



机器学习

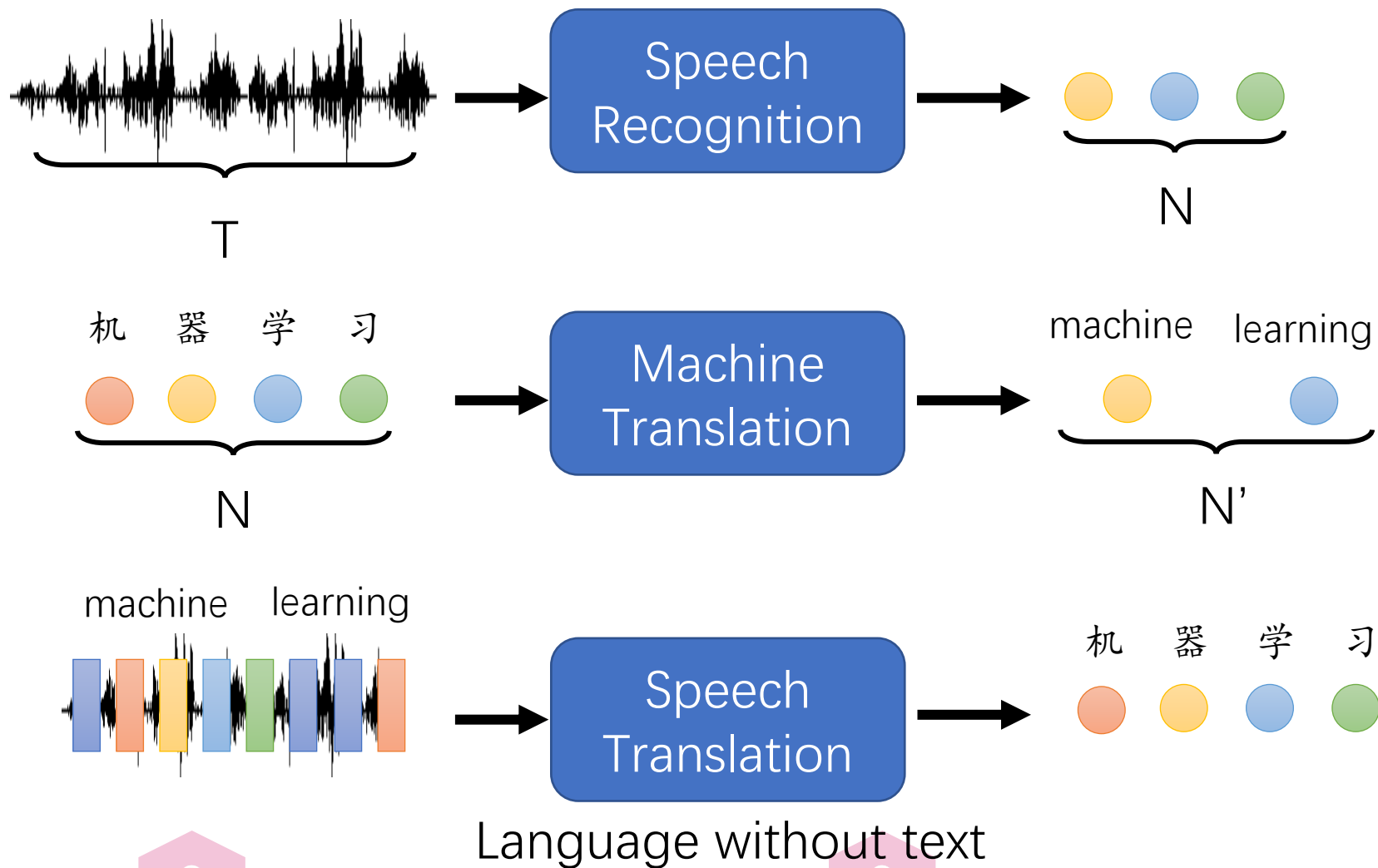
苏州大学计算机科学与技术学院

自然语言处理实验室

主讲：周夏冰

邮箱：zhouxiabing@suda.edu.cn

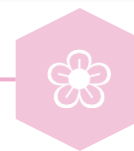
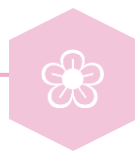
Sequence-to-sequence (Seq2seq)





01

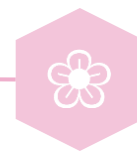
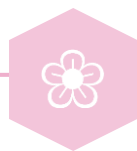
自注意力机制





注意力机制

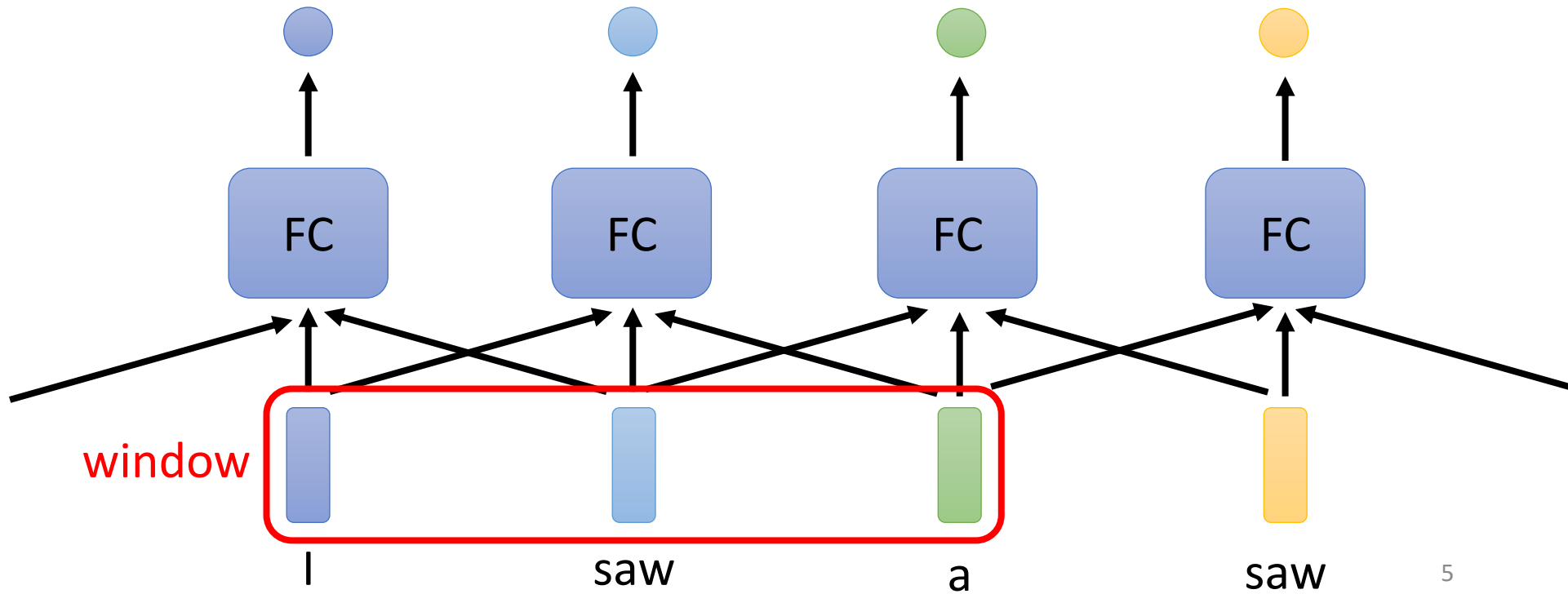
- 注意力机制（Attention Mechanism）是一种模仿人类视觉和认知系统的方法
- 在处理输入数据时集中注意力于相关的部分

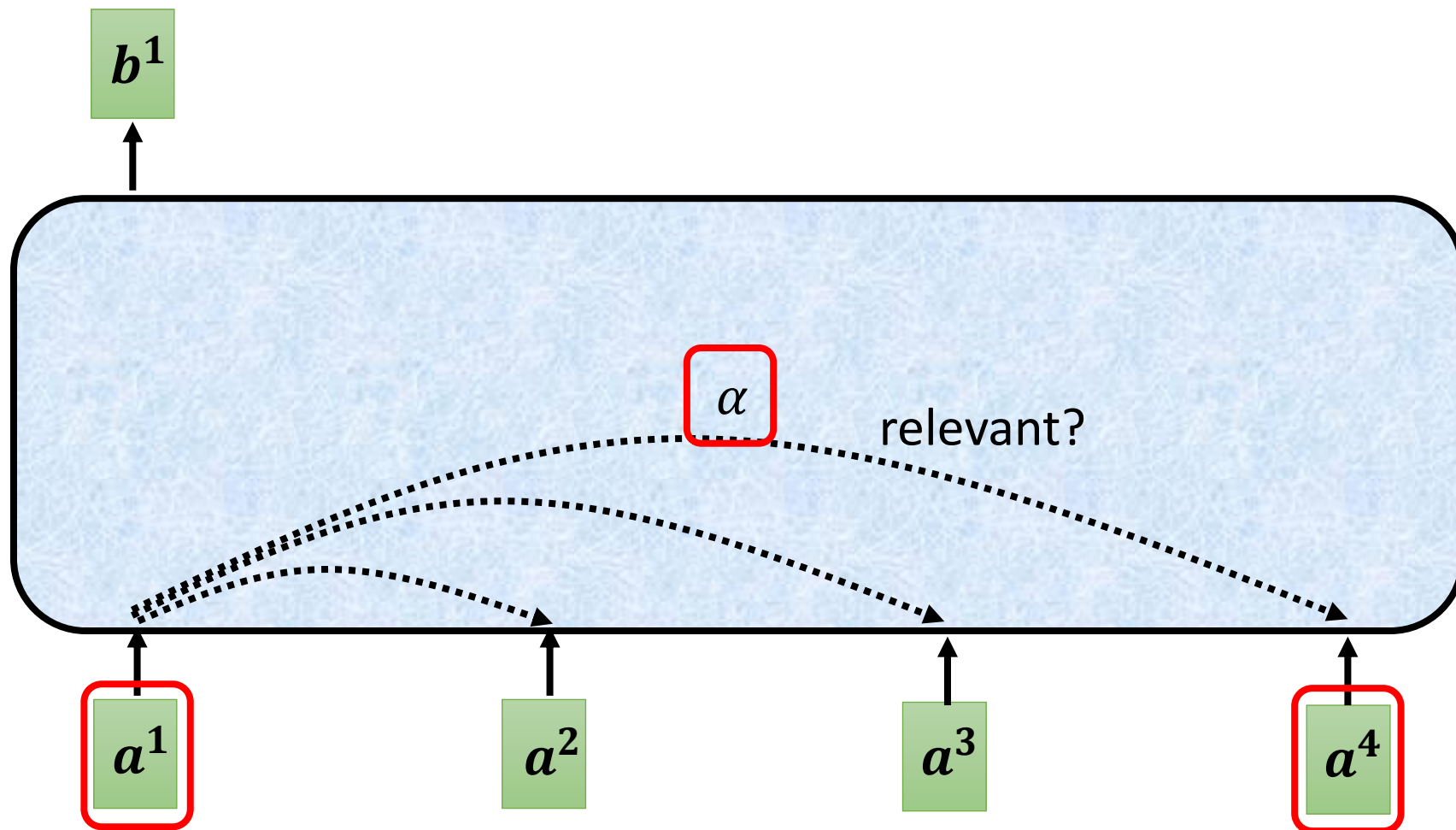


Sequence labeling

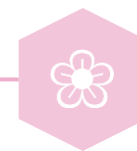
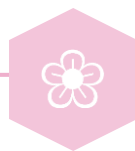
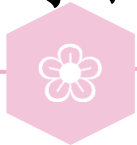


FC Fully-connected

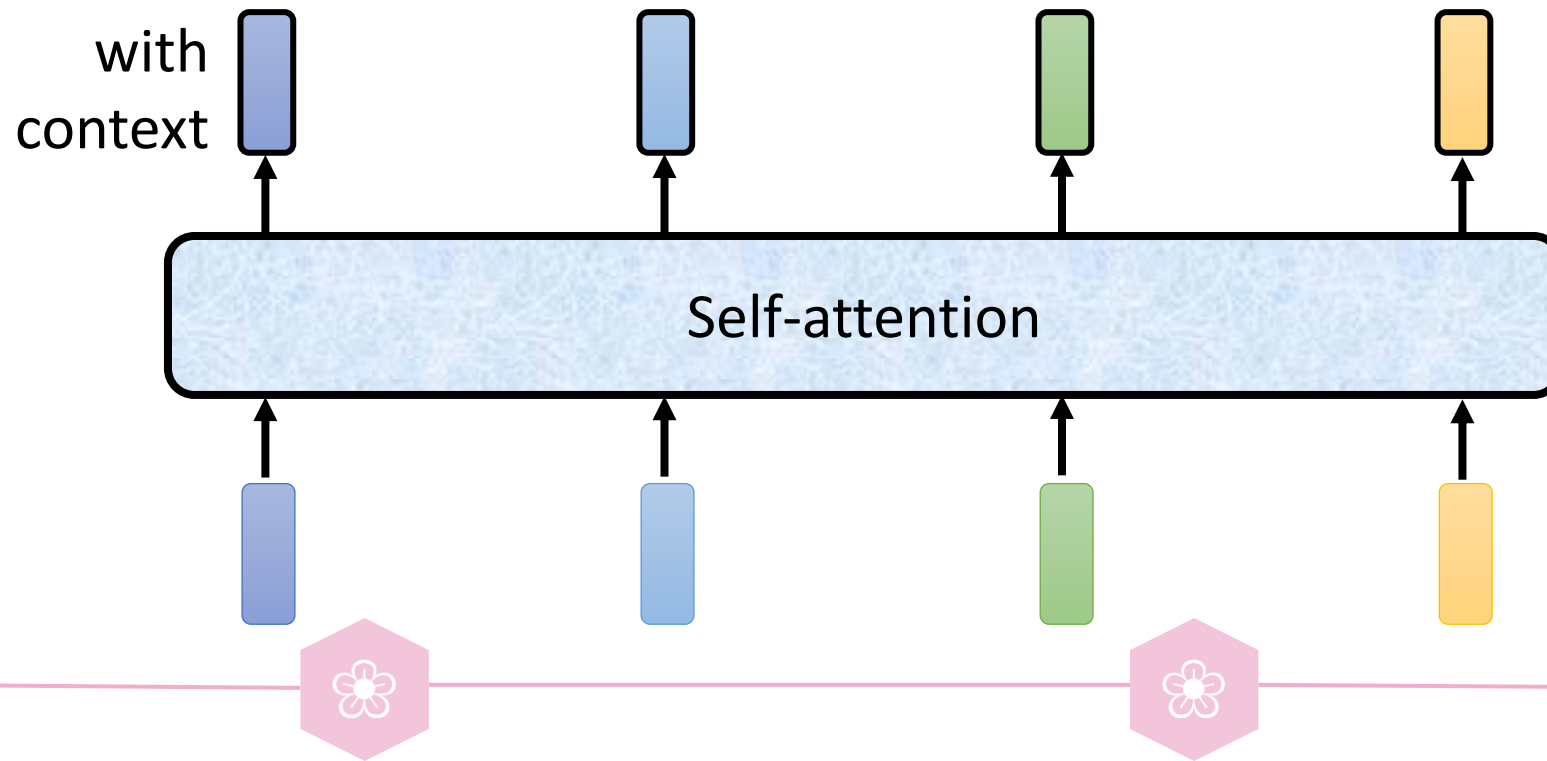




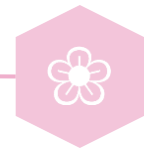
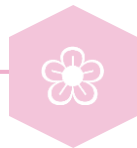
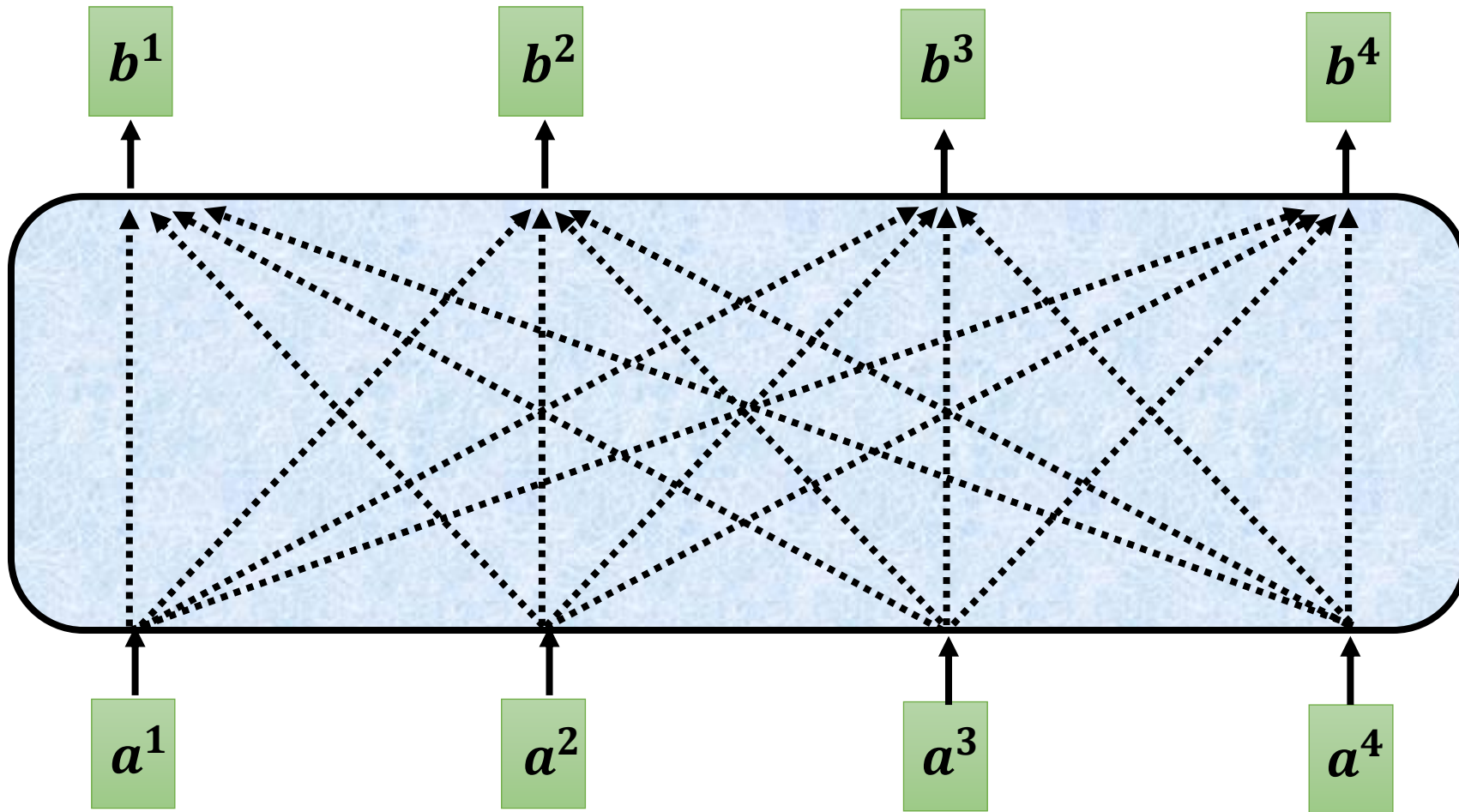
寻找输入中相关的信息



Self-attention

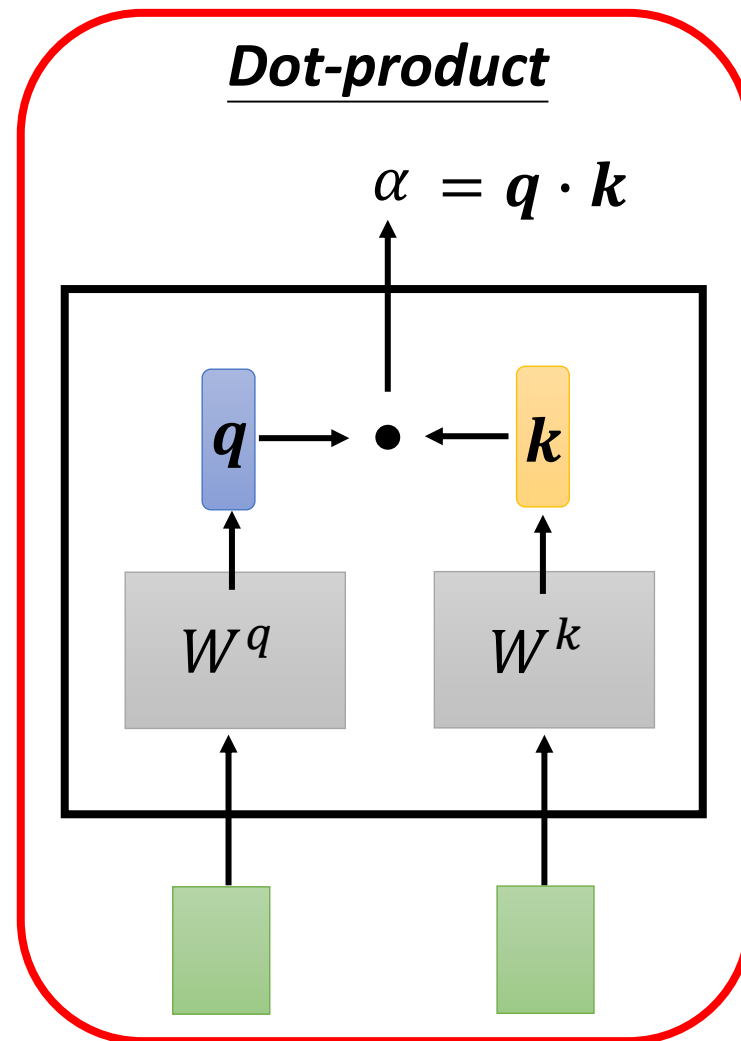


Self-attention

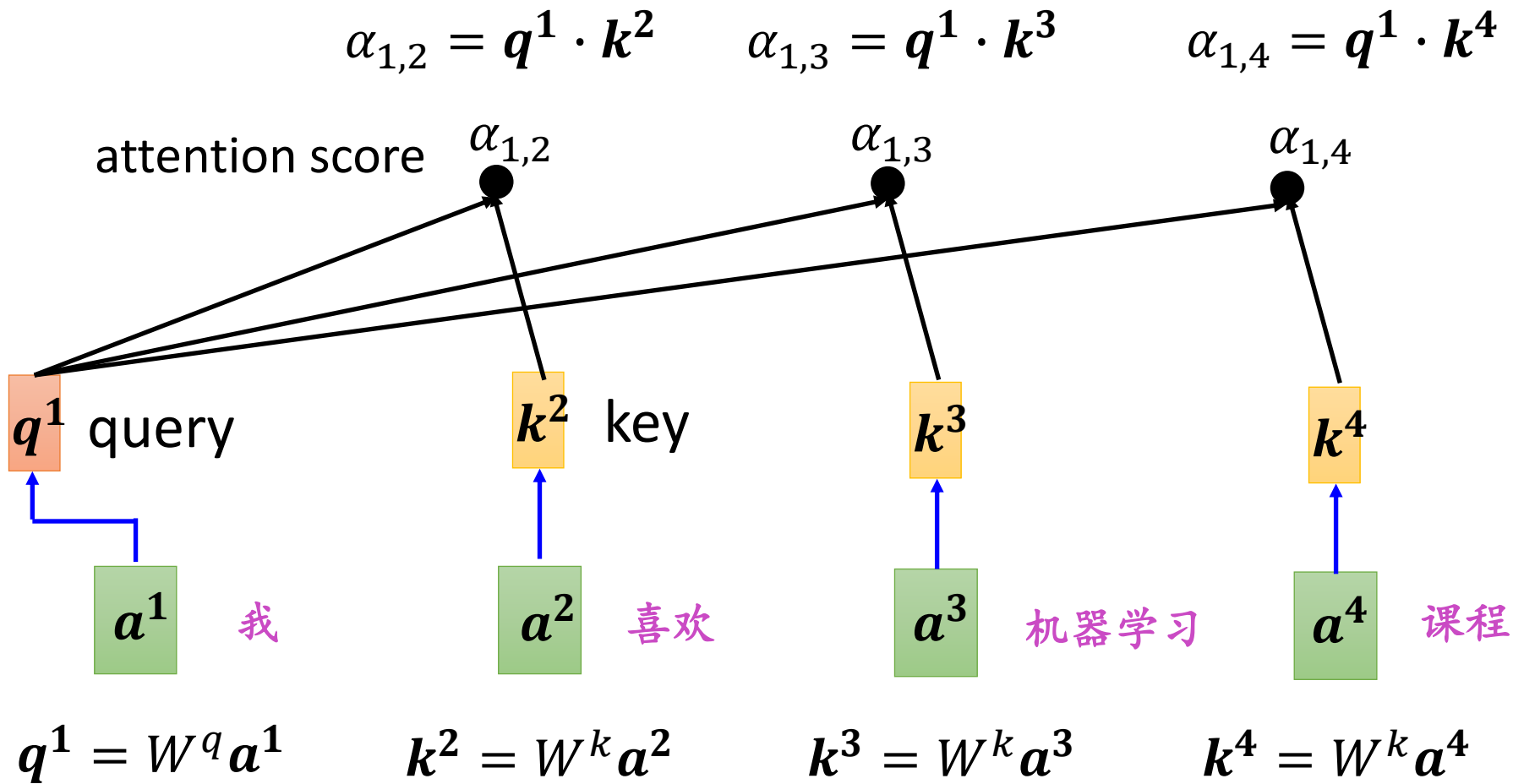


Self-attention

- Key, Query, Value
- Query: 查询向量
- Key: 键, 用来计算和Query相关性
- Value: 值, 用来获取信息的向量

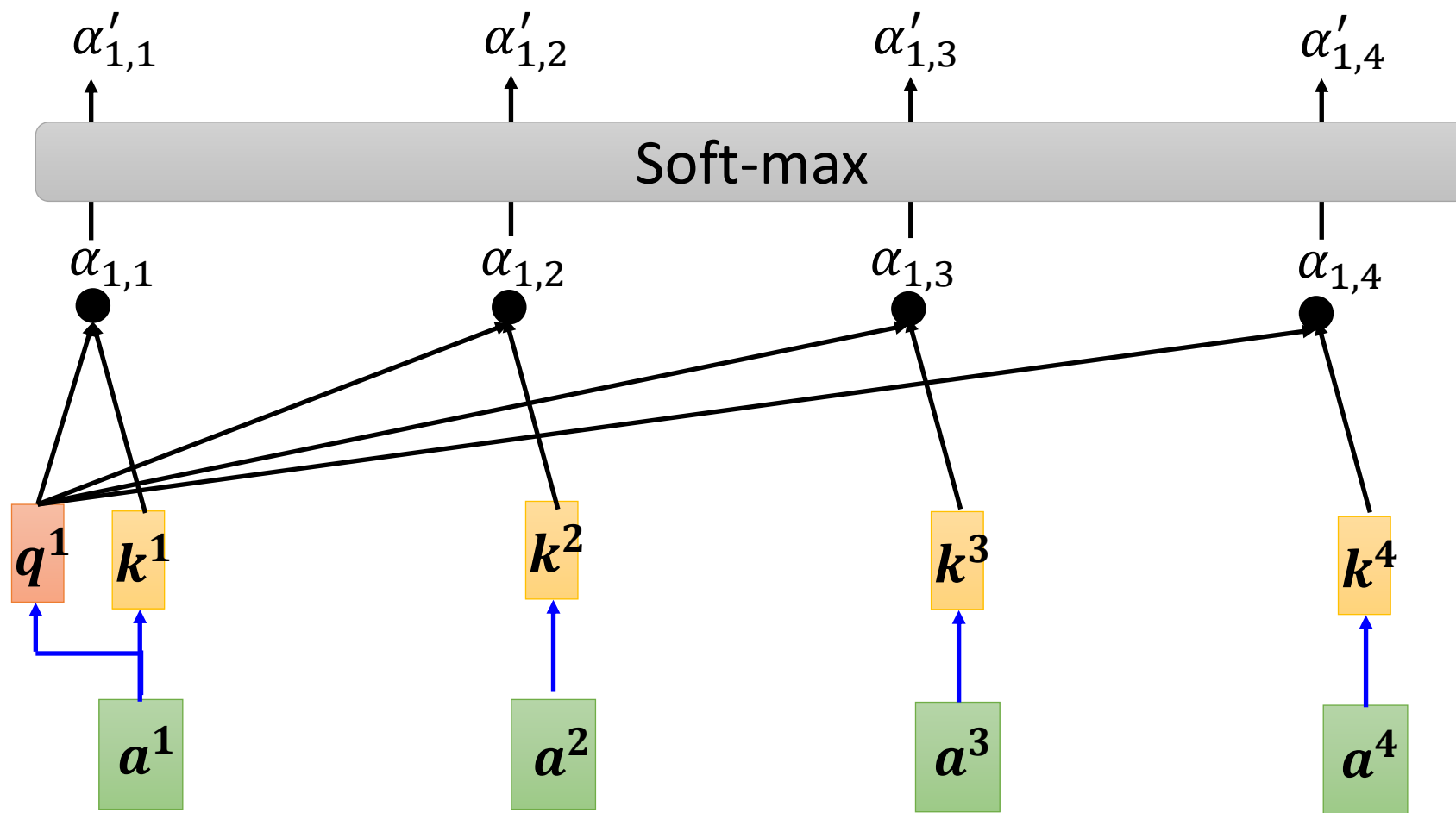


Self-attention



Self-attention

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



$$q^1 = W^q a^1$$

$$k^2 = W^k a^2$$

$$k^3 = W^k a^3$$

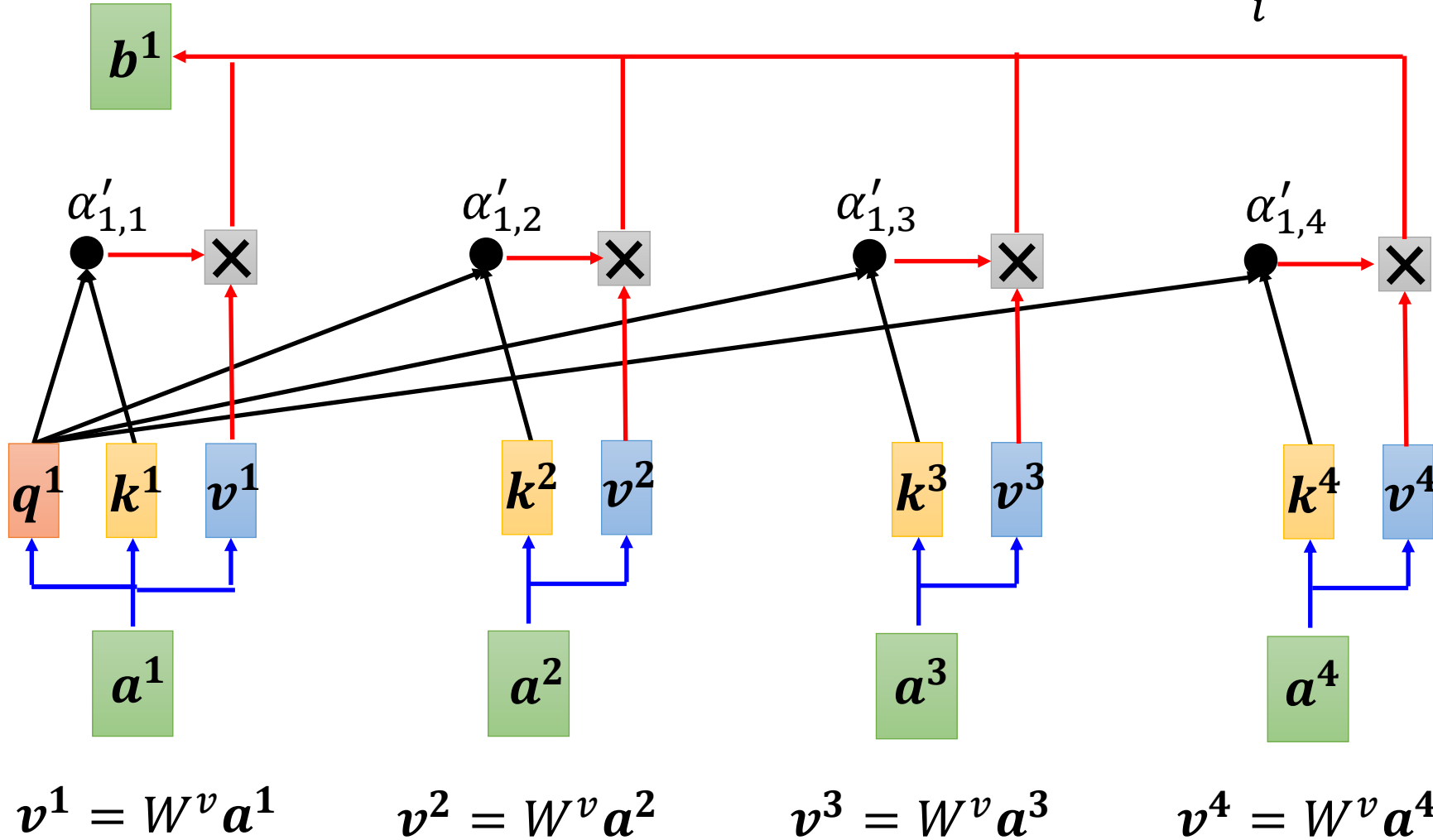
$$k^4 = W^k a^4$$

$$k^1 = W^k a^1$$



Self-attention

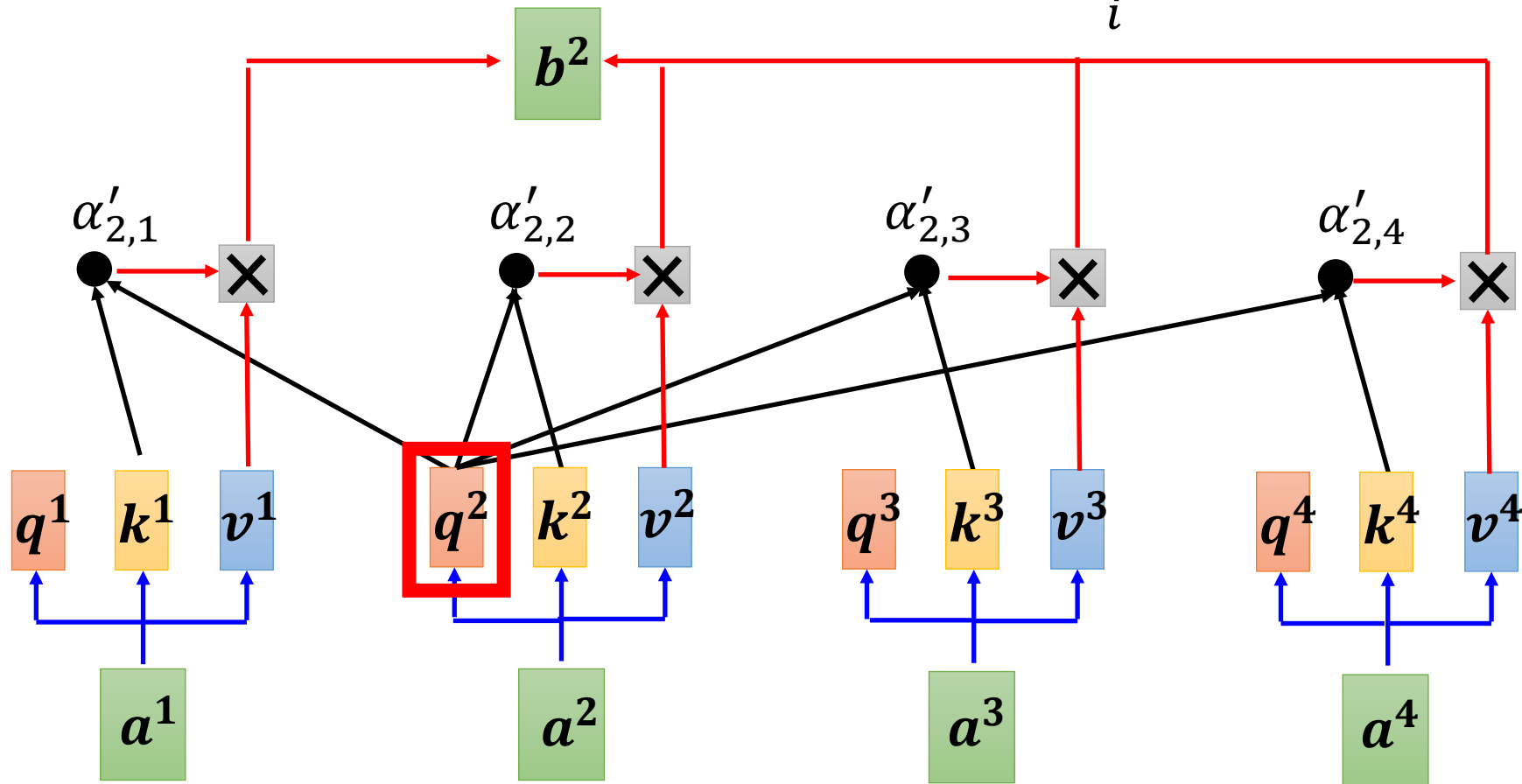
$$b^1 = \sum_i \alpha'_{1,i} v^i$$



根据权重抽取信息

Self-attention

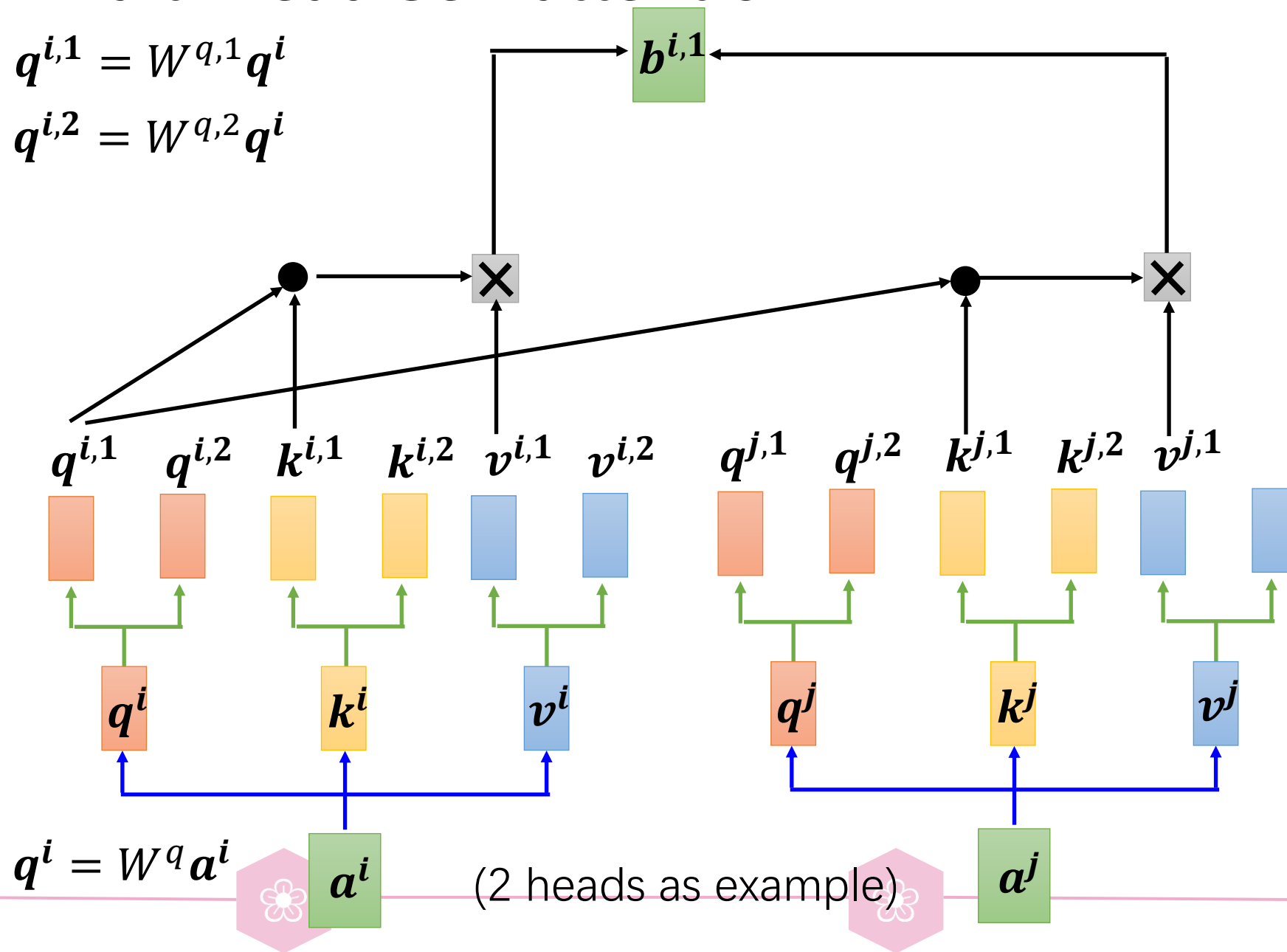
$$b^2 = \sum_i \alpha'_{2,i} v^i$$



Multi-head Self-attention

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

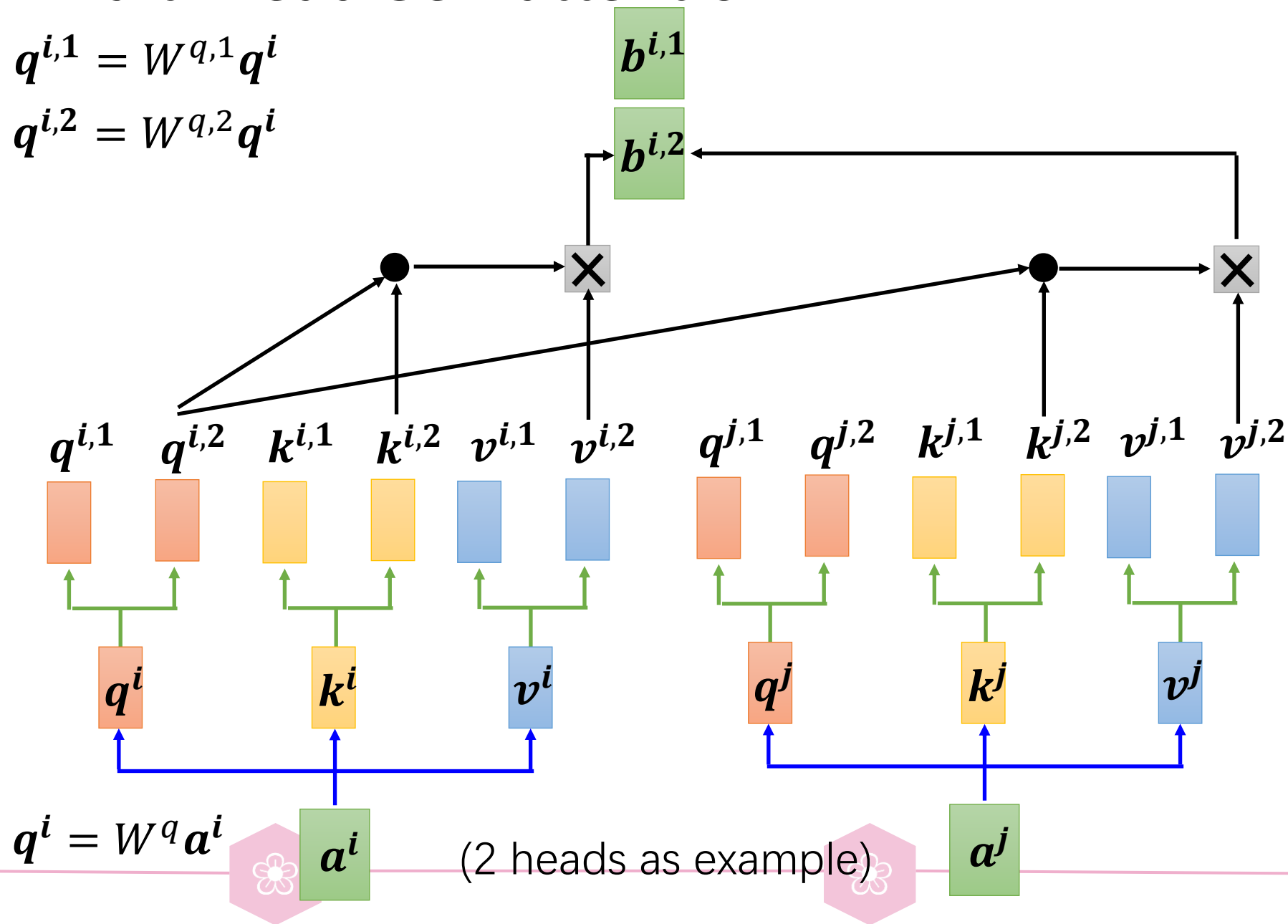
$$q^{i,2} = W^{q,2} q^i$$



Multi-head Self-attention

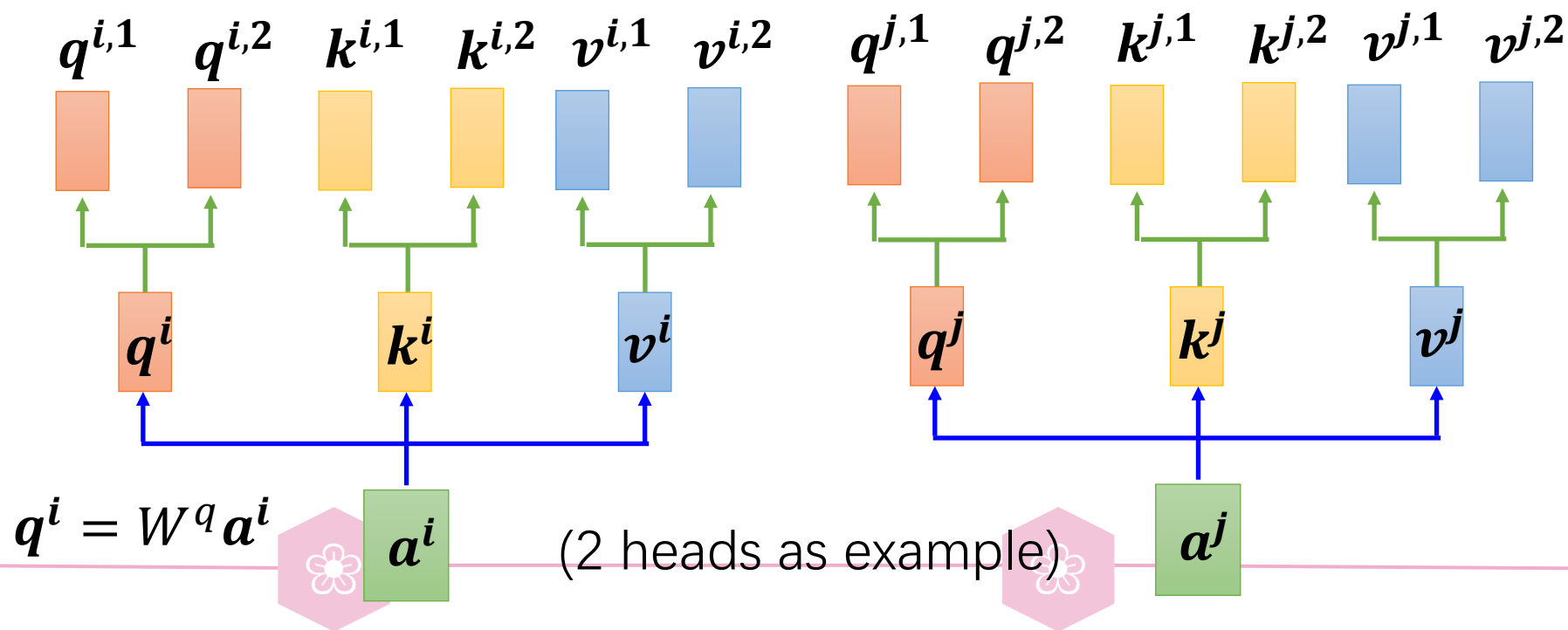
$$q^{i,1} = W^{q,1} q^i$$

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

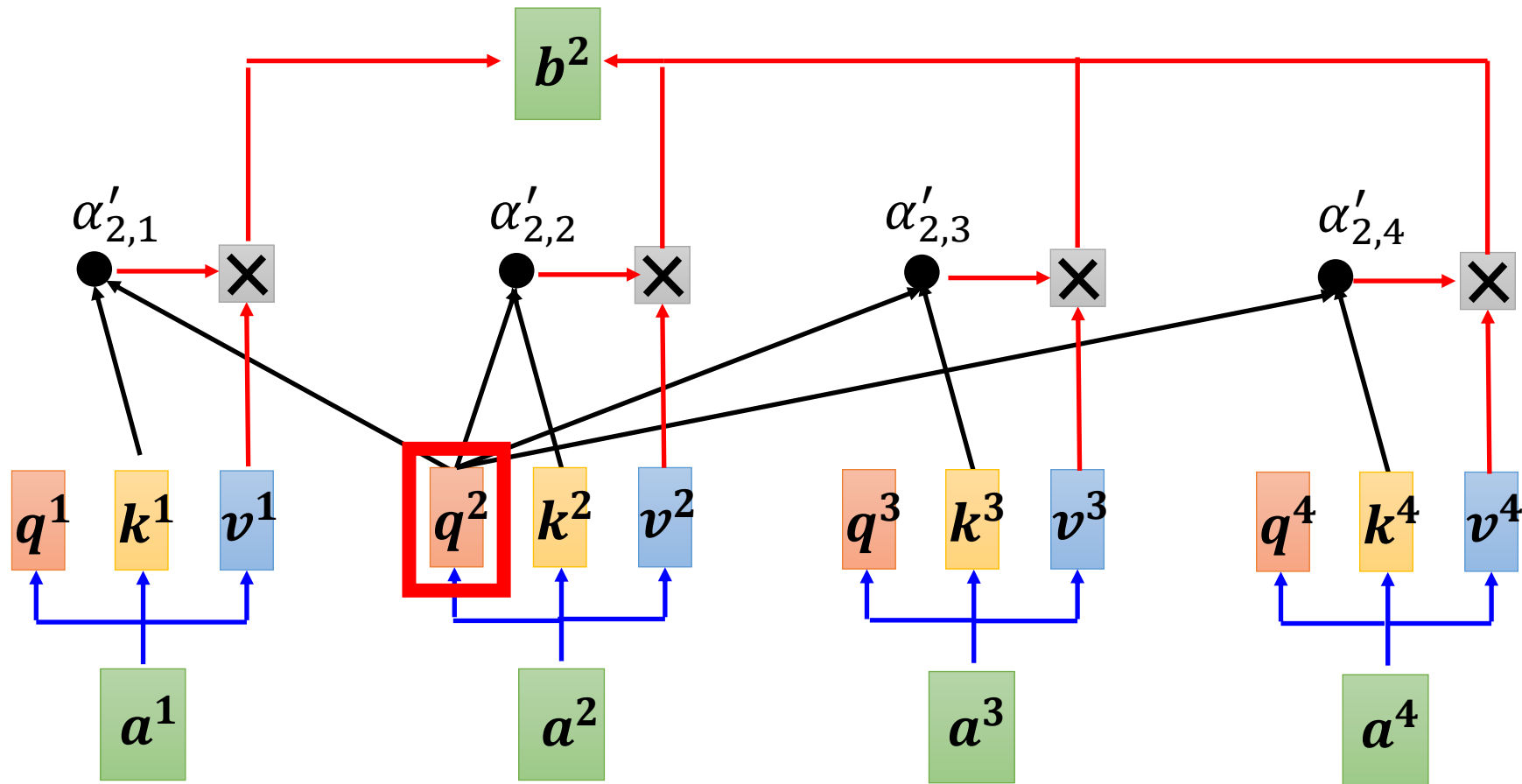


Multi-head Self-attention

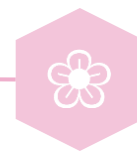
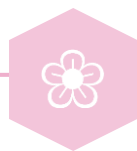
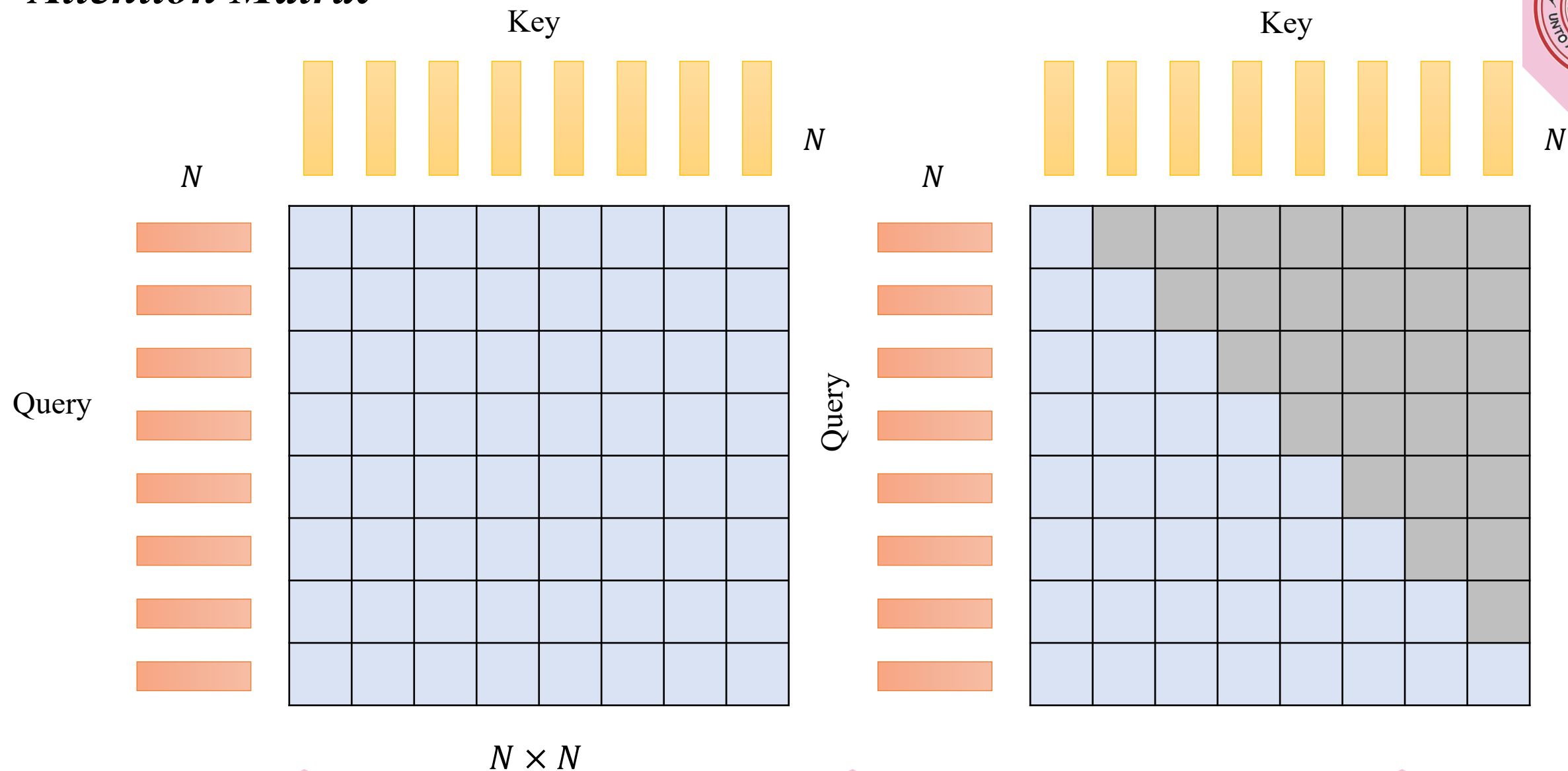
$$b^i = W^O \begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix}$$



Self-attention \rightarrow Masked Self-attention

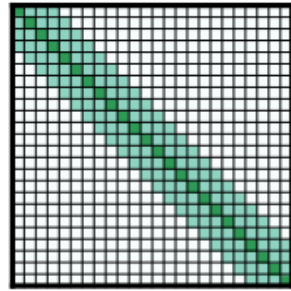


Attention Matrix

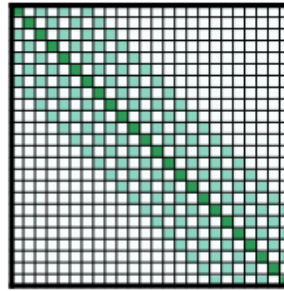


- Longformer

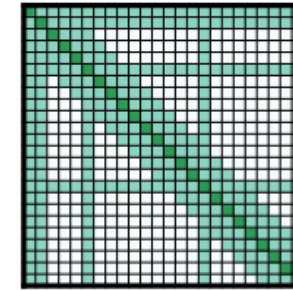
<https://arxiv.org/abs/2004.05150>



(b) Sliding window attention



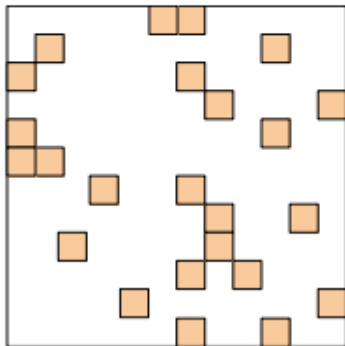
(c) Dilated sliding window



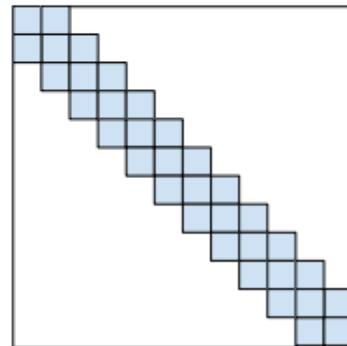
(d) Global+sliding window

- Big Bird

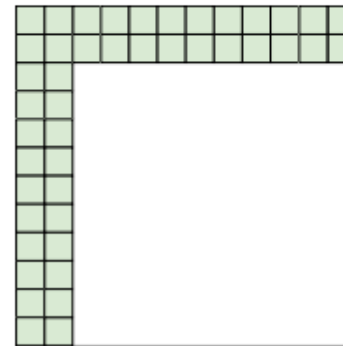
<https://arxiv.org/abs/2007.14062>



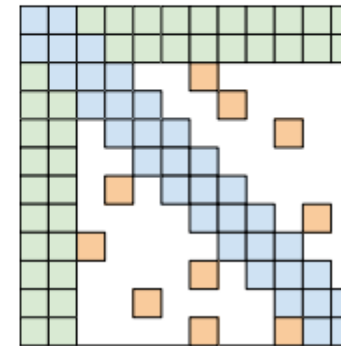
(a) Random attention



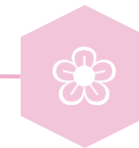
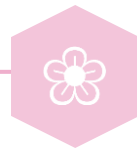
(b) Window attention



(c) Global Attention

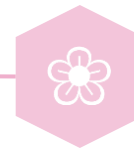
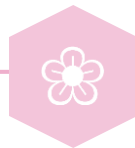
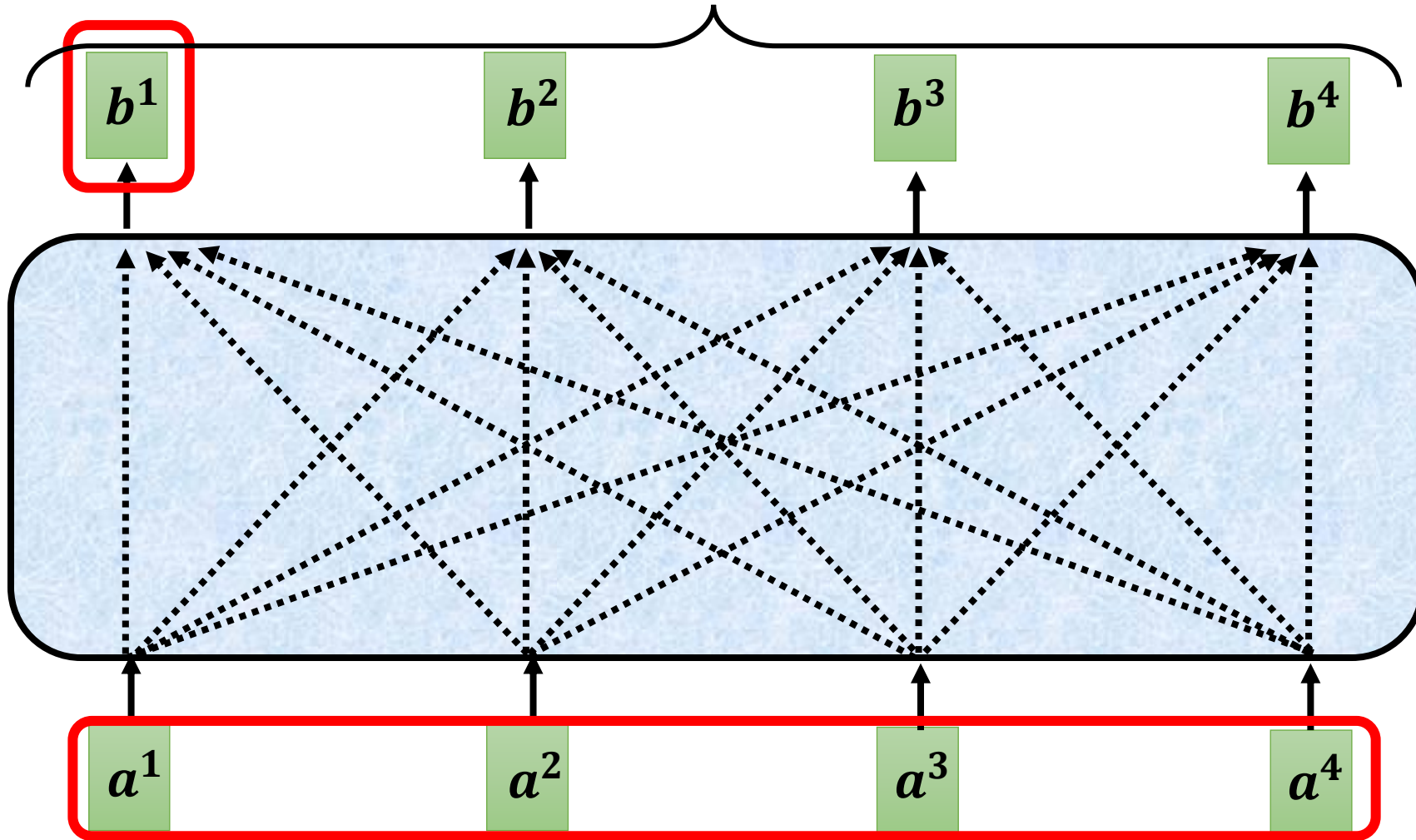


(d) BIGBIRD



Self-attention

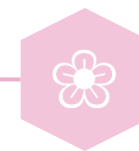
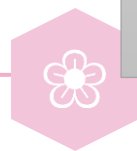
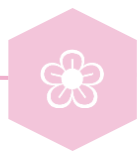
parallel



$$\begin{aligned} Q &= W^q I \\ K &= W^k I \\ V &= W^v I \end{aligned}$$

$$A' \leftarrow A = K^T Q$$

$$0 = V A'$$

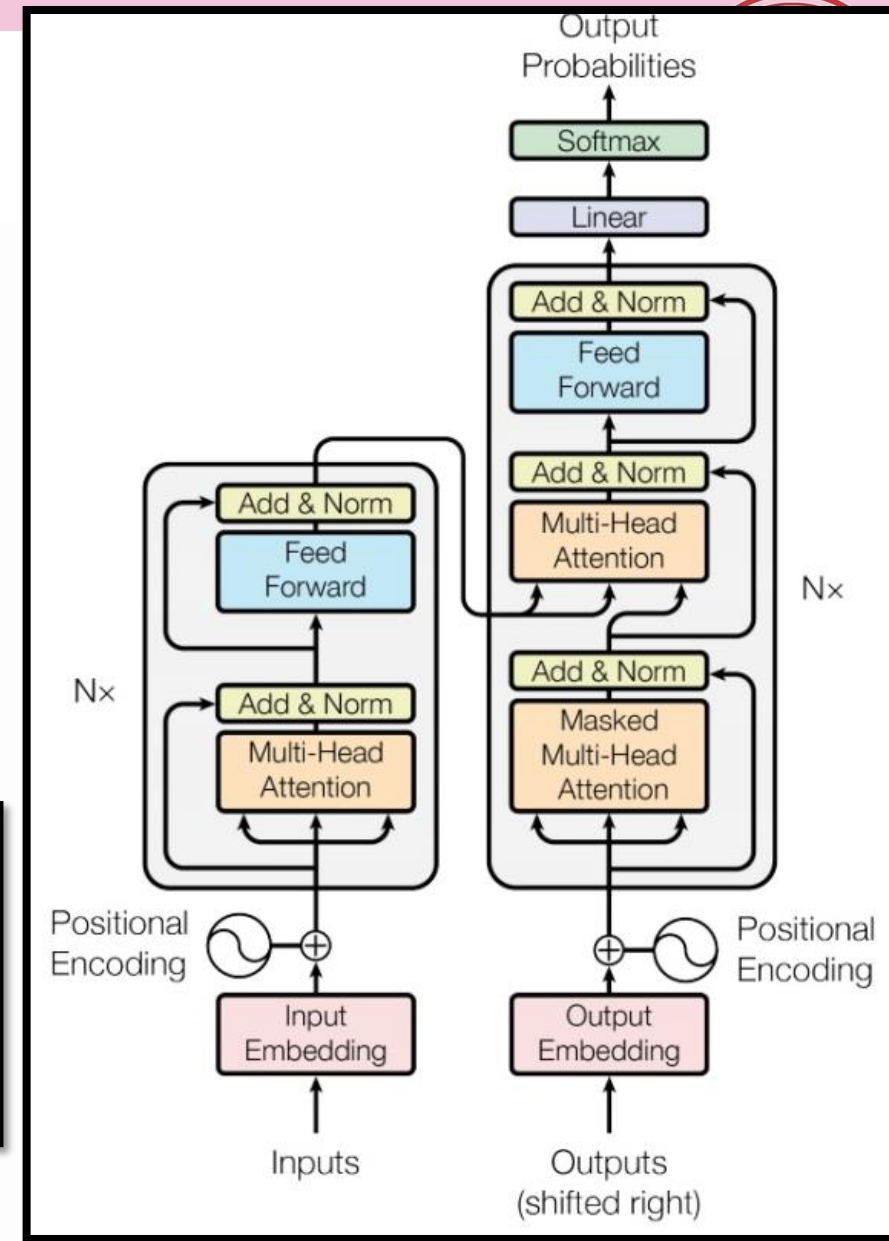
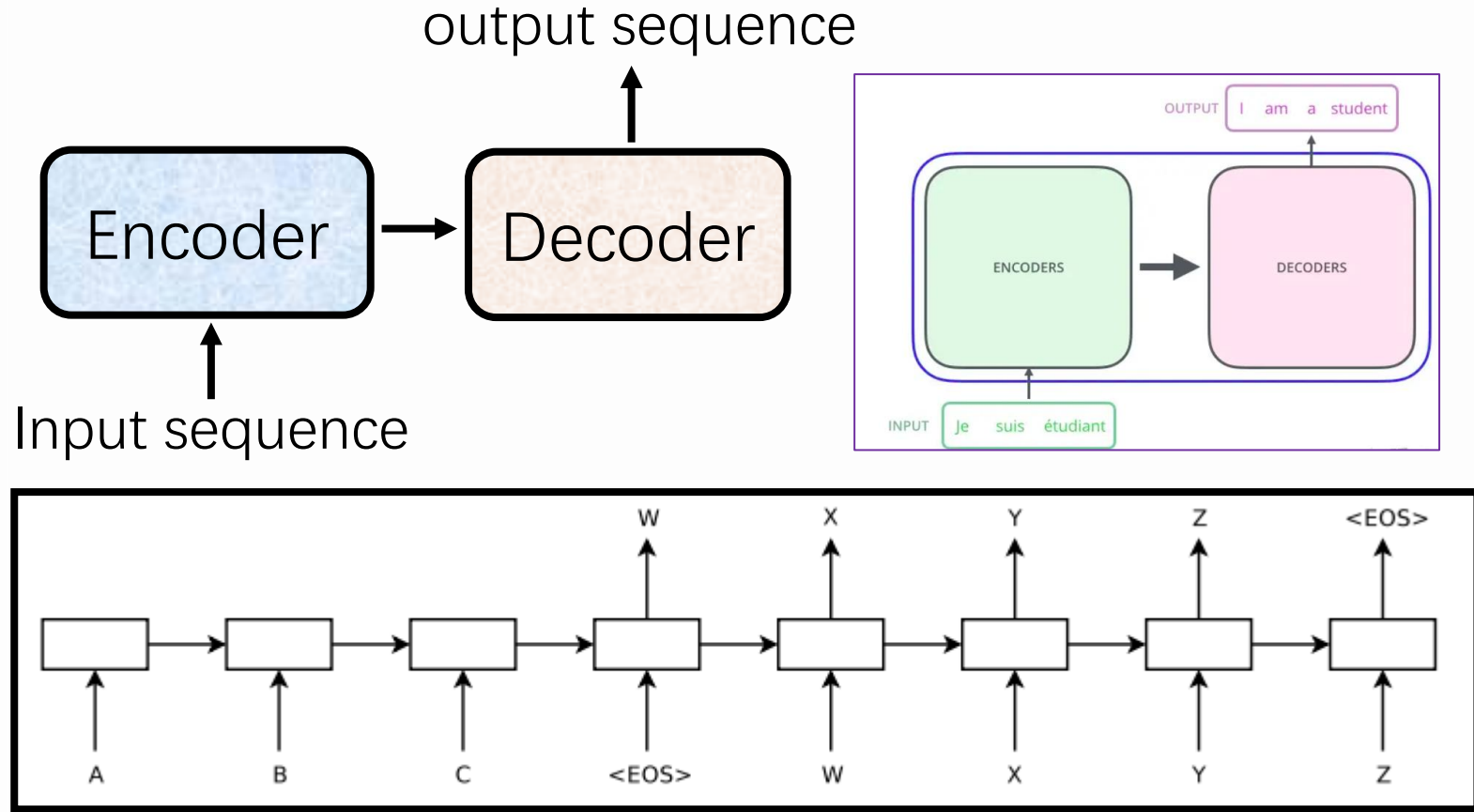


02

Transformer

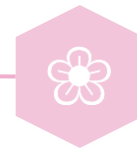
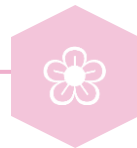
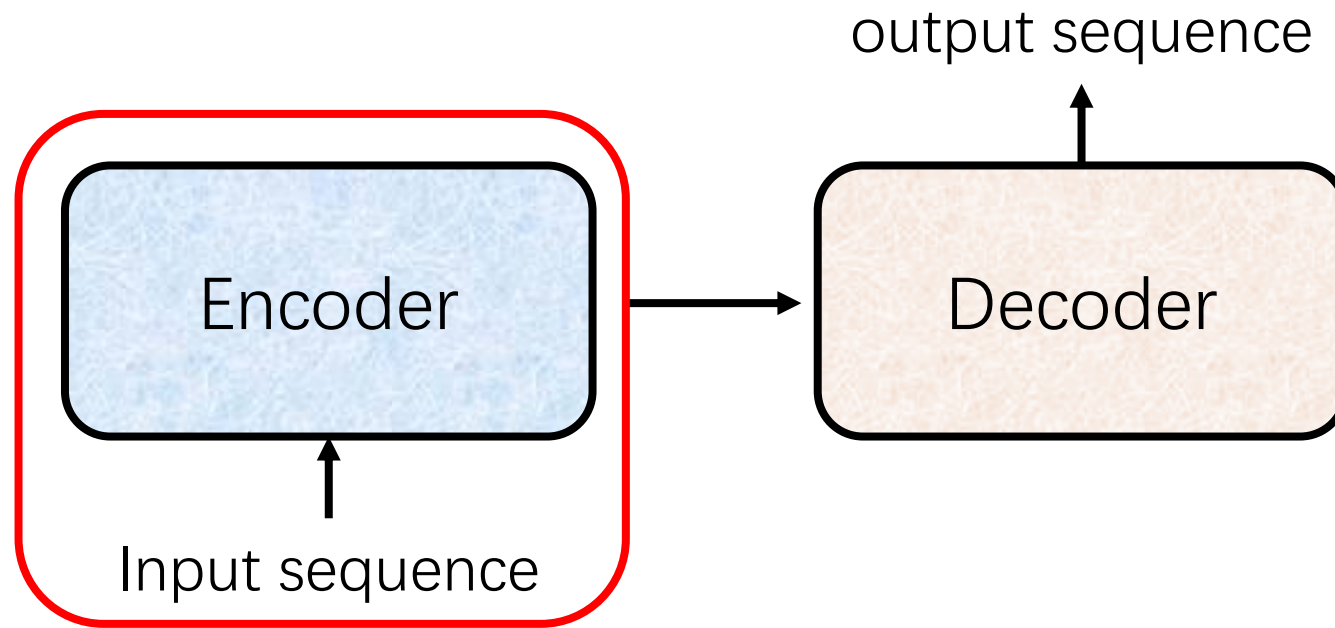


Seq2seq



Transformer

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

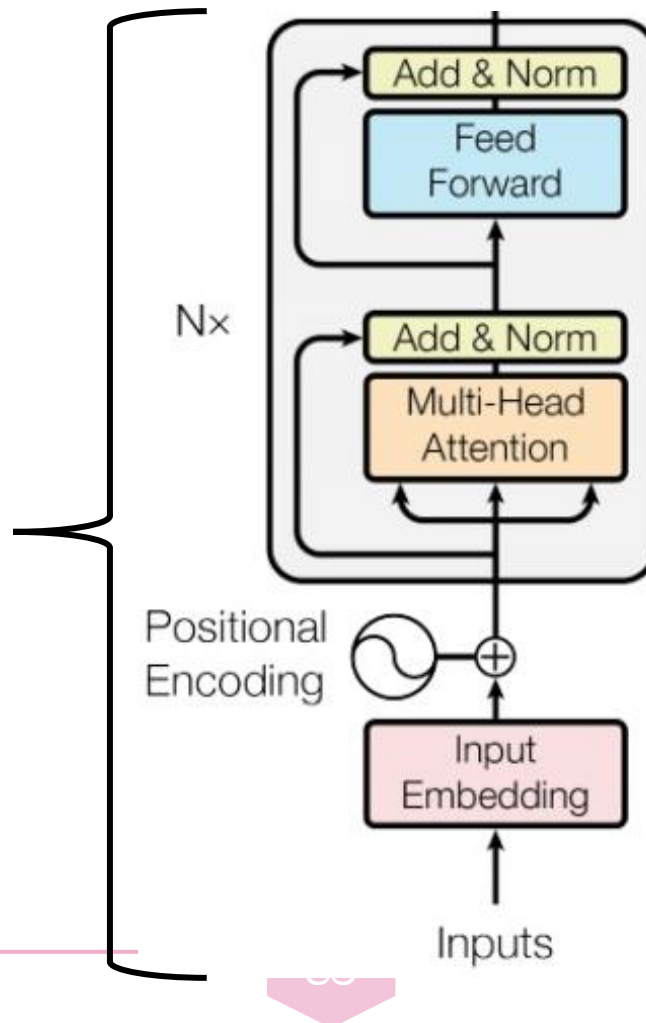
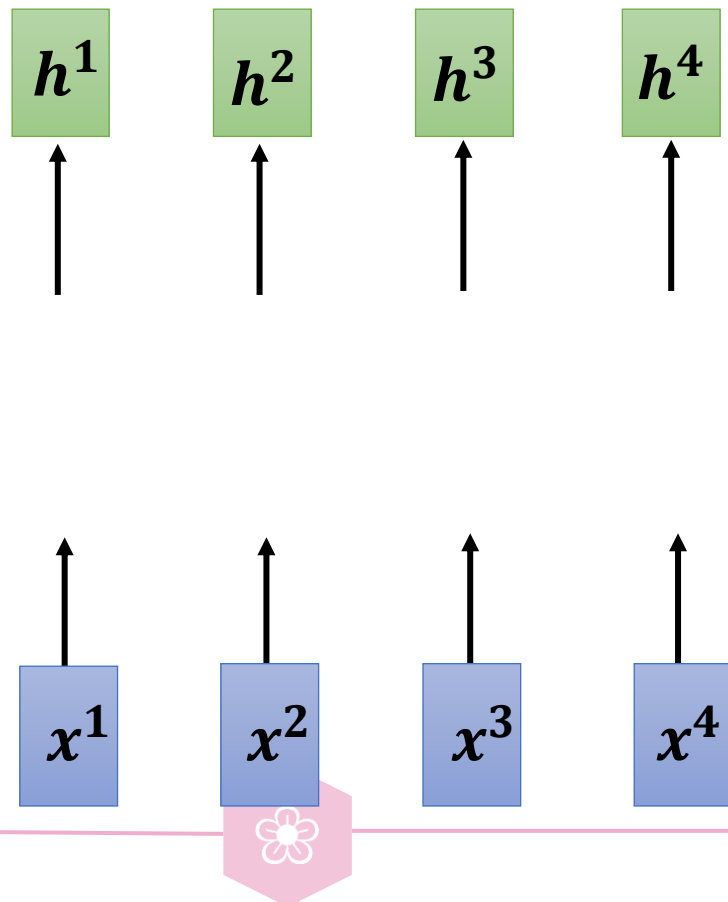


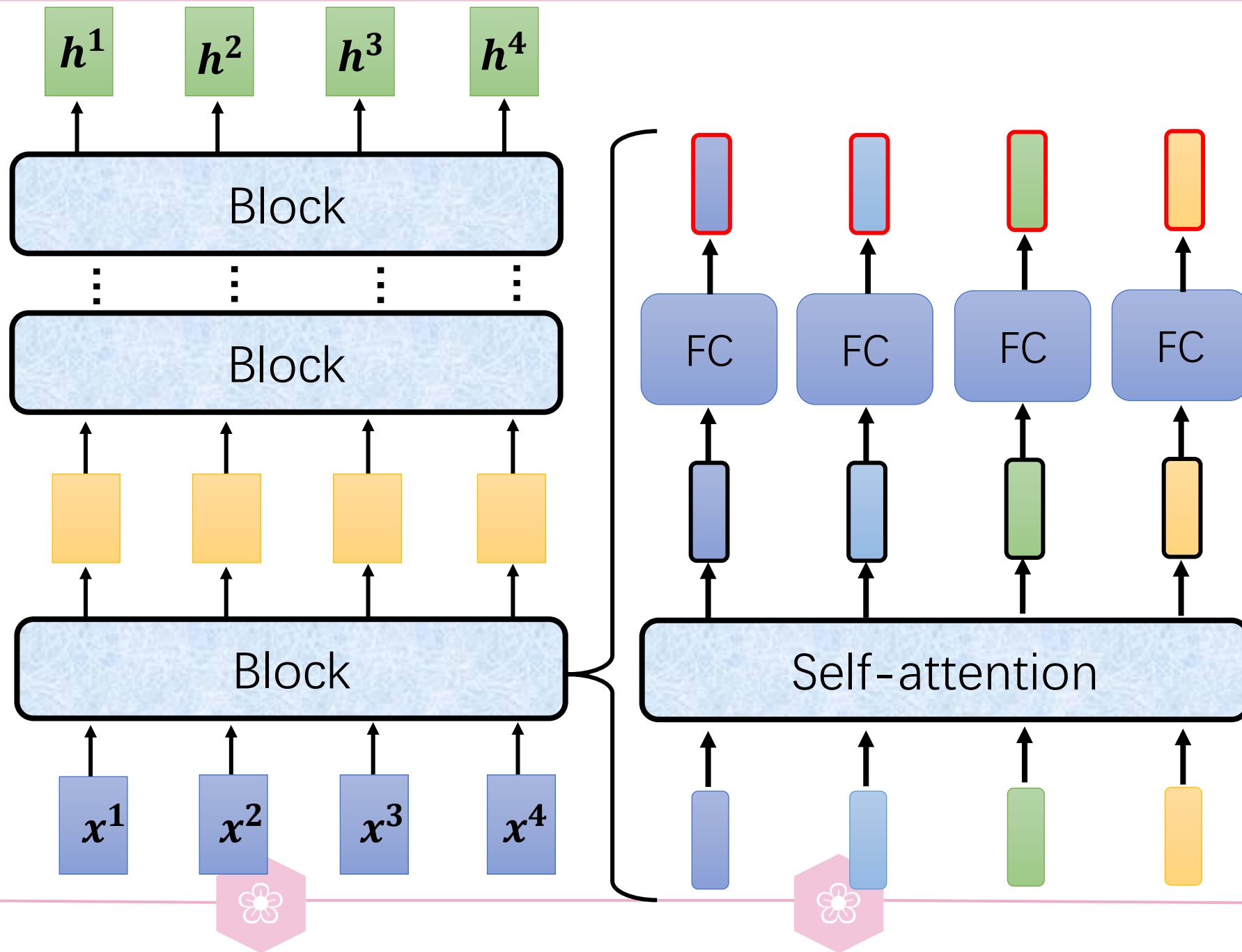
Encoder

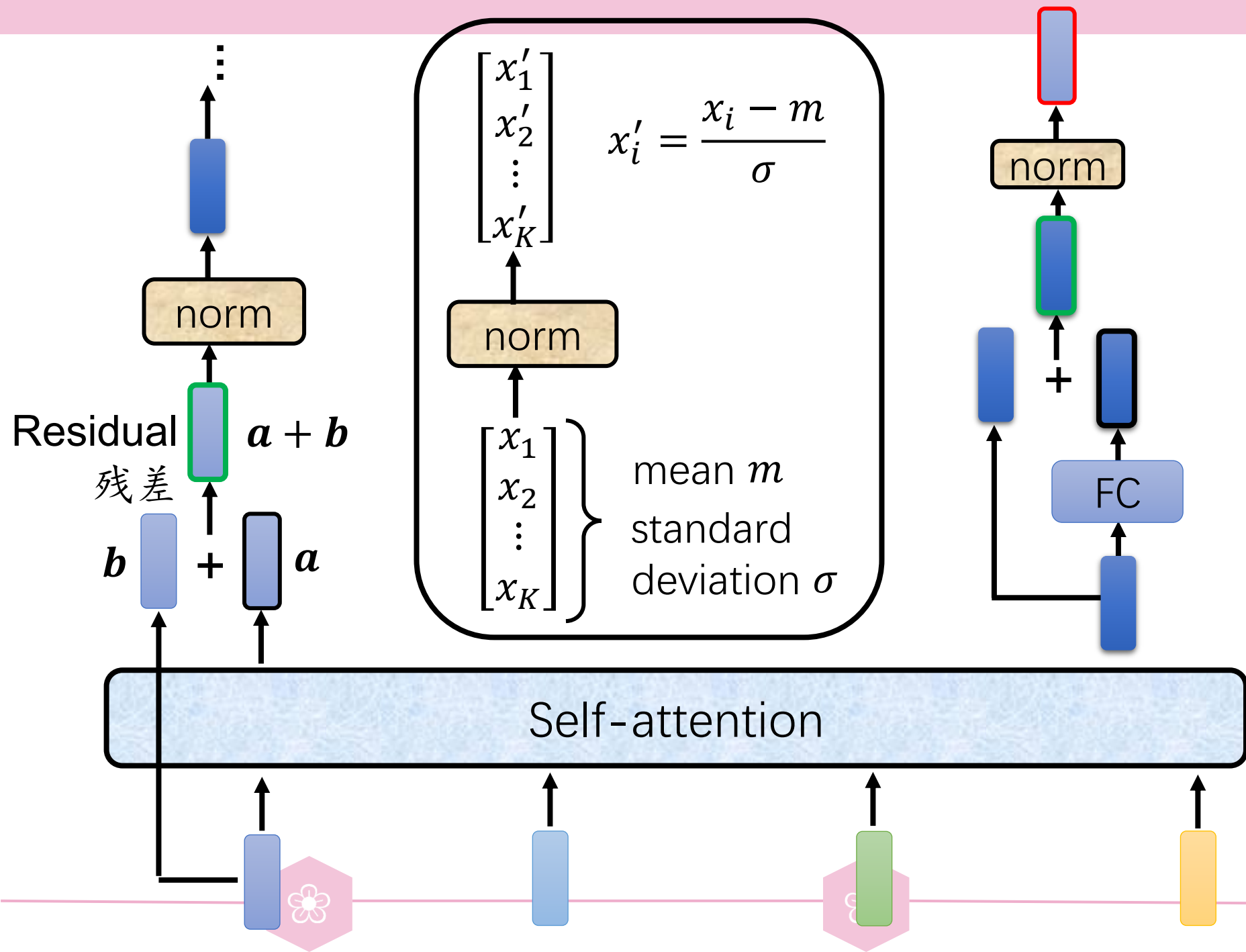


You can use **RNN** or **CNN**.

Transformer's Encoder







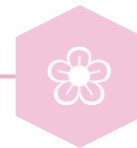
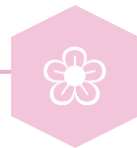
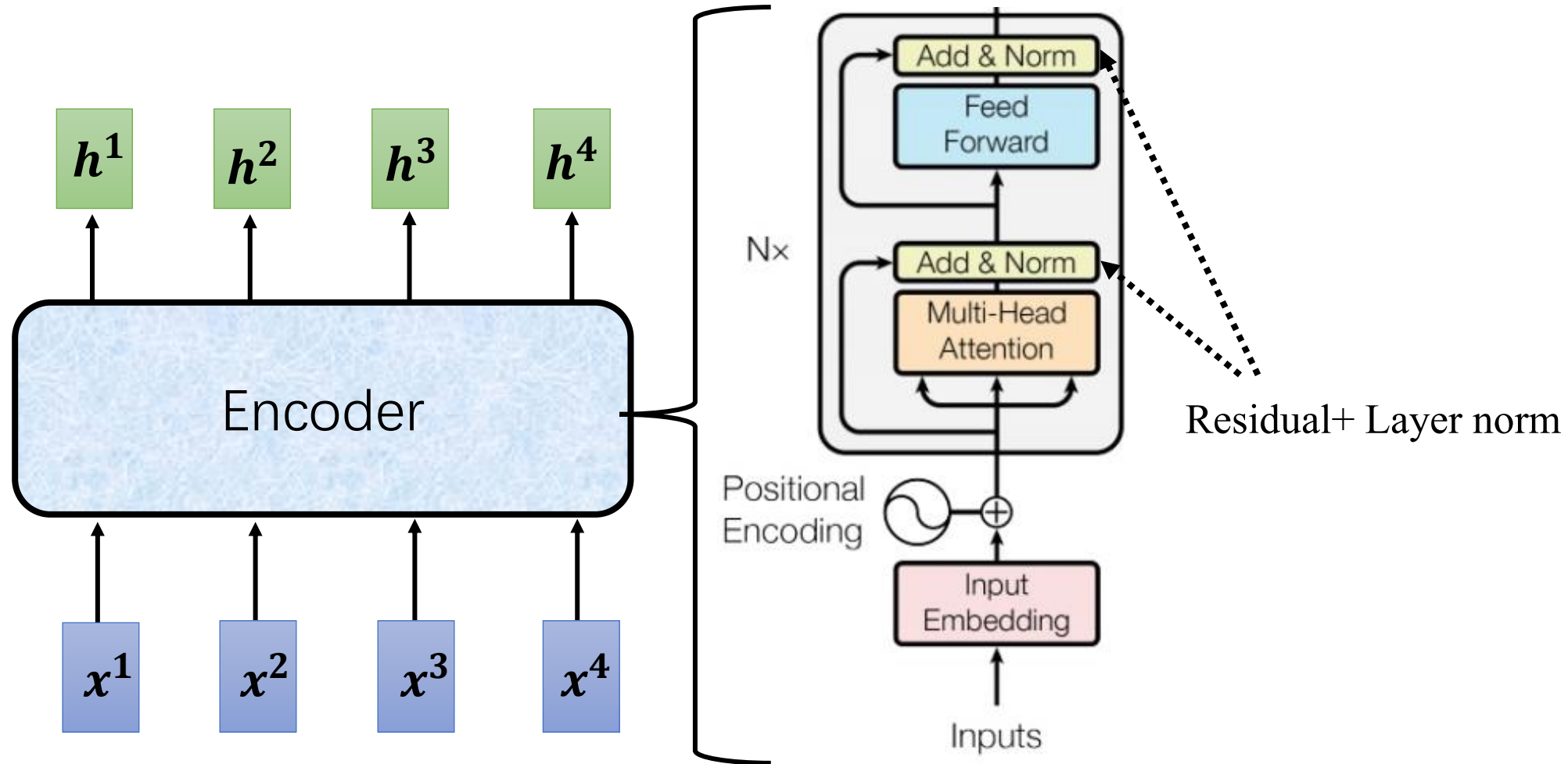
残差：
稳定训练、帮助梯度传播、避免退化问题

Norm：
稳定训练、加速收敛

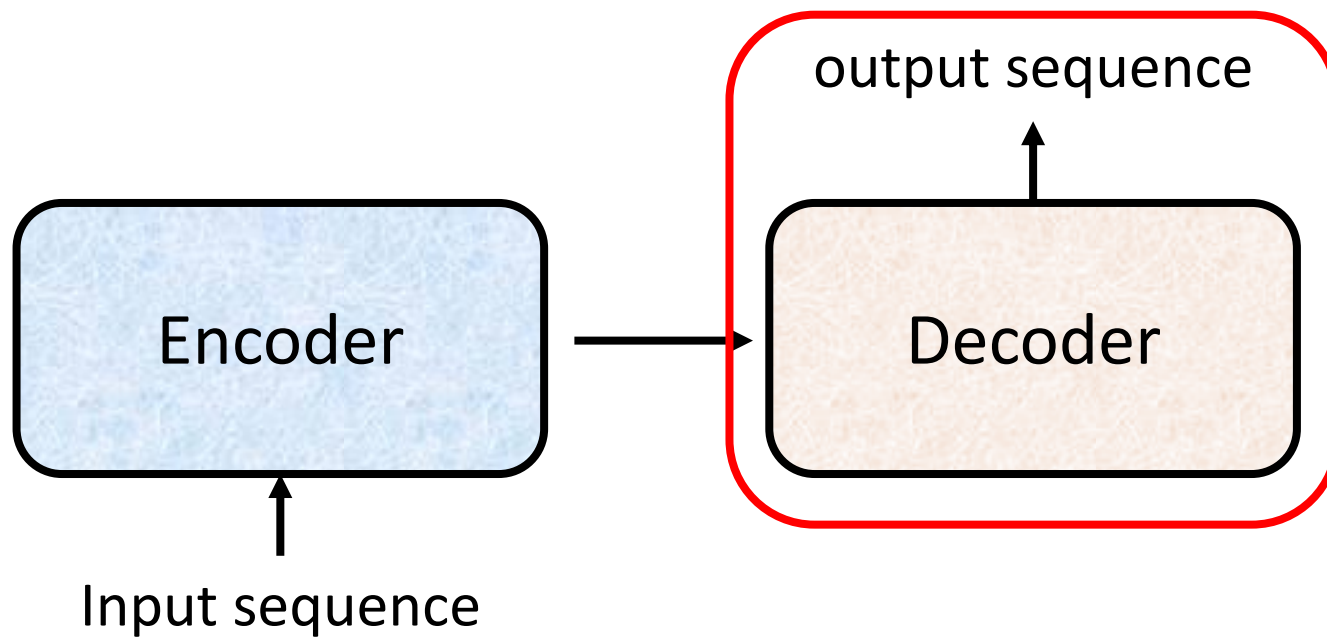
LayerNorm: 单个样本所有特征维度归一化

BatchNorm: 一个batch中同一通道的所有样本归一化 (CNN)

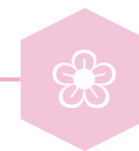
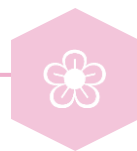
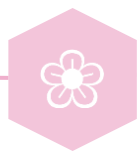
Encoder



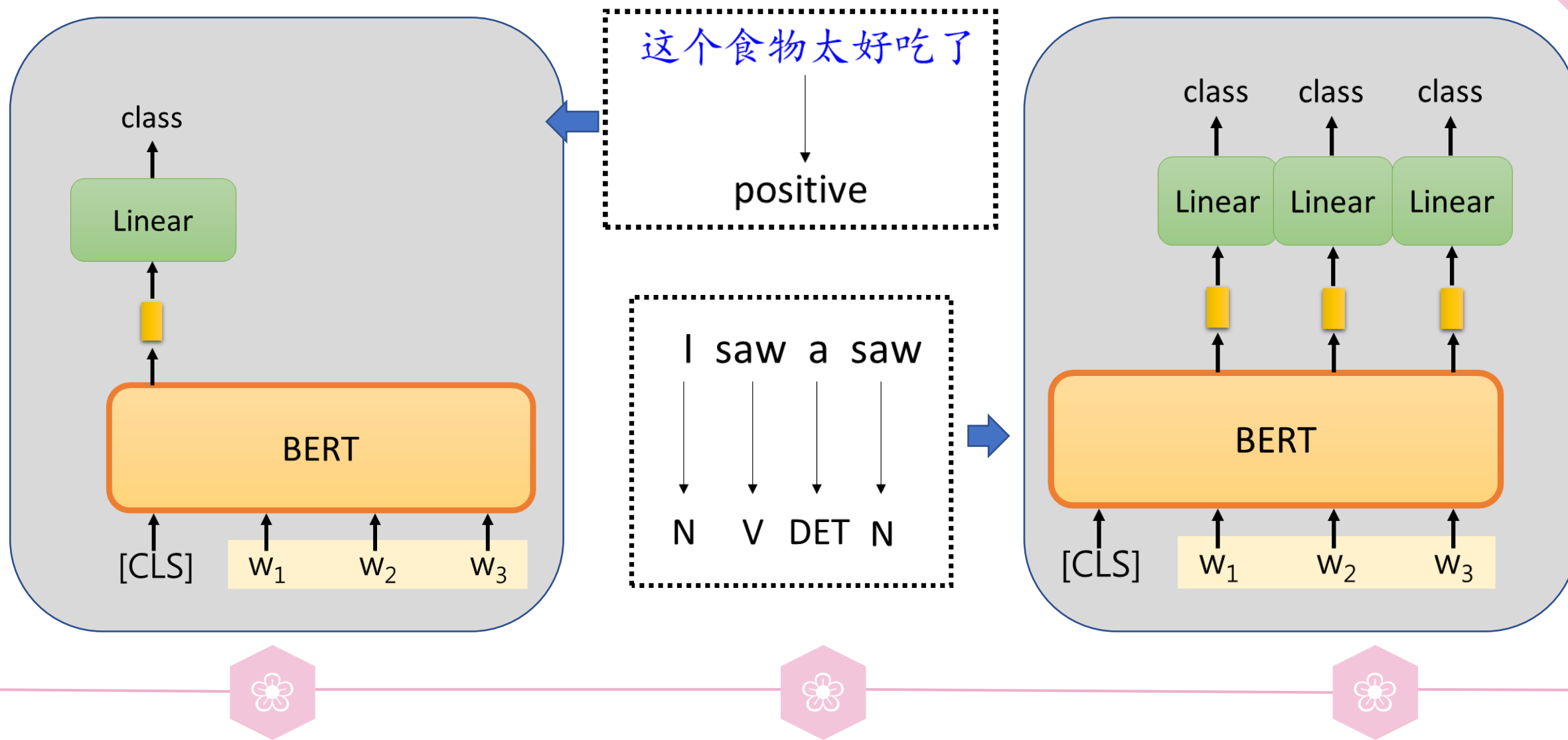
Decoder



Autoregressive (AT)



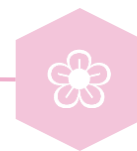
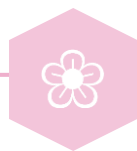
Bert



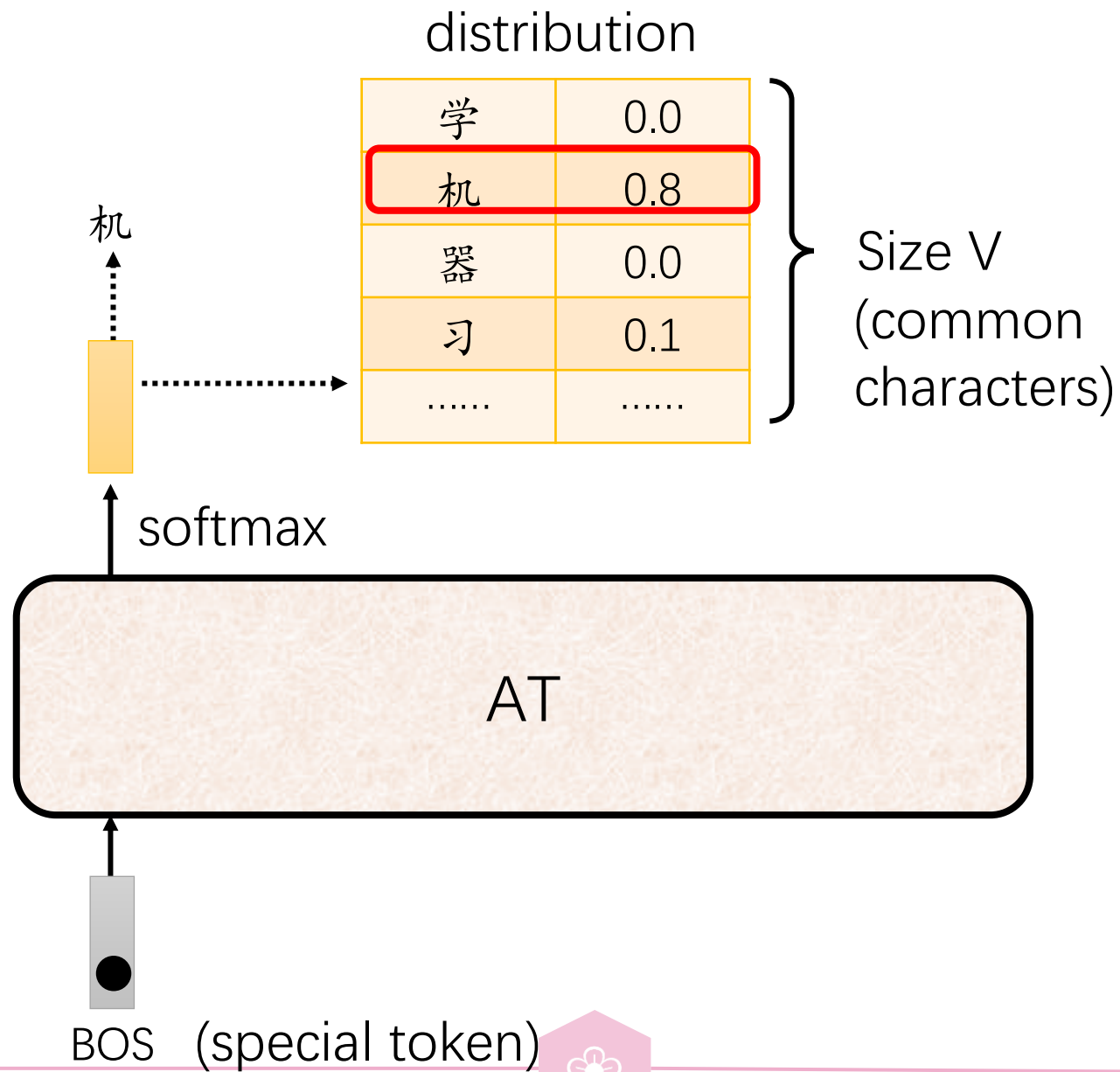


自回归Autoregressive

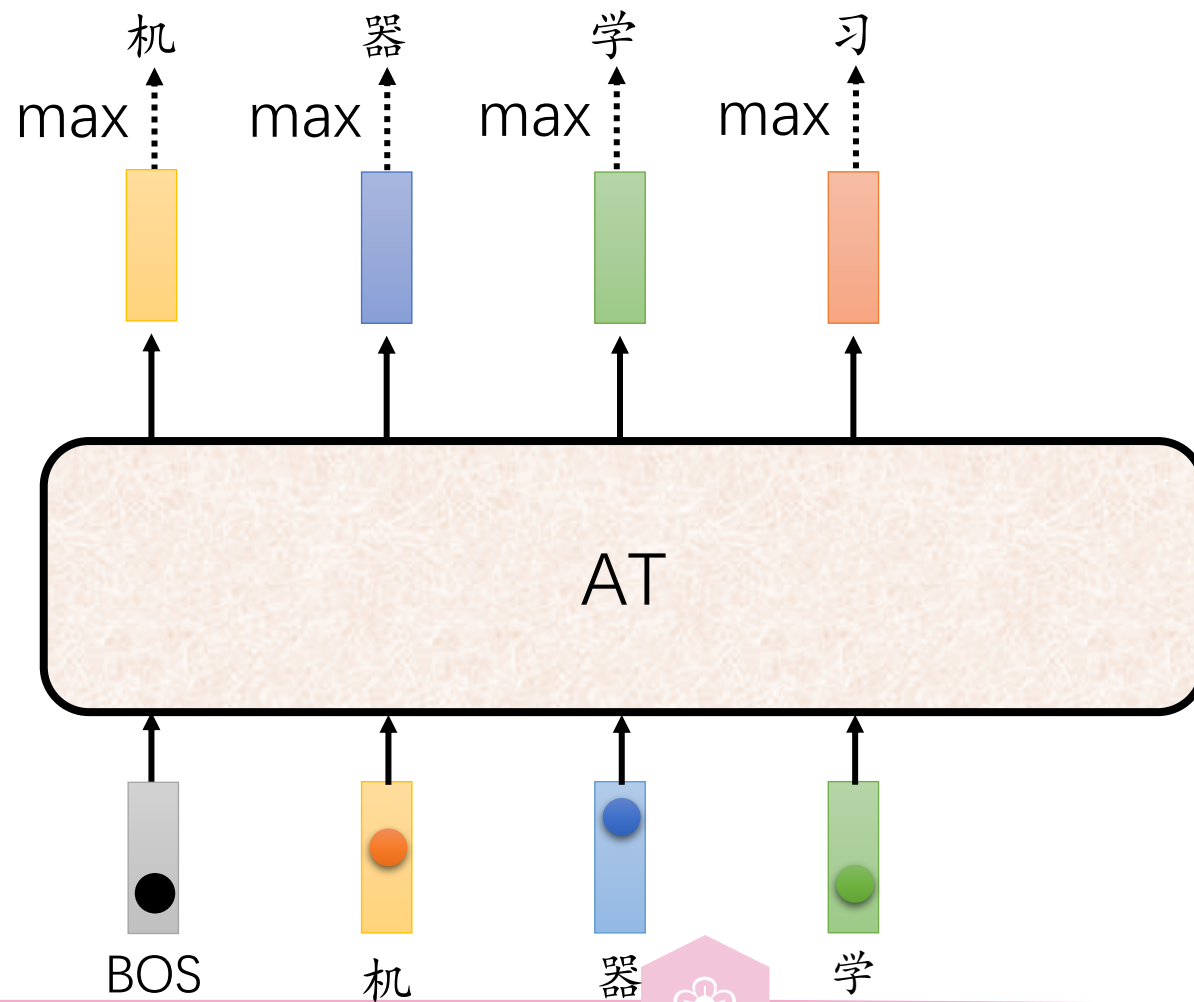
- 将预测对象按照时间顺序排列起来，构成一个所谓的时间序列
 - 从过去预测未来，依次生成每一个输出值
- 目标函数： $-\log \sum_{t=1}^T p(y_t | y_{<t}, x)$



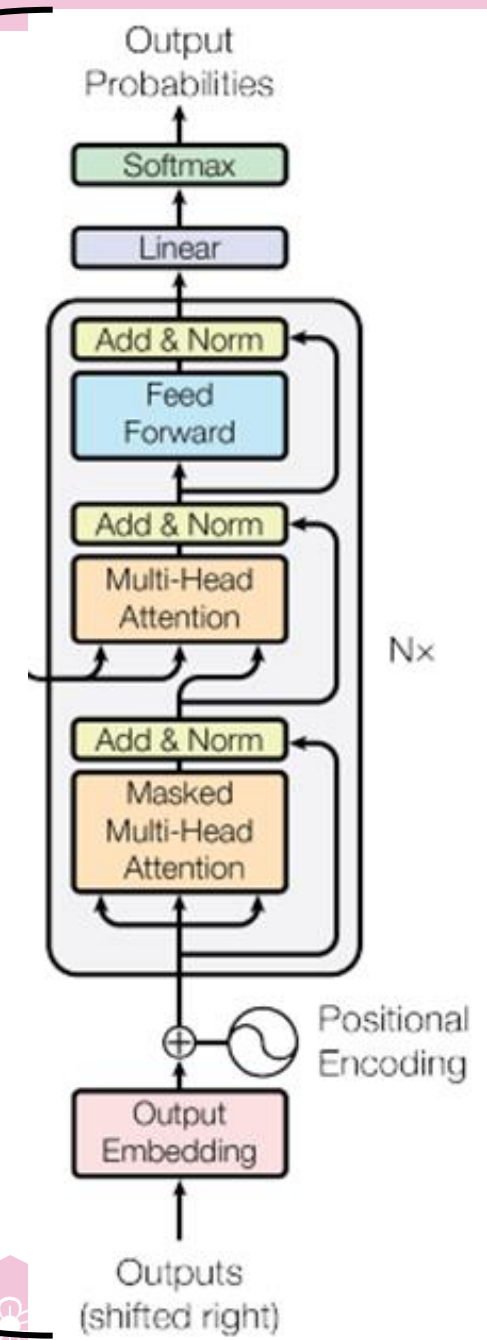
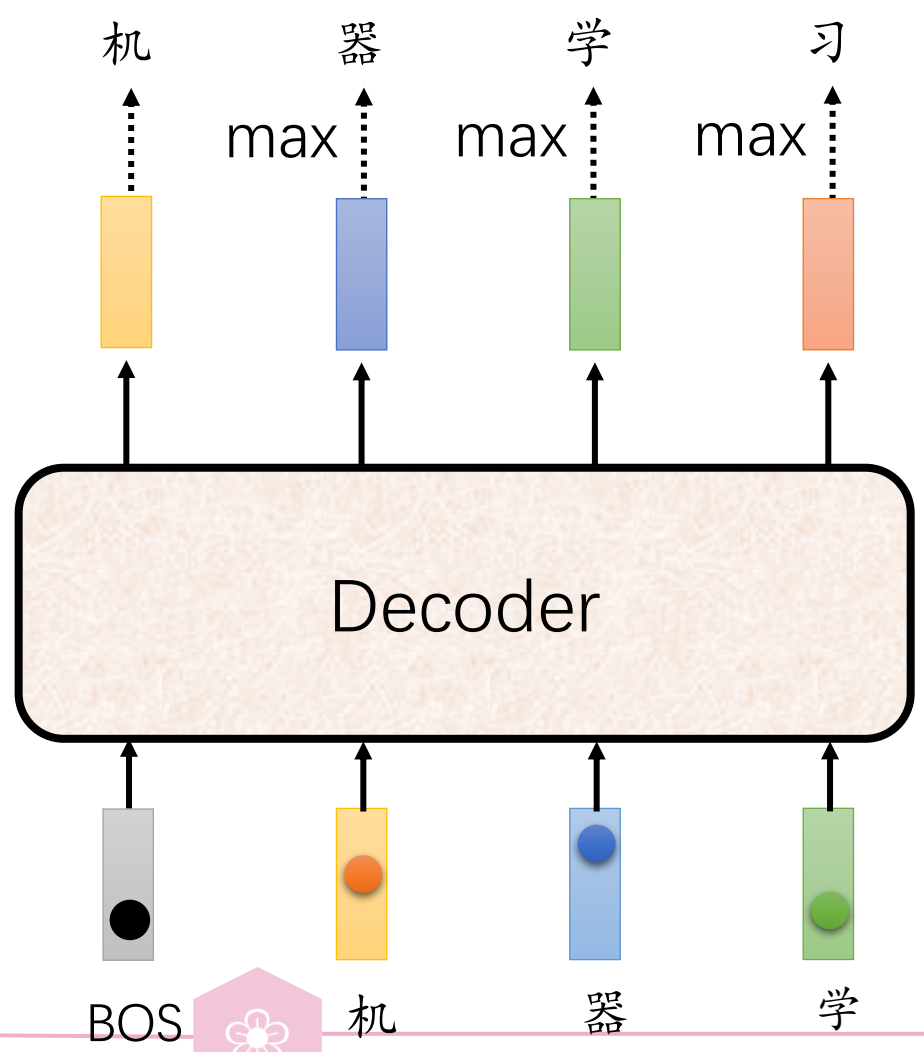
AT



AT



Decoder

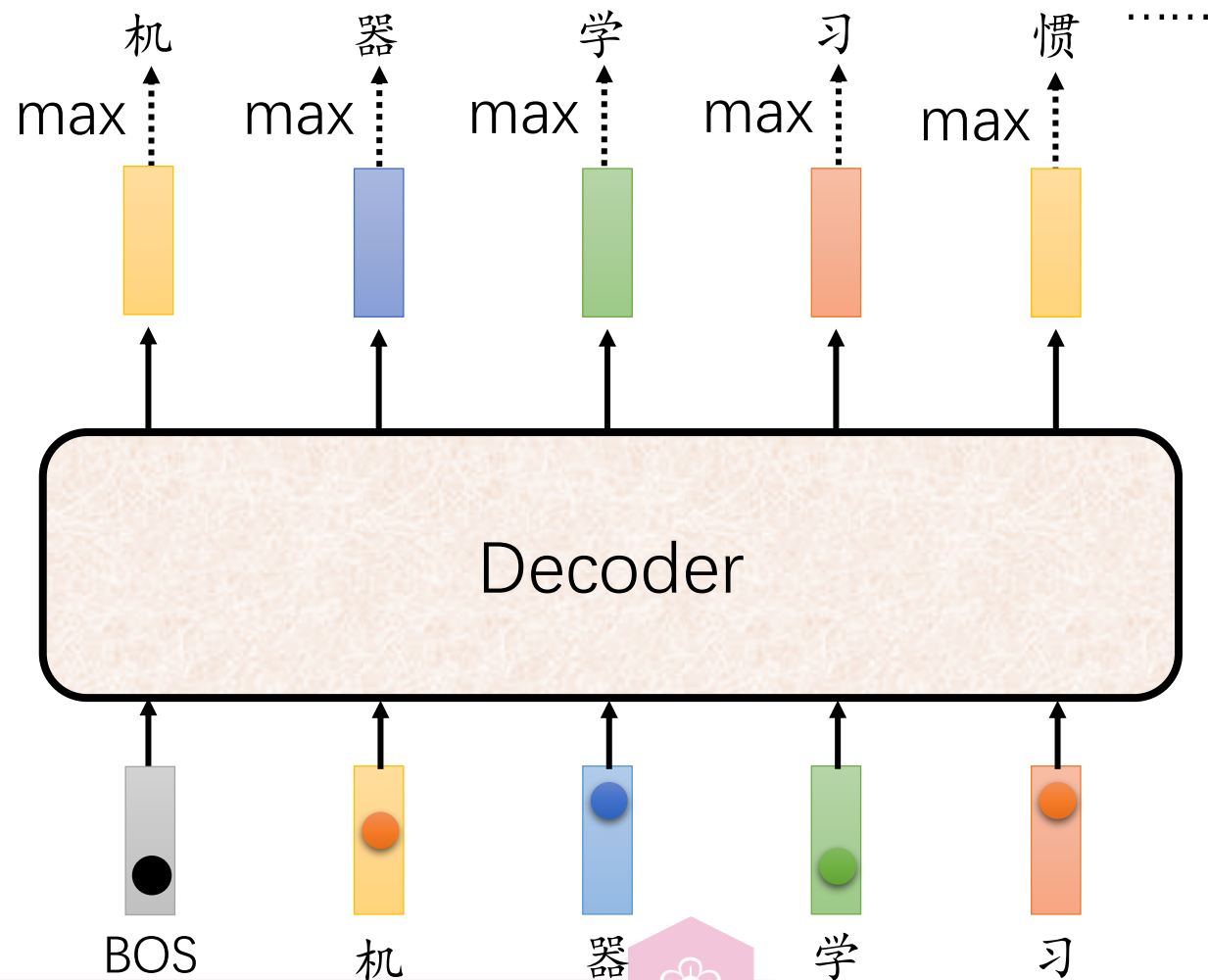


Decoder

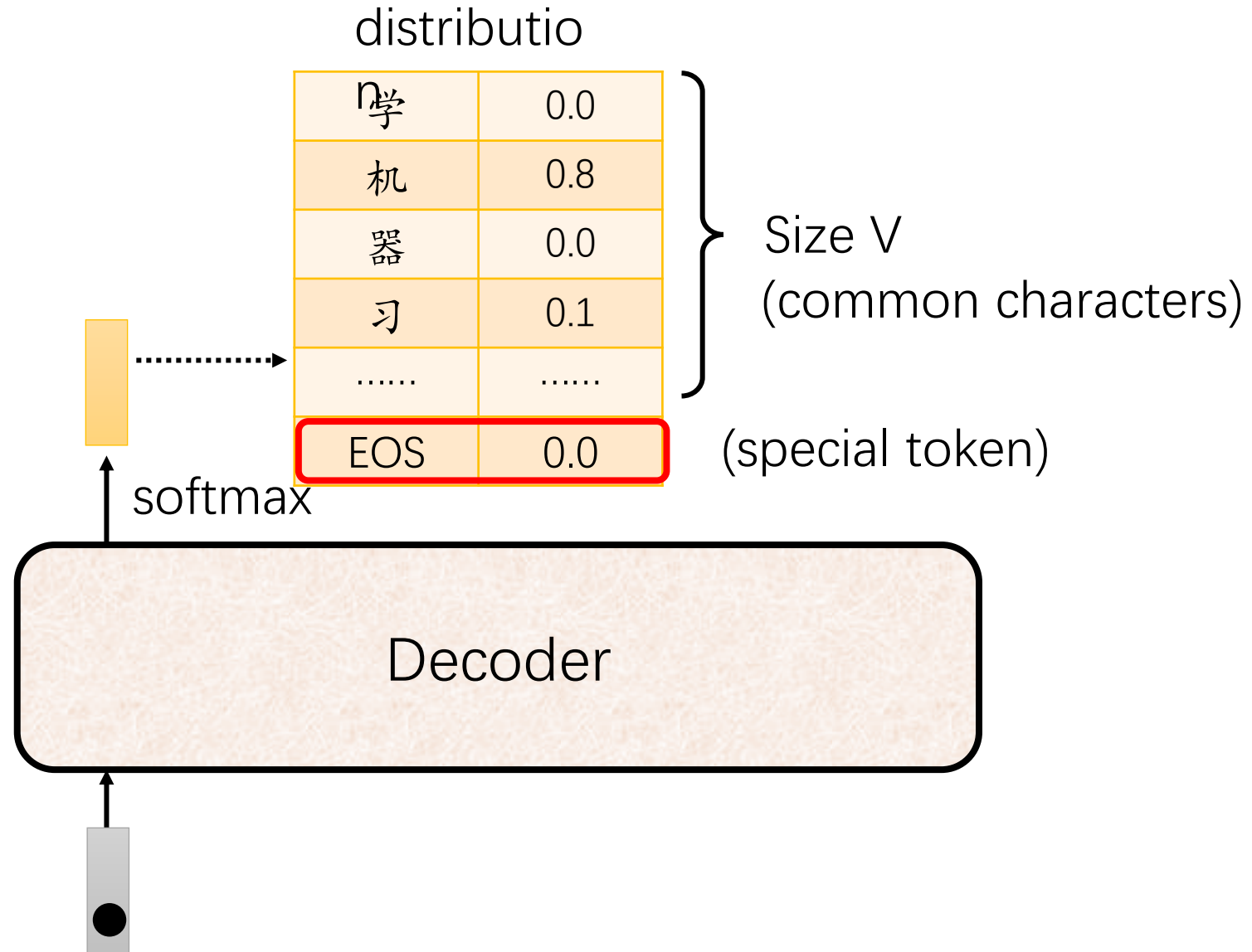


由于真实的输出序列长度未知.

Never stop!

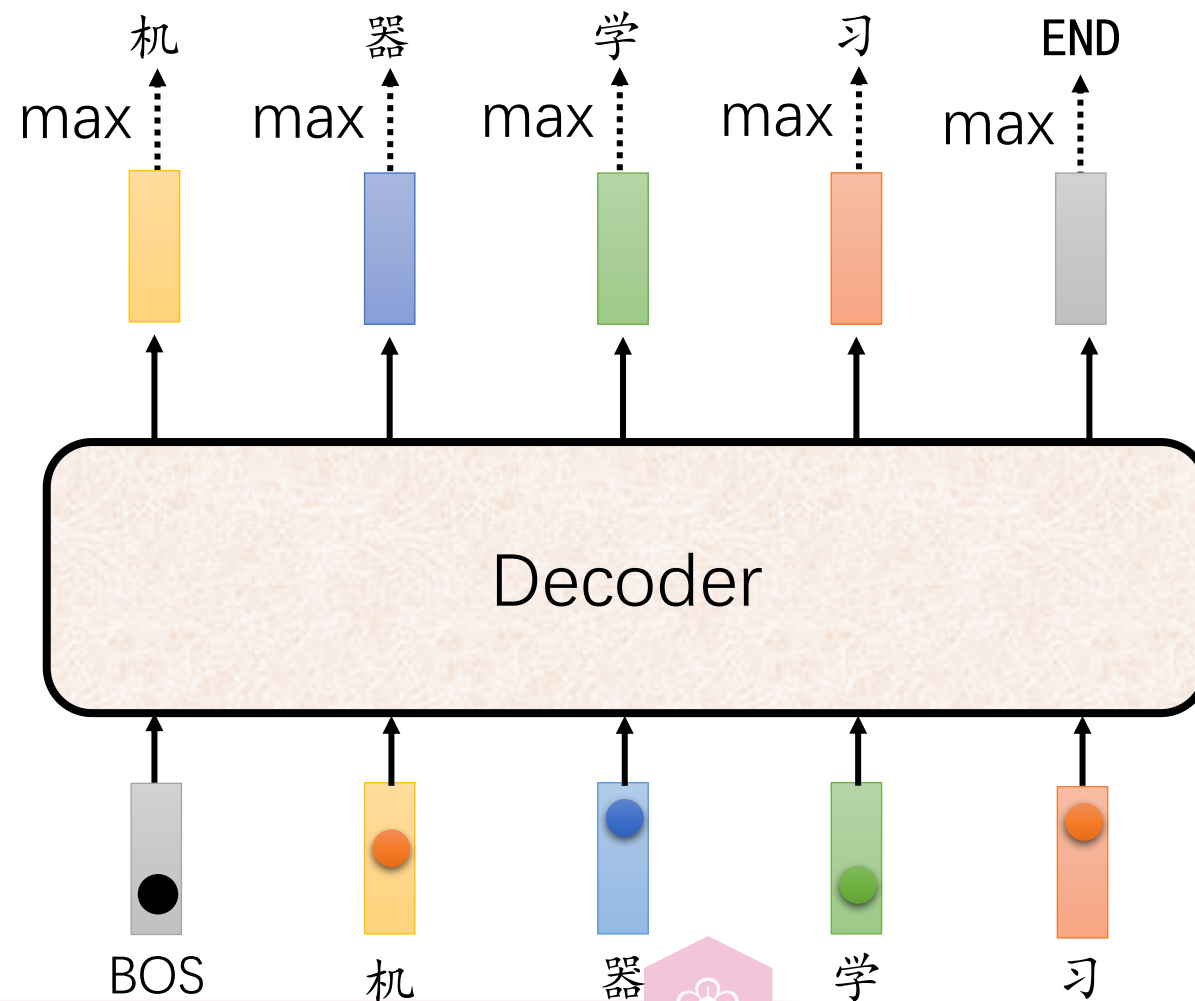


Decoder



Decoder

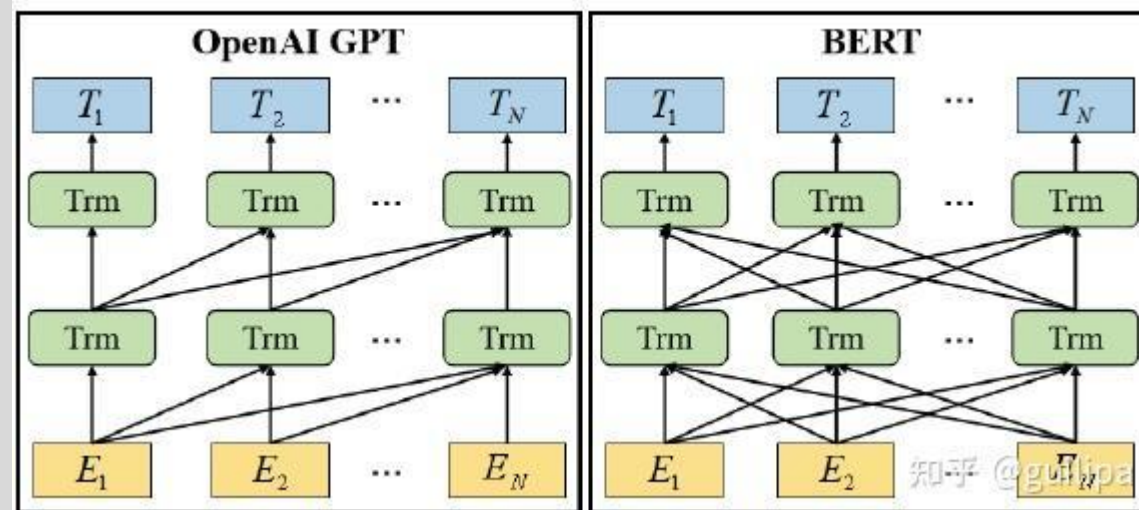
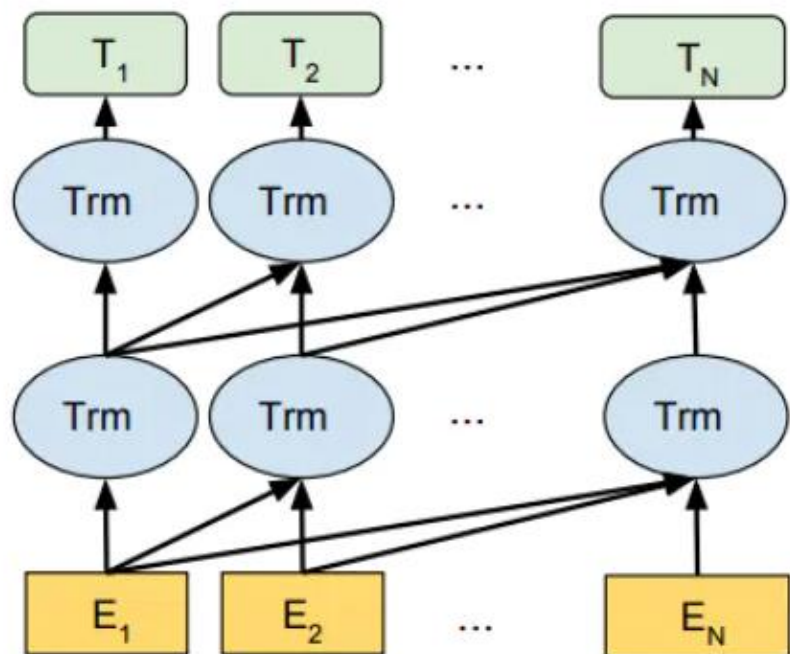
Stop at here!



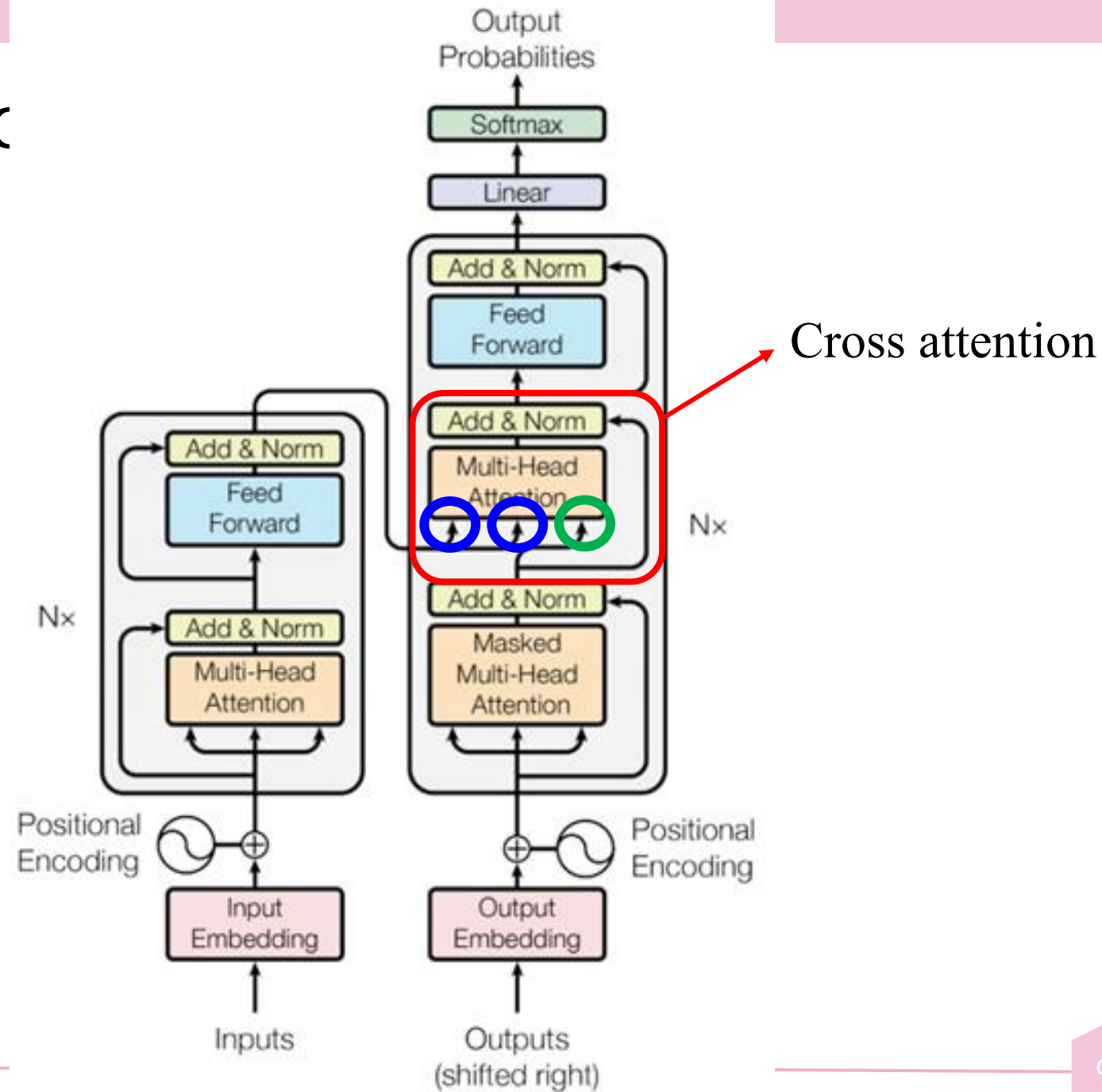
GPT



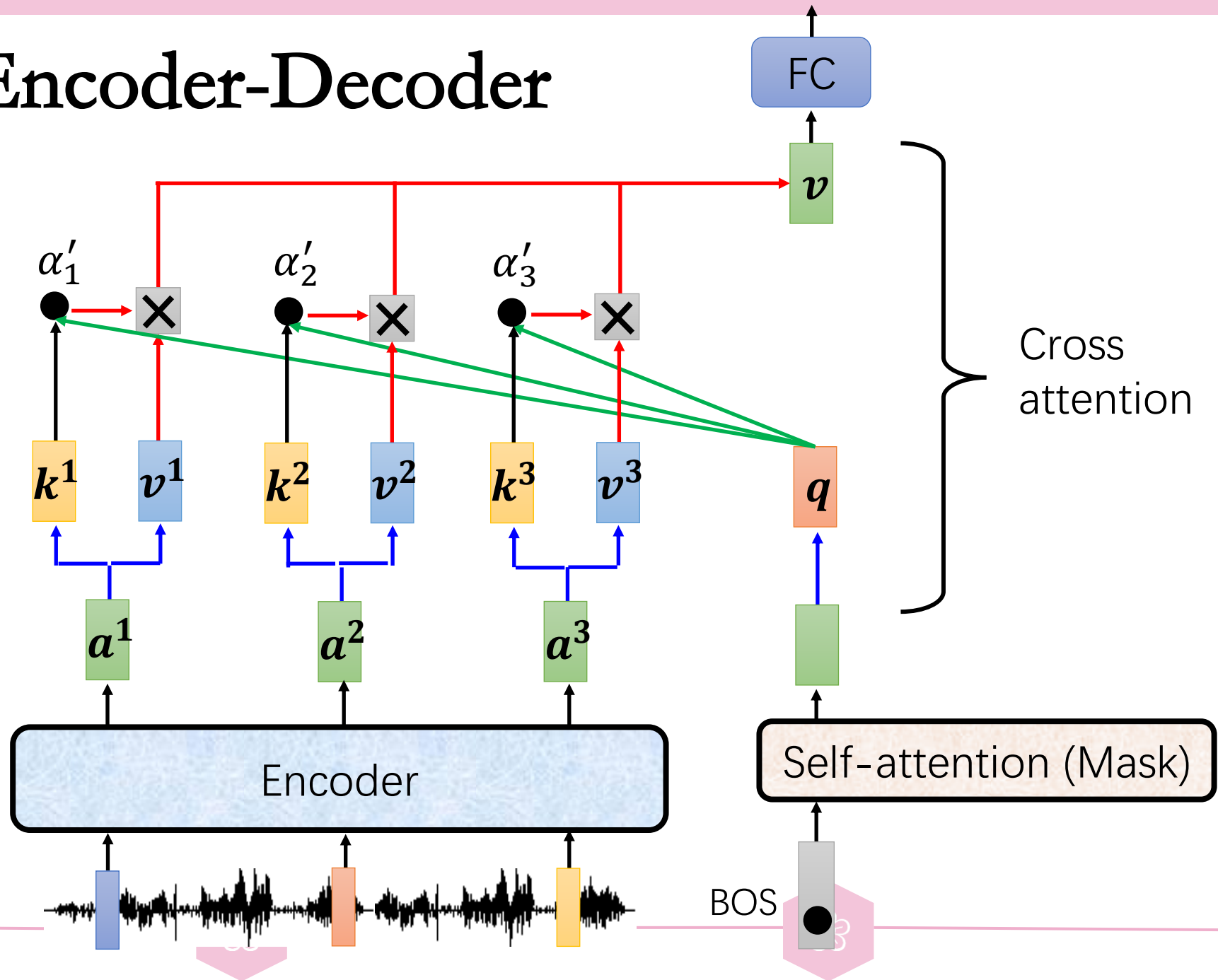
OpenAI GPT



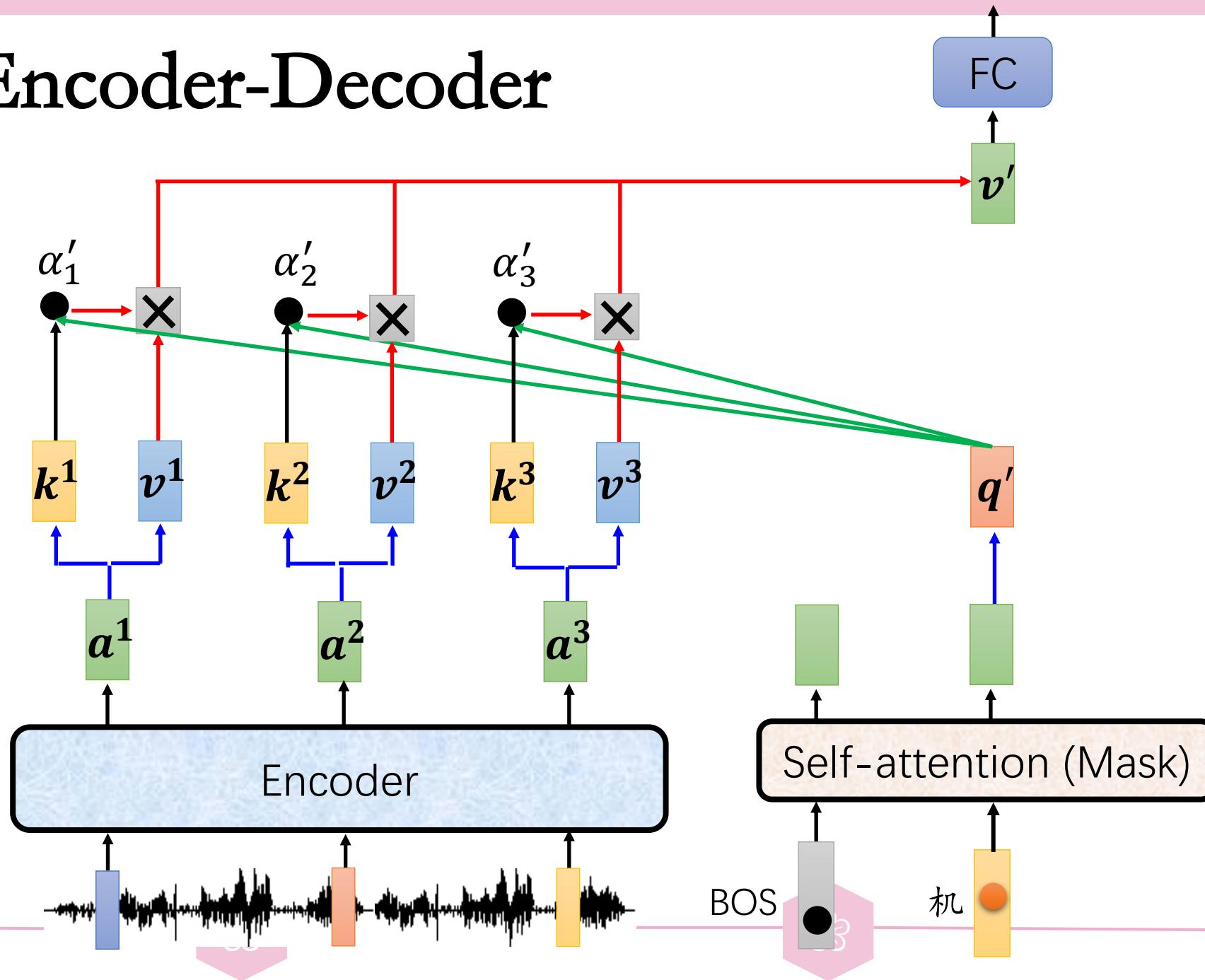
Encoder-Decoder



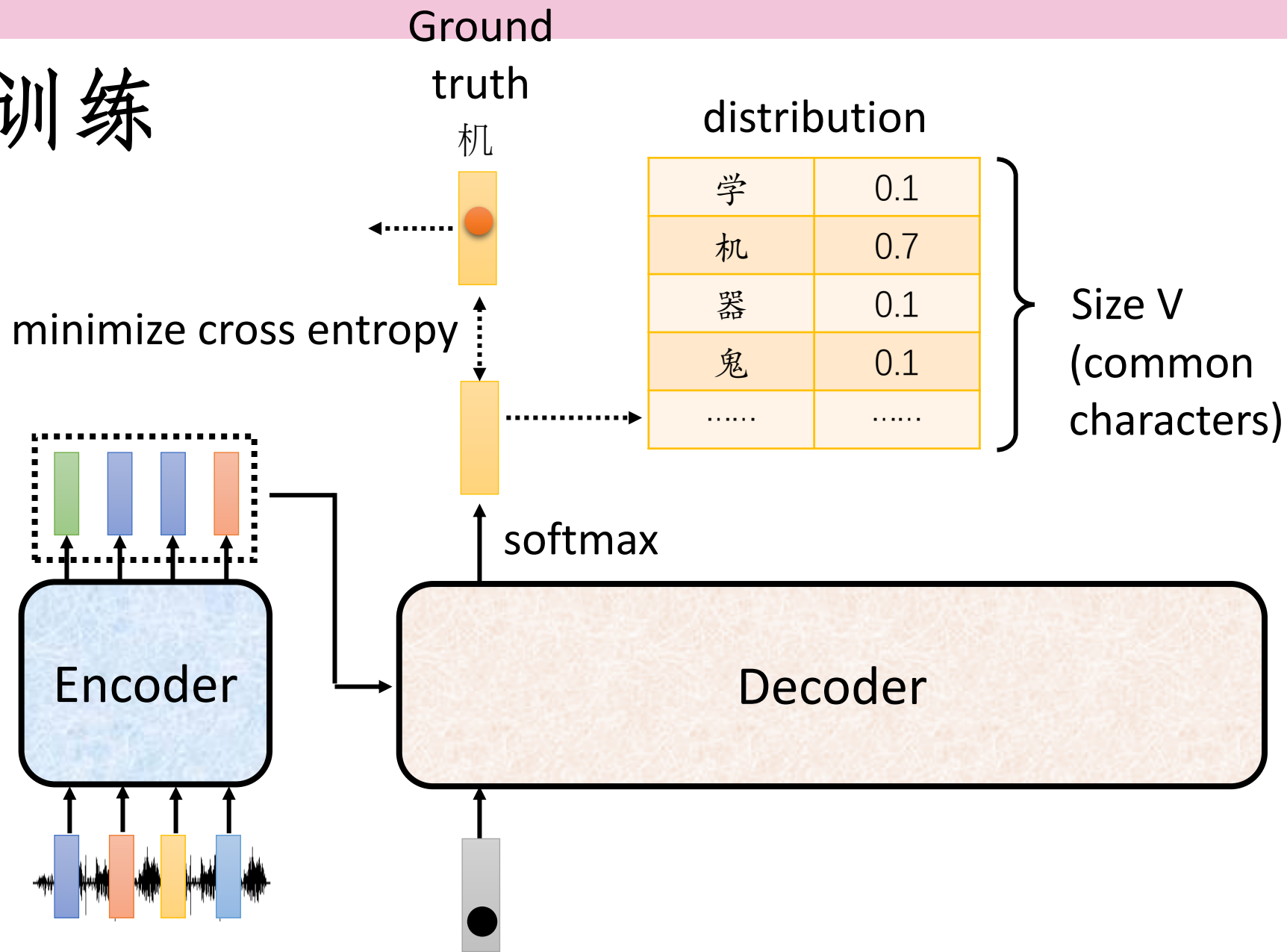
Encoder-Decoder



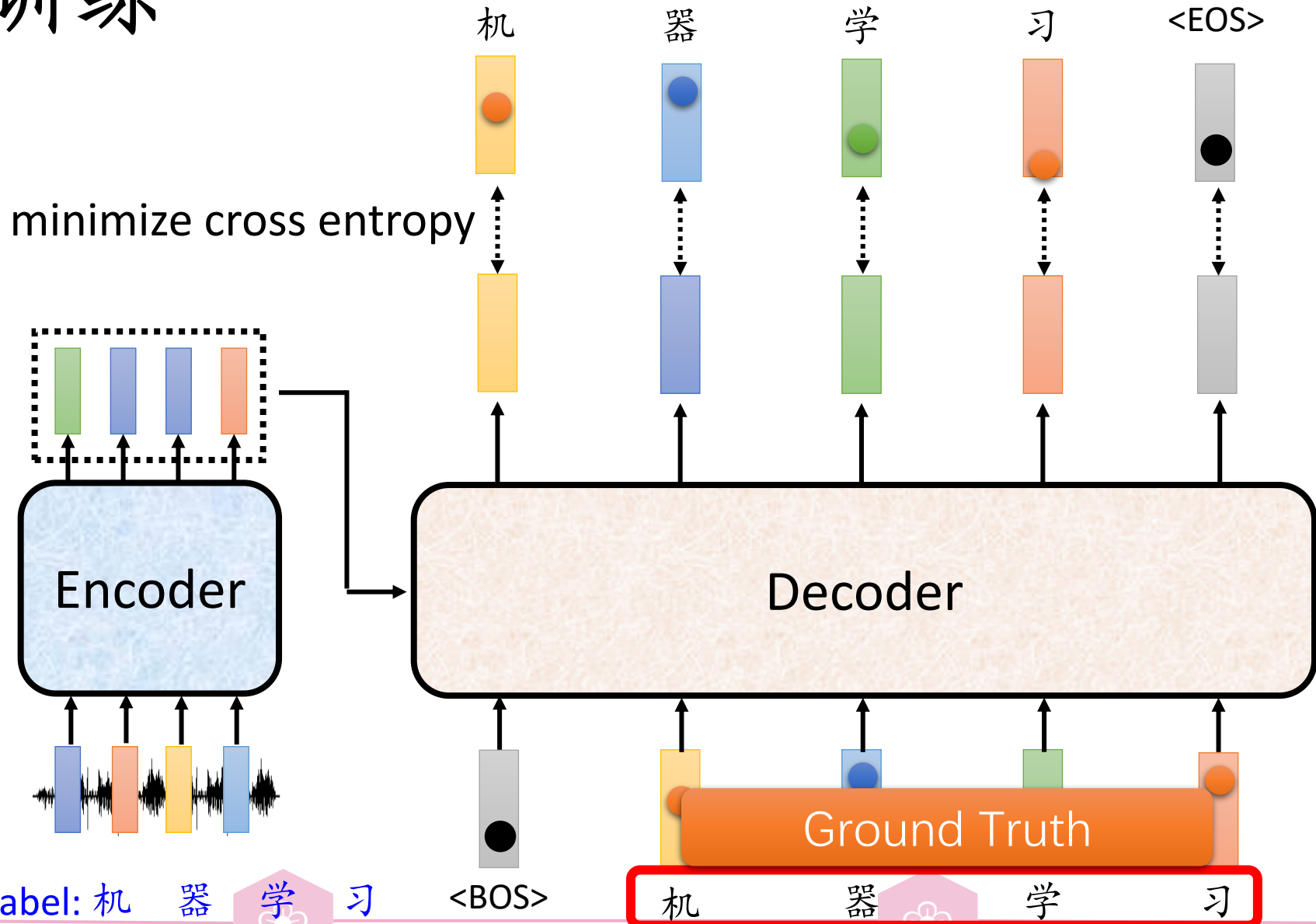
Encoder-Decoder



训练

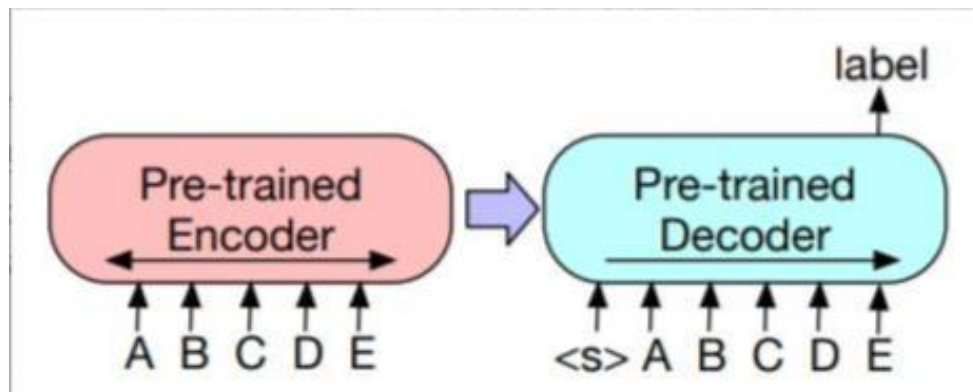


训练

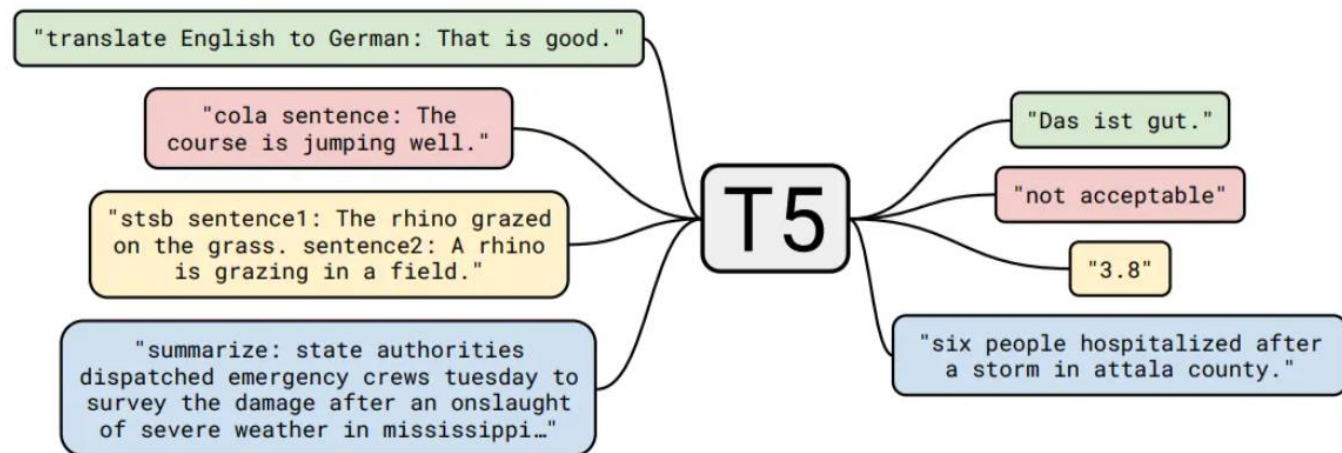


Encoder-decoder架构

BART



T5



UniLM

