

Part 1: Encode Layer

$$X = [t_1, t_2, t_3, \dots]$$

$$E' = \text{Wemb}(X_i)$$

$$E' = E' \sqrt{d}$$

$$X^0 = E' + PE \left\{ \begin{array}{l} PE_{pos, 2i} = \sin\left(\frac{pos}{10^{\frac{2i}{d}}}\right) \\ PE_{pos, 2i+1} = \cos\left(\frac{pos}{10^{\frac{2i}{d}}}\right) \end{array} \right\} i \in [0, \dots, \frac{d}{2}-1]$$

$$Q = X^0 W^Q$$

$$K = X^0 W^K$$

$$V = X^0 W^V$$

$$\text{score} = \frac{Q_i K_j^T}{\sqrt{d_k}}$$

mask to $-\infty$

$$\text{head}_j = \text{Attention}(Q, K, V) = \text{Softmax}(\text{scores}) \cdot V$$

$$\Psi = \text{concat}(\text{head}_1, \dots, \text{head}_h)$$

$$\text{MHA}(H) = \Psi \cdot W^O$$

Add initial X^0 and normalize:

$$\text{residue}(x) = X^0 + \text{MHA}(H)$$

$$\begin{bmatrix} pos \\ 10^{\frac{2i}{d}} \end{bmatrix}$$

$$PE = [\cos \theta, \sin \theta, \dots]$$

single head

$$H^1 = \gamma \odot \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta \quad \text{layernorm}$$

$$\mu = \frac{1}{d} \sum x_i, \quad \sigma^2 = \frac{1}{d} \sum (x_i - \mu)^2$$

$\gamma, \beta \Rightarrow$ learnable params

$$\text{FFN}(x) = \phi(x W_1 + B_1) W_2 + B_2$$

$$H^1 = \text{LayerNorm}(H^1 + \text{FFN}(H^1))$$

$$\phi = x \cdot \Phi(x)$$

$$\Phi(x) = \frac{1}{2} \left[1 + \tanh\left(\frac{x}{\sqrt{2}}\right) \right]$$

GELU

SO FAR: Encode = Embed(x) + PE \rightarrow Scores \rightarrow Single head \rightarrow MHA \rightarrow add & LayerNorm \rightarrow FFN \rightarrow A&L

$$H^l = \text{Encode}(H^{l-1}), \quad H^0 = X^0$$

Part 2: Decode Layer (Analogous to Encoder Layer)

$$T_0 = \text{Embed}(Y) \sqrt{d} + PE$$

$$T_1 = \text{LayerNorm}(T_0 + \text{MHA}_{\text{masked}}(T_0))$$

Cross attention: T_L from iter 2, T_L from iter 2

$$T_2 = \text{LayerNorm}(T_1 + \text{MHA}_{\text{cross}}(Q = T_1, K \Rightarrow H^l, V \Rightarrow H^l))$$

$$T^3 = \text{LayerNorm}(T_2 + \text{FFN}(T_2))$$

$$T^{l+1} = T_3 \Rightarrow \text{Stack decoder layer} \Rightarrow T^l = \text{Decode}(T^{l-1})$$

Part 3: Prediction

$$\text{logits} = T^L W_{\text{out}} + B_{\text{out}}$$

$$P(\text{next token} | \text{context}) = \text{softmax}(\text{logits})$$