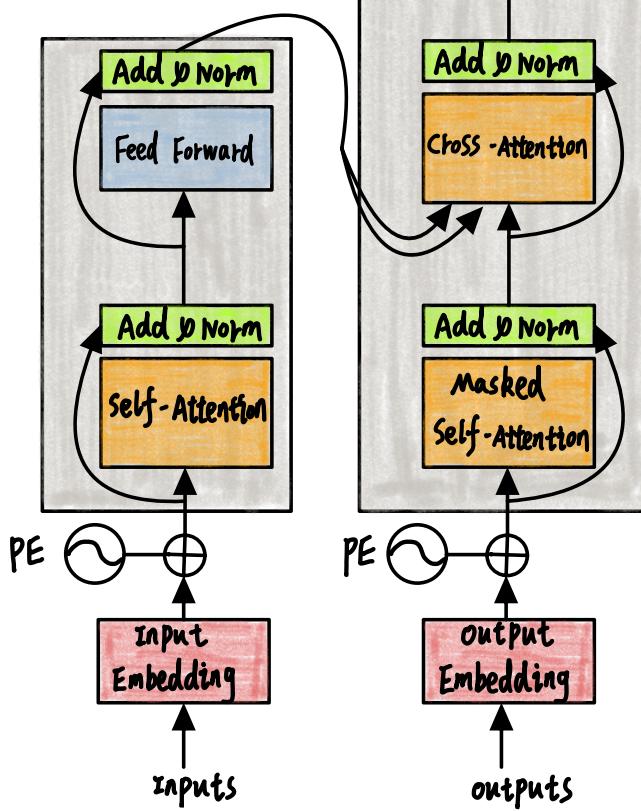


## Transformer



output probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Cross -Attention

Add & Norm

Masked  
Self-Attention

output  
Embedding

Input  
Embedding

PE

outputs

$$(\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^n) = E_1(o^1, o^2, \dots, o^n)$$

$$\tilde{X} = \sum_{(m,n)} O$$

$$(x^1, x^2, \dots, x^n) = (\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^n) + (p^1, p^2, \dots, p^n)$$

$$x = \tilde{x} + p$$

$$Q_i = W_{q,i} X$$

$$K_i = W_{k,i} X$$

$$V_i = W_{v,i} X$$

$$A_i = K_i^T \cdot Q_i$$

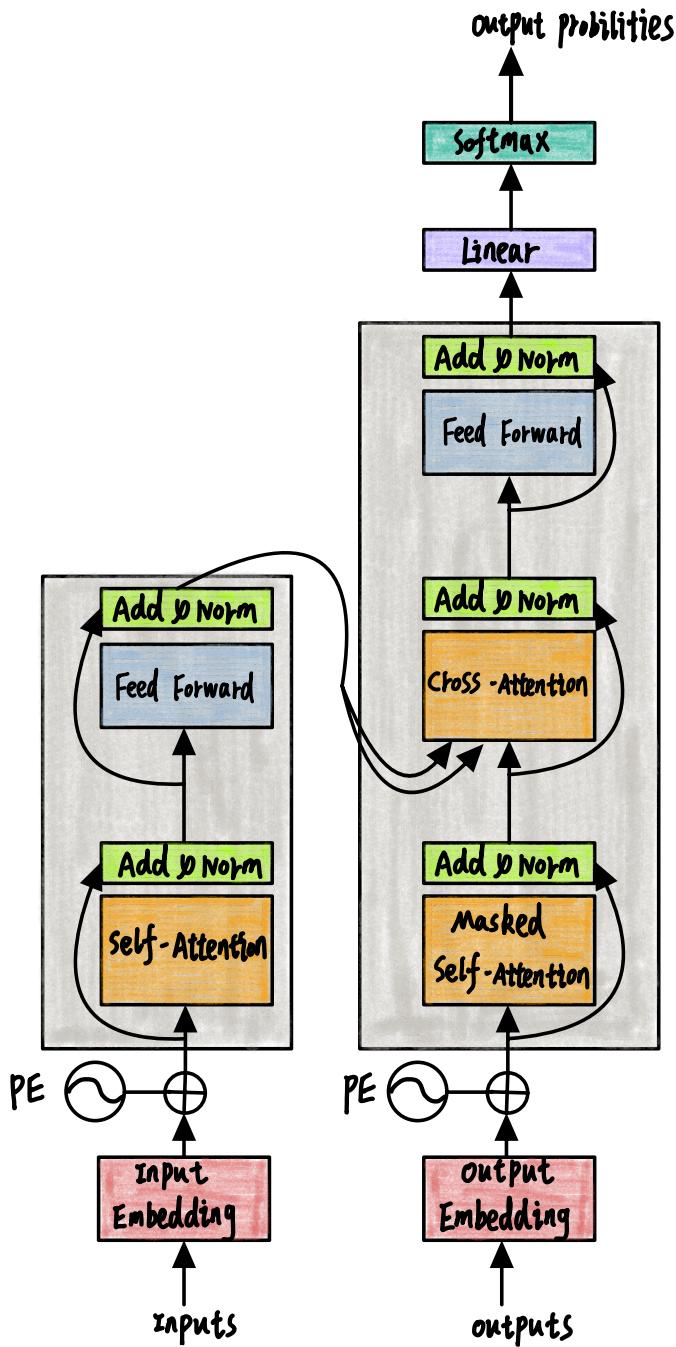
$$A'_i = \text{Softmax}(A_i)$$

$$Y_i = V_i A'_i$$

$$Y'_i = \text{Norm}(Y_i + X)$$

$$Y''_i = W_E Y'_i$$

$$S = \text{Norm}(Y''_i + Y'_i)$$



$$(\tilde{z}^1, \tilde{z}^2, \dots, \tilde{z}^t, \dots, \tilde{z}^n) = E_2(u^1, u^2, \dots, u^t, \dots, u^n)$$

$$\tilde{Z} = \sum_{(m,n)} \tilde{z}^m$$

$$(z^1, z^2, \dots, z^t, \dots, z^n) = (\tilde{z}^1, \tilde{z}^2, \dots, \tilde{z}^n) + (p^1, p^2, \dots, p^t, \dots, p^n)$$

$$Z = \tilde{Z} + P$$

$$Q_2 = W_{q2} Z$$

$$K_2 = W_{k2} Z$$

$$V_2 = W_{v2} Z$$

$$A_2 = K_2^T \cdot Q_2$$

$$A_2 \stackrel{\text{某些特例}}{\neq} \text{Score} \rightarrow -\infty$$

$$A'_2 = \text{Softmax}(A_2)$$

$$Y_2 = V_2 A'_2$$

$$R = \text{Norm}(Y_2 + Z)$$

$$Q_3 = W_{q3} R$$

$$K_3 = W_{k3} S$$

$$V_3 = W_{v3} S$$

$$A_3 = K_3^T \cdot Q_3$$

$$A'_3 = \text{Softmax}(A_3)$$

$$Y_3 = V_3 A'_3$$

$$Y'_3 = \text{Norm}(Y_3 + R)$$

$$Y''_3 = W_D Y'_3$$

$$Y'''_3 = \text{Norm}(Y''_3 + Y'_3)$$

$$L = W_L Y'''_3 + b$$

$$U' = \text{softmax}(L)$$

怎么回归？

将多层输入和输出连接

线性

if  $L_2^Q$  multi-head?

