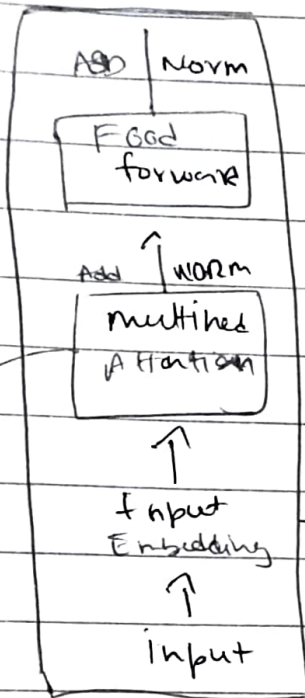
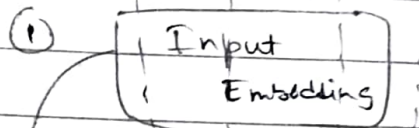


# Transformers (LSTMs are DEAD)

## # Encoders



Hi How ARE you



## ② Positional Encoding Odd/Even

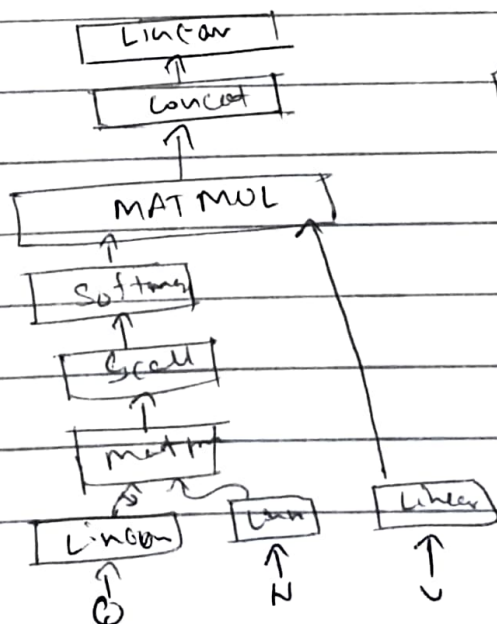


$$\text{Even} = \sin \left( \frac{\text{Pos}}{10000^{2i/d_{\text{model}}}} \right)$$

$$\text{Odd} = \cos \left( \frac{\text{Pos}}{10000^{2i/d_{\text{model}}}} \right)$$

## Multihead Attention

## -> Self Attention



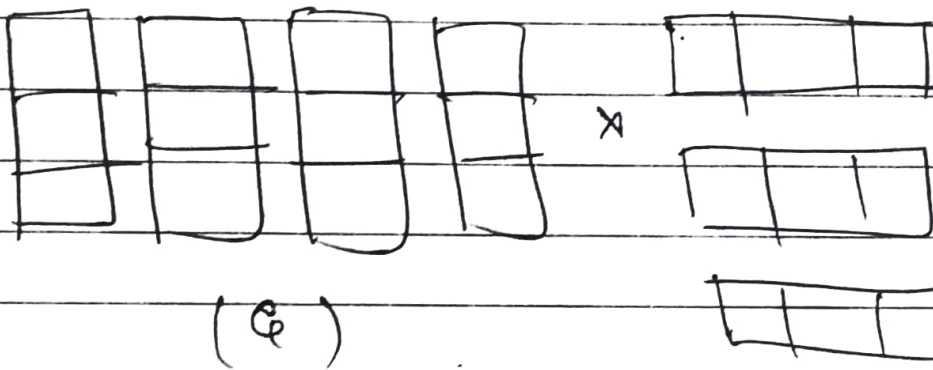
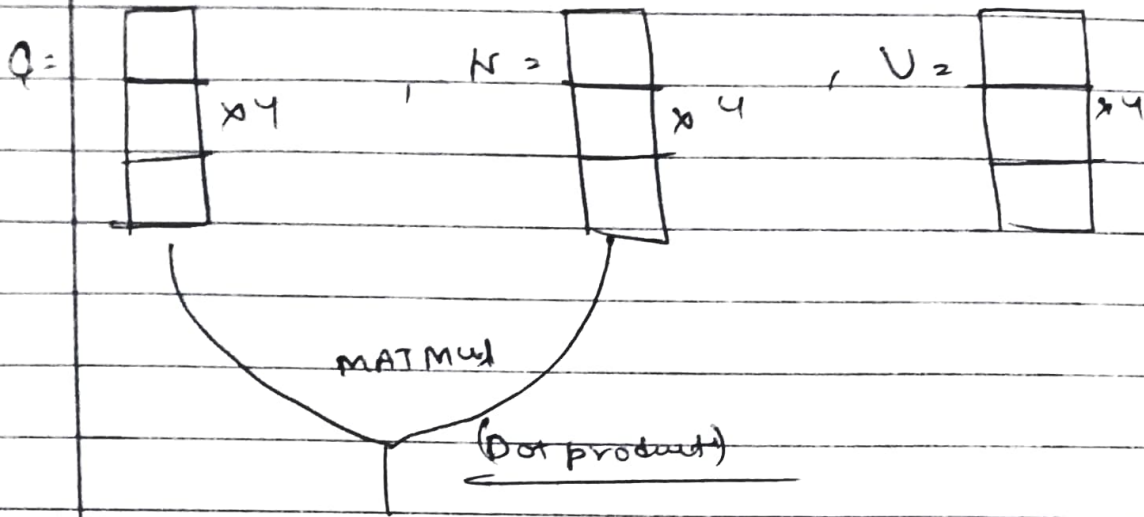
Self Att. helps model to associate each word with every word in the vocabulary Input.

Q, K, V { Youtube example

Hi    How    Are    You

Page :

Date :



$(K)^T$

hi	48	29	10	12
How	27	89	31	67
Are	10	31	91	54
You	12	67	54	92
	hi	How	Are	You

(Score)

$\Rightarrow$  Basically Score matrix is a kind of confusion matrix which tells the correlation b/w each input vector or word . .

Higher the score in score matrix higher is the focus.

⇒ Now scores are scaled by dividing them with dimension of key

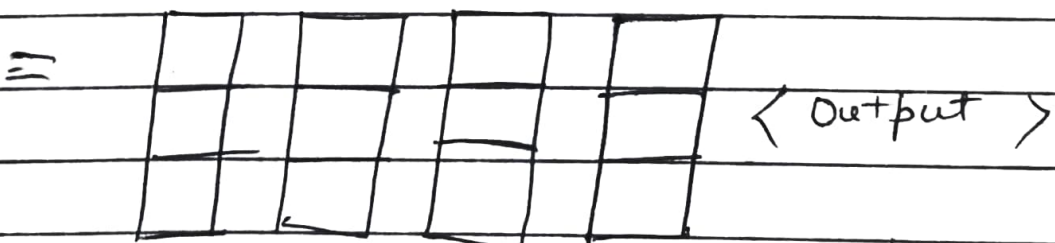
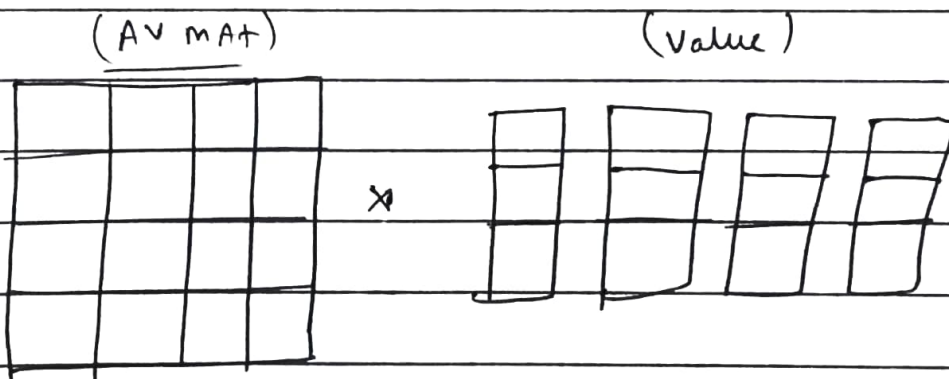
⇒ Apply softmax to scaled scores.

$$\text{softmax}(z_{ij}) = \frac{e^{z_{ij}}}{\sum_{j=1}^n e^{z_{ij}}}$$

By softmax high scores are relatively probabalized higher and low scores are diminished.

⇒ now score matrix becomes Attentioned volume matrix (AV MAT)

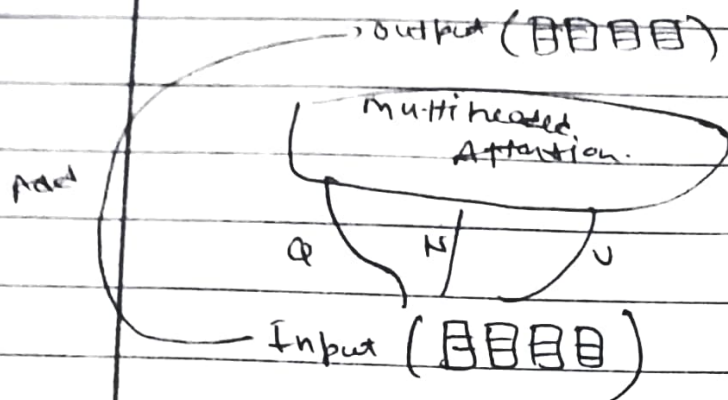
⇒ Multiply AV matrix with value.



{ This vector gives the metric of how each word should attend to other's words. }

4) 7) Residual Connect / Layer Norm

and Point Feed forward.



# Now Input and output are added and passed through norm layer to get values easy to be calculated.

# Now it is passed in in a linear layer with ReLU activation and At last is normalised.

x ————— x End of the Encoder layer ————— x

Note :- We can use Encoder much more times to further encode the information and thus it learns further more about the features and gets better.



# <DECODERS>

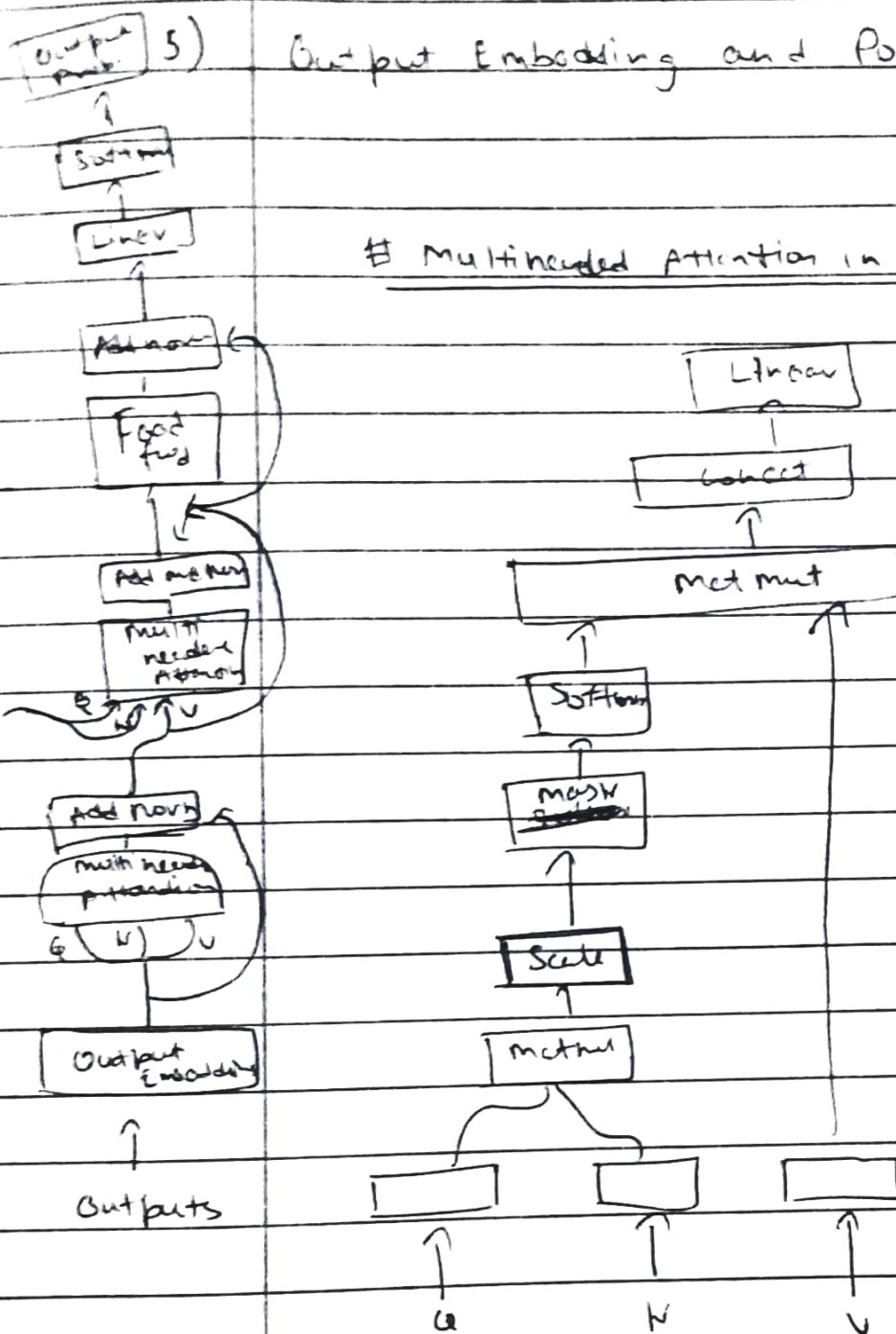
Page

Date

5)

## Output Embedding and Position Encoding

### # Multitasked Attention in Decoders



Note:

The mask layer is here because

De-coder is generating a sequence so it is putting the word at the end of the output what is the significance of the future coming value acc. to the matrix.



So every word has access to scores of itself and words before it.

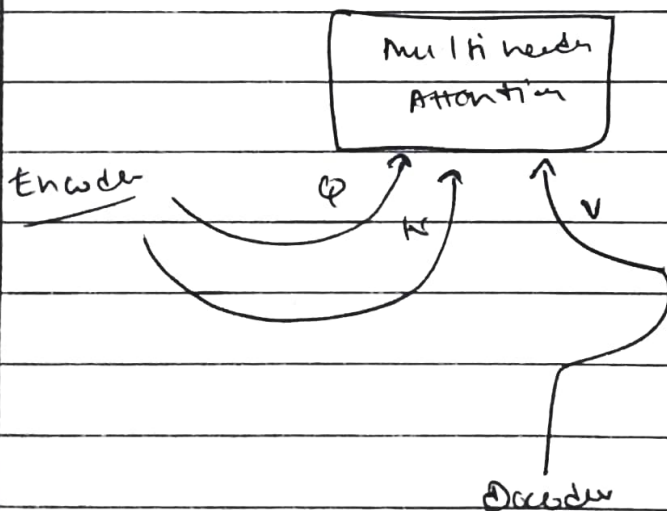
$$\begin{bmatrix} 0.1 & 0.6 & 0.3 \\ 0.3 & 0.7 & 0.2 \\ 0.4 & 0.4 & 0.2 \end{bmatrix} + \begin{bmatrix} 0 & -5H & -5H \\ 0 & 0 & -5H \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.1 & -5H & -5H \\ 0.3 & 0.8 & -5H \\ 0.4 & 0.4 & 0.2 \end{bmatrix}$$

softmax

Now taking the softmax of mask matrix will give us

$$\begin{bmatrix} 0.1 & 0 & 0 \\ 0.8 & 0.3 & 0 \\ 0.1 & 0.2 & 0.2 \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} 0.1 & 0 & 0 \\ 0.8 & 0.3 & 0 \\ 0.1 & 0.2 & 0.2 \end{bmatrix}} \right\} \text{Remember it}$$

6.) Second Multi-Headed Attention { check b/w encoder & k and decoder v. }



Here the Decoder and Encoder provided weights are matched and new output matrix is given and in last linear and softmax function is used to get the probability of every word.