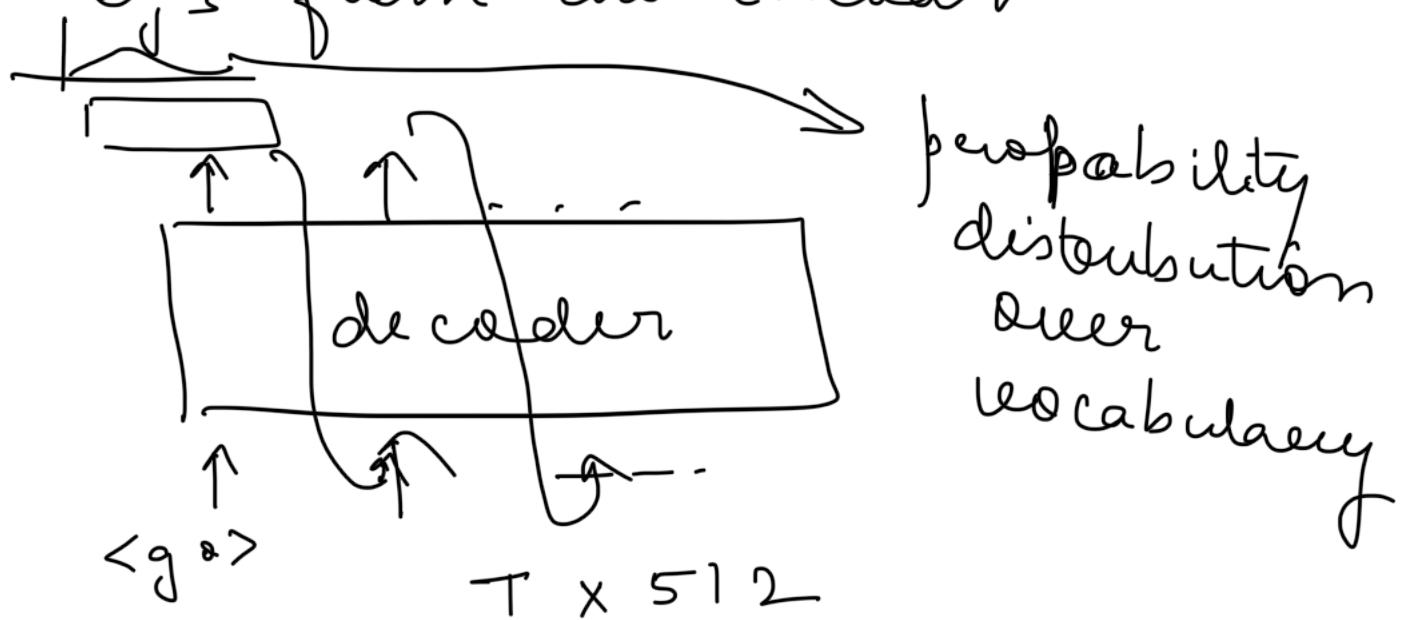


# LLM

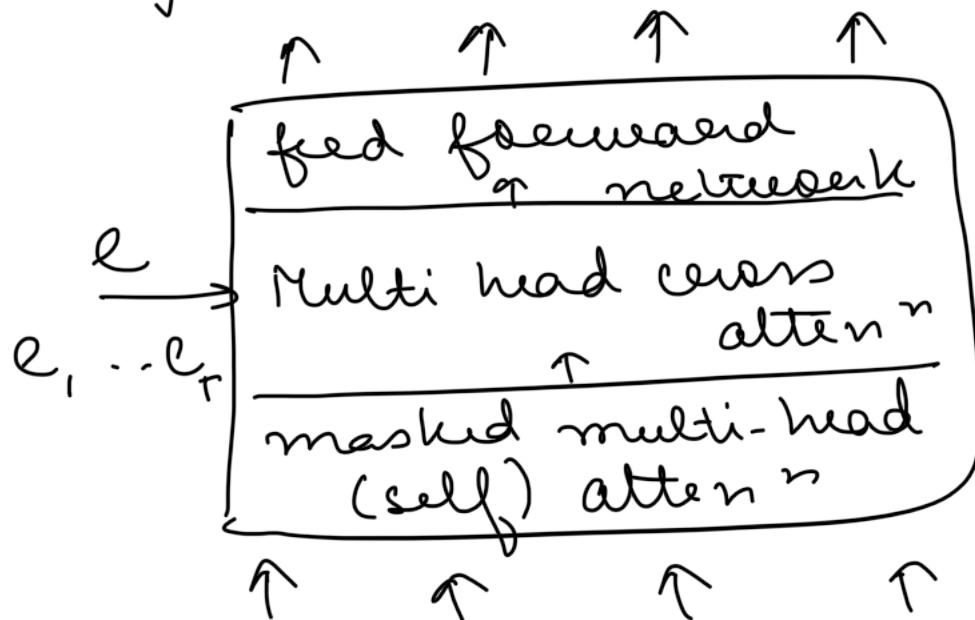
## Week-2

### 1 Teacher Forcing & Masked Atten<sup>n</sup>

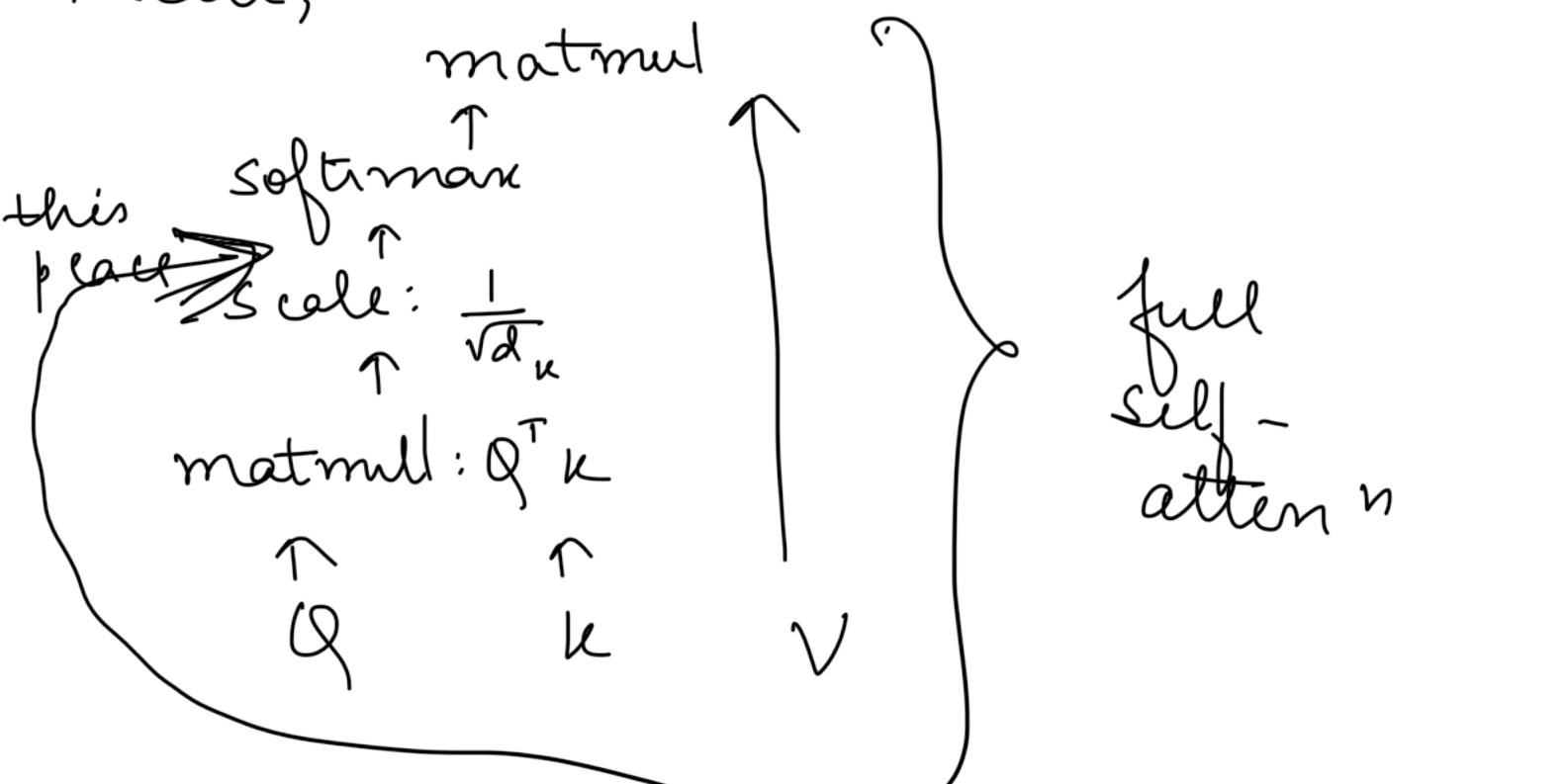
- decoder → we had received the  $e_{i,s}$  from the encoder



- again decoder is a stack of  $N$  layers. Each layer has 3 sub-layers.



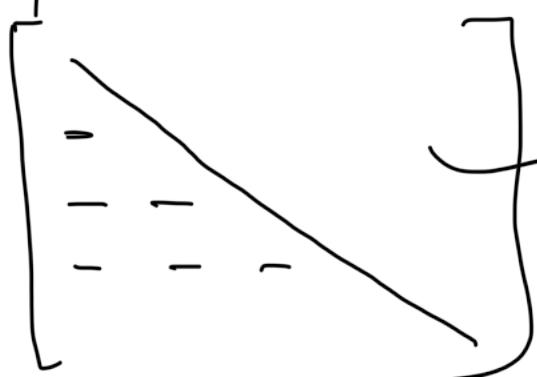
- if you see there is problem, ~~there~~ is a possibility that your first word itself is wrong. Then further, thing will be surely garbage.
- Again, you have both eng. & tamil labelled data, feed it. Then you don't need any encoder.
- feed actual outputs (teacher forcing)
- do not feed everything, do masking on other objects.
- Recall,



"if you want  
masking"

- masking to implement the teacher forcing approach during training.

Remember, never use teacher forcing during inference. Then, the decoder directly acts as a auto-regressor.



it does masking on the other words.

$$Q_i = W_{Q_i} H, K_i = W_{K_i} H, V_i = W_{V_i} H$$

$$Q_i^T K = A$$

Assign 0 wts.  $a_{ij} = 0$  for the value vectors  $v_j$  to be masked in seq.

$$\begin{bmatrix} -A \\ - \end{bmatrix}_{T \times T} \Rightarrow \begin{bmatrix} -\infty \\ 0, M \end{bmatrix}_{T \times T}$$

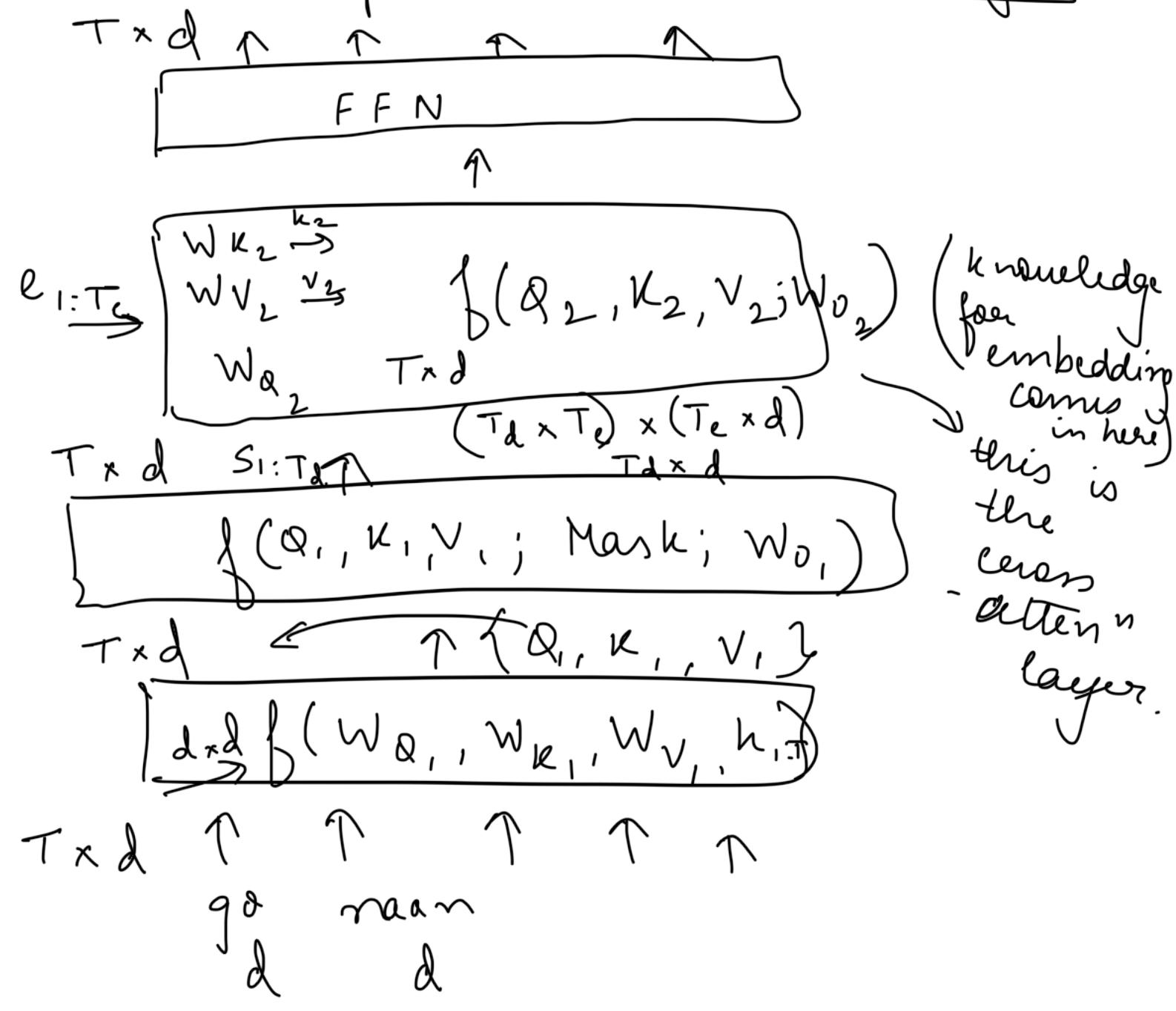
$$Z = \text{softmax}(A + M) V^T$$

$$A + M$$

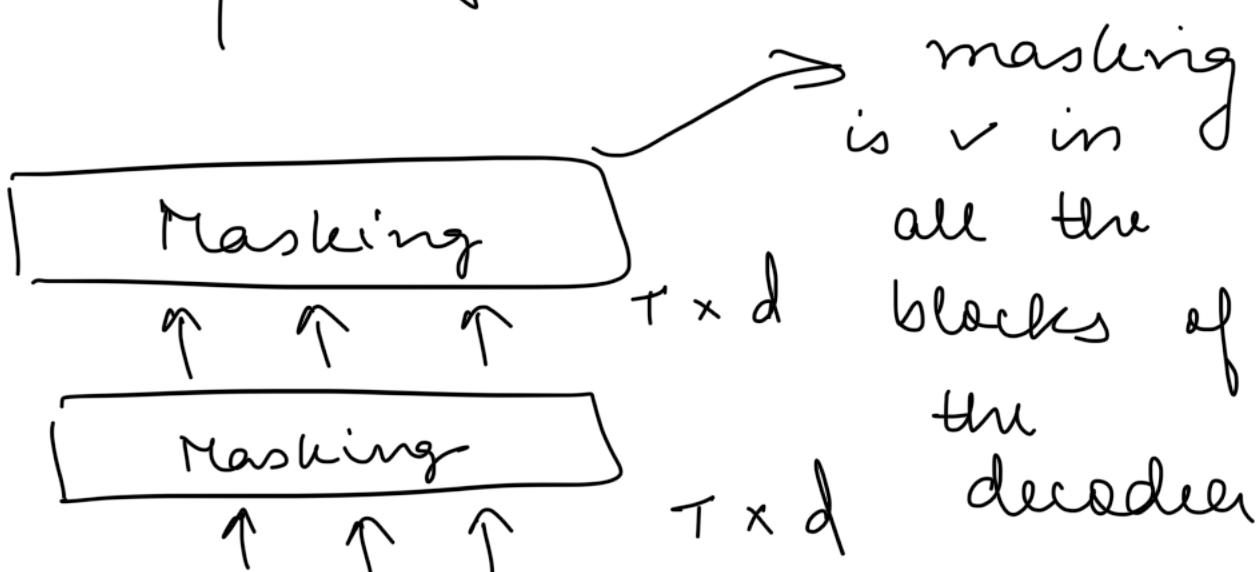
$$\downarrow$$

$$\rightarrow \begin{bmatrix} q_1 k_1 - \infty & - \infty \\ q_2 k_1, q_2 k_2 - \infty \\ \vdots & \vdots \\ \vdots & \vdots \end{bmatrix}_T$$

## L2 Zooming into the decoder layer



- similarly it goes



Softmax

calculates  
the distribu<sup>n</sup>  
over words.

Linear  $W_D$



all the  
layers of  
decoder

these  
are  $d$   
dimensional  
vectors

$$[w_D] \begin{bmatrix} \quad \\ \quad \end{bmatrix}_{\|v\| \times d} = \begin{bmatrix} \quad \\ \quad \end{bmatrix}_{d \times 1} \quad \|v\| \times 1$$

This is how the output layer  
will be. This calcula<sup>n</sup> is meant  
for linear  $W_D$  layer.

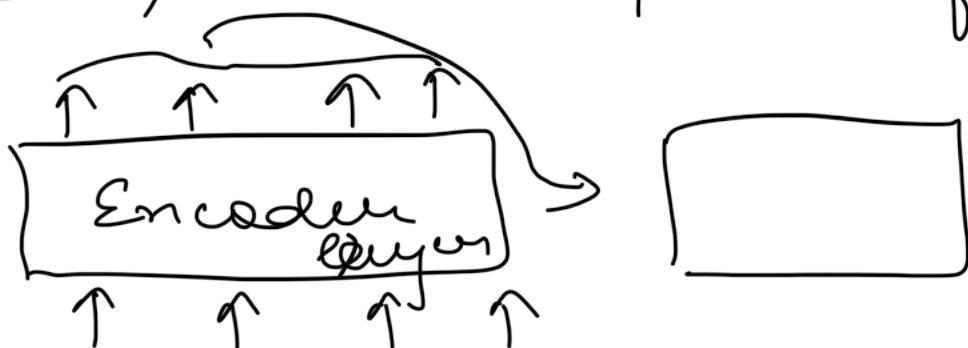
But we don't need a  $\{\|v\| \times 1\}$   
kind of a vector. We want  
probability distribu<sup>n</sup>. . . , I'll use  
softmax at the last layer.

## L3 Positional Encoding

- The pos<sup>n</sup> of words in a sentence was encoded in the hidden states of RNN-based seq-to-seq. models.

$$z_i = \sum_{i=1}^4 (q_i^\top k_i) v_i$$

- Currently, the words does not have any sense of positions. But relative position of words is needed.
- Output from the self-atten<sup>n</sup> is "permuta<sup>n</sup>" invariant. So it is necessary to encode posi<sup>n</sup> info.



- how do we do this?

$$\begin{array}{l} p_0 \\ i=0 \\ i=1 \\ \vdots \end{array} \rightarrow \begin{array}{c} h_0 \in \mathbb{R}_{5 \times 2} \\ + \\ h'_0 \in \mathbb{R}_{5 \times 2} \\ | \end{array}$$

it could be  
a constant  
vector (all  
elements are of  
constant value j)

- for  $p_j$ )
- these nos. are typically very large  
for initialization.
- Can we use one-hot encoding?
- $\hookrightarrow j = 0, 1, \dots, T$
- $\hookrightarrow$  may loose info., problem of sparse matrix.
- learn all embeddings for all possible positions?
  - these methods are not suitable for dynamic sentence lens.
  - in one hot-coding, all 1s will have a 'just same' distance of  $\sqrt{2}$ .  
But, this should not be so.

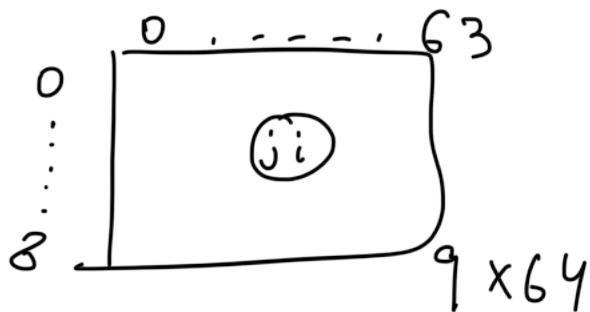
## Sinusoidal Encoding

- embed a unique pattern of features for each pos<sup>n</sup>  $j$  and the model will learn to attend by the relative pos<sup>n</sup>
- but how do we generate such features?
- $$PE(j, i) = \begin{cases} \sin\left(\frac{j}{10000^{2i/d_{\text{model}}}}\right) & i = 0 \dots 25 \\ \cos\left(\frac{j}{10000^{2(i+1)/d_{\text{model}}}}\right) & i = 0 \dots 25 \end{cases}$$

$d_{model} = \text{dim of transformer model} = 512$

- for a fixed pos<sup>n</sup>  $j$ , we use sin func<sup>n</sup> when  $i$  is even, & cos if  $i$  is odd

$$j = 0, \dots, 8 \quad i = 0, \dots, 63$$



- some obs. -

1. dist. increases on the left. & right of 0.

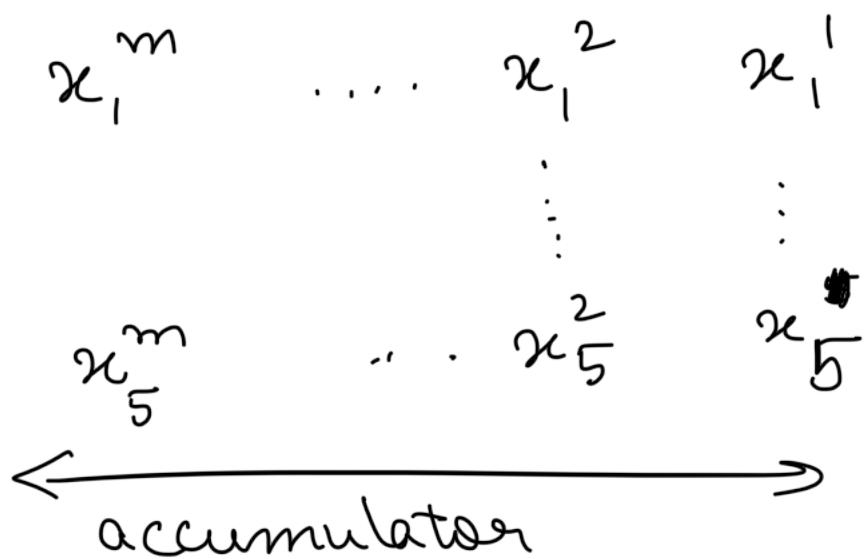
2. sym. ~~at~~ center of sentence

- PE func<sup>n</sup> solves our problem.

## L5 Batch Normaliza<sup>n</sup>

- roughly, we can say transformer architecture is composed of attention layers & hidden layers.
- in such a network, how will the gradient flow will take place?  
by using residual connections
- how to speed up?  $\rightarrow$  Normaliza<sup>n</sup>

- Batch Normalization at  $l^{th}$  layer  
 $x_i^j \rightarrow$  activation of  $i^{th}$  neuron for  $j^{th}$  training sample.  
 accumulator with  $l^{th}$  layer stores the activation of batch inputs



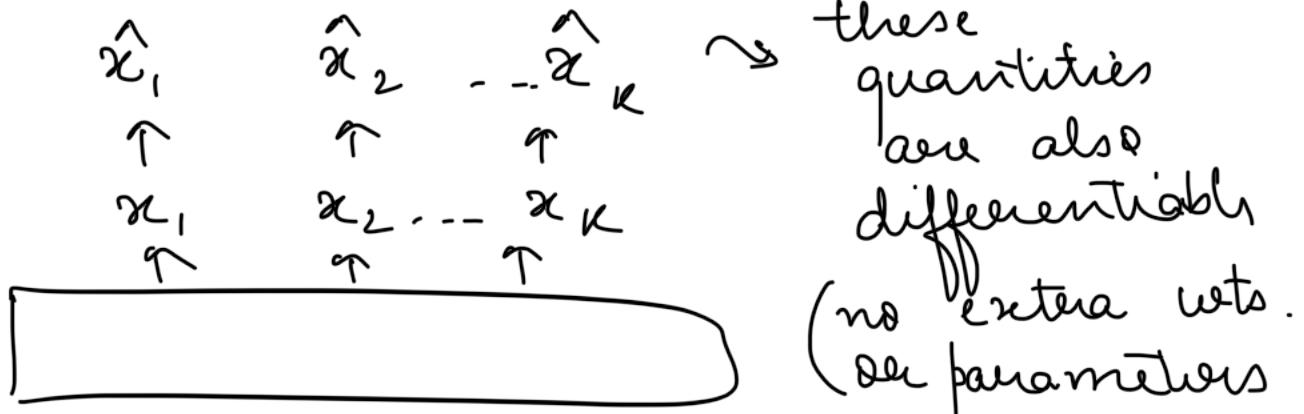
- it is generally known that we do standardization on the features while in ML or DL.

$$\mu_{i=2} = \frac{1}{m} \sum_{j=1}^m x_2^j$$

$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_2^j - \mu_i)^2$$

$$\hat{x}_i = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

} these are  
the  
standardized  
values.



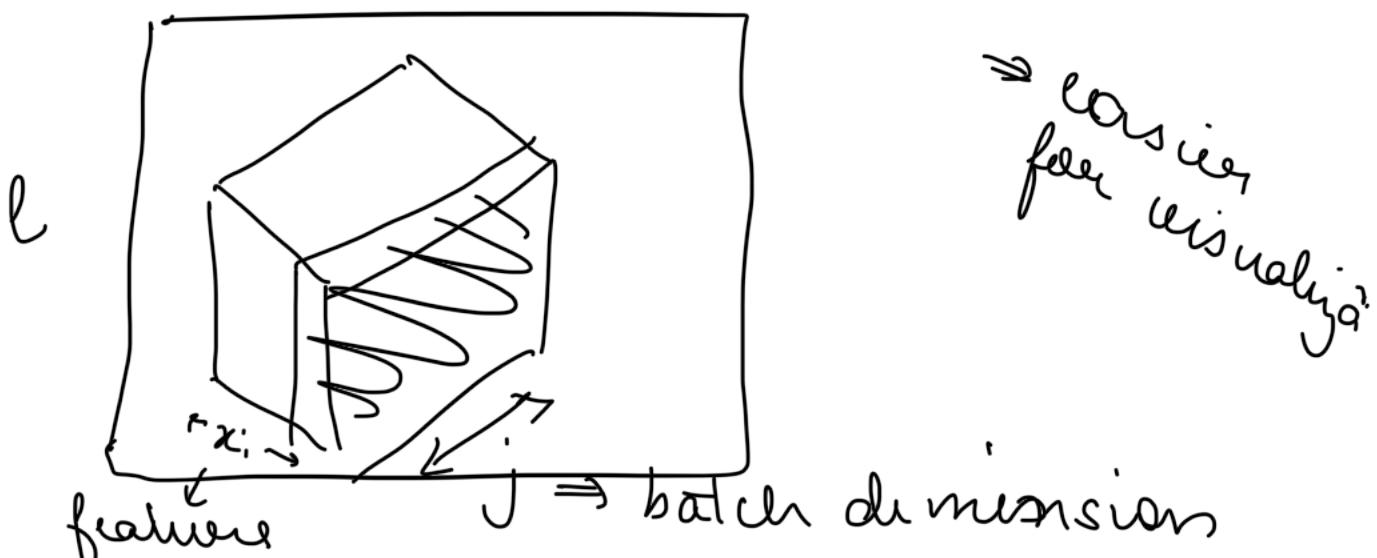
- we are calculating the  $\mu$  &  $\sigma$  from a batch ( $m$ ).  
Why? we just passed 32 or 64 samples  $\Rightarrow$  hence, batch normalizes.
- why did you standardize?

$$\hat{y}_i = \gamma \hat{x}_i + \beta \rightarrow \mu$$

new params

$\sigma \rightarrow$  then standadized  
(older one)

- there 3 variables  $\rightarrow l, i, j$



## L<sup>l</sup> layer Normalization

- yes we can apply batch normalization to transformers. The accuracy depends on the size of m.
- Naive use of BN leads to performance degradation in NLP tasks.
- layer normalization at l<sup>th</sup> layer.

$$D \rightarrow x_1 \rightarrow$$

$$D \rightarrow x_2 \rightarrow$$

$$D \rightarrow x_3 \rightarrow$$

$$D \rightarrow x_4 \rightarrow$$

$$D \rightarrow x_5 \rightarrow$$

$$\sum \rightarrow \frac{1}{H} \rightarrow \mu_l = \frac{1}{H} \sum_{i=1}^H x_i$$

take the avg. across outputs of hidden units in the layer. ∴, the normalization is indep. to no. of samples in a batch.

$$\sigma_l = \sqrt{\frac{1}{H} \sum_{i=1}^H (x_i - \mu_l)^2}$$

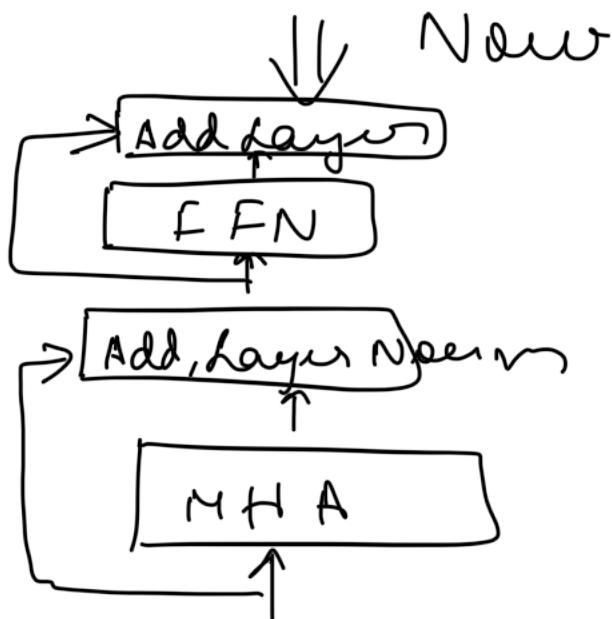
$$\hat{x}_i = \frac{x_i - \mu_l}{\sqrt{\sigma_l^2 + \epsilon}}, \quad \forall i$$

$$\hat{y}_i = \gamma \hat{x}_i + \beta$$

# Encoder



MHA  
↓  
multi-head  
atten<sup>n</sup>



add a residual  
connect<sup>"</sup> & layer  
norm block after  
every MHA, FF N,  
Cross atten<sup>n</sup>  
block.

$$\hat{x} = x + \hat{h} \Rightarrow \text{in the add & normalize layer.}$$

- similarly, we have these connections in the Decoder as well.