

# RECURRENT NEURAL NETWORK

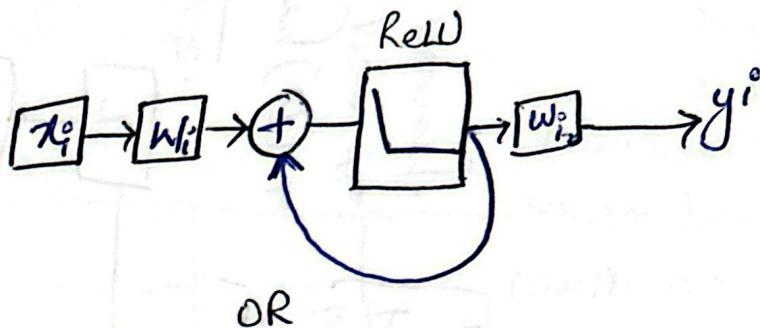
→ is a class of neural network designed to model sequential data. It has a memory mechanism mean its output depend on current state plus previous state

Hidden state update.

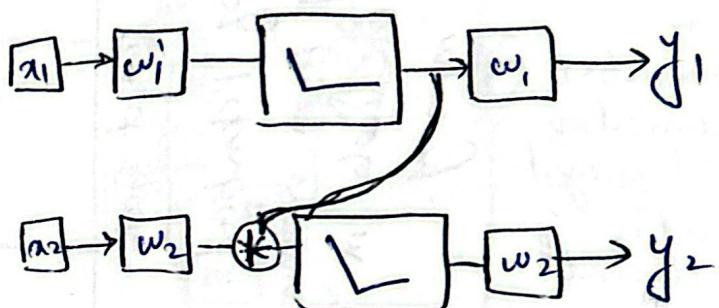
$$h_t = f(W_h h_{t-1} + b_h)$$

Output.

$$y = g(W_y h_t + b_y)$$



OR



\* The LSTM / GRU provides a solution.

→ Vanishing gradients :-

low steps cause reach max  
steps without finding optimum

→ Exploding gradient  $> 1$  : too  
big steps cause divert from  
optimum.

→ Difficulty with long-term  
dependencies.

captures short term context

struggles over long sequences.

→ Training Insufficiency

sequential in nature  
can't use GPU etc.

## RNN Pro's

- o same set of weights & biases are used.
- o good for time dependent tasks (speech, text ch)
- o can process sequences of differential length
- o can be configured for different Input/Output mappings

# LONG - SHORT TERM MEMORY

cell-state  $\rightarrow$  long term memory  
hidden state  $\rightarrow$  short term memory

Sigmoid.  
Take  $x$  and feed it to  $y$  in byw  
 $\text{Tanh}$   
Take  $x \rightarrow y (-1, 1)$

- ① LSTM learns long term dependencies and resolve RNN issues.

can remember information and forget about irrelevant

still sequential mean no GPU,

more complex than RNN

Memory intensive due + gates

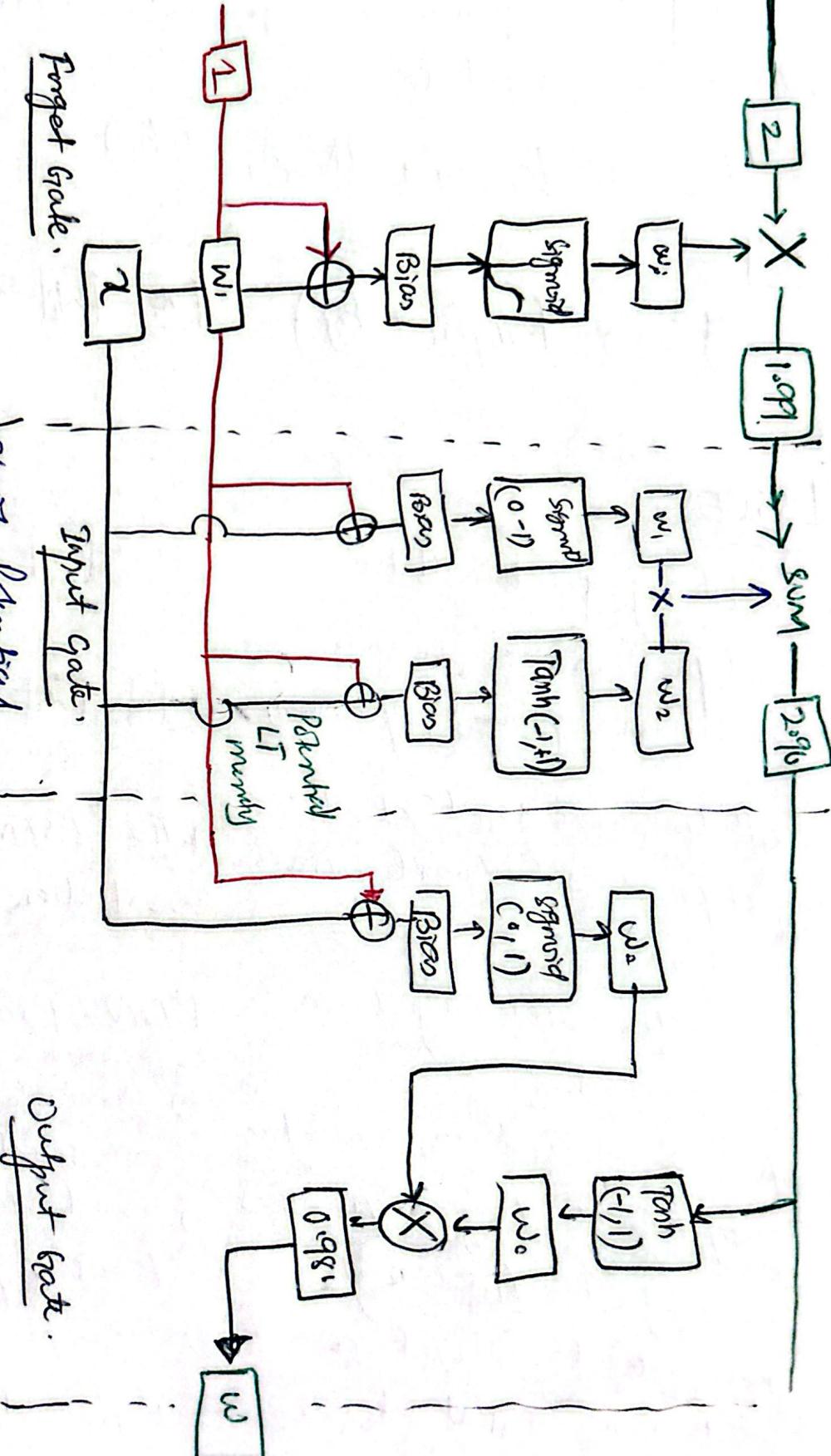
% of Memory ( $L_T$ ) to remember from the long Term memory.

% of Practical long term memory to include in  $L_T$  memory.

Forget Gate.

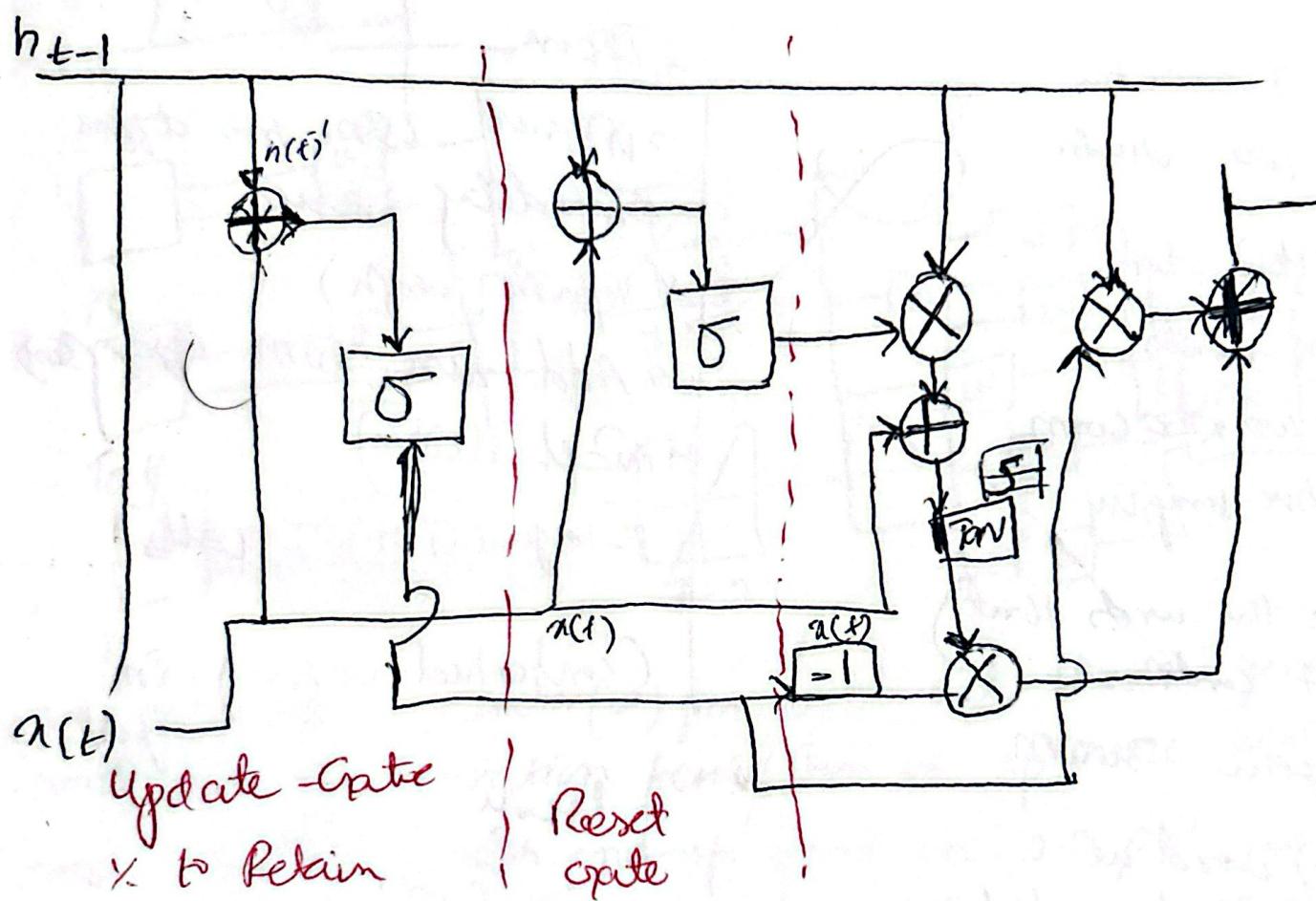
Input Gate.

Output Gate.



## Gated Recurrent Unit (GRU)

- up of RNN designed to avoid gradient issues and learn long-term dependencies in sequential data. → similar to LSTM but faster. (best for small sequence) LSTM for long.
- combines long & short term memory
- Update → How much past memory to keep
- Reset → How much to forget



(deux)

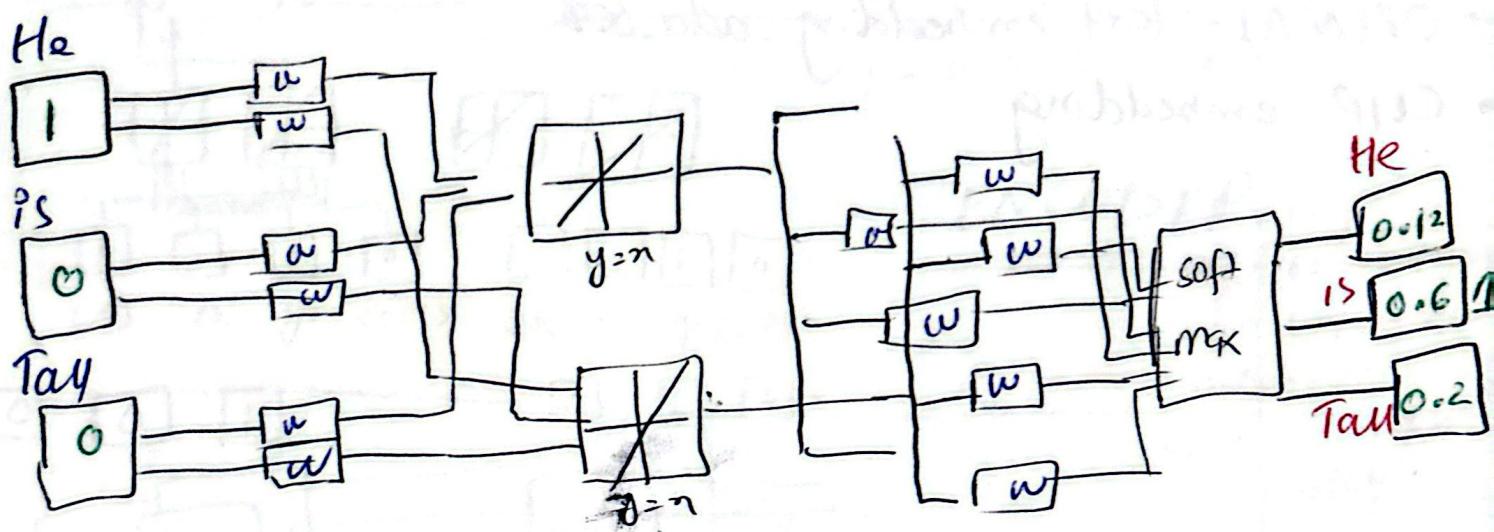
## WORD - EMBEDDINGS

Word embedding is a vector representation of word in a continuous vector space.

### Pros:-

- ) Dense and low-dimensional representation
- ) Captures semantic/syntactic meaning.
- ) Enable transfer learning, better generalization vs one hot vector

Example by NN:



→ We associate 2 numbers ( $w$ ) per word but 100 in general and 100 of activation fn's. Then we optimize weight and similar words end-up near in 2-D plot. we associate various to capture the context (+ve, -ve, etc) word 2 vec.

Usually, we train over whole Wikipedia etc

### 1) CBOW:

We give surrounding words to predict middle ones mean predict current word given context words.

### 2) Skip-gram :

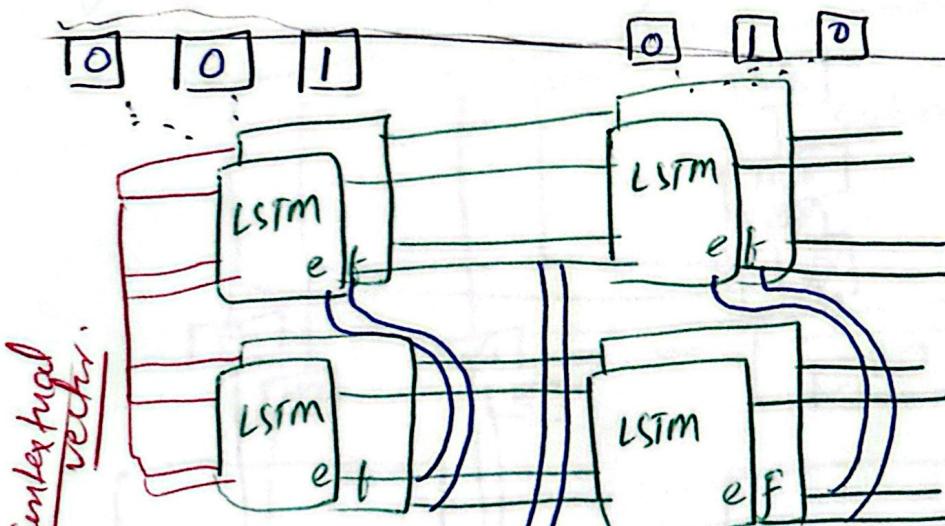
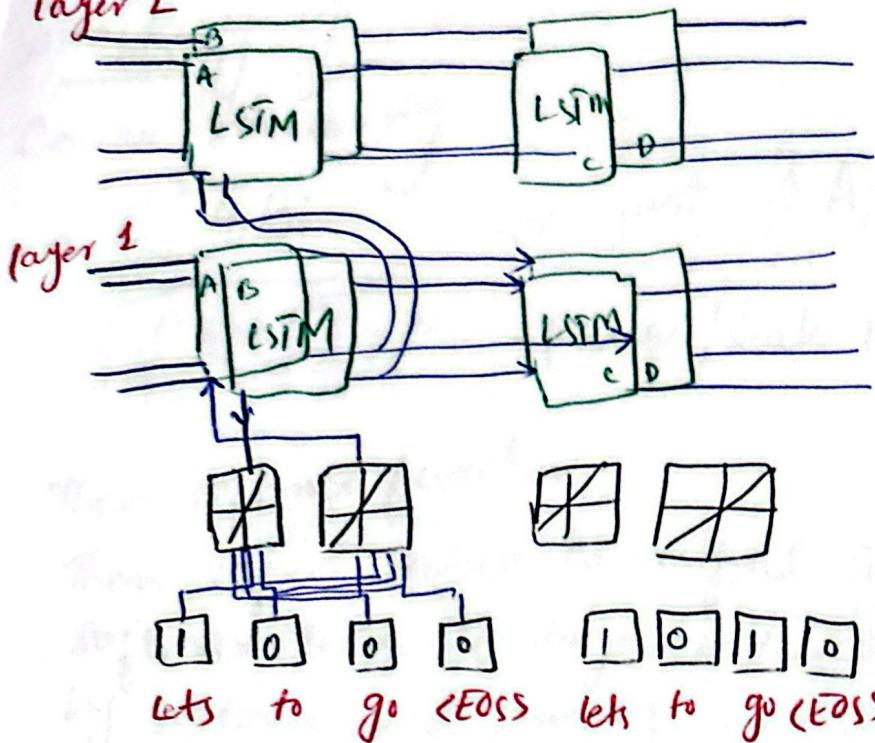
vise-versa

## SEQ-2-SEQ: (Encoder - Decoder)

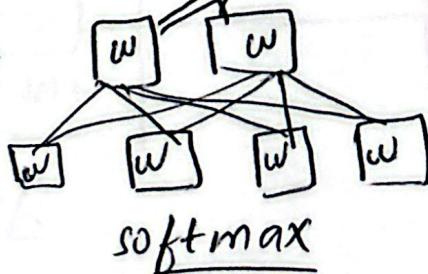
is a RNN design to transform one sequence into another.

Eg: for Translation.

layer 2



FC-layer



1 0 0 0  
Vamos ir y <EOS>

0 0 1 1  
Nunca p y <EOS>  
(cont'd.)

## TEACHER-FORCE(ing)

contextual  
vector

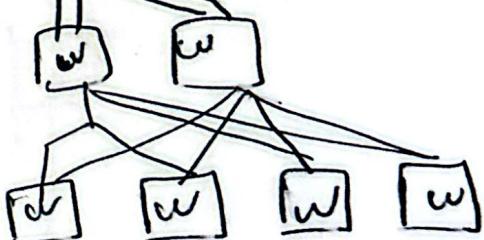
- 2 cells - 2-layer LSTM
- each cell have some weight/prob
- but differed in layer

## ENCODER, Paper:

input	160,000 tokens
output	80,000 tokens
embed	1000
LSTM	1000 nodes 2 layers

## DECODERs

384 m  
weight/  
biases



softmax

# ATTENTION FOR NEURAL NETWORKS :-

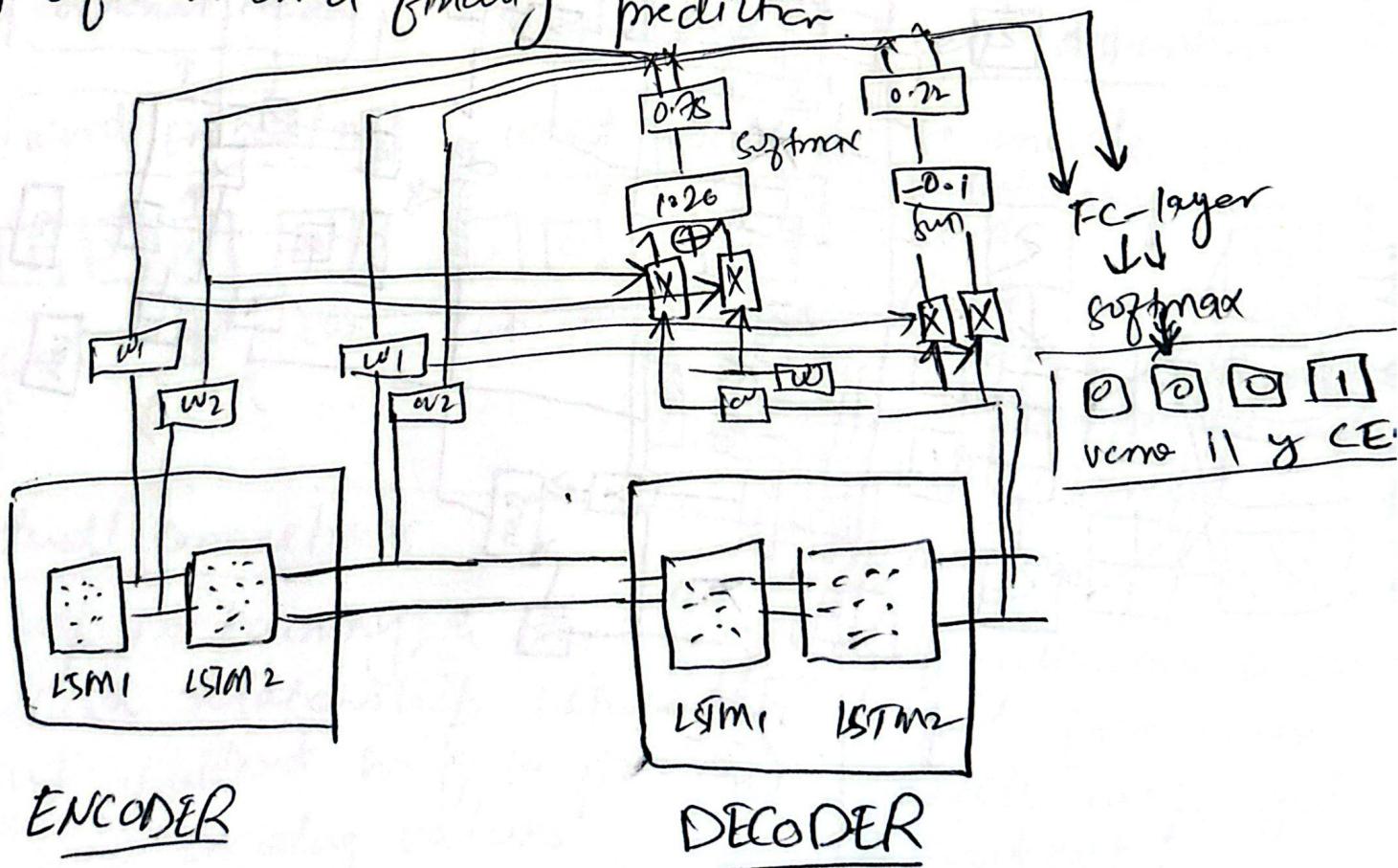
- adding an additional path for each input value so that each step of decoder can directly access these values similarity of series.

Cosine-Similarity

$$\Rightarrow \frac{\sum_{i=1}^n A_i \cdot B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} \quad \text{or just } \sum_{i=1}^n A_i \cdot B_i \quad \begin{cases} \text{dot product} \\ \text{range / scale to } (-1, 1) \end{cases}$$

Then softmax function.

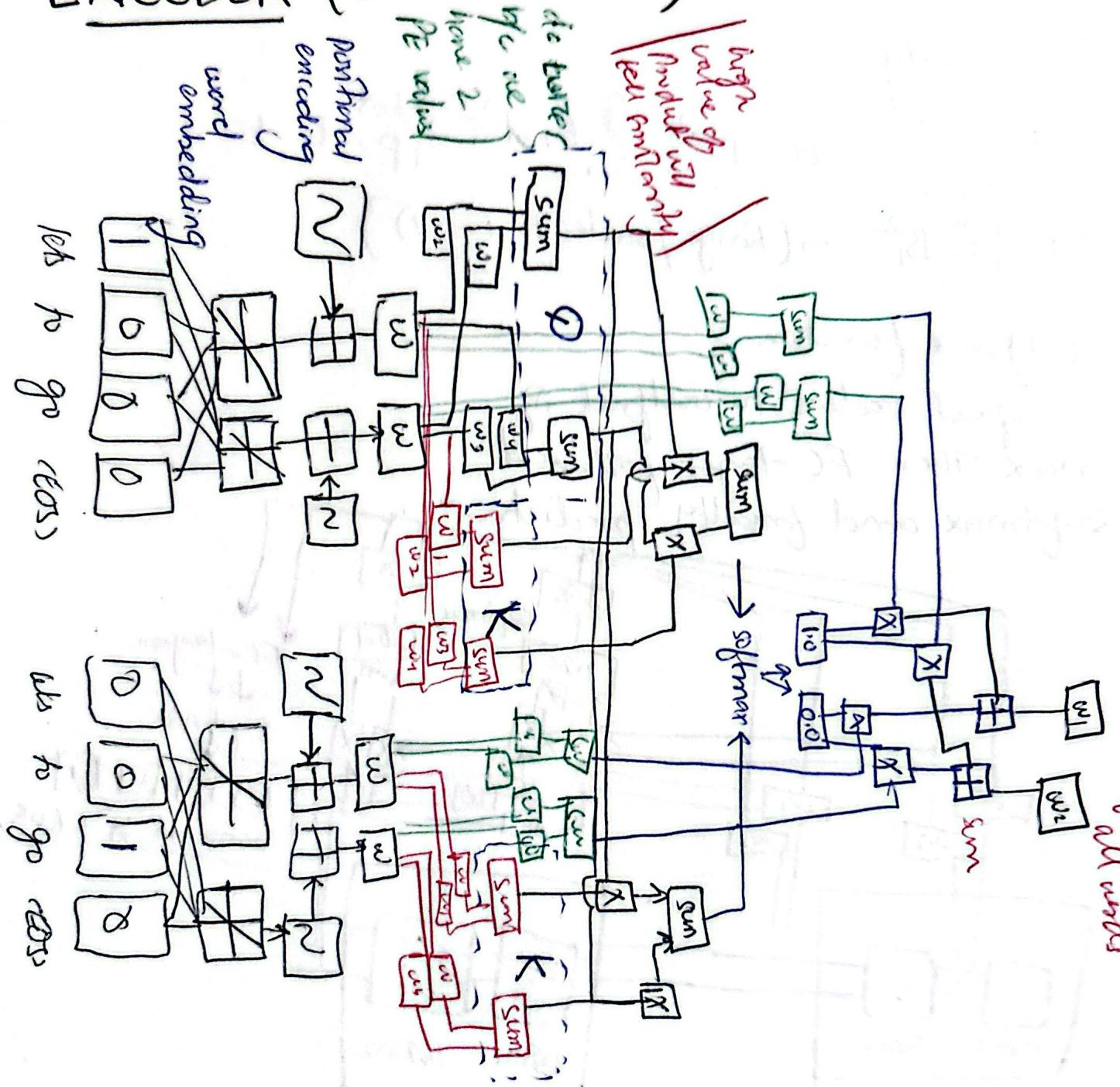
Then input value to output of softmax then FC-layer followed by softmax and finally prediction.



# TRANSFORMER ARCHITECTURE

By Statquest. (examples: lets go → varmos) <sup>trans</sup>

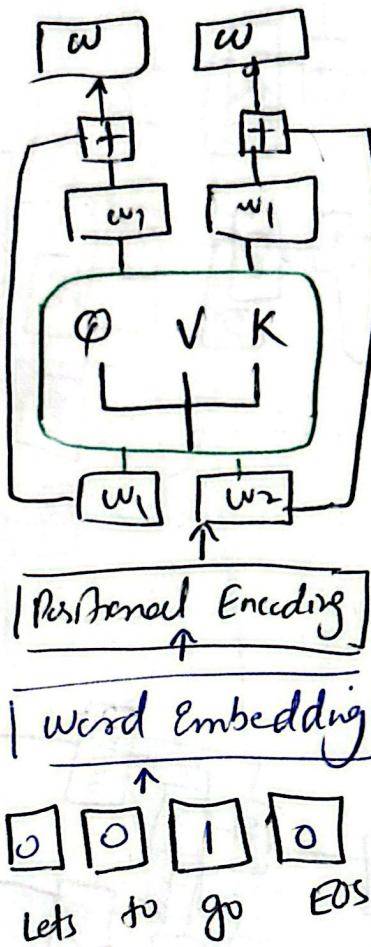
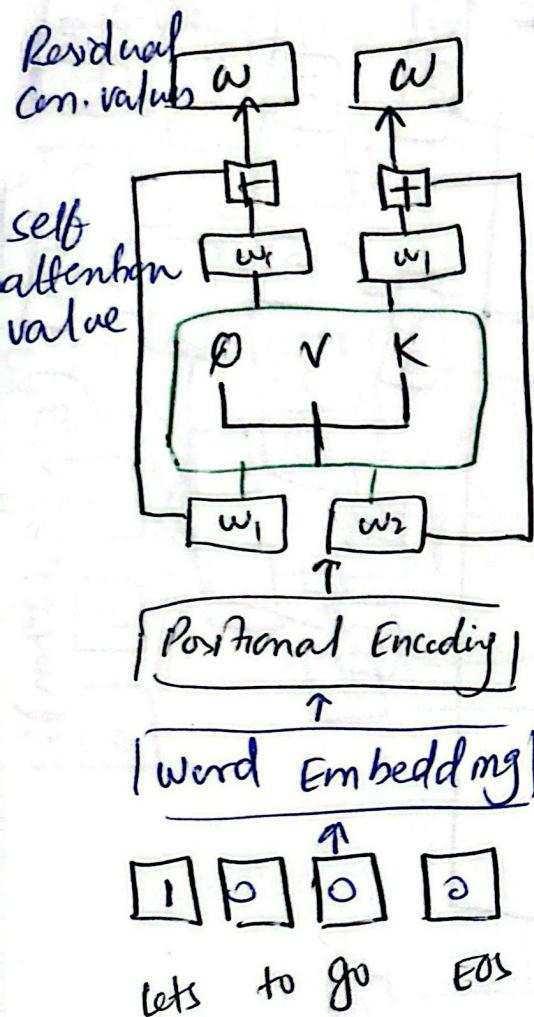
ENCODER (SELF-ATTENTION) for each words



similarly  
for all words

## 3- ARCHITECTURE (CONTINUED)

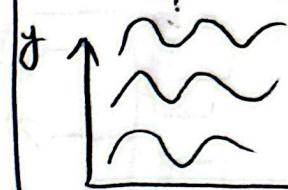
### ENCODER:



### Residual connections

To help in training, & to establish relationship between input & output (without having to preserve) encoding values

FOR Positional Encoding.



for  $n$  (words), we pick  $y$  value and add them to encode. (To keep track of the order of words)

### Self Attention -

To encode the relationship among words.

(used same weights  $Q, V, K$  to get self attention for all words)

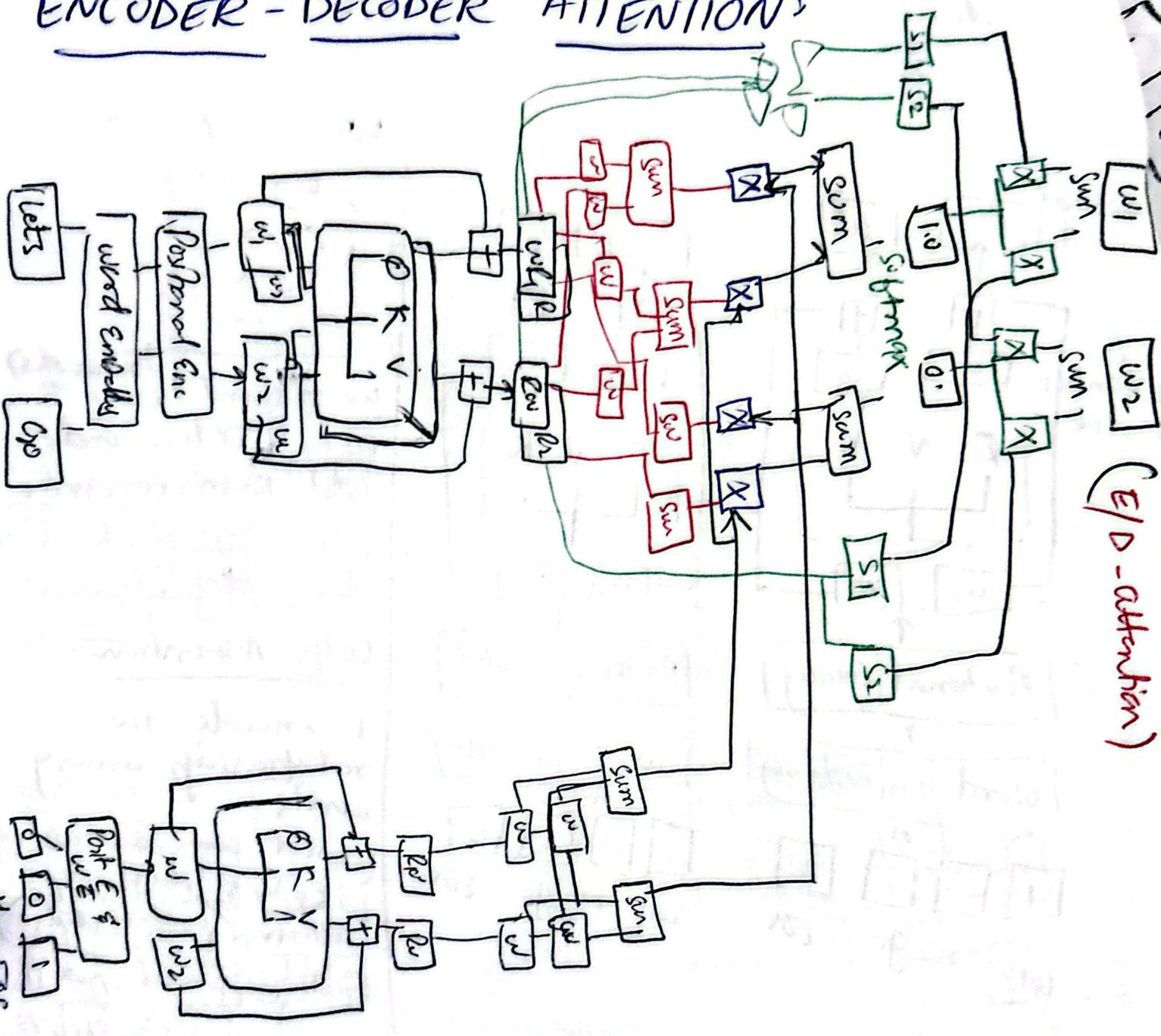
But different for  $St$

(Many stack = multi head)

Self attention have input from other words help get context.

# DECODER

## ENCODER - DECODER ATTENTION



(E/D - attention)

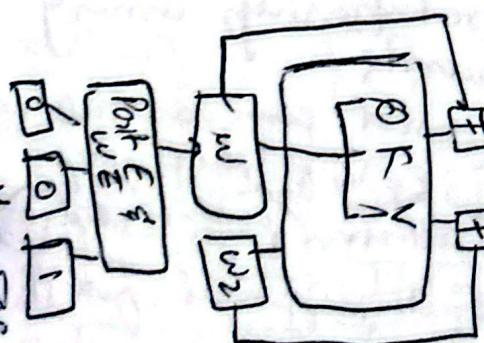
## Decoder

- Starts at (EOS).
- Decoder  $\Phi, K, V$  are different.

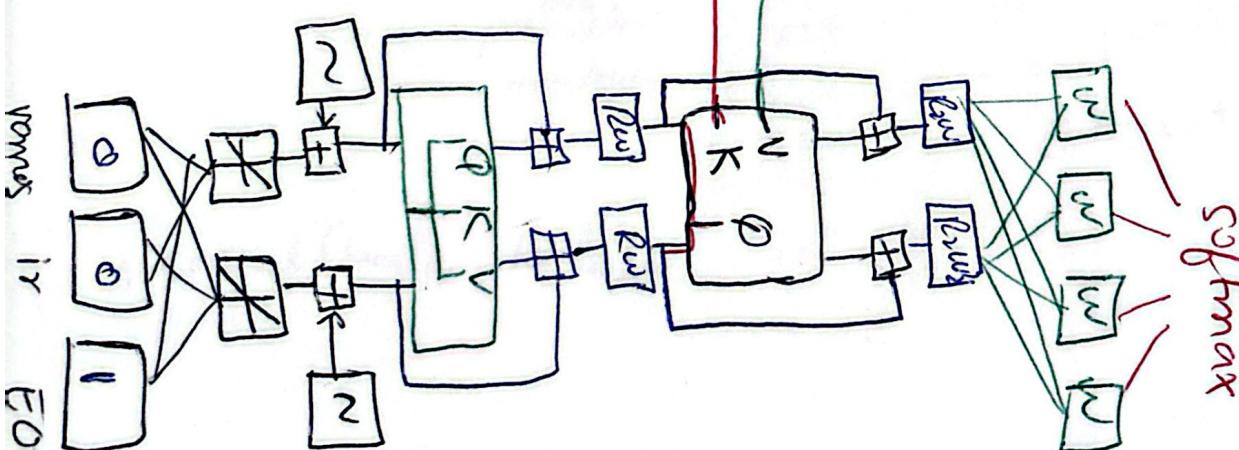
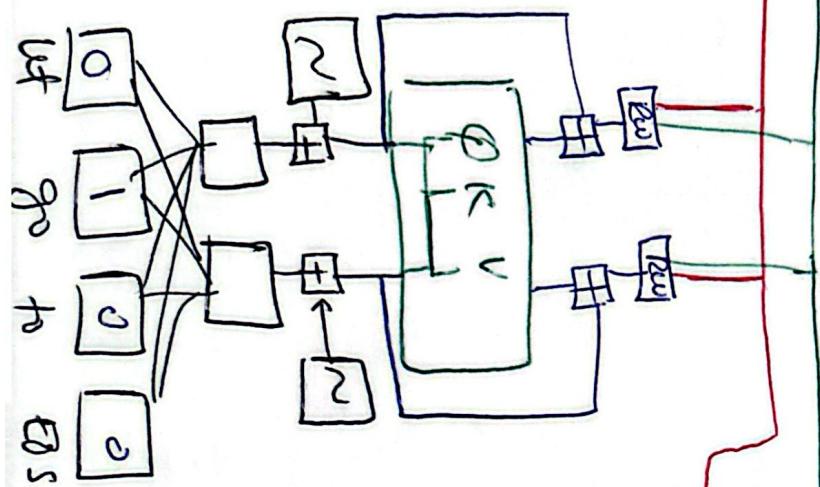
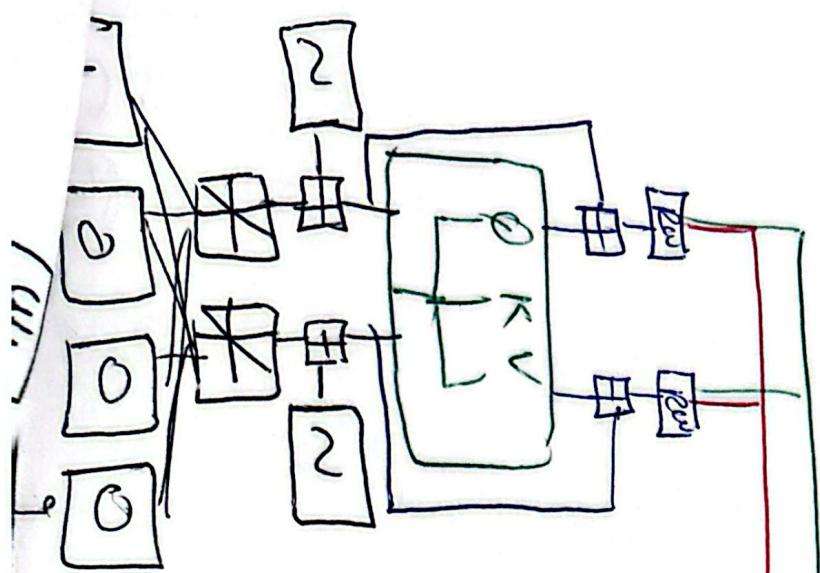
- E/D allows to keep track of significant words of an input and output.

- $E/D$ ,  $\Phi, K, V$  are also different but can be stacked.

lets  
word endibly  
Op



CONTINUED :-



next input to copy of Decol  
1 0 0  
Vemos or EOS  
softmax  
Fc-layer

10/16  
IDS  $\rightarrow$  Embedding 512 (learned)  $\rightarrow$  PE (fixed not learned)

PEs ~~odd~~ =  $\sin \frac{\text{pos}}{1000^{2i}}$  even  $\frac{\text{odd}}{\text{dimModel}} = \cos \cdot -$  (i = vector i)

why  $\sin/\cos$ ?  
our model can recognize as continuous  
and see the pattern  
self-attention. (each word to other in a sentence)

Attention =  $(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{dk}} \right) V$

eg = 6 word sentence  $dk = 512$

$6 \times [Q]^{512}, [K]^{512}, [V]^{512}$  /  $\sqrt{512}$  then softmax  $[6 \times 6]$

each value represent a product of 1st token w/ 2nd token  
each row sum to 1.

$[6 \times 6] \times [V]_{6,512} = [6, 512] \rightarrow$  each row capture many,  
position & interaction w/ other words.

same dimension  
as we  
started

# STATQUEST

Affusion is all you need !

## (1) Word Embedding

1) Positional encoding (To keep track of word order)  
 we add these values in embedding even if it similar, we'll get unique values.

[0.2]

+

[0.9] + [0.2]

=

[0.11]

even if we exchange the position of words the positional encoding value will be same, but result will change



for  $n$ , we words ( $x$ )  
get  $y$

## 3) Transformer (dependency)

by - self attention. (similarity with others and itself) for all some weight of  $(\phi, \theta, v)$

Residual  
help fixing,

Self attention  
have right for  
other and  
help context.

8-self attention

multi-heads (and,  
are needed  
in paragraph etc.)  
with own  
weights

encode  
Word to no,  
position of words,  
selection ong words,  
easy to turn in parallel  
E-D attention to keep  
track of things relative  $\text{IP} \rightarrow \text{OP}$

similarity =  $\frac{\text{DT Product}}{\# \text{Enthely values per dt}}$

~~to establish R/W~~  
~~not used~~ without also  
preserve the inherent PE  
info

Feb 21.

1) Permutation invariance (if we ignore PE) { mean if we change position of words their value will remain exact, just changes position}

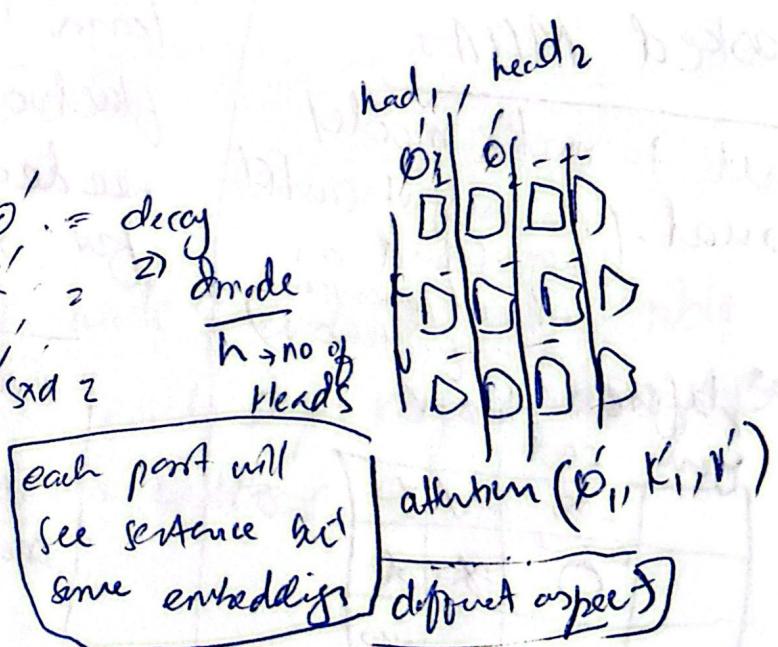
2) self attention needs no parameter (except embeddings)

3) we expect value along diagonal to be max (as its dot with itself) and if we want word not to interact, before softmax we replace it with  $\text{softmax}(x)$ .

## Multi-head attention

Encoder:  $(6 \times 512)$   
input

$$\begin{array}{l} \text{input} \\ \boxed{\begin{array}{l} Q \times w^Q \\ K \times w^K \\ V \times w^V \end{array}} = \begin{array}{l} Q' \\ K' \\ V' \end{array} \end{array} \quad \begin{array}{l} \text{parameters} \\ \text{matrices} \end{array}$$



$\Rightarrow \text{concat}(\text{head}_1, \dots, \text{head}_n) W^O$

$$\begin{array}{c} \text{seq, } h \times d_V \\ \text{seq, } h \times d_{\text{model}} \\ 6 \times 512 \end{array} \times \begin{array}{c} W_O \\ h \times d_{\text{model}} \\ h \times d_{\text{model}} \end{array}$$

$$\Rightarrow \begin{array}{c} \text{MHA} \\ \text{seq, } r \text{ d_mdl} \end{array}$$

→ OK → softmax - visualize attention.

Q, K, V:

Add&Norm:

Layer Normalization (we calculate mean & variance of each item independently) & replace by

$$\hat{x}_j = \frac{x_j - \bar{x}_j}{\sqrt{\sigma_j^2 + \epsilon}} \quad (\text{to range } (0-1))$$

Decoder:

(cross-MHA.)  $\xrightarrow{K^T_{\text{Query}}$  Encoder}

Masked MHAs

Want to make model causal. (the output should only depend on current/previous, not future)

so ~~before softmax~~

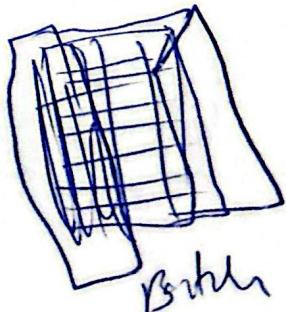
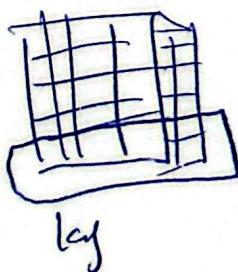
	you	cat
you	-∞	-∞
cat	cancel	0
is		
run		∞

we also multiply by (gamma) & add beta. ( $\gamma$  &  $\beta$  are learnable) to learn these value and introduce fluctuations (amplified or not) when needed,  $(0-1)$  may be too restrictive for network.

Batch vs layer

layer Norm, treating each item independently

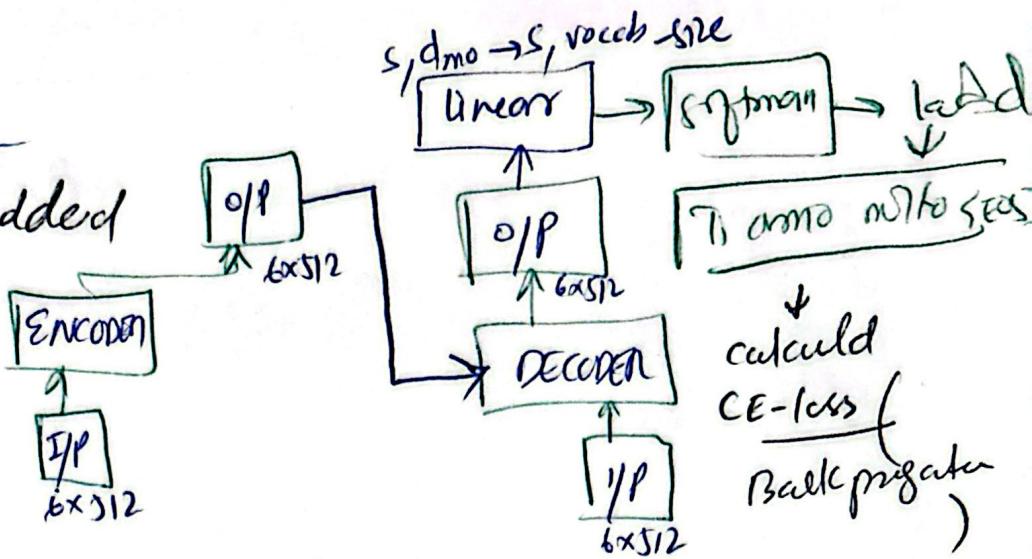
• Batch = using for same premises



## Inference & Training

Decoder input is padded

<SOS> --> <EOS>  
<SOS> -----



This all happen in all 1-time step  
while RNN - N-step for N-length sequence

## Inference:

Some → (output = logits)

But not in Greedy, by token-token ( $t_i = 1$ ,  $a_{m^0} = 2$ )

Greedy: at every time step we select word with max softmax value

## Beam Search:

at each step select top  $B$  words & evaluate all possible next words for each of them and at each step keeping the top  $B$  most probable sequences. This performs better.