



学术分享

Transformer入门

北京理工大学宇航学院

李嘉鑫

◆ Transformer的前世今生

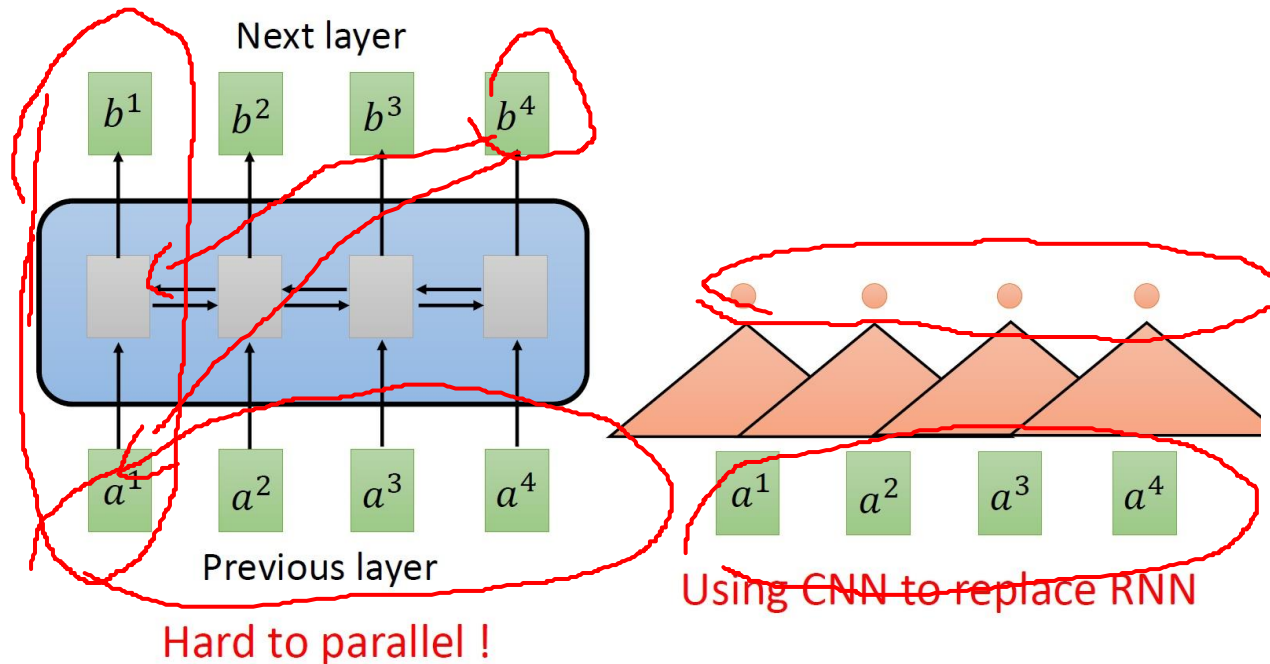
Transformer是Google团队在2017年提出的一种NLP经典模型。

1. 采用**Self-Attention**机制
2. **不采用**RNN的顺序结构，能够并行化训练，且拥有全局信息

◆ RNN与CNN

网络结构

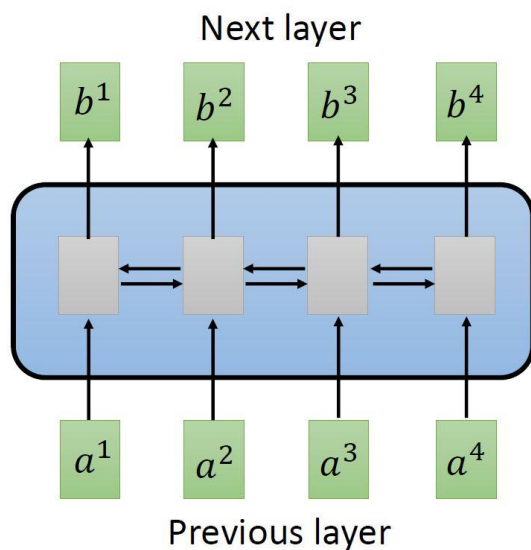
Sequence



◆ RNN与CNN

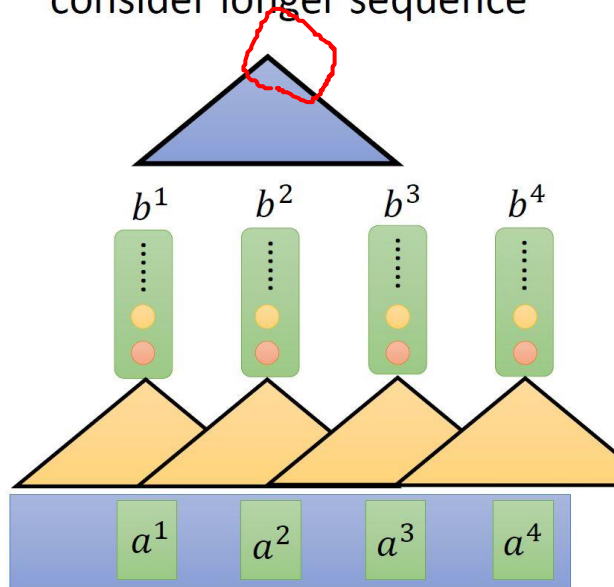
网络结构

Sequence



Hard to parallel

Filters in higher layer can consider longer sequence

Using CNN to replace RNN
(CNN can parallel)

◆ Self-attention

网络结构

Self-attention<https://arxiv.org/abs/1706.03762> q : query (to match others)

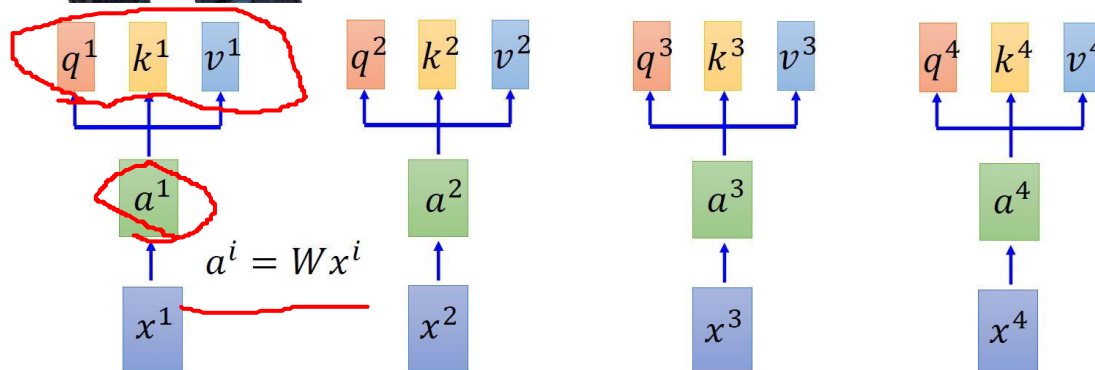
$$q^i = W^q a^i$$

 k : key (to be matched)

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

 v : information to be extracted

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

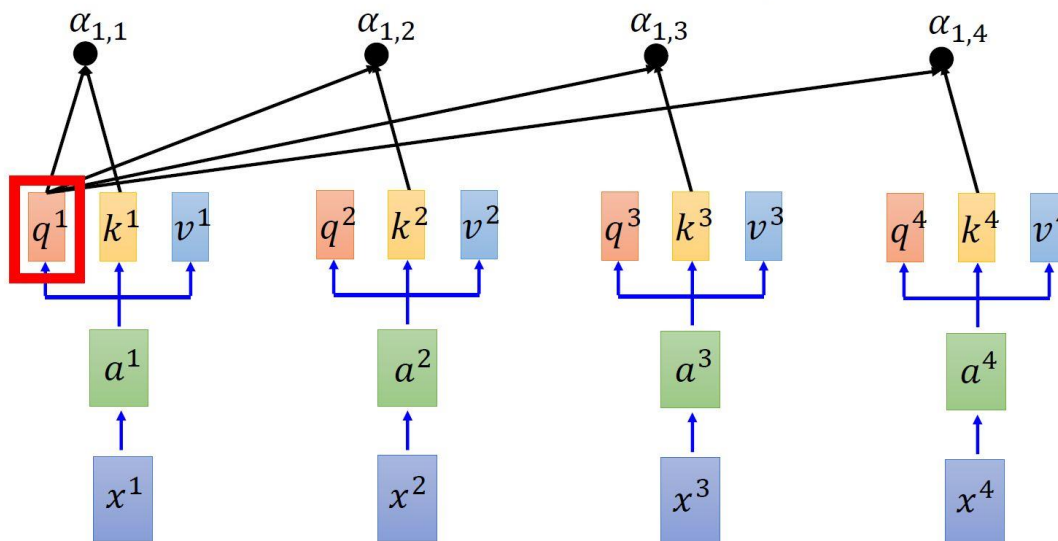


◆ Self-attention

网络结构

Self-attention拿每個 query q 去對每個 key k 做 attention d is the dim of q and k

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

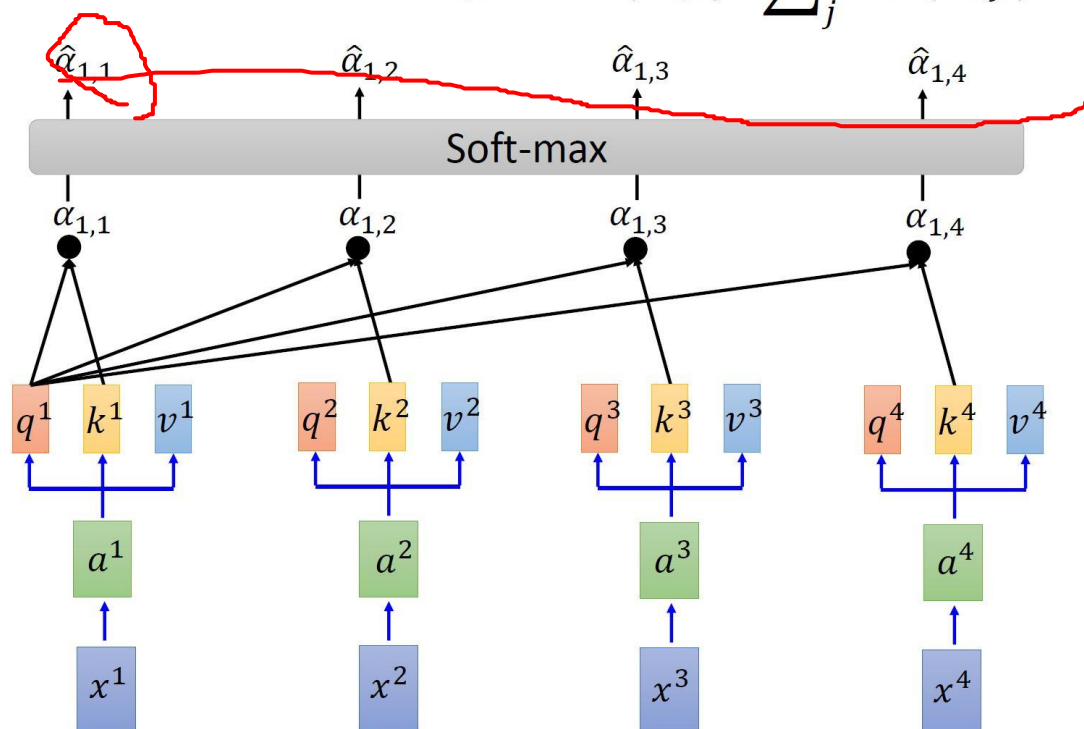


◆ Self-attention

网络结构

Self-attention

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



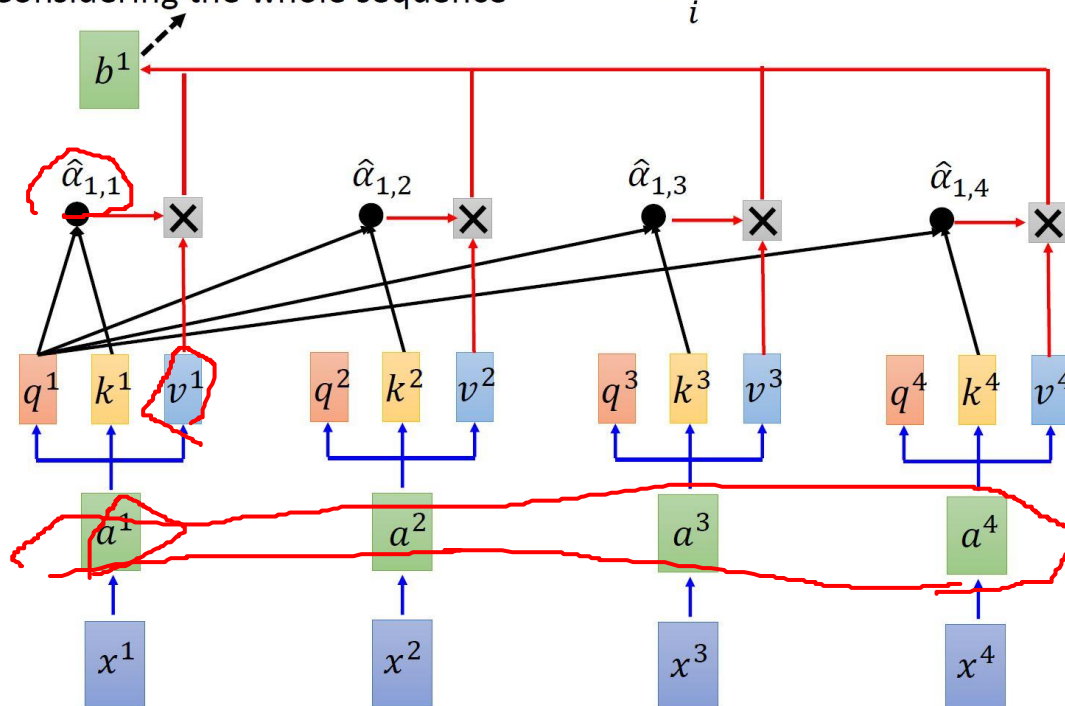
◆ Self-attention

网络结构

Self-attention

Considering the whole sequence

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



◆ Self-attention

网络结构

Self-attention

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

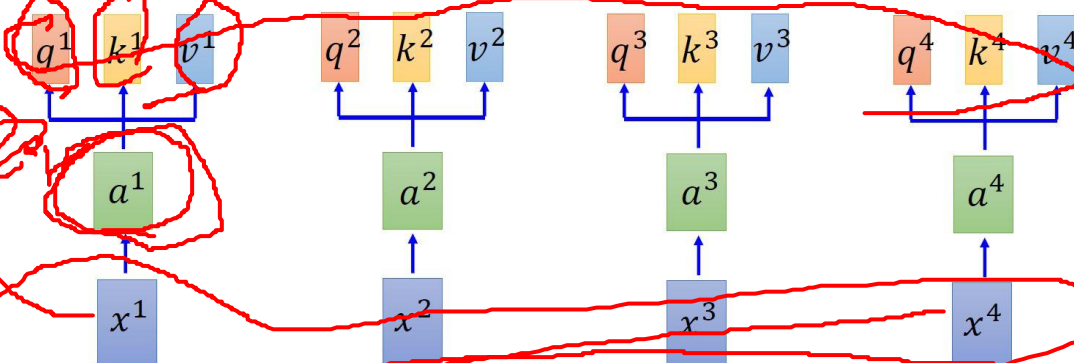
$$k^i = W^k a^i$$

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

$$\begin{matrix} q^1 & q^2 & q^3 & q^4 \\ Q \end{matrix} = \begin{matrix} W^q \\ \text{I} \end{matrix} \begin{matrix} a^1 & a^2 & a^3 & a^4 \\ \text{I} \end{matrix}$$

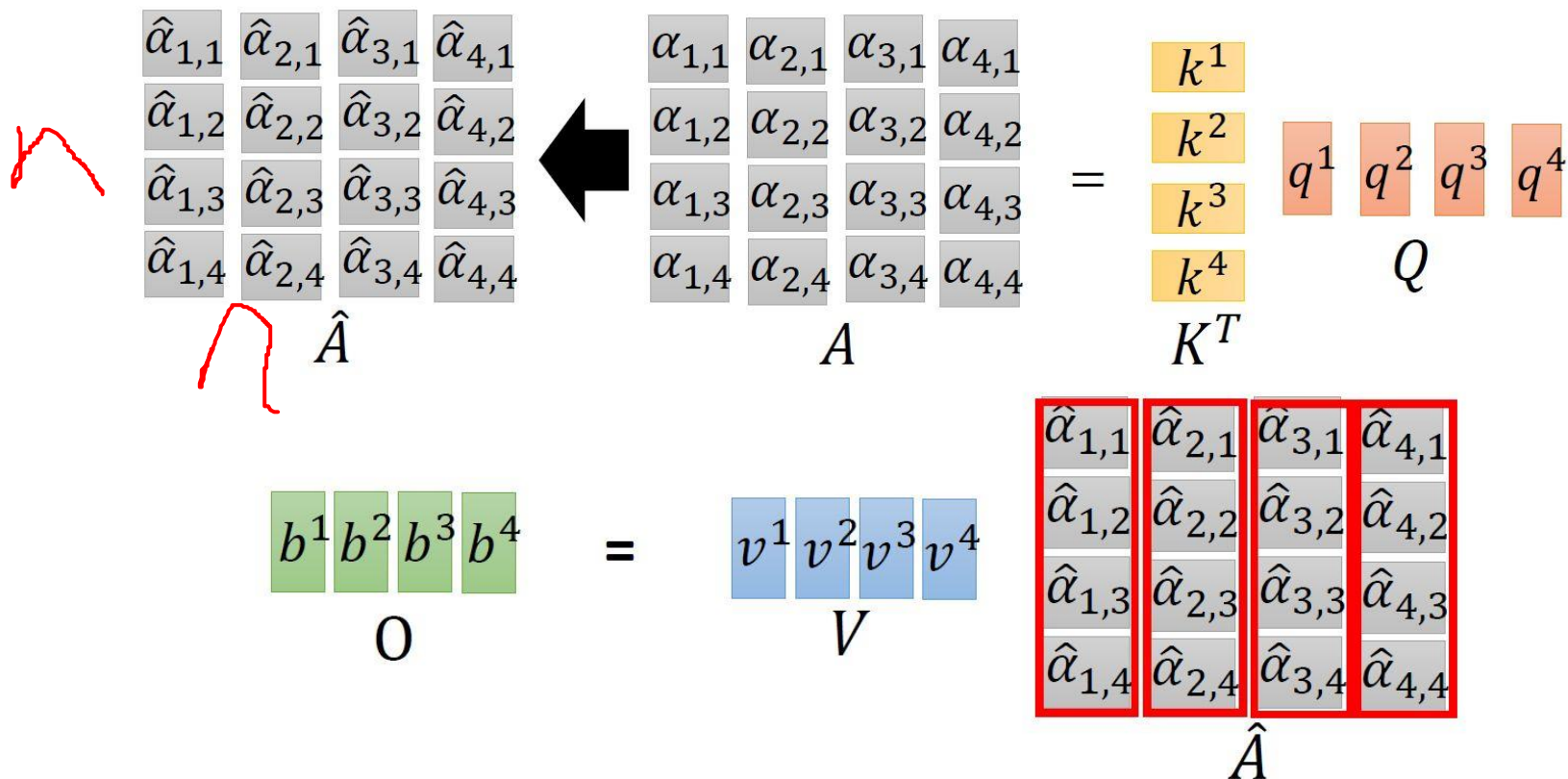
$$\begin{matrix} k^1 & k^2 & k^3 & k^4 \\ K \end{matrix} = \begin{matrix} W^k \\ \text{I} \end{matrix} \begin{matrix} a^1 & a^2 & a^3 & a^4 \\ \text{I} \end{matrix}$$

$$\begin{matrix} v^1 & v^2 & v^3 & v^4 \\ V \end{matrix} = \begin{matrix} W^v \\ \text{I} \end{matrix} \begin{matrix} a^1 & a^2 & a^3 & a^4 \\ \text{I} \end{matrix}$$



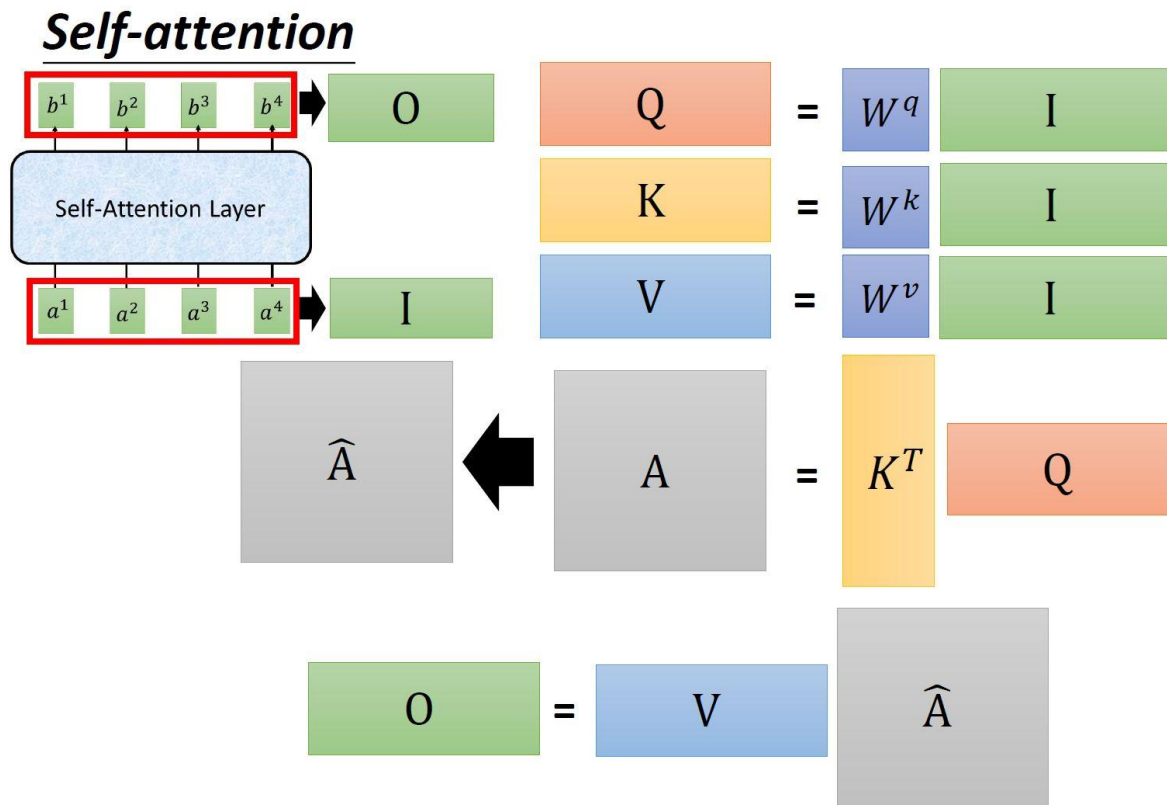
◆ Self-attention

网络结构



◆ Self-attention

网络结构

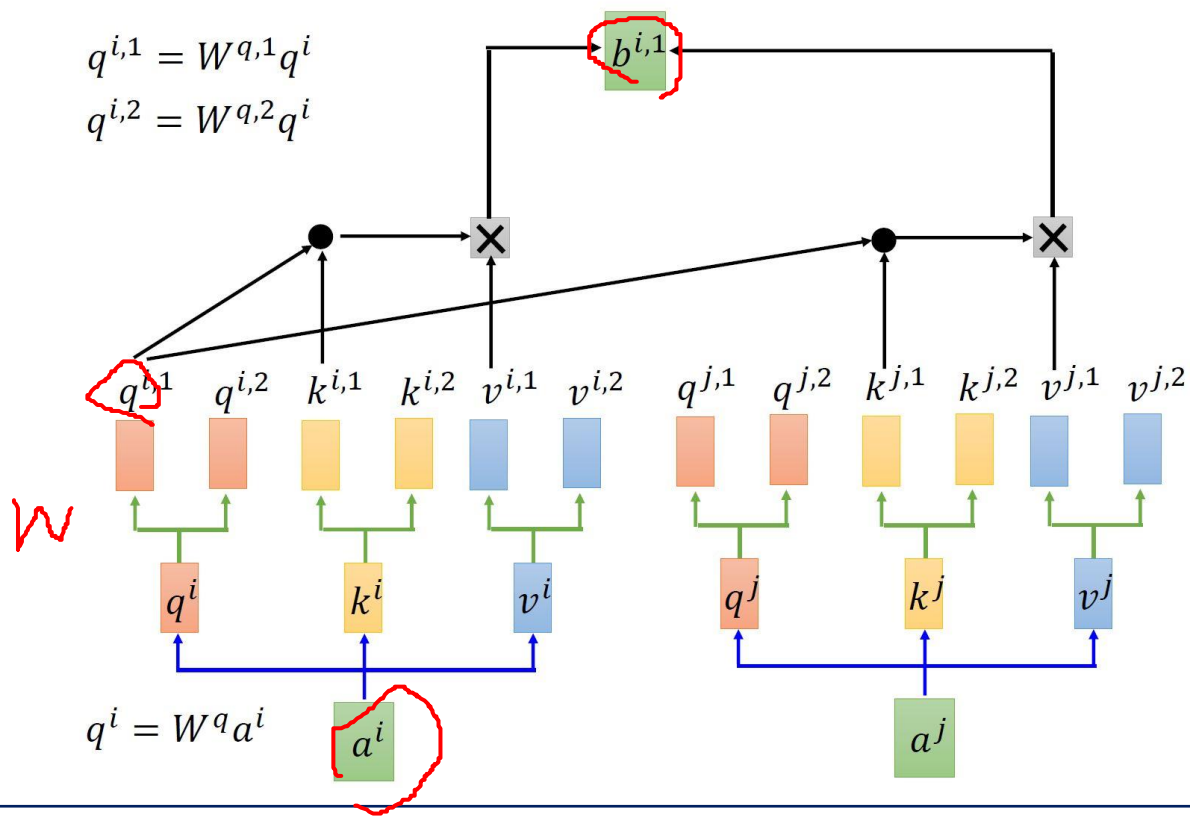


◆ Multi-head Self-attention

网络结构

Multi-head Self-attention

(2 heads as example)



◆ Multi-head Self-attention

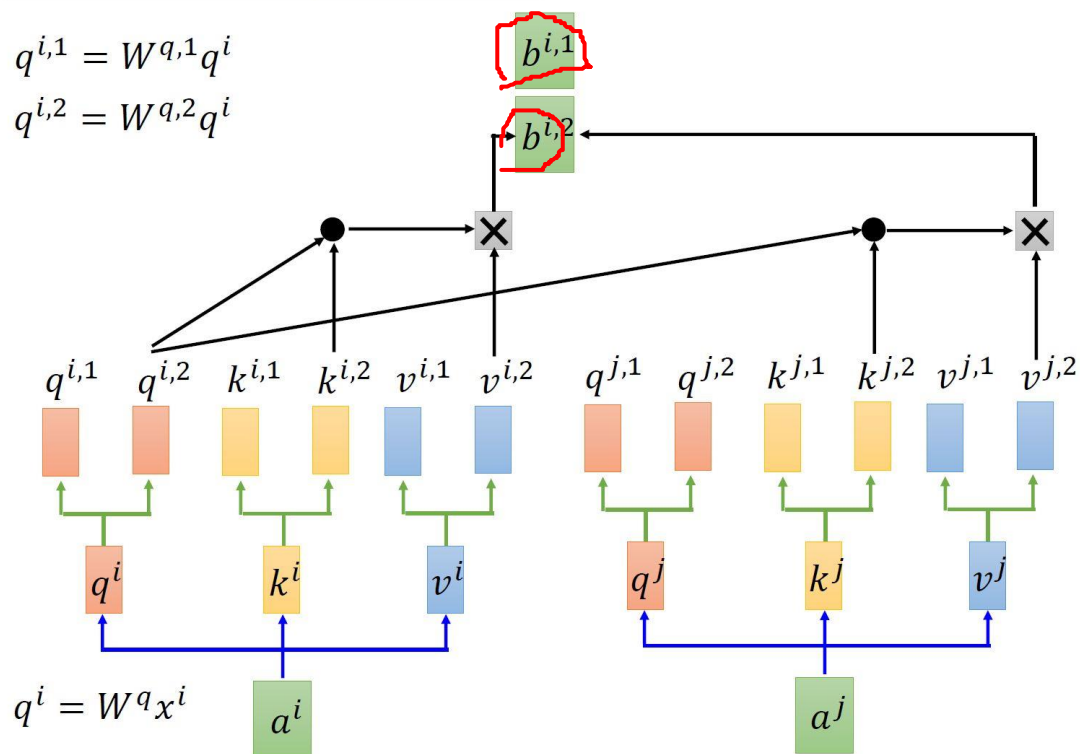
网络结构

Multi-head Self-attention

(2 heads as example)

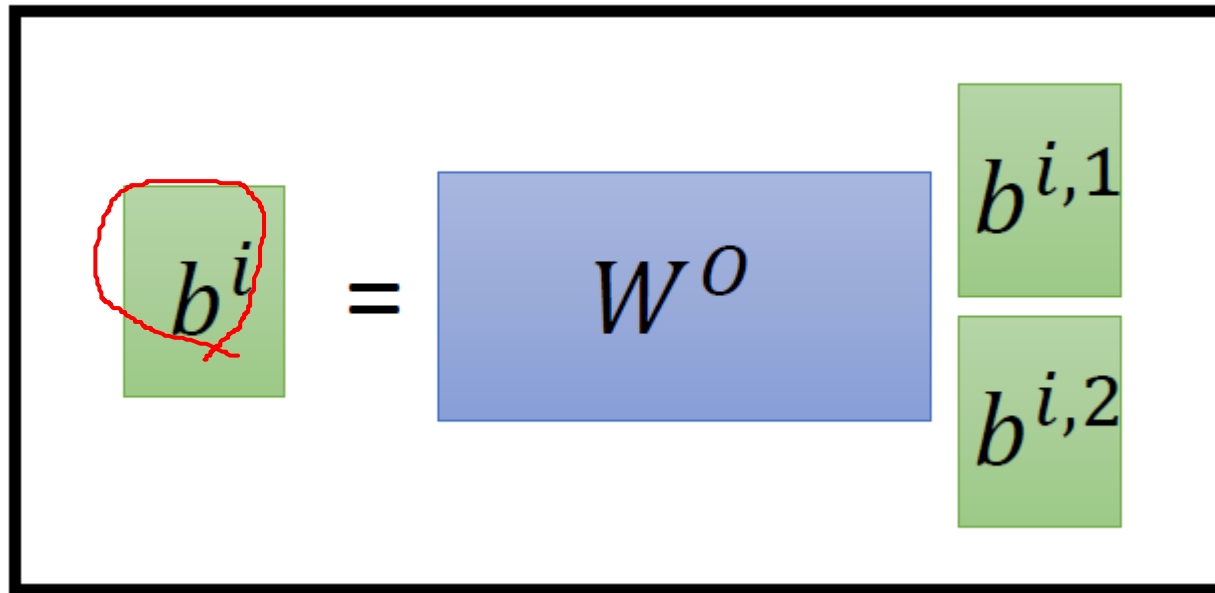
$$q^{i,1} = W^{q,1} q^i$$

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



◆ Multi-head Self-attention

网络结构

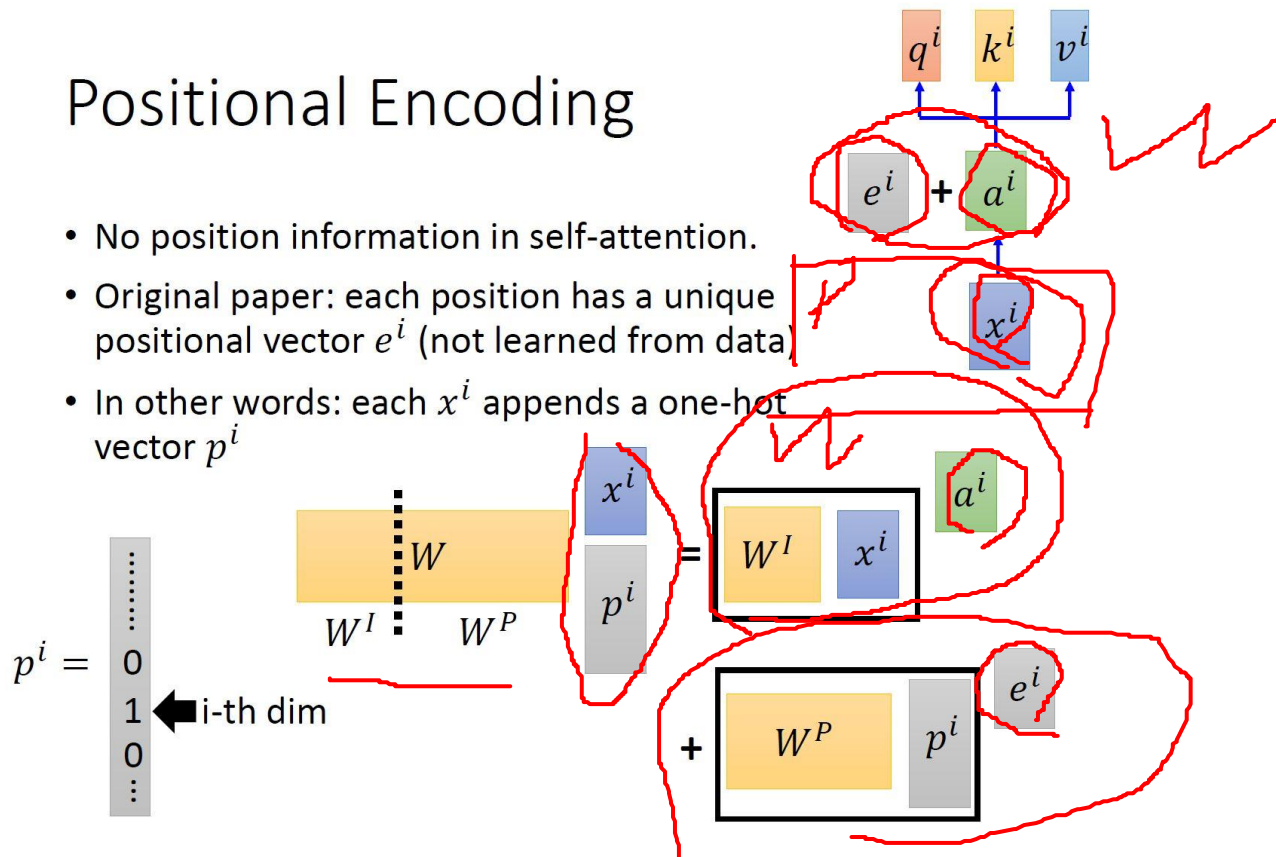


◆ Position Encoding

网络结构

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



◆ Position Encoding

网络结构

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

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

$$\cos(\alpha + \beta) = \cos(\alpha)\cos(\beta) - \sin(\alpha)\sin(\beta)$$

$$\sin(\alpha + \beta) = \sin(\alpha)\cos(\beta) + \cos(\alpha)\sin(\beta)$$

◆ Position Encoding

网络结构

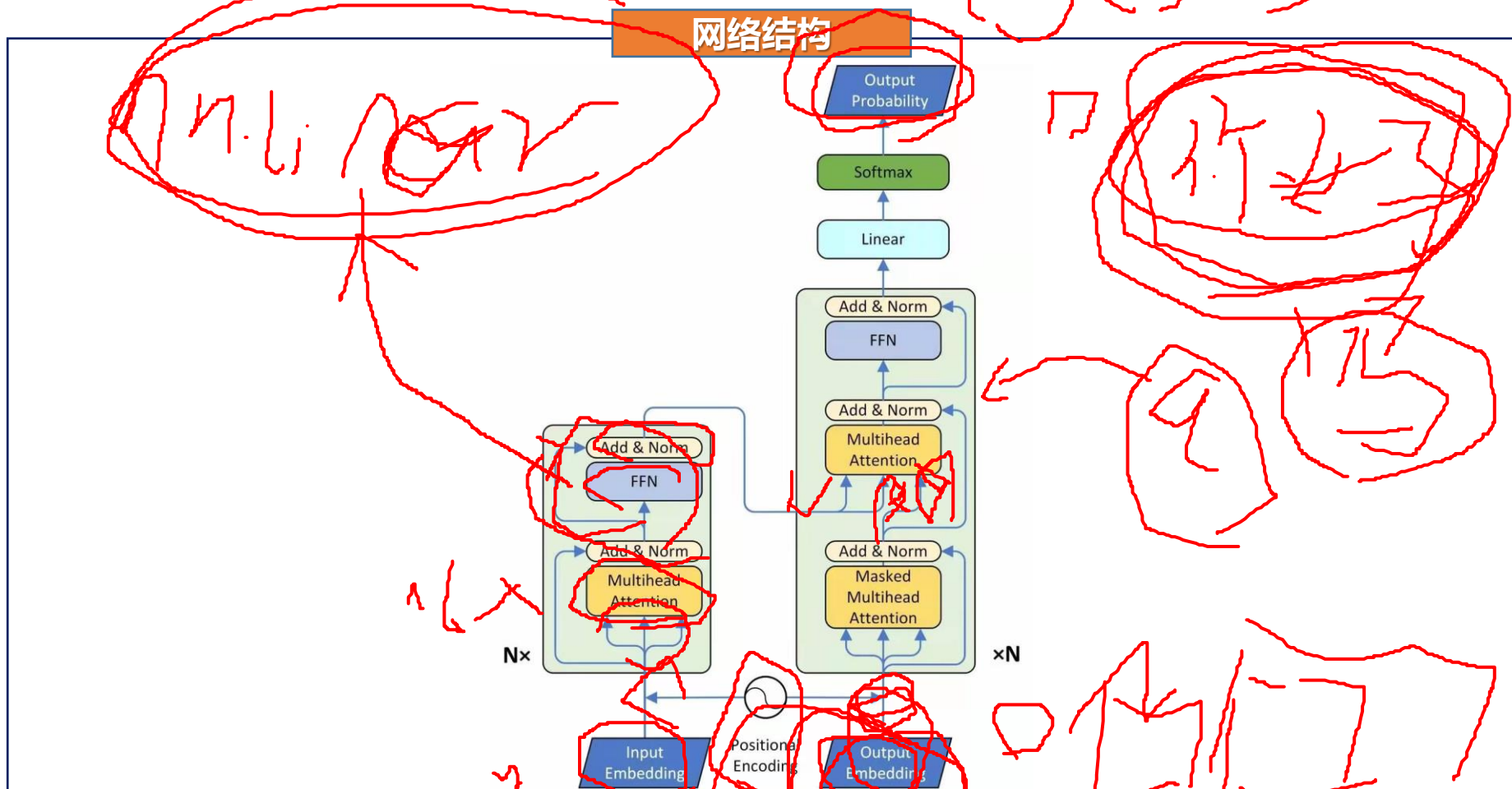
$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

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

$$\cos(\alpha + \beta) = \cos(\alpha)\cos(\beta) - \sin(\alpha)\sin(\beta)$$

$$\sin(\alpha + \beta) = \sin(\alpha)\cos(\beta) + \cos(\alpha)\sin(\beta)$$

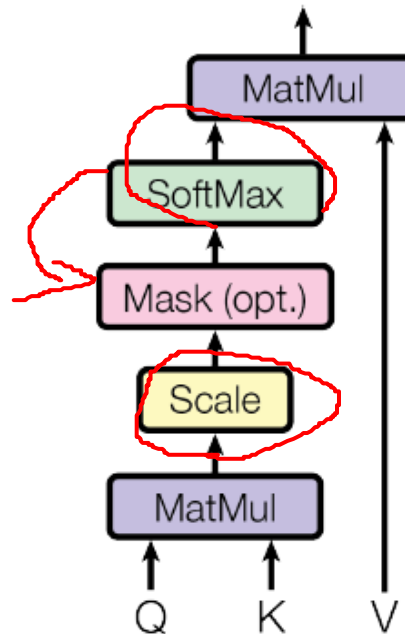
Transformer



◆ Transformer

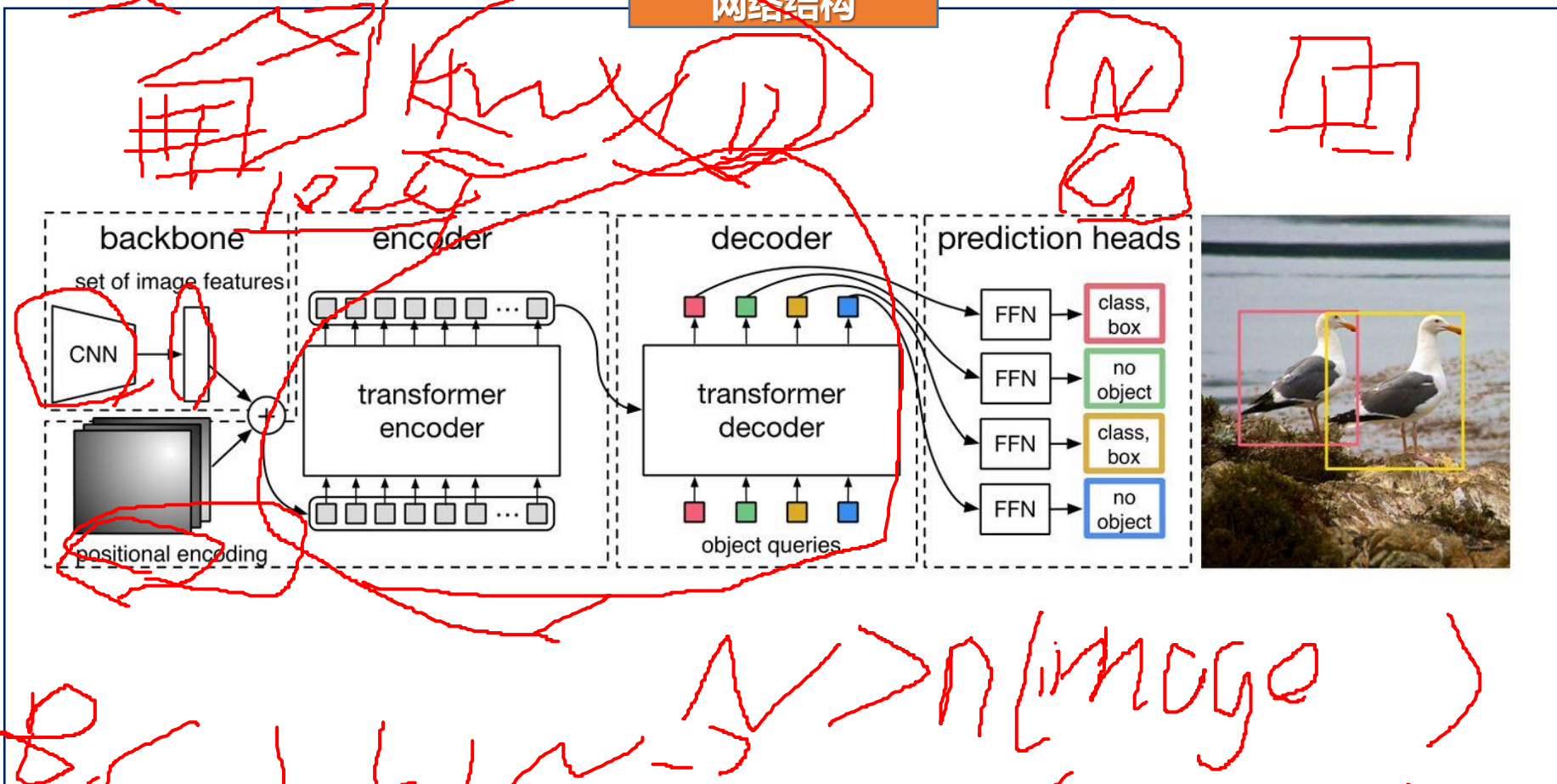
网络结构

Scaled Dot-Product Attention



DETR

网络结构



◆ DETR

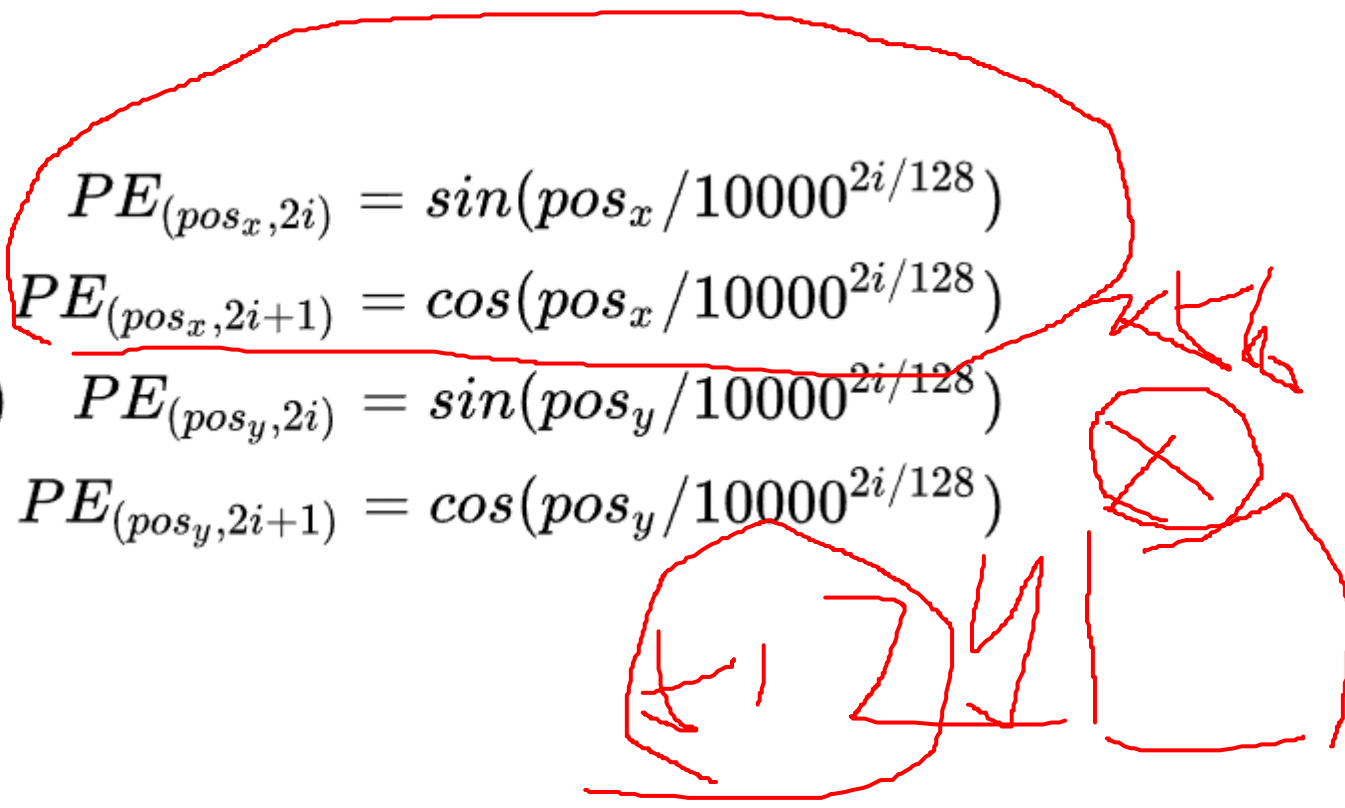
网络结构

a) $PE_{(pos_x, 2i)} = \sin(pos_x / 10000^{2i/128})$

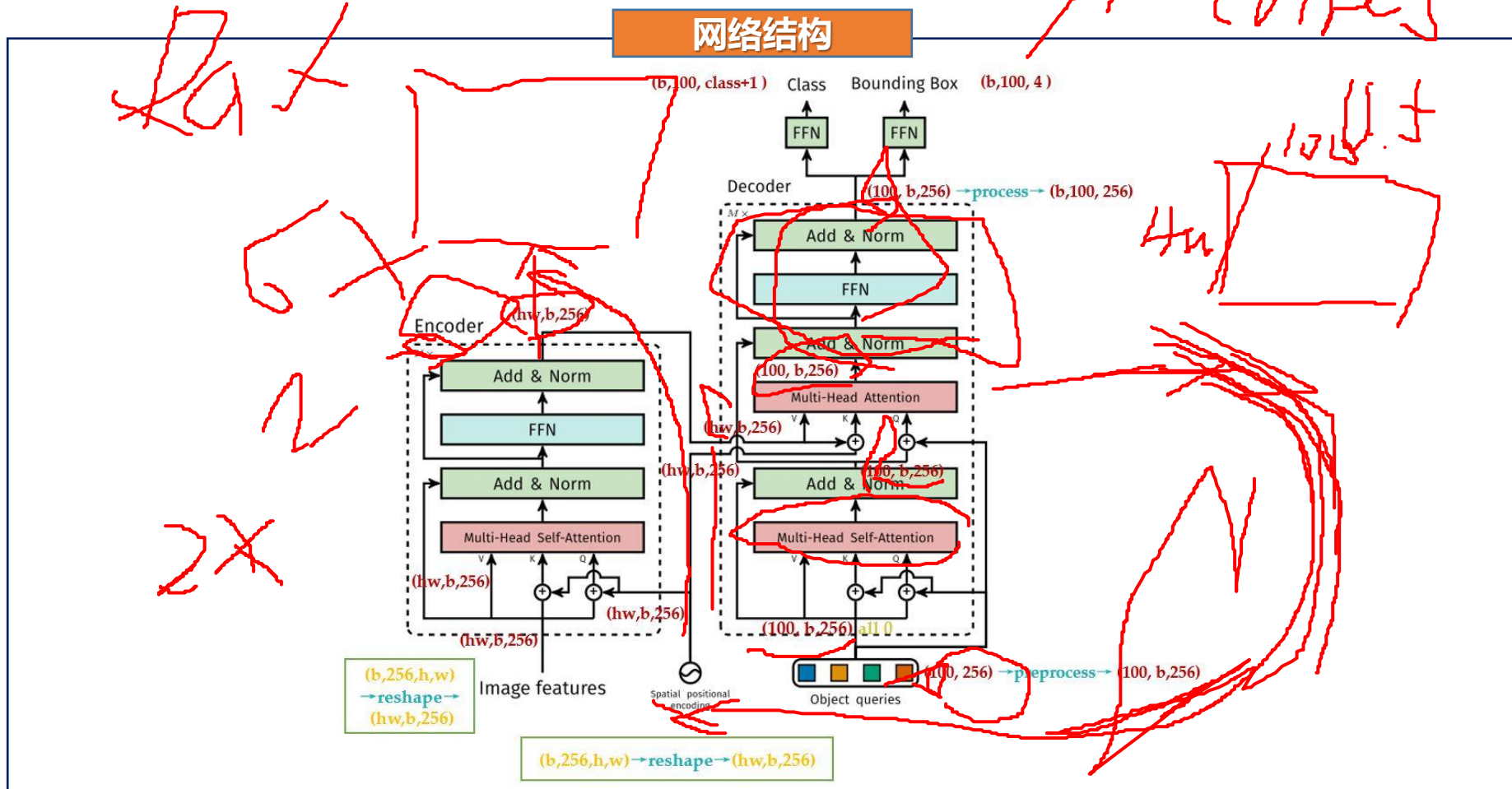
b) $PE_{(pos_x, 2i+1)} = \cos(pos_x / 10000^{2i/128})$

c) $PE_{(pos_y, 2i)} = \sin(pos_y / 10000^{2i/128})$

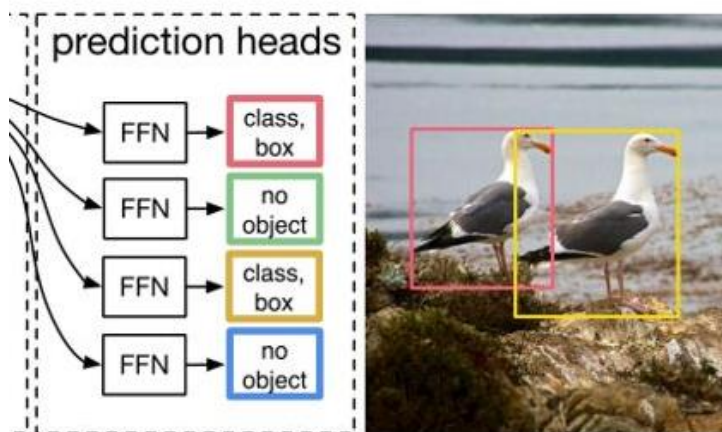
d) $PE_{(pos_y, 2i+1)} = \cos(pos_y / 10000^{2i/128})$



◆ DETR



◆ DETR的损失函数



$$L_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^N \left[-\log \hat{p}_{\hat{\sigma}(i)}(c_i) + 1_{\{c_i \neq \emptyset\}} L_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]$$

$$L_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) = \lambda_{\text{iou}} L_{\text{iou}}(b_i, \hat{b}_{\hat{\sigma}(i)}) + \lambda_{\text{L1}} \|b_i - \hat{b}_{\hat{\sigma}(i)}\|_1, \text{ where } \lambda_{\text{iou}}, \lambda_{\text{L1}} \in \mathbb{R}$$