

Agenda

1. GRU
2. Transformer Recap End to End
3. Transformer in training
4. BERT
5. Masked Language Modelling

TRANSFORMER

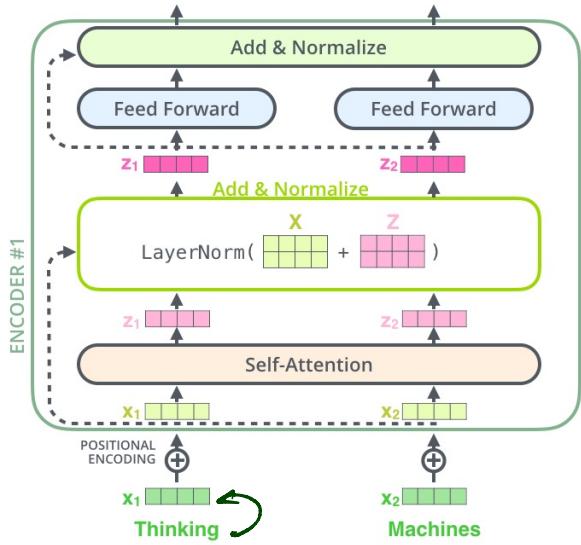
- End - To - End

<https://jalammar.github.io/illustrated-transformer/>

Inference

"I love Pizza"

"He encontró la Pizza"

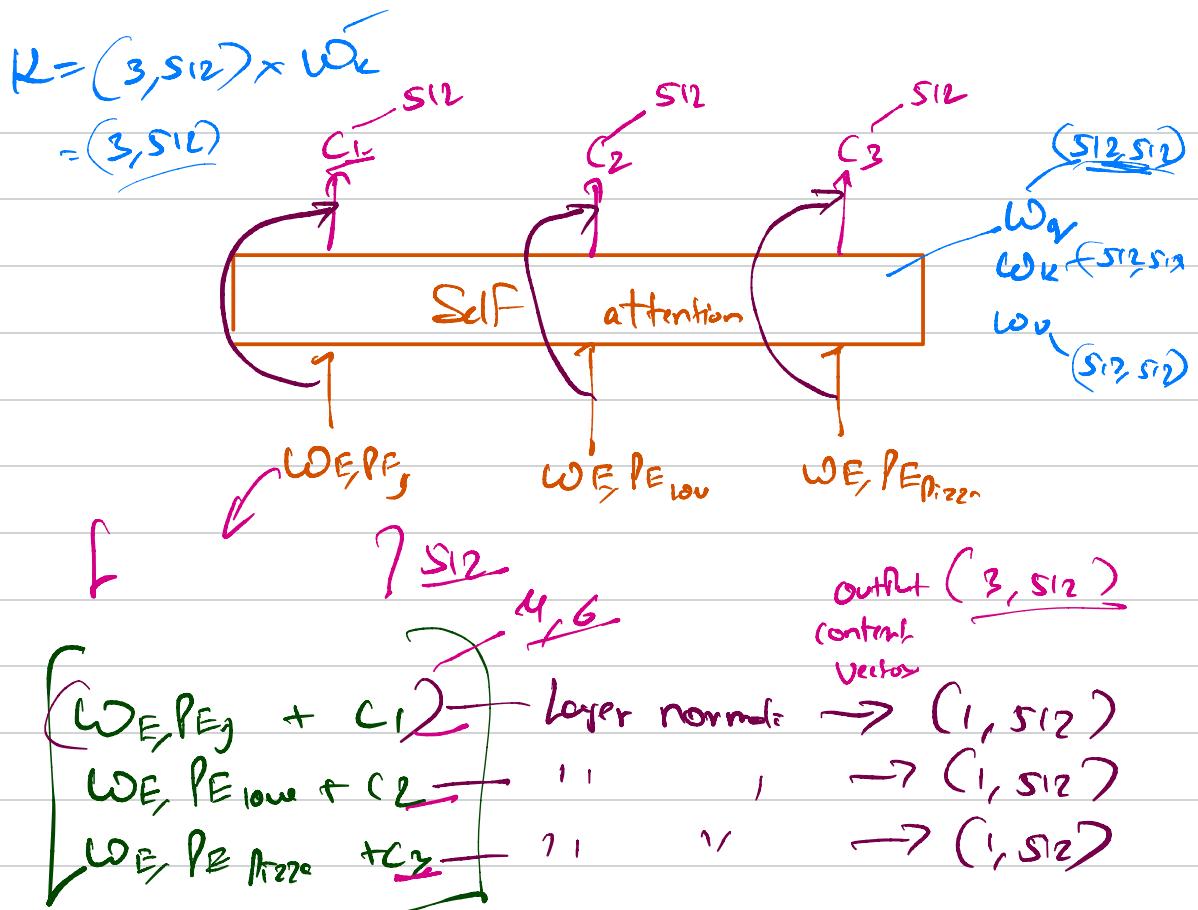


$$\textcircled{1} \quad (3, 512)$$

Diagram illustrating the word embedding calculation:

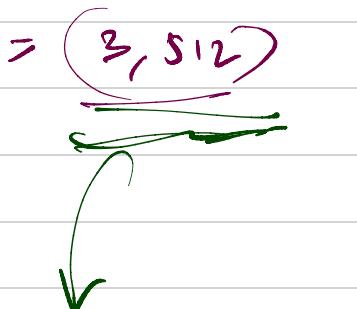
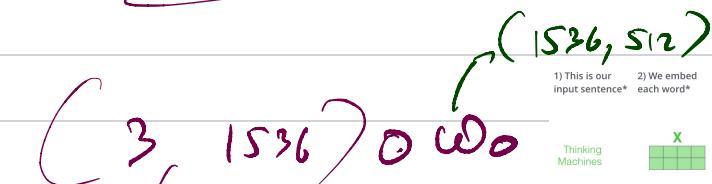
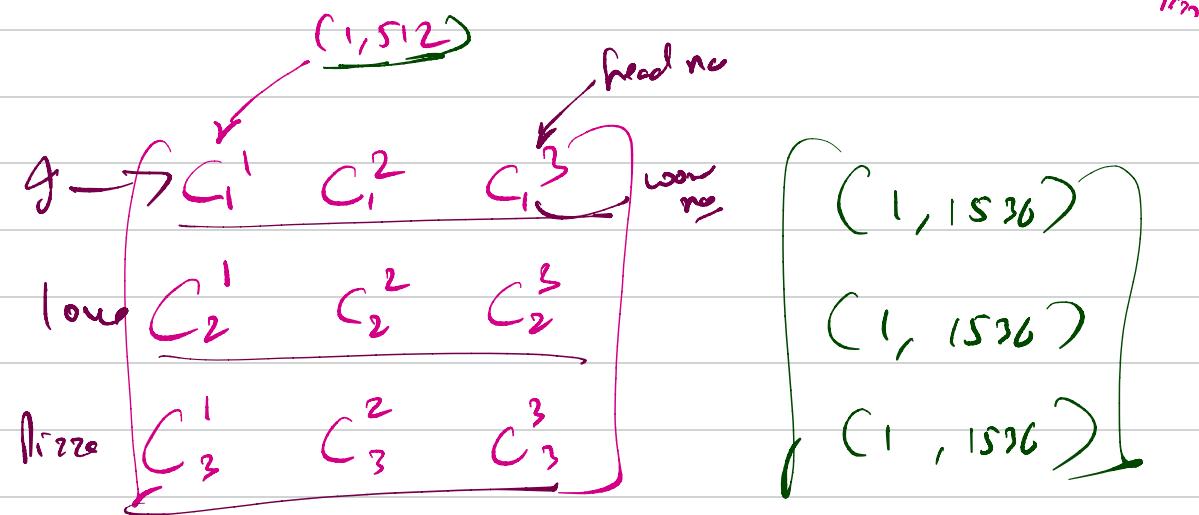
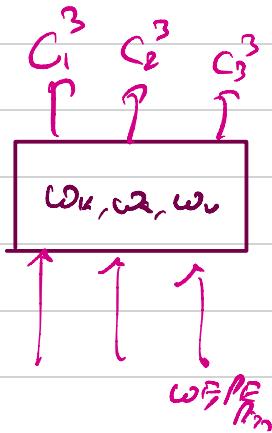
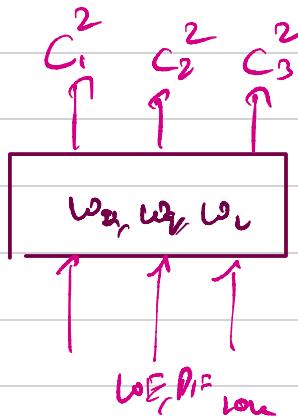
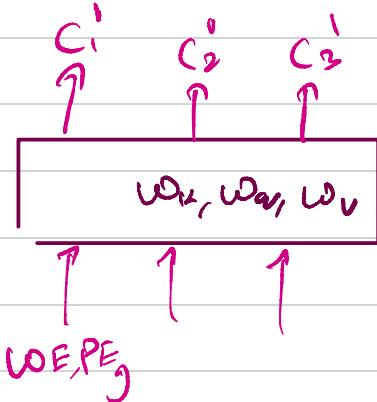
$$\left[\begin{array}{c} \text{Word}_1 \\ \text{Word}_2 \\ \text{Word}_3 \end{array} \right] + \left[\begin{array}{c} \text{PE}_{w_1} \\ \text{PE}_{w_2} \\ \text{PE}_{w_3} \end{array} \right] = \left[\begin{array}{c} \omega_E, \text{PE}_{\text{Thinking}} \\ \omega_E, \text{PE}_{\text{Machines}} \\ \omega_E, \text{PE}_{\text{Pizza}} \end{array} \right]$$

Annotations: $S12$ above the first column, $S12$ above the second column, and $(3, S12)$ below the result.



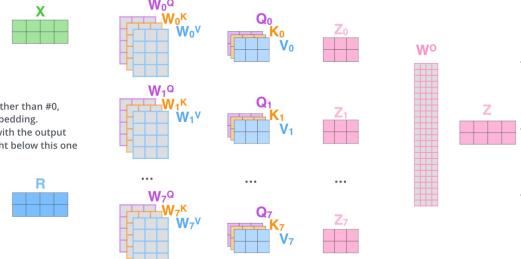
$$\begin{cases}
 \frac{1-2}{G_1}, \frac{2-2}{G_2}, \frac{3-2}{G_2} \\
 \frac{4-5}{G_2}, \frac{5-5}{G_2}, \frac{6-5}{G_2} \\
 \frac{7-8}{G_3}, \frac{8-8}{G_7}, \frac{9-8}{G_7}
 \end{cases}$$

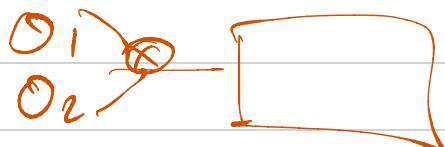
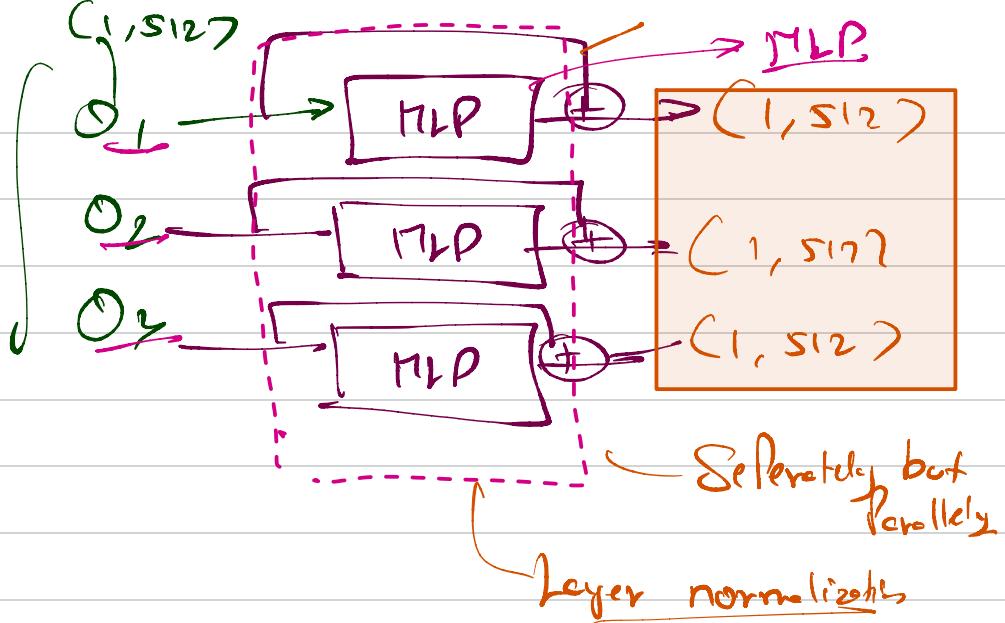
Assume head = 3



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

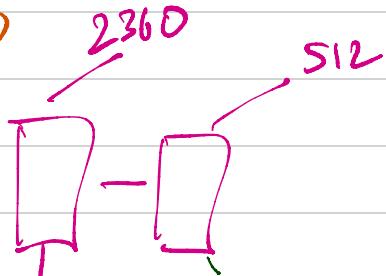
* In all encoders other than #0, we don't need embedding.
We start directly with the output of the encoder right below this one





$$\text{Output} = (1, \text{SIZ})$$

$$(2, \text{SIZ}) \times (512, 2360)$$



$$O_1 \rightarrow (1, \text{SIZ})$$

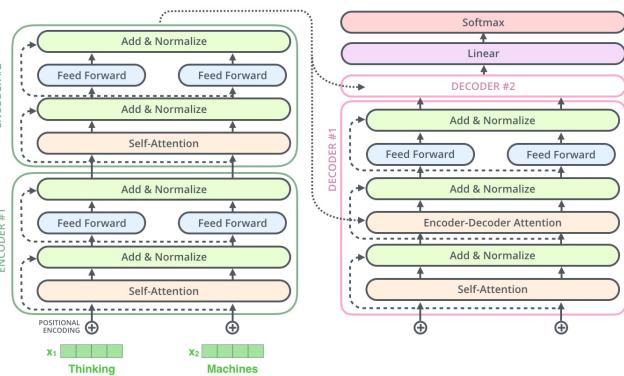
$$(512, 2360) \xrightarrow{\quad} (2360, \text{SIZ})$$

$$= 1, 2360$$

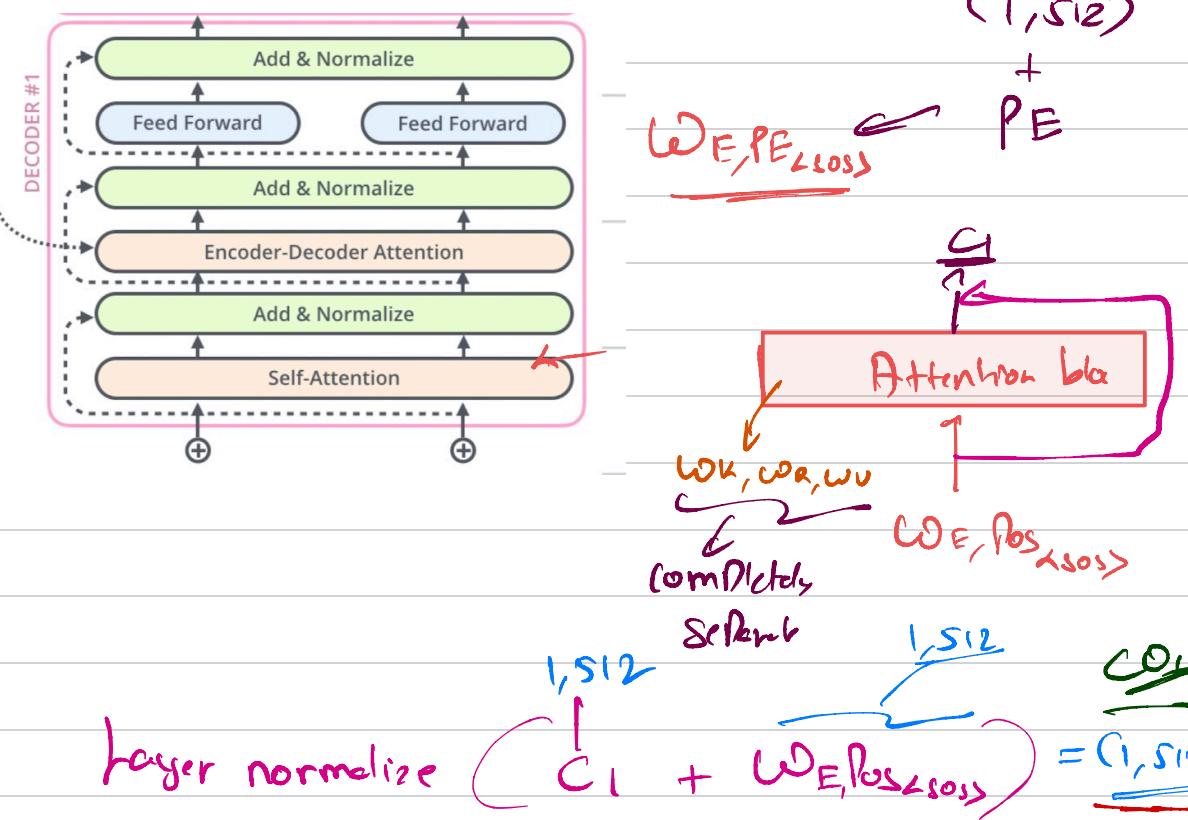
$$\rightarrow (1, \text{SIZ})$$

Decoder - Inference \rightarrow

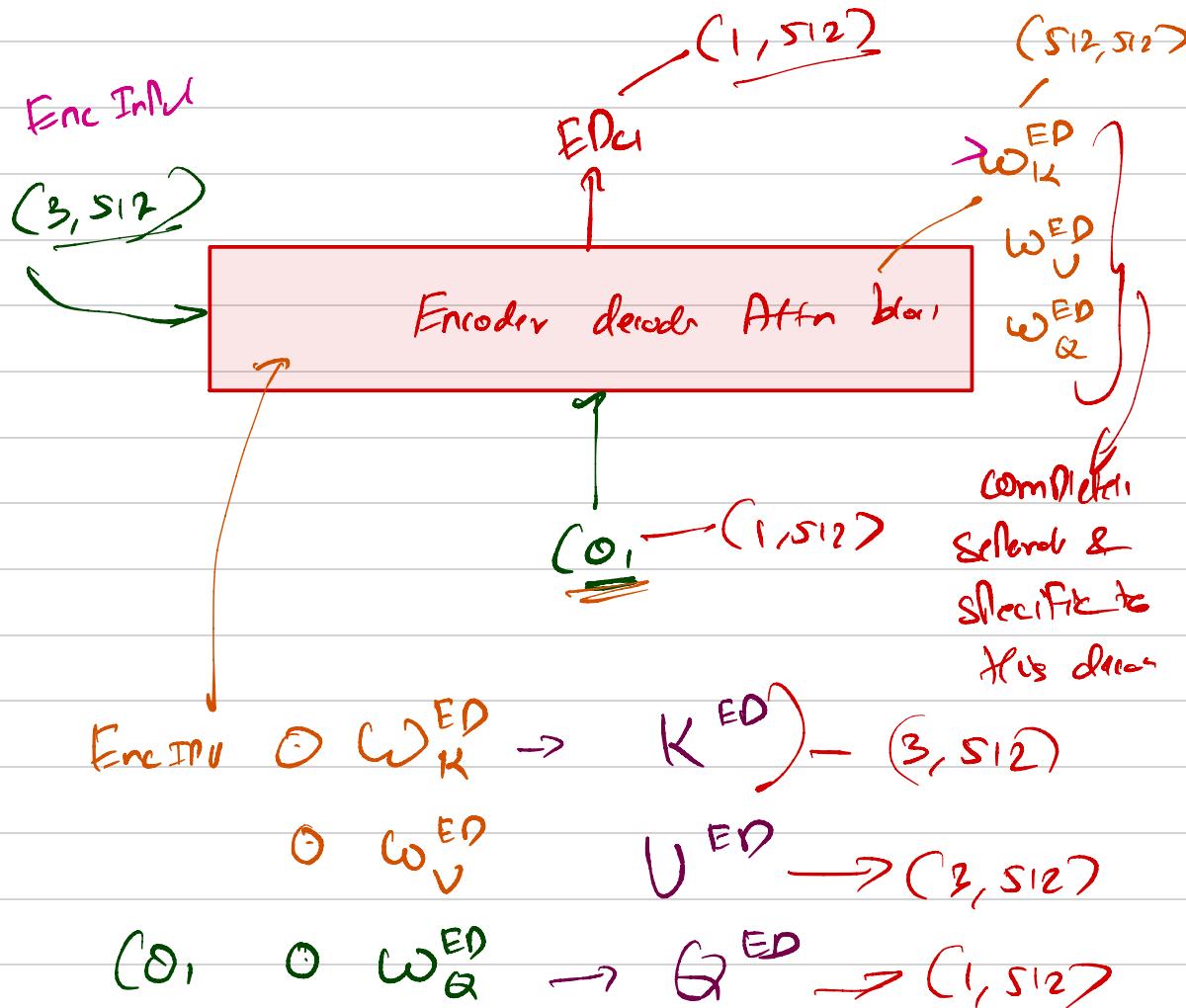
SINCE ITEND



At $t=1$ $\text{<SOS>} \rightarrow$ Word Embedding



Let's talk about encoder decoder attention block.

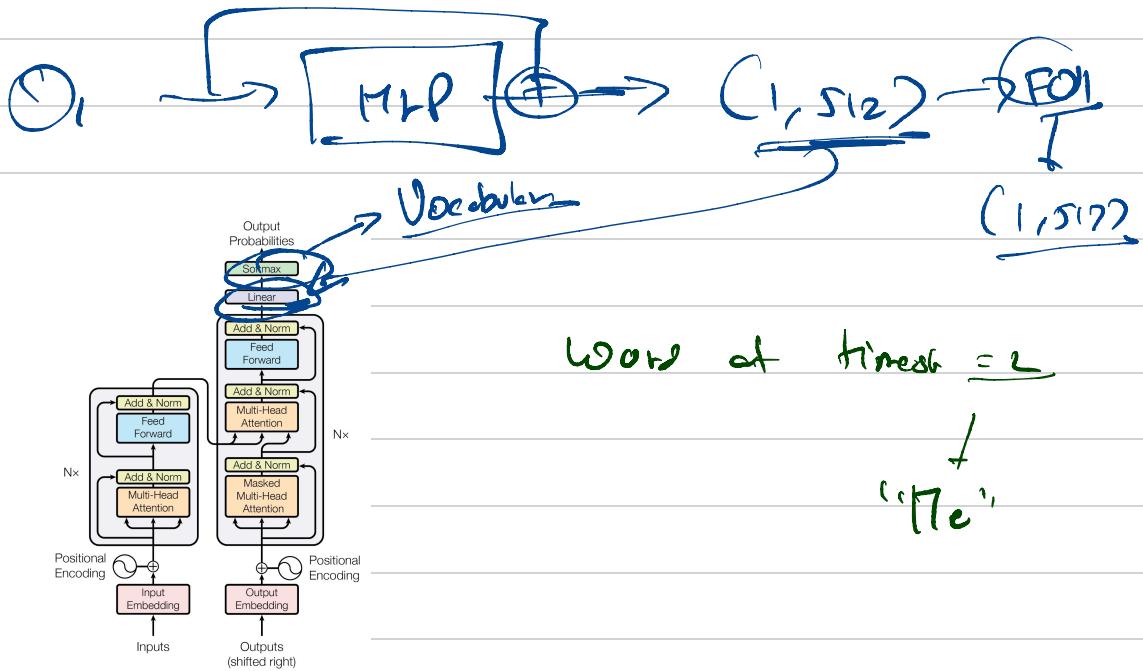


What we will do, is compute similarity of each word of query (which for us is just one word) with respect to each word in key-matrix

Compute similarity, which will be 3 numbers, corresponding to how similar each query word is wrt to each key word,

then $s_1 * v_1 + s_2 * v_2 + s_3 * v_3$; where s_1, s_2 and s_3 are softmax out.

Layer Norm ($ED_{cl} + CO_1$) $\rightarrow \underline{O_1}$



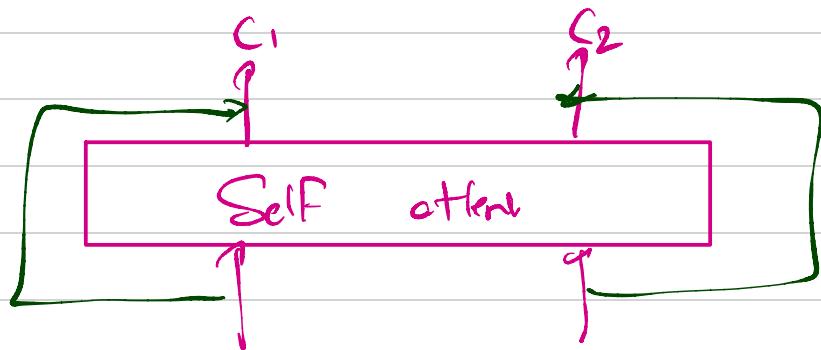
At timestep 2: $\text{<SOS>} \& \text{"Me"}$

Word \downarrow
 $\omega_E + \hat{\rho}_E$

$\omega_E, \rho_E_{\text{SOS}}$

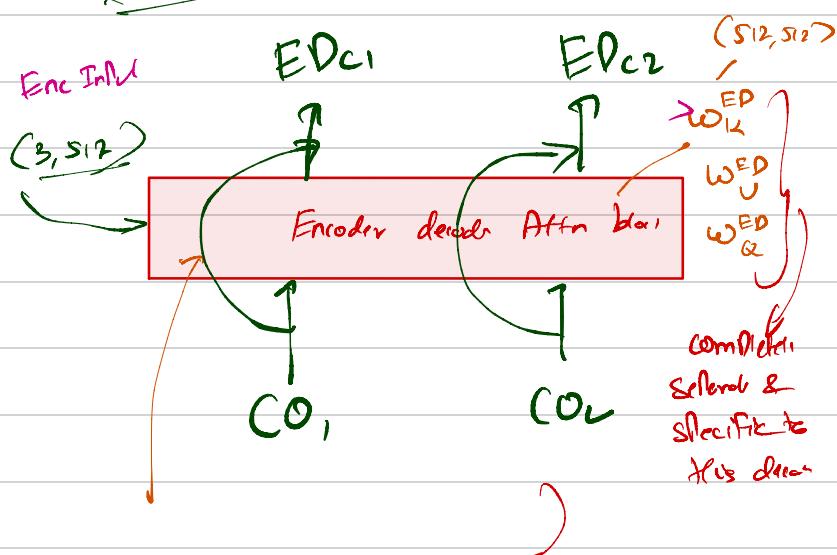
word \downarrow
 $\hat{\rho}_E$

$\omega_E, \rho_E_{\text{Me}}$

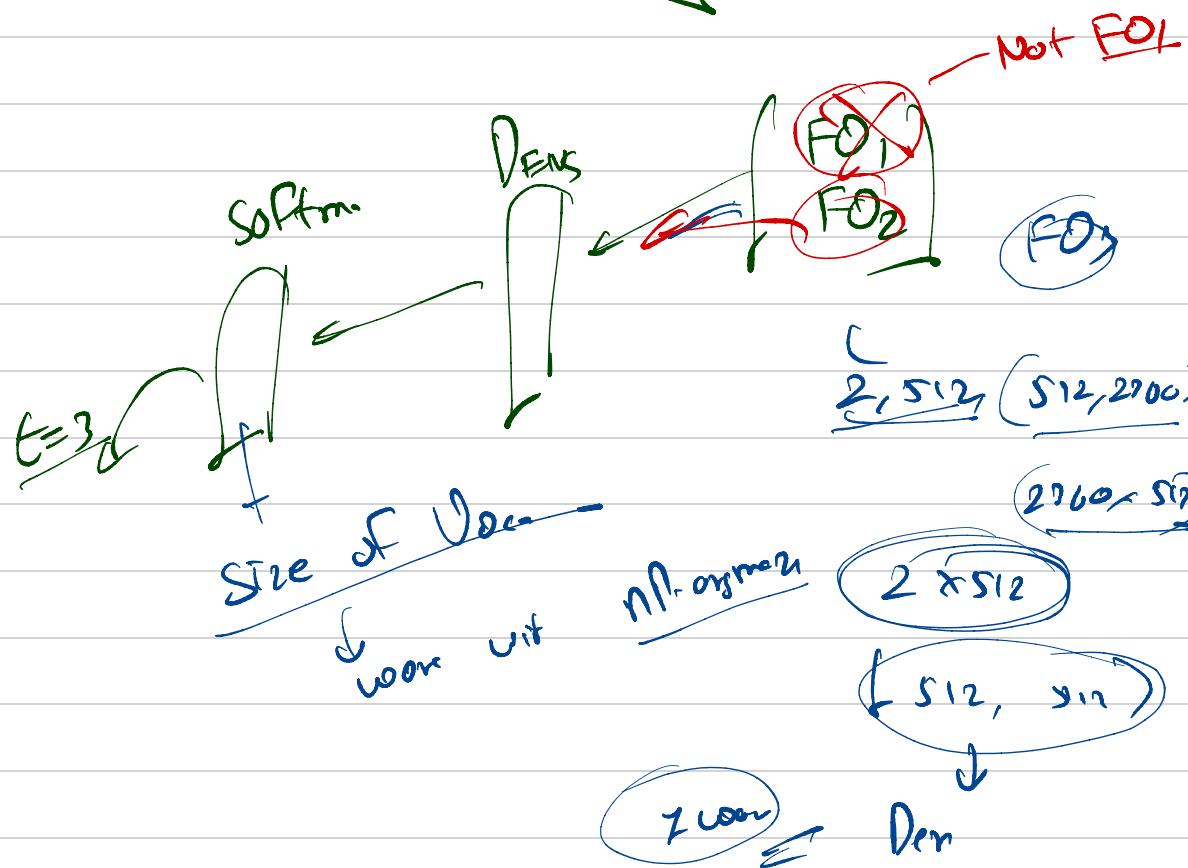
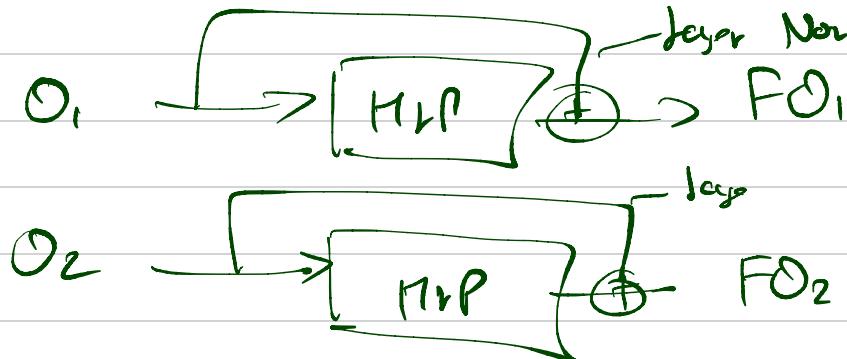


$$\text{Layer norm} \left(C_1 + \text{WE, PE, loss} \atop C_2 + \text{WE, PE, m} \right) = \begin{pmatrix} CO_1 \\ CO_2 \end{pmatrix} \quad (1, 512)$$

Diagram showing the result of layer normalization on the concatenated outputs of the multi-head attention mechanism. The result is a vector of shape $(1, 512)$, represented as $\begin{pmatrix} CO_1 \\ CO_2 \end{pmatrix}$.



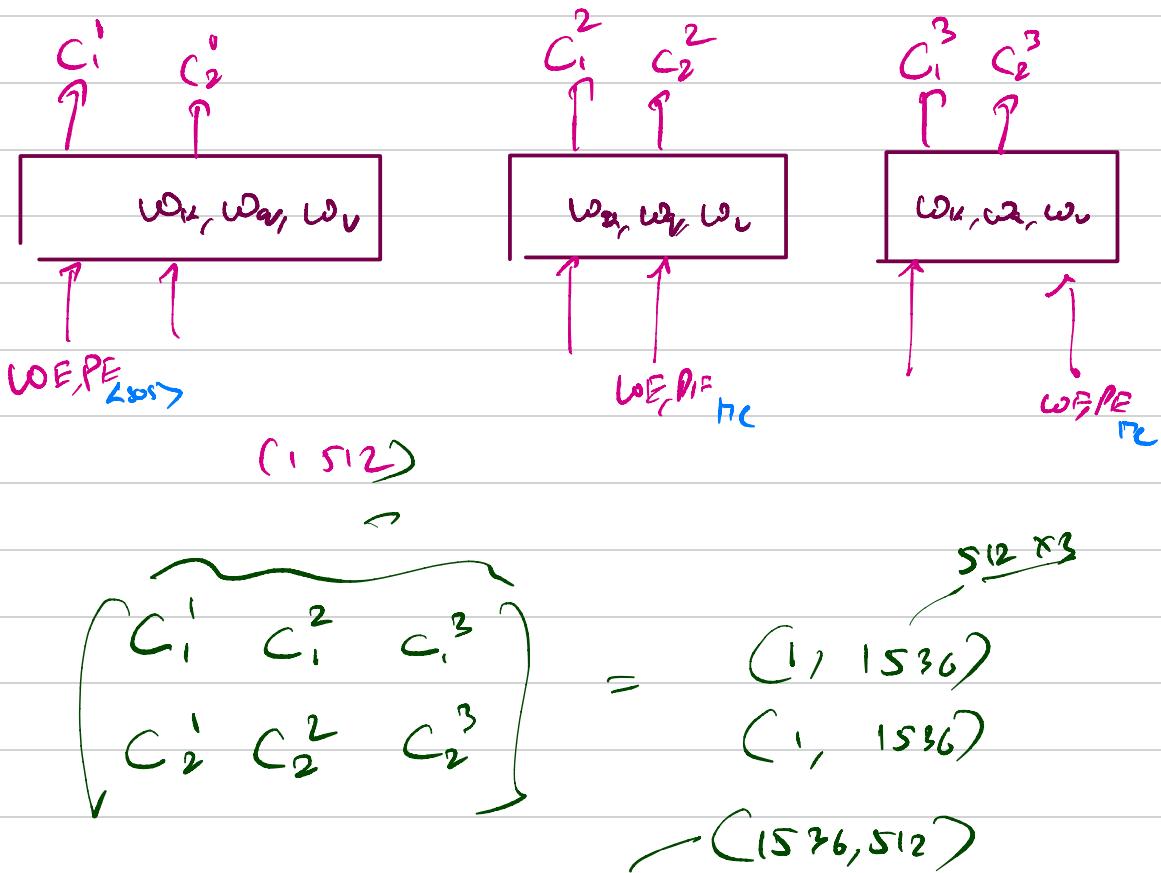
$$\text{Layer Norm} \quad \left(\begin{matrix} C_1 + \text{EDC}_1 \\ C_2 + \text{EDC}_2 \end{matrix} \right) = \left(\begin{matrix} O_1 \\ O_2 \end{matrix} \right)$$



Multi head decoder

$t=2$, L sos , "He"

Assume head = 3



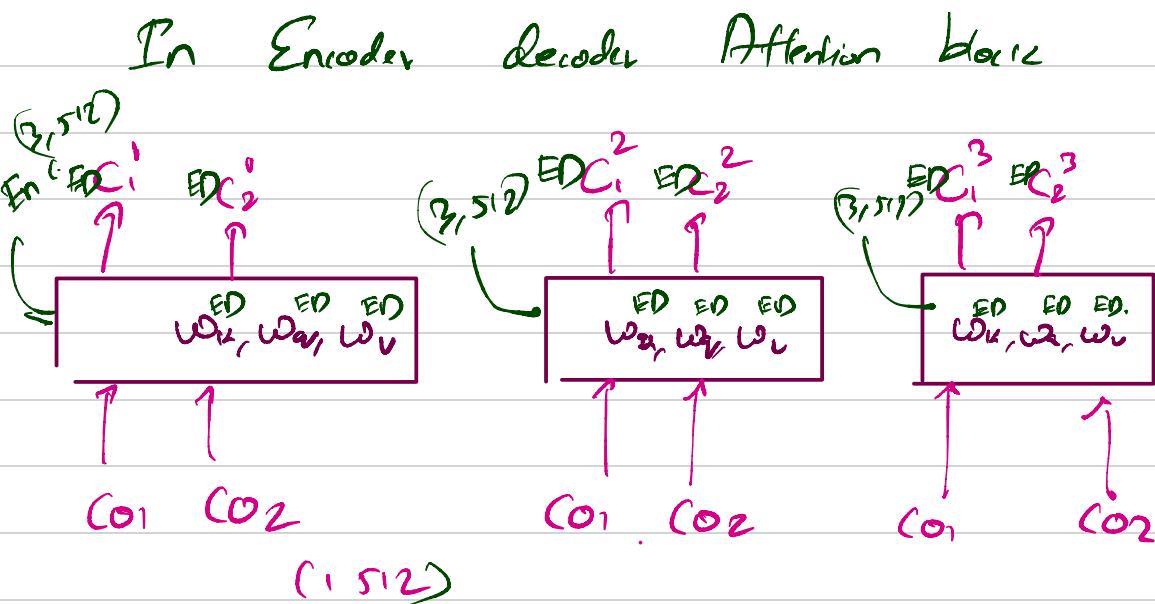
$$(2, 1536) \times W_0^D = \frac{1, 512 - CC_1}{1, 512 - CC_2}$$

$$\text{LayerNorm} \left(\frac{\underline{CC}_1 + \underline{\omega_{E, PE_{SOS}}}}{\underline{CC}_2 + \underline{\omega_{E, PE_{nc}}}} \right) = \underline{(2, 512)}$$

(C_1)
 (C_2)

In multi head decoder, like multi head encoder, attention blocks are separate, but feed forward layers are same.

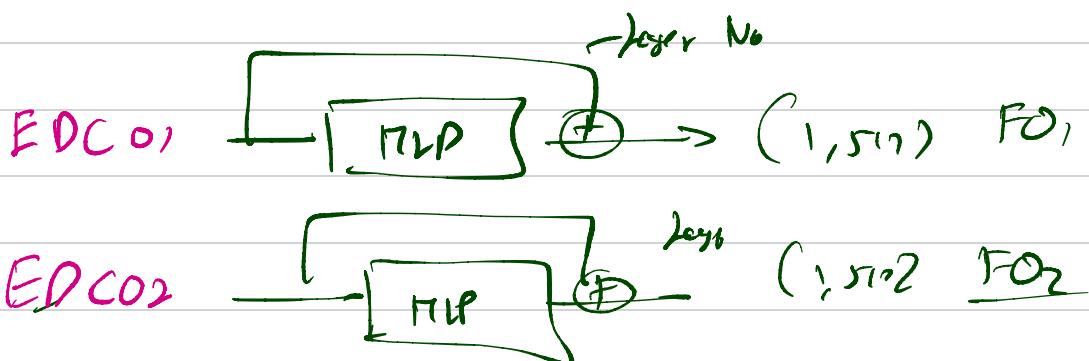
In decoder, for each head, the self attention block and encoder decoder attention blocks are separate, feed forward layers are same



$$\begin{array}{c}
 \overbrace{\text{EDC}_1^1 \text{ EDC}_1^2 \text{ EDC}_1^3}^{\text{S12}} = (1, 1536) \\
 \text{EDC}_2^1 \text{ EDC}_2^2 \text{ EDC}_2^3 = (1, 1536) \\
 (1, 1536) \times \omega_0^{\text{ED}} = \underbrace{(1, 1536)}_{(1536, 512)} \xrightarrow{\text{S12} \rightarrow \text{EDCC}_1} \\
 \xrightarrow{1, 512 \rightarrow \text{EDCC}_2}
 \end{array}$$

Layer Norm

$$\begin{aligned}
 & (\text{EDCC}_1 + \text{CO}_1) \rightarrow \text{EDCO}_1 \\
 & (\text{EDCC}_2 + \text{CO}_2) \rightarrow \text{EDCO}_2
 \end{aligned}$$



Until you find LEOSS

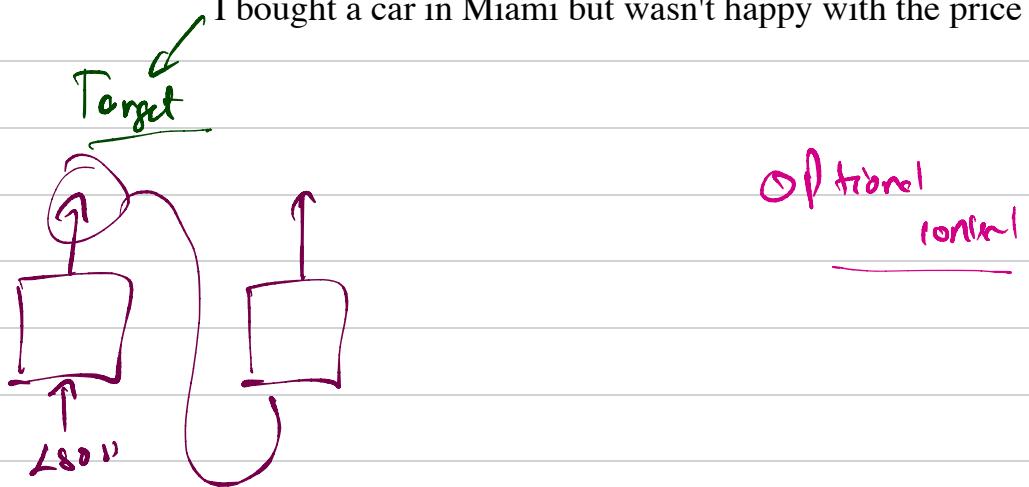
LSTM + Attention → Scale it up
PE

How training is done

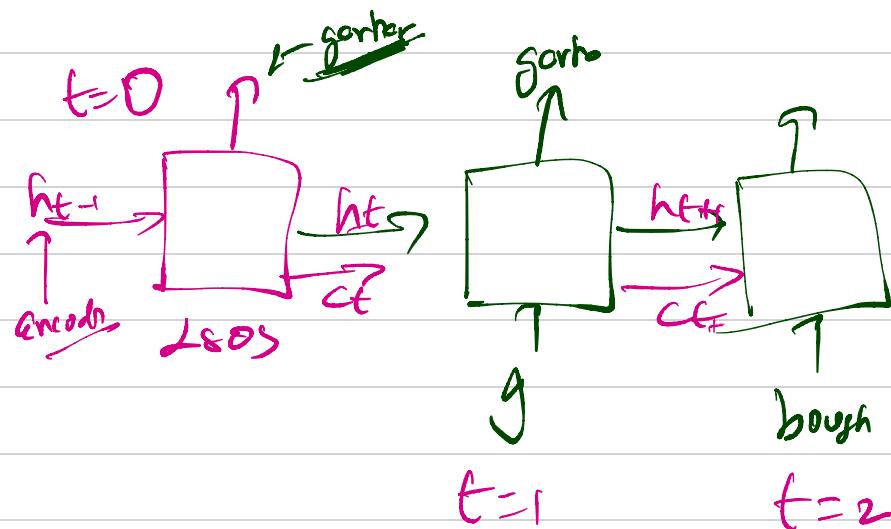
Input

Ich habe in Miami ein Auto gekauft, war aber mit dem Preis nicht zufrieden

I bought a car in Miami but wasn't happy with the price



Teacher Forcing \rightarrow LSTM



TRANSFORMER

Input
Ich habe in Miami ein Auto gekauft, war aber mit dem Preis nicht zufrieden

I bought a car in Miami but wasn't happy with the price

Target
○

Optional

Input \rightarrow 14 word_dim = 512

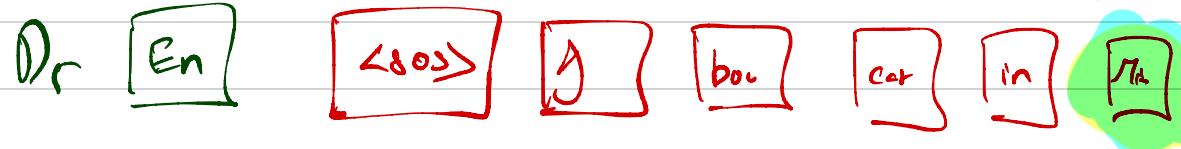
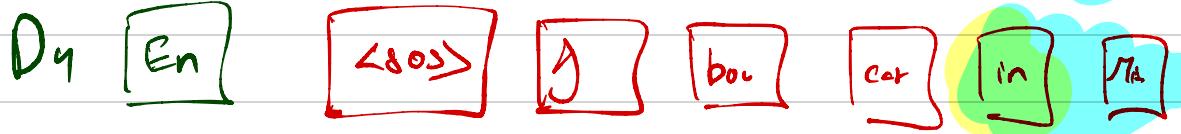
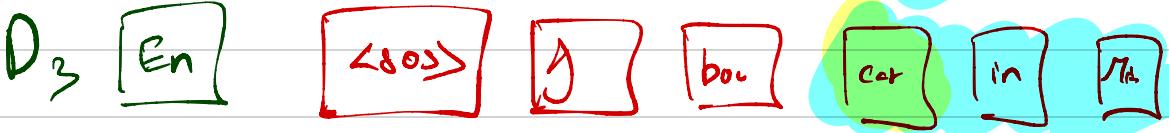
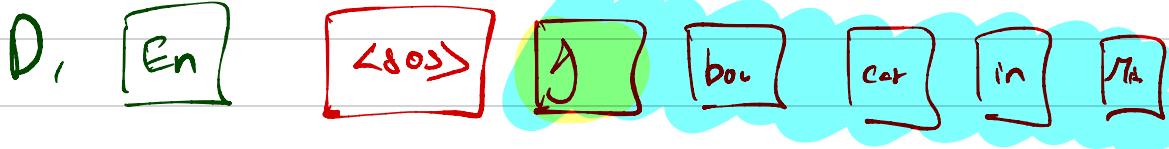
14, 512

Mod/L self alter

Enc-outfit

↑
(u, s(2))

Scoring at my decou



D₁

D --

D₅

↓
1011

↓
202

↓
102

grade W.r.t for

avg per cent

