

碎语

1. 本人先看了下面的视频
2. 阅读了《Attention is all you need》发现文章内容和视频讲的基本一样，细节论文不清楚，。
3. 所以，又看了一遍视频。

视频笔记

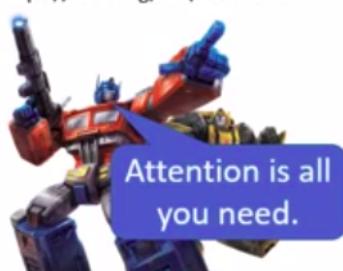
<https://www.bilibili.com/video/av56239558?from=search&seid=4166606460205828626>

self-attention

得到: q, k, v

Self-attention

<https://arxiv.org/abs/1706.03762>



q : query (to match others)

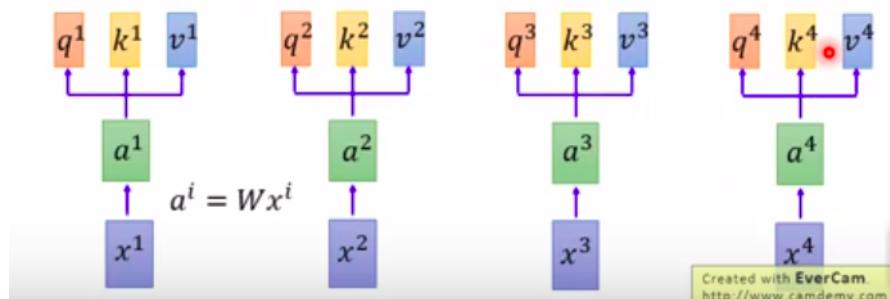
$$q^i = W^q a^i$$

k : key (to be matched)

$$k^i = W^k a^i$$

v : information to be extracted

$$v^i = W^v a^i$$

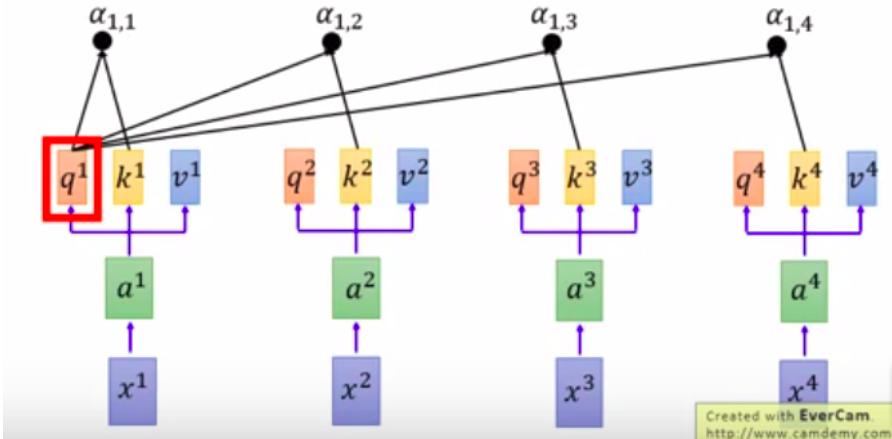


接下来，计算attention

Self-attention

拿每個 query q 去對每個 key k 做 attention

$$\text{Scaled Dot-Product Attention: } \alpha_{1,i} = \underbrace{q^1 \cdot k^i}_{\text{dot product}} / \sqrt{d}$$



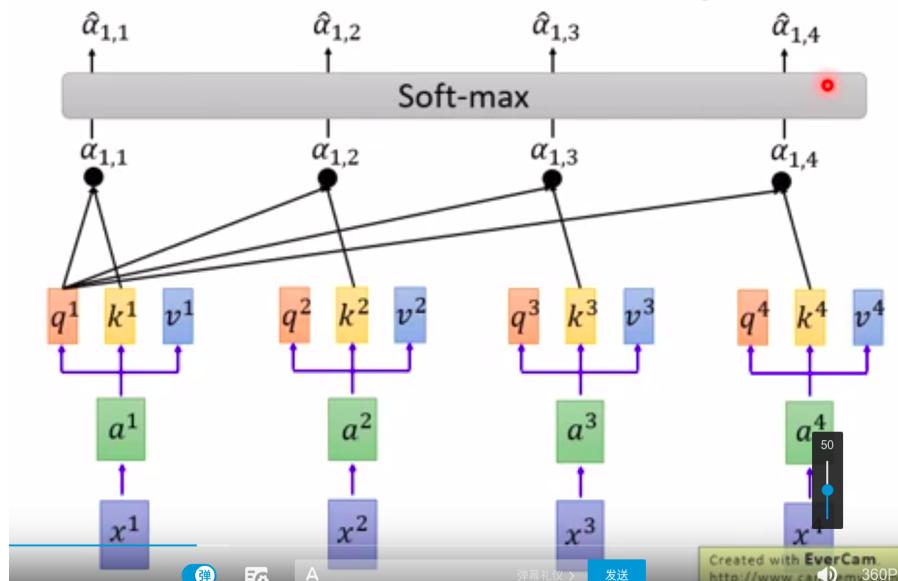
Created with EverCam,
<http://www.camdemey.com>

dot product: 点积 (attention is all you need 文章中, 用此方法计算相似性)。

为什么除以根号d: 因为点积是相乘求和, 维度越大, 相应数值越大。

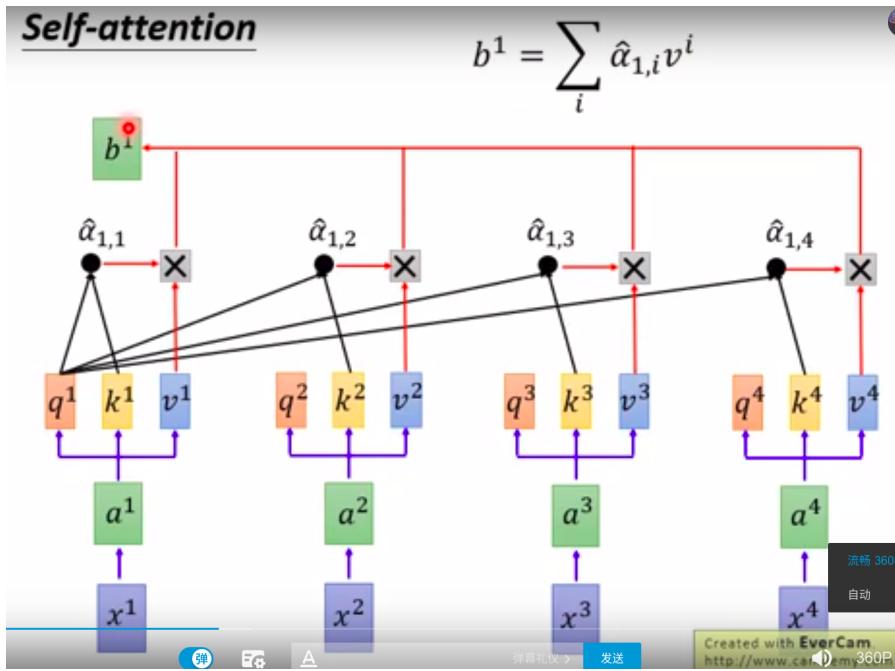
接下来, 数值变概率 (softmax)

$$\text{Self-attention} \quad \hat{\alpha}_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$

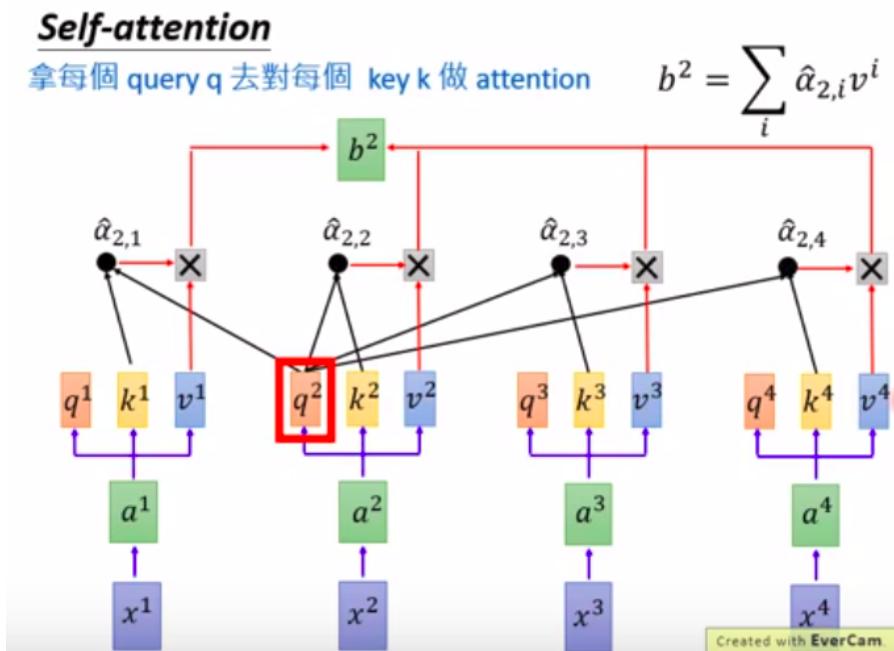


Created with EverCam
<http://www.camdemey.com> 360P

第一个输出: (这个输出, 用了全部输出的信息)



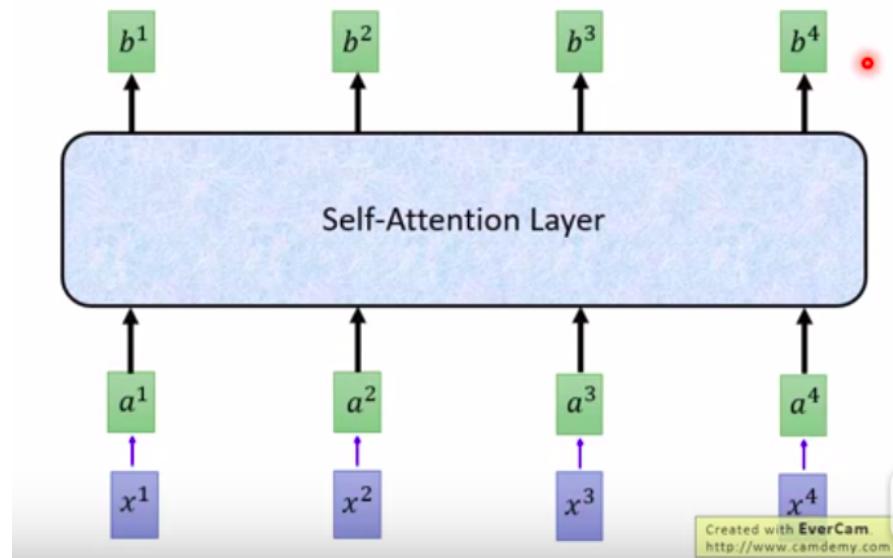
第二个输出:



整体展示:

Self-attention

b^1, b^2, b^3, b^4 can be parallelly computed.



矩阵运算平行化

Self-attention

$$q^1 \boxed{q^2} \boxed{q^3} \boxed{q^4} = \boxed{W^q} \boxed{a^1 a^2 a^3 a^4}$$

Q I

$$q^i = W^q a^i$$

$$\boxed{k^1} \boxed{k^2} \boxed{k^3} \boxed{k^4} = \boxed{W^k} \boxed{a^1 a^2 a^3 a^4}$$

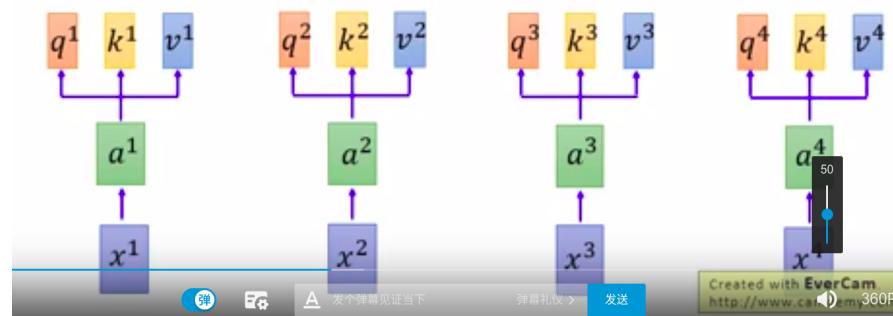
K I

$$k^i = W^k a^i$$

$$v^i = W^v a^i$$

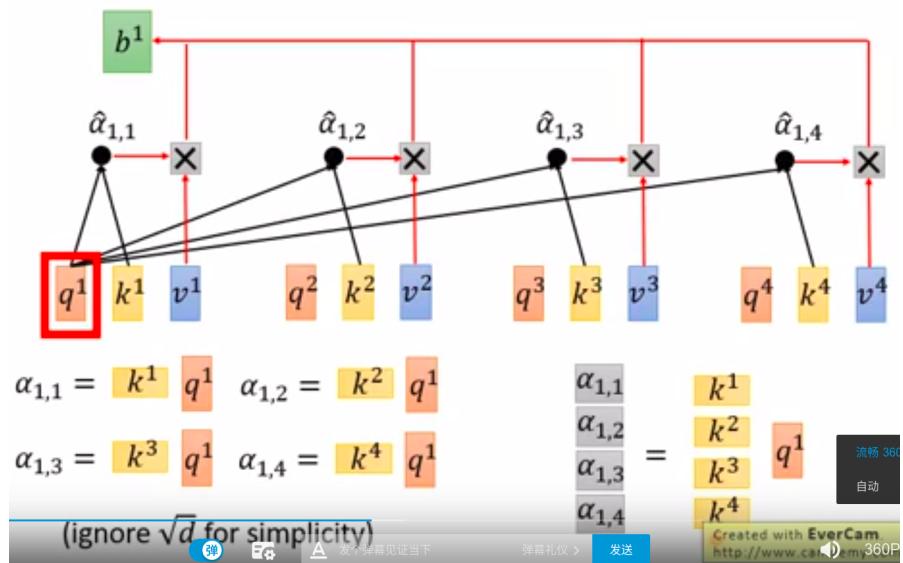
$$\boxed{v^1} \boxed{v^2} \boxed{v^3} \boxed{v^4} = \boxed{W^v} \boxed{a^1 a^2 a^3 a^4}$$

V I

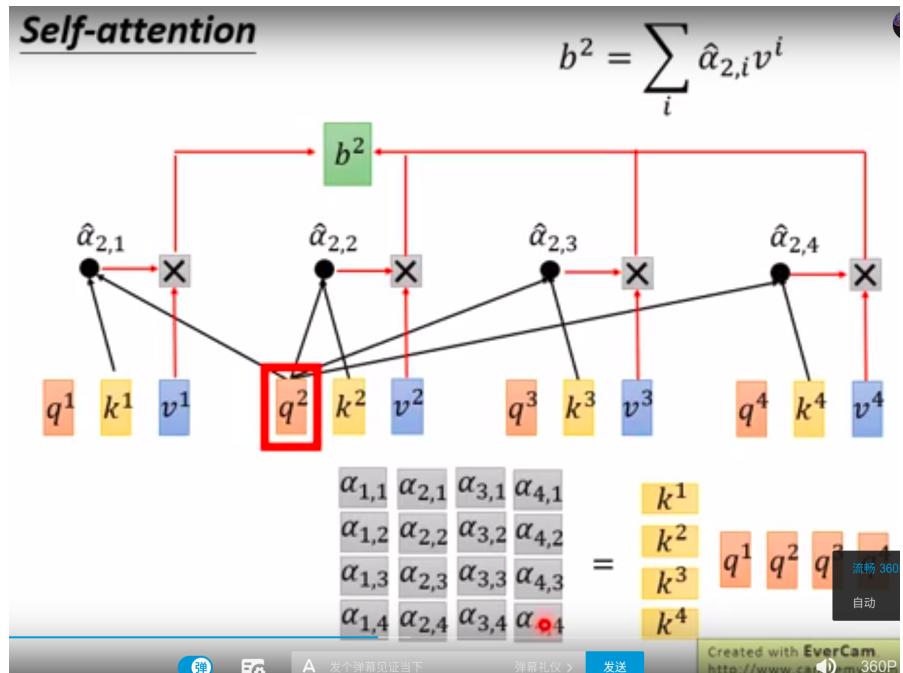


一个

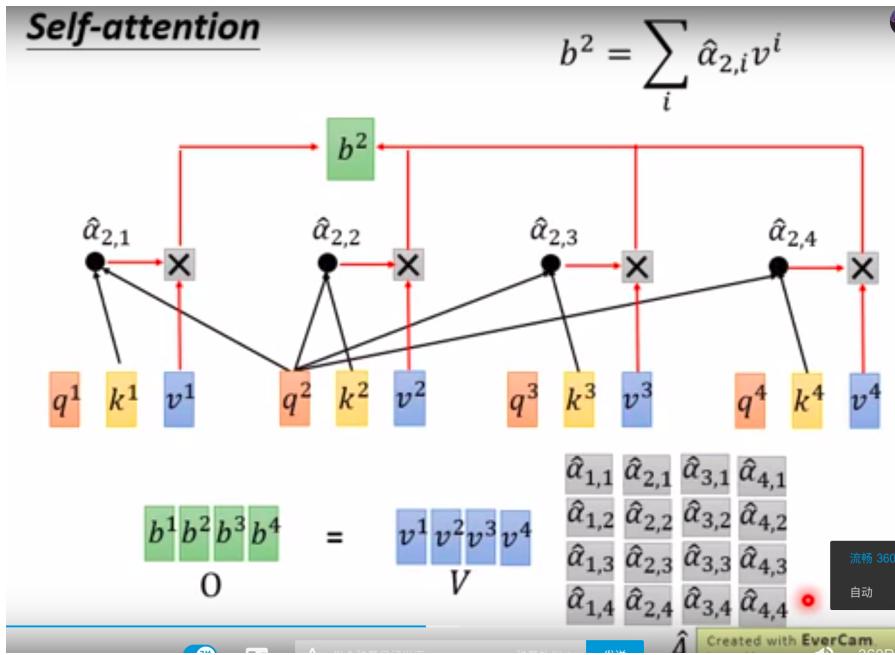
Self-attention



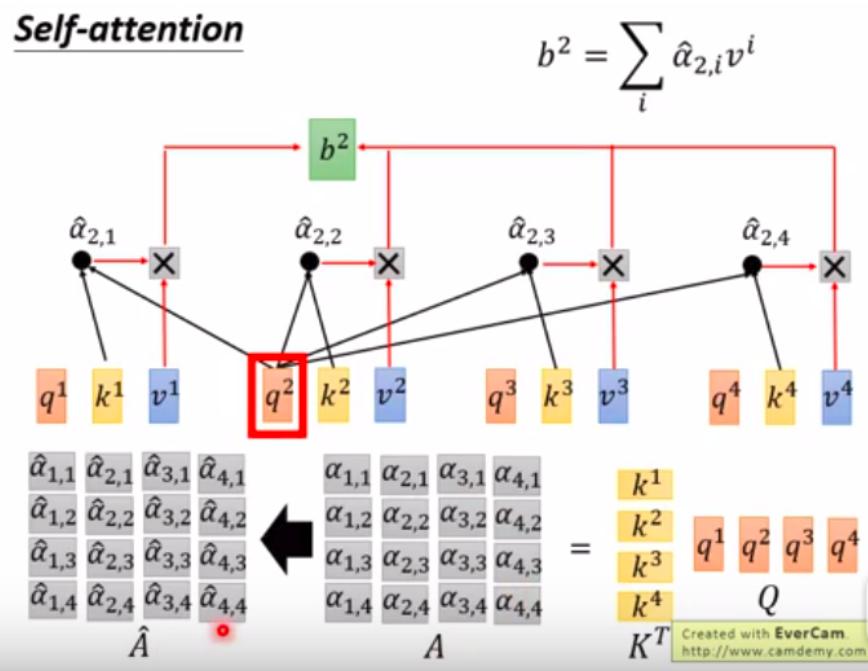
第二个



整个过程，可以转化为上面



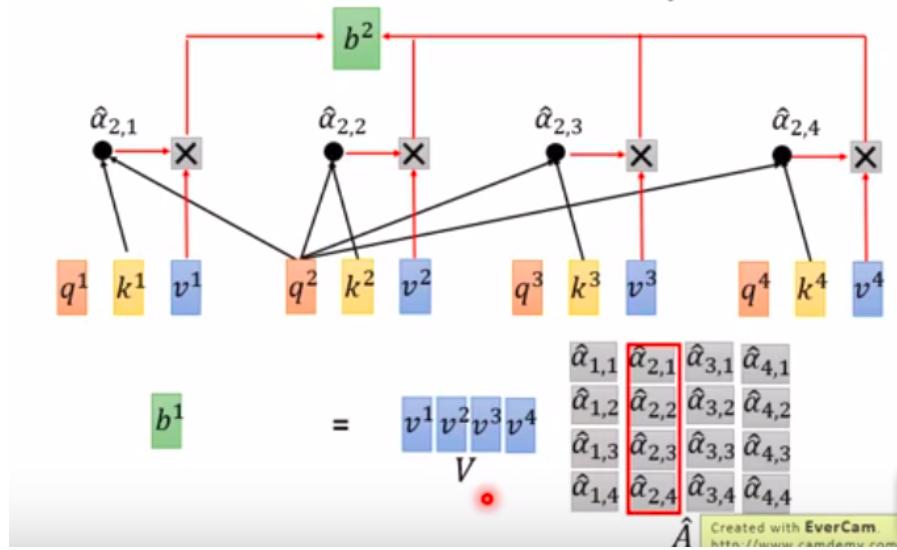
softmax



输出求解

Self-attention

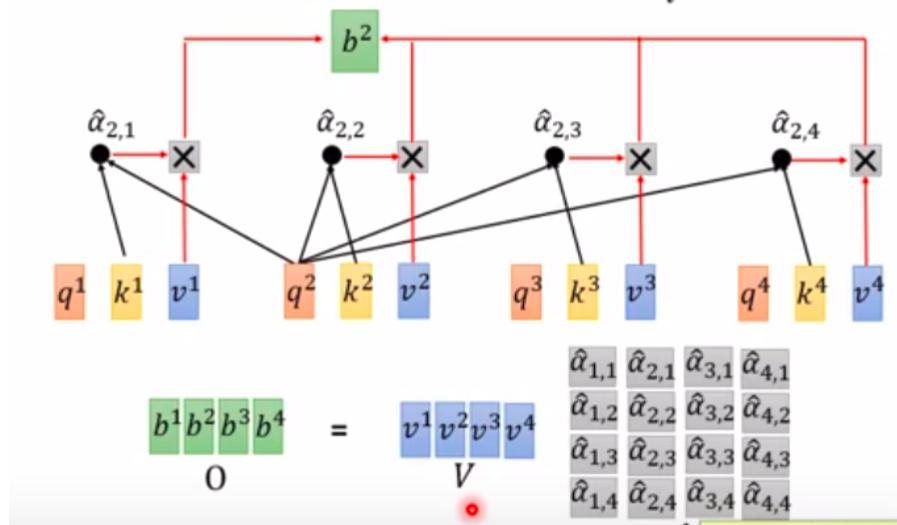
$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$



所有输出

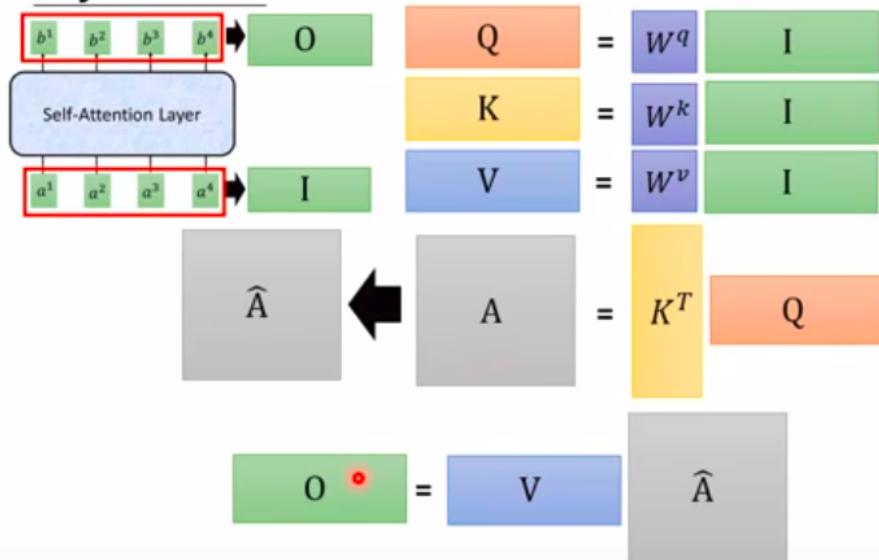
Self-attention

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$



换个形式再看一遍

Self-attention



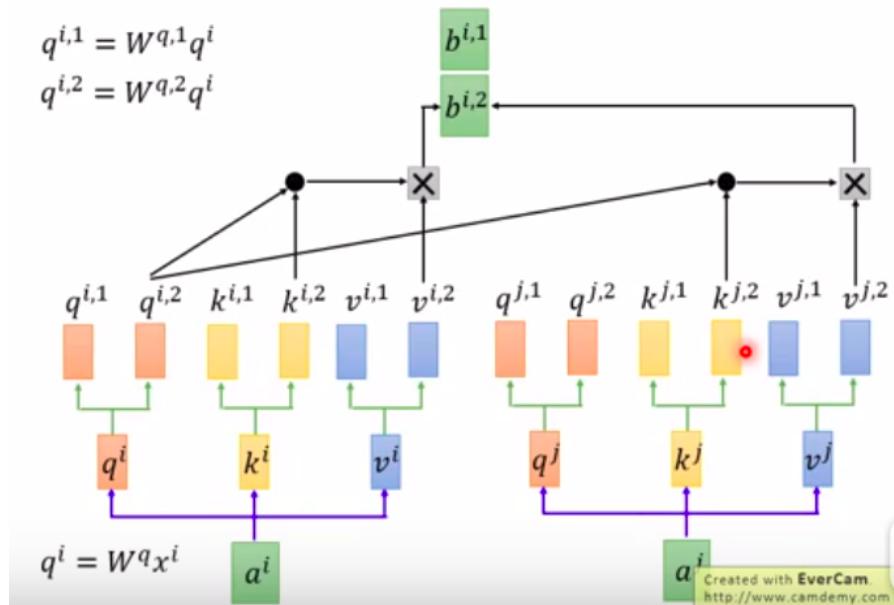
反正就是一堆矩阵乘法 · 用 GPU 可以加速

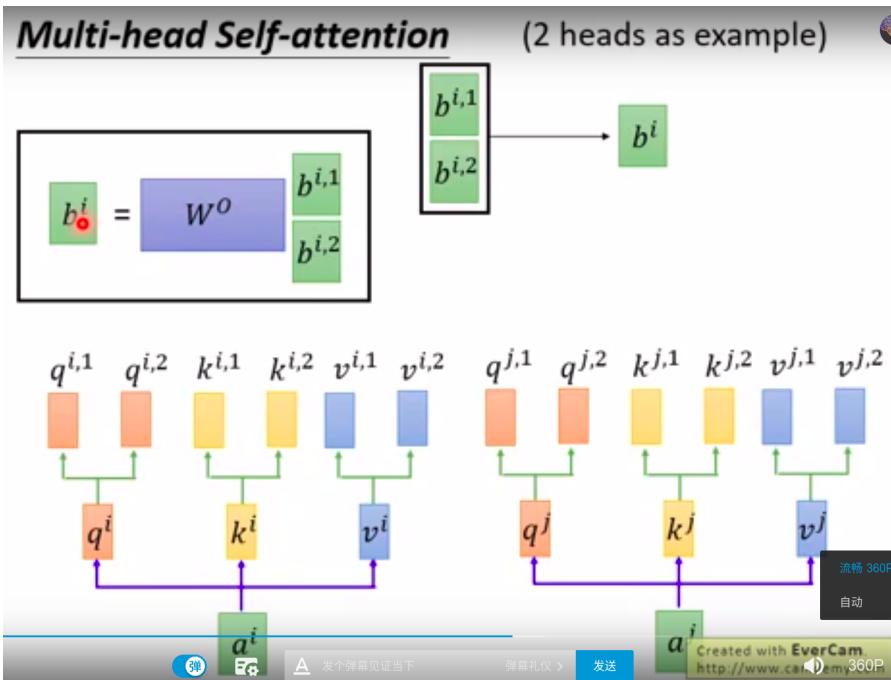
Created with EverCam.
http://www.camdemmy.com

重点: GPU 加速矩阵乘法

Multi-head self-attention

Multi-head Self-attention (2 heads as example)



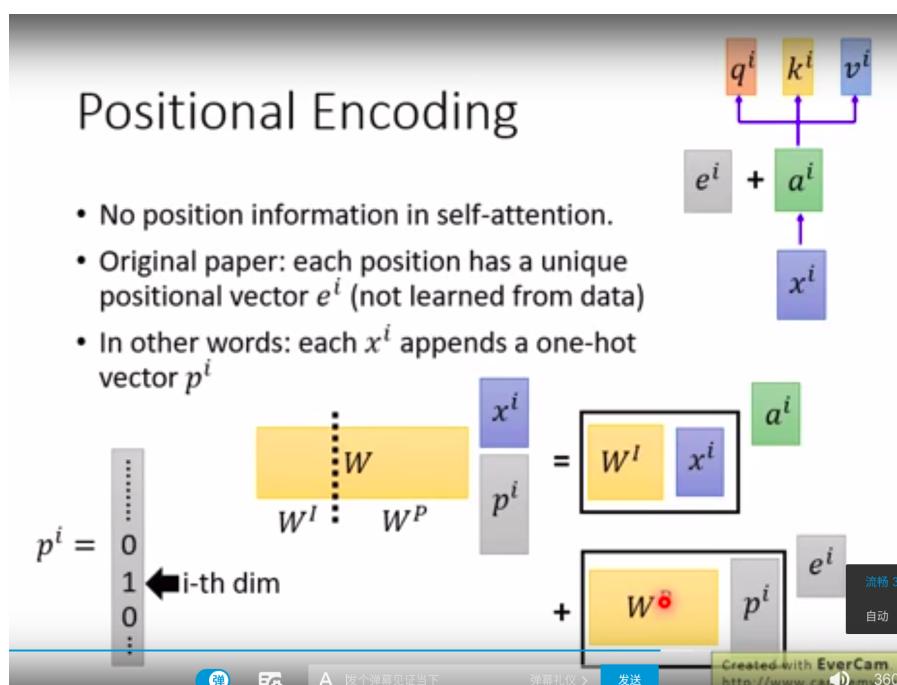


多头有什么好处？

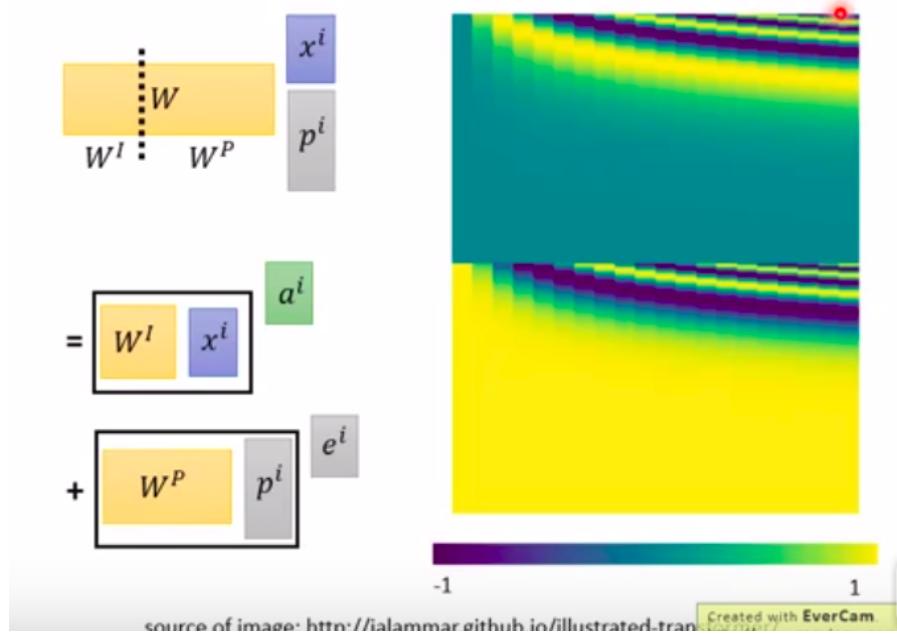
不同的头，关注可以不同，多个关注点。论文里说，避免了原来 attention 的平均化。

Position Encoding

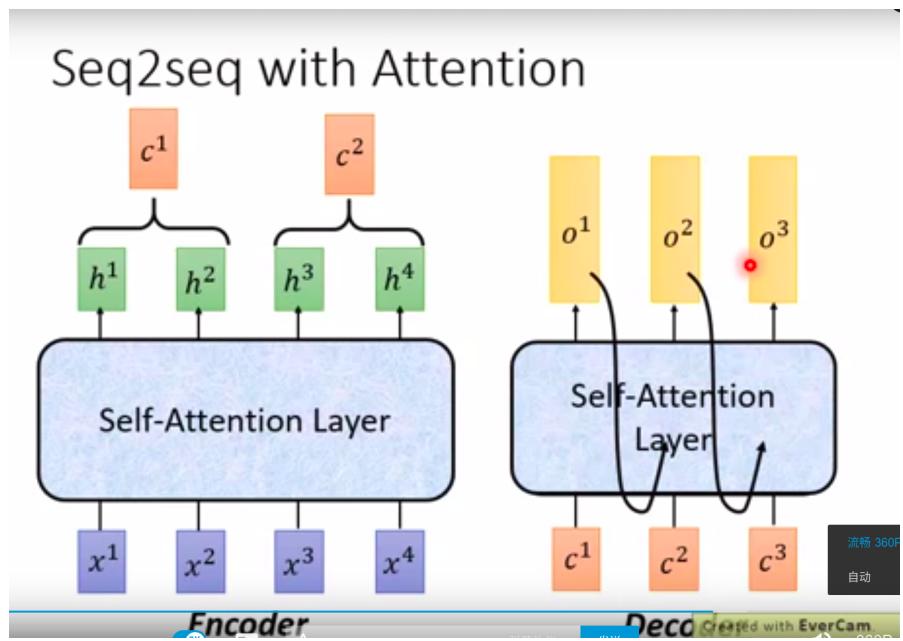
注意，原有self-attention对位置不敏感，没有位置信息



1. ei 是人根据式子人工设的。
2. 为什么 ei 和 ai 是相加, 不是concat进 x 里? 看上面公式计算, 计算后他们等价。
3. 有人试过, learn Wp , 但是并不好, 所以有人, 通过式子人工手设 wp 。如下图:

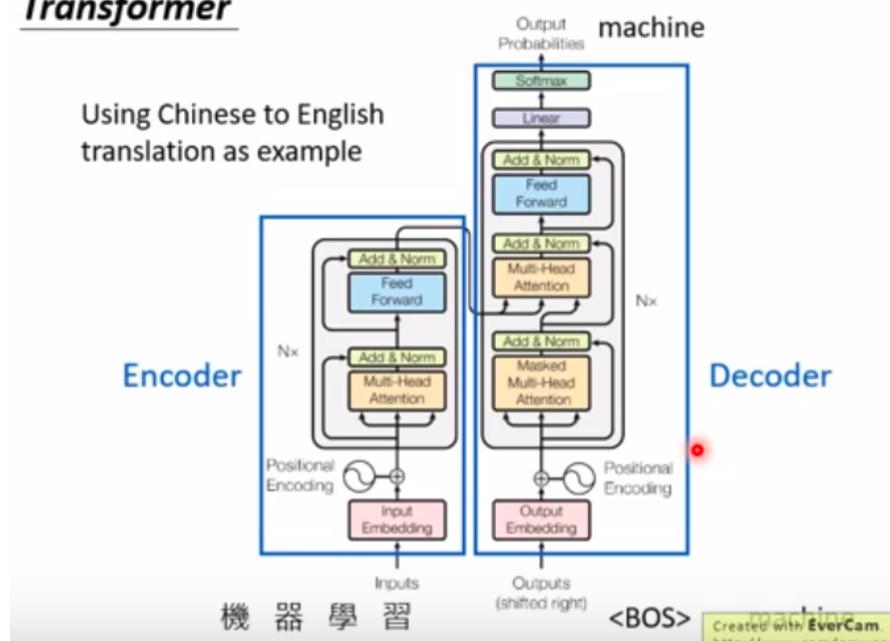


self-attention 在seq2seq

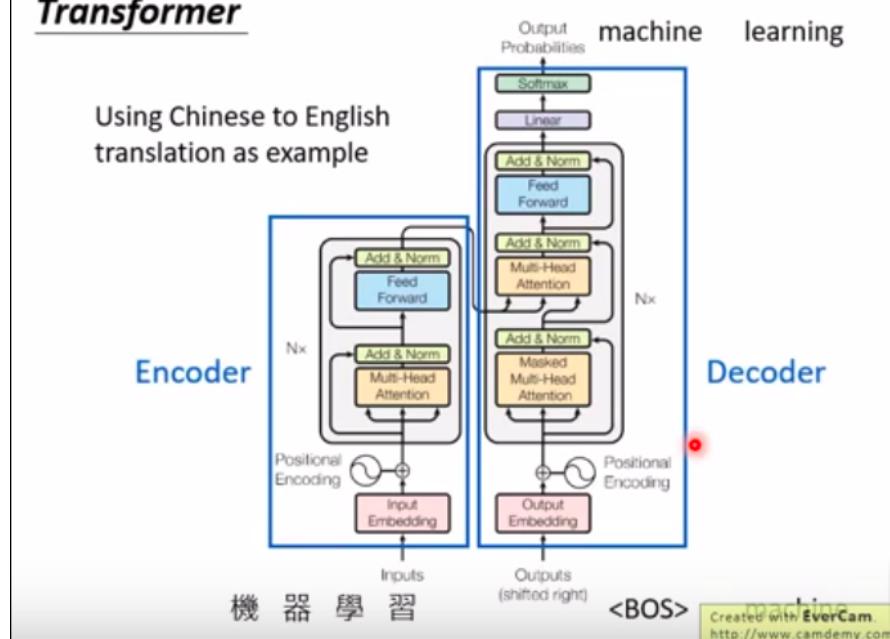


Transformer

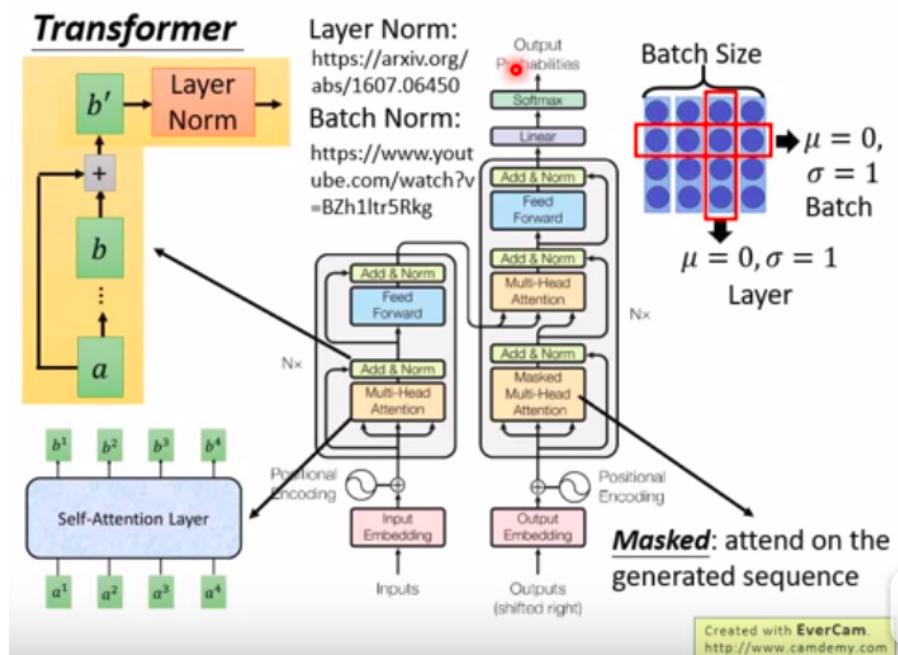
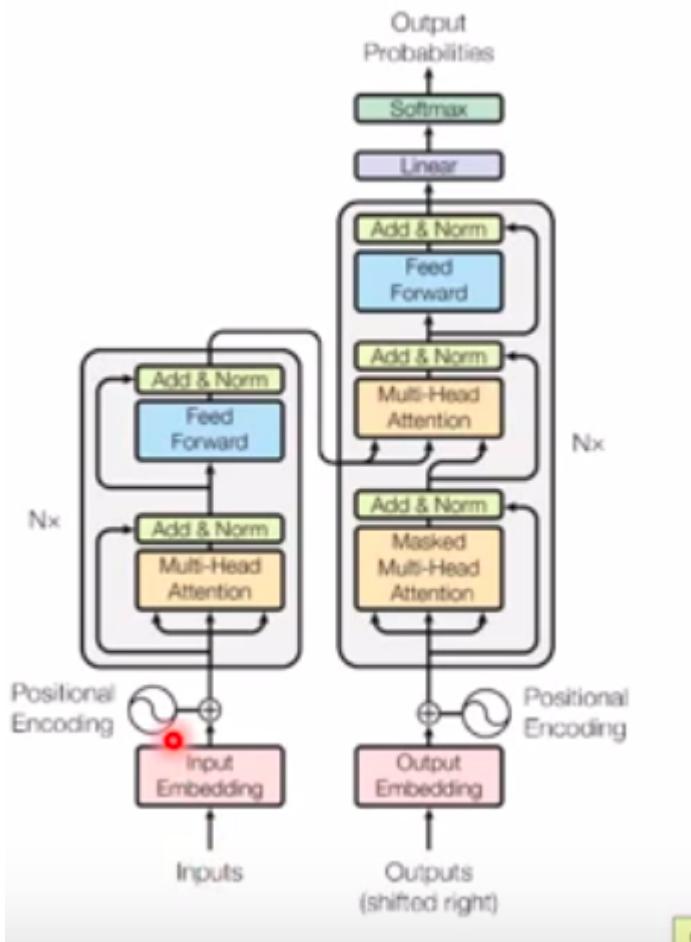
Transformer



Transformer



分析layer



encoder

1. $input \rightarrow embedding = vctor | + position \rightarrow$ 最后输出
2. 进入灰色块 (循环N次)

3. 灰色第一层 (*multi-head*) : 输入一个sequence得到另一个等长的sequence。
4. 下一层: (如图左上) ADD: *multi_input + mulit_output*
—> *layernorm*.
5. *batch_norm*和*layer_norm*区别: (右上角) : *batch*横向
(注意, 横向是*batch_size*) , 整个batch里同一维度做
norm, *layer*相反, 不考虑batch, 一个样本内做*norm*, 一般
*layer_norm*搭配RNN。
6. 下一层: *feed Forward*.
7. 最后 ADD & Norm

decoder

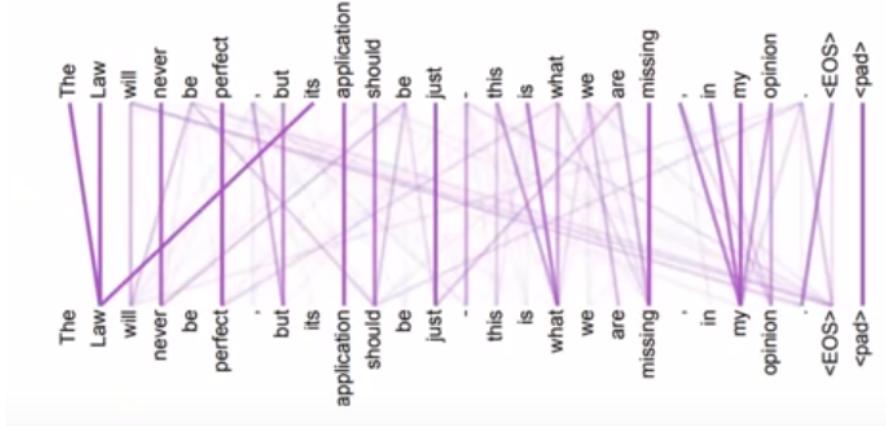
1. *input->* 前一个timestep产生的output。
2. *embeddding+postion*进入灰色block。

1. 灰色第一层: *Mask_multi,mask*表示只attention产生出来的句子。
2. *Add & Norm*
3. 下一层*multi_attention*对的是encoder的输出。
4. *Add & Norm*
5. *feed*
6. *Add & Norm*

1. *linner*
2. *softmax*

Attention visualization

Attention Visualization

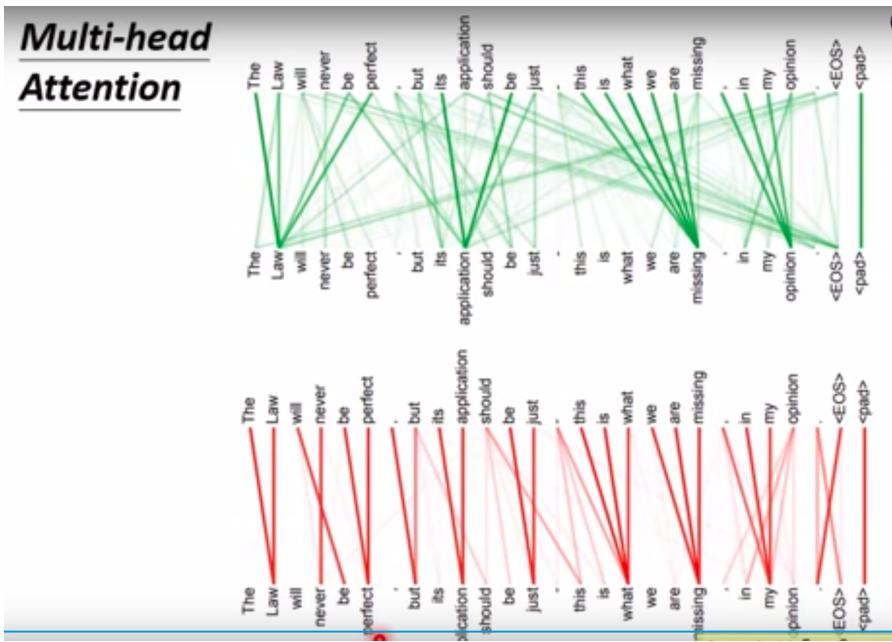


Attention Visualization



The encoder self-attention distribution for the word "it" from the 5th to the 6th layer of a Transformer trained on English to French translation (one of eight attention heads).

Multi-head Attention



上面那个是：全局attention

下面那个是：局部attention

Transformer 例子

第一个(wiki)

摘要：

自动写一个wiki

Example Application

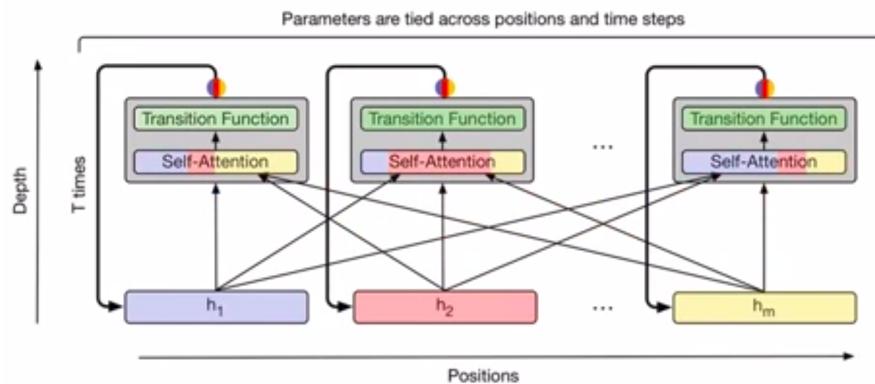
- If you can use seq2seq, you can use transformer.



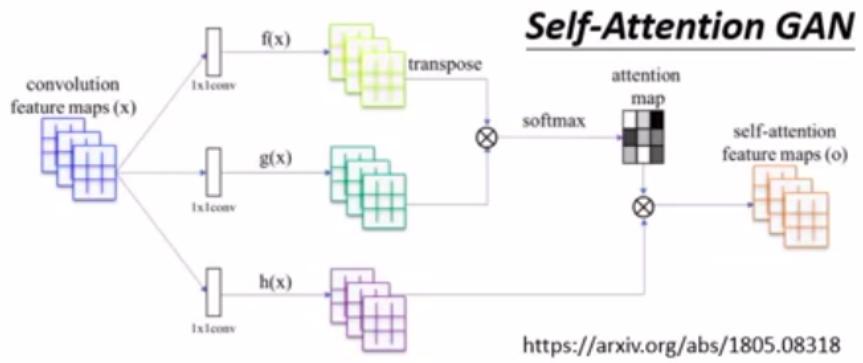
Dataset	Input	Output	# examples
Gigaword (Graff & Cieri, 2003)	10^1	10^1	10^6
CNN/DailyMail (Nallapati et al., 2016)	$10^2\text{--}10^3$	10^1	10^5
WikiSum (ours)	$10^2\text{--}10^6$	$10^1\text{--}10^3$	10^6

第二个 (universal Transformer)

Universal Transformer



第三个 (影像)



<https://arxiv.org/abs/1805.08318>

