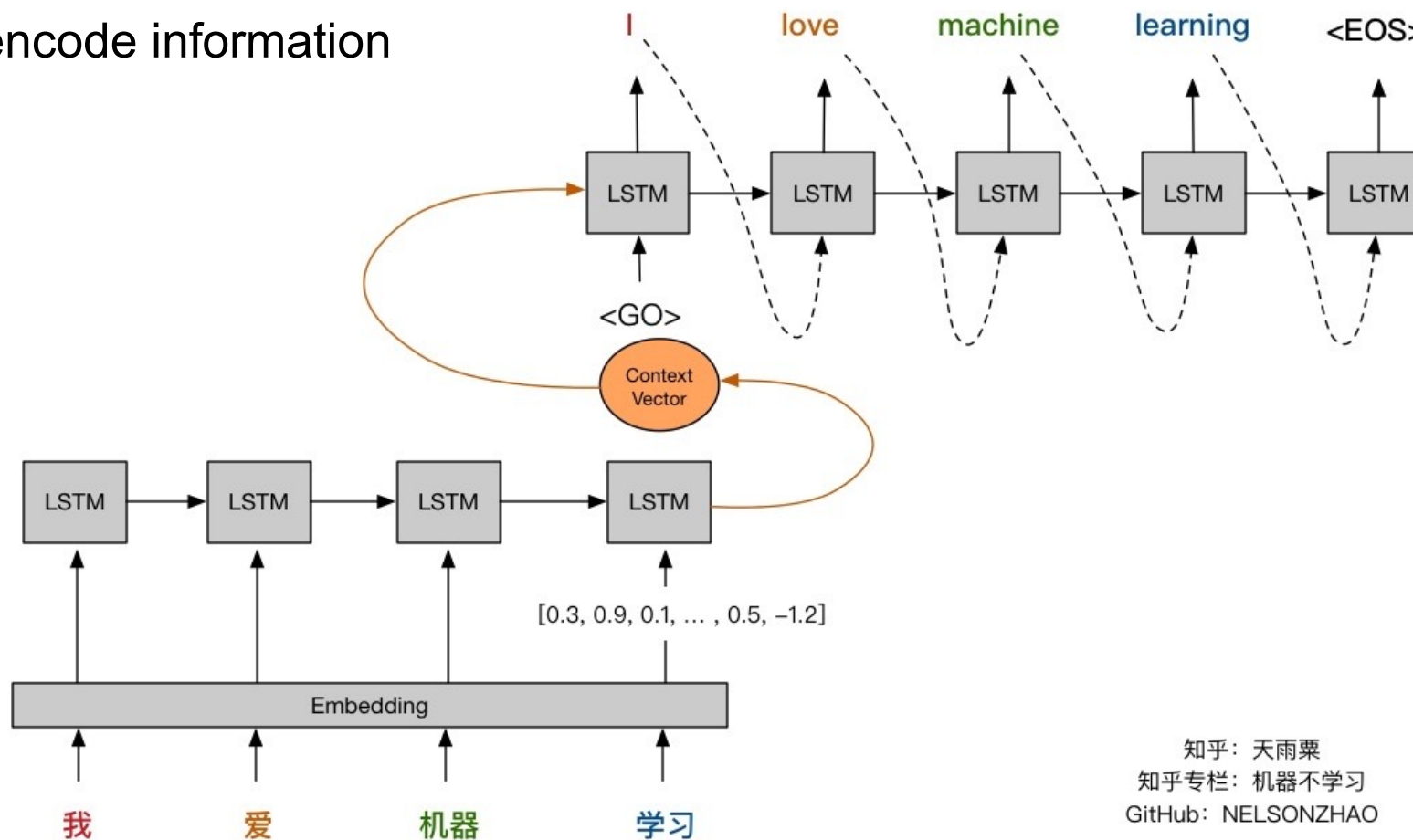
Kun Yuan

Center for Machine Learning Research @ Peking University

Oct. 24, 2023

Traditional seq2seq model

RNN is hard to encode information that is far away



知乎：天雨粟
知乎专栏：机器不学习
GitHub：NELSONZHAO

知乎 @天雨粟

How to capture the most valuable information from a pool of candidate?

Consider a pool of candidate information $D = (k_1, v_1), (k_2, v_2), \dots, (k_m, v_m)$

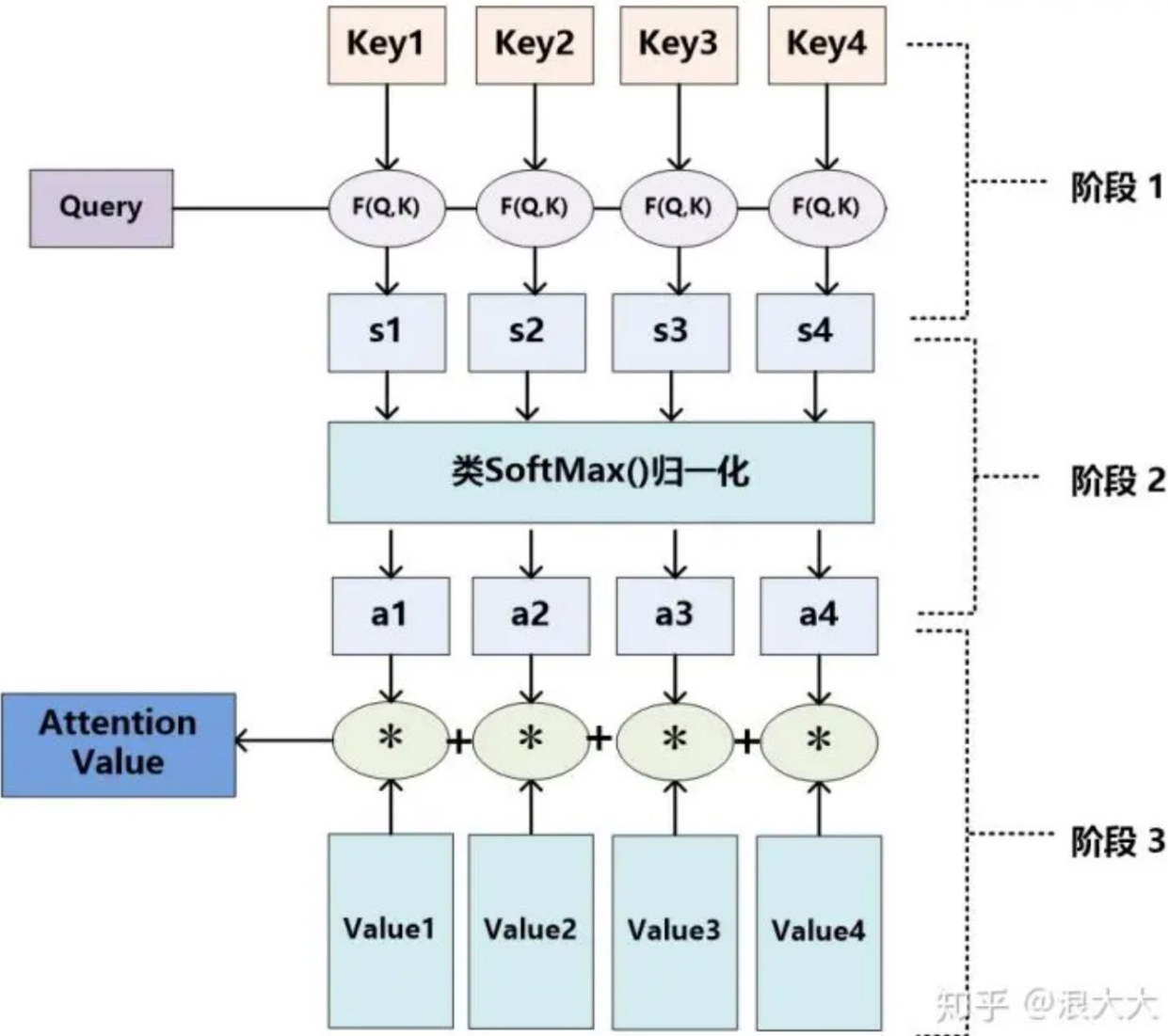
Given a query q , we can capture the most valuable information as follows

$$\text{Attention}(\mathbf{q}, \mathcal{D}) \stackrel{\text{def}}{=} \sum_{i=1}^m \alpha(\mathbf{q}, \mathbf{k}_i) \mathbf{v}_i,$$

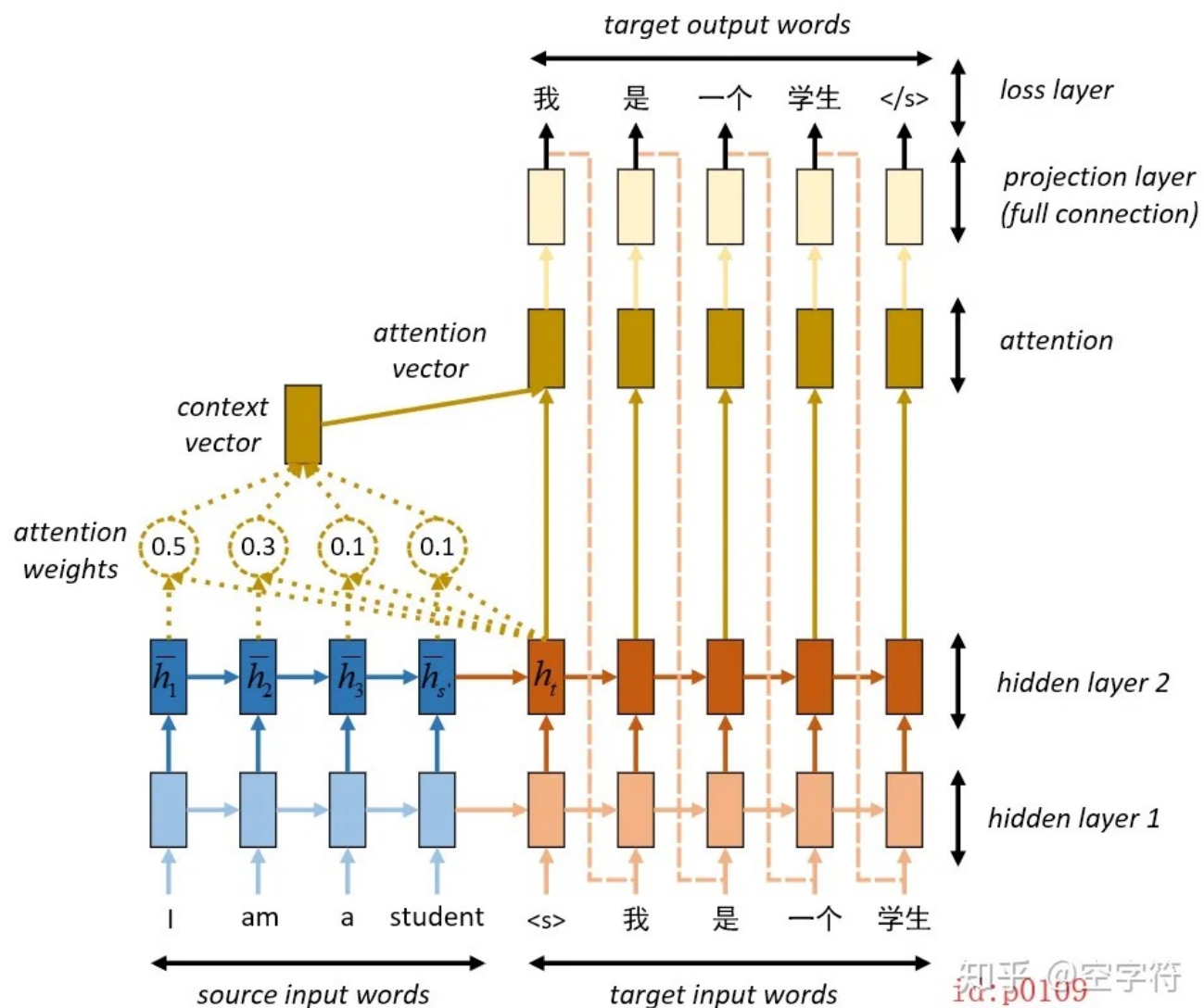
where weight α is to evaluate how close the query q is to key k_i

$$\alpha(\mathbf{q}, \mathbf{k}_i) = \text{softmax}(a(\mathbf{q}, \mathbf{k}_i)) = \frac{\exp(\mathbf{q}^\top \mathbf{k}_i / \sqrt{d})}{\sum_{j=1} \exp(\mathbf{q}^\top \mathbf{k}_j / \sqrt{d})}.$$

Attention



Seq2Seq with attention



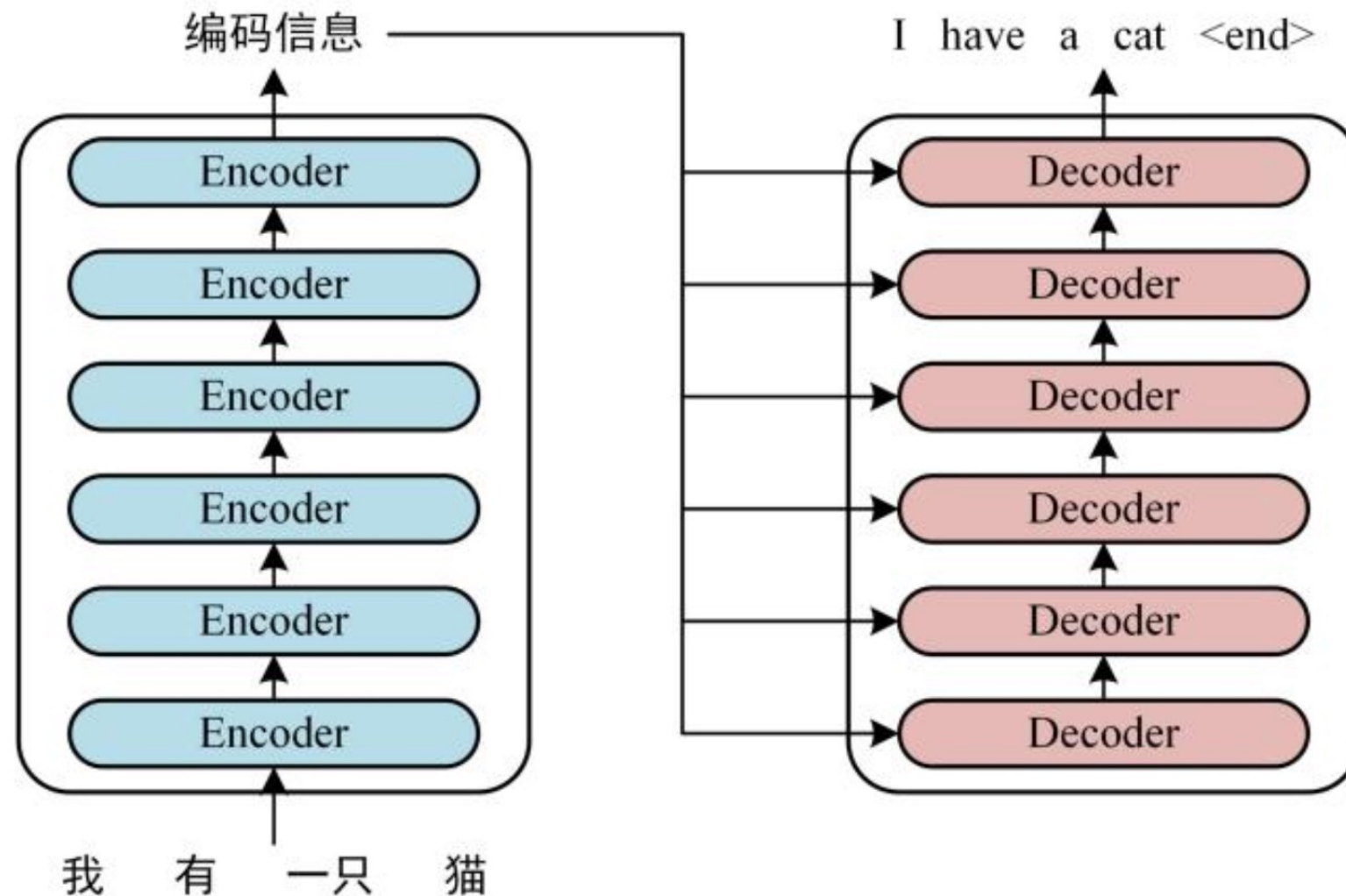
How to get representation h ?

One way is RNN

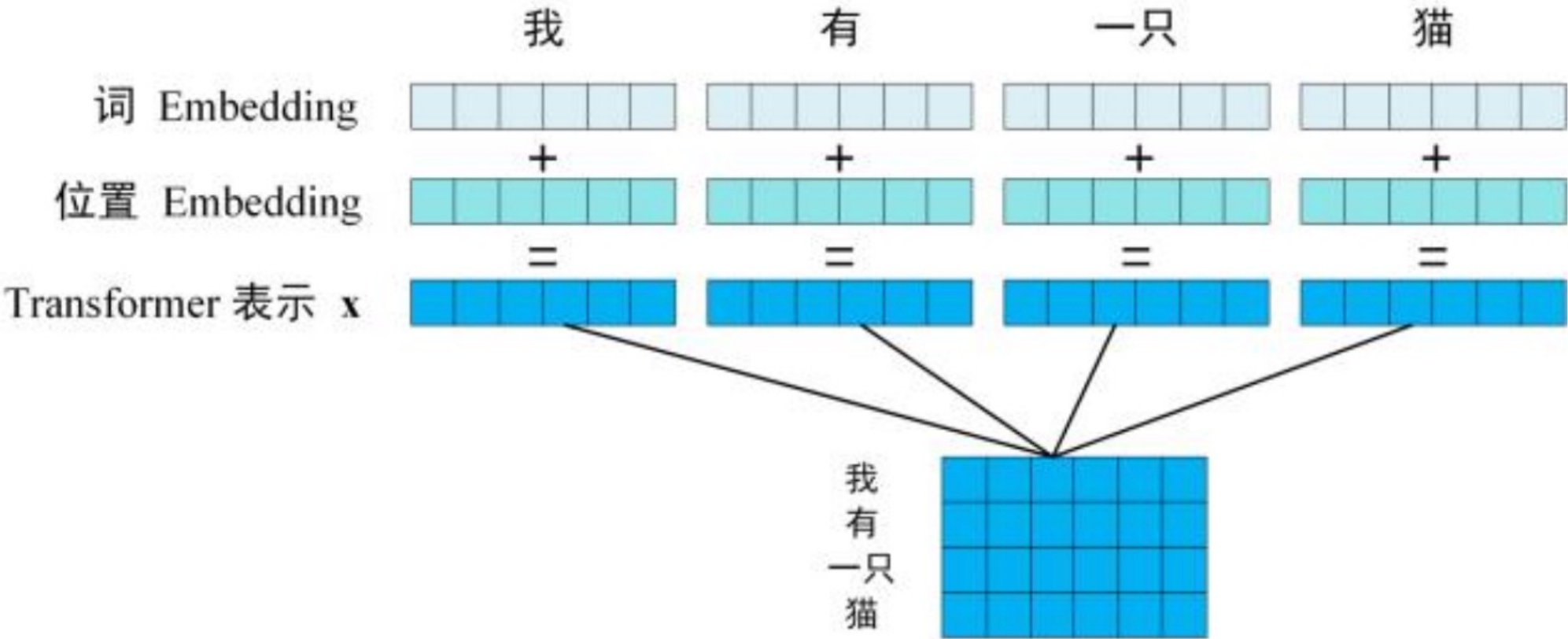
The other is Transformer!

Transformer overview

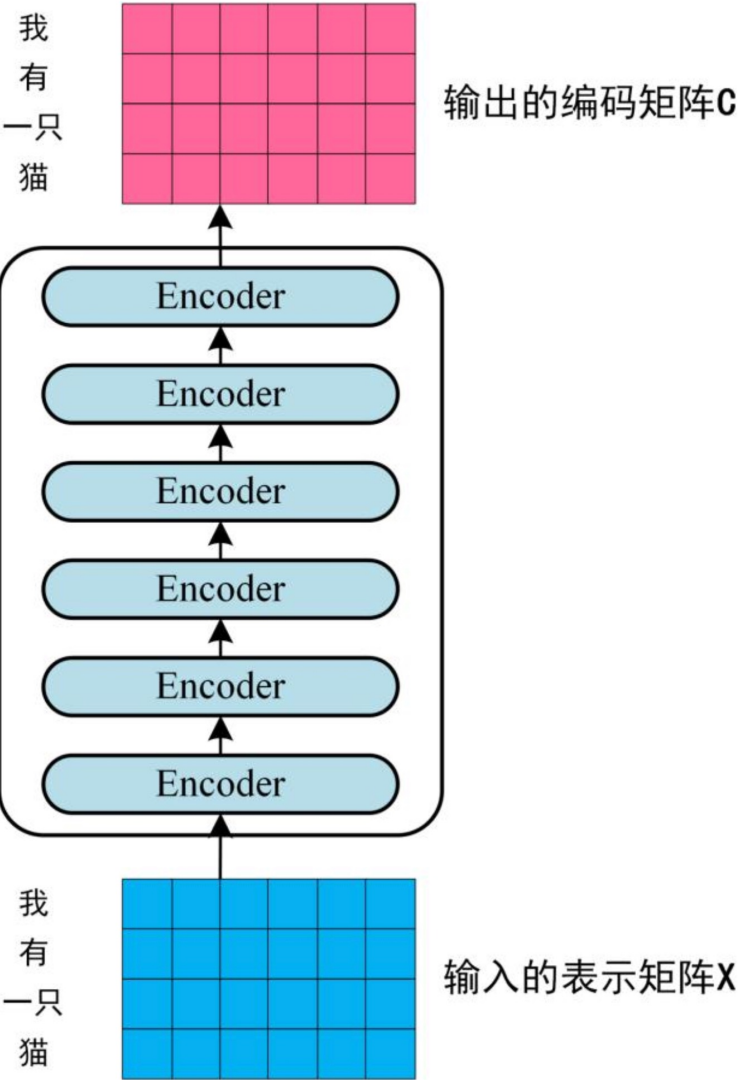
Transformer also
utilizes seq2seq



Step 1: represent input information

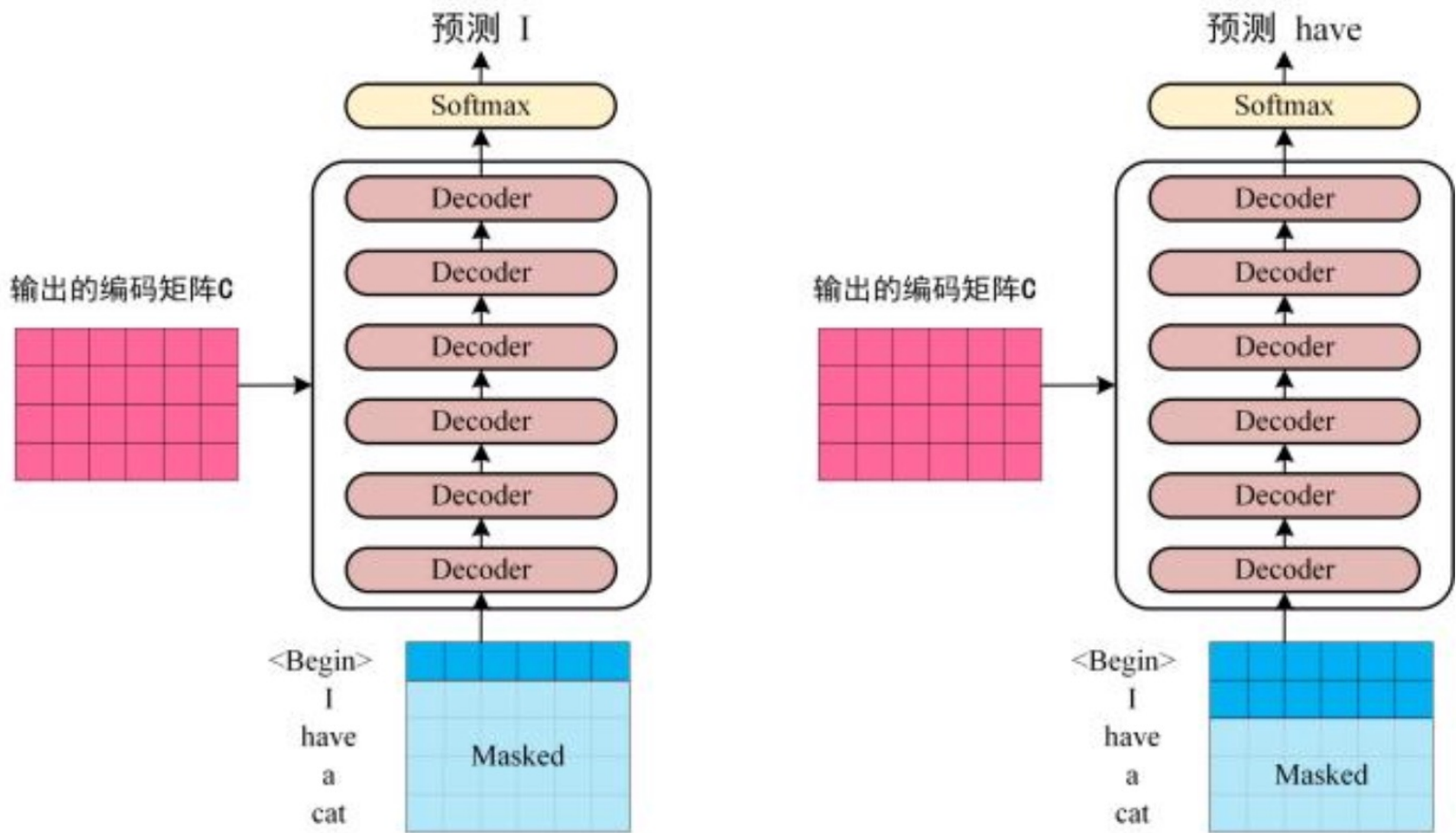


Step 2: Encode input



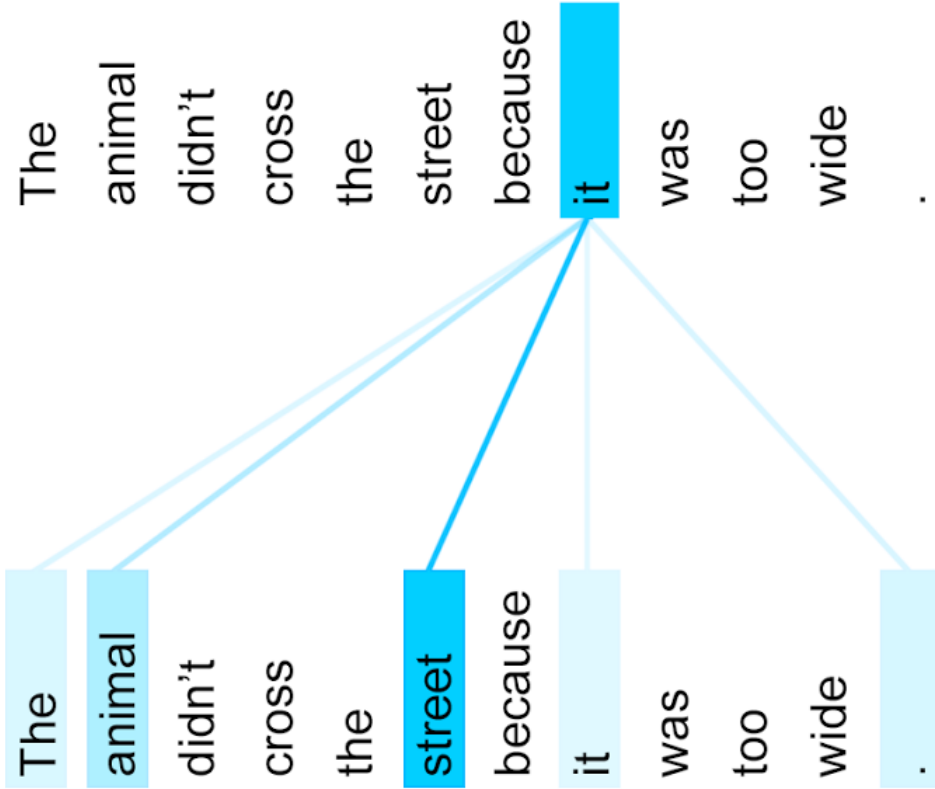
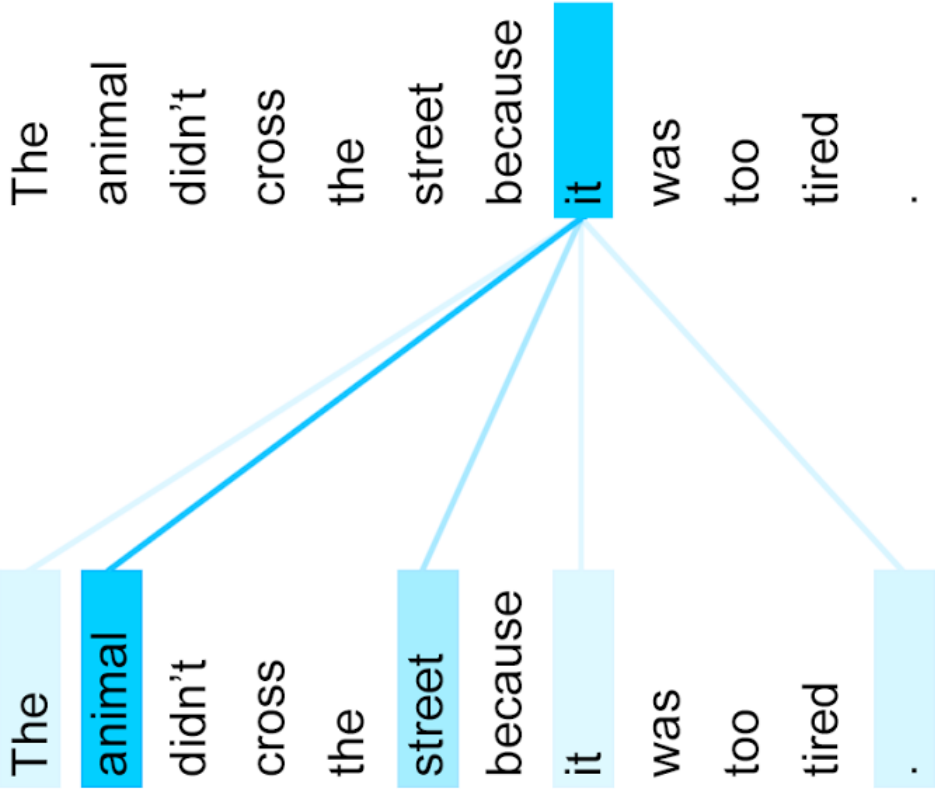
Transformer Encoder 编码句子信息

Step 3: Decode and predict

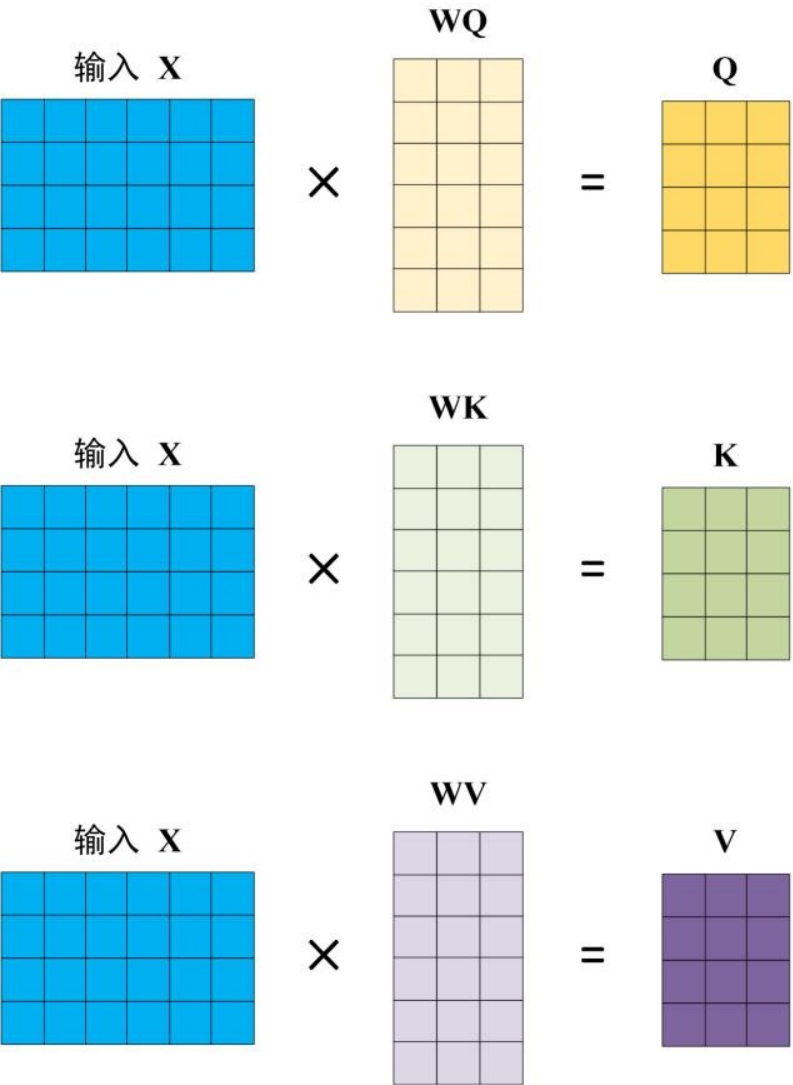
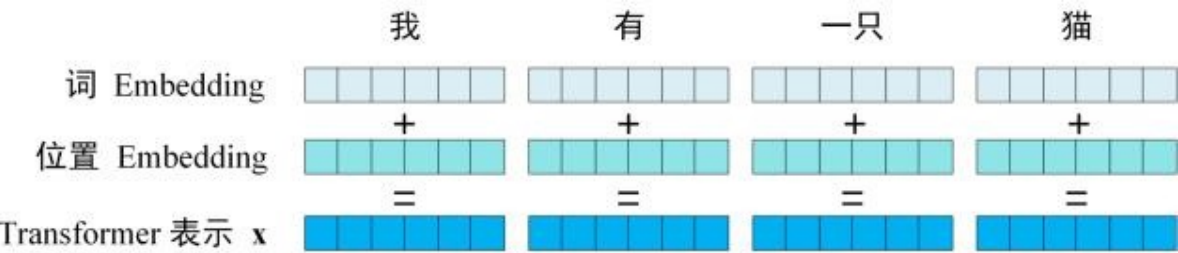


Transofrmer Decoder 预测

Encode: self-attention

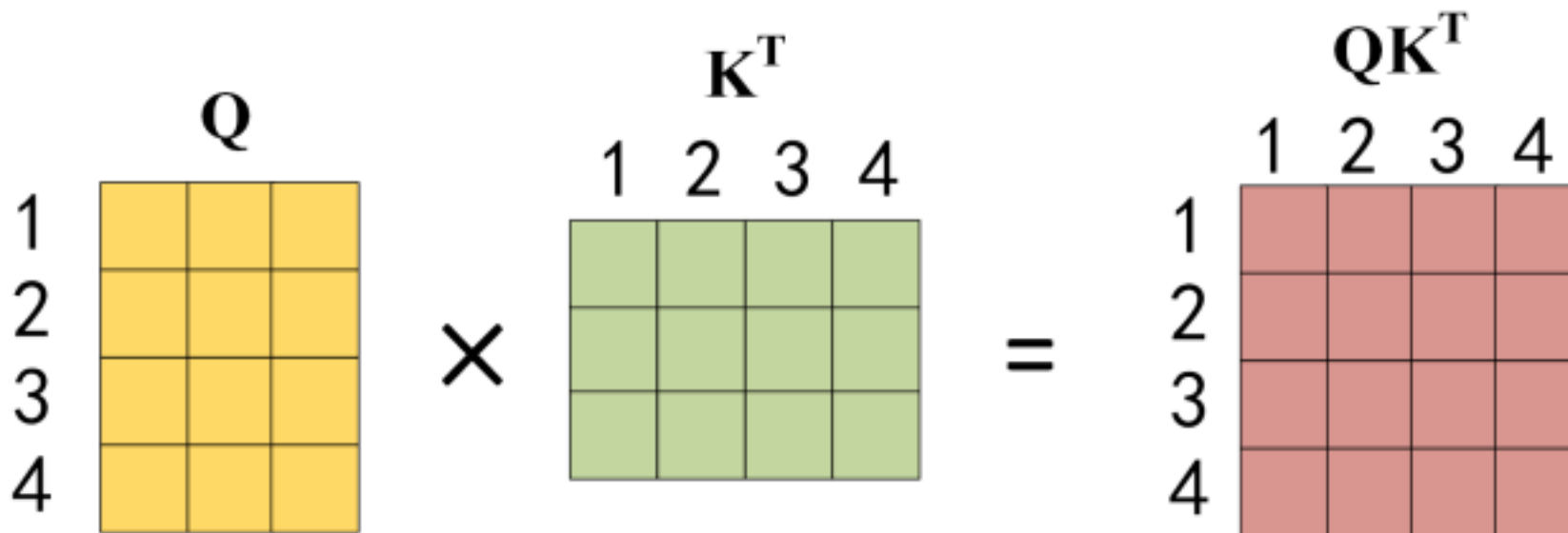


Encode: self-attention



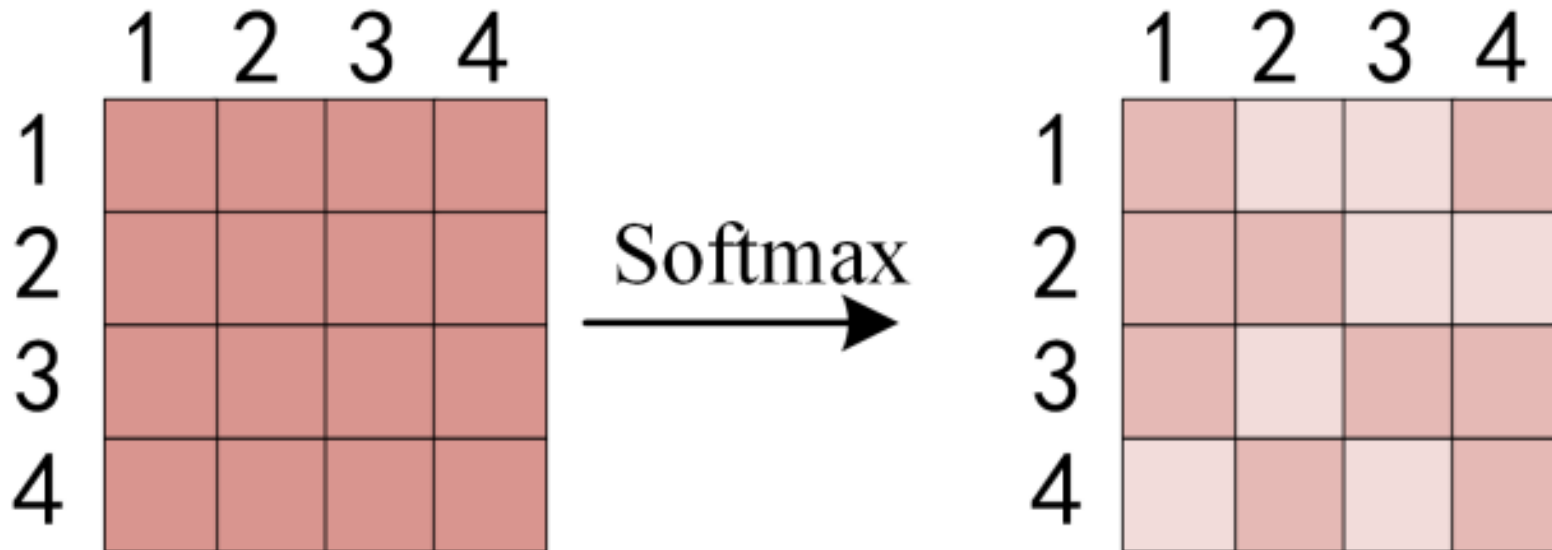
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

d_k 是 Q, K 矩阵的列数，即向量维度



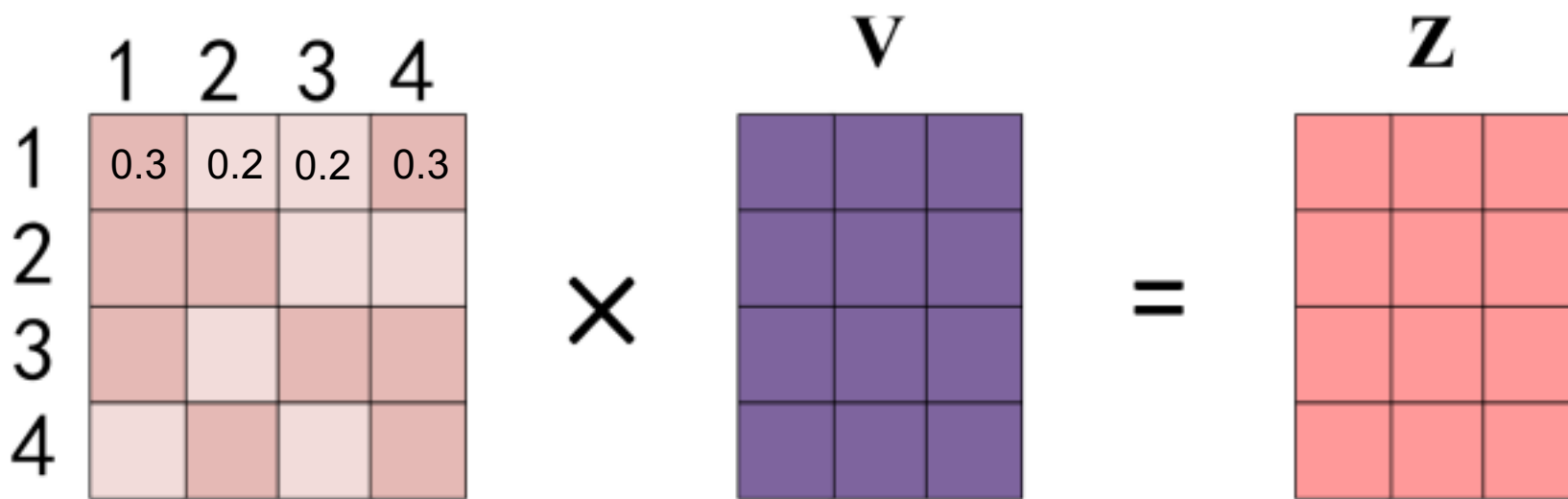
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

d_k 是 Q, K 矩阵的列数，即向量维度

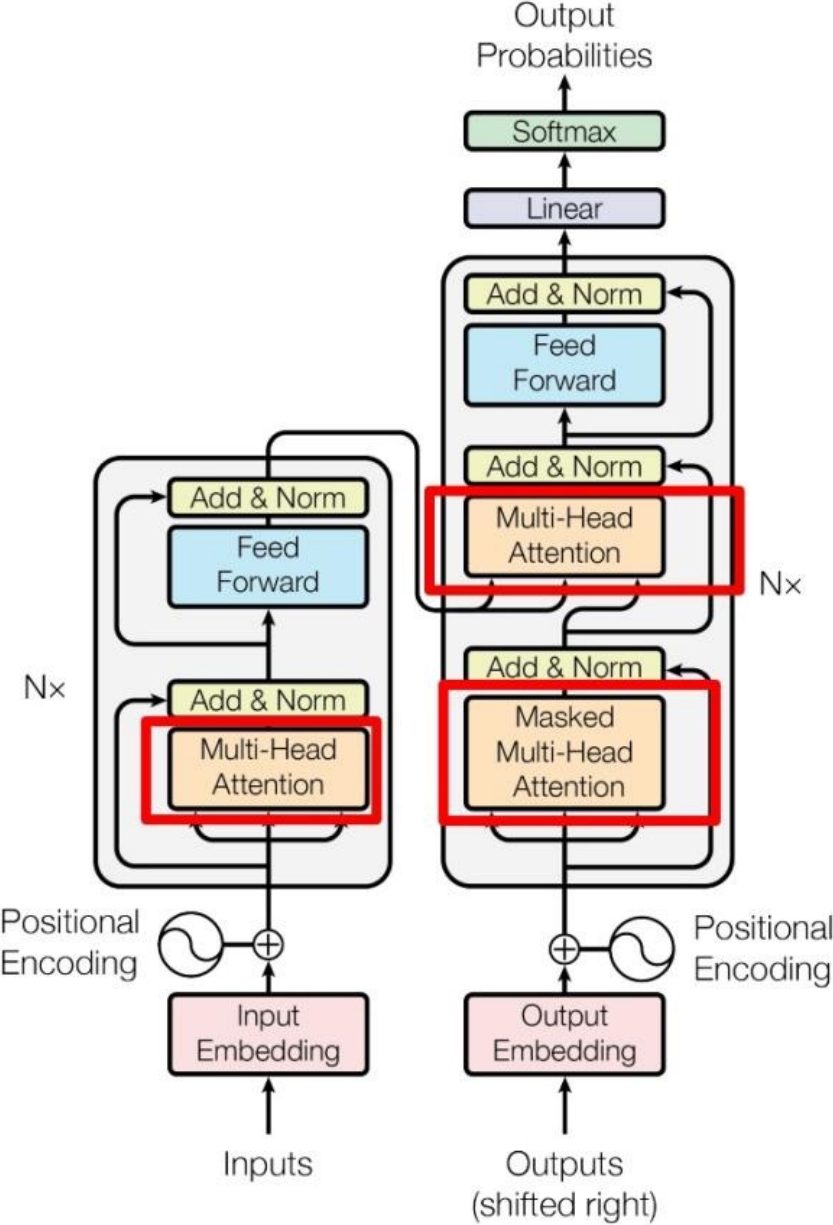


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

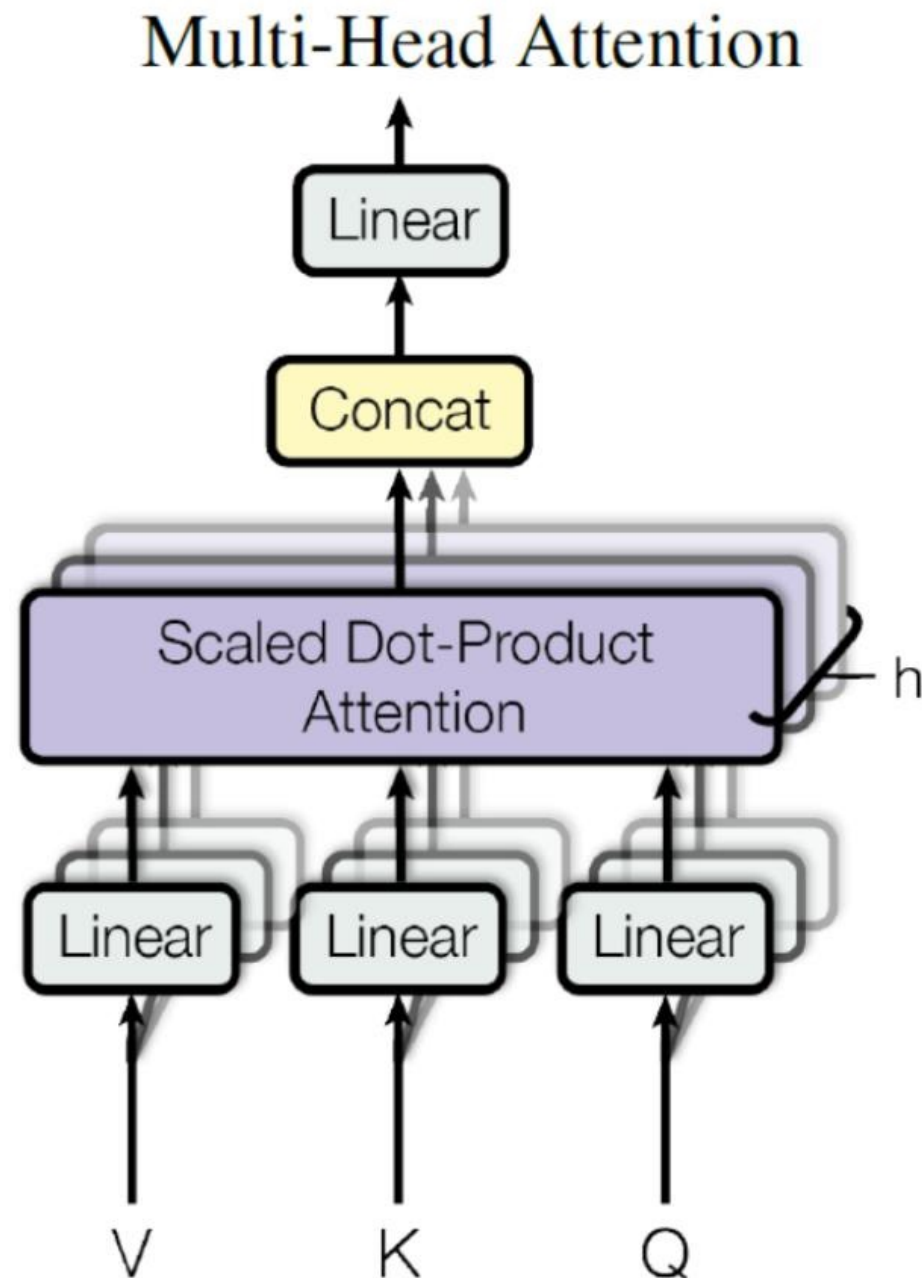
d_k 是 Q, K 矩阵的列数，即向量维度



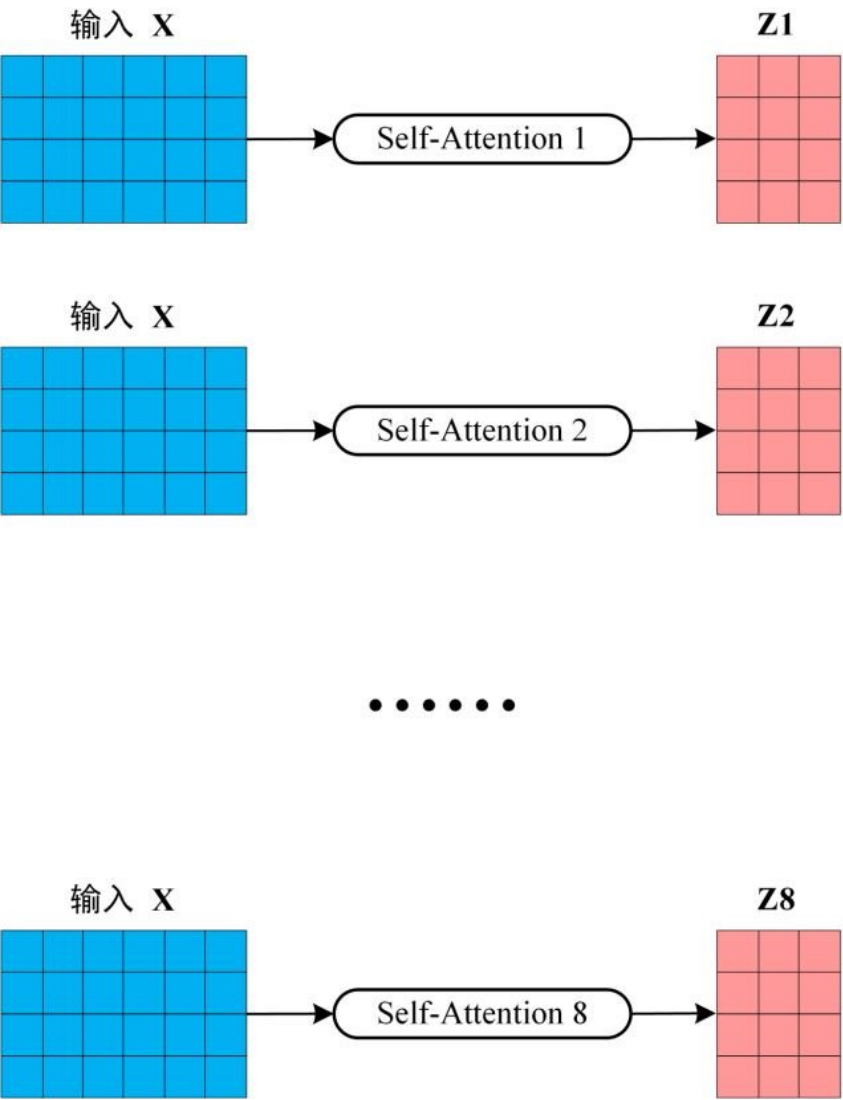
Encode: multi-head attention



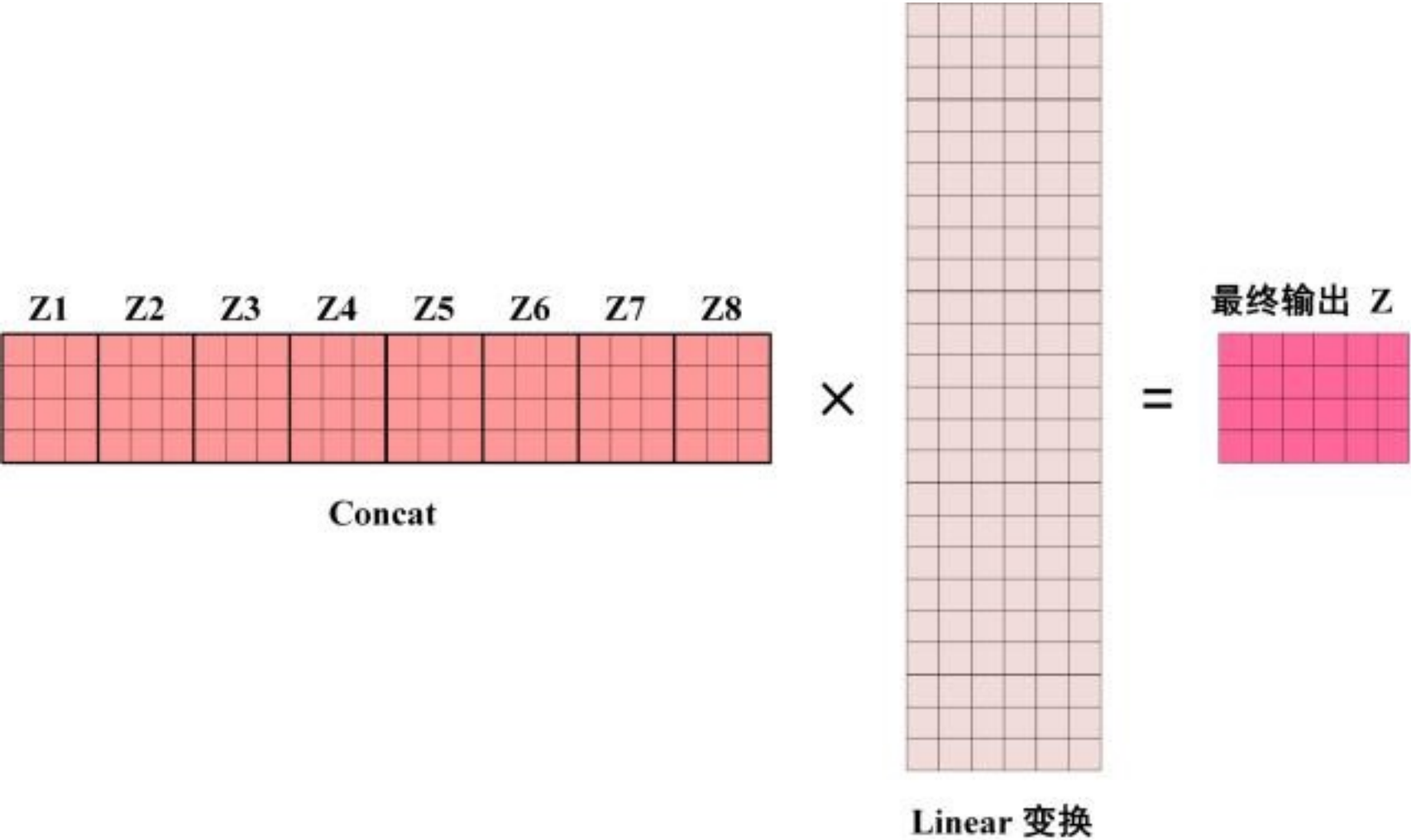
Encode: multi-head attention



Encode: multi-head attention



Encode: multi-head attention

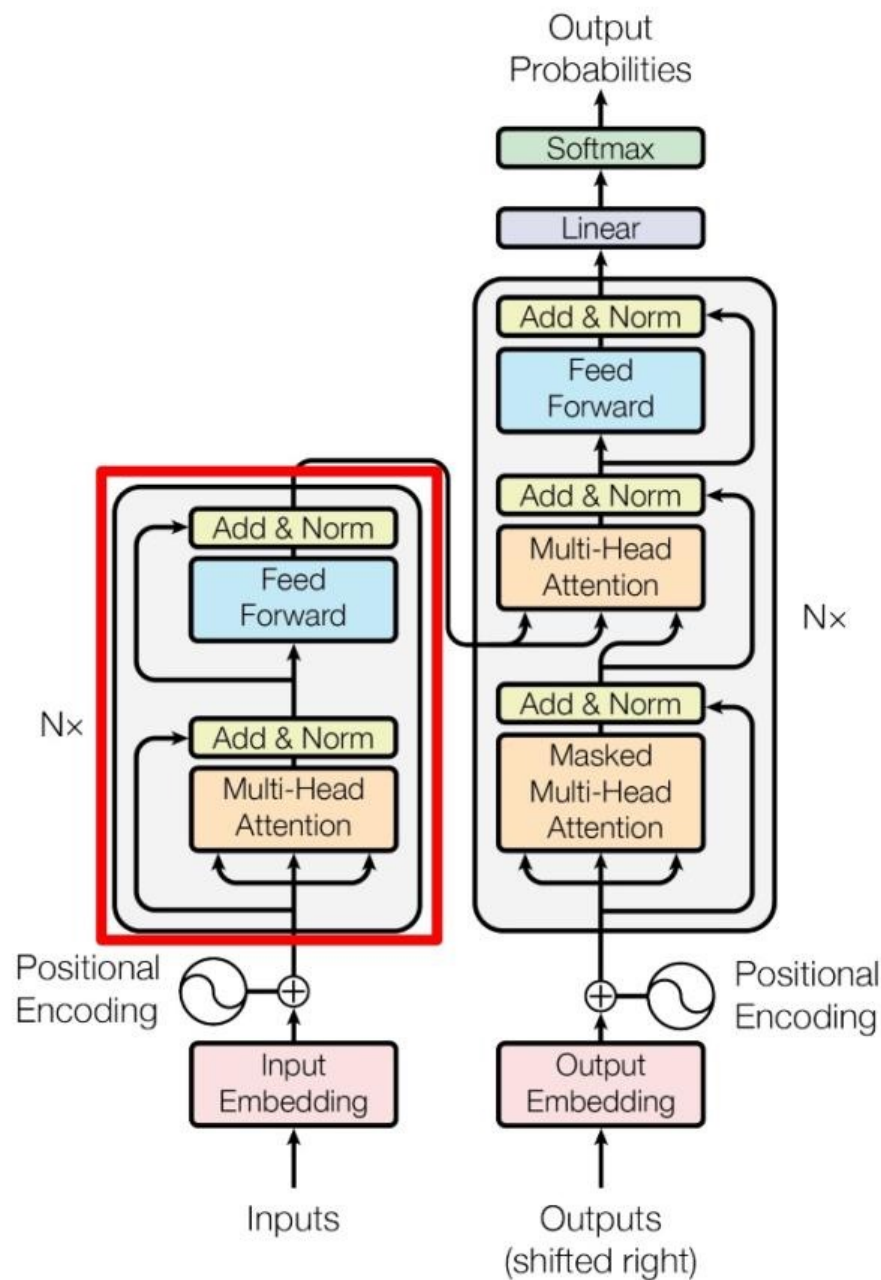


Encode: add and norm

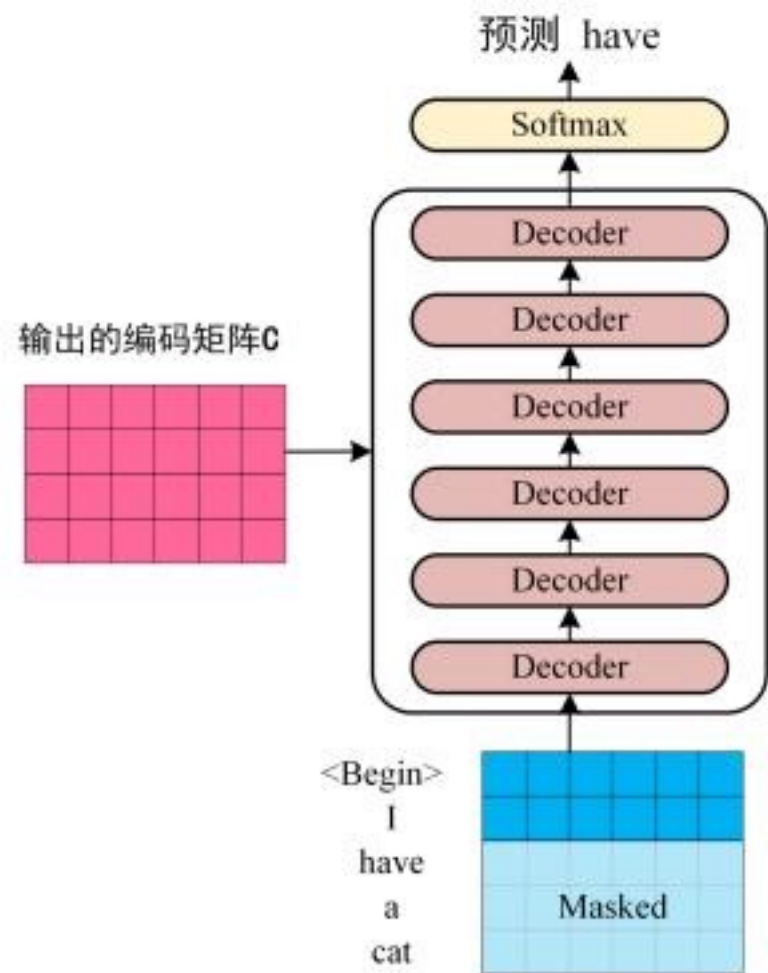
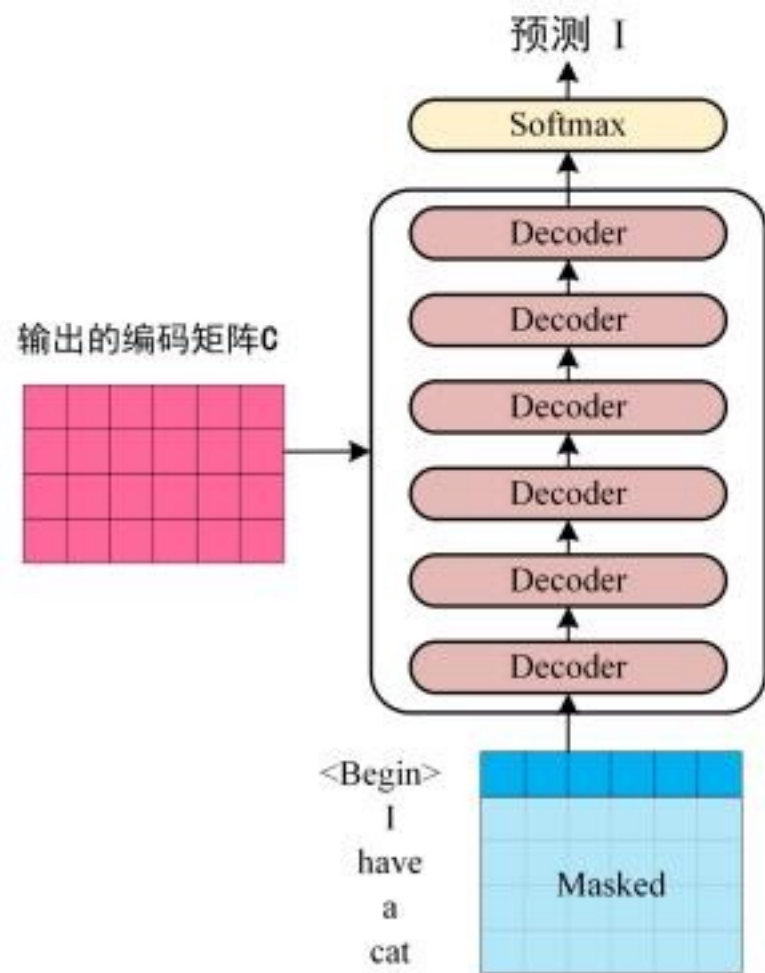
$\text{LayerNorm}(X + \text{MultiHeadAttention}(X))$

$\text{LayerNorm}(X + \text{FeedForward}(X))$

残差链接，解决深层训练问题



Decoder: masked self-attention

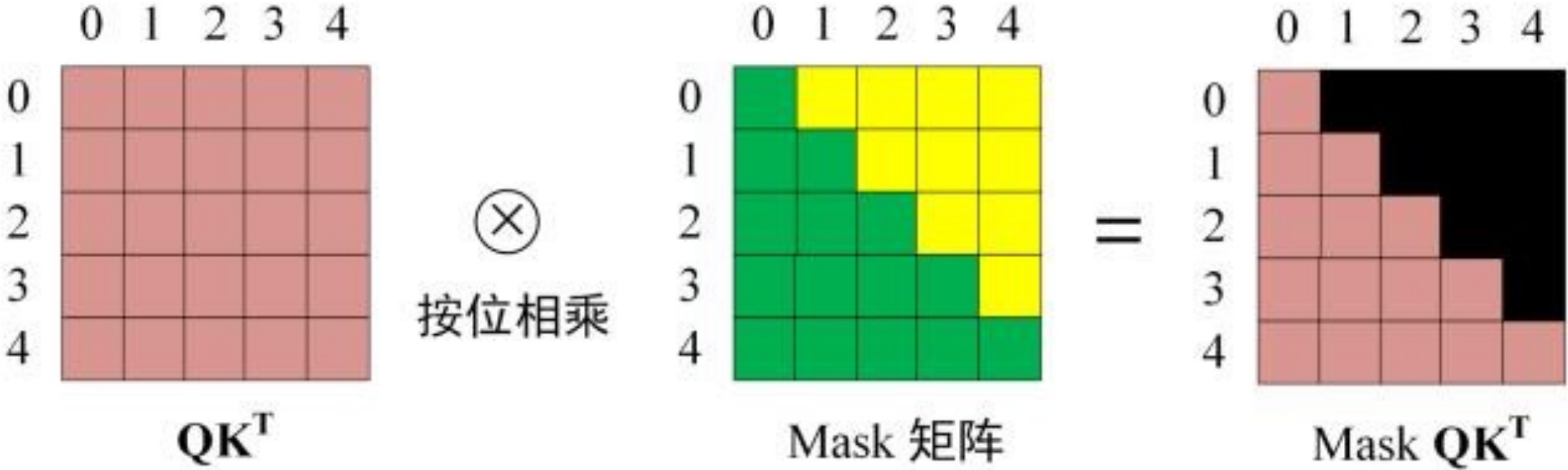


Decoder: masked self-attention

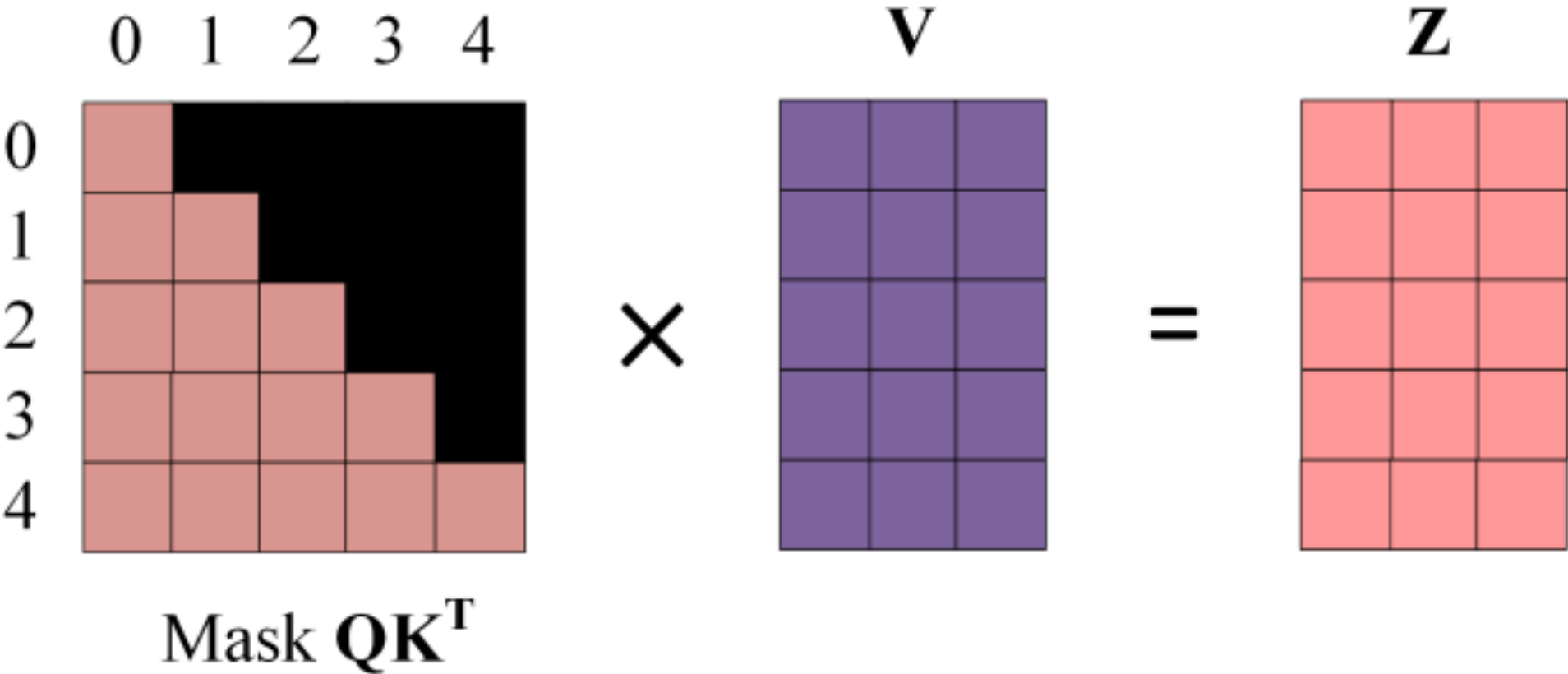
$$\begin{array}{c} \mathbf{Q} \\ \begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \\ 4 \end{array} \end{array} \times \begin{array}{c} \mathbf{K}^T \\ \begin{array}{c} 0 \ 1 \ 2 \ 3 \ 4 \end{array} \end{array} = \begin{array}{c} \mathbf{QK}^T \\ \begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \\ 4 \end{array} \end{array}$$

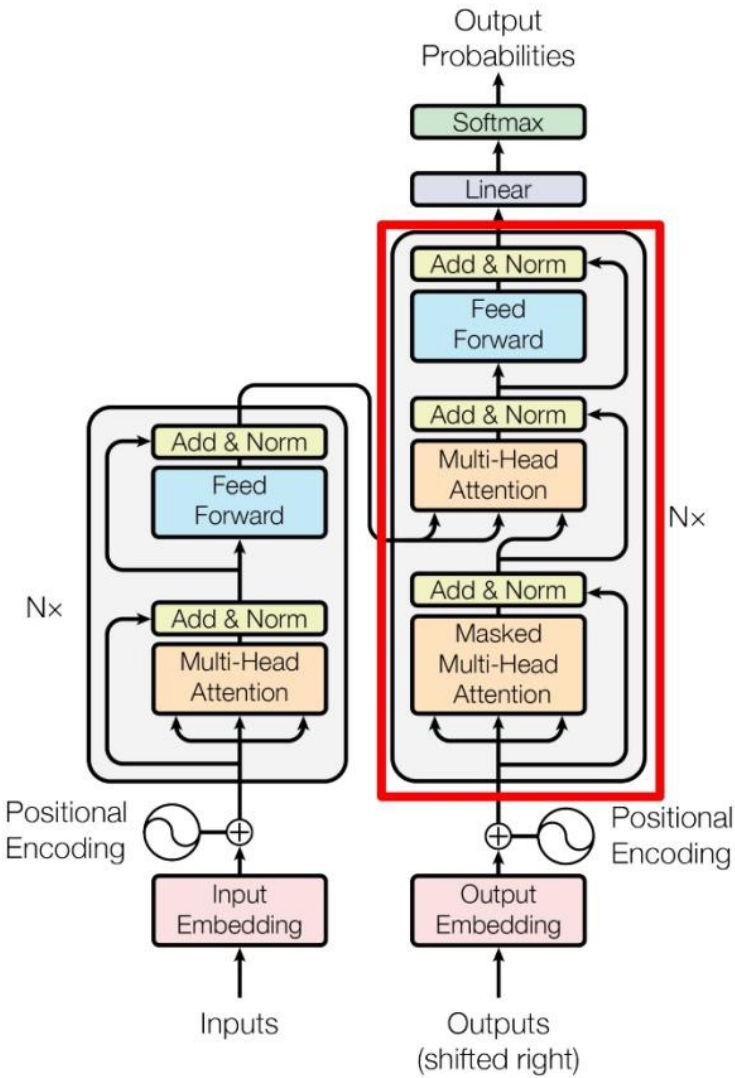
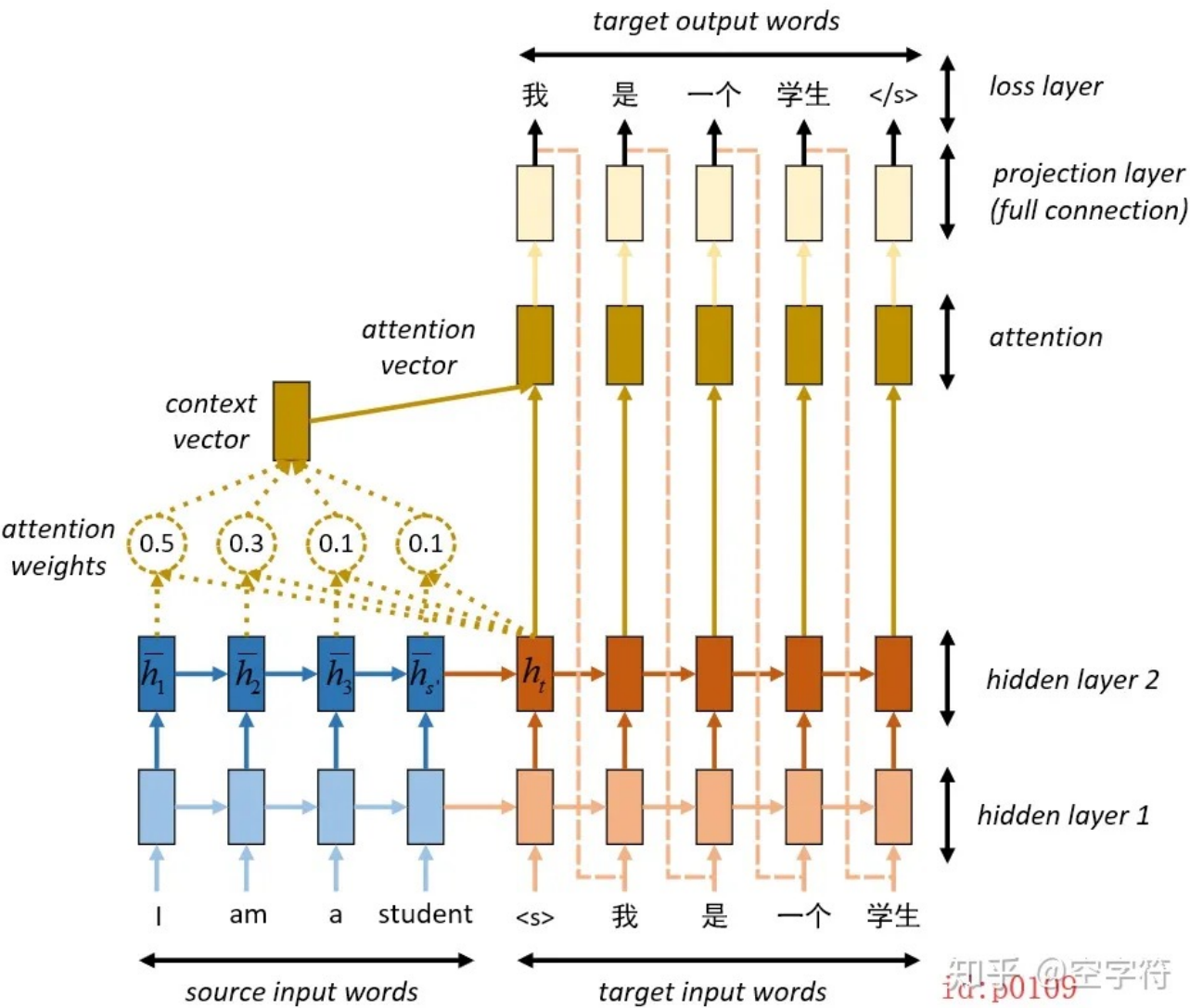
The diagram illustrates the matrix multiplication of \mathbf{Q} and \mathbf{K}^T to produce \mathbf{QK}^T . Matrix \mathbf{Q} is a 5x3 yellow grid, \mathbf{K}^T is a 5x5 green grid, and the result \mathbf{QK}^T is a 5x5 red grid. All grids have indices 0 to 4 on both axes.

Decoder: masked self-attention

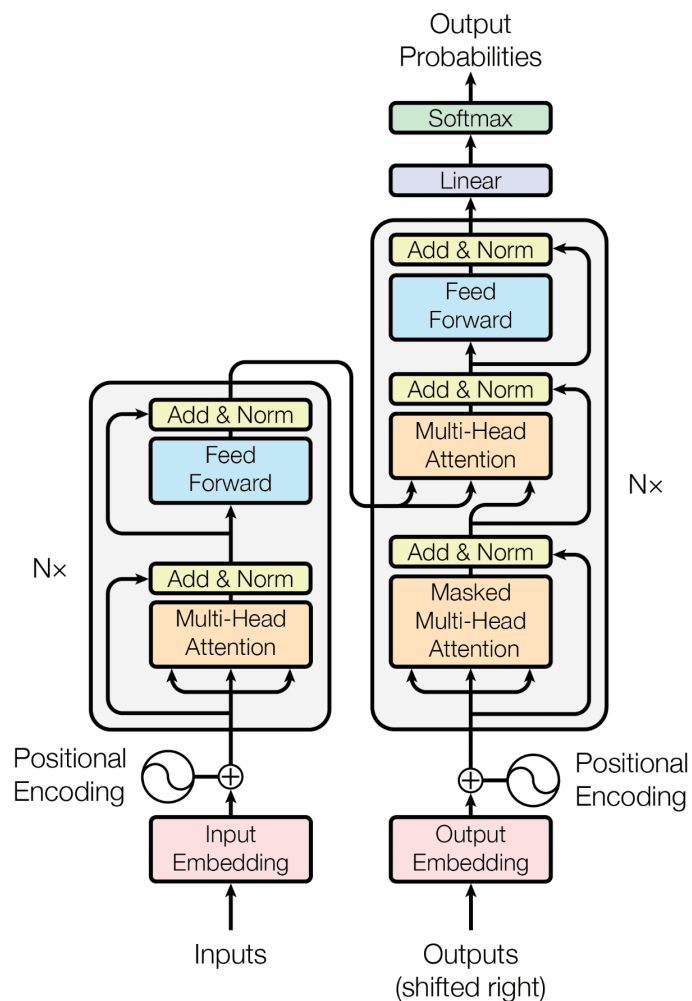


Decoder: masked self-attention





Summary



Transformer没有时序结构，容易并行

需要引入位置的embedding

内存开销巨大

Reference

<https://zhuanlan.zhihu.com/p/338817680>