# Transformer

$$\text{output} = \text{Attention} \left( \underset{\text{"query"}}{q}, \underset{\text{"keys"}}{k}, \underset{\text{"values"}}{v} \right)$$

$$d_t = \frac{\exp\left( q^T k_t / \sqrt{\dim(q)} \right)}{\sum_{j=1}^{T} \exp\left( q^T k_j / \sqrt{\dim(q)} \right)}$$

$$\text{output} = \sum_{j=1}^{T} d_j \cdot v_j$$

$$q = h_{k-1}$$
$$k = \{ h_1, h_2, \ldots, h_T \}$$
$$v = \{ h_1, h_2, \ldots, h_T \}$$

$$Q \in \mathbb{R}^{\tilde{T} \times H}, \quad K \in \mathbb{R}^{T \times H}, \quad V \in \mathbb{R}^{T \times H}$$

$$O = \text{Attention}(Q, K, V)$$

$$\iff O_k = \text{Att}(q_k, K, V) \quad k = 1, \dots, \tilde{T}$$

$$\underset{\underset{\mathbb{R}^{\tilde{T} \times H}}{\uparrow}}{O} = \underset{\underset{\text{rowwise}}{\uparrow}}{\text{Softmax}} \left( \underset{\underset{\mathbb{R}^{\tilde{T} \times T}}{\uparrow}}{Q \cdot K^T / \sqrt{H}} \right) \cdot \underset{\underset{\mathbb{R}^{T \times H}}{\uparrow}}{V}$$

# Multi Head Attention

$$\text{head}_i = \text{Attention}\left(Q \cdot W_i^Q, \; K \cdot W_i^K, \; V \cdot W_i^V\right)$$

$\mathbb{R}^{\tilde{T} \times H_{head}}$     $\tilde{T} \times H$   $H \times H_{head}$   $T \times H$   $H \times H_{head}$   $T \times H$   $H \times H_{head}$

$$\left\{ W_i^Q, W_i^K, W_i^V \right\}_{i=1}^{n} - \text{trainable params}$$

$$O = \text{MultiHeadAttention} = \text{Concat}\left(\text{head}_1, \text{head}_2, \dots, \text{head}_n\right) \cdot W^O$$

$\mathbb{R}^{\tilde{T} \times H}$        $\tilde{T} \times H_{head}$   $\tilde{T} \times H_{head}$      $n \cdot H_{head} \times H$

# Transformer encoder

$$x_t, \quad t = 1, \ldots, T$$

$$h_t^0 = emb(x_t)$$

$$\text{for } \ell = 1, 2, \ldots, L :$$

$$O^\ell = LayerNorm\left(H^{\ell-1} + MHAtt(H^{\ell-1}, H^{\ell-1}, H^{\ell-1})\right)$$

$$H^\ell = LayerNorm\left(O^\ell + g(O^\ell W_1 + b_1) W_2 + b_2\right)$$

$$\underset{T \times H}{} \quad \underset{H \times H}{} \quad \underset{H \times 1}{} \quad \underset{H \times H}{} \quad \underset{H \times 1}{}$$

Add & LayerNorm

block 1

$MHAtt.(E, E, E)$

$emb(x_1)$  $emb(x_2)$  $emb(x_T)$

$posemb(1)$  $posemb(2)$  $posemb(T)$

$x_1$  $x_2$  $x_T$

- - - -

$$h_{t+1} = g\left(\underset{H \times H}{W^h} \cdot h_t + \underset{H \times L}{W^x} \cdot x_{t+1} + b\right)$$
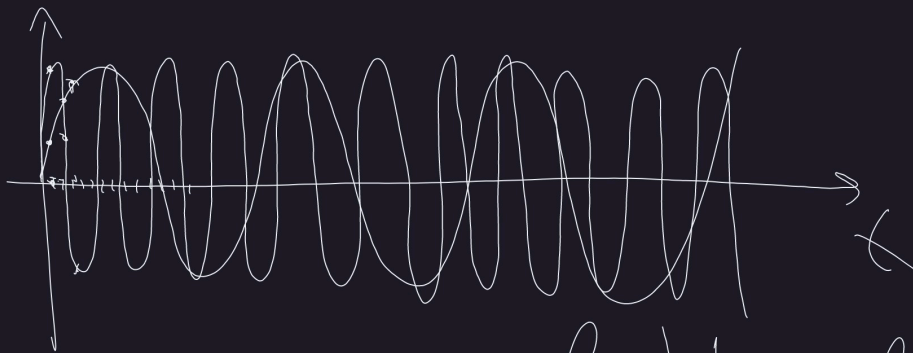
Comp. complexity :

| RNN | MH SelfAtt. |
|---|---|
| $O(T \cdot H \cdot H)$ | $O\left(n\left(T \cdot H \cdot H_{head} + T \cdot T \cdot H_{head}\right)\right)$ |
| | $= O\left(T \cdot H \cdot H + T \cdot T \cdot H\right)$ |
| | $n H_{head} = H$ |

$T \quad H \quad H_{head} \quad n$

$$posemb_{t,2i} = \sin(f_i \cdot t)$$
$$posemb_{t,2i+1} = \cos(f_i \cdot t)$$

$$f_i = 10000^{\frac{emb\_dim}{2i}}$$



$$posemb_{t+k} = W_k \cdot posemb_t \Leftarrow$$

$$posemb_{t+k,2i} = \sin(f_i(t+k)) =$$
$$= \sin(f_i \cdot t)\cos(f_i \cdot k) +$$
$$+ \cos(f_i \cdot t) \cdot \sin(f_i \cdot k) =$$
$$= posemb_{t,2i} \cdot \cos(f_i \cdot k) +$$
$$+ posemb_{t,2i+1} \cdot \sin(f_i \cdot t)$$

Transformer decoder

block 2 → ... → block L

block L

MH Att. $(\tilde{O}, H^L, H^L)$

Masked MH Att. $(\tilde{E}, \tilde{E}, \tilde{E})$

$emb(\tilde{x}_1) + posemb(1)$

$emb(\tilde{x}_T) + posemb(\tilde{T})$

$\tilde{x}_1$   $\tilde{x}_2$  ----  $\tilde{x}_T$

SOS

$\tilde{x}_t, \quad t = 1, \ldots, T$

$h_t^0 = emb(\tilde{x}_t) + posemb(t)$

for $\ell = 1, 2, \ldots, L$:

$\tilde{O}^\ell = LayerNorm\left( \tilde{H}^{\ell-1} + \right.$

$\left. + MaskedMHAtt.(\tilde{H}^{\ell-1}, \tilde{H}^{\ell-1}, \tilde{H}^{\ell-1}) \right)$

$\tilde{\tilde{O}}^\ell = LayerNorm\left( \tilde{O}^\ell + \right.$

$\left. + MHAtt.(\tilde{O}^\ell, H^L, H^L) \right)$

$\tilde{H}^\ell = LayerNorm\left( \tilde{\tilde{O}}^\ell + \right.$

$\left. + g(\tilde{\tilde{O}}^\ell W_1 + b_1) W_2 + b_2 \right)$

$$\tilde{O} = \text{Masked Attention}\left(\tilde{H}, \tilde{H}, \tilde{H}\right) =$$

$$= \text{Softmax}\left(\tilde{H} \cdot \tilde{H}^T \Big/ \sqrt{d\text{im}} + m\right) \cdot \tilde{H}$$