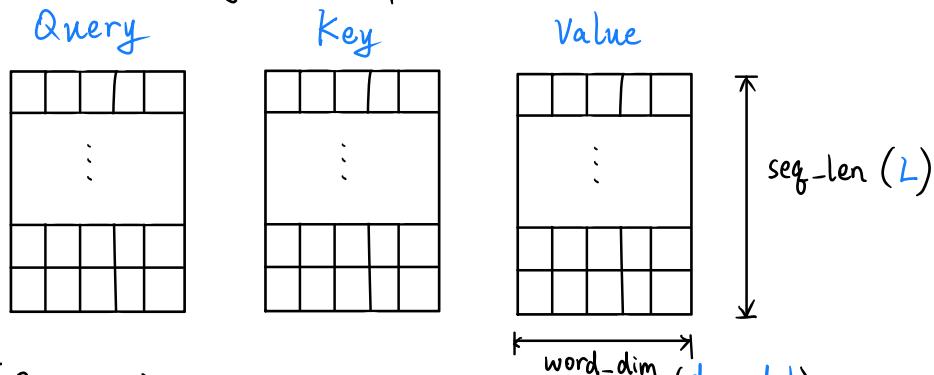


Attention : 因的是计算两个序列的相似度是多矩阵(分数), 也就是 Q 与 V 的相似度  
(权重)

每一行代表一个词(向量), Query 是你要算哪些词对 key-value 的注意力分数



def attention(Q, K, V):

(1)  $Q \cdot K^T$ : 每个要被查询的 q 对于每一个词 k 的关系, 得到 Q 与 K 的关系矩阵

(2)  $\text{softmax}(Q \cdot K^T)$ : 将关系矩阵归一化成和为 1 的权重矩阵

(3)  $\text{softmax}(Q \cdot K^T) V$ : 权重矩阵与真值 V 做乘法得到“分数”

(4)  $\text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$ :  $d_K$  为  $Q \cdot K^T$  的维度 ( $d_K \cdot d_K$ ),  $d_K$  高的话 softmax 会放大值之间的差异

Self - Attention : 只关注自身序列的注意力分数

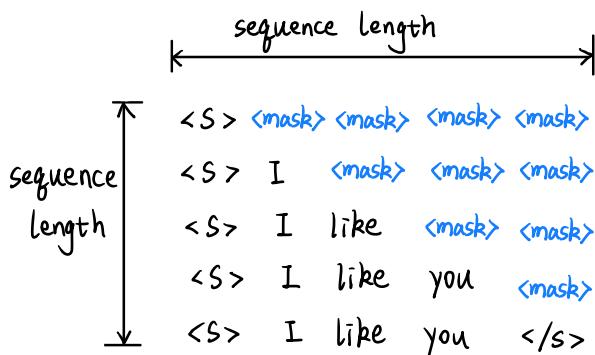
def self-attention = attention(x, x, x)

Mask Self - Attention : 盖住某些 token 不让 model 学习 (只看到历史信息而看不到未来)

① 对任意文字进行补全

② 并行学习

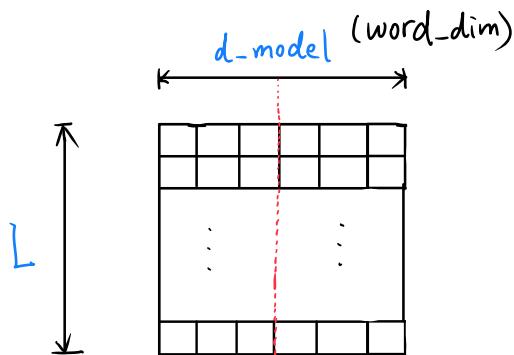
③ 一般为上三角矩阵



计算注意力分数时, score += mask 即  
忽略上三角部分的注意力分数

MultiHead - Attention : 每个头学习文字的不同语义特征

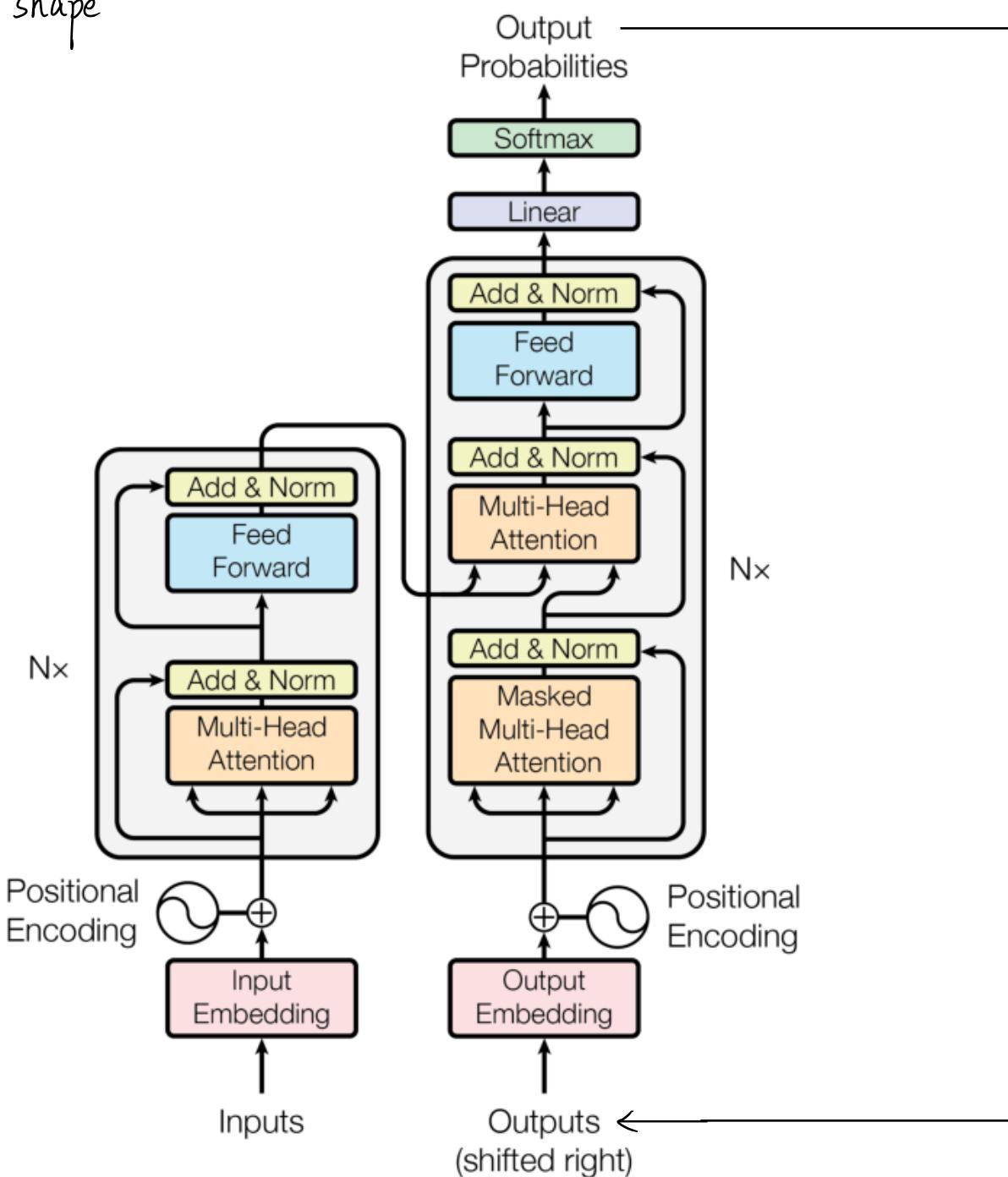
$\text{MultiHead}(Q, K, V) = \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_n)W^O$   
 , where  $\text{head}_n = \text{Attention}(QW_n^Q, KW_n^K, VW_n^V)$



(仍能表达词与词之间的关系,但仅停留在保留的  $k$  个维度上)

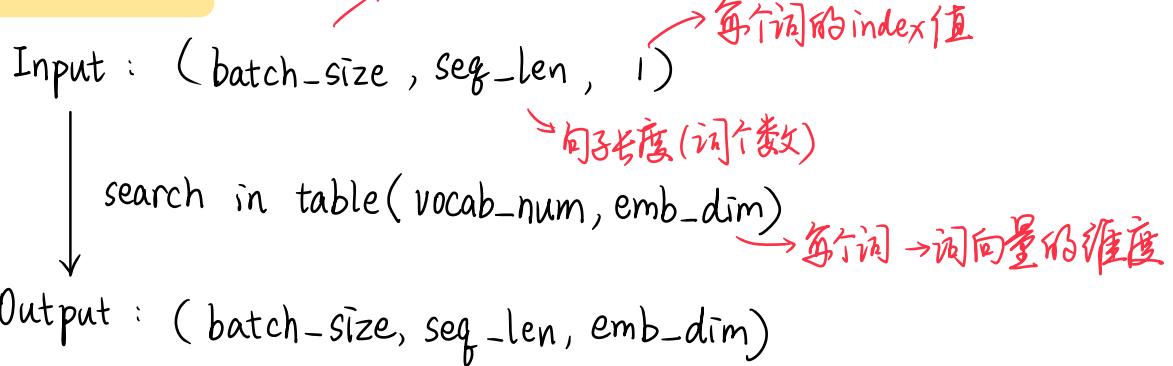
剪的是每个词的维度,所以每个 head 中的注意力分数仍是关于所有词的

code with shape

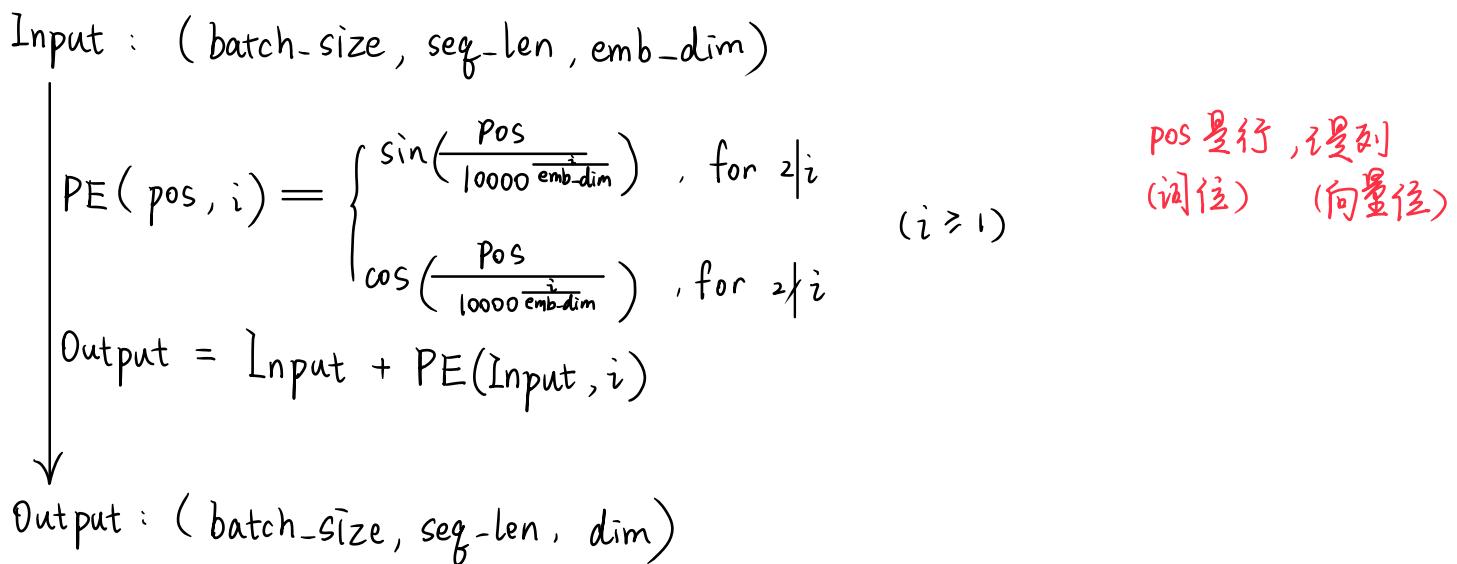


Embedding ↗:

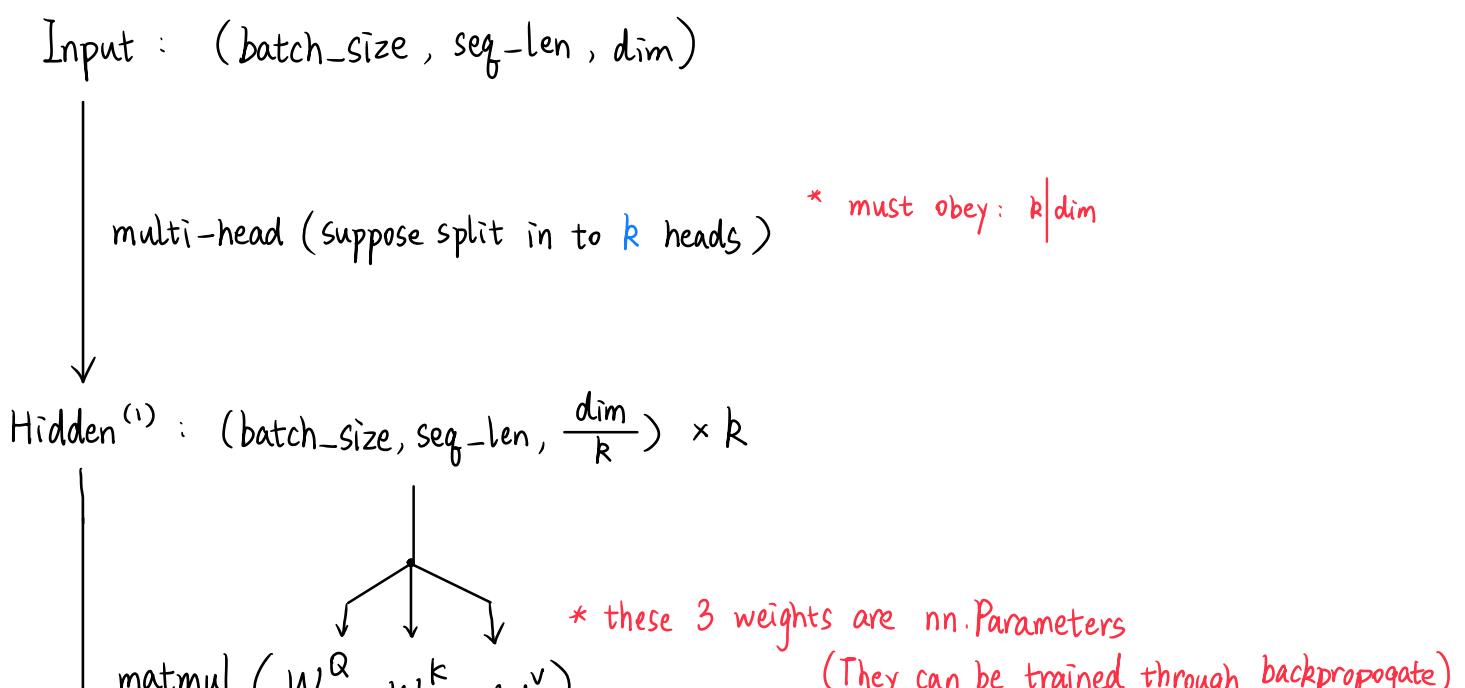
→ 一个 batch 的句子

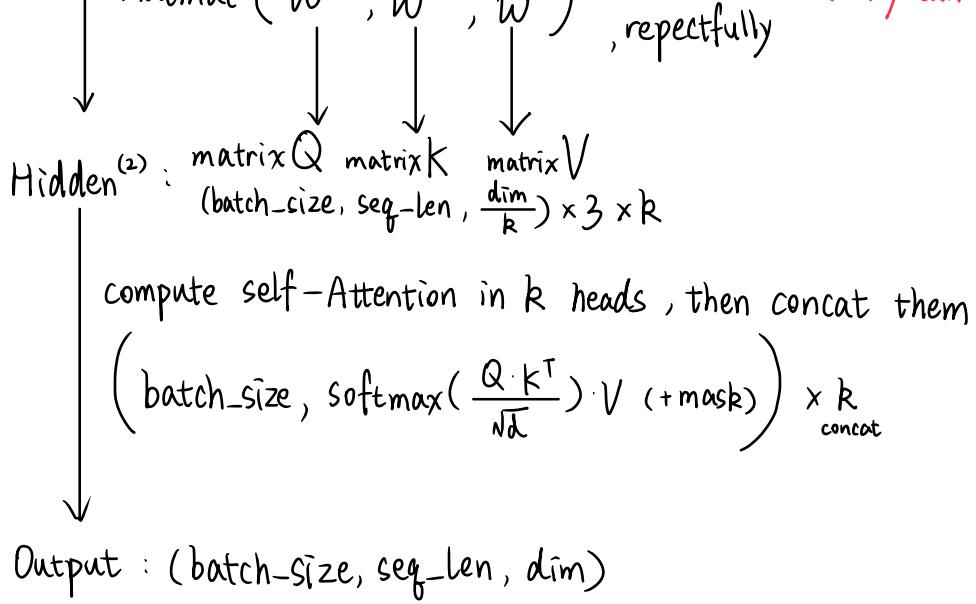


Positional Encoding : 提供时序信息

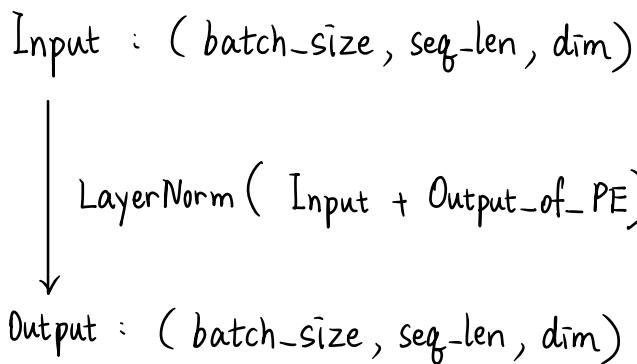


Self-Attention  $\tilde{P}_E$  (mask/None mask) :

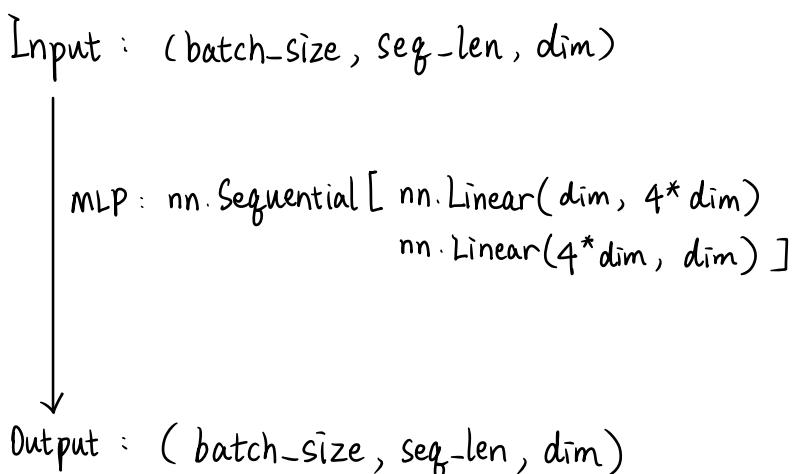




Add & Norm 



feed forward 



Cross-Attention 

$Q$  is Decoder's output, while  $K-V$  is Encoder's output



