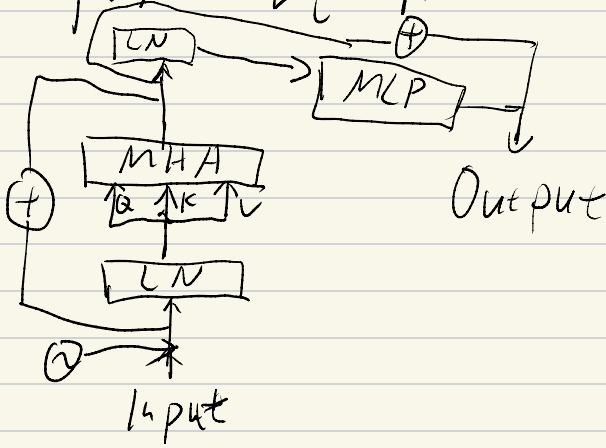


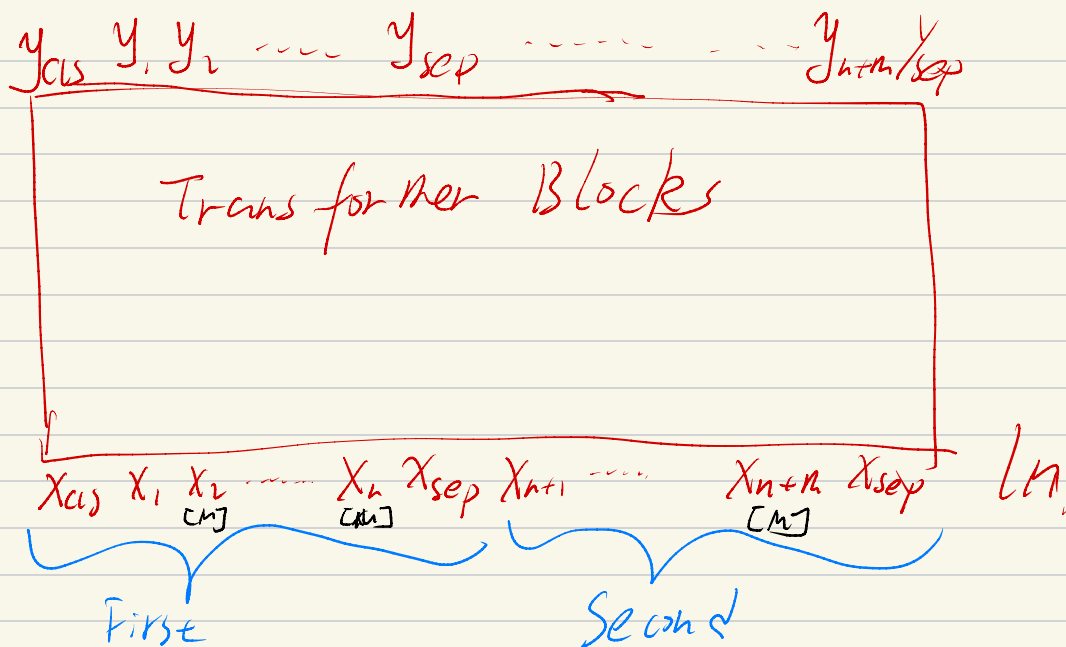
25.02.18

Transformer Block



- ① MHA: Parallelization
- ② ResNet / LN: Better Optimization
- ③ MLP: Non-Linearity, Information Compression
- ④ Position Encoding, Sequence Info
- ⑤ Future-Sequence Masked Attention

Out



MLM: For all $[m], [r], [n]$
we use y_2 to predict

$$z_2 = \text{Softmax}(A_{MLM} \cdot y_2 + b_{MLM})$$

$$\text{Loss: } -\log(z_2) \quad \begin{matrix} \downarrow \\ \mathbb{R}^{|\mathcal{V}| \cdot d_{model}} \\ \text{Actual token of } x_2 \end{matrix}$$

$$\text{NSP: } \text{Sigmoid}(A_{NSP} \cdot y_{cls} + b_{cls})$$

$$\text{Loss: } \text{Cross-Entropy} \quad \begin{matrix} \downarrow \\ \mathbb{R}^{d_{model}} \end{matrix}$$