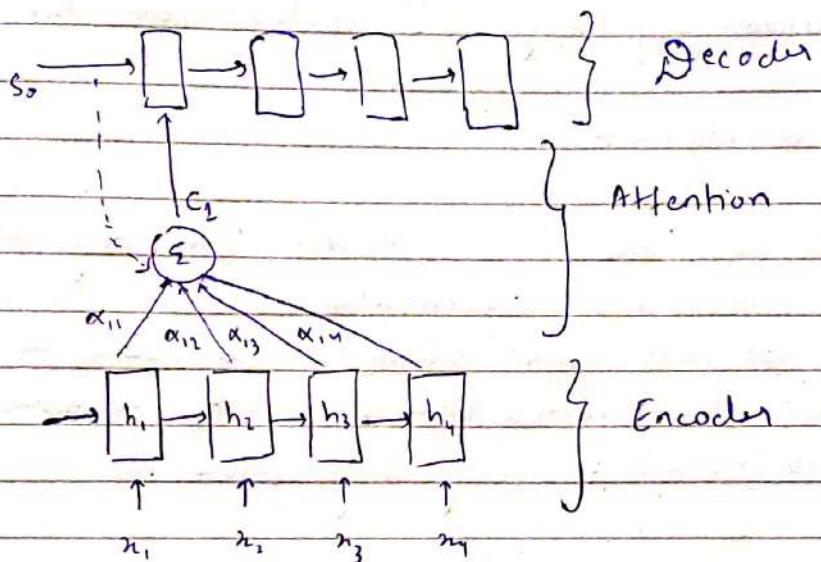


# TRANSFORMERS

\* Introduction to transformers

→ Limitations of Attention based RNNs



As we already know following are some formulas

- $c_t = \sum_{i=1}^n \alpha_{ti} h_i$  context vector for output  $y_t$

$n$ , number of words

$t$  = time step of for decoder

- $\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{i=1}^n \exp(e_{ti})} = \text{softmax}(e_{ti})$

- $e_{ti} = V_{att} \tanh(W_{att} s_{t-1} + U_{att} h_i)$

(Note - notations are a little changed from earlier please manage)

Ishan Modi

Now, what we currently do is as follows

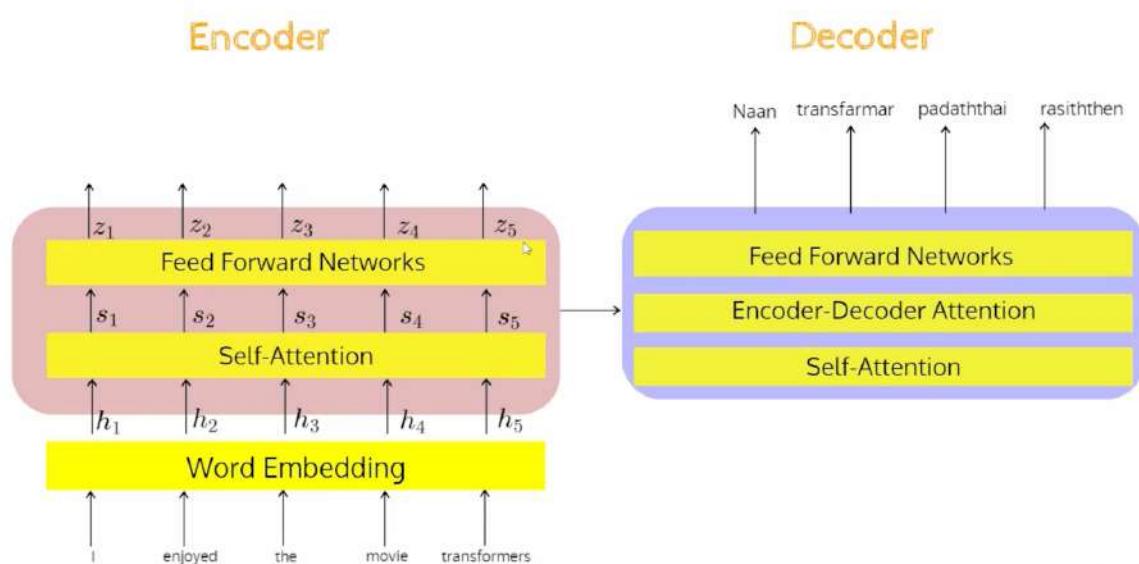
- Compute all the  $h_i$ 's sequentially by iterating through encoder
- Once we have  $h_i$ 's we can compute all ~~eg's~~  $c_t$ 's with each  $c_t$  parallelly
- Then sequentially get  $y_t$ 's (outputs) from the decoder

What can be improved,

- Initially we have all  $x_i$ 's (input) then why do we want to compute  $h_i$ 's sequentially
- So for a particular time stamp in the decoder we already have the previous time step in the decoder, so we want to parallelize every thing below it.

### \* Attention is all you Need (Encoder)

→ Following the above mentioned idea, we want parallelism in the encoder



Ishan Modi

Initially, we had the following formula for attention

$$f_{att} = \text{V}_{att} \tanh \left( \text{U}_{att} s_{t-1} + \text{W}_{att} h_i \right)$$

$\downarrow$                                      $\downarrow$   
 decoder timestep                        encoder timestep

Now, what we have done is

$$f_{att} = f_{att}(h_i, h_j)$$

$\downarrow$                                      $\downarrow$   
 encoder timestep                        encoder timestep

since both are encoder for attention computation  
 this is called Self Attention

Note → we have also removed the recurrent connections  
 that were previously there.

Because of the above mentioned changes parallelism  
 is not only possible per time step but also across  
 multiple timestep.

e.g. for  $h_1 \rightarrow \alpha_{11}, \alpha_{12}, \dots, \alpha_{1n}$

$\left. \begin{matrix} h_2 \rightarrow \alpha_{21}, \alpha_{22}, \dots, \alpha_{2n} \\ \vdots \\ h_n \rightarrow \alpha_{n1}, \alpha_{n2}, \dots, \alpha_{nn} \end{matrix} \right\}$

we are able to compute because it no longer depends on decoder timestep

we are able to compute because all encoder words embeddings are available from get go

Ishan Modi

## \* Self Attention

→ Previously, we had

$$f_{att} = V_{att} \tanh \left( \underbrace{U_{att} s_{t-1}}_{d \times d} + \underbrace{W_{att} h_t}_{d \times d} \right)$$

Here we have 3 vectors ( $v_{AT}$ ,  $s_{f-1}$ ,  $h_i$ )

Following the same idea,

$$f_{att} = f_{att}(h_i, h_j)$$

Step 1 lets say we have 3 matrix ( $w_\alpha$ ,  $w_k$ ,  $w_v$ )

Query matrix	Key matrix	Value matrix
$(d \times d)$	$(d \times d)$	$(d \times d)$

for  $h_1$ ,

$$h_1^T \begin{bmatrix} \longrightarrow \end{bmatrix} \begin{bmatrix} \quad \end{bmatrix} = \begin{bmatrix} \longrightarrow \end{bmatrix} q_1^T$$

$w_\alpha$

ChanModi

$$h_1^T \begin{bmatrix} \longrightarrow \end{bmatrix} \begin{bmatrix} \quad \end{bmatrix} = \begin{bmatrix} \longrightarrow \end{bmatrix} \\ w_k$$

$$h_2^T \begin{bmatrix} \longrightarrow \end{bmatrix} \begin{bmatrix} \quad \end{bmatrix} = \begin{bmatrix} \longrightarrow \end{bmatrix} \\ w_v$$

$\therefore$  for  $h_1$  we get  $(q_1, k_1, v_1)$  by matrix multiplication

similarly, for all  $h_i$  we get  $(q_i, k_i, v_i)$  using same matrix

Step 2

Now for each query we want the importance / weight of each key

$$q_1^T \begin{bmatrix} \longrightarrow \end{bmatrix} \begin{bmatrix} | & | & | \end{bmatrix} = \begin{bmatrix} \longrightarrow \end{bmatrix} \\ k_1 \quad k_2 \dots k_n$$

similarly you can do it for all  $q$ 's by keeping  $K$  matrix the same

Step 3

Do a softmax on  $q^T K$  to make probabilistic distribution

Step 4

Now,  $V_i$ 's are representations of  $h_i$ 's and we want to pick from  $V_i$ 's based on the probability distribution we obtained in above step

Ishan Modi

$$\begin{bmatrix} \dots \\ q_1^T k \end{bmatrix} \begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix} v_i = \begin{bmatrix} \dots \end{bmatrix}$$

representation  
for  $h_i$   
after self attn

Let's formulate the above for  $h_1$ ,

$$e_1 = [q_1 \cdot h_1, q_1 \cdot h_2, \dots, q_1 \cdot h_n]$$

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

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

→ Now, we also need to prove that we can compute all  $(z_1, \dots, z_n)$  for  $(h_1, \dots, h_n)$  in one go.

Step 1

$$\begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_n \end{bmatrix} \begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix} = \begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix} \begin{bmatrix} q_1 \\ q_2 \\ q_3 \\ \vdots \\ q_n \end{bmatrix}$$

$H^T$        $W_a$        $Q^T$

$$\begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix} \begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix} = \begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix}$$

$H^T$        $W_k$        $K^T$

$$\begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix} \begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix} = \begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix}$$

$H^T$        $W_v$        $V^T$

Ishan Modi

Step 2

$$\begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$$

$Q^T$        $K$        $Q^T K$

Step 3

$$\text{Softmax}\left(\frac{Q^T K}{\sqrt{d_k}}\right) = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$$

Step 4

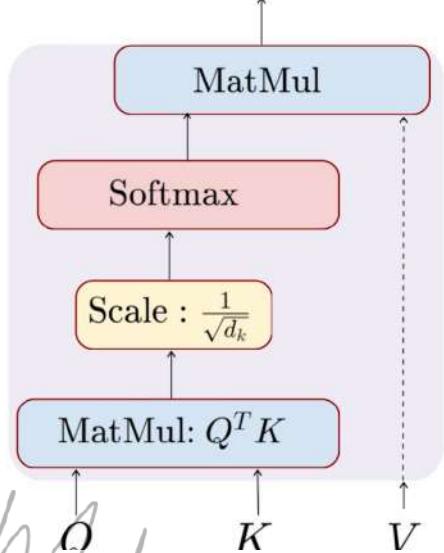
$$\begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n$$

$\text{Softmax}\left(\frac{Q^T K}{\sqrt{d_k}}\right)$        $V^T$        $\text{Self Attn}^T$

So final formulation of self Attn<sup>T</sup>

$$= \text{Softmax}\left(\frac{Q^T K}{\sqrt{d_k}}\right) V^T$$

→ (is the scaling factor generally equal to dimension of key)

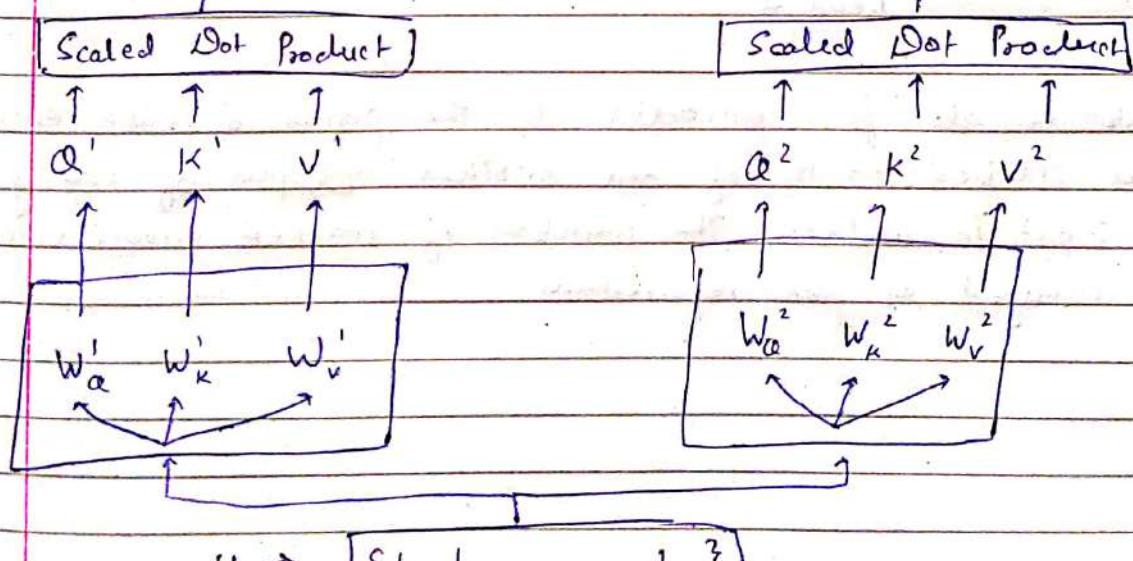
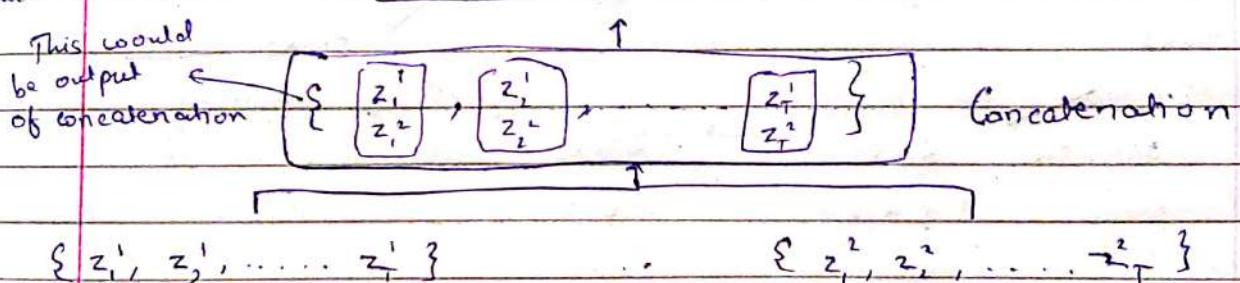
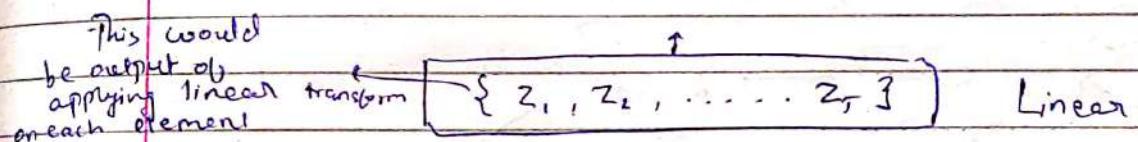


Ishan Modi

## \* Multi-Headed Attention

→ The idea is to use multiple ( $W_a$ ,  $W_k$ ,  $W_v$ ) matrices and generate multiple representations

→ Motivation is the way we use multiple filters to create multiple feature maps, where each feature map is created using one filter. At the end we concat all feature maps for further processing.

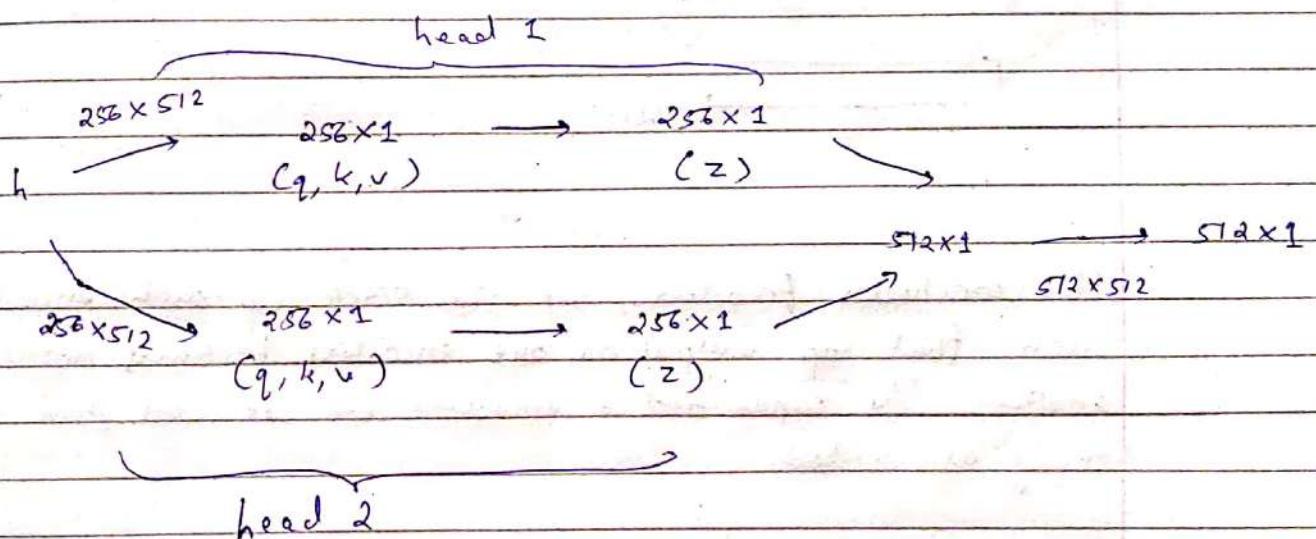


Here we have only shown 2 heads, you can have as many heads as needed.

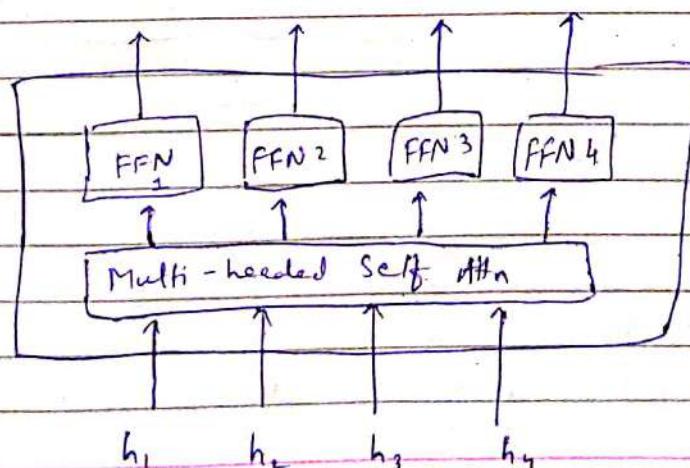
Ishan Modi

In practice, we adjust the dimensions of  $(W_q, W_k, W_v)$  across all heads such that query, key & value vectors are less in dimension than  $h$ 's (input).

Thus each head will generate  $z$ 's which will be less than  $h$ 's in dimension. But when we concatenate the  $z$ 's across all heads the dimension becomes on par or even more than  $h$ 's which is then scaled by linear layer.



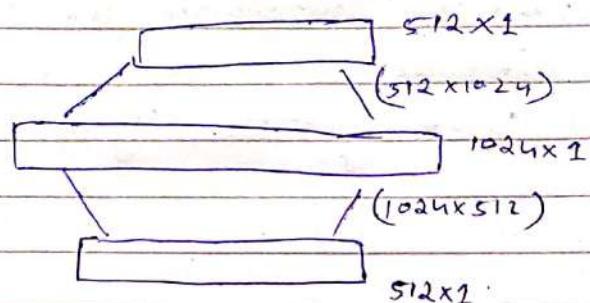
## \* Feed Forward Network (FFN)



Ishan Modi

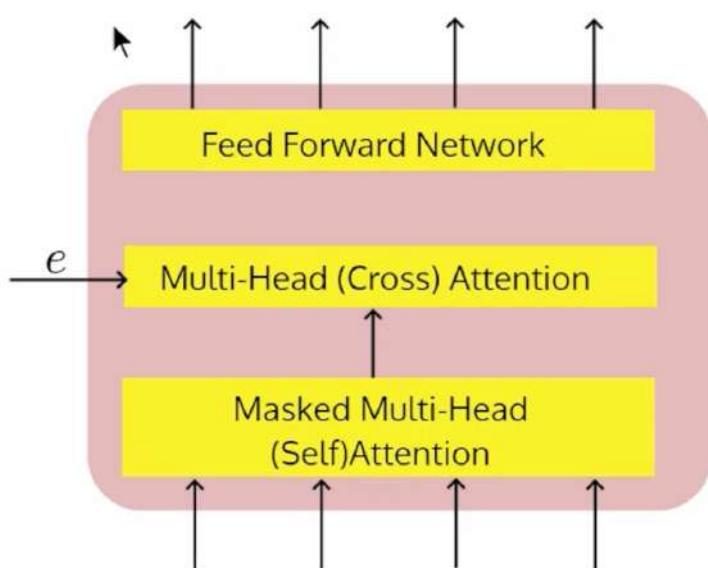
Unlike linear layer block in multiheaded Attn, which had only one projection to the required dimension, FFN, contain hidden layers.

### FFN



This concludes Encoders, we can stack up such encoders such that ~~imp~~ output of one encoder becomes input to another. The paper had 6 encoders but we can pick as many as needed.

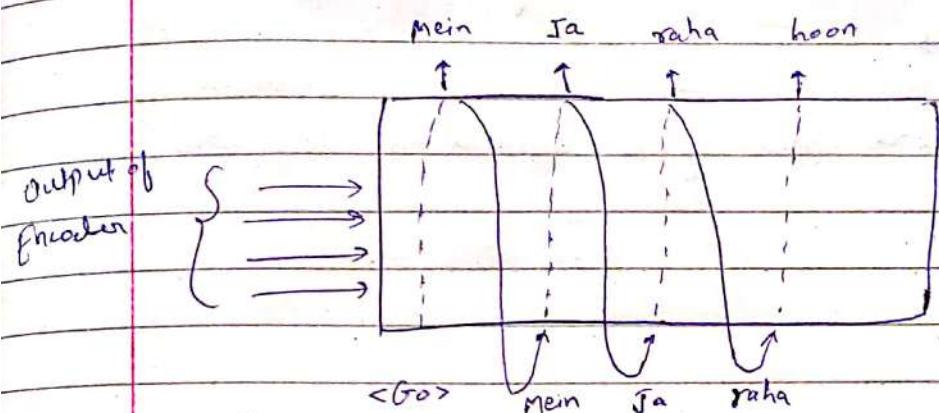
### \* Decoder



Ishan Modi

## \* Teacher Forcing

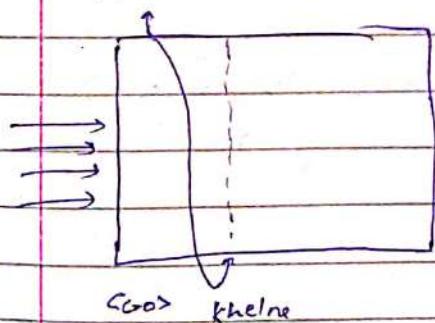
→ As we can see, decoder are auto-regressive, meaning output at time step one is given as an input along with all the previous inputs at timestep 2 and so on..



original sentence - I am going

Now, what if the word at the first timestamp is predicted wrong ?

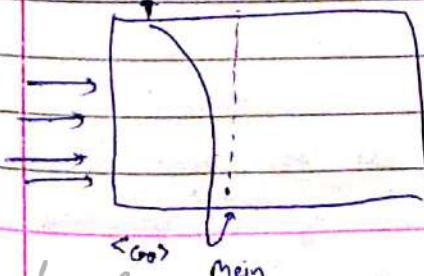
khelne



This could mess up all the subsequent predictions.

A solution to this would be

khelne



Though it predicted "khelne" wrongly we give the correct thing "mein" in the input. We do this for the first few epochs forcing the correct word as input & later stop forcing once network has learnt relevant amount.

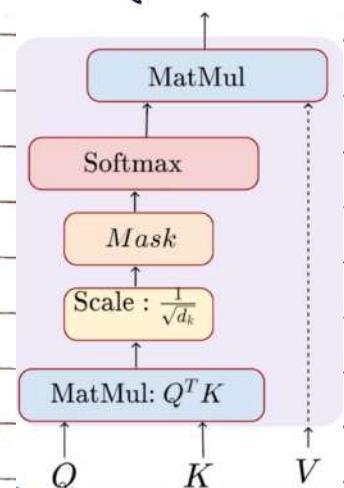
Ishan Modi

## \* Masked - Multi-Head Self Attention

→ As we know, we cannot pass all the words at a given timestamp to the decoder

At a given timestamp we can only pass all the previously predicted words.

Following is where masking is applied in self-attn block



How ~~does~~ does masking work,

$$= \begin{bmatrix} q_1 k_1 & q_1 k_2 & q_1 k_3 \\ q_2 k_1 & q_2 k_2 & q_2 k_3 \\ q_3 k_1 & q_3 k_2 & q_3 k_3 \end{bmatrix} + \begin{bmatrix} 0 & -\infty & -\infty \\ 0 & 0 & -\infty \\ 0 & 0 & -\infty \end{bmatrix} \quad \text{Masking}$$

$$= \begin{bmatrix} q_1 k_1 & -\infty & -\infty \\ q_2 k_1 & q_2 k_2 & -\infty \\ q_3 k_1 & q_3 k_2 & q_3 k_3 \end{bmatrix}$$

$$\text{Softmax} \left( \begin{bmatrix} q_1 k_1 & -\infty & -\infty \\ q_2 k_1 & q_2 k_2 & -\infty \\ q_3 k_1 & q_3 k_2 & q_3 k_3 \end{bmatrix} \right) = \begin{bmatrix} S_{11} & 0 & 0 \\ S_{21} & S_{22} & 0 \\ S_{31} & S_{32} & S_{33} \end{bmatrix} \quad \text{Softmax}$$

Basically we are adding  $(-\infty)$  to all the next words in the sequence so that output after softmax becomes 0 so that no importance is given to next word in sequence

Ishan Modi

Note → Masking also happens in encoder,

token length → dimension of vector representing each word  
 sequence length → total no of tokens allowed.

So for an encoder if the number of tokens are less than sequence length, padding is done

These new tokens generated using padding should be ignored at the time of self attention, this is done using masking

## \* Cross Attention

→ As we know from earlier, query is the word for which we are trying to get representation. So the query part in cross attention comes from decoder & key, value part comes from encoder

→ Following is formula for masking - multi-head attention  
 & Encoder representation

- Decoder representation =  $\text{Softmax}(\mathbf{Q} \mathbf{k}^T + \mathbf{M}) \mathbf{v} = (\mathbf{D}_i)$

- Encoder representation = output from encoder. =  $(\mathbf{E}_i)$

$$\mathbf{D}_i \rightarrow \mathbf{D}_i \cdot \mathbf{W}_Q \rightarrow \mathbf{Q}$$

$T \times d$        $T \times d$        $d \times d$        $T \times d$

$$\mathbf{E}_i \rightarrow \mathbf{E}_i \cdot \mathbf{W}_K \rightarrow \mathbf{K}$$

$T \times d$        $T \times d$        $d \times d$        $T \times d$

$$\mathbf{E}_i \rightarrow \mathbf{E}_i \cdot \mathbf{W}_V \rightarrow \mathbf{V}$$

$T \times d$        $d \times d$        $T \times d$

After this steps  
 similar to self  
 attention are followed

Ishan Modi

Note =  $D_i$  &  $E_i$  can also have different dimension

$$\text{eg } D_i = T_1 \times d$$

$$E_i = T_2 \times d$$

but at the end of decoder we need output similar to input ie  $\underline{T_1 \times d}$ . If we do all calculations

with  $D_i \rightarrow T_1 \times d$  &  $E_i \rightarrow T_2 \times d$ , The decoder block will give us  $T_1 \times d$  as the output because

$$T_1 \times T_2 \cdot T_2 \times d \rightarrow T_1 \times d$$

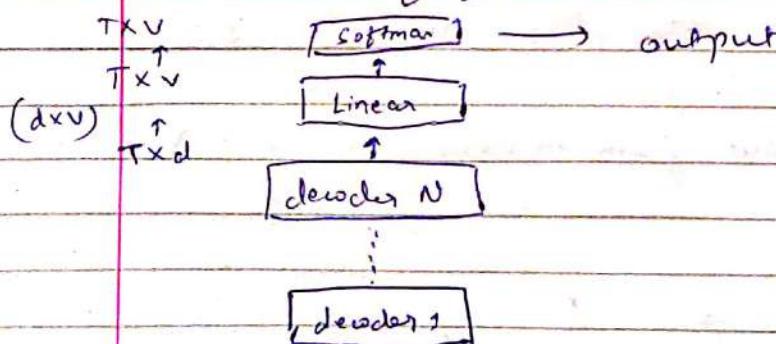
$$(Q^T) \cdot V \rightarrow (\text{output})$$

\* Feed Forward Network (Similar to Encoder; Masked)  
 attention output is feed to FFN

\* Linear & Softmax

→ We can stack up as many decoders on top one another, original paper had 6.

→ At the end of final decoder we have linear layer + softmax



So the linear layer has  $d \times v$  weights & it projects to  $v$  dimension & softmax gives probability distribution

$$T \times d \cdot d \times v \rightarrow T \times v \xrightarrow{\text{softmax}} T \times v$$

decoder      linear      output

probability distribution in  
each row

Note - till now we have discussed attention with the following formula

$$A^T = \text{softmax} \left( \frac{\alpha^T K}{\sqrt{d_k}} \right) V^T$$

Here following are the  $\alpha$ ,  $K$ ,  $V$  matrices

$$\alpha = \begin{bmatrix} | & | & | & | \\ k_1 & k_2 & k_3 & \dots & k_n \end{bmatrix}$$

Now  $K = \begin{bmatrix} | & | & | & | \\ v_1 & v_2 & v_3 & \dots & v_n \end{bmatrix}$

$$V = \begin{bmatrix} | & | & | & | \\ a_1 & a_2 & a_3 & \dots & a_n \end{bmatrix}$$

But if we use the above formula can also be written as,

$$A^T = \text{softmax} \left( \frac{\alpha K^T}{\sqrt{d_k}} \right) V$$

This is possible with the following  $\alpha$ ,  $K$ ,  $V$  matrices

ShanMod

$\Omega =$ 

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

fundamental frequencies

values in terms of fundamental frequencies are given by

 $K = \text{diag}$ 

$$\begin{bmatrix} k_1 & & \\ k_2 & \ddots & \\ \vdots & & \\ k_n & & \end{bmatrix} \quad \text{and} \quad N = \text{diag} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

 $A =$ 

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

plot of radius of gyration with respect to axis, with successive segments of beamwise areas change in mode shape

last writing indicates that bending angle is in terms of angle because it shows sum of bending of each segment is linear sum of individual bending angles at larger end

and if there is no load then bending moment is zero and so will be the deflection curves and problems omitted

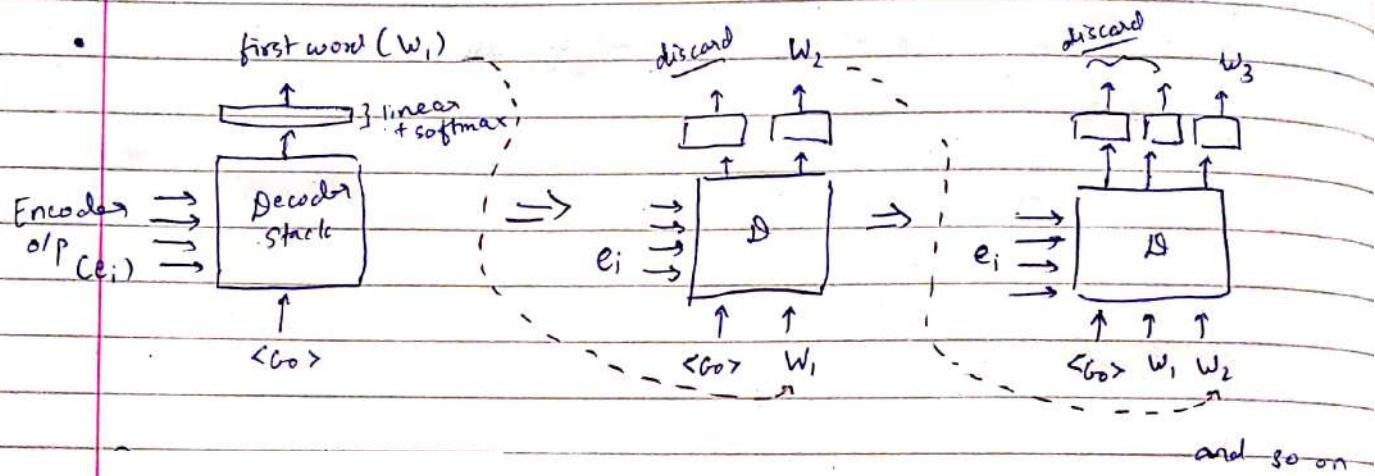
indicates that there is also a gap between two segments which is filled with transition function having two boundary conditions at both ends

Ishan Modi

Note - Lets discuss how decoder actually work.

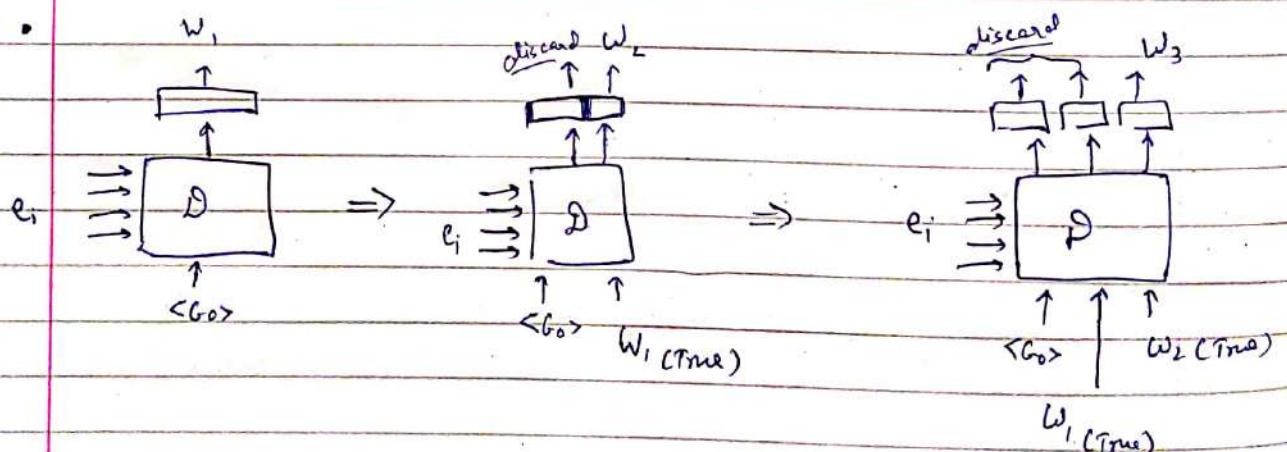
CASE - 1

Training phase without teacher forcing



CASE - 2

Training phase with teacher forcing

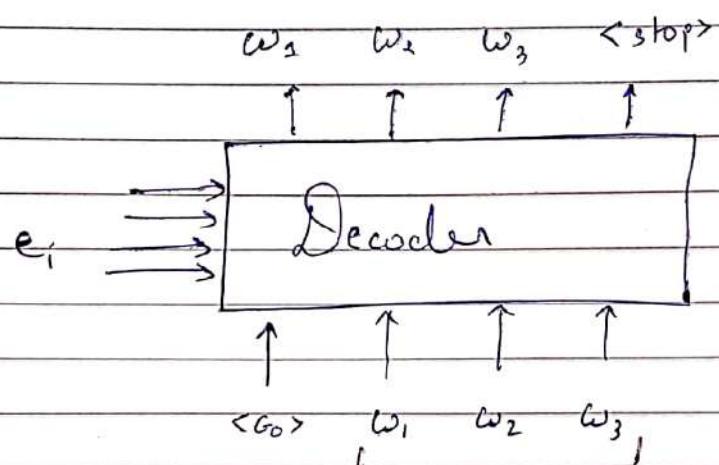


Here we are not using the predicted labels/tokens at all in next iteration

Ishan Modi

Since we are not using the predicted token at each iteration & we have the entire input during training, we can run the decoder in parallel during training under the teacher forcing scheme.

Though teacher forcing is faster, it has to be used along with normal training, so model learns to deal with mistakes at inference.



These come from the training data & are same irrespective of what is predicted

Ishan Modi

- Other points to remem' remember

- Self attention matrix  $\rightarrow (T+1) \times (T+1)$  because sequence starts with  $\langle \text{bos} \rangle$  & ends with  $\langle \text{eos} \rangle$
- For every word <sup>entire</sup>, decoder has to be iterated
- Masking is required in subsequent decoder blocks in a decoder stack irrespective of teacher forcing.

## \* Positional Encodings

- Motivation - Currently, up until now we don't have any positional information in the words. If we switch words in input given to encoder, it won't affect output.
- Following are our initial thoughts on how to get position into token

- $w_1 \quad w_2 \quad w_3 \quad \text{and so on}$  } This is not feasible  
 $+ \quad + \quad +$   
 $[1, 1, \dots] \quad [2, 2, 2, 2, \dots] \quad [3, 3, \dots]$   
 because if sequence length is 256 then you will end up adding 256 to last token & 256 is huge when compared to random initialization
- $w_1 \quad w_2 \quad w_3 \quad \text{and so on}$  } use one-hot  
 $+ \quad + \quad +$   
 $[1, 0, \dots] \quad [0, 1, \dots] \quad [0, 0, 1, \dots]$

The notion of distance is lost, eg  $\text{dist}(w_1, w_2) = \sqrt{2}$   
 distance is calculated between one-hot encoding of  $w_1$  &  $w_2$

now, the  $\text{dist}(w_1, w_3) = 5_2$ , essentially distance between  $w_1$  &  $w_3$  should be more.

- $w_1$        $w_2$        $w_3$       +      and so on      } positional encoding  
 $[....]$        $[....]$        $[....]$       } are learnable  
 $\alpha_1$        $\alpha_2$        $\alpha_3$       } parameters

This also has an issue, if we have encountered only max sequence length of 40 during training, then only 40 positional encodings will be learned. At inference if a sentence has more than 40 tokens we won't have any learned ~~each~~ encodings.

## \* Sinusoidal Embedding

$$\rightarrow \text{PE}_{(j,i)} = \begin{cases} \sin\left(\frac{j}{10000} 2\pi / d_{\text{model}}\right), & i = 0, 2, \dots, d_{\text{model}} \\ \cos\left(\frac{j}{10000} 2\pi / d_{\text{model}}\right), & i = 1, 3, \dots, d_{\text{model}} \end{cases}$$

Here,

1           6  
 2           5  
 3           4  
 .  
 .  
 .  
 5           6

j = position of token in sequence

$i$  = value inside token embedding

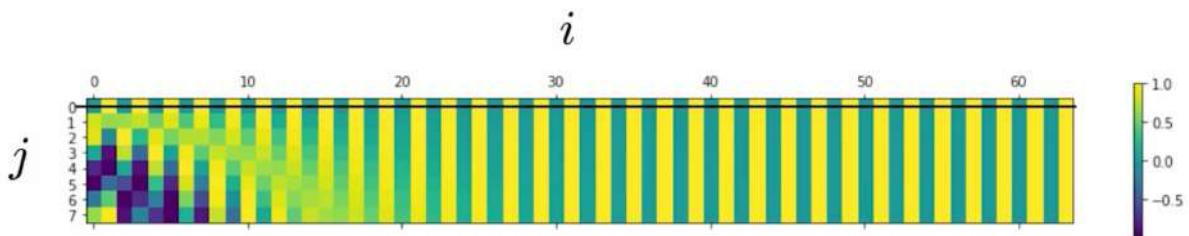
$d_{model}$  = length of token embedding

oldmodel = ;

Ishan Modi

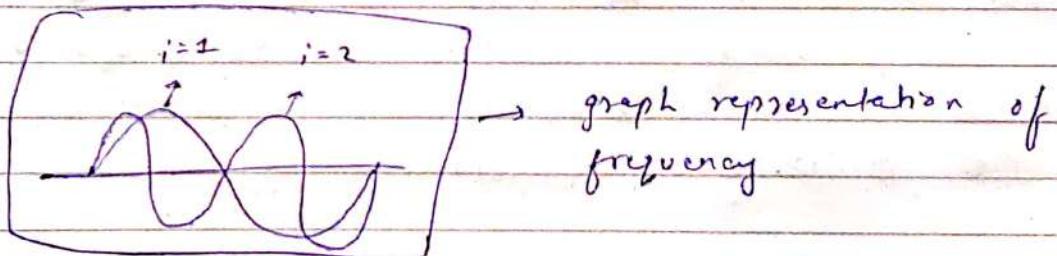
Note - If "i" is even we use sin function & if it is odd we use cosine function.

→ Following is a diagram



→ In the following function,

$\sin\left(\frac{j}{10000} \cdot 2i/\text{modulo}\right)$ , for a fixed j as the value of i increases the frequency decreases



→ By using sinusoidal function,

- We are creating a dependence between  $j \& i$ , every value at  $i$  positions in the embedding depend on  $j$ .
- Notion of distance is also there. Eg - 32 word will be close to 31st word & 33rd word & distance increases as we move in either direction

Ishan Modi

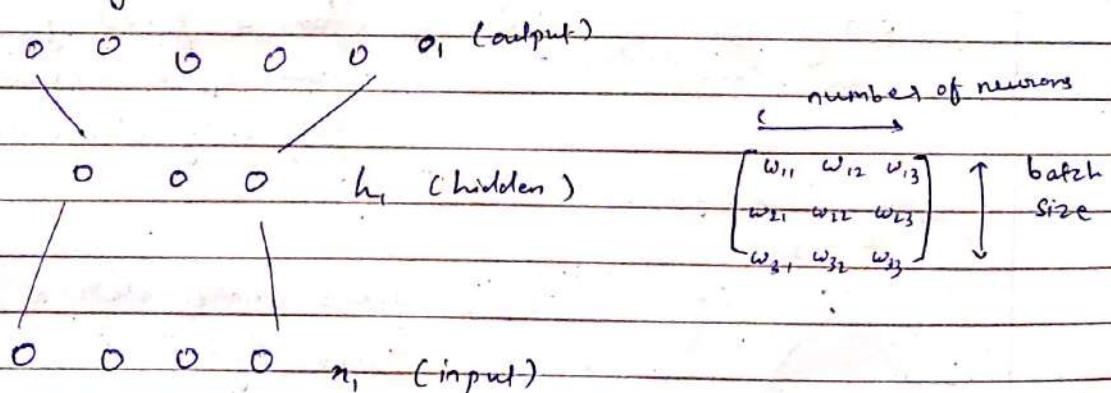
## \* Residual Connection & Layer Normalization

### → Residual Connection

Since there is a stack of encoder & a stack of decoder in transformer, and individually these encoder & decoder also contain attention layers, there is a chance of vanishing gradient. To deal with that we use residual connections.

### → Layer Normalization

- It is very similar to Batch Norm



Let consider the above Neural Network.

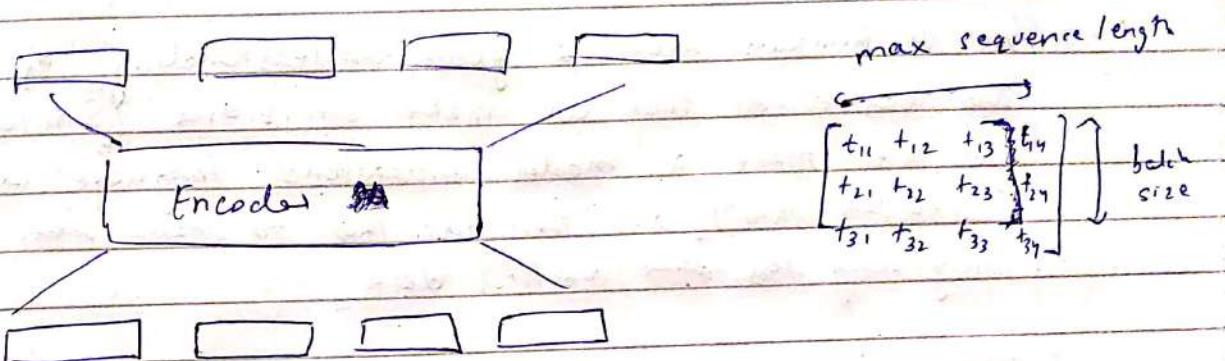
- For Batch normalization,

all the values within the batch are normalized, so each column will be normalized individually (normalization takes place across column). i.e.  $(w_{11}, w_{21}, w_{31})$  are normalized.

- For Layer normalization,

all the values at the layer are normalized, so the each row will be normalized individually. i.e (normalization happens across rows) ie - ( $w_{11}, w_{12}, w_{13}$  are normalized)

Now, how does layernorm work in transformers

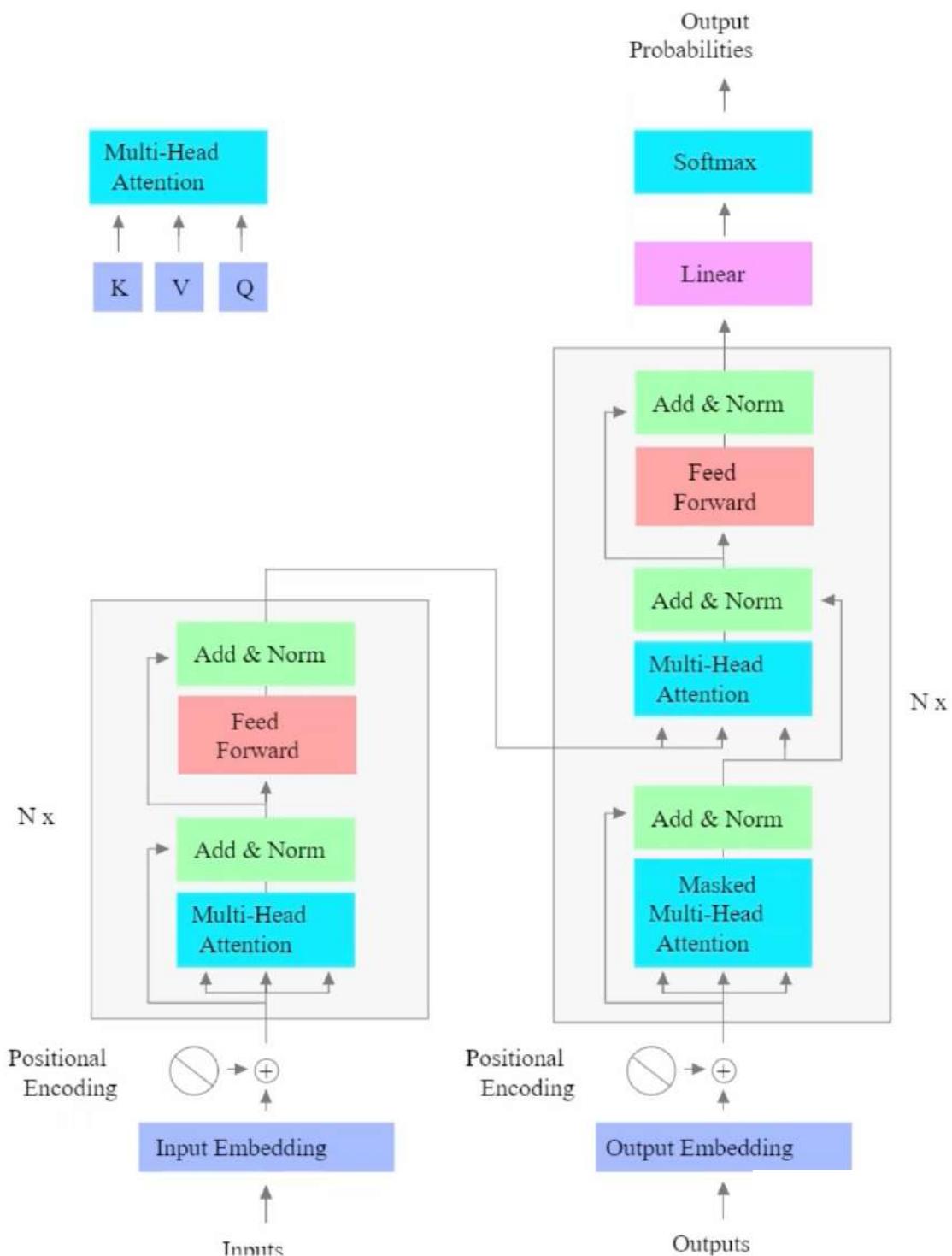


so layernorm is calculated across each row, also each token is dmodel (dimensional). Thus the layer norm will be calculated for the following

$$(t_{111} + t_{112} + \dots + t_{11d}) + (t_{121} + t_{122} + t_{123} + \dots + t_{12d}) \\ + (t_{131} + t_{132} + \dots + t_{13d}) + (t_{141} + t_{142} + \dots + t_{14d})$$

similar for other rows, as well.

# Transformer Architecture



Ishan Modi

~~Decoder- Only - Models~~

## \* Language Modelling

- Labelled data is always limited, what if we can use un-labelled data for language modelling
- We as humans have a good understanding of language and based on that we make decisions. Similarly can we make a model understand language using raw-text (unlabelled data) and then use that to make decisions using very less data (labelled data) -

→ pretraining      Supervised -  $\leftarrow$  Finetuning

- What is language modelling ?

Let say, we have  $\{x_1, x_2, \dots, x_n\} \in V$ , where  $V$  is the vocabulary &  $x_i$ 's are all the words

We can construct a lot of different sentences using these words. Now when we are reading/refering some text the probability of some sentences occurring is higher because they are more common than the other sentences.

The probability of <sup>a</sup>sentence appearing can be given as follows.

$$f: (x_1, x_2, x_3, \dots, x_n) \rightarrow [0, 1]$$

function

some probabilistic function suggesting words will occur in the fixed sequence

Ishan Modi

Now,

$$f = p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_1, x_2) \cdot \dots \cdot p(x_T | x_1, x_2, \dots, x_{T-1})$$

probability distribution over  
vocab that gives first  
word (marginal probability)

probability distribution over the  
vocab that gives second word  
given previous (first) word  
(conditional probability)

All other are also conditional probability

$$= p(x_1) \cdot \prod_{i=2}^T p(x_i | x_1, \dots, x_{i-1})$$

→ How do we link pre-training, with the above formula  
of language modeling?

So, let's pick a task of predicting next token in a sequence

Here you don't need labelled data, from raw-text available  
you can give the model  $T$  tokens and tell it to predict  
 $T+1^{th}$  token.

Also, each term in the above formula, can be calculated  
by the previously seen transformer architecture.

Following diagram is showing, the equation & its relation with  
transformers,

Ishan Modi

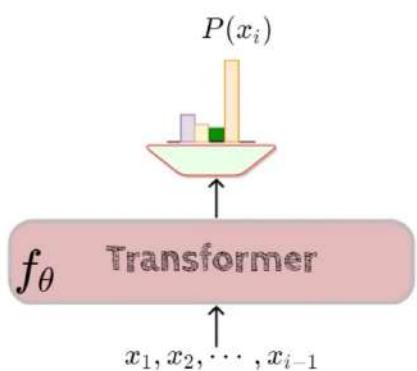
## Causal Language Modelling (CLM)

$$\begin{aligned} P(x_1, x_2, \dots, x_T) &= \prod_{i=1}^T P(x_i | x_1, \dots, x_{i-1}) \\ &= P(x_1)P(x_2 | x_1)P(x_3 | x_2, x_1) \cdots P(x_T | x_{T-1}, \dots, x_1) \end{aligned}$$

We are looking for  $f_\theta$  such that

$$P(x_i | x_1, \dots, x_{i-1}) = f_\theta(x_i | x_1, \dots, x_{i-1})$$

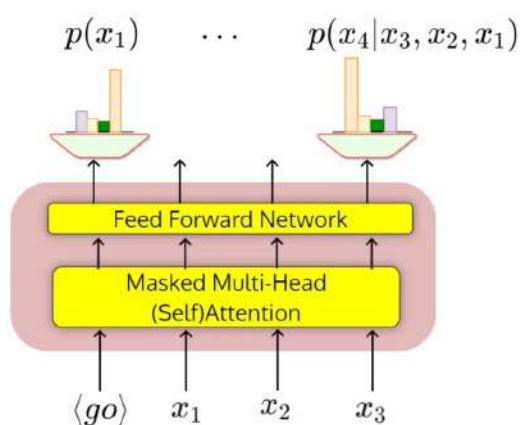
Can  $f_\theta$  be a transformer?



The transformer architecture can be,

- Encoder Only
- Decoder Only
- Encoder - Decoder

## \* Decoder Only Models

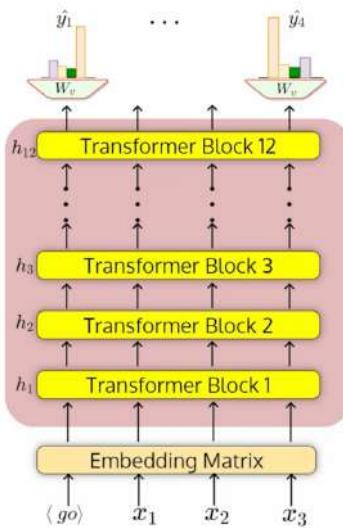


This is a basic diagram for decoder-only model and here the joint probability for the sequence is

$$f = p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) p(x_4 | x_1, x_2, x_3)$$

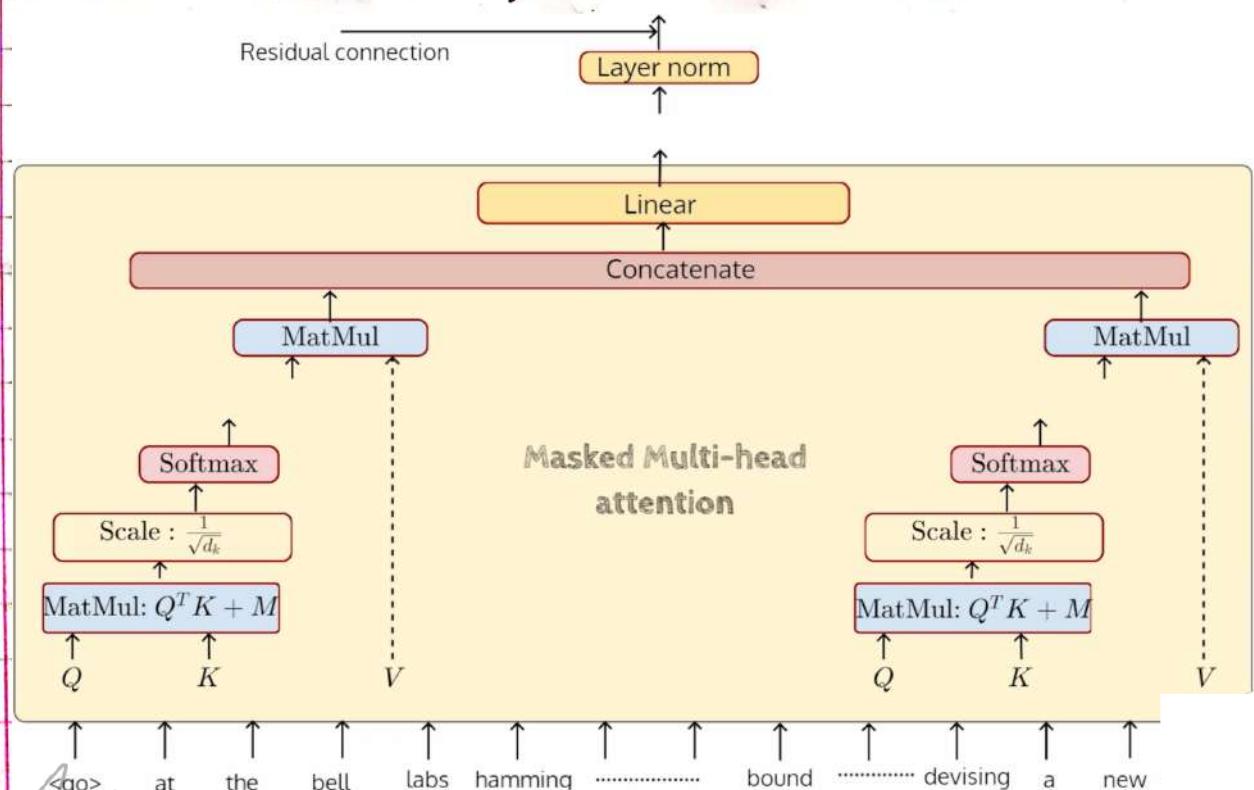
# \* Generative Pretrained Transformers (GPT)

Note - The transformer block in the figure below is similar to the one from previous section i.e- decoder only models



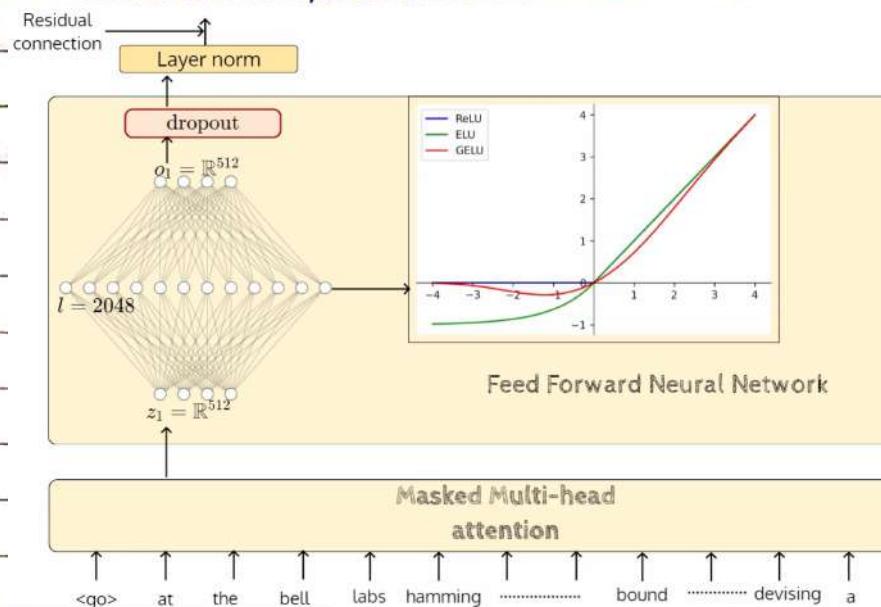
→ Now, in the previous section we also saw that there are two main components of decoder only models, following are the big diagrams of each component w.r.t GPT architecture

- Masked Multi-Head Self Attention



Ishan Modi

## • Feed Forward Network



→ Let us now discuss the model configurations of GPT

(1) Tokenizer - Byte Pair Encoding

(2) Vocab Size - 40478

(3) Embedding Dimension - 768

(4) Decoder Layer (Transformer Blocks) - 12

(5) Context Size (Sequence length) - 512

(6) Attention heads - 12

(7) FFN hidden Dimension -  $768 \times 4 \Rightarrow 3072$

→ Let us now count the total parameters in the model

(1) Input Embedding Layer

$$\begin{aligned}
 \text{Token Embeddings} &= 40478 \times 768 \\
 &= 31 \times 10^6 \\
 &= 31 \text{ M}
 \end{aligned}$$

- Positional Embeddings =  $512 \times 768$   
 $= 0.3 \times 10^6$   
 $= 0.3 \text{ M}$

Total = 31.3 M

## (2) Attention layers

Per Block  $W_Q, W_K, W_V$  all have equal number of parameters / dimension

$$\therefore W_Q = W_K = W_V = (768 \times 64)$$

Per Head

So in one head we have

$$3 \times (768 \times 64) = 147 \times 10^3$$

Each layer has 12 heads

$$12 \times 147 \times 10^3 = 1.7 \text{ M}$$

## Linear layer

$$768 \times 768 = 0.6 \text{ M}$$

$$\text{Total} = 1.7 \text{ M} + 0.6 \text{ M} = 2.3 \text{ M}$$

Each GPT has 12 blocks

$$\text{Total Total} = 12 \times 2.3 = 27.6 \text{ M}$$

Ishan Modi

### (3) Feed Forward layers (FFN)

~~Per Block~~

$$2 \times (768 \times 3072) + \underbrace{3072 + 768}_{\text{bias}} \\ = 4.7 \times 10^6 = 4.7 \text{ M}$$

GPT has 12 blocks

$$\text{Total} = 12 \times 4.7 = 56.4 \text{ M}$$

→ Combining (1), (2), (3) → we have 117 M parameters in GPT model

Note - Final layer needs to generate a distribution over the vocab, therefore given 768 dim embedding it needs to generate 40478 dim output. Same is applied for all of the sequence length.

Thus we need  $768 \times 40478$  matrix for the transformation and it is shared, because token embeddings also have same sized matrix.

## \* Pretraining & Finetuning

~~Pretraining~~

Let discuss the pre-training configurations

- (1) Batch size - 64
- (2) Sequence length - 512
- (3) hidden dimension/embedding length - 768
- (4) Loss -

$$- \sum_{\substack{\text{inputs} \\ \text{sequence}}} \sum_{\substack{\text{tokens} \\ \text{in sequence}}} (\log \hat{y}_i)$$

Ishan Modi

(5) Optimizer - Adam with cosine learning rate scheduler  
(6) Strategy - Teacher Forcing.

## Tine tuning

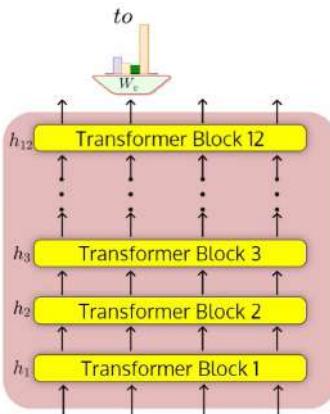
Note - Fine-Tuning can be done in the following ways

- Freezing all the pre-trained weights in the network, add a linear layer head / small neural network, and change only these added weights during backpropagation for required task
  - Freezing the bottom k-layers, others as above
  - Don't freeze the pre-trained model at all, others as above

→ Some use cases of finetuning can be

- Sentiment Analysis
  - Question / Answering
  - Textual Entailment etc

Text Generation is also a usecase, where given a sequence of tokens you can generate new tokens

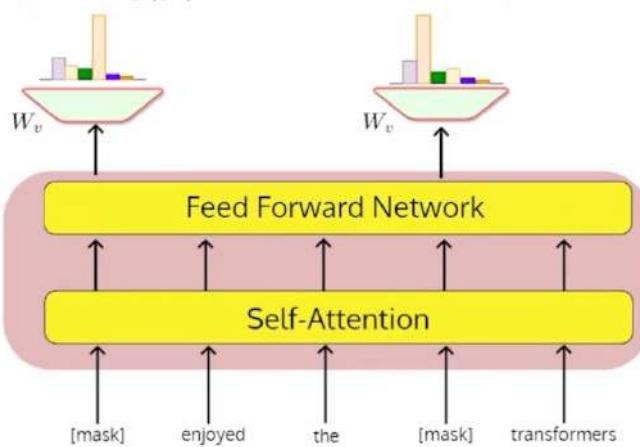


The only problem here is that this seems to be deterministic i.e. - get same o/p everytime

Ishan Modi

## \* Encoder - Only - Models

- Why do we need encoder only models?  
decoder only models are auto-regressive & they look at past tokens to predict future tokens
- Future & Past tokens both can contribute to the prediction of the current token.



- Now, we used causal language modeling for decoder-only model. In case of encoder-only-models we use something called as masked language modeling

$$P(y_i | n_1, n_2, \dots, n_i = [\text{mask}], \dots, n_j)$$

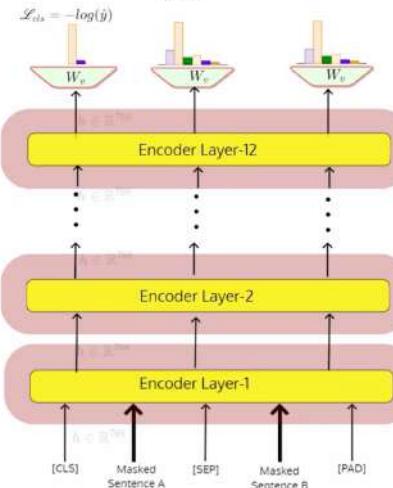
Here we will provide  $n_j$ 's as the input & will mask ( $n_j$  where  $i=j$ ) & will try to predict  $y_i$  for the mask ( $n_i$ ) given the  $n_j$ 's.

## \* BERT (Bidirectional Encoder Representations of Transformers)

→ Continuing from the previous section, the question is how does masking work internally?

Answer is, you introduce a new token called [MASK] → (has its own embedding) into the vocabulary. Then you pass it through the network and try to learn the actual word through back propagation

$$\mathcal{L} = \frac{1}{N} \sum_{y_i \in M} -\log(\hat{y}_i) + \mathcal{L}_{cls}$$



→ Following is the model configuration of BERT

- (1) Tokenizer - Word Piece
- (2) Vocab Size - 30522
- (3) Embedding Dimension - 768
- (4) Encoder layers - 12
- (5) Sequence Length - 512
- (6) Attention heads per layer - 12
- (7) FFN layer Dimension -  $768 \times 4 = 3072$

→ Let's now count the total number of parameters of the model

### (1) Embedding layers

- Token Embeddings =  $30522 \times 768 = 23 \text{ M}$

- Segment Embeddings =  $2 \times 768 = 1536$

- Positional Embeddings =  $512 \times 768 = 0.4 \text{ M}$

$$\text{Total} = 23.4 \text{ M}$$

We will learn more in pre-training section

### (2) Attention layers

~~per block~~  $w_k = w_q = w_v = 768 \times 64$  ~~per head~~

$\frac{1}{3}$  of them =  $768 \times 64 \times 3$

$$\text{for 12 heads} = 768 \times 64 \times 3 \times 12 = 1.7 \text{ M}$$

Also, Linear layer =  $768 \times 768 = 0.6 \text{ M}$

$$\begin{aligned} \text{Total 12 blocks} &= 2.3 \times 12 \\ &= 27.6 \text{ M} \end{aligned}$$

(3) FFN Layer =  $(2 \times 3072 \times 768) + 3072 + 768$

$$= 4.7 \text{ M}$$

$$\begin{aligned} \text{For total 12 blocks,} &= 4.7 \times 12 \\ &= 56.4 \text{ M} \end{aligned}$$

Ishan Modi

Adding, all of them, (1), (2), (3) it is around 110 M

## \* Pretraining

### → Masked Language Modeling (MLM)

We know from previous discussion, how masked language modelling works. But we don't know how many tokens to mask and some other details

- 15% of the words are masked from the input sequence
- 80% of the total masked token are replaced with [MASK] embedding
- 10% of the total masked tokens are replace with random word embeddings
- 10% are unchanged

Eg - Let's say sequence length is 200 tokens

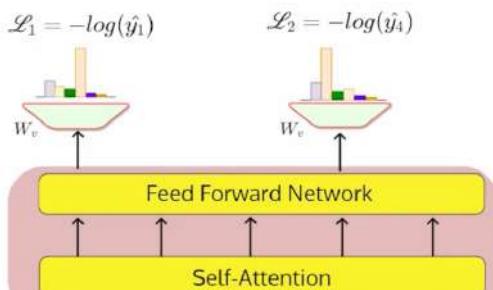
$\therefore 15\% \text{ of } 200 \rightarrow 30 \text{ words are selected for masking}$

$80\% \text{ of } 30 \rightarrow 24 \text{ words are replaced with [MASK]}$

$10\% \text{ of } 30 \rightarrow 3 \text{ words } " " " \text{ random word}$

$10\% \text{ of } 30 \rightarrow 3 \text{ words remain unchanged}$

$$\mathcal{L} = \frac{1}{N} \sum_{y_i \in M} -\log(\hat{y}_i)$$



## → Next Sentence Prediction (NSP)

- Introduce [CLS] token that denotes beginning of input sequence
- Introduce [SEP] token that denotes sentence break
- Segment Embedding were introduced, dim  $\Rightarrow 768 \times 2$

$$\text{ie } 768 \times 1 = SE_A \quad 4 \quad 768 \times 1 = SE_B$$

Denotes that it is the first sentence

Denotes that it is the 2nd sentence

- 50% of the sentences had correct sequence & for other 50% of the sentences 2<sup>nd</sup> sentence was picked randomly from the corpus.

