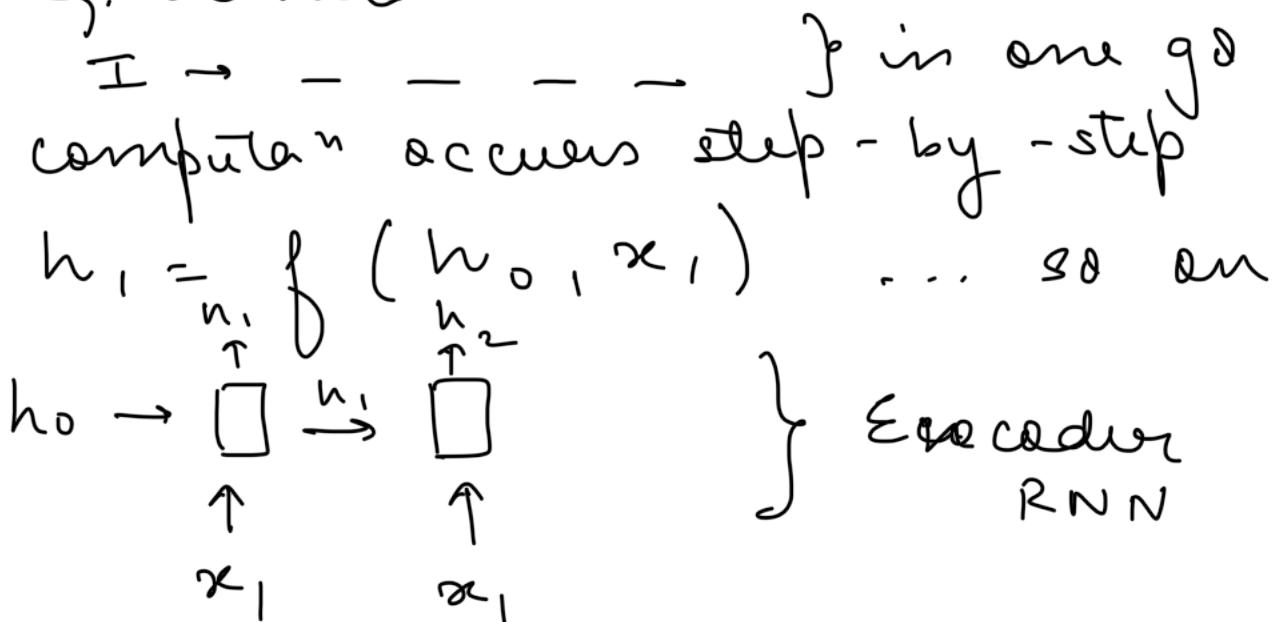


# LLM Week-1

## ↳ Introduction to Transformer Architecture

- previously RNNs were used, now transformers are more used.
- then came attention based building block.
- RNN based models

e.g. translation



why calculating  $h_2$  only after  $h_1$ ?  
→ to know from the context of the sentence. Contextual representation of every word is imp.

→ cons → can't do parallel processing.  
wasting of resources.

I want the bidirectional, not recurrent in nature.

Decoder RNN  $\rightarrow$  also gives the sequential output.

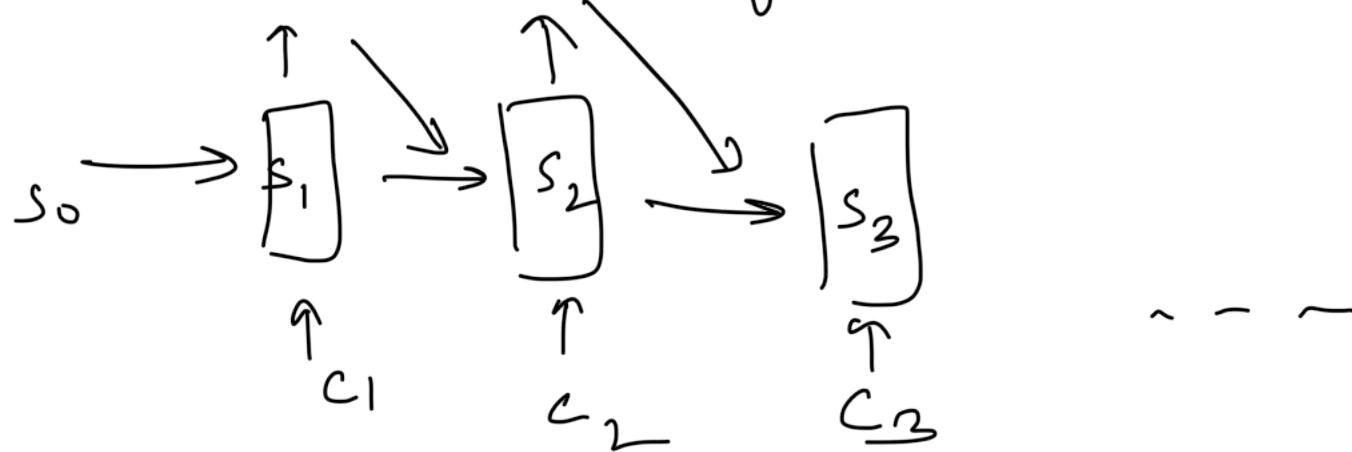
We can speed up the encoder part.

- The final state is called the a concept vector, thought vector, or content vector see annotation.
- However, we still don't care about the alignment of words b/w source sentence & target sentence. Therefore, we will be needing a attention mechanism to help us resolve this issue.
- I compute the RNN encoding. I have the vectors  $h_1, \dots$  with me. We can drop the encoder block.

$$c_t = \sum_{i=1}^n \alpha_{t;i} h_i ; \text{ context vector for output}$$

$n \Rightarrow$  no. of words

$t \Rightarrow$  time step for decoder



- we generally look at a heatmap, when dealing with transformers.

$$\begin{aligned}\alpha_{t_i} &= \text{align}(y_t, h_i) \\ &= \frac{\exp(\text{score}(s_{t-1}, h_i))}{\sum_{i=1}^n \exp(\text{score}(s_{t-1}, h_i))}\end{aligned}$$

Ques: Can  $\alpha_{t_i}$  be computed parallelly for all  $i$  at time step  $t$ ?

Yes  $\Rightarrow h_i$  is available for all  $i$  &  $s_{t-1}$  is also available at time step  $t$ .

Attention can be parallelized

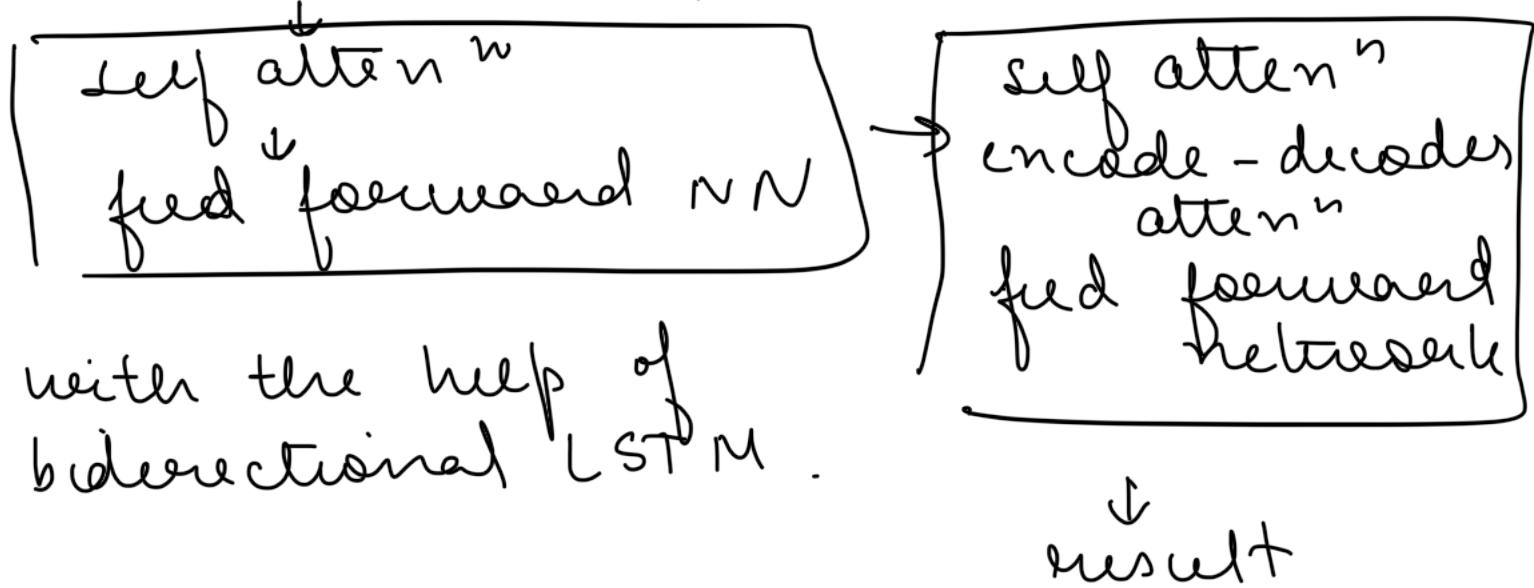
Ques: Can  $h_i$  for all  $i$  can be computed parallelly?  $\Rightarrow$  No

- We need parallel mechanism at same time to get rid of vanishing gradient problem.

L2 Attention is all you need

- We have recurrent connection  $\Rightarrow$  the RNN + attention model.

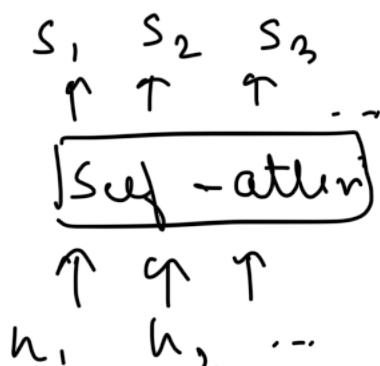
- transformer architecture  
word embedding:



- with the help of  
bidirectional LSTM.

- self atten<sup>n</sup>  
takes word embeddings and creates  
the  $h_i$ 's. The atten<sup>n</sup> requires a  
pair of vectors as input.

$$\alpha_{t-i} = \text{ATT}(s_{t-1}, h_i)$$



$$s_2 = \sum_{j=1}^{T=5} \alpha_{2j} h_j$$

as you can see there  
is no recurrence here.

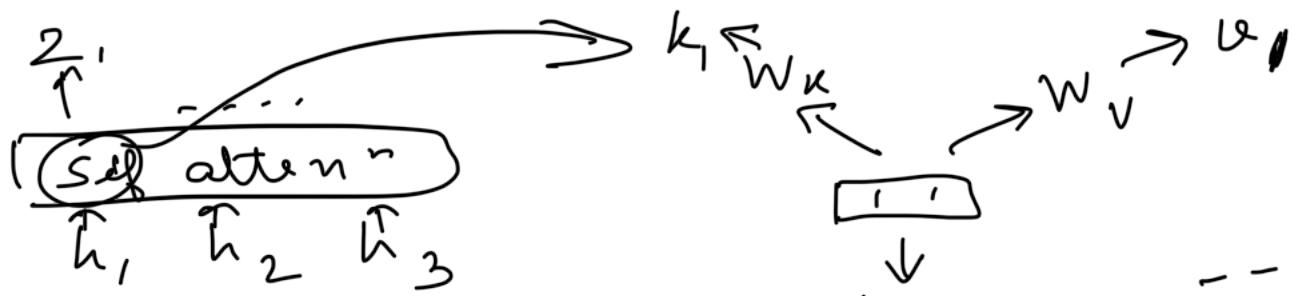
- atten<sup>n</sup> is needed because "it" should  
refer to more atten<sup>n</sup> to the  
contextual word.

⇒ Self-Atten<sup>n</sup> (Atten<sup>n</sup> mechanism)

- given a word, we want to compute the relational score b/w the word & rest of the words in the sentence. It will be a  $T \times T$  matrix.
- Now both the vectors ( $s_i$ ) & ( $h_j$ ) for all  $i, j$  are available for all the time. (whereas it was not so in the seq2seq model.)
- Score func<sup>n</sup> of seq2seq  
 $\text{score}(s_{t-1}, h_j) = W_a^T \tanh (U_{\text{att}} s_{t-1} + W_{\text{att}} h_j)$   
 3 vectors  $\rightarrow (s, h, v)$ , 2 linear linear transformation, 1 non-linearity finally the dot-product.  
 However the input the self-attention module is only the word embedding  $h_j$  for all  $j$ .  
 We need matrix transformation.

### ↳ Self-Attention

- transformation of matrices  
 $q_f = R^d = R^{d \times d} h_j \cdot v_j, k_j$   
 $(W_Q)$



score ( $s_{t-1}, h_j$ )  
score ( $q_1, k_j$ )

variables  
finid

func" is a dot product

$$\alpha_{1,j} = \text{softmax}(e_{1,j})$$

$$e_1 = [q_1 \cdot k_1, q_1 \cdot k_2, \dots, q_1 \cdot k_5]$$

$$z_1 = \sum_{j=1}^5 \alpha_{1,j} v_j$$

- Can we vectorize all the computation and find the outputs in one go.

$$\begin{bmatrix} q_1 & q_2 & \dots & q_r \end{bmatrix} = W_Q \begin{bmatrix} h_1 & h_2 & \dots & h_r \end{bmatrix}$$

$$Q \in \mathbb{R}^{64 \times T}$$

$\downarrow$   
 $d_1 \times T$

$W(64) \times \underline{64}$   
 $d_1 \times d$

$[h_T]$   
 $\frac{64}{d} \times (T)$   
 $d \times T$

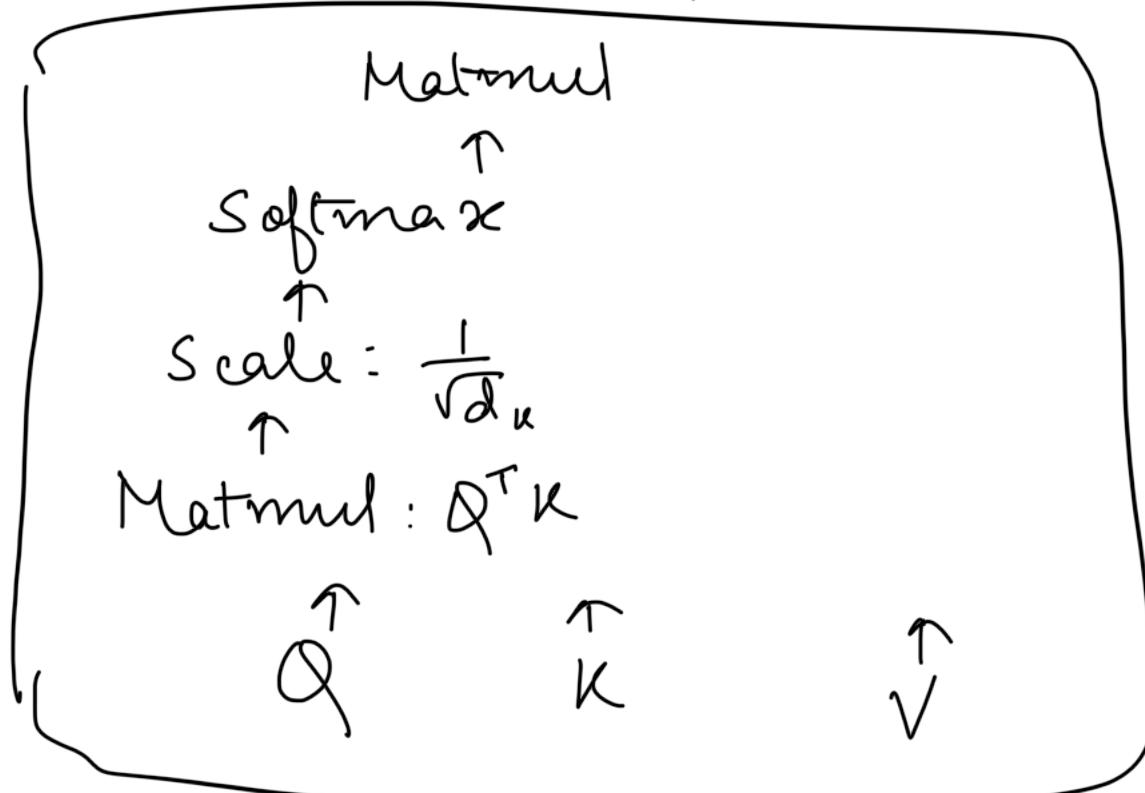
similarly it goes for K and V.

$$\begin{aligned}
 Z &= [\gamma_1, \gamma_2, \dots] \\
 &= \text{softmax} \left( \frac{Q^T K}{\sqrt{d_K}} \right) V^T
 \end{aligned}
 \quad [ ]_{T \times T}$$

$$\begin{aligned}
 Q &\rightarrow T \times 64 \\
 K &\rightarrow 64 \times T \\
 V &\rightarrow d \times T
 \end{aligned}
 \quad \left. \begin{array}{l} \text{atten"} \text{ matrix} \\ \vdots \end{array} \right\}_{T \times d}$$

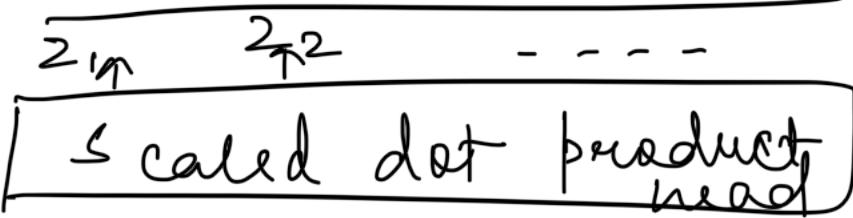
$d_K$  is the dim of key vector

$d_K$  scales the values of  $Q^T K$ .  
 ↳ scaled dot product!



→ self atten" block

## L4 Multi Headed Atten<sup>n</sup>



$$\begin{array}{ccc} Q & K & V \\ \uparrow & \uparrow & \uparrow \\ w_Q T & w_K T & w_V T \end{array}$$

$$H = \{ h_1, h_2, \dots, h_T \}$$

they take  
care of  
contextual  
info.

this is one 'block', we can have 2 or multiple head blocks.

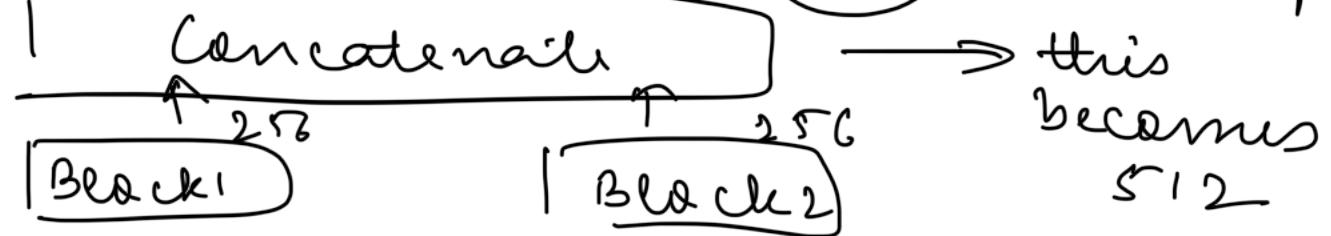
But what is the significance of multiple blocks?

Just like CNN, we have multiple filters that would capture diff. input points.

Similarly, if there are atten<sup>n</sup> blocks, first will capture 'it' with 'animal' and the second will capture the info. in 'was'.

In the same way, diff. atten<sup>n</sup> block we calculate  $\alpha_s$ , the block will

learn more contextual info.

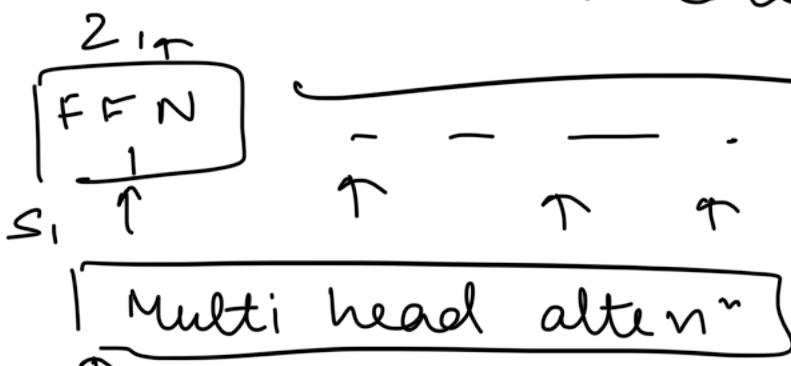


this will be the output.

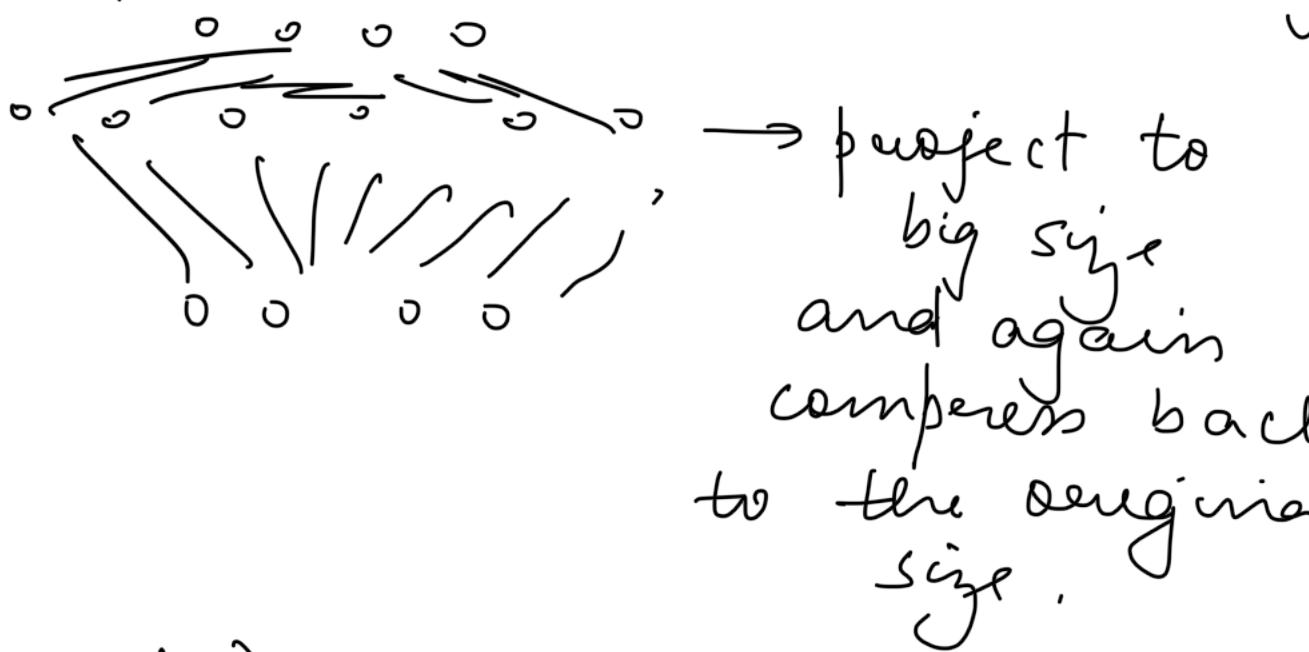
this becomes 512

$^256$   
 $Q$   
 $K$   
 $V$   
 $512 \times 256$   
 $W_Q$   $W_K$   $W_V$   
 $512, \dots$

if there are  $\sigma$  blocks then they will 64 size, then again concatenates gives 512.



acts as a project layer



$$FFN(z) = \max(0, W_1 z + b_1) W_2 + b_2$$

FFN  
multi head

→ one encoder  
block, there will  
many blocks.

- Please note, parallelism is horizontally and not vertically. So basically in layer 1 of encoder that computation will occur parallelly. But we will waiting for the step-by-step process layer after layer.