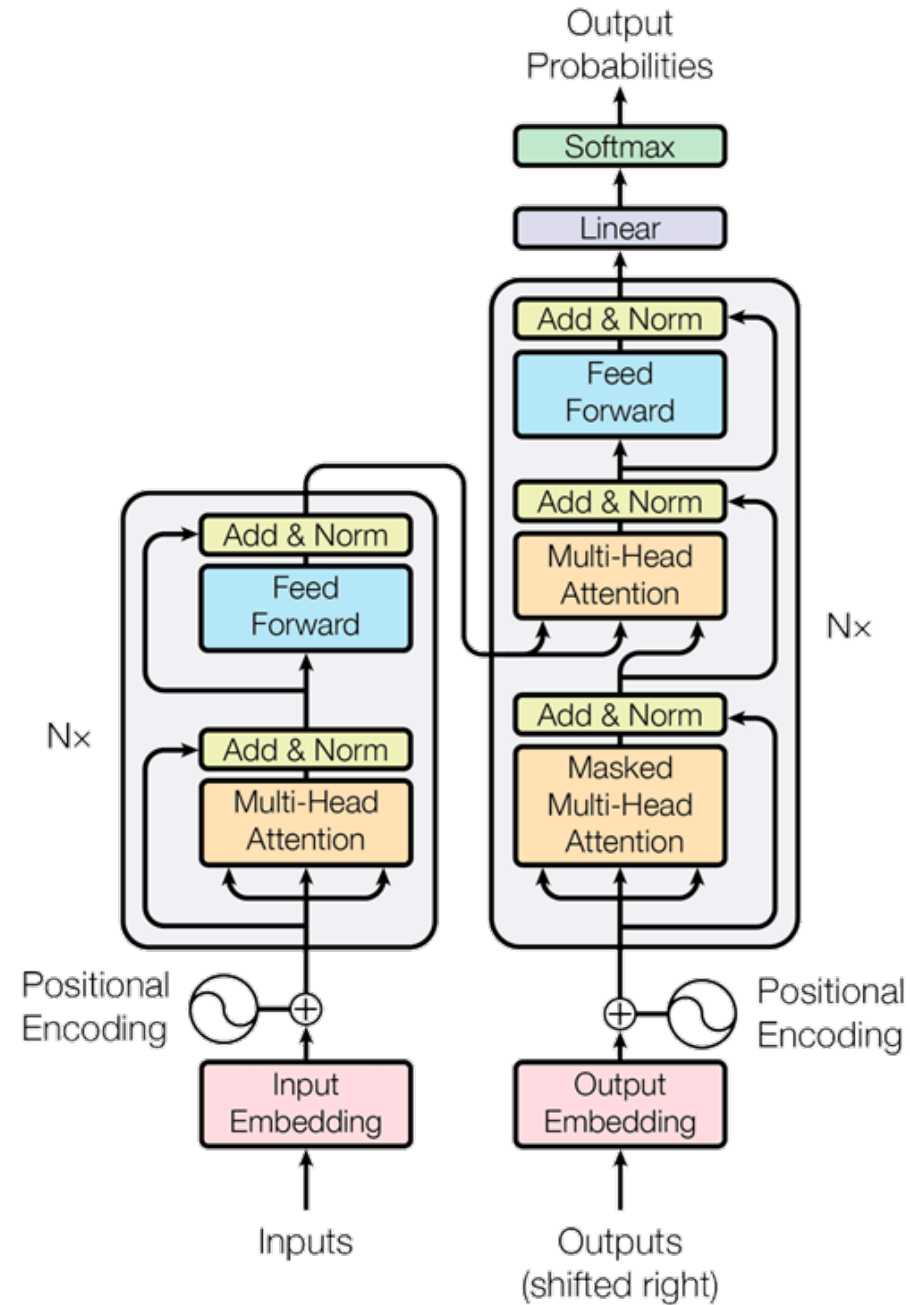


# Transformer: Encoder Block

Ki Hyun Kim

[nlp.with.deep.learning@gmail.com](mailto:nlp.with.deep.learning@gmail.com)

# Transformer



# Equations

- Given Dataset,

$$\mathcal{D} = \{x^i, y^i\}_{i=1}^N$$
$$x^i = \{x_1^i, \dots, x_m^i\} \text{ and } y^i = \{y_0^i, y_1^i, \dots, y_n^i\},$$

where  $y_0 = \langle \text{BOS} \rangle$  and  $y_n = \langle \text{EOS} \rangle$ .

- What we want is

$$\hat{y}_{1:n} = f(x_{1:m} : \theta)$$

# Equations

- $Q$ ,  $K$  and  $V$  are from previous layer:
  - Residual connections and Layer Normalizations are used.

$$\begin{aligned}h_{0,1:m} &= \text{emb}(x_{1:m}) + \text{pos}(1, m) \\ \tilde{h}_{i,1:m}^{\text{enc}} &= \text{LayerNorm}(\text{Multihead}_i(Q, K, V) + h_{i-1,1:m}^{\text{enc}}), \\ &\text{where } Q = K = V = h_{i-1,1:m}^{\text{enc}}.\end{aligned}$$

$$\begin{aligned}\text{FFN}(h_{i,t}) &= \text{ReLU}(h_{i,t} \cdot W_i^1) \cdot W_i^2 \\ &\text{where } W_i^1 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}} \text{ and } W_i^2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}.\end{aligned}$$

$$h_{i,1:m}^{\text{enc}} = \text{LayerNorm}([\text{FFN}(\tilde{h}_{i,1}^{\text{enc}}); \dots; \text{FFN}(\tilde{h}_{i,m}^{\text{enc}})] + \tilde{h}_{i,1:m})$$

# Equations

- Encoder is stack of encoder blocks:

$$\begin{aligned} h_{\ell_{\text{enc}},1:m}^{\text{enc}} &= \text{Block}_{\text{enc}}(h_{\ell_{\text{enc}}-1,1:m}^{\text{enc}}) \\ &\dots \\ h_{1,1:m}^{\text{enc}} &= \text{Block}_{\text{enc}}(h_{0,1:m}^{\text{enc}}) \end{aligned}$$

# Summary

- Encoder는 self-attention으로 구성되어 있음
  - $Q, K, V$ 는 이전 레이어의 출력 값 - 즉, 같은 값
- Seq2seq의 attention과 달리,  $Q$ 도 모든 time-step을 동시에 연산
  - 빠르지만 굉장히 메모리를 많이 먹게 됨
- Residual connection으로 인해 깊은 네트워크 구성 가능
  - Big LM의 토대 마련