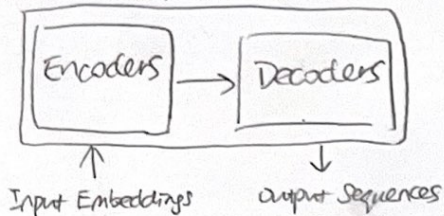
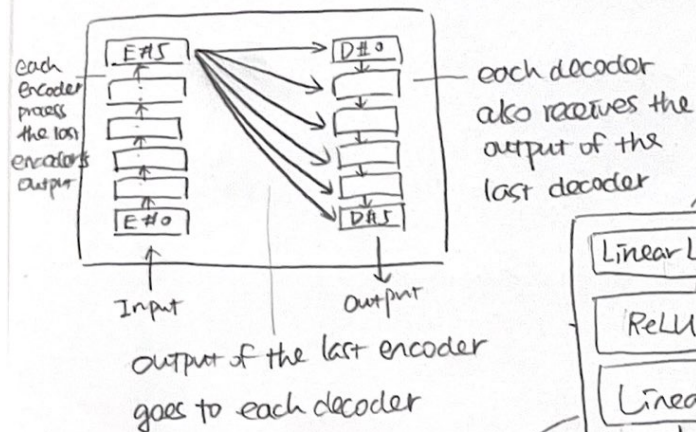


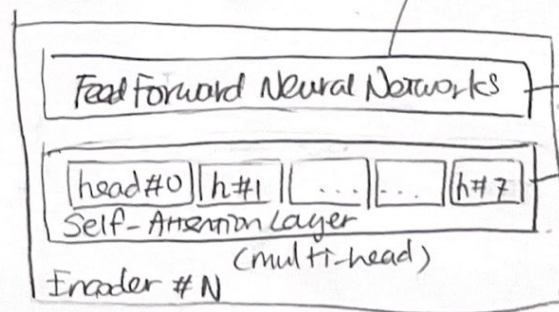
1. Transformer Architecture Encoder



① Same no of encoders & decoders stacked (6 in the original paper)



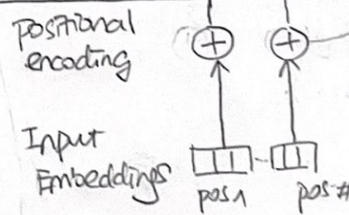
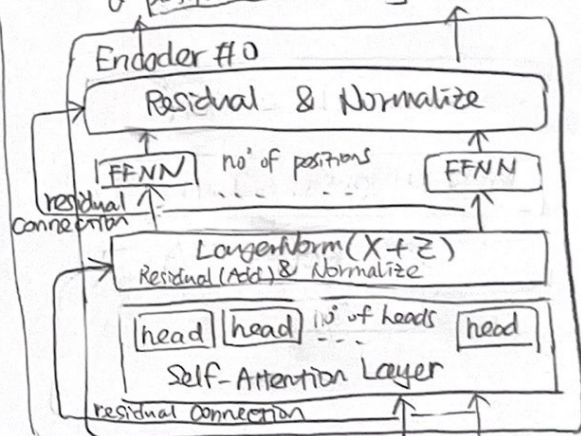
② Encoder layers



②-① Same FFNN applied to each position's Self-attention layer output
parallel processing all positions independently
calculate all positions inter-dependently
all heads are different

②-②

- more detailed Encoder architecture with Residual Connection & Layer Norm & positional encoding



②-③ positional encoding

- gives the model info about the position of each token in the sequence (and relative positions)
- same dimension as embedding so they can be added together
- use sine and cosine functions of different frequency

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

i : vector index, $\in (0, \dots, \frac{d}{2}-1)$

d : dimensionality of embedding

- a way to generate unique encoding for each position
- can scale to sequence lengths that are bigger than max sequence length in training data