

# LSTM (Hochreiter & Schmidhuber)

[ Long Short Term Memory ]

ANN

RNN

→ Disadvantages ←

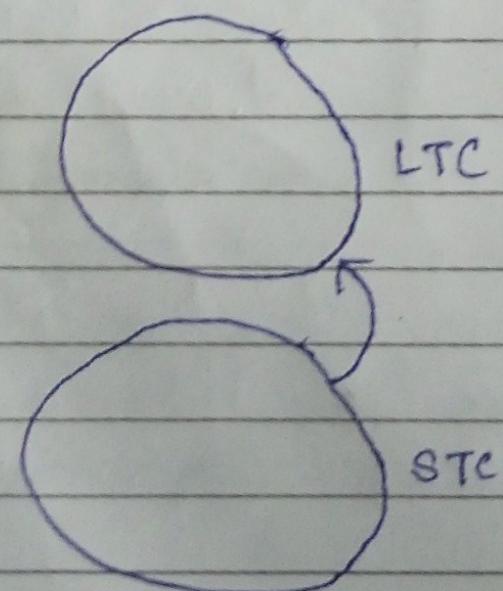
(i) It's not useful  
for sequential  
data.

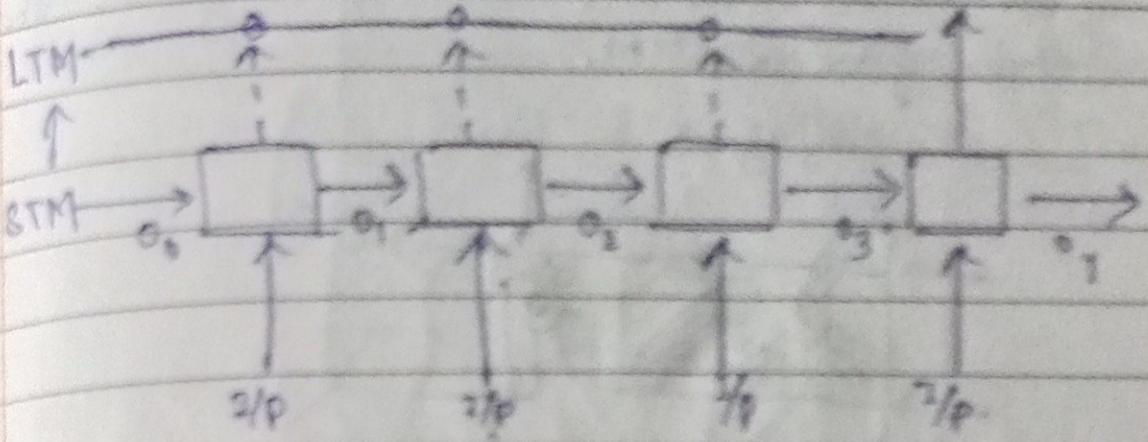
(ii) Hidden unit used to return  
feedback that propagate  
with other nodes. where  
time stamp plays a vital  
role. When the chain become

Vanishing  
Gradient

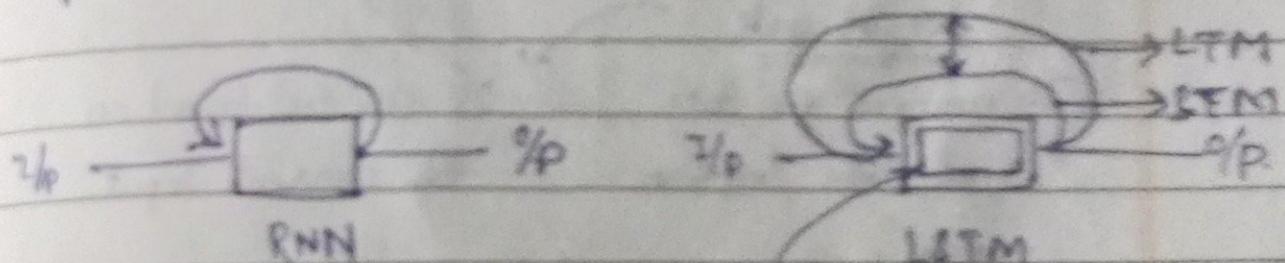
long, then the machine (or)  
forgot the old data.

When we deriving a story percephone, process  
word by word, that maintain two context.

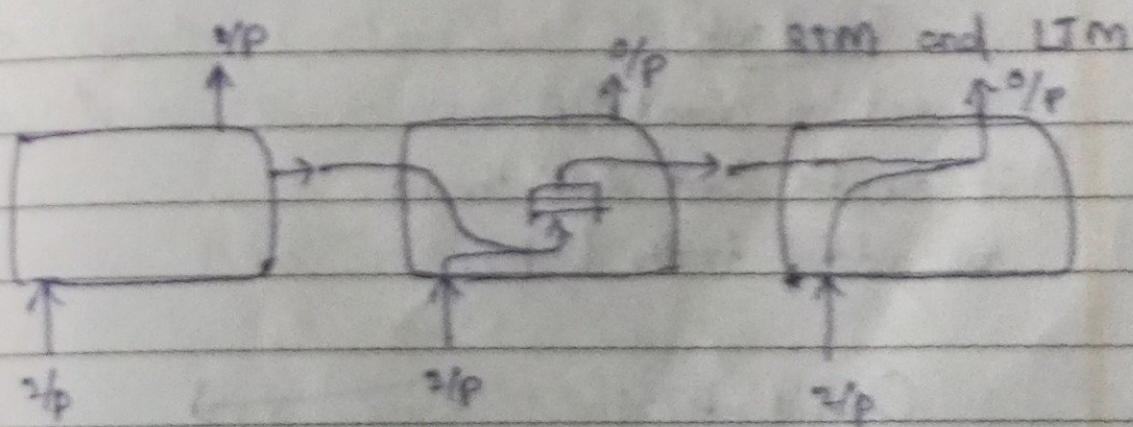




If you find any datablock or perception that is important in long term then we move that data from STM to LTM.



Communication between

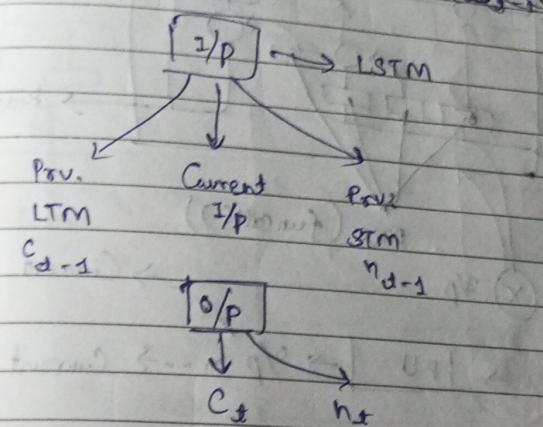


(RNN) architecture

Where we can find a Short Term memory is there which is there to

How they are working →

During input we provide previous LTM + previous time step i.e.  $(c_{t-1})$  and STM as  $(n_{t-1})$  (at time  $t$ )



Processing →

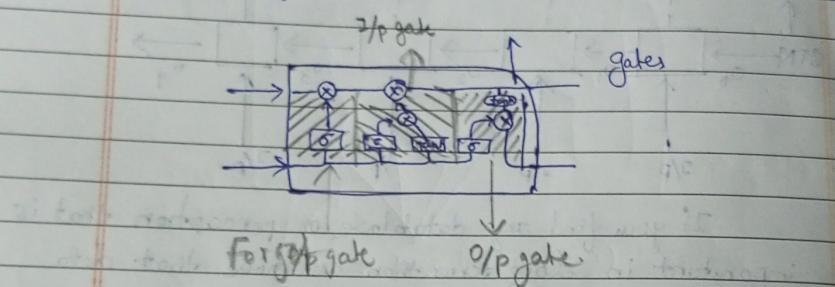
- Updation of Long term memory
- Create short term memory.

Now I get it →



What is LSTM.

The architecture of LSTM is called as gates, there are three gates in LSTM.



Forget gate  
It removes the data from the long term memory.

If decided based on the current Input.

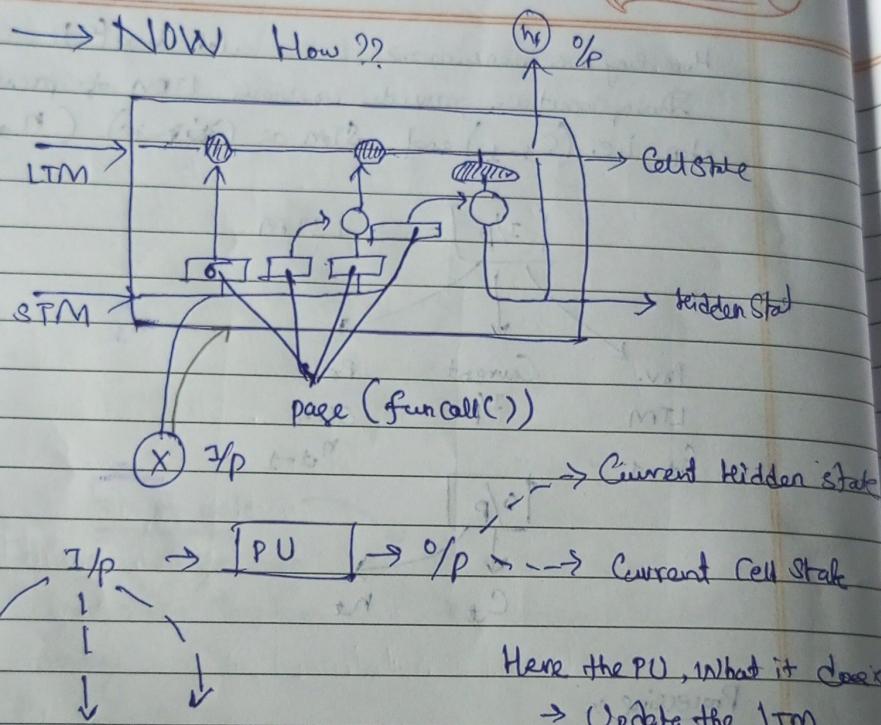
Input gate  
Based on current I/P  
it decides what data should be added to the LTM.

Output gate :-  
Based on current input it returns the data based on long term memory.

It not only returns output but also Create a short term memory.

fun call( )

Date \_\_\_\_\_  
Page \_\_\_\_\_



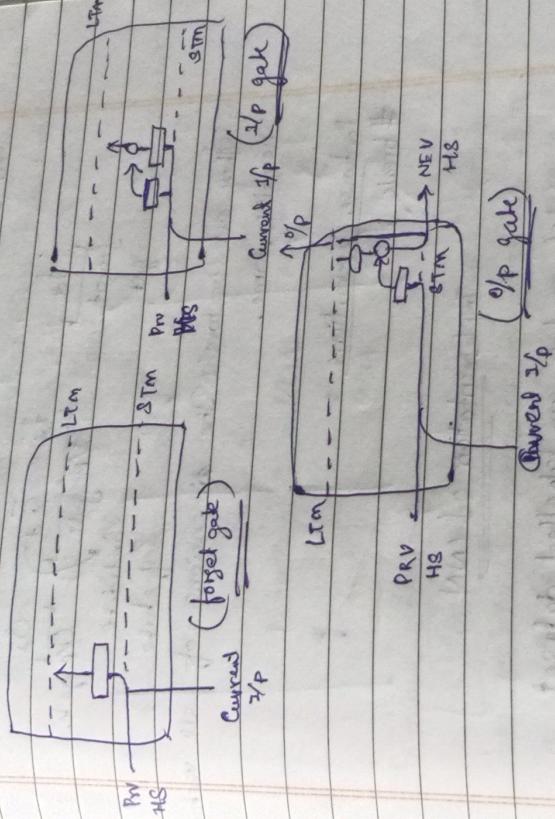
Present hidden state → Present Cell State → PU for the current time step ( $x_t$ ) → Calculate the hidden state / STM

While going from  $C_t$  →  $C_{t+1}$ .

(i) Based on the current input  $x_t$ , we are going to decide which data is to remove.

(ii) Based on the current input  $x_t$ , we are going to add the necessary data.

The Gates (LSTM architecture) is based on three gates in below figure.  $\Rightarrow$



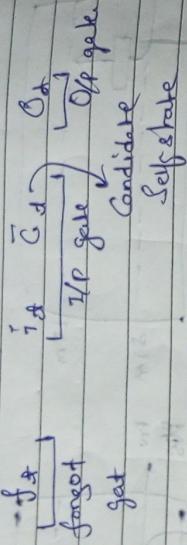
Forget gate need to remove something / data from cell state. If gate is there to add something to the cell state. If it is used to calculate the hidden & cell for the next time step.

Note:-

(Cell State) (Hidden State)	The dimensions of C & H.S will be same.
C H.S	Same

Vectors

$X_i$  is also vector / list of numbers. It's dimension depends on the data. Let's say more 4 gather.



They are vectors. The dimensions of all vectors remain same.

Pointwise operations :-

```
#break;
```

### Implementation of CNN

- We need a unlabelled dataset
- We make a dataset and fit the data with the dataset to form a dataset frame
- Split the data, train and test
- Convert train, test → train generator and test generator. These are the method which generates batch of data.

Here we modify some hyperparameter i.e  
image-size, batch-size, specify the column

→ Creation of Model.

Model is created using CovAD, unpelling 3D Dropout, Flatten, Dense, activation, Batch Normalisation.

→ Training of the model.

→ using the train-generator and test-generator we train & validate our data.

→ Visualize the accuracy by training loss and accuracy plot.

→ testing the model.

♦

Pointwise operations →

The pointwise operations are execute in between two vectors. The pointwise operators can be a

→  $\oplus$

→  $\otimes$  ←

→  $\tanh$

(x) : let a vector  $C_{3 \times 3} = [4 \ 5 \ 6]^{T} 1 \times 3$

$f_{3 \times 3} = [1 \ 2 \ 3]^{T} 1 \times 3$

%

↓

$\otimes \rightarrow [4 \ 10 \ 18]$

$\oplus \rightarrow [5 \ 7 \ 9]$

$[\tanh(x) \ \tanh(y) \ \tanh(z)]$

(+) :

(tan) :

## (Go to) function()

In this page

$\delta \rightarrow$  represents a neural network layer  
and it's layer containing number of neurons  
and each of them containing with an activation function  
(Sigmoid activation function).

TYPE	Activation
$3, 4, 5 \boxed{\delta} \rightarrow$ ANN	Sigmoid
$3, 4, 5 \boxed{\delta} \rightarrow$ ANN	Sigmoid
$3, 4, 5 \boxed{\tanh} \rightarrow$ ANN	tanh

Containing Let's consider a hyperparameter  
Sam No of neurons

e.g. No. of neuron.

$\rightarrow$  user will provide it setting up  $\delta$

And the hyperparameter will be user defined

universal one if there is a change triggered  
occurred in the 1st layer then it will  
affect all activation function.

$$[x_1, x_2, \dots, x_n] \leftrightarrow (\delta)$$

$$[x_1, x_2, \dots, x_n] \leftrightarrow (\delta)$$

$$[x_1, x_2, \dots, x_n] \leftrightarrow (\delta)$$

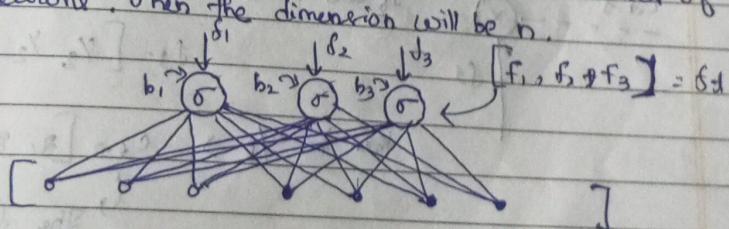
## The forget gate:

forget gate is a container / get three i/p as arguments i.e.  $c_{t-1}, h_{t-1}, x_t$

If give a o/p i.e. will remove information from cell state.

Here [n dimension] of  $h_{t-1}$  and  $c_{t-1}$ , where the no. of neurone will be n. i.e.

dim  
the product of  $h_{t-1}, c_{t-1}$  dependent on the number of neurone. if there are n numbers of neurone. Then the dimension will be n.



So there are all weight and each node containing one bias, i.e. 3 biases are there.

$$\text{Here dim, } W_f = (7 \times 1)$$

$$h_{t-1} = (3 \times 7)$$

$$b_f = (3 \times 1)$$

$$f_f = (3 \times 1)$$

$$f_t = \sigma(w_s [h_{t-1}, x_t] + b_s)$$

$3 \times 9 \quad 3 \times 1 \quad 1$

$3 \times 1$

$3 \times 2$

again there is a pointwise operator  $\circ$

$$f_s \circ g_{t-1}$$

$\times$

How this sign is really removing data from long term memory.

If we take  $c_{t-1} = [4, 5, 6]$

$$d_s = [y_1, y_2, y_3]$$

then the  $\circ$   $\rightarrow c_t = [2, 2.5, 3]$

The long term memory is cleared.

When we use  $[0, 0, 0]$  then the long term memory will be empty.

(ex. 1)

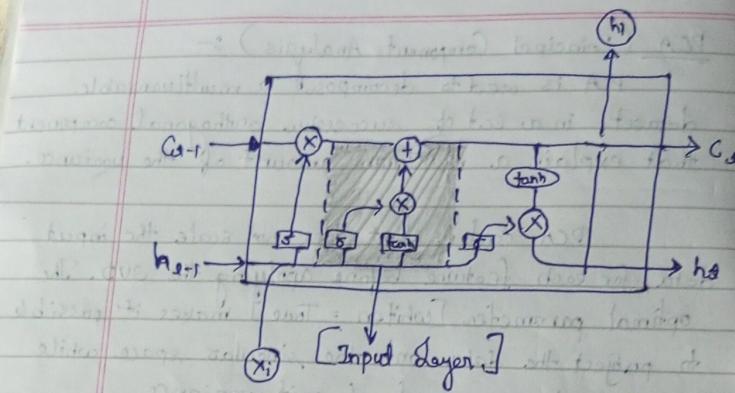
(ex. 2)

Contd..... (LSTM)

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_



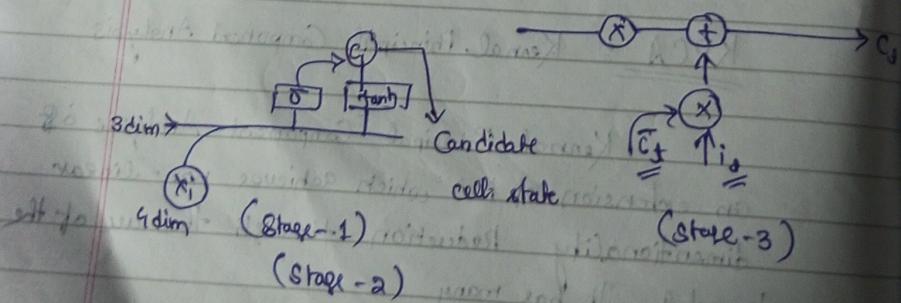
When we input data into the layer it checks whether it's important or not. If it is, then it will move to the cell state.

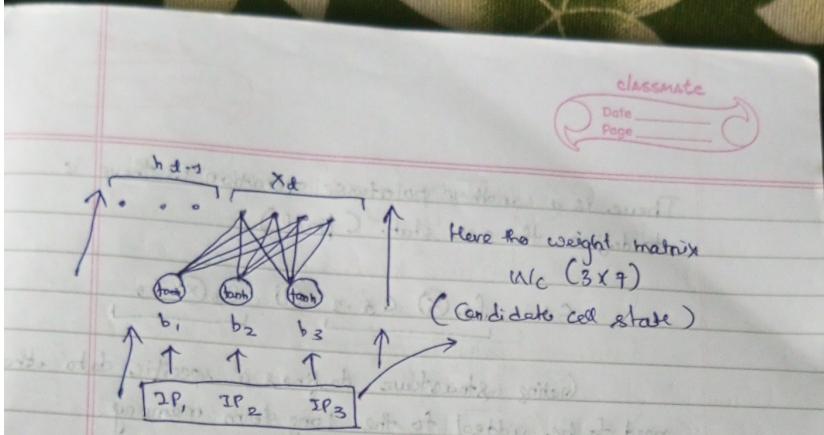
It's working under 3 stages -

1) Candidate cell state ( $\tilde{C}_t$ ) is formed.

2) ( $i_t$ ) decides which cell is / cell state is going to add to the Long term memory.

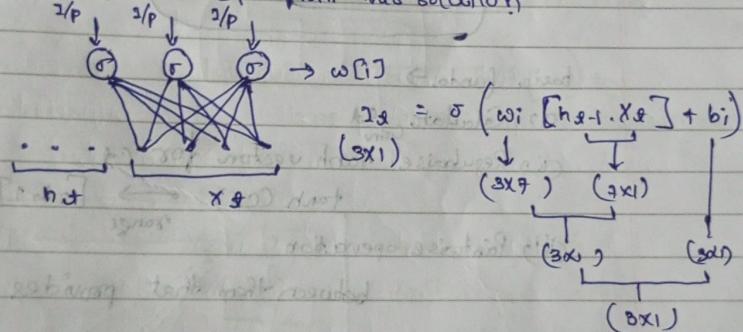
3) ( $c_t$ ) finds out to get/update the Long term context.





$$\bar{c}_d = \tanh(w_c [h_{d-1} \cdot x_d] + b_c)$$

Now we find out the o/p then we will add the o/p.  
which is returned from this solution.



There is a pointwise operation in between  $\bar{c}_d$  and  $i_g$ . See to the graph.

pointwise,

$$i_g \otimes \bar{c}_d \rightarrow \tilde{c}_d \quad (\text{filtered candidate cell state})$$

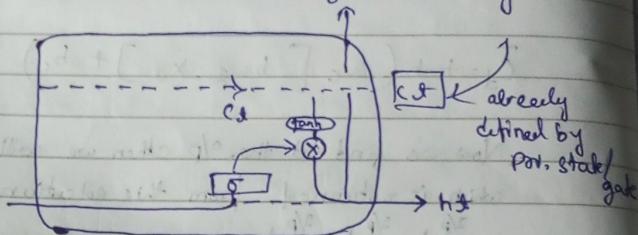
$\begin{matrix} (3 \times 1) & (3 \times 1) & (3 \times 1) \\ \downarrow & \downarrow & \downarrow \\ (2 \times 1) & & \end{matrix}$

$$(i_g \otimes \bar{c}_d) \oplus o_g = c_d$$

There is another pointwise operation is there is to define the cell state (final)

$$c_t = f_t \underbrace{\otimes c_{t-1}}_{\text{gate}} + \underbrace{\bar{c}_t \otimes i_t}_{\text{gate}}$$

Gating structure defines a specific data that need to be added to the long term memory.

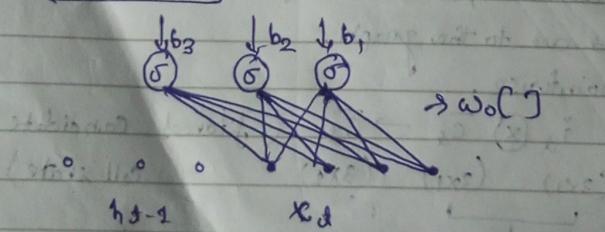


basic funda  $\rightarrow$

- (i)  $\rightarrow$  Calculate  $c_t$
- (ii) Regularise tanh vector for  $c_t$   
 $\tanh(c_t) \rightarrow [-1, 1]$
- (iii) Pointwise operation

between them that provides  $h_t$

$$\tanh(c_t)$$



$$o_t = \sigma(w_o [h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma \left( w_o [h_{t-1}, x_t] + b_o \right)$$

$3 \times 7$

$7 \times 1$

$$\sigma (3 \times 1) = 3 \times 1$$

$$to \ find \ h_t \rightarrow we \ need \ [o_t \otimes \ tanh(c_t)]$$

$(3 \times 1) \quad (3 \times 1)$

$\downarrow$   
 $(2 \times 1)$

then  $h_t$  containing  $(3 \times 1)$  /  $h_t$  value is found by  
us and we find that  $h_t$  containing  $(3 \times 1)$  matrix.

## (LSTM ARCHITECTURE) ↪

THE OVERALL THEORY BEHIND THIS

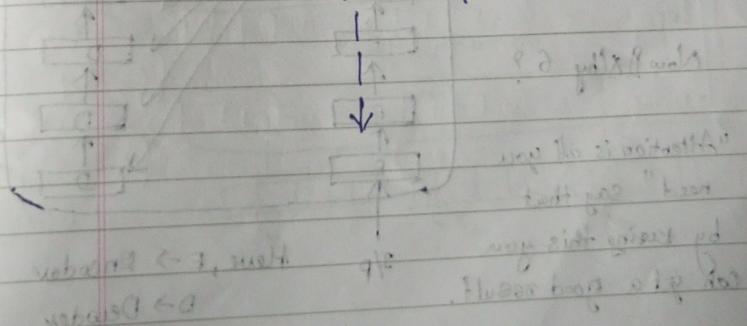
NOW

LET'S

MOVE

TO

TRANSFORMERS



classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$o_t = \sigma \left( w_o [h_{t-1}, x_t] + b_o \right)$$

$3 \times 7$        $7 \times 1$   
 $3 \times 1$        $3 \times 1$   
 $\sigma (3 \times 1) = 3 \times 1$

to find  $h_t \rightarrow$  we need  $[o_t \otimes \tanh(c_t)]$

$$(3 \times 1) \quad (3 \times 1)$$

$(3 \times 1)$

then  $h_t$  containing  $(3 \times 1)$  /  $h_t$  value is found by us and we find that  $h_t$  containing  $(3 \times 1)$  matrix.

(LSTM ARCHITECTURE) ↫

THE OVERALL THEORY BEHIND THIS