

◦ Recap: Attention mechanism

◦ Transformer: Architecture, attention, I/O relation

◦ Attention: $J(i,j) = \sum_{a,b \in N_k(i,j)} \text{softmax}_{a,b}(\langle q(i,j), k(a,b) \rangle) \cdot v(a,b)$.

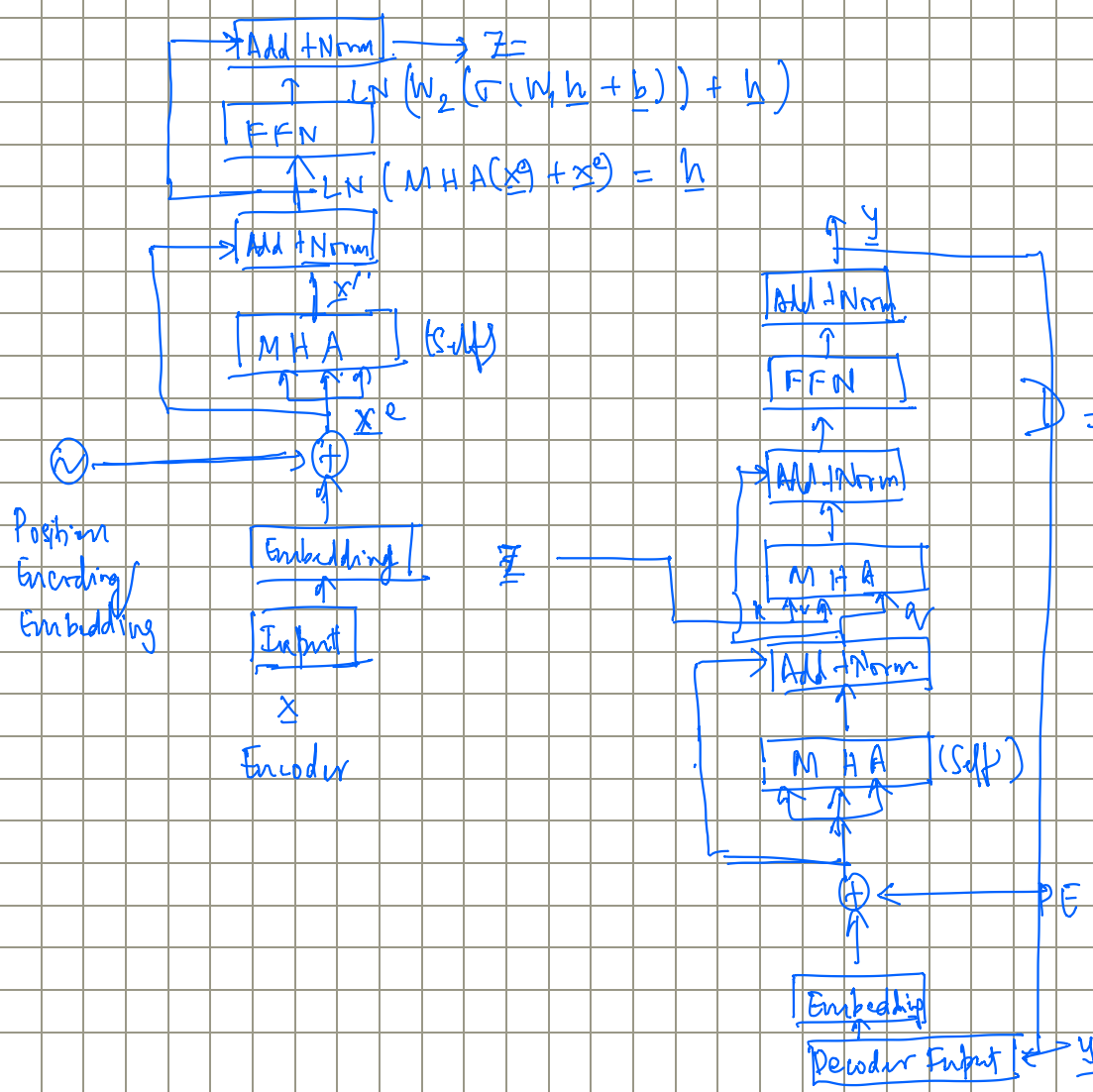
- $q(i,j) = W_q z(i,j)$ query
 - $k(i,j) = W_k z(i,j)$ key
 - $v(i,j) = W_v z(i,j)$ value
- } Self-attention

- $z(i,j) \in \mathbb{R}^e$, $q(i,j), k(i,j), v(i,j) \in \mathbb{R}^d$

- $W_q, W_k, W_v \in \mathbb{R}^{d \times e}$

- $N_k(i,j)$ is a neighbourhood around (i,j) of size $k \times k$

◦ Transformer architecture (Vaswani et al. 2017)



* Lookup GELU activation

$$\text{GELU}(u) = u \cdot \phi(u)$$

$$\phi(u) = \text{CDF}(\mathcal{N}(u, 1))$$

$$= \left\{ (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \right. \\ \left. \dots (x^{(T)}, y^{(T)}) \right\}$$

time