

Multi-Head Self-Attention（从空间角度解释为什么做多头）

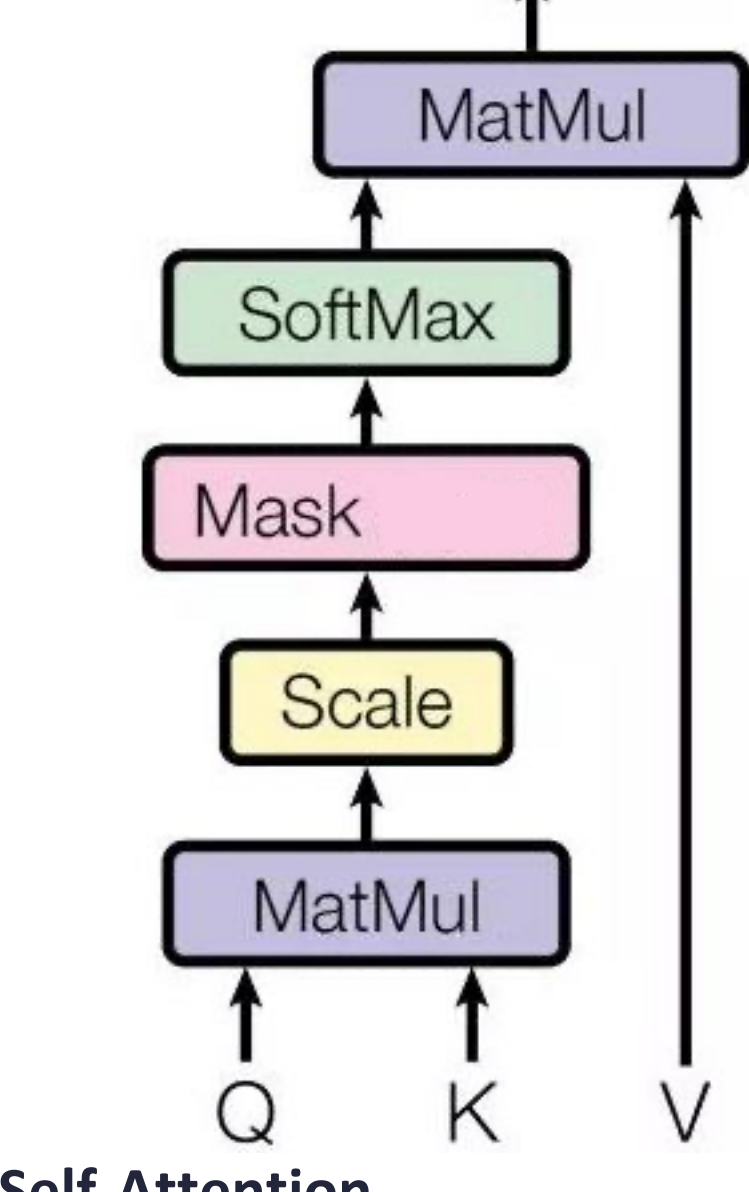
2022年9月3日 星期六 17:30

随笔 - 832 文章 - 0 评论 - 313 阅读 - 151万

上节课回顾

0: 40

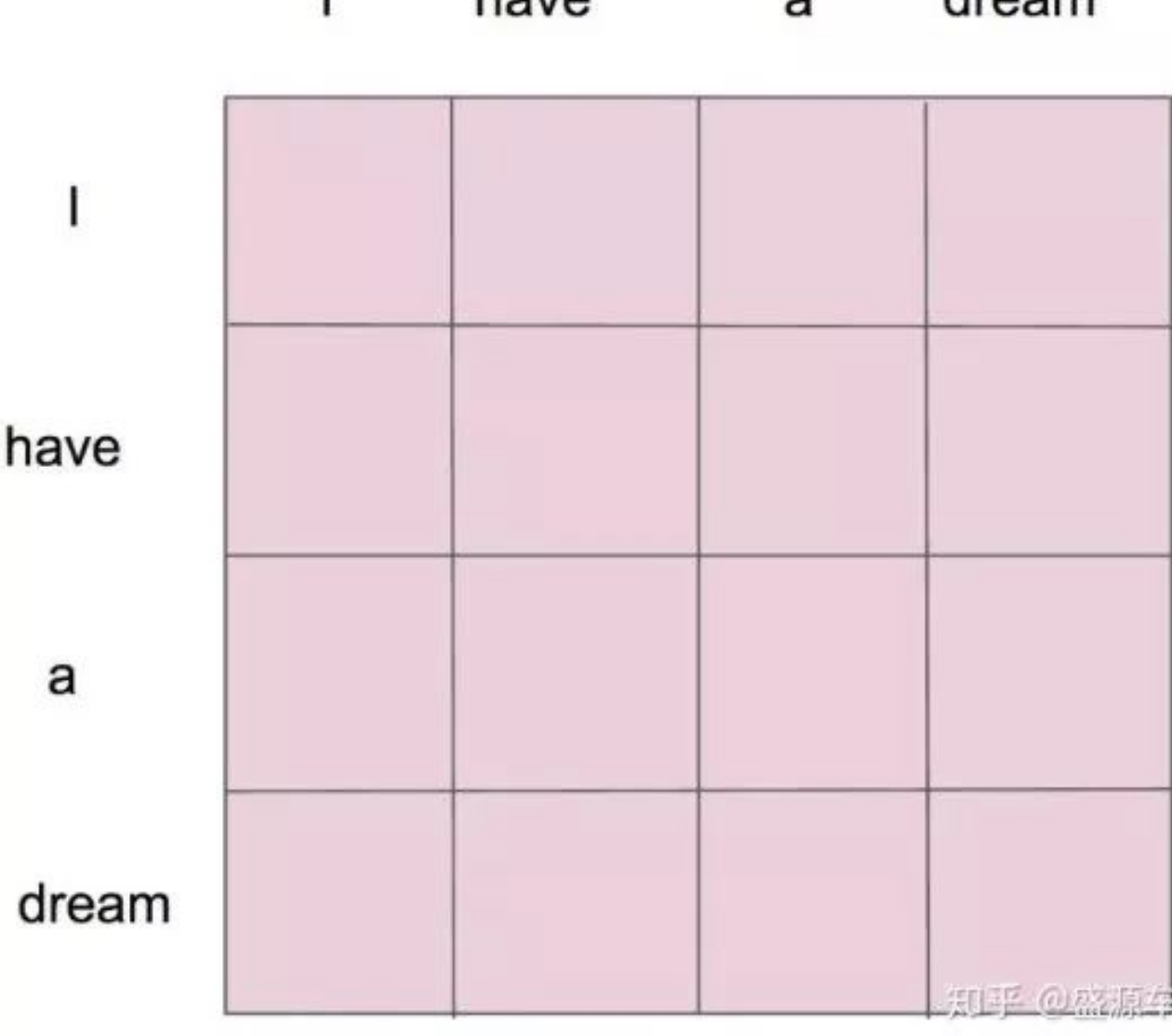
Attention



Self-Attention

Self-Attention 其实是 Attention 的一个具体做法

给定一个 X，通过自注意力模型，得到一个 Z，这个 Z 就是对 X 的新的表征（词向量），Z 这个词向量相比较 X 拥有了句法特征和语义特征



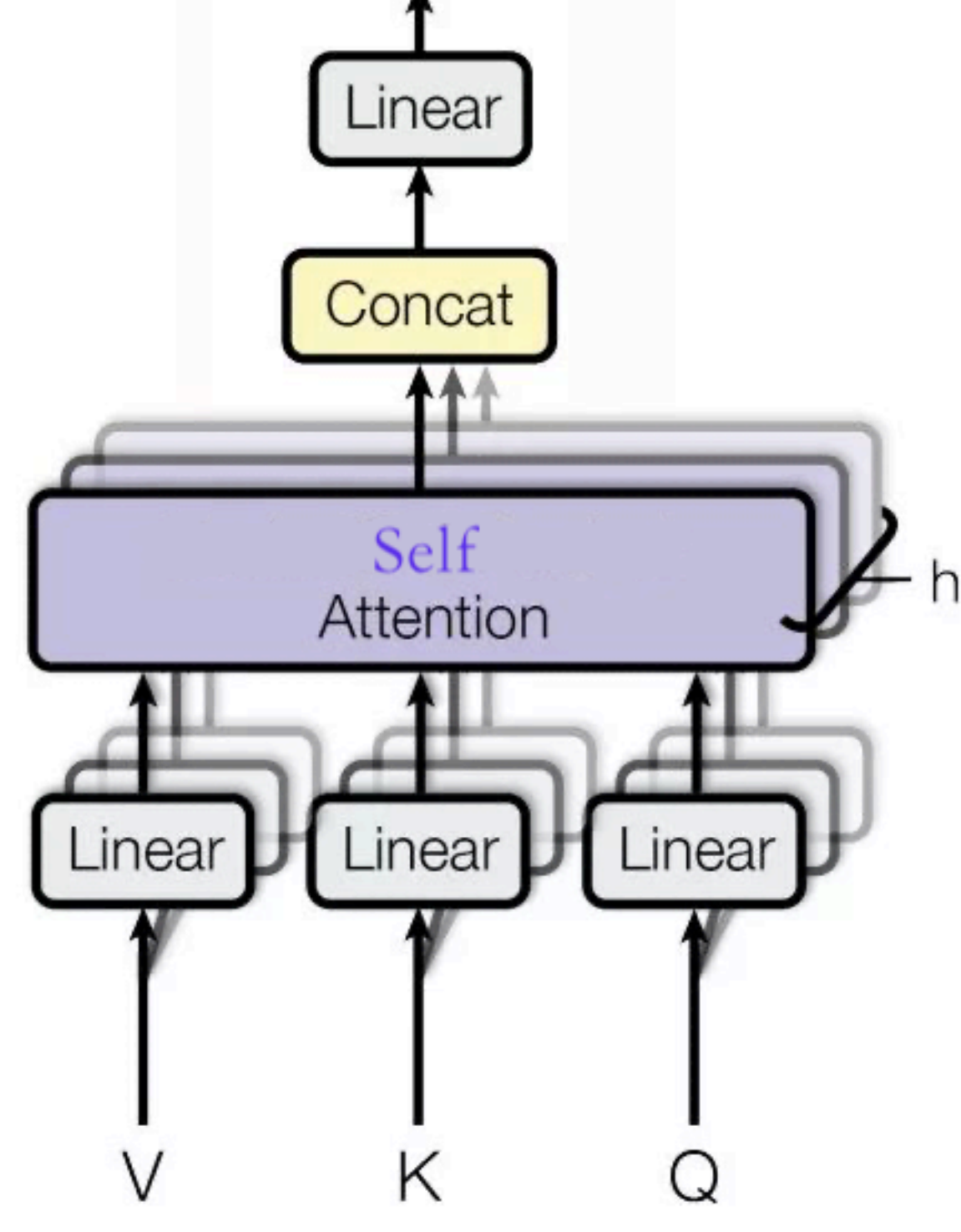
Multi-Head Self-Attention（多头自注意力）

Z 相比较 X 有了提升，通过 Multi-Head Self-Attention，得到的 Z' 相比较 Z 又有了进一步提升

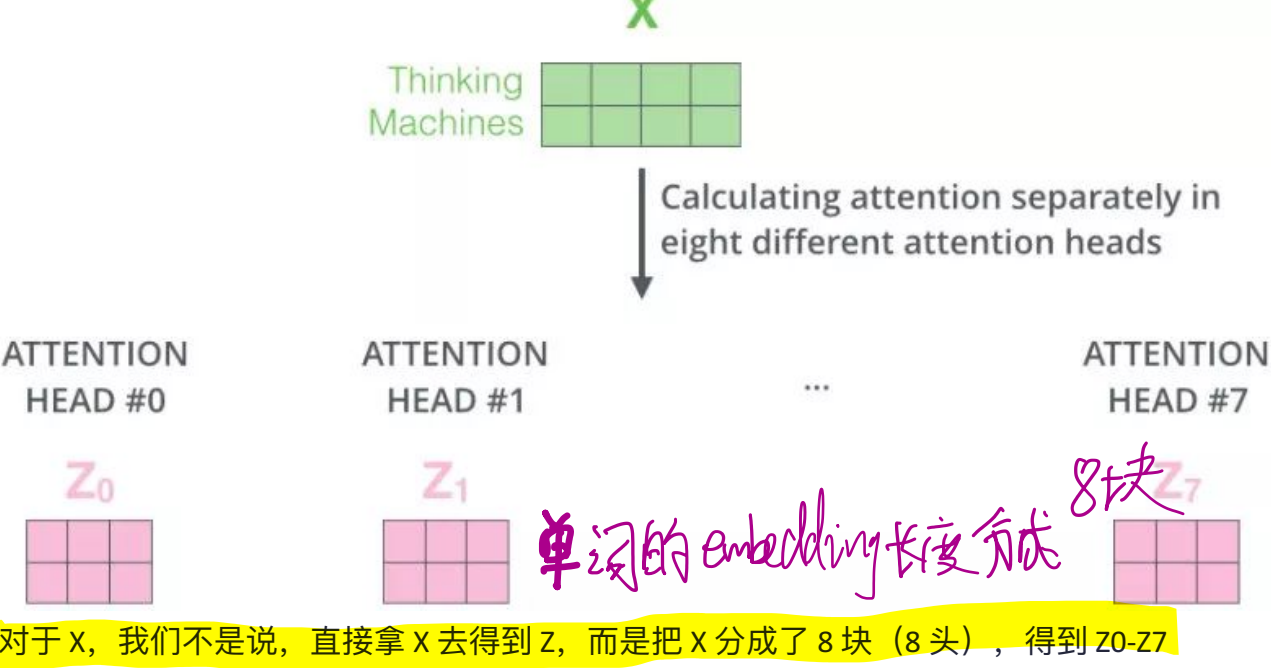
多头自注意力，问题来了，多头是什么，多头的个数用 h 表示，一般 $h=8$ ，我们通常使用的是 8 头自注意力

什么是多头

Multi-Head Attention



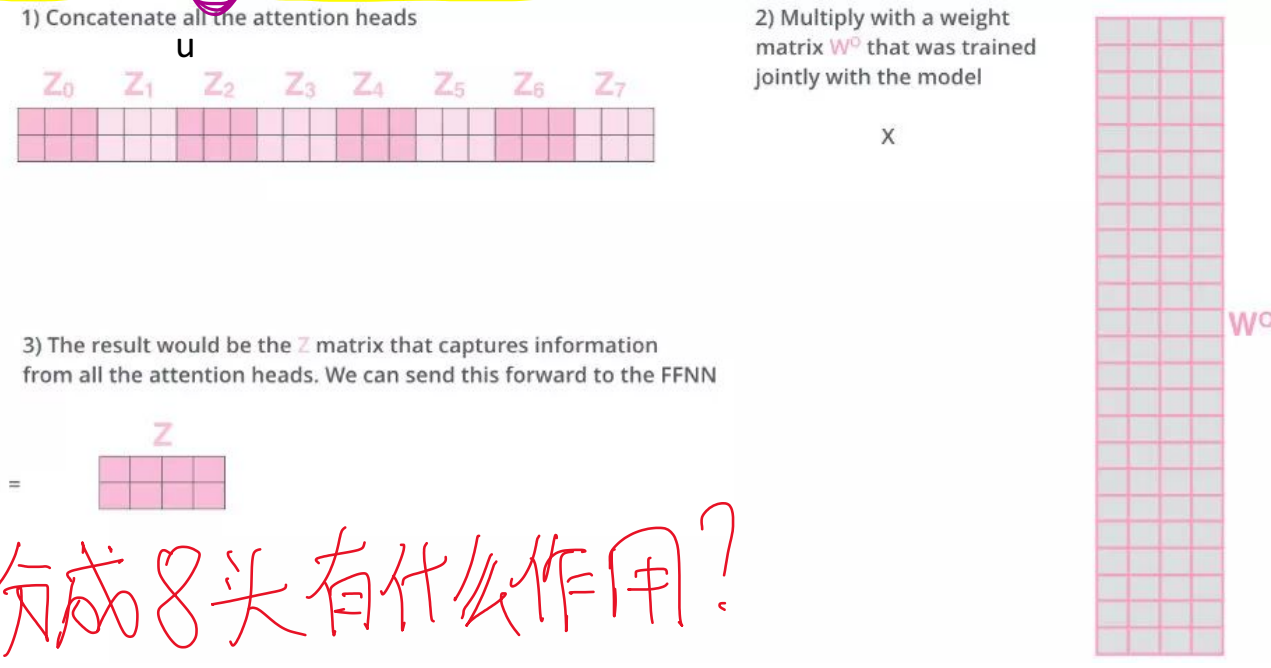
如何多头 1



对于 X，我们不是说，直接拿 X 去得到 Z，而是把 X 分成了 8 块（8 头），得到 Z0-Z7。

如何多头 2

然后把 Z0-Z7 拼接起来，再做一次线性变换（改变维度）得到 Z



分成 8 头有什么作用？

机器学习的本质是什么： $y = \sigma(wx+b)$ ，在做一件什么事情，非线性变换（把一个看起来不合理的东西，通过某个手段（训练模型），让这个东西变得合理）

非线性变换的本质又是什么？改变空间上的位置坐标，任何一个点都可以在维度空间上找到，通过某个手段，让一个不合理的点（位置不合理），变得合理

这就是词向量的本质 词向量的本质

one-hot 编码 (0101010)

word2vec (11, 222, 33)

emlo (15, 3, 2)

attention (124, 2, 32)

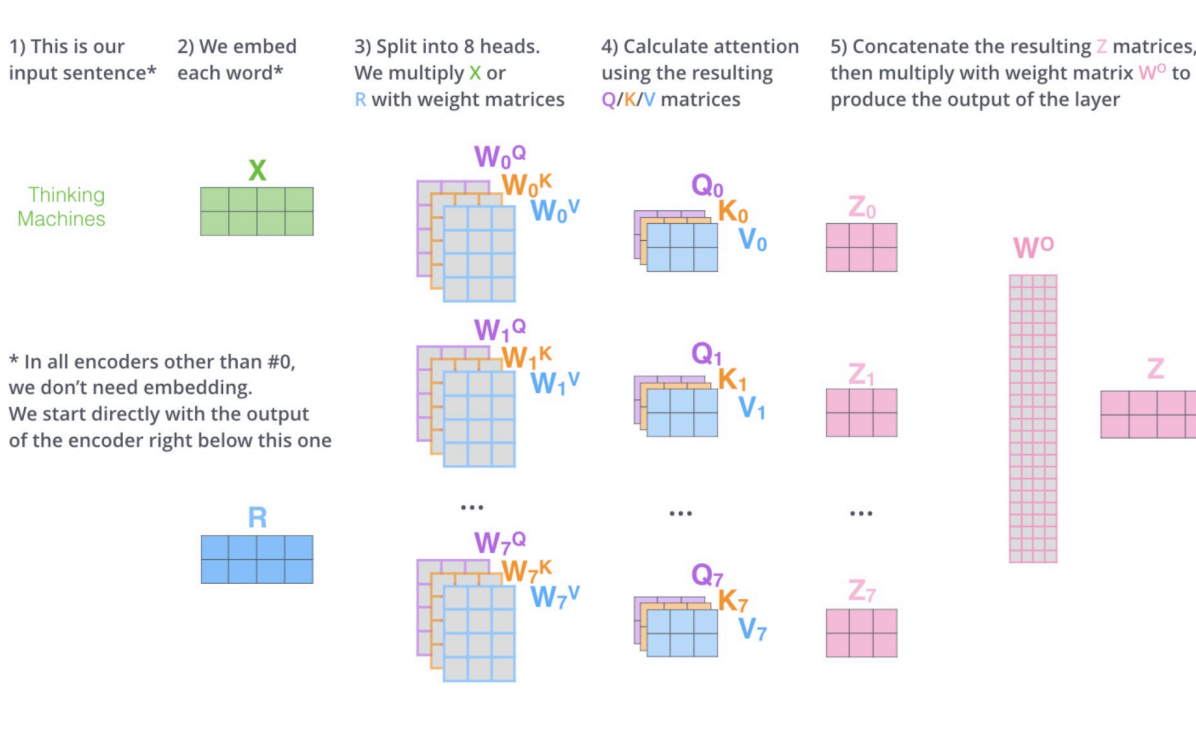
multi-head attention (1231, 23, 3)，把 X 切分成 8 块（8 个子空间），这样一个原先在一个位置上的 X，去了空间上 8 个位置，通过对 8 个点进行寻找，找到更合适的位置

词向量的大小是 512

假设你的任务，视频向量是 5120, 80

对计算机的性能提出了要求

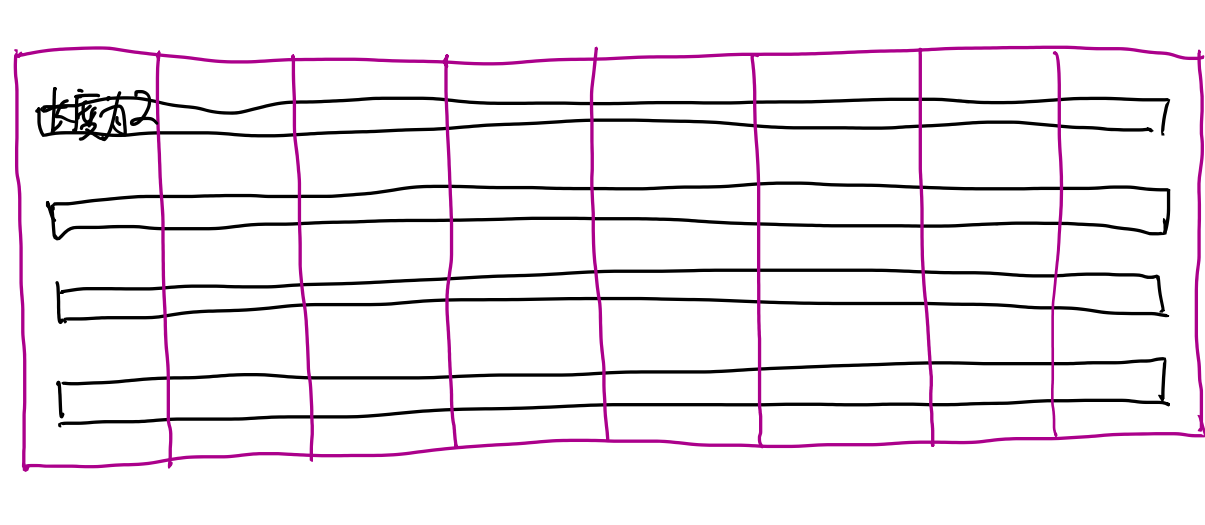
多头流程图



* In all encoders other than #0, We start directly with the output of the encoder right below this one

分割(下面是斜率丰例子)

假设，batch=3，词嵌入长度为16。



$[3, 4, 16]$

$\downarrow \text{view}(\text{batch}, -1, \text{head}, \text{dim-v})$

$[3, 4, 8, 2]$

$\downarrow \text{transpose}(-2, -1)$

$[3, 8, 4, 2]$

\downarrow 多头，self-attention

$[3, 8, 4, 2]$

\downarrow (拼接多头结果) $\text{view}(\text{batch}, -1, \text{head} \times \text{dim-v})$

$[3, 4, 16]$. Multi-Head-self Attention 输出