

# Transformer

## Language Modeling (MLE)

$$p_{\theta}(y_{1:L} \mid x_{1:T}) = \prod_{i=1}^L p_{\theta}(y_i \mid y_{<i}; \text{enc}(x_{1:T}))$$

## Attention

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

$$\begin{aligned} Q &= (q_1, q_2, \dots, q_N)^{\top} \in \mathbb{R}^{N \times d_q} & W_Q &\in \mathbb{R}^{d_X \times d_q} & X &\in \mathbb{R}^{N \times d_X} \\ K &= (k_1, k_2, \dots, k_M)^{\top} \in \mathbb{R}^{M \times d_q} & W_K &\in \mathbb{R}^{d_X \times d_q} & QK^{\top} &\in \mathbb{R}^{N \times M} \\ V &= (v_1, v_2, \dots, v_N)^{\top} \in \mathbb{R}^{M \times d_v} & W_V &\in \mathbb{R}^{d_X \times d_v} & Y &\in \mathbb{R}^{N \times d_v} \end{aligned}$$

Soft Attention:  $a = \text{Softmax}$ , Hard Attention:  $a = \text{argmax}$

## Complexity

Time:  $(m \times n) \cdot (n \times p) : \mathcal{O}(mnp)$ . Space:  $(m \times n) \cdot (n \times p) : \mathcal{O}(mp)$

Time:  $QK^{\top} : \mathcal{O}(NMd_q), a(QK^{\top}) : \mathcal{O}(NM), AV : \mathcal{O}(NMd_v)$  All:  $\mathcal{O}(NMd_q + NMd_v)$

Space:  $\mathcal{O}(MN + Nd_v)$

## MHA

$$\begin{aligned} \text{Multihead}(Q, K, V; a) &= \text{concat}(\text{head}_1, \text{head}_2, \dots) W_O \\ \text{head}_i(Q, K, V; a) &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V; a); h = |\text{head}| \\ W_i^Q &\in \mathbb{R}^{d_q \times \tilde{d}_q} & W_i^K &\in \mathbb{R}^{d_q \times \tilde{d}_q} & W_i^V &\in \mathbb{R}^{d_v \times \tilde{d}_v} & W_O &\in \mathbb{R}^{h\tilde{d}_v \times d_o} \\ Q^* &: h \times N \times \tilde{d}_q & K^* &: h \times M \times \tilde{d}_q & V^* &: h \times M \times \tilde{d}_v & H &: h \times N \times \tilde{d}_v \end{aligned}$$

$h = 1$ : Time:  $\mathcal{O}(NM\tilde{d}_q + NM\tilde{d}_v)$ . Space:  $\mathcal{O}(MN + N\tilde{d}_v)$

Time:  $\mathcal{O}(h(d_q\tilde{d}_q(M + N) + d_v\tilde{d}_vM) + h(NM\tilde{d}_q + NM\tilde{d}_v) + Nh\tilde{d}_vd_o)$

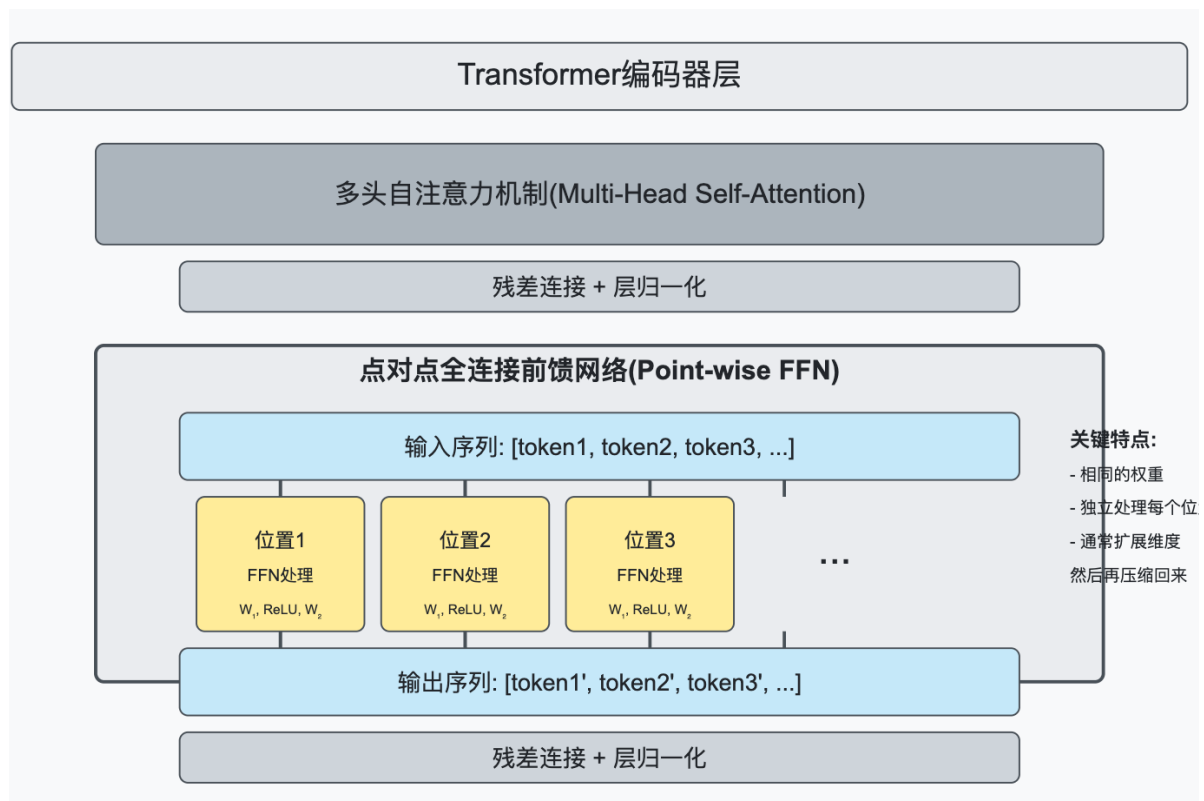
Projection + Attention + Output Projection

Space:  $\mathcal{O}(h(\tilde{d}_q(M + N) + \tilde{d}_vM) + h(MN + N\tilde{d}_v) + Nd_o)$

## Layer Norm

$$\begin{aligned} Y &= \{y_1, y_2, \dots, y_N\}, \quad \text{where } y_i = XW + b \\ y_i &\leftarrow \frac{y_i - \mu}{\sigma}, \quad \text{where } \mu = \mathbb{E}Y, \sigma = \sqrt{\mathbb{E}_{y \in Y}[y - \mu]^2} \\ \text{Add \& Norm}(X) &= \text{LN}(X + \text{Sublayer}(X)) \end{aligned}$$

## Point-wise FFN (PFFN)



假设输入  $X: \text{Batch} \times \text{Context} \times \text{Embedding} : B \times C \times E$

$$\text{FFN } f(X) = \text{ReLU}((XW_1 + b_1)W_2 + b_2) : E \rightarrow V$$

$$(B \times C \times E) \xrightarrow{\text{PFFN}: E \rightarrow V} (B \times C \times V)$$