

NEURAL NETWORKS AND DEEP LEARNING

Transformer

http://speech.ee.ntu.edu.tw/~tlkagk/courses_ML19.html

李宏毅Transformer教程：

http://speech.ee.ntu.edu.tw/~tlkagk/courses_ML19.html

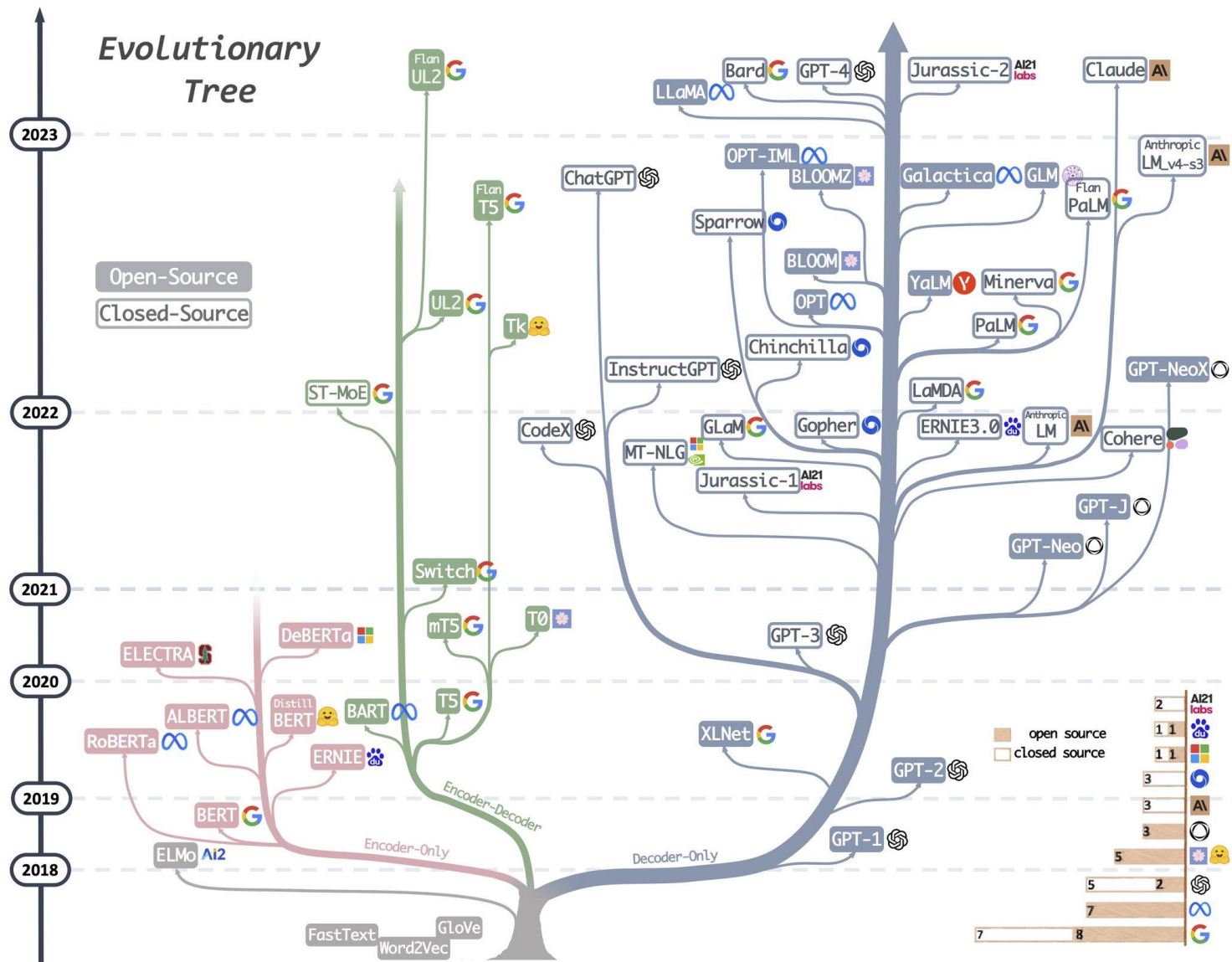
Transformer Blog:

https://slds-lmu.github.io/seminar_nlp_ss20/attention-and-self-attention-for-nlp.html

<https://lilianweng.github.io/posts/2018-06-24-attention/>

<http://jalammar.github.io/illustrated-transformer/>

LLM in NLP



LLM in Time-series

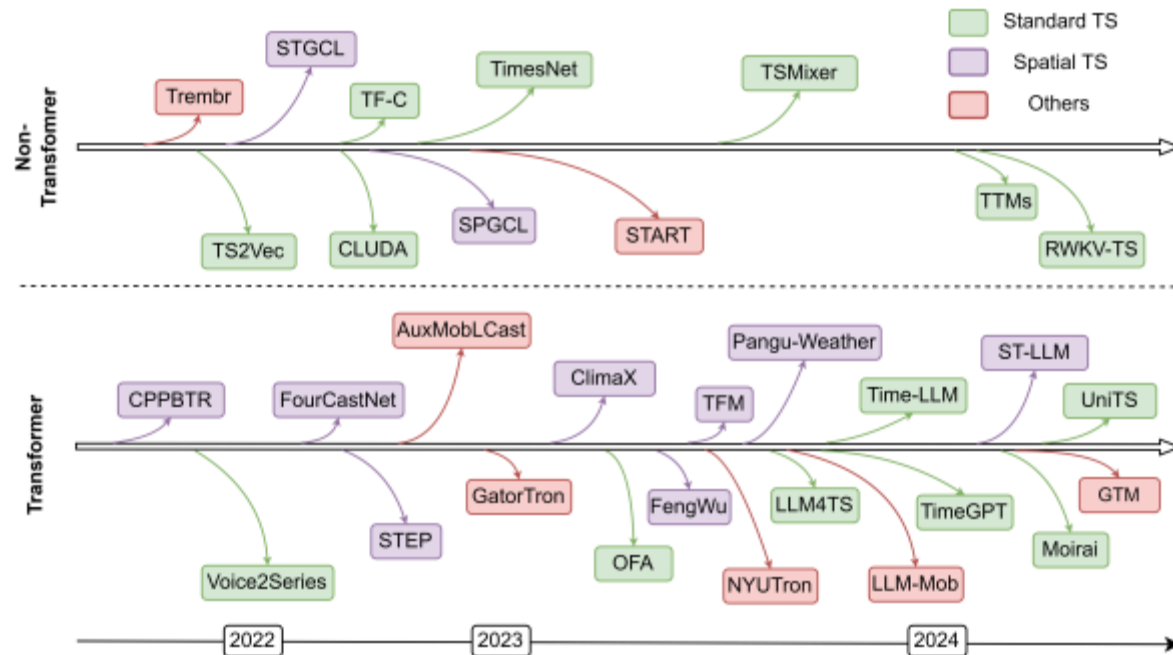


Figure 1: Roadmaps of representative TSFMs.

<https://arxiv.org/pdf/2403.14735>

Transformer in Time-series

- Vanilla Transformer, NeurIPS 2017
- Informer, AAAI 2021
- Autoformer, NeurIPS 2021
- FEDformer, ICML 2022

... (Transformer variants)

- Dlinear, AAAI 2023

... (non-Transformer model)

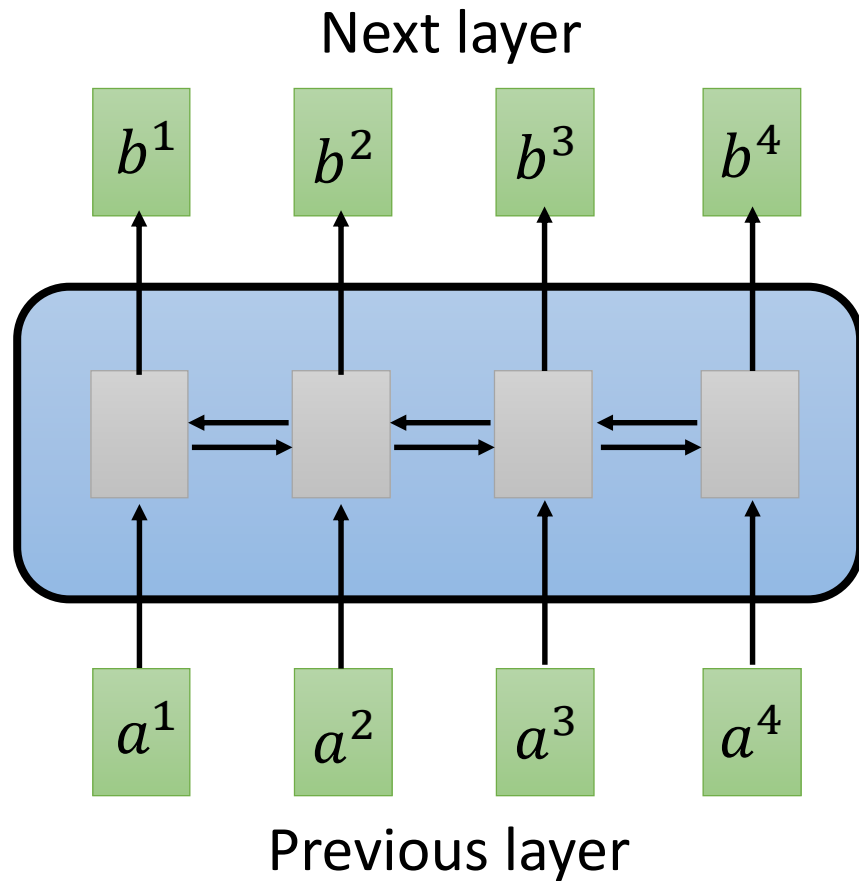
- PatchTST, ICLR 2023

... (Patching based Transformer)

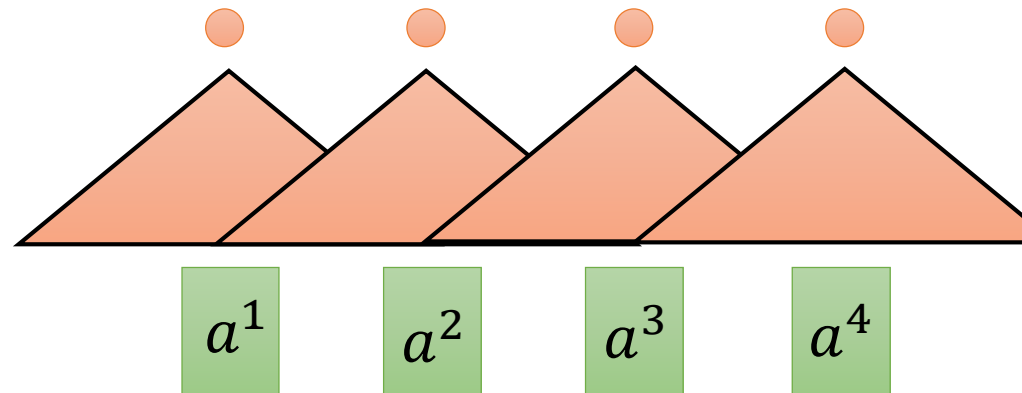
- Time-series Foundation model

... (TSFM)

Sequence

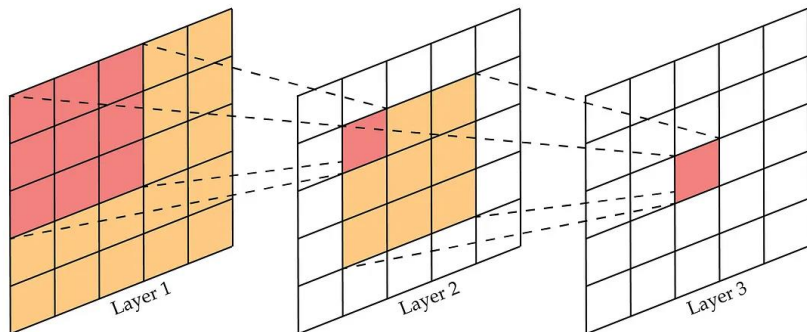


Hard to parallel !



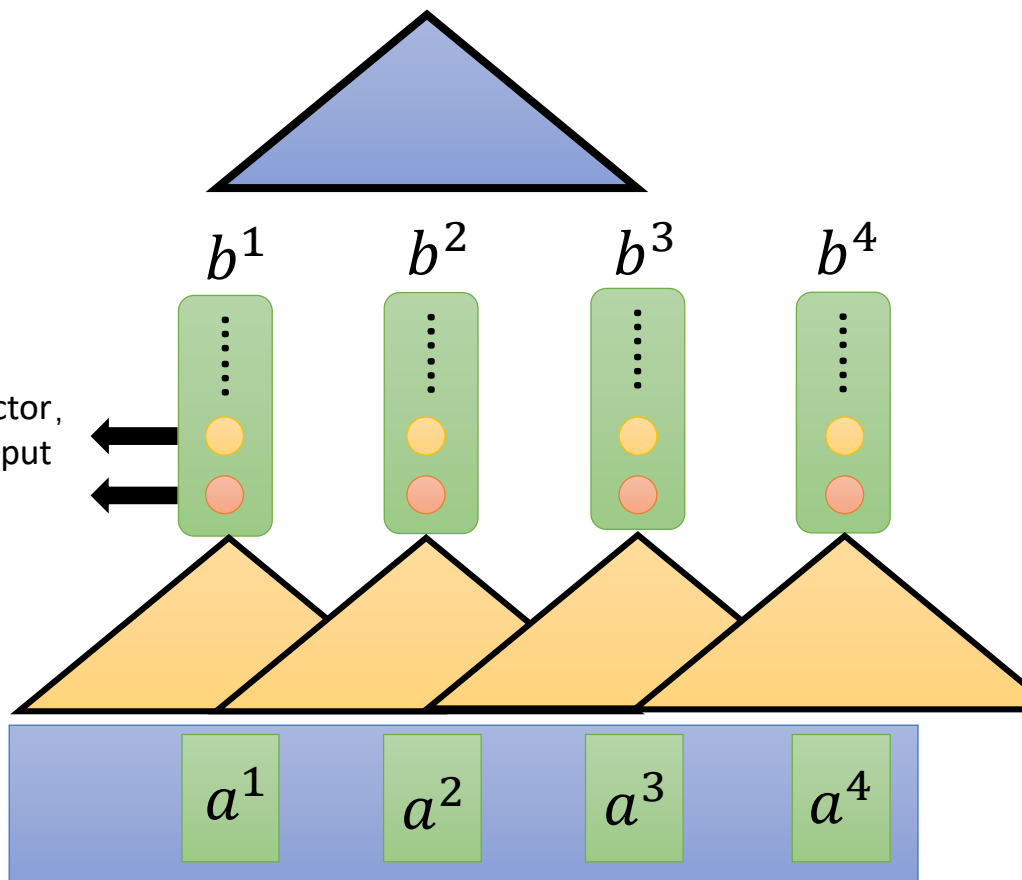
Using CNN to replace RNN
CNN can parallel !

Receptive Field in Convolutional Networks



Filters in higher layer can consider longer sequence

不同卷积核得到的hidden vector,
有多少个卷积核对应 output
hidden vector有多少channel

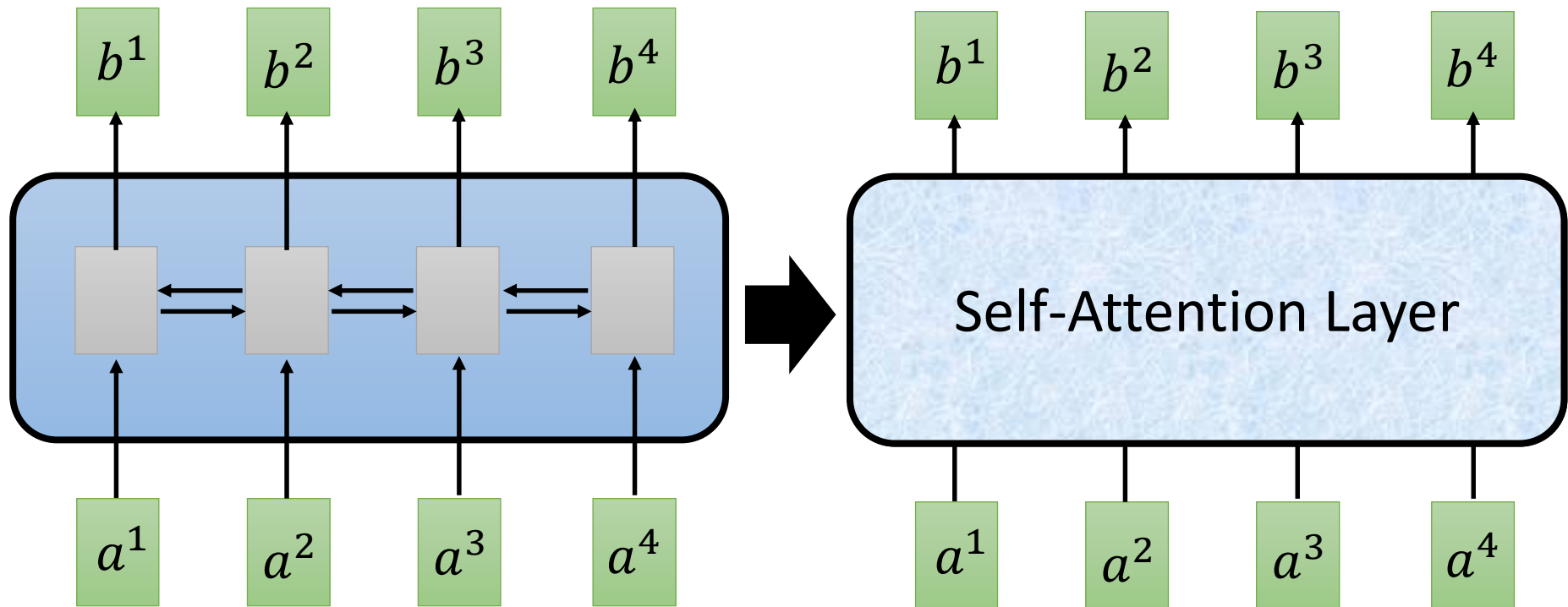


虽然higher layer filters能够捕捉longer sequence, 但lower layer filters始终只能捕捉shorter (局部) sequence

Using CNN to replace RNN
(CNN can parallel)

Self-Attention

b^i is obtained based on the whole input sequence. 与双向RNN相同
 b^1, b^2, b^3, b^4 can be parallelly computed. 与RNN不同



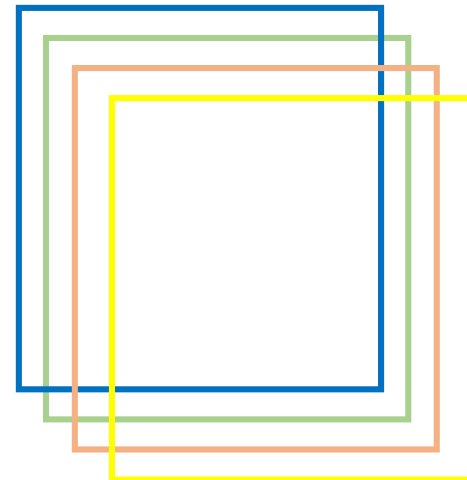
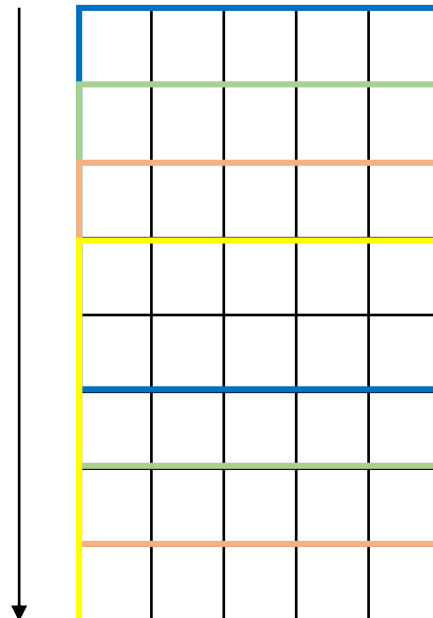
Anyway, you can try to replace any thing that has been done by RNN with self-attention.

Self-attention

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



Time



4x5x5

- 4: batch size
- 5: sequence length
- 5: feature dimension

NIPS papers
[https://papers.neurips.cc/paper/7181-attention...](https://papers.neurips.cc/paper/7181-attention-is-all-you-need) PDF

Attention is All you Need
by A Vaswani Cited by 140492 — We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely.
11 pages

对于time-series数据而言，attention一般作用在sequence维度（第二个维度）上

Self-attention

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



q : query (to match others)

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

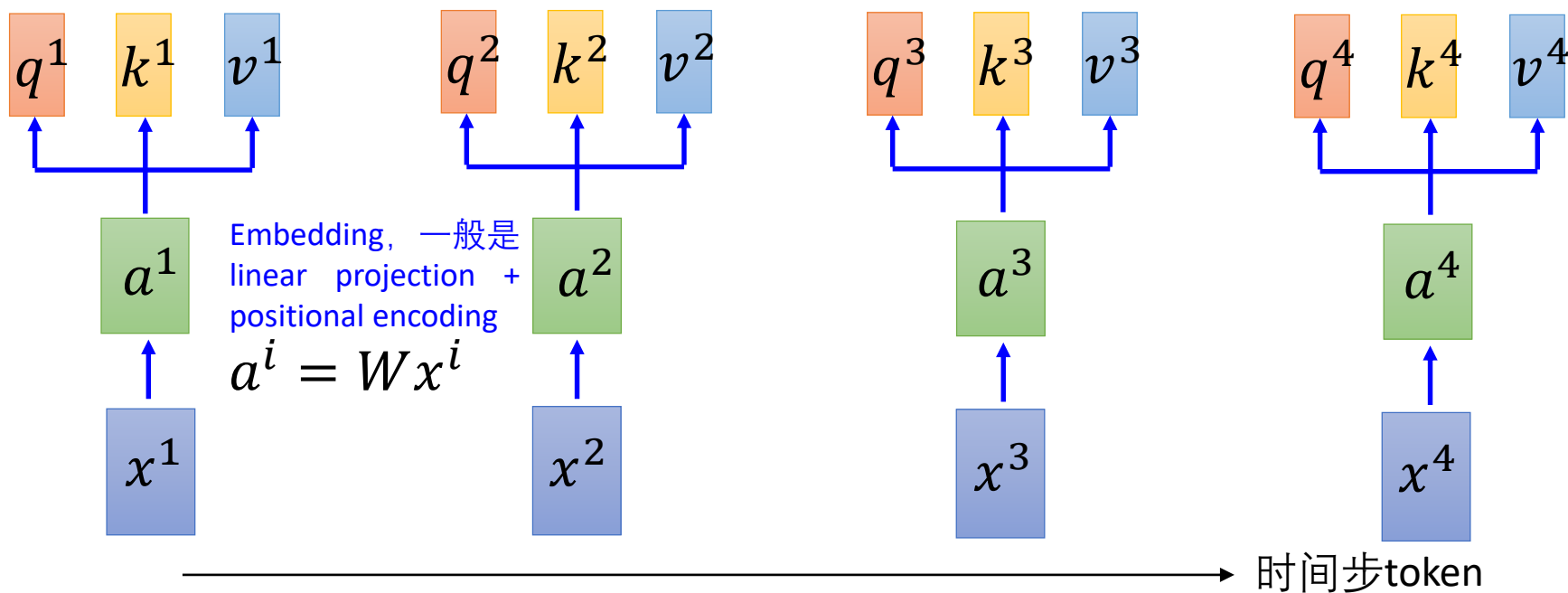
不同时间步embedding对应的 W^q, W^k, W^v 是同一个，换句话说qkv的projection parameter matrix在不同时间步上是shared

k : key (to be matched)

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

v : information to be extracted

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



q : query (to match others)

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

k : key (to be matched)

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

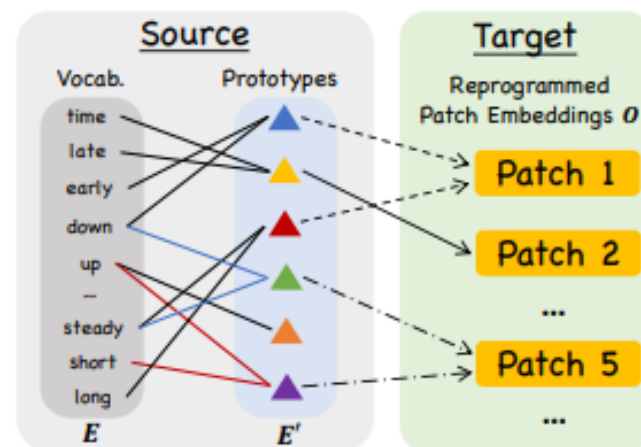
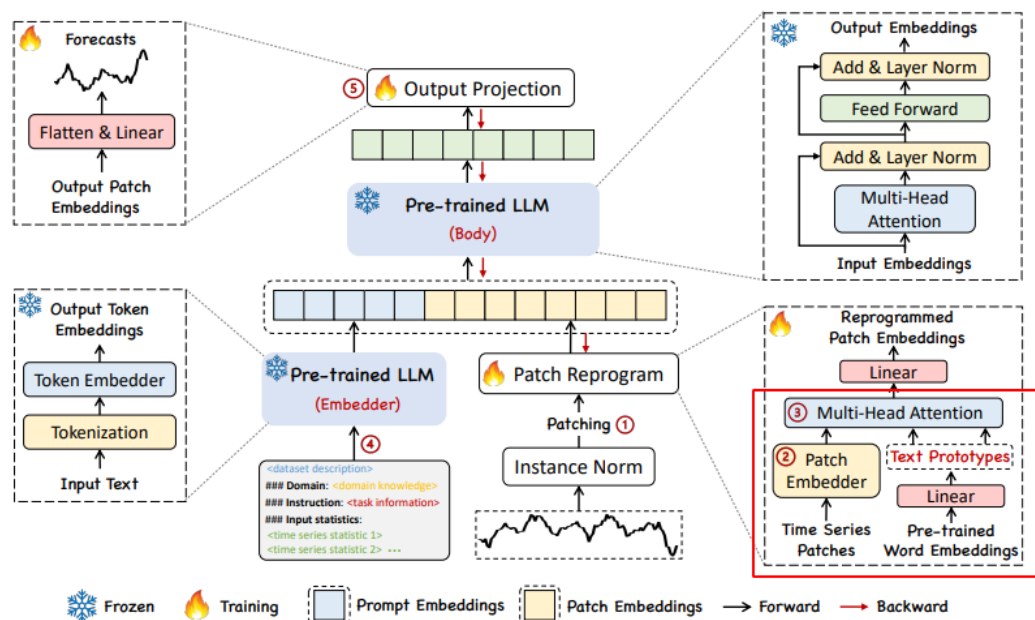
v : information to be extracted

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

Time-LLM: Time Series Forecasting by Reprogramming Large Language Models, ICLR 2024

Paper: <https://openreview.net/pdf?id=Unb5CVPtat>

Blog: <https://zhuanlan.zhihu.com/p/676256783>



(a) Patch Reprogramming

Pretrained LLM (embedder)是在文本语料上预训练的，具有文本modality上的知识，但现在我们需要的是time-series modality上做forecasting任务，因此需要以 **time-series data (embedding)** 作为 **query** 去 **match text embedding** (text embedding是被matched)，然后再从**text embedding**里面**extract**信息，相当于将text中和当前time-series embedding相似度高的部分（例如up, short)提取出来

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

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

刚才介绍的是cross-attention，涉及两个data modality之间的信息交互，self-attention就是自己和自己算attention，qkv来自于同一个data modality

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

Self-attention

qk之间做dot product (element-wise相乘相加), 会随着元素个数, 即维度d的增大而使得variance增大, 因此用根号d做scaled

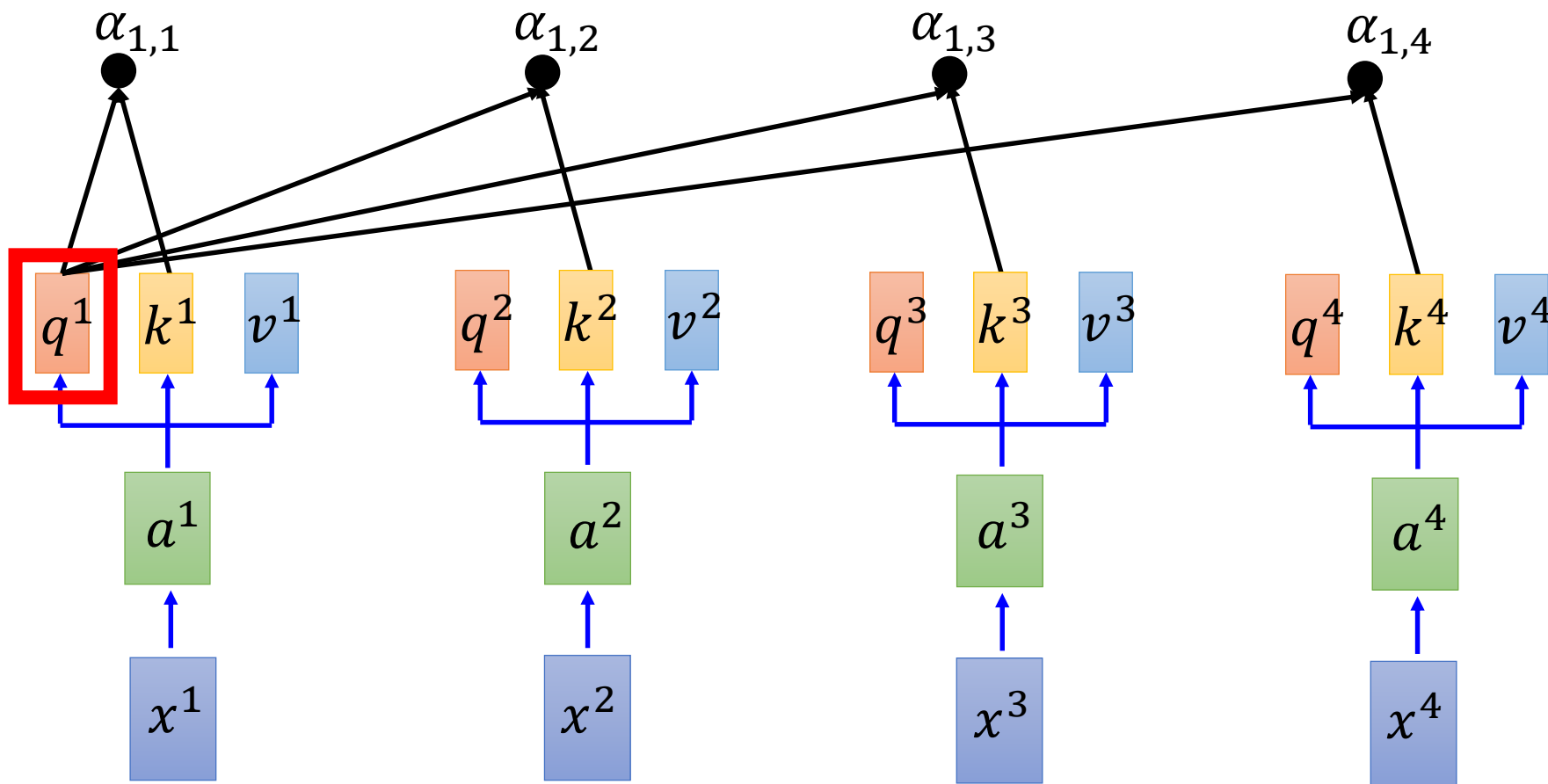
拿每個 query q 去對每個 key k 做 attention

d is the dim of q and k

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

dot product

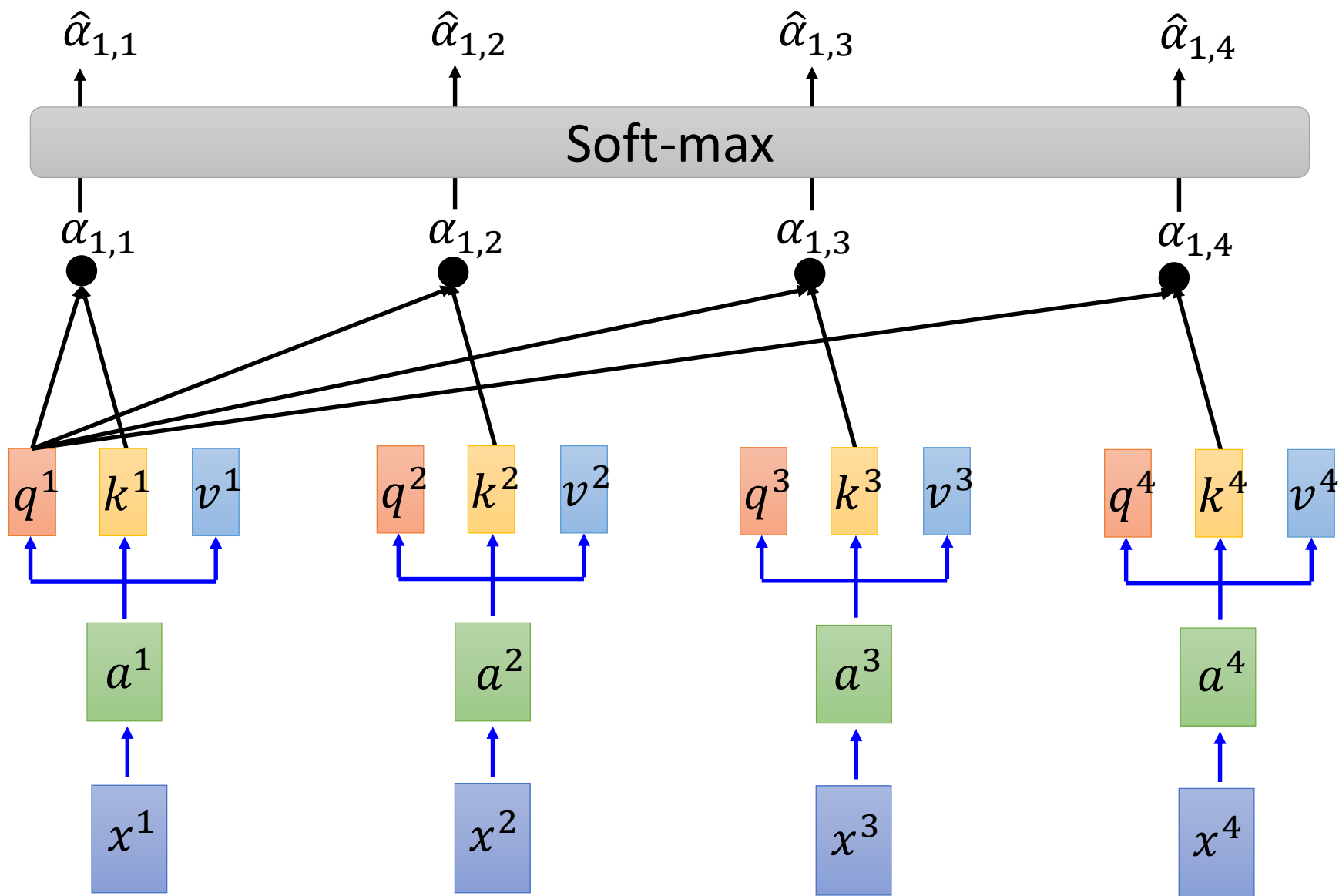
q^1 和 k^1 之间的
attention score



Self-attention

Softmax相当于把attention score归一化

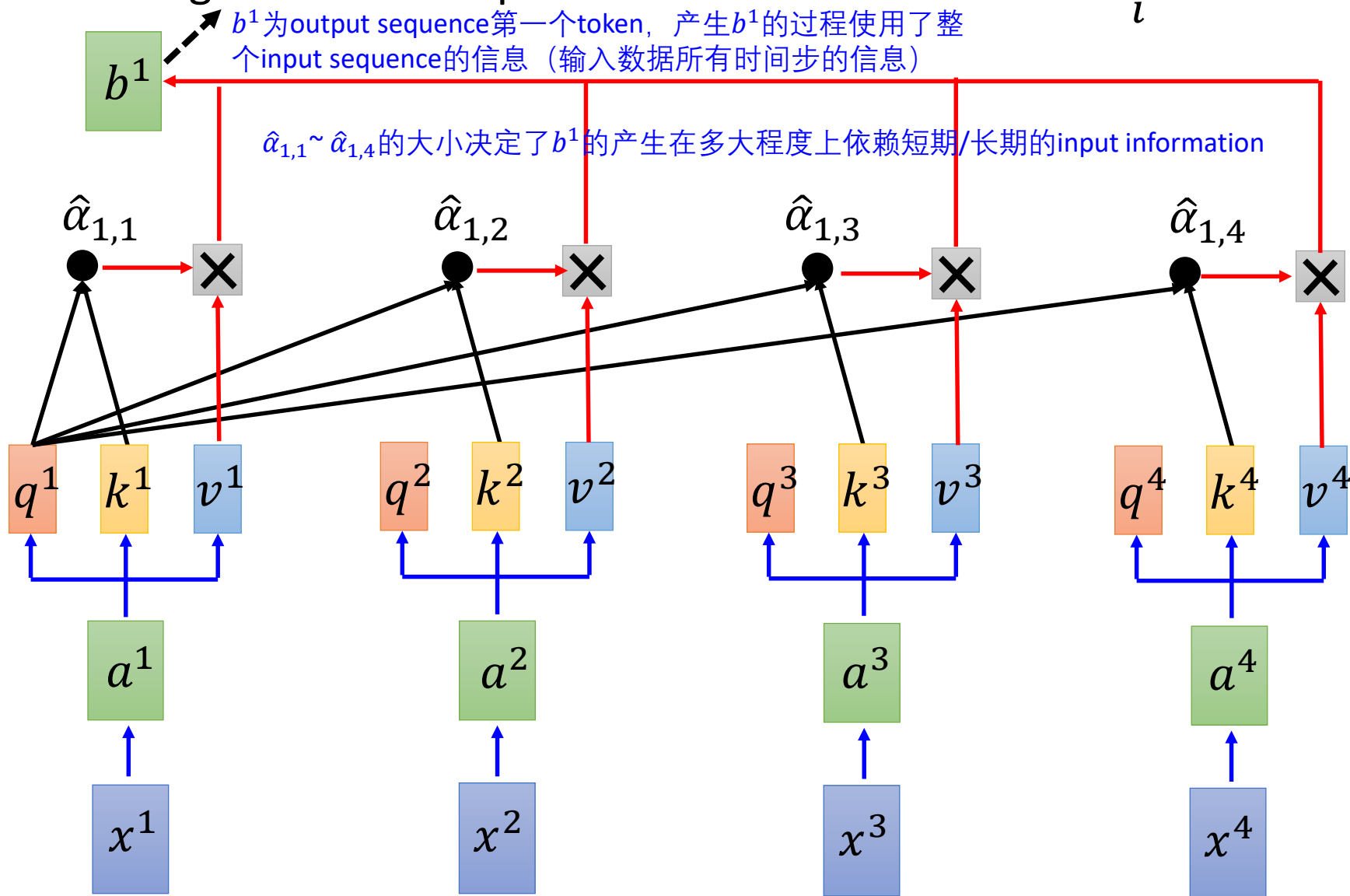
$$\hat{\alpha}_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



Self-attention

Considering the whole sequence

$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$

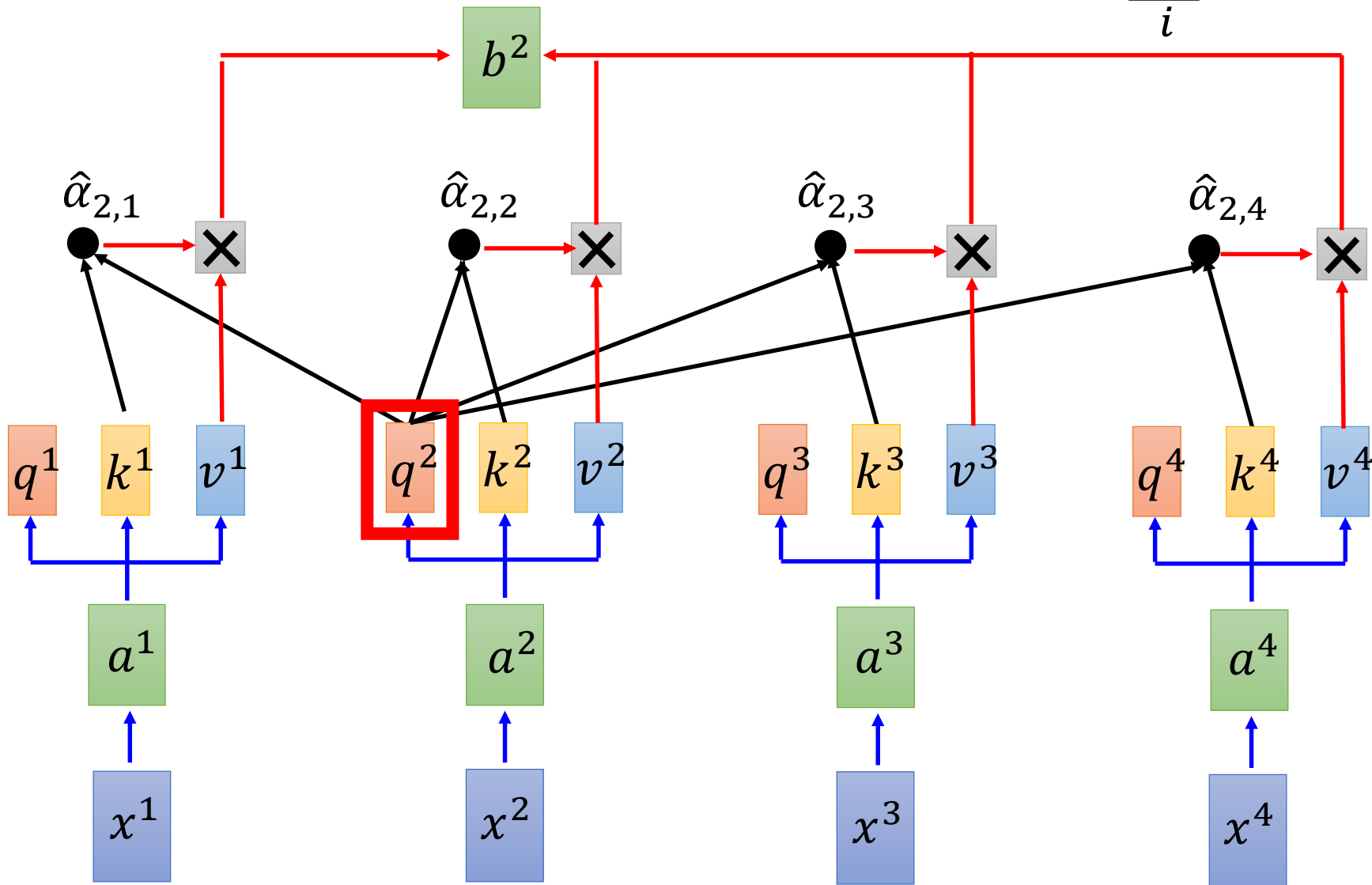


Self-attention

$b^1, b^2 \dots$ 的运算是parallel的, b^2 并不需要像RNN那样等 b^1 计算完成

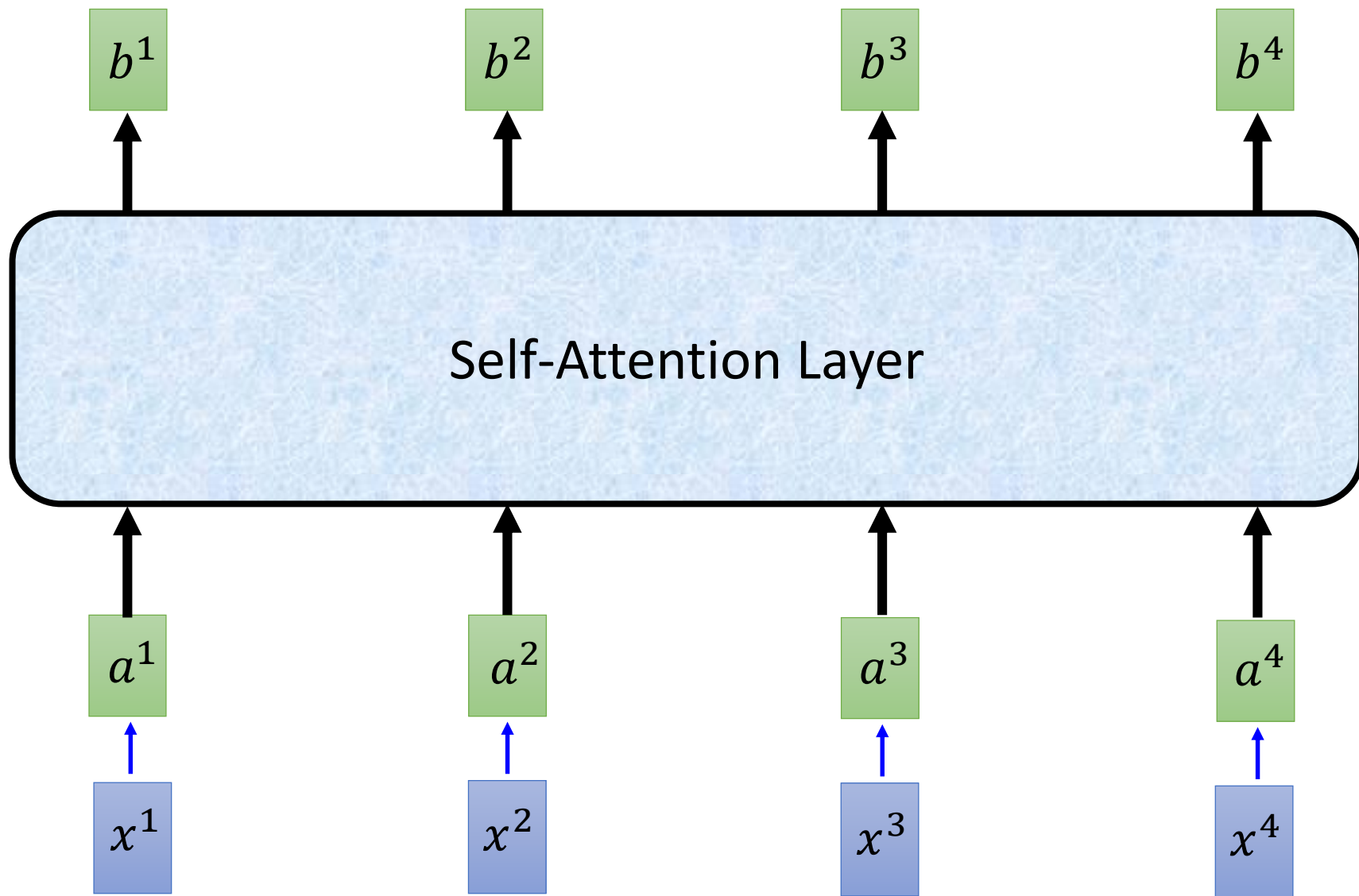
拿每個 query q 去對每個 key k 做 attention

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

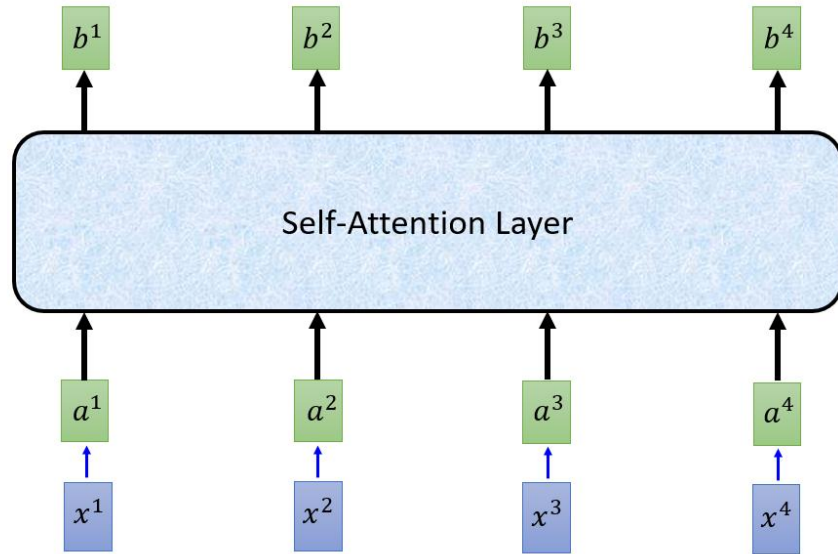


Self-attention

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



Self-attention



Attention 本质上在计算 token 之间的相似度 (与 cosine similarity 的计算类似), 然后用归一化的相似度去加权得到 output token

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

$$\text{cosine similarity} = S_C(A, B) := \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \cdot \sqrt{\sum_{i=1}^n B_i^2}},$$

Self-attention

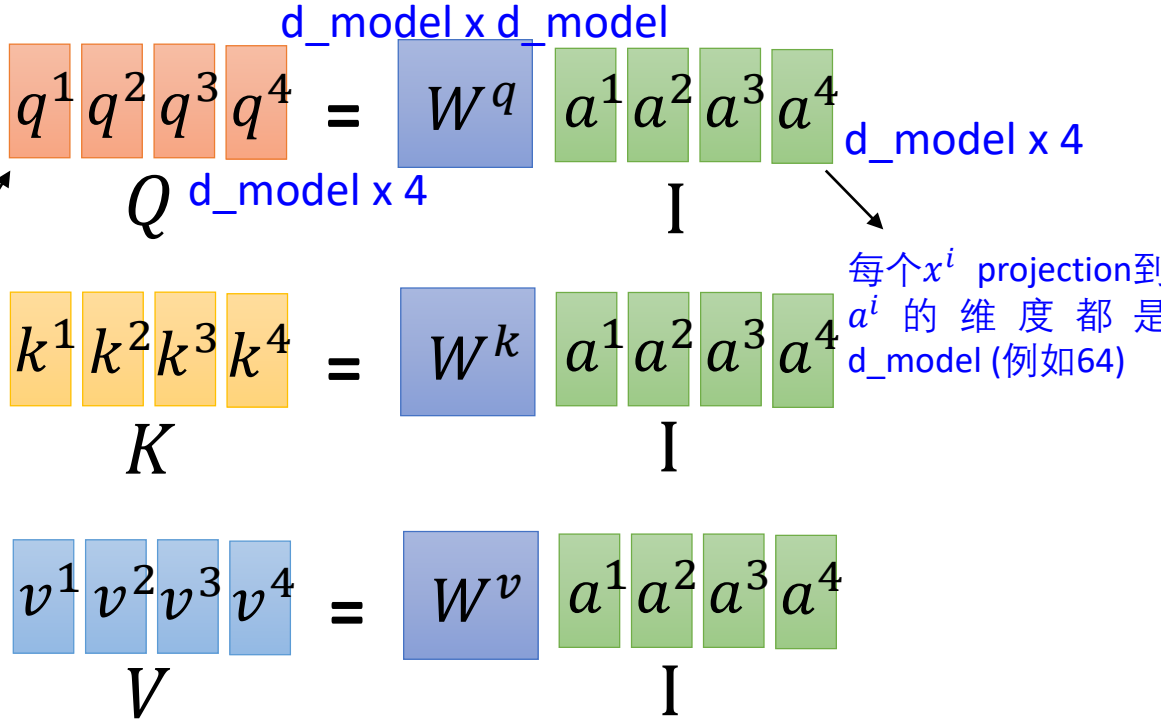
用矩阵的思维来理解
self-attention

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

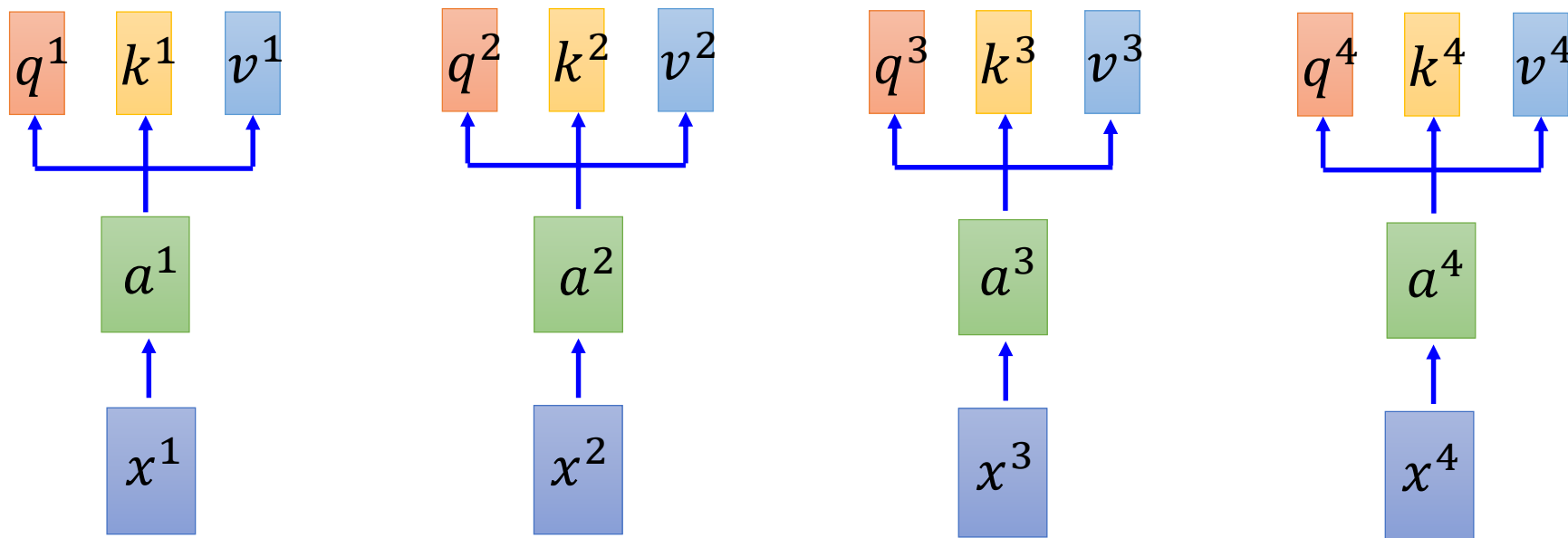
$$k^i = W^k a^i$$

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

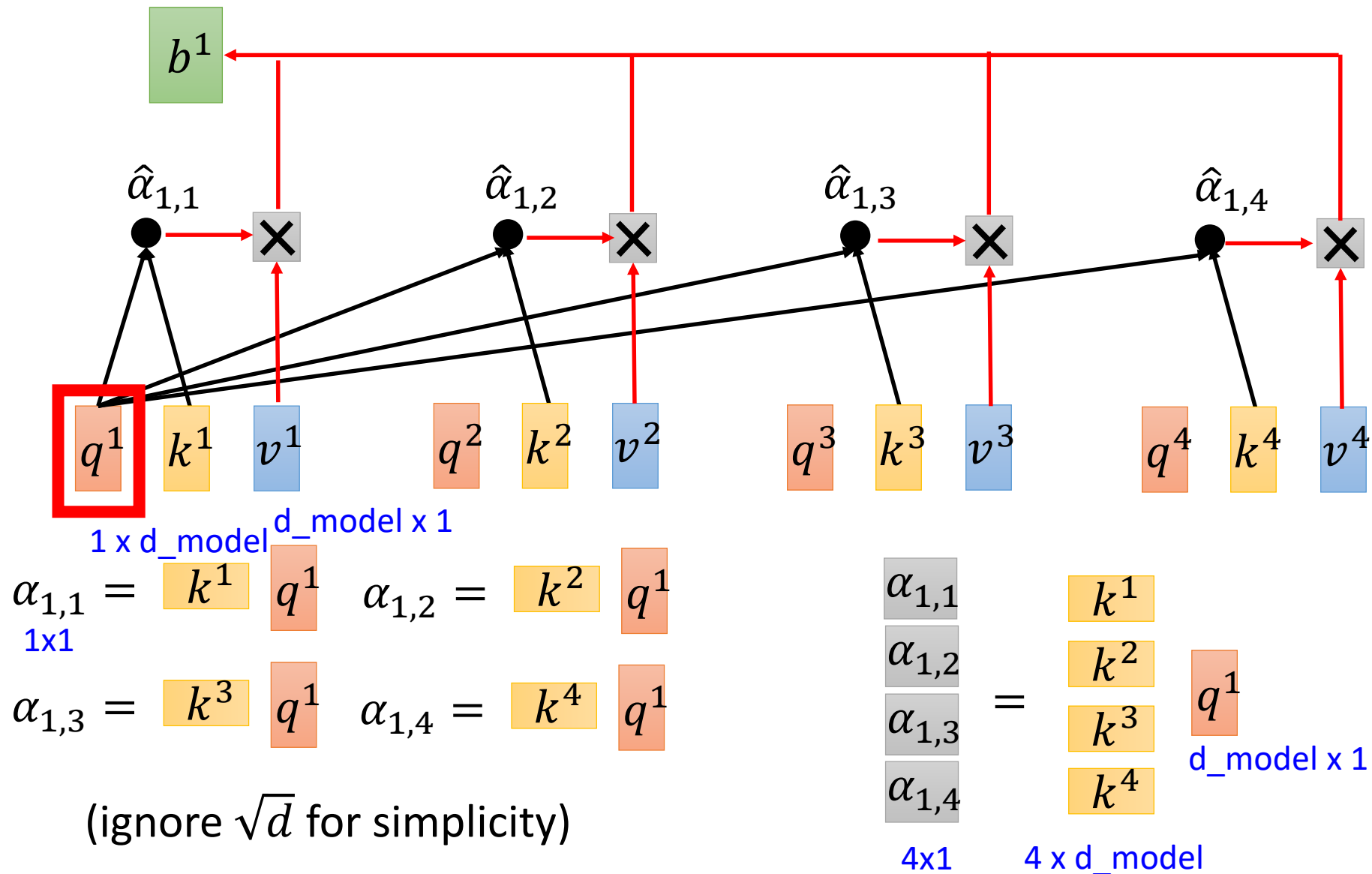
拼起来



每个都是 d_{model} 维度的一维向量



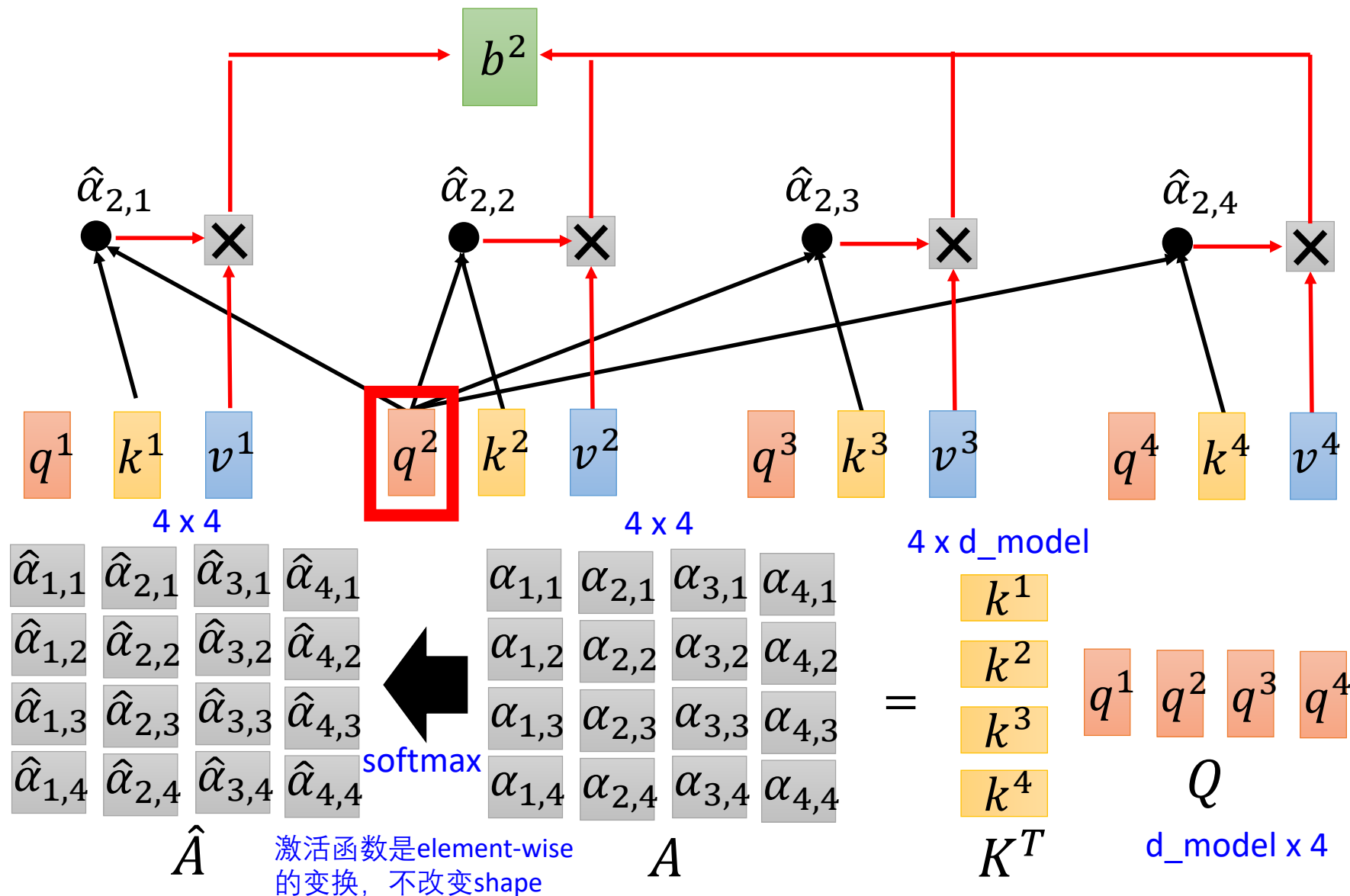
Self-attention



Self-attention

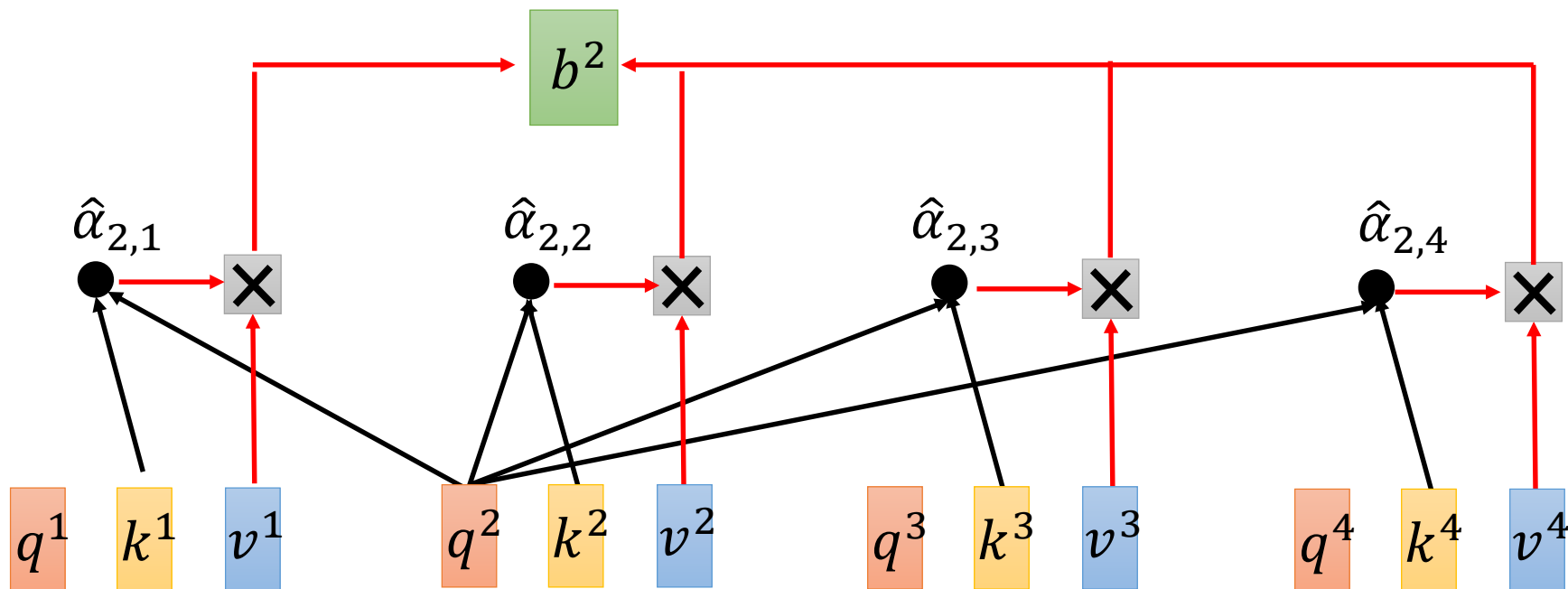
重复之前的操作

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



Self-attention

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



$d_{\text{model}} \times 4$

$b^1 b^2 b^3 b^4$

O

$d_{\text{model}} \times 4$

$v^1 v^2 v^3 v^4$

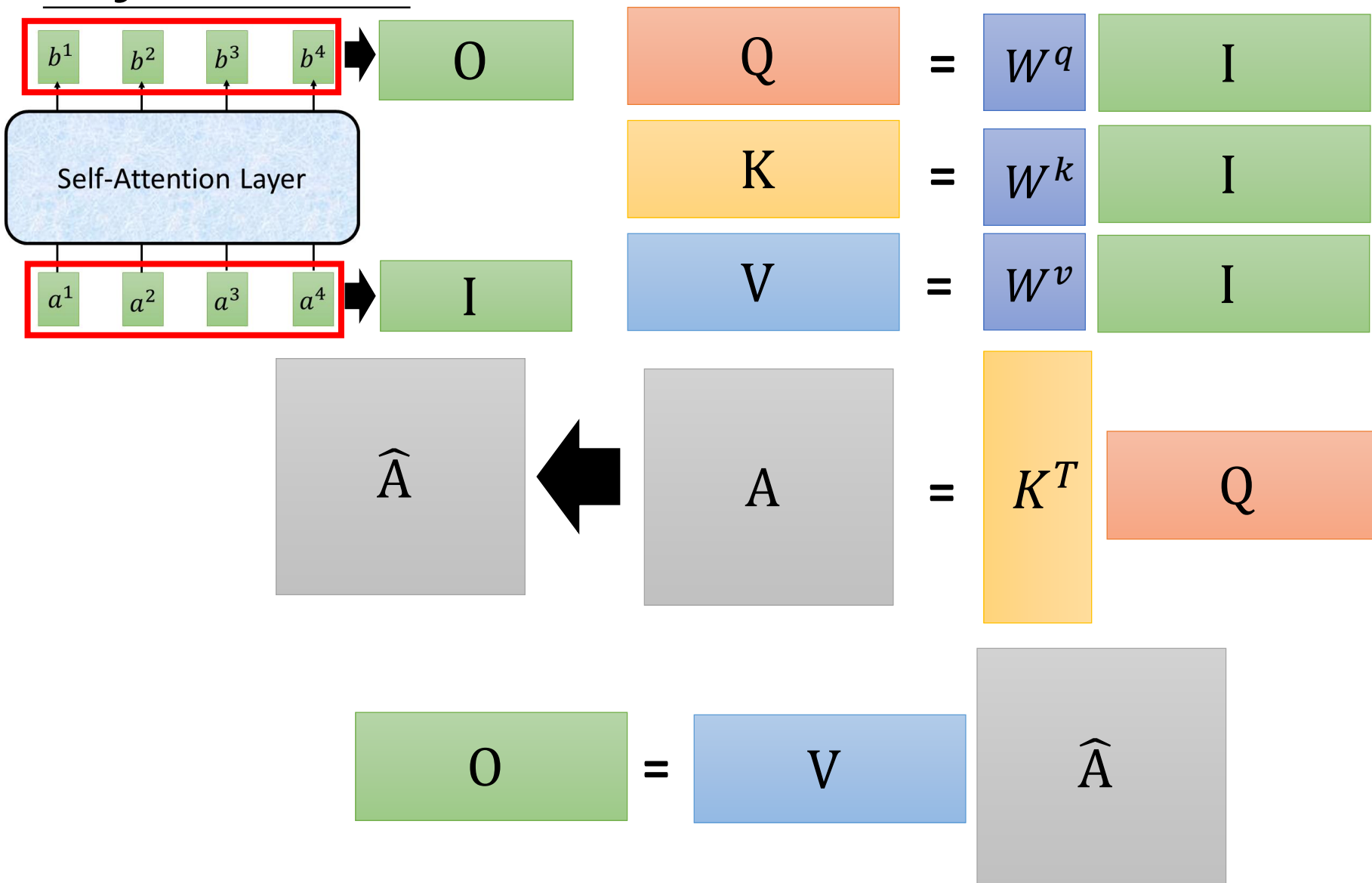
V

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

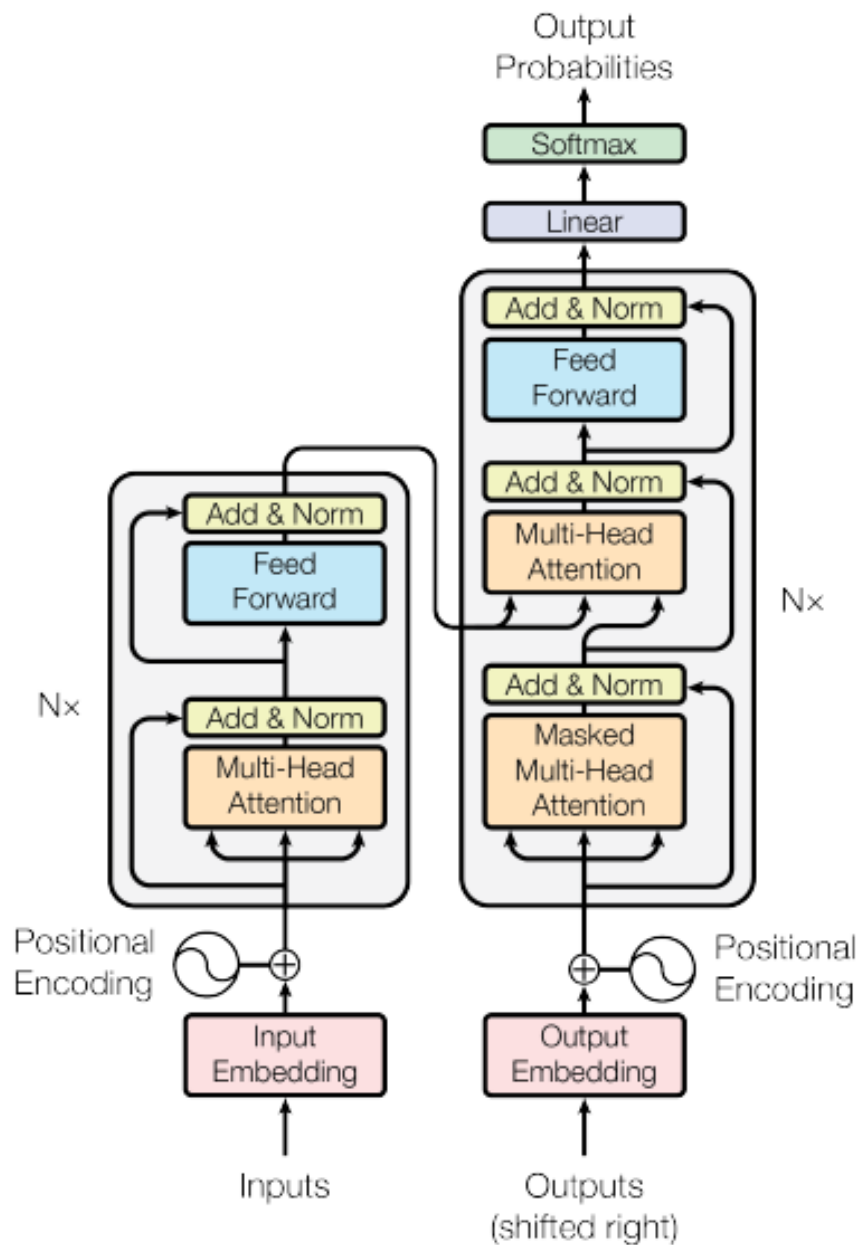
\hat{A}

4×4

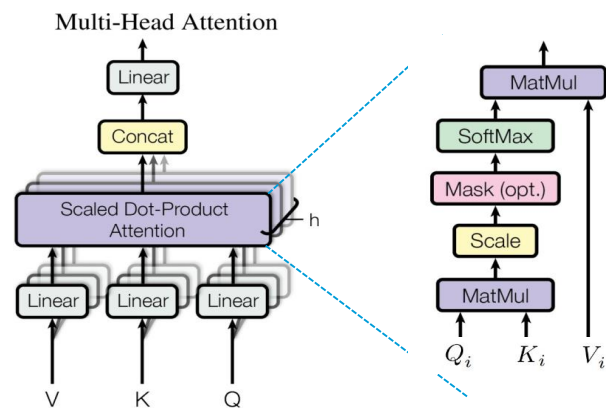
Self-attention



反正就是一堆矩陣乘法，用 GPU 可以加速



- Multi-head attention
- Masked multi-head attention
- Attention中的permutation invariant问题
- Transformer整体结构解读
- Transformer/Attention变体

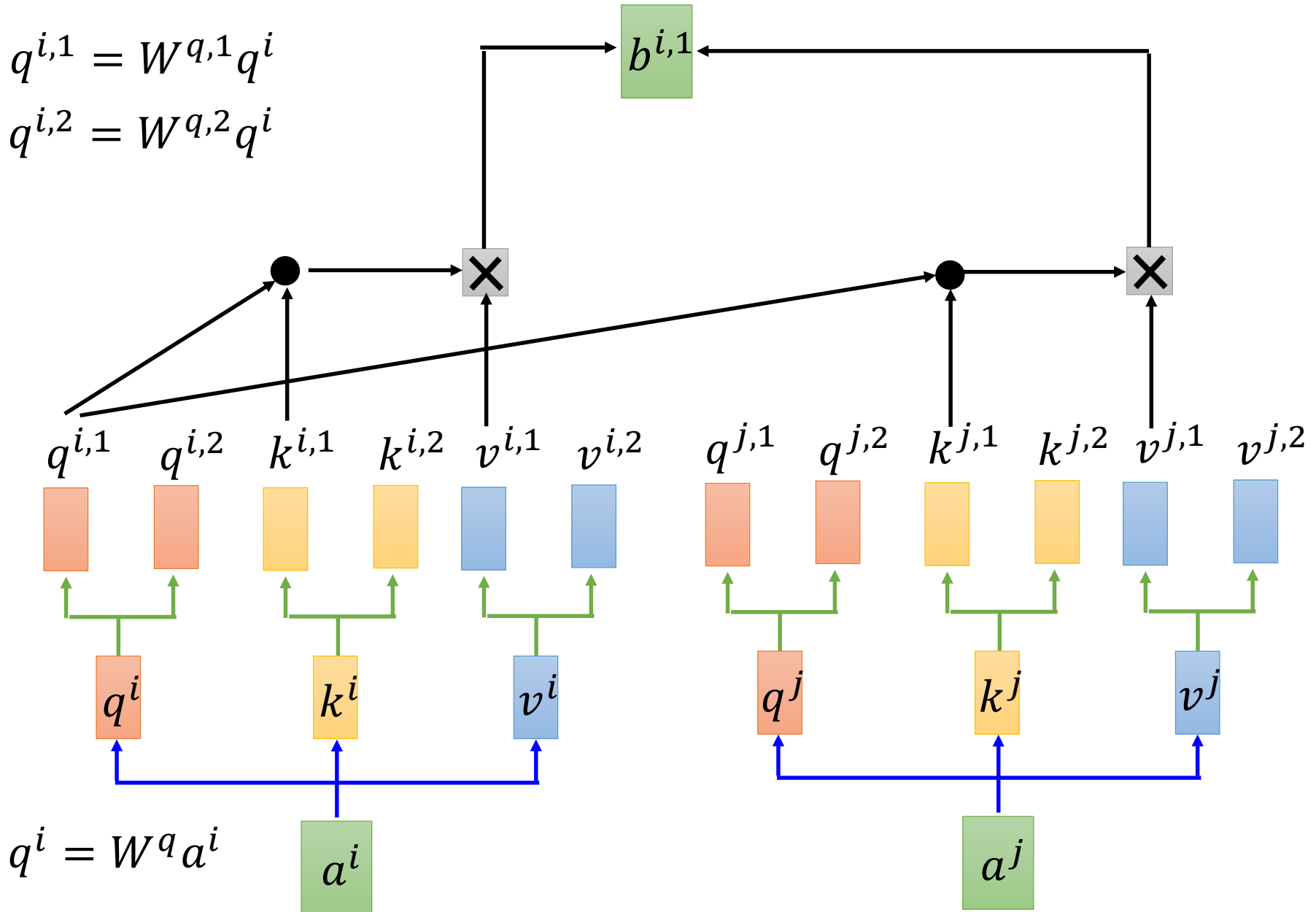


Multi-head Self-attention

(2 heads as example)

$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$

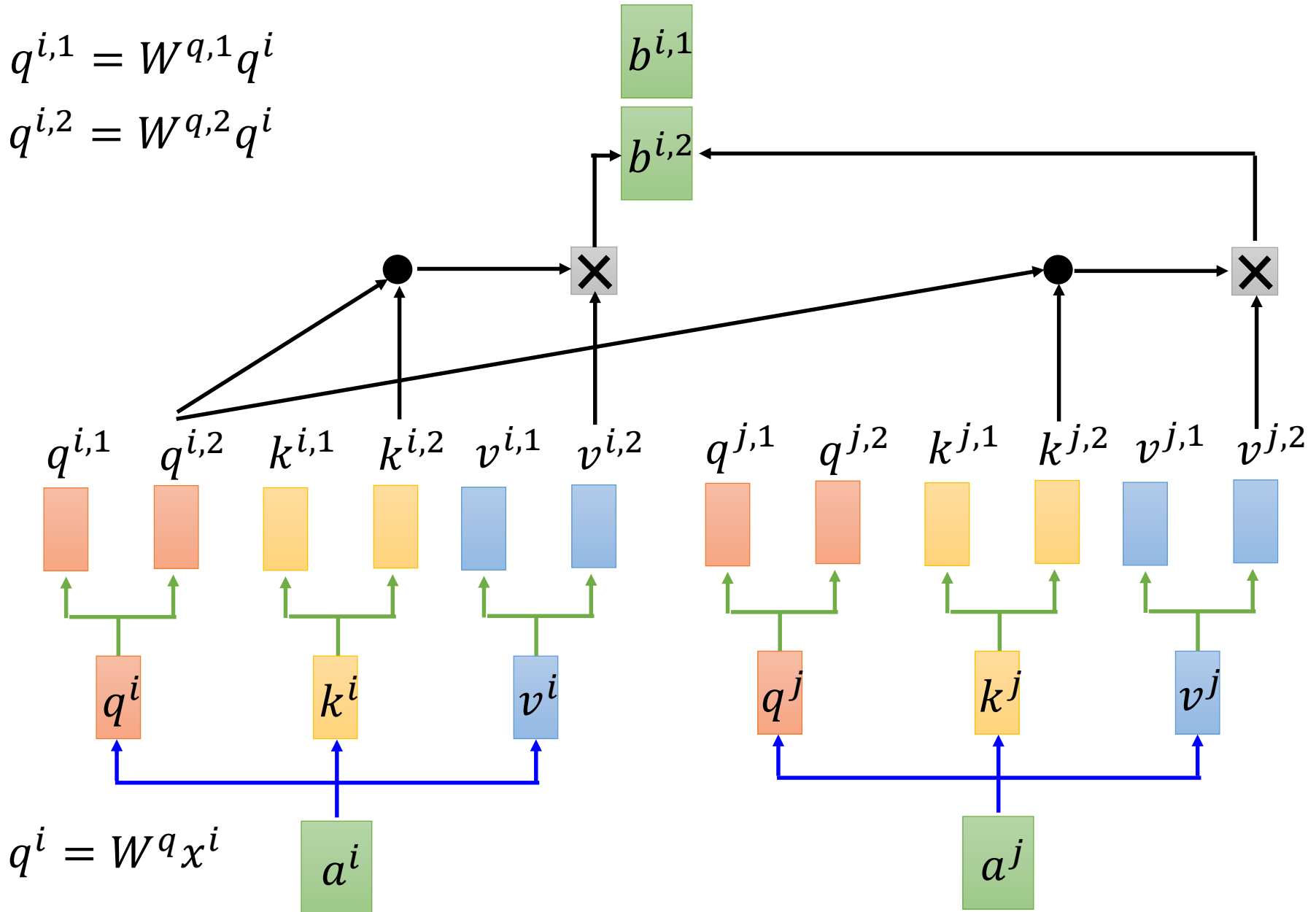


Multi-head Self-attention

(2 heads as example)

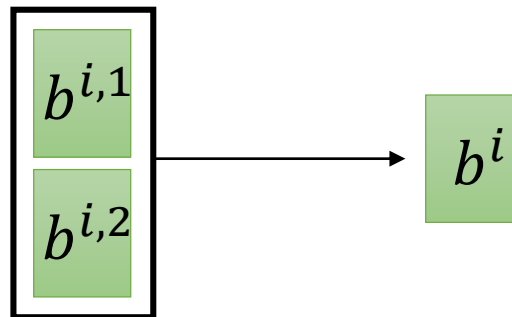
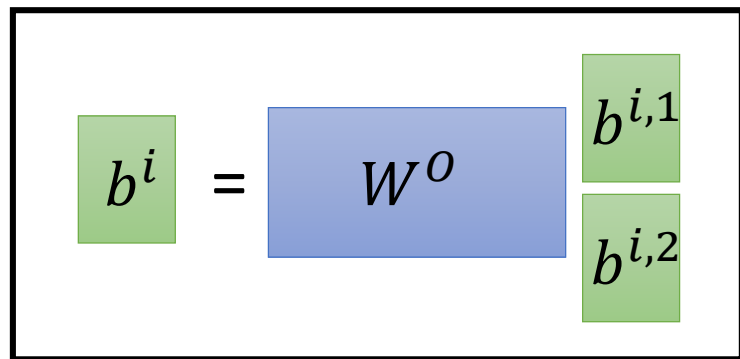
$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$

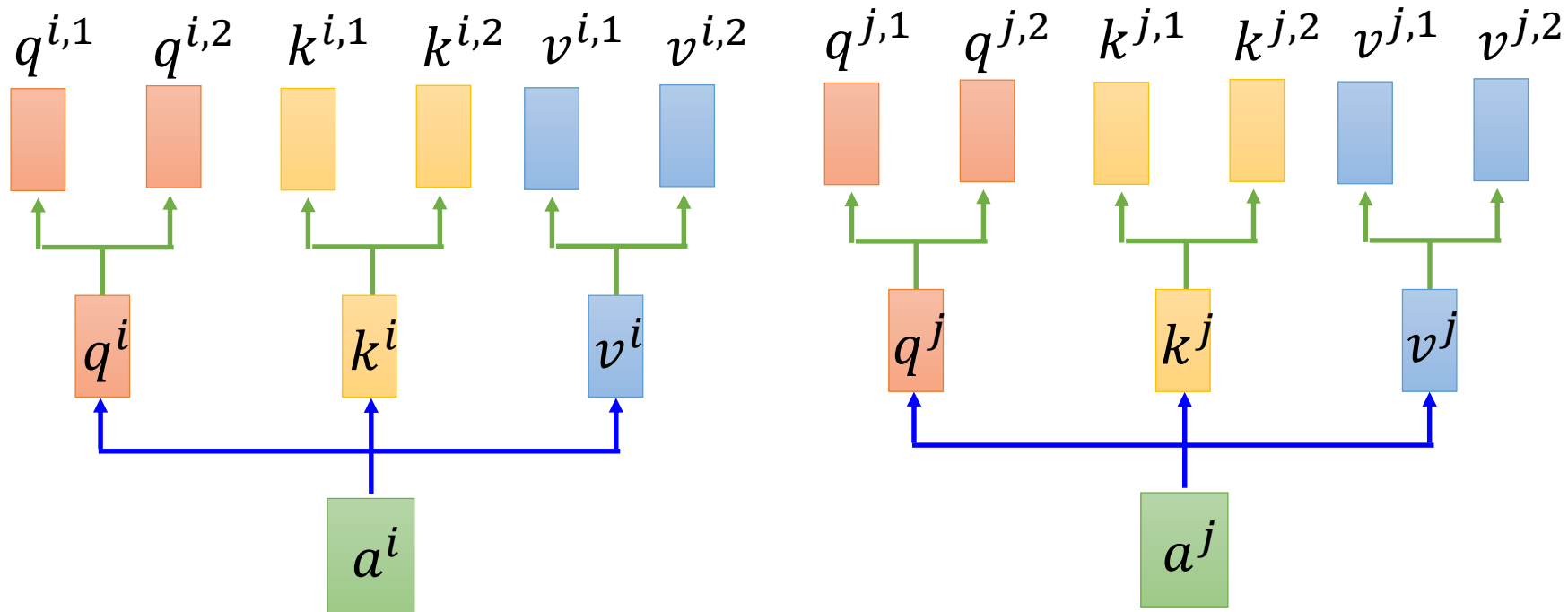


Multi-head Self-attention

(2 heads as example)



Again, Parallel



Multi-head Self-attention

(2 heads as example)

```
# multi-head attention layer
class AttentionLayer(nn.Module):
    def __init__(self, configs,
                  attention: nn.Module): 可以嵌套多种attention
        super(AttentionLayer, self).__init__()
        d_model = configs['model']['d_model']
        n_heads = configs['model']['n_heads']
        seq_len = configs['model']['seq_len']
        self.configs = configs

        d_queries = d_model // n_heads
        d_keys = d_model // n_heads
        d_values = d_model // n_heads

        self.attention = attention
        self.query_projection = nn.Linear(d_model, d_queries * n_heads)
        self.key_projection = nn.Linear(d_model, d_keys * n_heads)
        self.value_projection = nn.Linear(d_model, d_values * n_heads)
        self.out_projection = nn.Linear(d_values * n_heads, d_model)
        self.n_heads = n_heads

    def forward(self, q, k, v, attn_mask=None):
        B, L, _ = q.size()
        _, S, _ = k.size()
        H = self.n_heads

        q = self.query_projection(q).view(B, H, L, -1)
        k = self.key_projection(k).view(B, H, S, -1)
        v = self.value_projection(v).view(B, H, S, -1)

        out, attn = self.attention(q, k, v)
        out = out.view(B, L, -1) # concat multi-head attention calculated results

        return out, attn
```

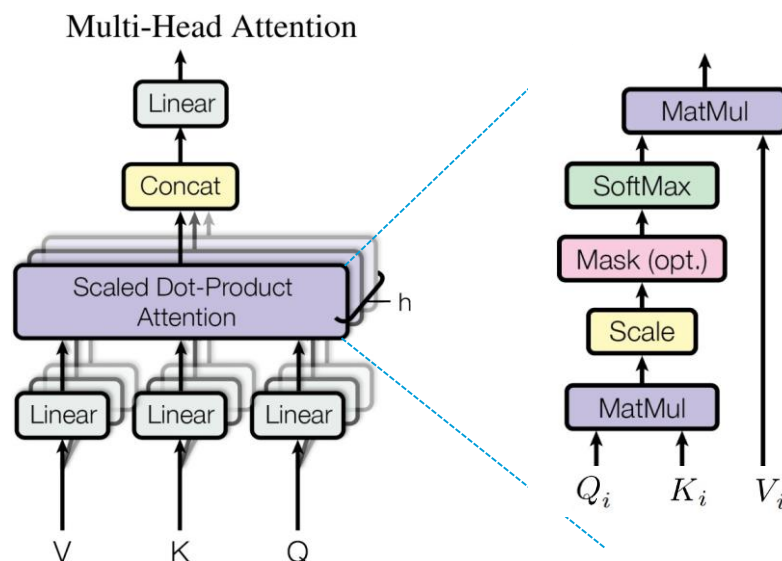
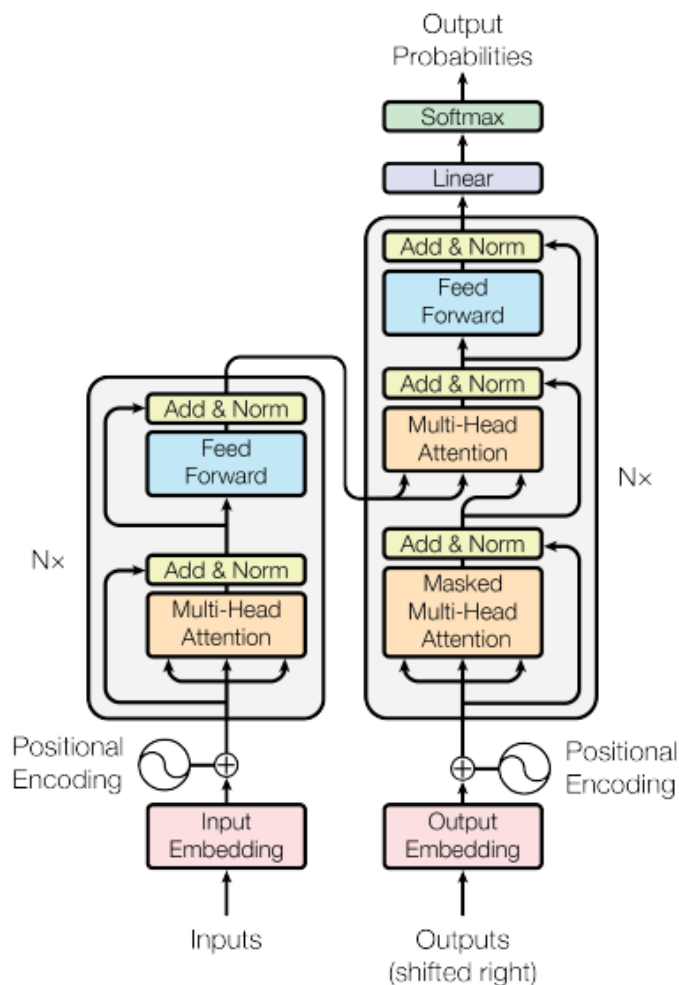
input_size 维的输入先通过 input embedding (linear projection) + positional encoding 到 d_model 维，然后通过 n_heads 等分

This multi-dimensionality allows the attention mechanism to jointly attend to different information from different representation at different positions

从矩阵角度来看multi-head只是reshape了一下，把原来的特征空间等分维多个特征子空间，然后再不同特征子空间的角度，关注序列中的不同信息。

即在每个特征子空间上做self-attention，每个特征子空间关注到的序列信息不同

Masked Multi-head Self-attention



NLP: 确保模型在生成每个词时，只能关注到当前和之前的词，而不能关注到未来的词；

Time-series: 确保模型只能关注到当前时间步及之前时间步的信息，而不用未来信息作弊

Masked Multi-head Self-attention

$$\begin{bmatrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{bmatrix} = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{bmatrix} q^1$$

$$\hat{A} \xrightarrow{\text{softmax}} A = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{bmatrix} K^T \quad q^1 \quad q^2 \quad q^3 \quad q^4 \quad Q$$

应该让主对角线下方的
attention score 变成0

Masked Multi-head Self-attention

```
class MaskGenerator():
    def __init__(self, batch_size, seq_length, device):
        # broadcast on the attention heads dimension
        mask_shape = [batch_size, 1, seq_length, seq_length]

        with torch.no_grad():
            self.mask = torch.triu(torch.ones(mask_shape, dtype=torch.bool), diagonal=1).to(device)

    @property
    def mask(self):
        return self.mask

# self-attention
class FullAttention(nn.Module):
    def __init__(self, configs):
        super(FullAttention, self).__init__()
        self.mask = configs['model']['mask']
        self.dropout = nn.Dropout(p=configs['model']['dropout'])

    def forward(self, q, k, v, attn_mask=None):
        # qkv的输入维度分别是batch size, attention heads, sequence length, embedding dim
        # b, h, d分别代表batch size, attention heads个数, 以及embedding维度
        # l和s分别是query和key的sequence length, 只是因为输出矩阵bhss是不合法的, 需要取两个变量
        # QK^T, lxdxdxs=lx s, 再考虑batch size和attention heads为bxbx1xs的张量
        attn_scores = torch.einsum("bhld,bhsd->bhls", q, k)

        if self.mask:
            attn_mask = MaskGenerator(batch_size=q.size(0), seq_length=q.size(2), device=q.device)
            attn_scores.masked_fill_(attn_mask.mask, -np.inf)

        # 只考虑最后两个维度的话1xs, 即softmax计算每个query对于所有key的attention scores
        attn_scores = self.dropout(torch.softmax(1. / sqrt(q.size(-1)) * attn_scores, dim=-1))
        weighted_values = torch.einsum("bhls,bhsd->bhld", attn_scores, v)

        # .contiguous()确保张量的数据在内存中是连续的,以提高后续操作的效率
        # The .contiguous() ensures the memory of the tensor is stored contiguously
        # which helps avoid potential issues during processing.
        return attn_scores, weighted_values.contiguous()
```

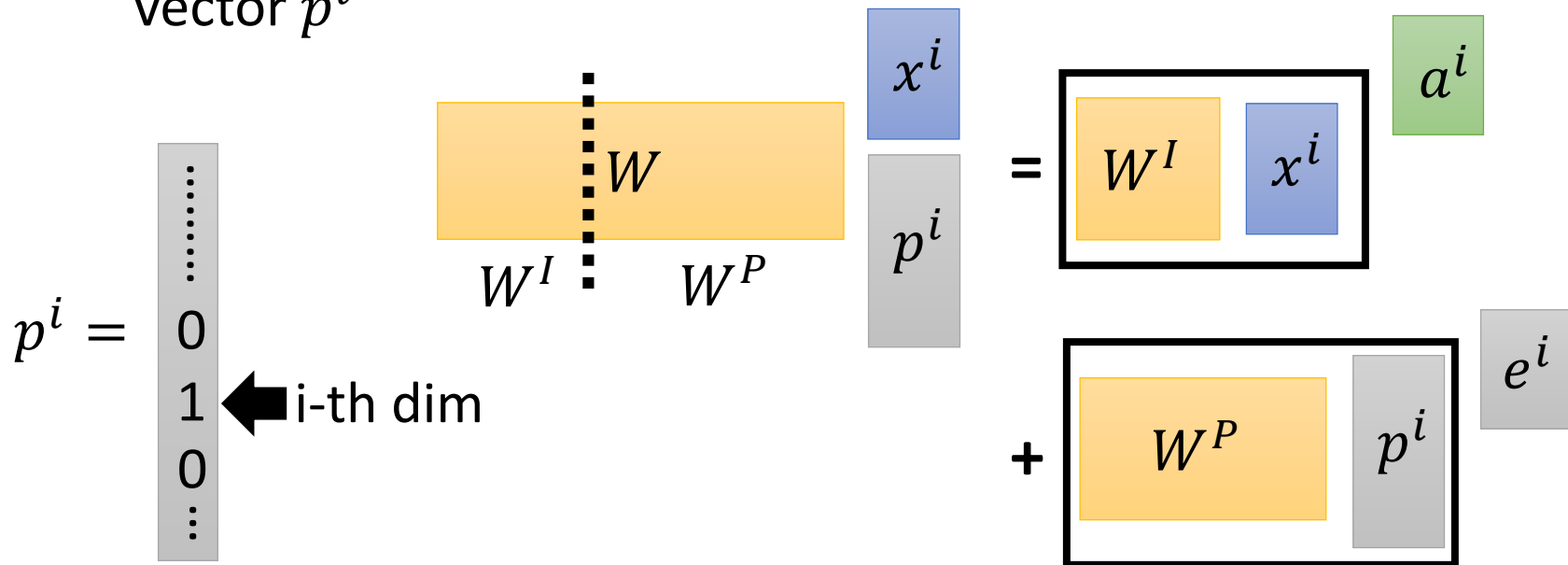
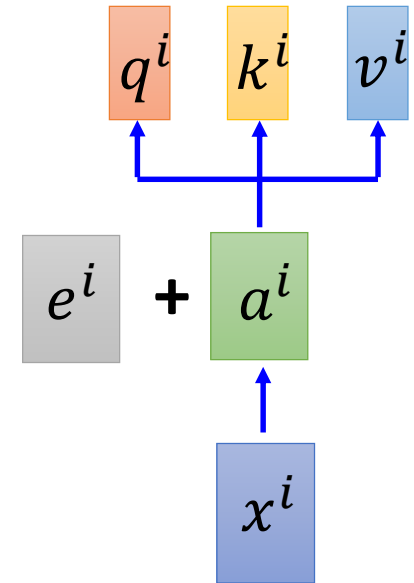
Mask那些位于主对角线
下方的元素

$$\hat{\alpha}_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$

让那些元素值变成-inf,
这样过 softmax 之后
probability就变成了0

Positional Encoding

- **No position information** in self-attention.
- Original paper: each position has a unique positional vector e^i (not learned from data)
- In other words: each x^i appends a one-hot vector p^i



Attention is all you need: <https://arxiv.org/pdf/1706.03762>

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

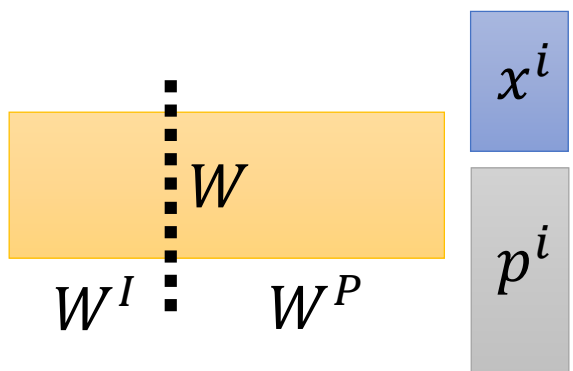
```
class PositionalEmbedding(nn.Module):
    def __init__(self, d_model, max_len=5000):
        super(PositionalEmbedding, self).__init__()
        # Compute the positional encodings once in log space.
        pe = torch.zeros(max_len, d_model).float()
        pe.requires_grad = False

        position = torch.arange(0, max_len).float().unsqueeze(1)
        div_term = (torch.arange(0, d_model, 2).float()
                    * -(math.log(10000.0) / d_model)).exp()

        pe[:, 0::2] = torch.sin(position * div_term) 2i
        pe[:, 1::2] = torch.cos(position * div_term) 2i+1

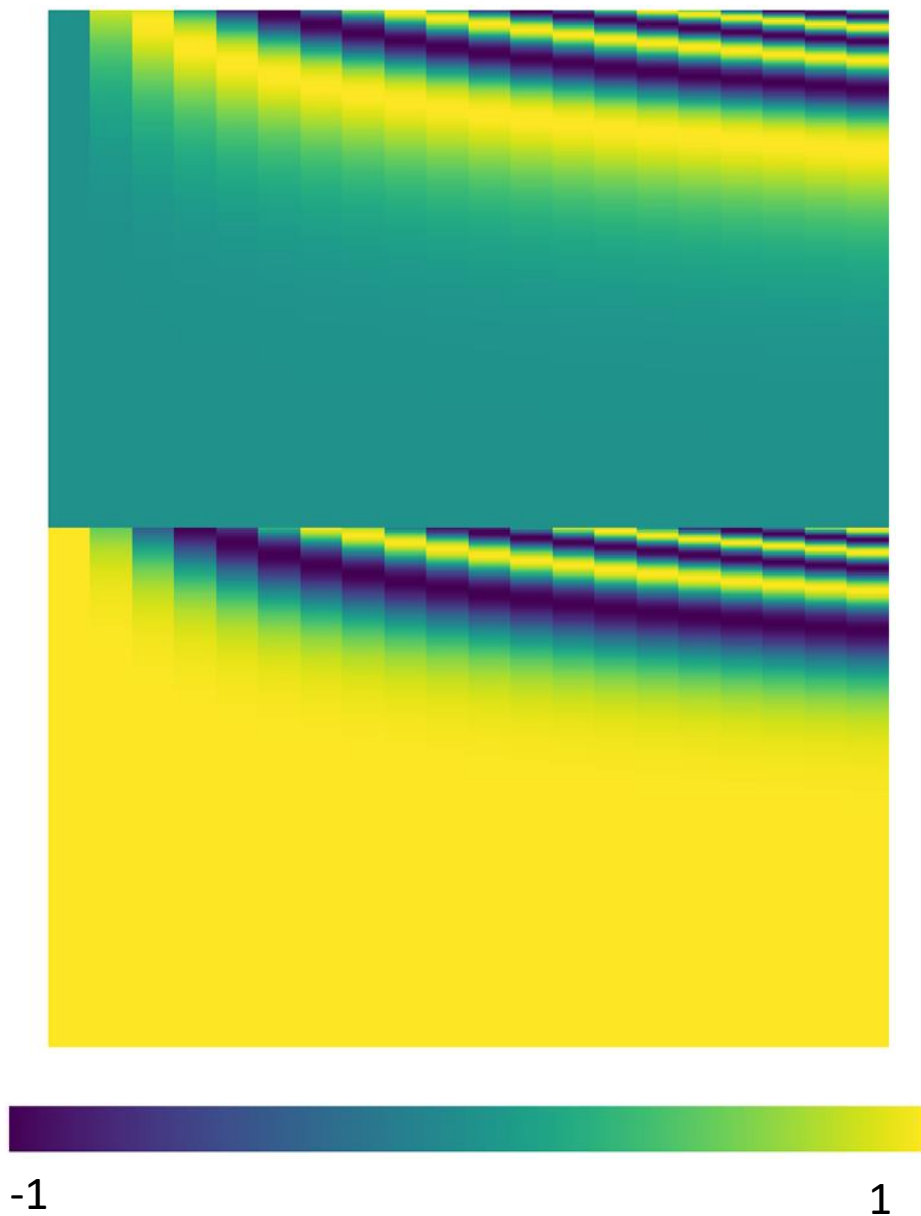
        pe = pe.unsqueeze(0)
        self.register_buffer('pe', pe)

    def forward(self, x):
        return self.pe[:, :x.size(1)]
```



$$= \begin{bmatrix} W^I & x^i \end{bmatrix} a^i + \begin{bmatrix} W^P & p^i \end{bmatrix} e^i$$

The equation shows the matrix W being multiplied by a vector a^i (green box) and added to the product of W^P and p^i (gray box). The vectors x^i and p^i are shown as inputs to the matrix multiplication.

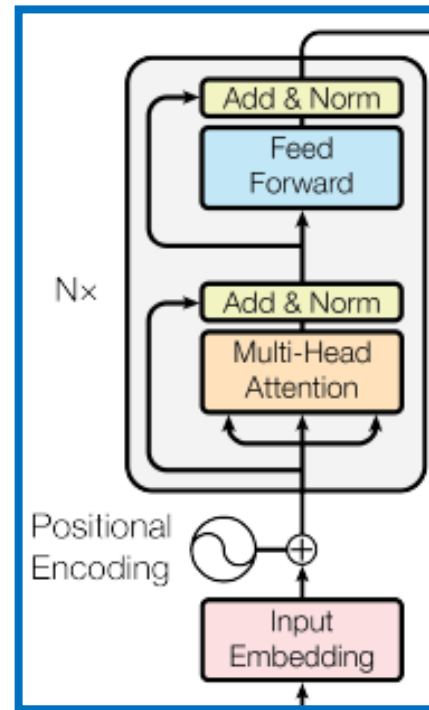


source of image: <http://jalammar.github.io/illustrated-transformer/>

Transformer

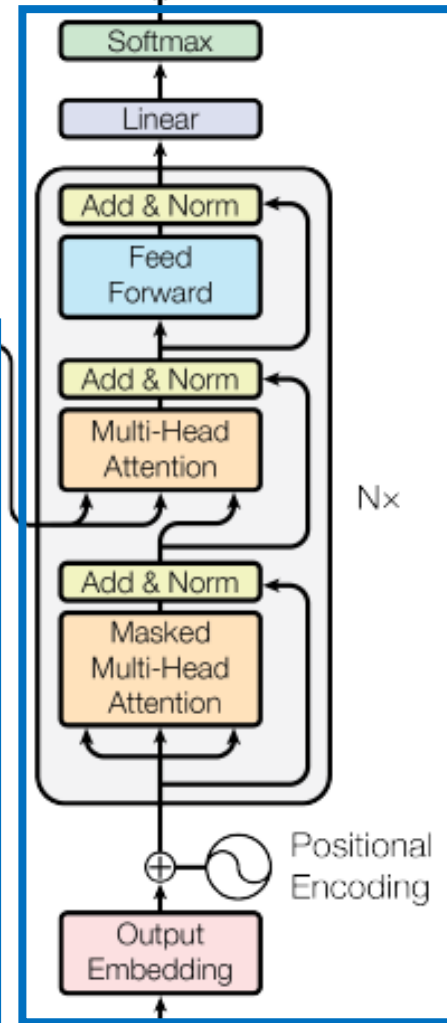
Using Chinese to English translation as example

Encoder



機器學習

machine learning



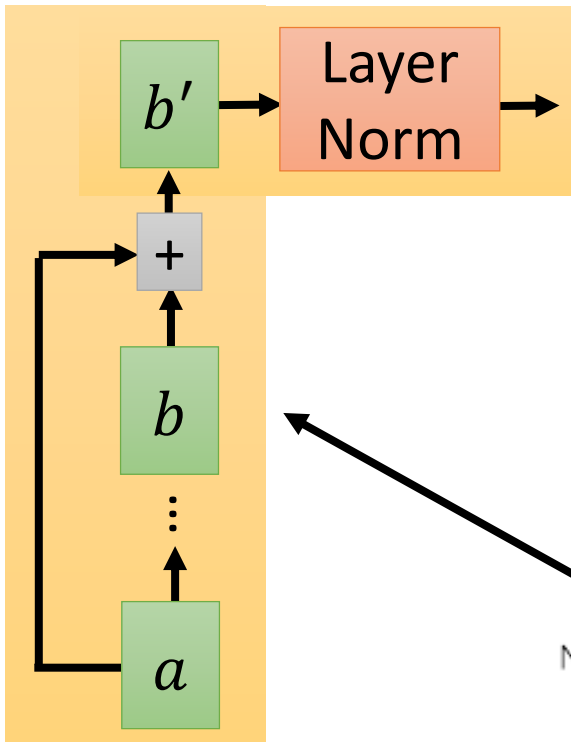
Decoder

Outputs
(shifted right)

<BOS>

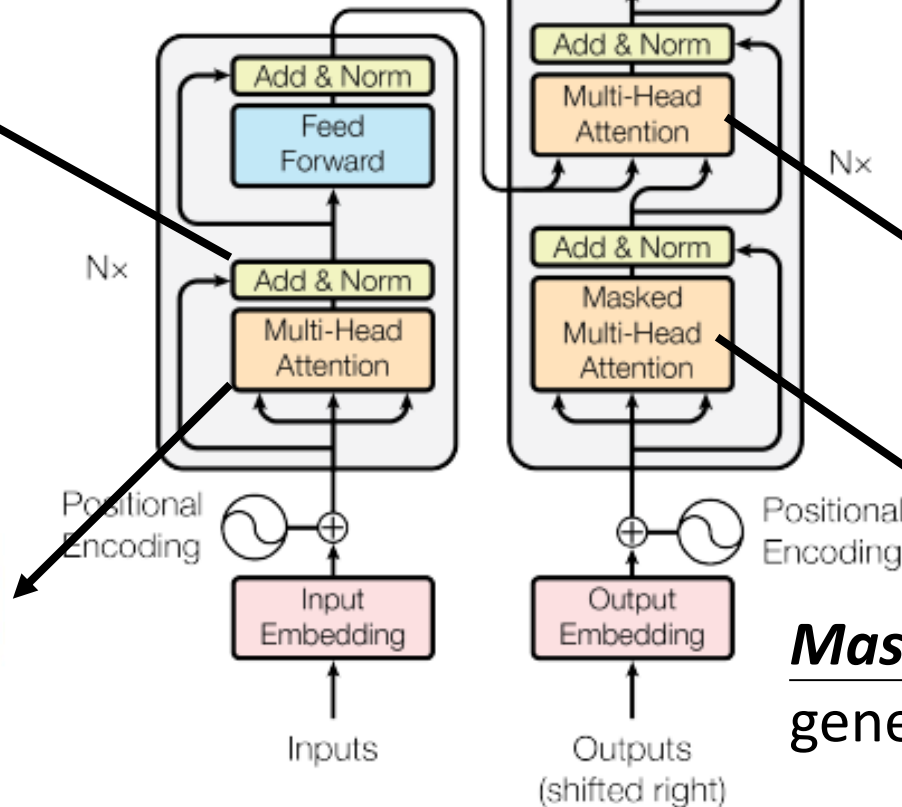
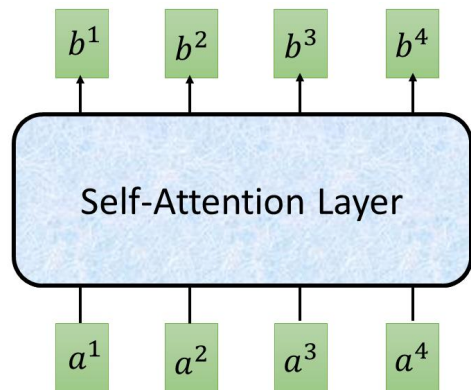
machine

Transformer

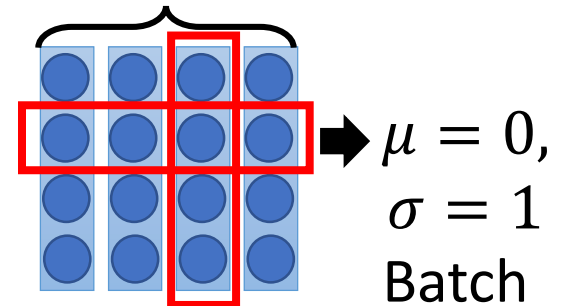


Layer Norm:
<https://arxiv.org/abs/1607.06450>

Batch Norm:
<https://www.youtube.com/watch?v=BZh1ltr5Rkg>



Batch Size



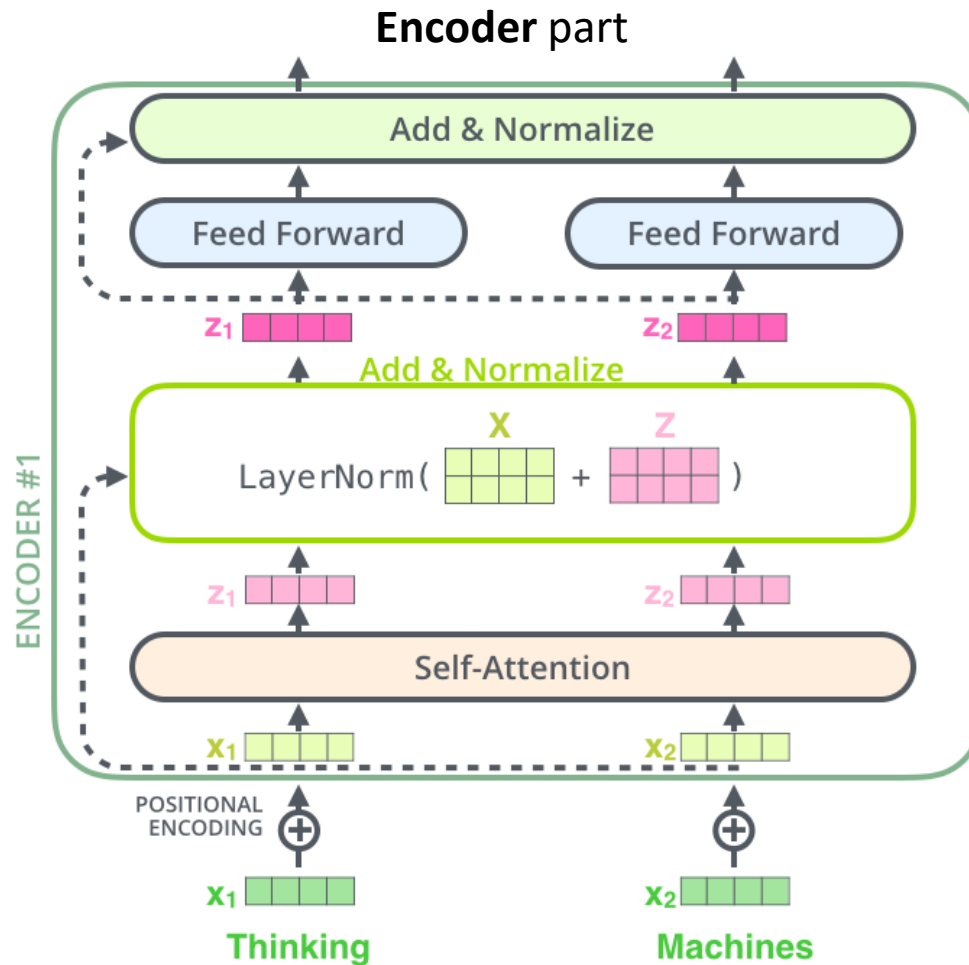
$\mu = 0, \sigma = 1$
Layer

attend on the
input sequence

Masked: attend on the
generated sequence

Transformer

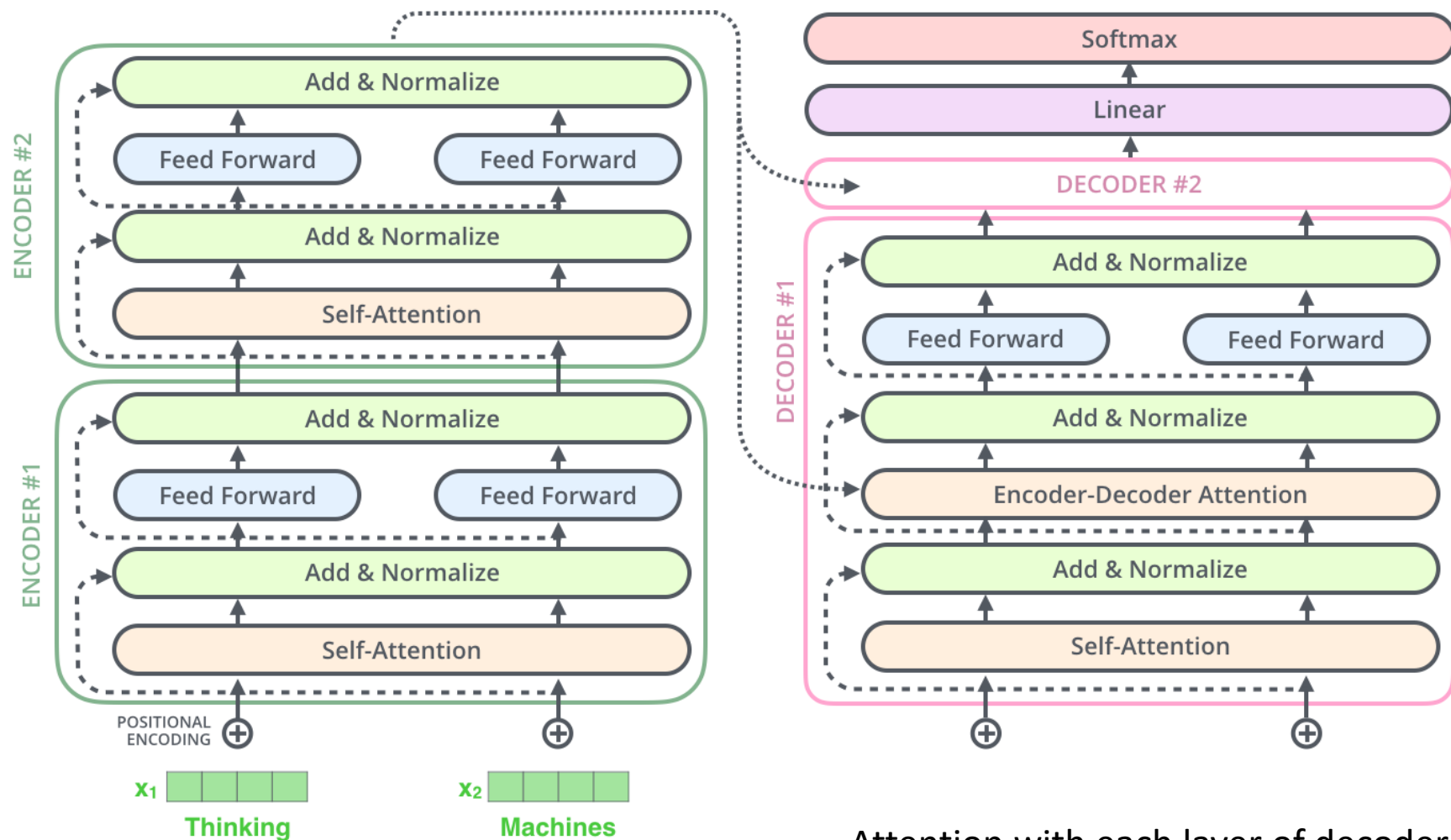
<http://jalammar.github.io/illustrated-transformer/>



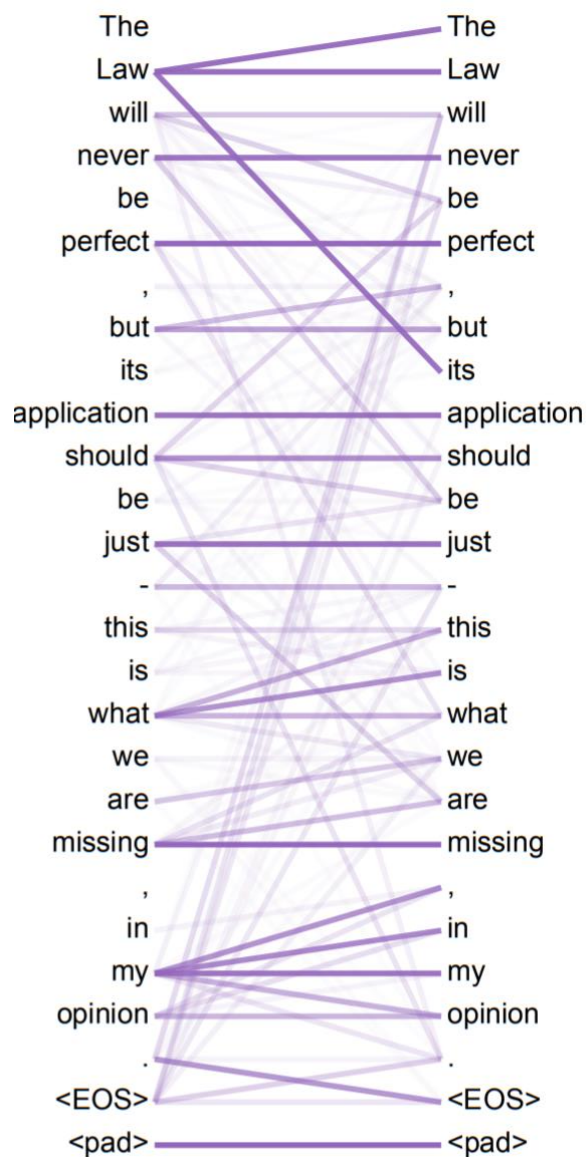
Transformer

Encoder可以有很多层

Encoder-Decoder



Attention Visualization



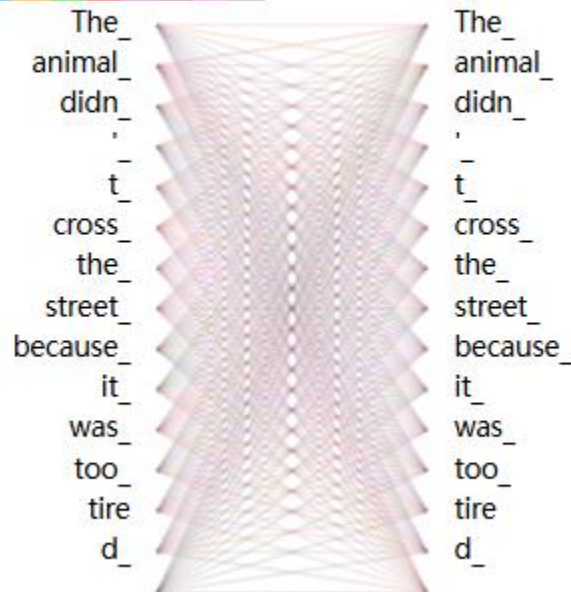
<https://arxiv.org/abs/1706.03762>

Attention Visualization

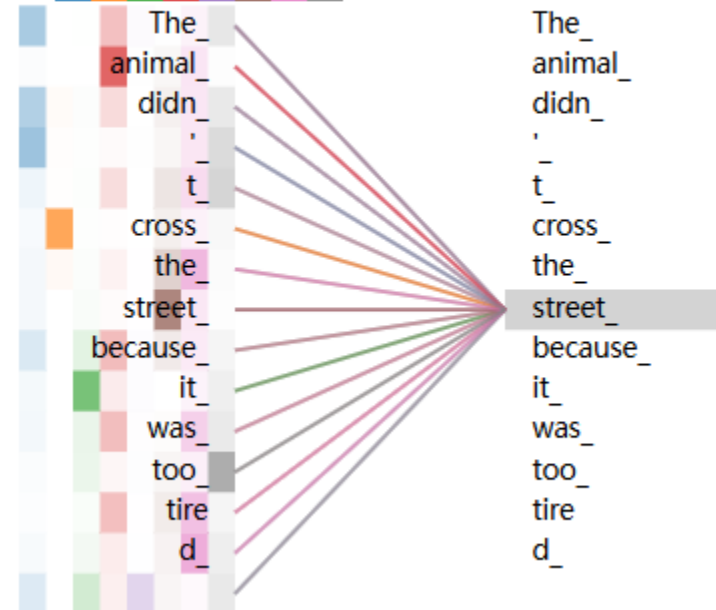
https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb#scrollTo=OJKU36QAfqOC

See `tf.nn.softmax_cross_entropy_with_logits_v2`.

Layer: 5 Attention: Input - Input



Layer: 5 Attention: Input - Input



Multi-head Attention

每个attention head学到的知识不同



Transformers in Time Series: A Survey

Qingsong Wen¹, Tian Zhou², Chaoli Zhang², Weiqi Chen², Ziqing Ma², Junchi Yan³, Liang Sun¹

¹DAMO Academy, Alibaba Group, Bellevue, USA

²DAMO Academy, Alibaba Group, Hangzhou, China

³Department of CSE, MoE Key Lab of Artificial Intelligence, Shanghai Jiao Tong University
{qingsong.wen, tian.zt, chaoli.zcl, jarvus.cwq, maziqing.mzq, liang.sun}@alibaba-inc.com,
yanjunchi@sjtu.edu.cn

