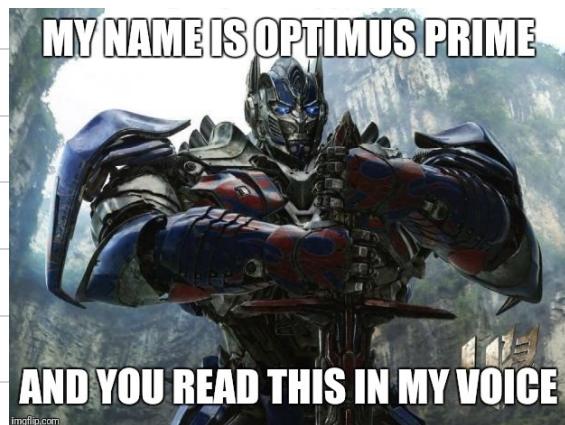


Agenda

1. Seq2Seq Learning
2. Neural Machine Translation
3. Transformer
 - a. Encoder
 - b. Decoder



Self attention.

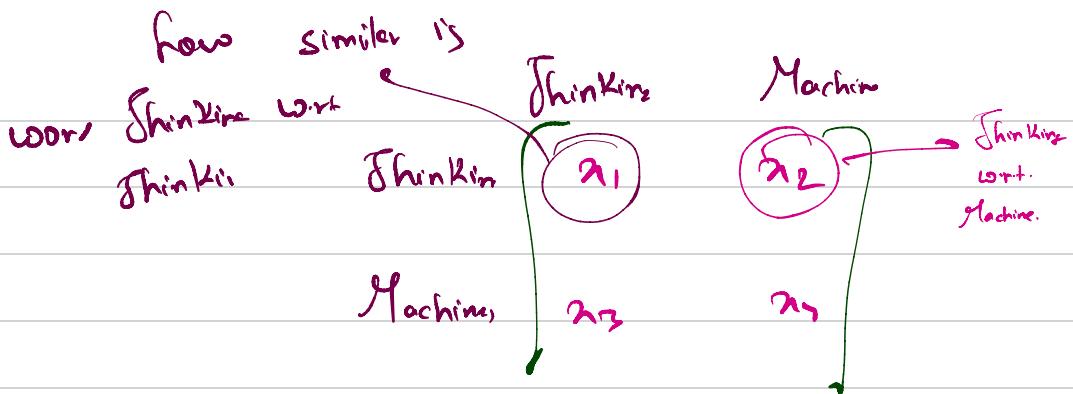
Input $\rightarrow (10, 512) \rightarrow W_q \rightarrow (10, 512)$
 $W_k \rightarrow (10, 512)$
 $W_v \rightarrow (10, 512)$

$W_1 = \text{Thinking} \rightarrow (1, 512)$
 $W_2 = \text{Machine} \rightarrow (1, 512)$ (2, 512)

Create $W(k, Q, V) \rightarrow (512, 512)$
 \downarrow
 $(512, 512)$

$K = (2, 512) \rightarrow \begin{cases} k_1 \\ k_2 \end{cases} \rightarrow (1, 512)$
 $Q = (2, 512) \rightarrow (1, 512)$
 $V = (2, 512)$

Attention Scores = $\text{Softmax} \left(\frac{k \cdot Q^T}{\sqrt{512}} \right) = []_{62 \times 2}$



$$G = \begin{pmatrix} Q_1 \\ Q_2 \end{pmatrix}$$

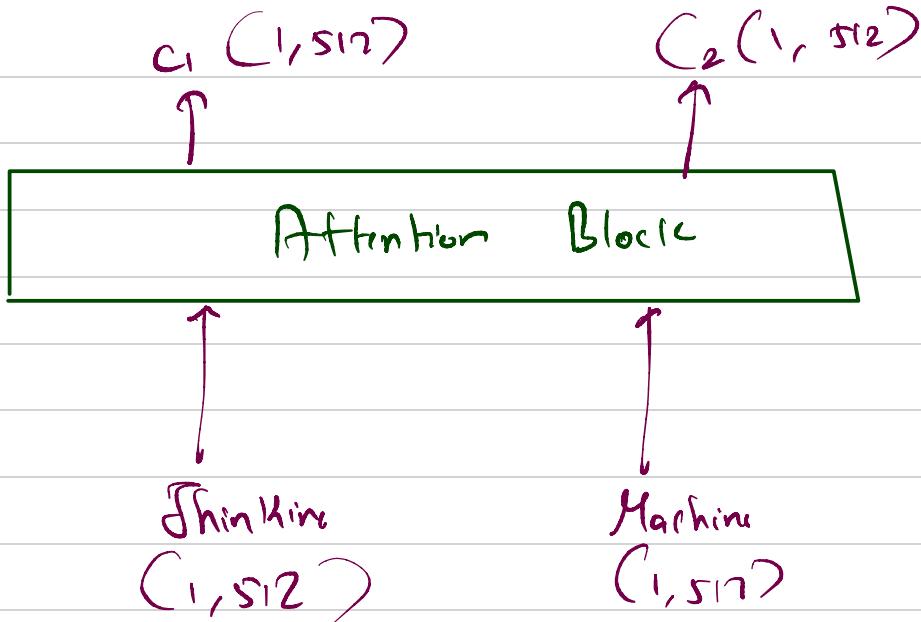
Thinking
Machine

Content Vector 1 = $U_1 \times \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = (1, s_{12})$

(Thinking)

Content Vector 2 = $U_2 \times \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = (1, s_{12})$

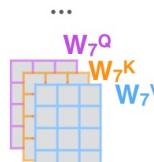
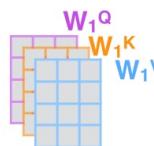
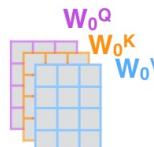
(Machine)



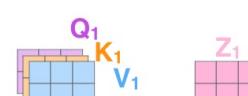
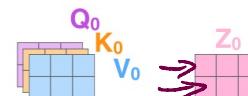
1) This is our input sentence* 2) We embed each word*



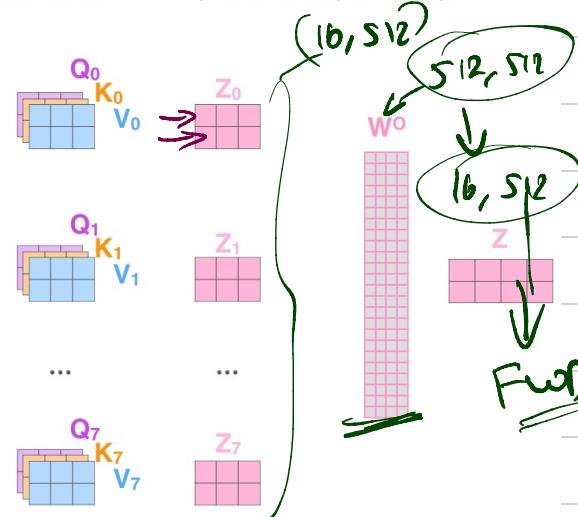
3) Split into 8 heads. We multiply X or R with weight matrices



4) Calculate attention using the resulting $Q/K/V$ matrices



5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



Output

1st attention block =



(2, 512)

)

)

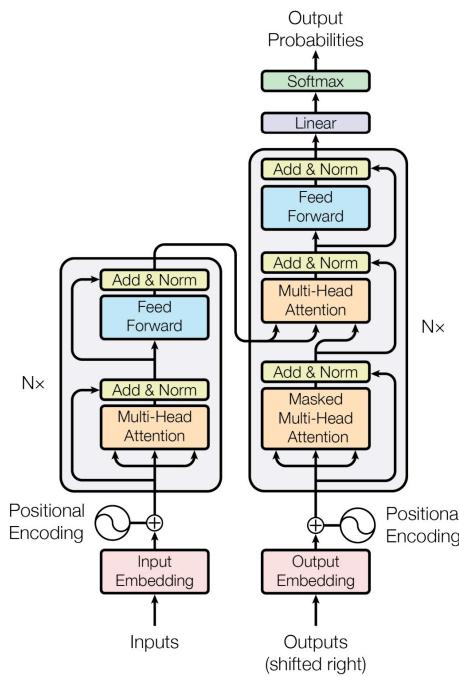
)

8th attention block =

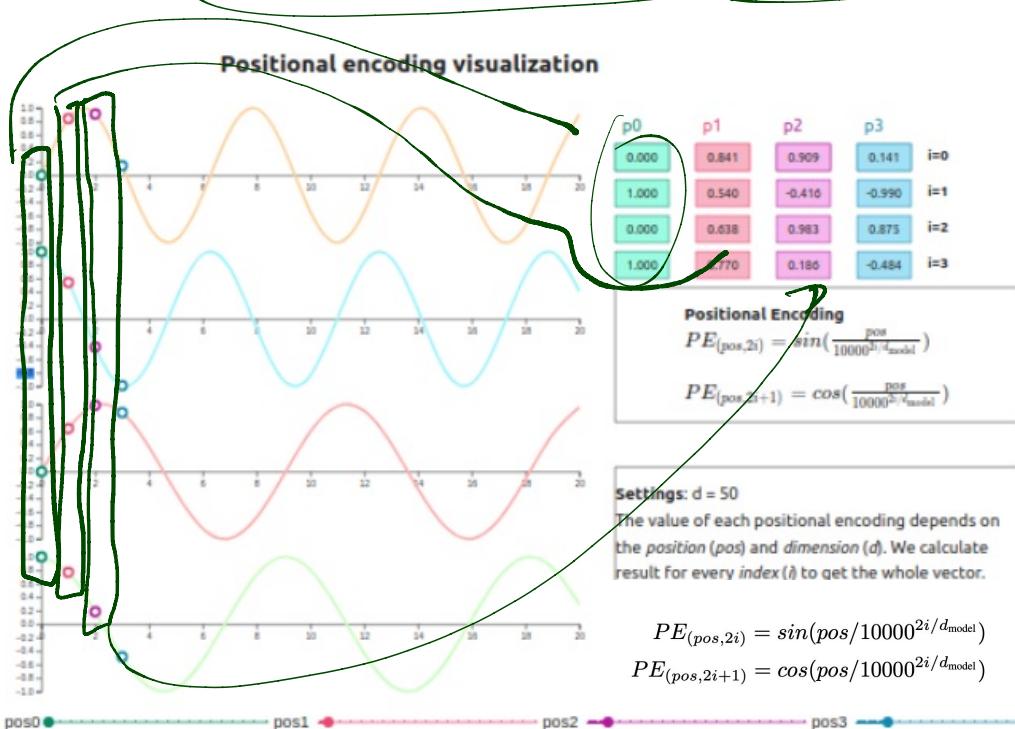
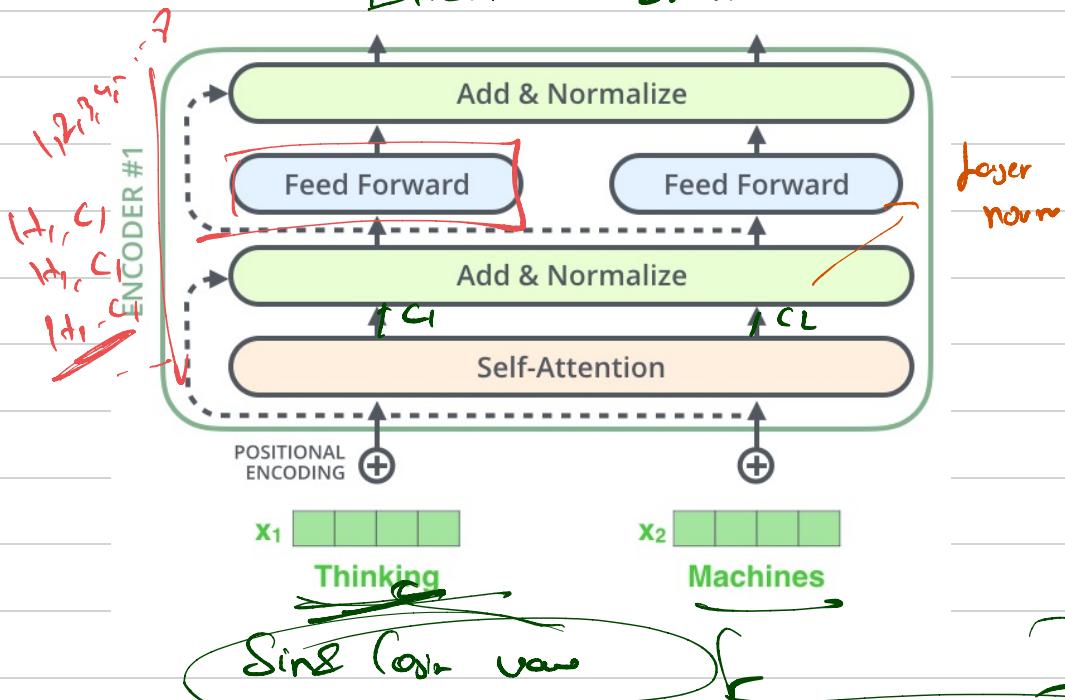


(2, 512)

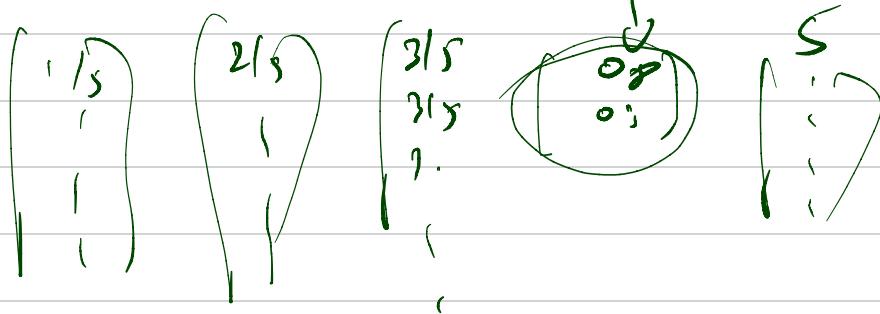
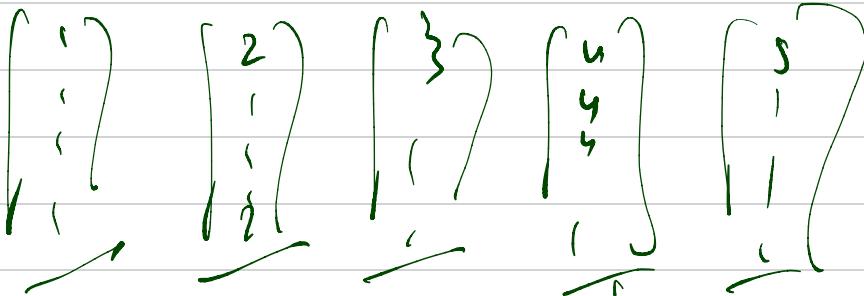
(16, 512)



Encoder Block

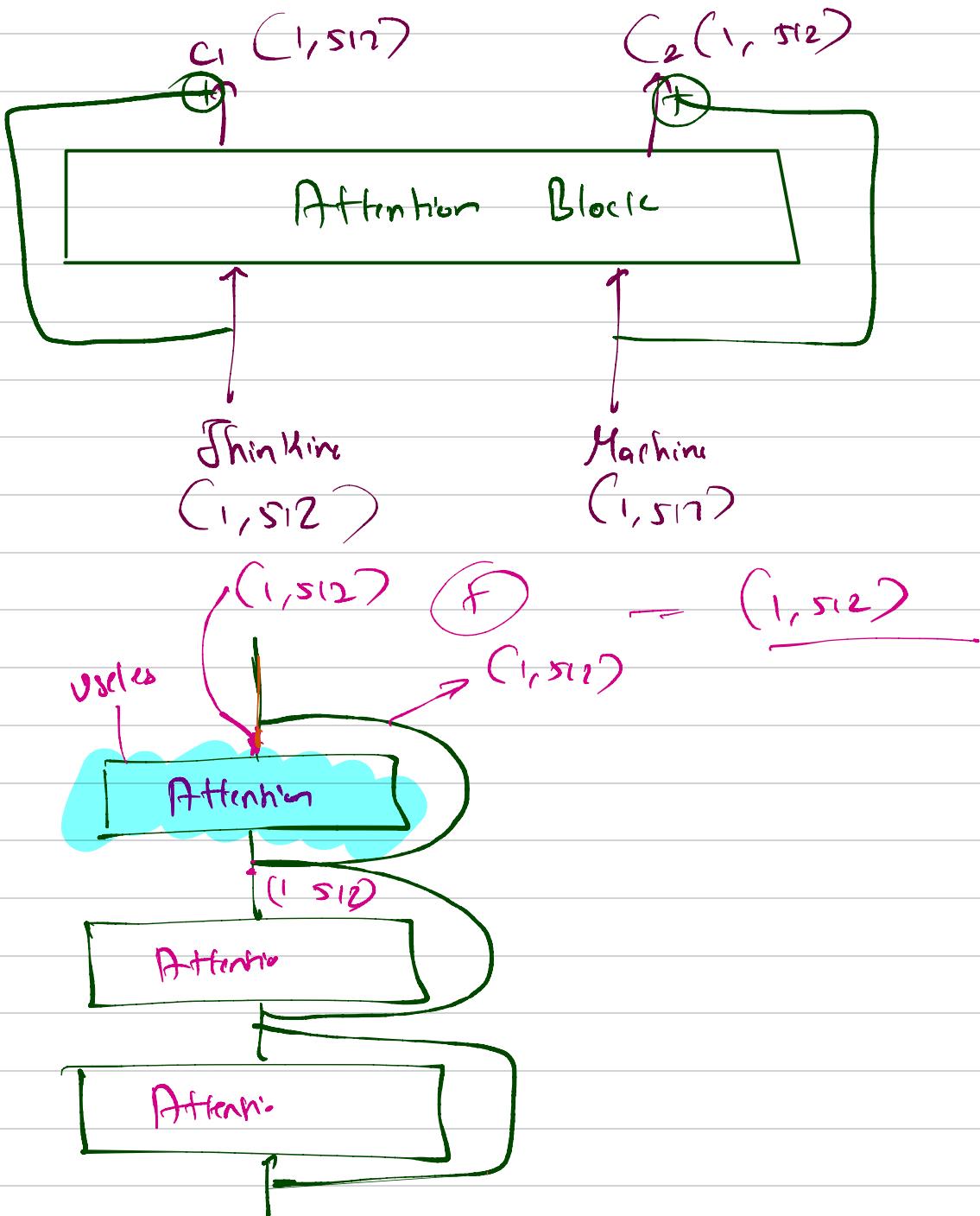


w_1 w_2 w_3 (w_4) w_5



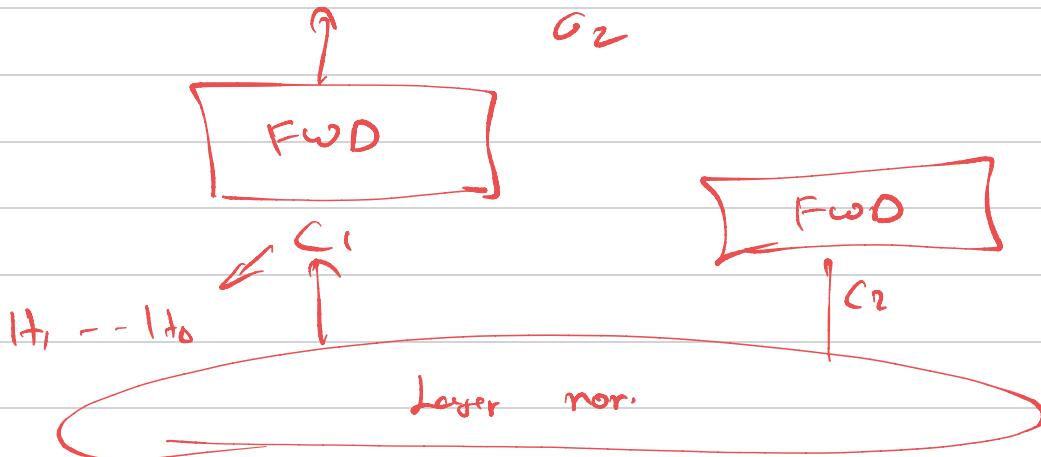
w_1 w_2 - - - w_4 w_5

$$\begin{pmatrix} 8/0 \\ 8/0 \\ 8 \\ 5 \end{pmatrix} - \begin{pmatrix} 0/8 \\ 0/8 \\ 0/8 \end{pmatrix}$$

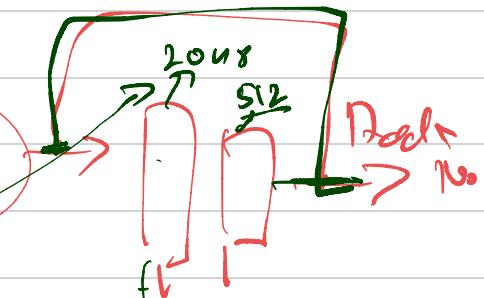
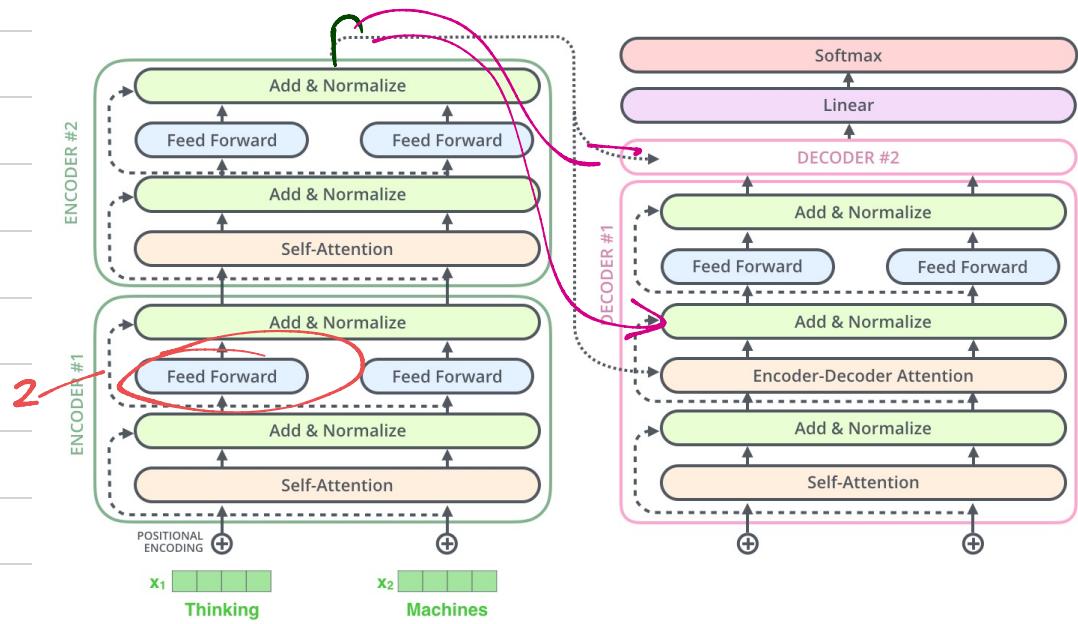


$$C_1 = \boxed{1 | 2 | 3 | 4 | 5} \quad \text{mean std} \leftarrow \frac{x - \bar{x}}{\sigma}$$

$$C_2 = \boxed{1 | 6 | 7 | 8 | 9 | 10 | 11} \quad \mu_{C_2} \quad \frac{x - \bar{x}_2}{\sigma_2}$$



$(2, \text{S12})$

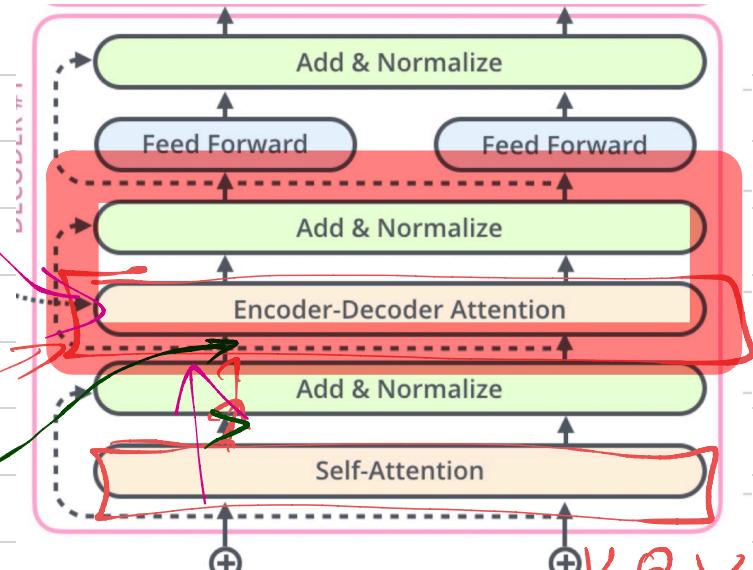


$(2, \text{S12})$

$\text{WD}_{KD} \rightarrow (2, \text{S12})$

$(\text{S12}, \text{S12})$

$(\text{S12}, \text{S12}) \xrightarrow{\text{WD}} (2, \text{S12})$

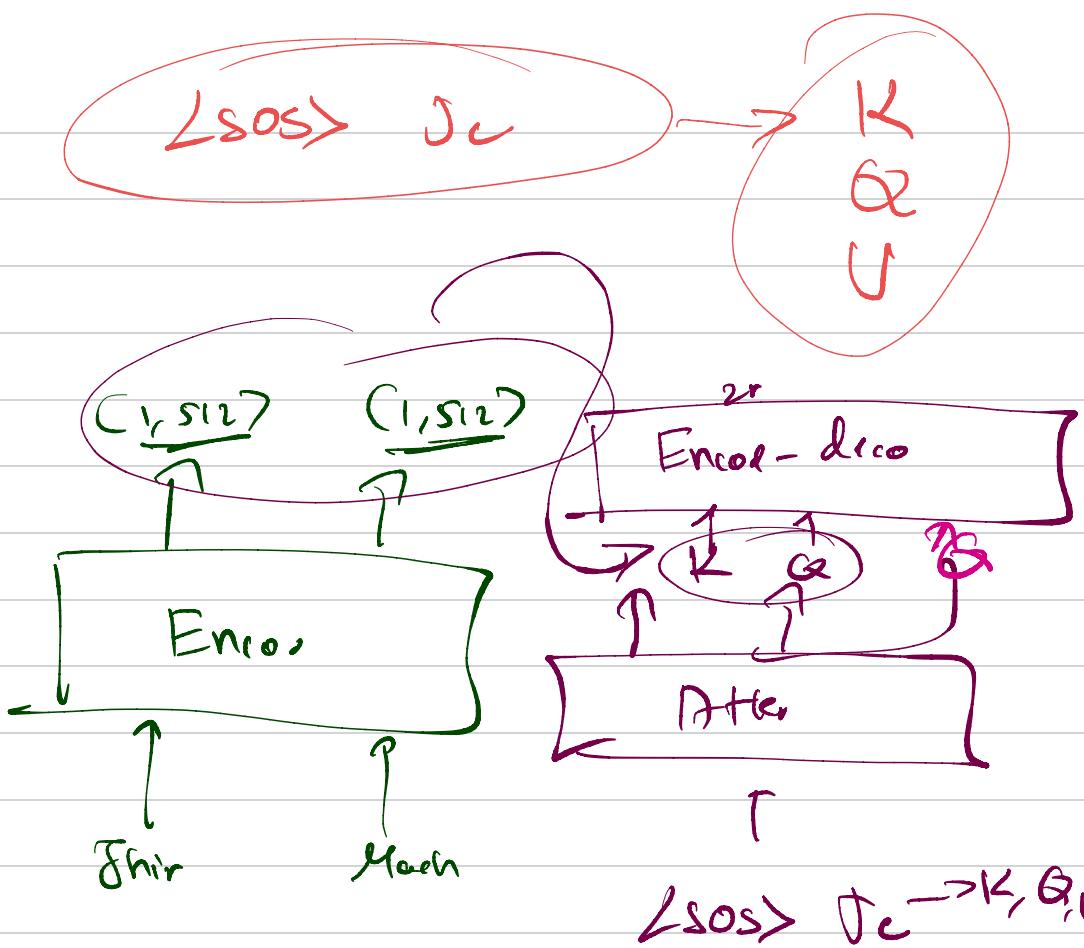


(SOS) → (1, 512)

g love apple → (1, 512)
 Je aime batata
 (250)

Je
 $t=2$

g in Frn



gnnPut → $(128K, s12)$

$$\begin{aligned}
 \omega_Q &= (s12, s12) & K &= (128K, s12) \\
 \omega_{12} &= (s12, s12) & Q &= (128L, s12) \\
 \omega_U &= (s12, s12) & U &= (128L, s12)
 \end{aligned}$$

