

Please note that in the practical we saw RNN and LSTM implementation using keras

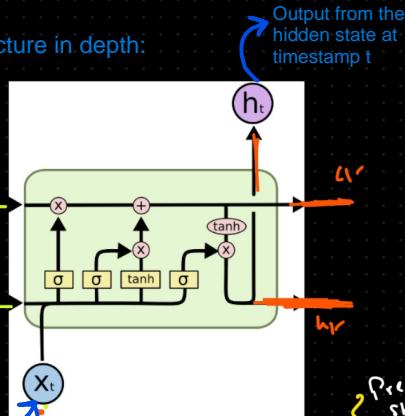
\Rightarrow LSTM

Decoding the LSTM architecture in depth:

Cell state(memory) from previous timestamp i.e 0

Output from hidden state from previous timestamp i.e 0

Input at t timestamp



Forward, Update, Output }

Forward invoke

current hidden state

Processing

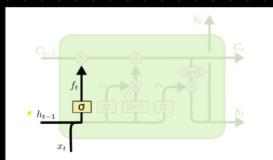
outPut

Current hidden state

- The LSTM cell does 3 processing
Forget, Update and Output the cell states by using 3 gates.

- It takes 3 input: Pre cell state, Pre hidden state and input from the current timestamp.

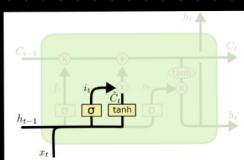
- It generates the 3 output: current hidden state, current cell state and third one is again a hidden state which is passed to the next LSTM cell as an input. Please note this 3rd output is utilized only if there exists another hidden state after the current cell.



Forget gate

To remove something from the C_t

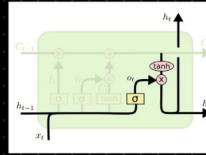
Cell state at t timestamp



Input gate

To add/update the C_t

Cell state at t timestamp

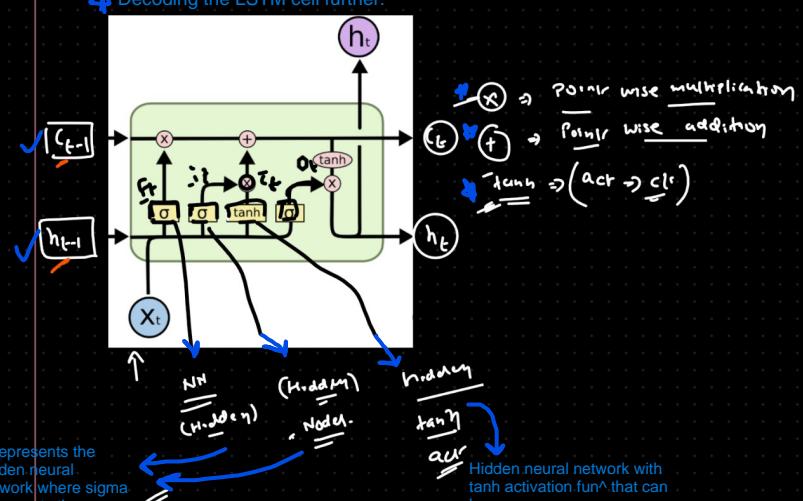


Output gate

Calculate the hidden
 $\{C_t / h_t\}$

Cell state and hidden state at timestamp

Decoding the LSTM cell further:



Hidden neural network with tanh activation fun^ that can have as many neurons as required

hidden vector
 C_t
Vector
 $[0.1 \ 0.2 \ 0.3]$

Important note on C_t and h_t :
Both of these are vector representation. Also, the size of both the vectors will be same.

$C_t = \{ \text{cell} \}$
 $h_t = \{ \text{hidden} \}$

same

- Represents the hidden neural network where sigma represents the sigmoid activation function (may be)

- Also, please note that the output from the first hidden NN is represented by F_t (where t represents the current timestamp) and output of second hidden NN is represented by it (Where t represents the current ts)

- Similarly, output of 3rd hidden NN is represented by

C_x
and for 4th hidden NN by
 O_t

$x_t \Rightarrow \text{1hot}_t$, on the current T_t which is also encoded vector representation of the textual data or any other sequence based data

Sentiment r

good bear better
bear bear better
work bear tool

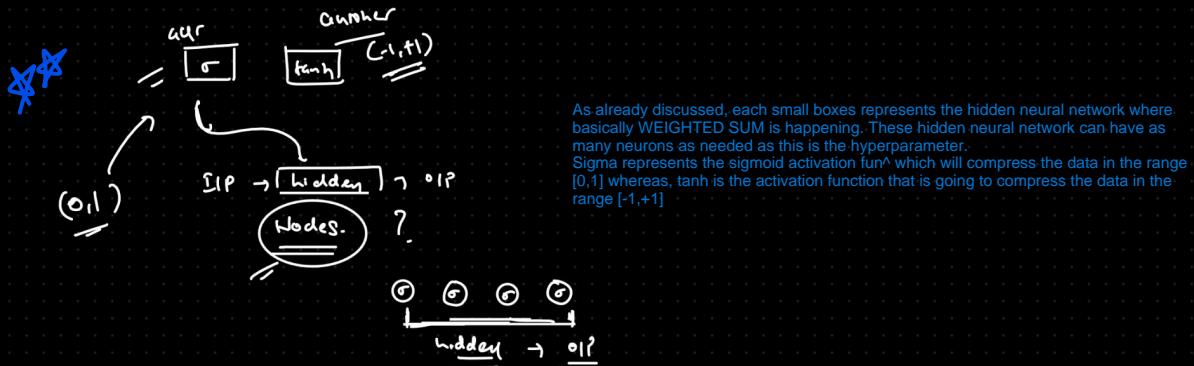
Vol

OME [good better better worse]

| | | |
|---|---|---|
| 0 | 1 | 0 |
|---|---|---|

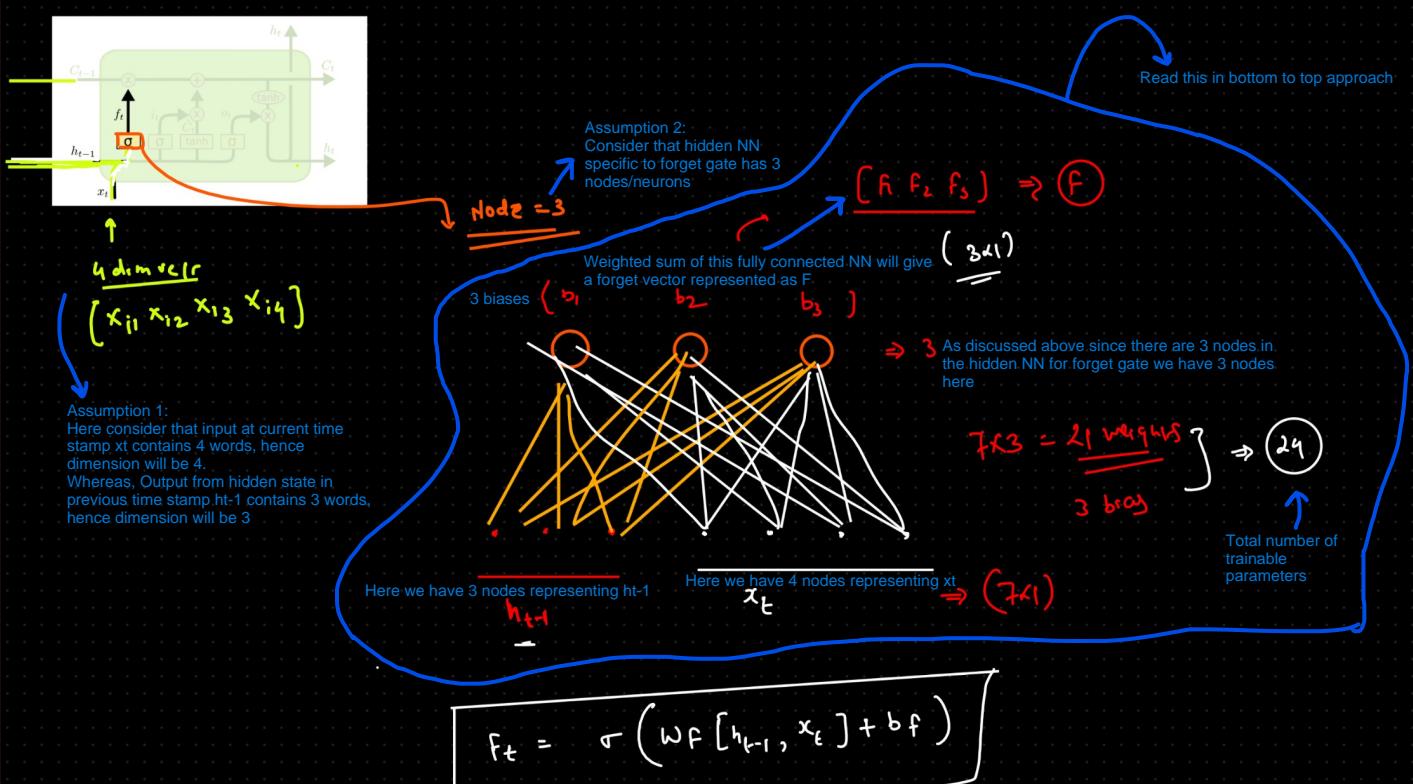
 $\left([1, 0, 0, 0] [0, 1, 0, 0] [0, 0, 1, 0] \right) \Rightarrow OME$
 $I_t \rightarrow x_{t+1} \rightarrow x_{t+2}$

word by word in lang



Decoding the Forget gate further:

Forget gate



Equation of forget gate vector. Where F_t represents the forget gate vector at t timestamp W_f represents the weight of forget gate, h_{t-1} represents the hidden state output at timestamp $t-1$, x_t represents the input at timestamp t and b_f represents the bias for the forget gate. This weighted sum is passed through the activation function represented by sigma

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Case where forget gate vector is a identity vector then in that case there will be no removal of the context being held by cell state in previous timestamp

$$[1 \ 1 \ 1] \Rightarrow [1 \ 1 \ 1]$$

removing \downarrow writing \uparrow

Understanding how forget gate at t time stamp (F_t) helps in removing the context being held in the previous state which is represented by cell state at $t-1$ timestamp (c_{t-1}).
 >>> This forget operation is implemented through point wise multiplication b/w F_t and C_{t-1}

$$[F_t \otimes C_{t-1}] \leftarrow$$

Case where forget gate vector is neither identity nor zero vector then in that case there will be some removal of the context being held by cell state in previous timestamp after applying pointwise multiplication

$$[1 \ 1 \ 1] \Rightarrow [1 \ 1 \ 1]$$

removing \downarrow writing \uparrow

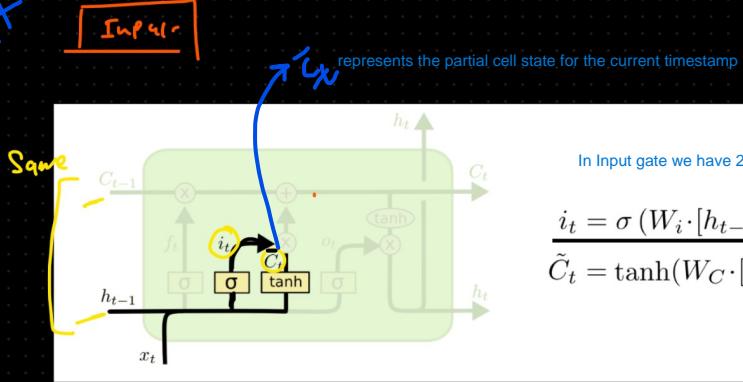
$$[4 \ 5 \ 6] \otimes [1 \ 1 \ 1] \Rightarrow [2, 2.5, 3]$$

$$[0 \ 0 \ 0] \otimes [1 \ 1 \ 1] \Rightarrow [0 \ 0 \ 0]$$

Case where forget gate vector is a 0 vector then in that case there will be 100% removal of the context being held by cell state in previous timestamp

Decoding the Input gate

Please note that both i_t and \tilde{C}_t are vectors of same size.
Similarly C_{t-1} and h_{t-1} represents the vector of same size.



In Input gate we have 2 equations:

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Try understanding these equations by observing the figure

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \rightarrow (3 \times 1)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \Rightarrow \text{Pointwise multiplication}$$

Firstly, we perform pointwise multiplication of these 2 equations

$$\tilde{C}_t = \begin{pmatrix} \cdot & \cdot & \cdot \end{pmatrix} \quad [T = \tanh]$$

$$\begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix} \quad \begin{matrix} T \\ T \\ T \end{matrix}$$

$$\begin{matrix} \times & \times & \times \\ 3 \times 1 & & 3 \times 1 \end{matrix} \quad \begin{matrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{matrix} \quad (7 \times 1)$$

$$7 \times 3 = 21 \rightarrow 21 \times 3 = 27$$

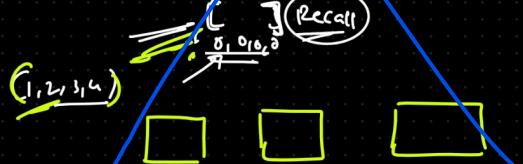
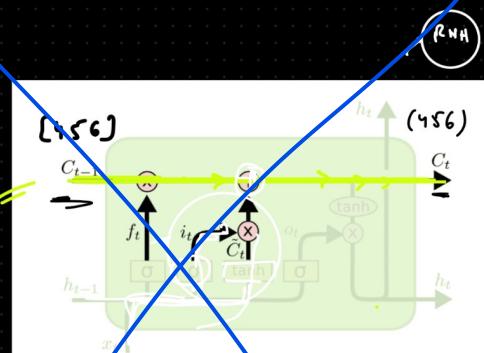
$$i_t \otimes \tilde{C}_t = \underline{\underline{\tilde{C}_t}}$$



Secondly, we will be performing the point wise addition operation b/w forget gate operation and input gate operation, whose outcome represents the cell state in t timestamp

$$C_t = \underbrace{f_t \otimes C_{t-1}}_{\text{forget gate operation}} + \underbrace{(i_t \otimes \tilde{C}_t)}_{\text{Input gate operation}} \Rightarrow \text{Current cell state } \downarrow t$$

To conclude unlike RNN, in LSTM we are able to forget the past context(fully or partially) using forget gate and add any new context to the cell state using input gate. This is how LSTM is able to hold long term context which was missing in the case of RNN.



~~$$C_{t-1} = [4 5 6] \quad f_t = [1 1]$$~~

~~$$i_t \otimes \tilde{C}_t = [0 0 0]$$~~

~~$$C_t = [4 5 6] \oplus [0 0 0]$$~~

~~$$= [4 5 6]$$~~

LSTM

~~Update~~

$[\quad]$

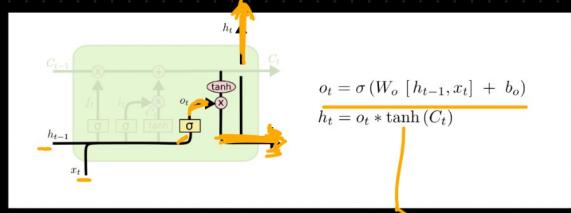
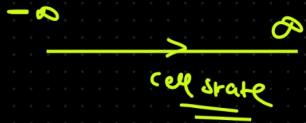


Decoding the Output gate:

Output gate

$$\left(\text{tanh}, \sigma_b + s \right) \circ$$

next hidden Stacked RNN



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

$$h_t = o_t * \tanh(c_t)$$

σ \Rightarrow hidden layer

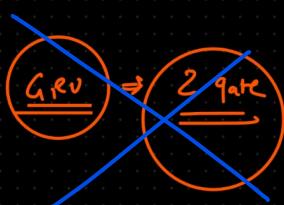
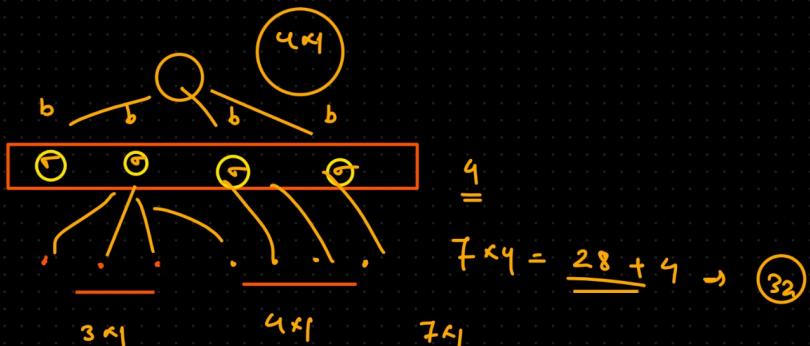
$$\Rightarrow h_t = o_t \otimes \tanh(c_t)$$

Here C_t denotes the cell state at t timestamp that we have updated in Input gate (Check C_t equation in input gate as discussed earlier for the better understanding)

x_t, h_{t-1}

$$\tau \left(W_0 [h_{t-1}, x_t] + b_0 \right)$$

(3×1) $\rightarrow (4 \times 1)$

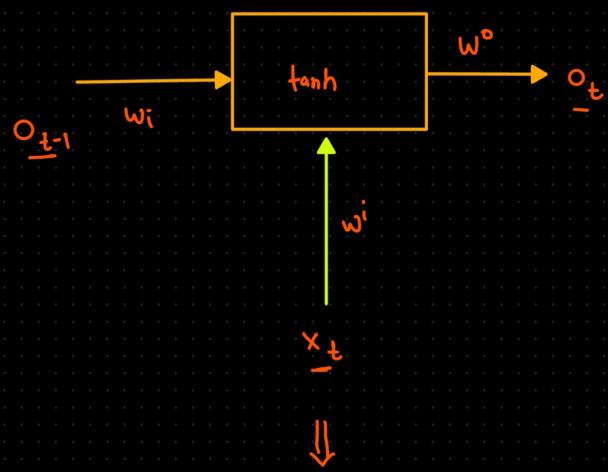




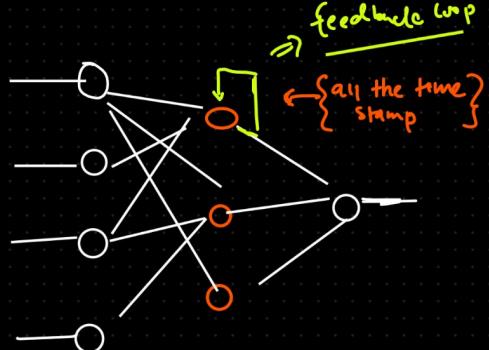
Deep RNN | Deep LSTM / Deep GRU

Understanding the basic RNN architecture:
Already discussed in the earlier session!

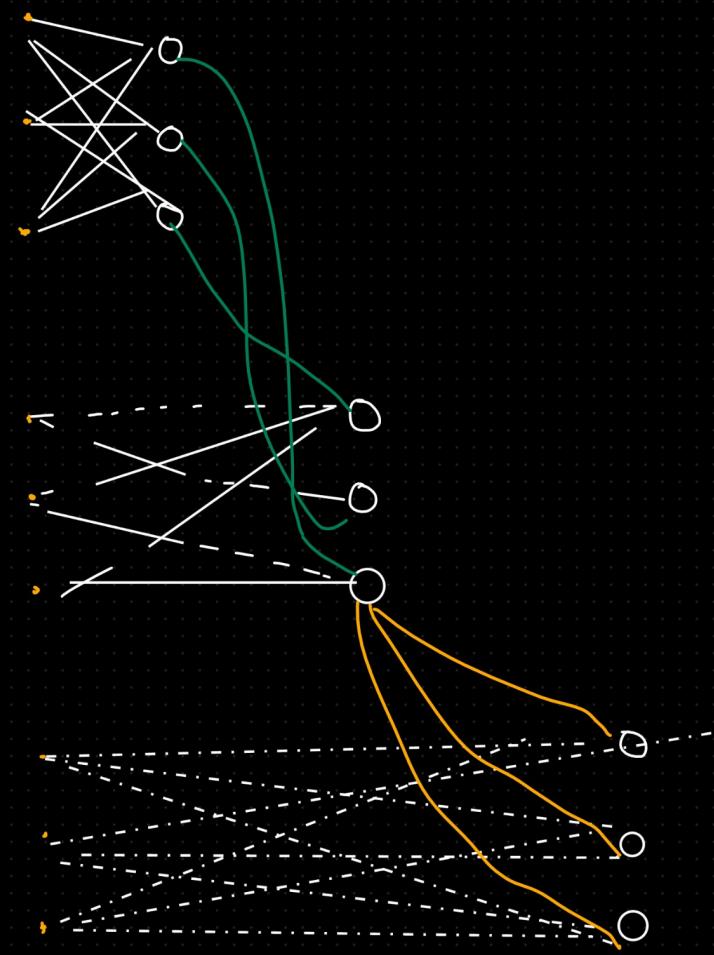
Feed back w/p



feed back w/p



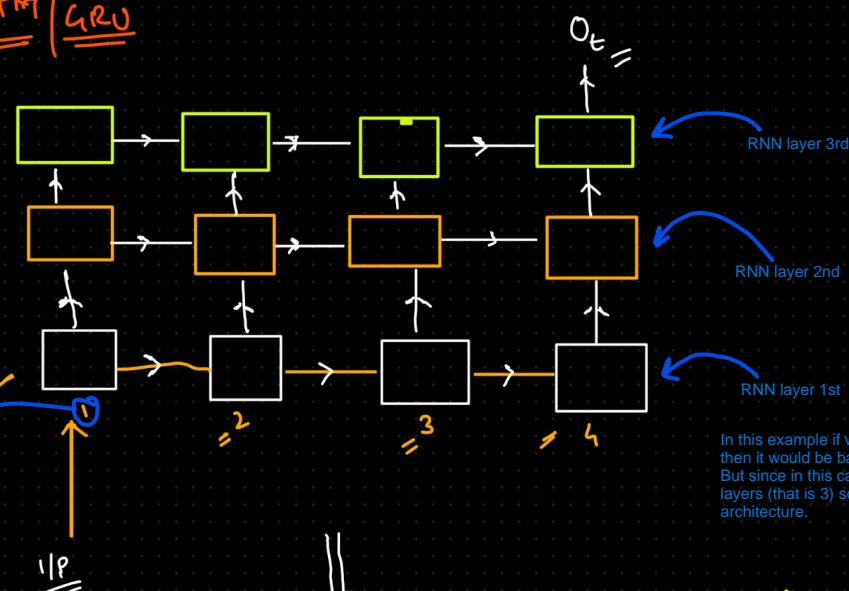
$t=1$



Deep RNN | LSTM | GRU

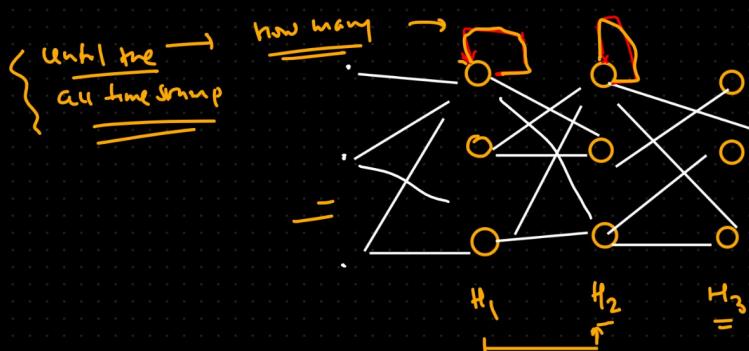
Hidden .

Deep means hidden layer more than 1. Therefore, we can say that Deep RNN is the RNN network where we have more than 1 hidden RNN layers.



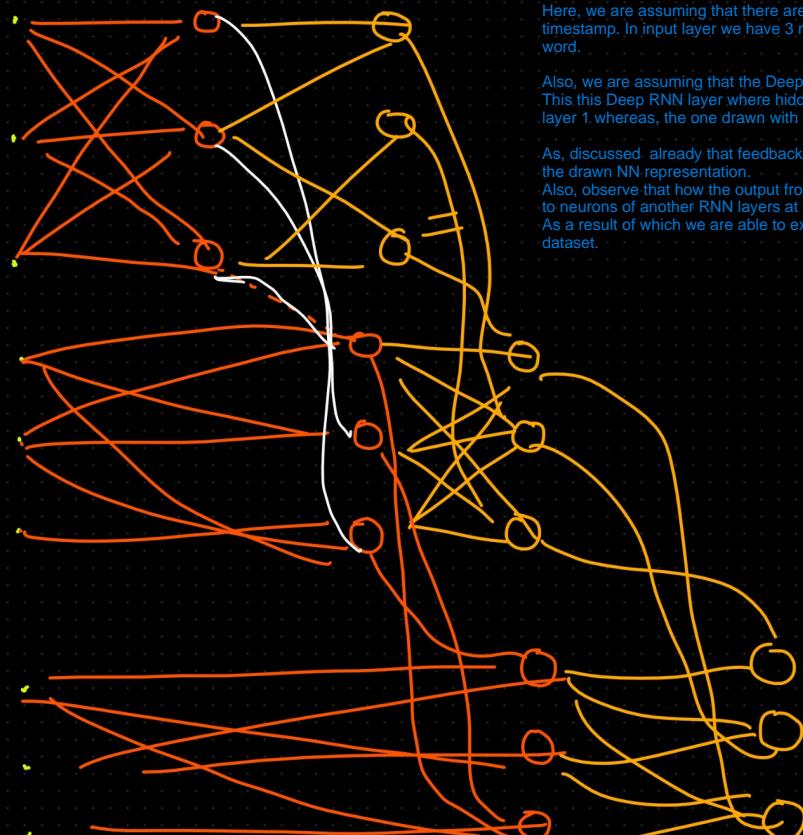
In this example if we were having only 1 RNN layer then it would be basic RNN architecture. But since in this case we are having multiple RNN layers (that is 3) so this is considered as Deep RNN architecture.

More than 1 hidden layer



Please note that here feedback loop will be passed to its own hidden layer. Meaning for layer 1 a feedback loop through a neuron will be passed to all the other neurons in the same hidden RNN layer 1 and not to neurons residing in the layer 2 or layer 3. Similar segregation happens in layer 2 and layer 3 as well.

Neural Network representation of Deep RNN:



Here, we are assuming that there are only 3 words in the sentence hence, we are provided with 3 timestamp. In input layer we have 3 nodes that will be capturing the vector representing the each word.

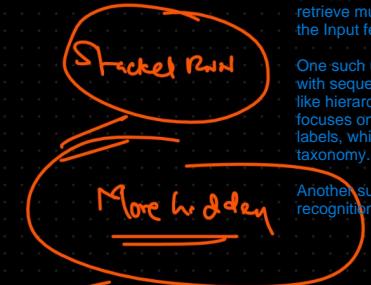
Also, we are assuming that the Deep RNN architecture consists of 2 RNN layer. This this Deep RNN layer where hidden nodes/neurons drawn using orange represents the RNN layer 1 whereas, the one drawn with yellow represents RNN Layer 2.

As, discussed already that feedback loop is passed to its own hidden layer. Observe the same in the drawn NN representation. Also, observe that how the output from nodes of hidden layers at each timestamp goes to neurons of another RNN layers at the same timestamp. This chain goes on for each timestamp. As a result of which we are able to extract more deeper context from the passed input sequential dataset.

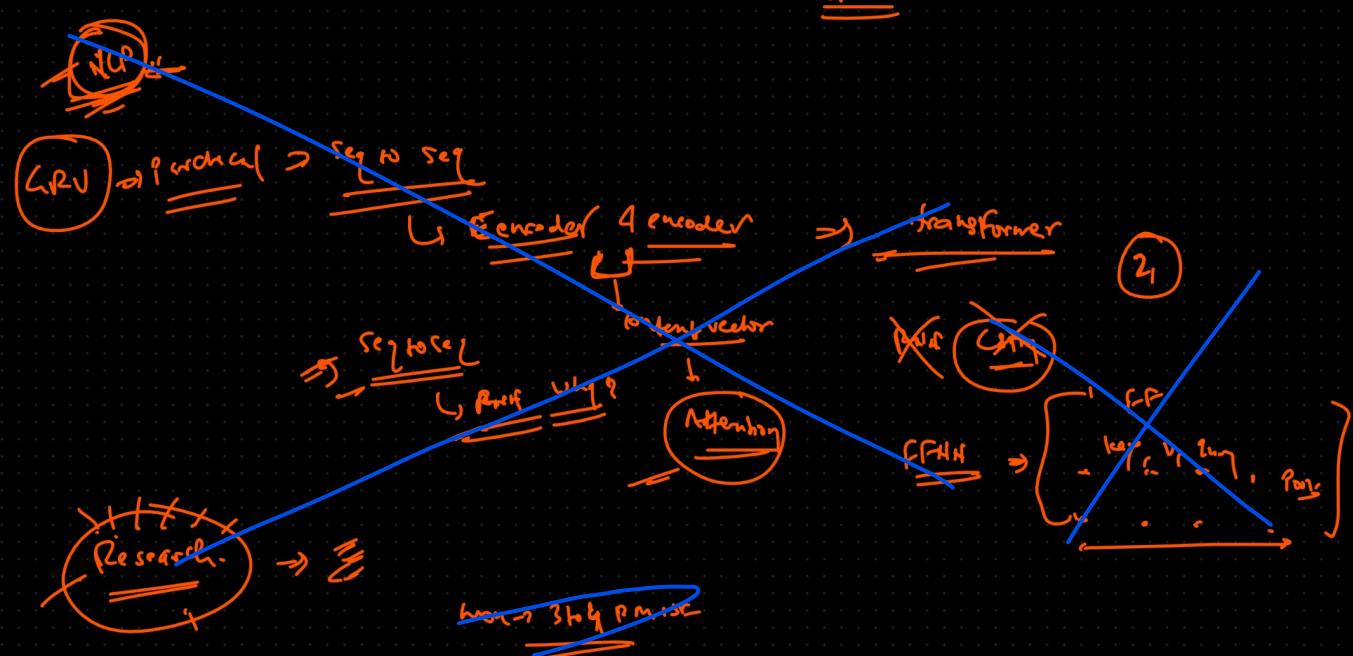
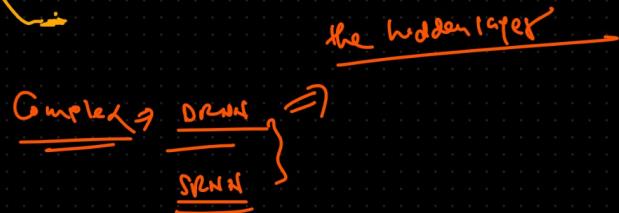
Please note that here we are using RNN cell. If we have used LSTM cell then it would be called Deep LSTM. Similarly, if we have used GRU cell then it would be called Deep GRU.

Why Deep RNN/Deep LSTM/Deep GRU:

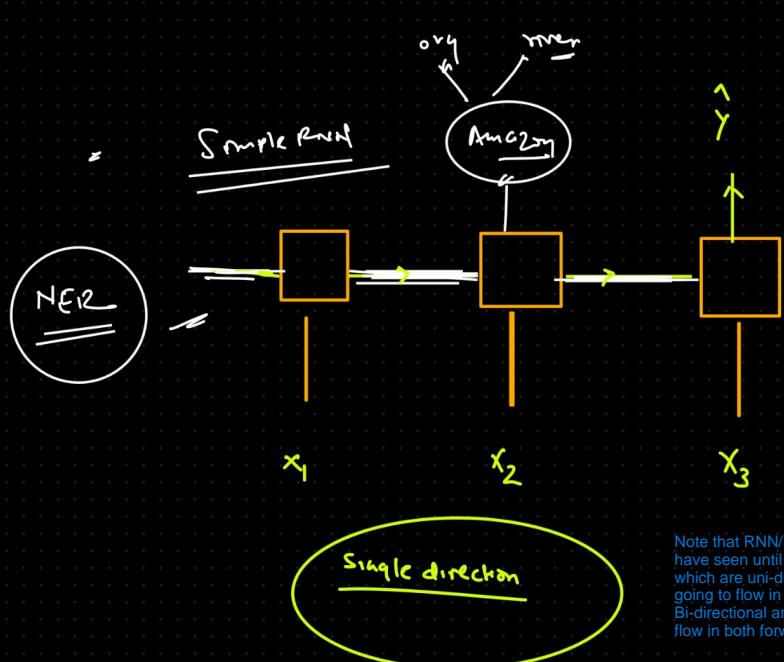
Why RNN/LSTM/GRU?



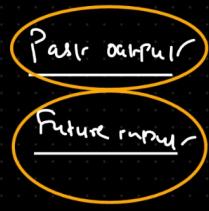
Stacked RNN/LSTM/GRU are same where also our main motive is to increase the hidden layers to capture complex patterns



Bi-directional RNN / Bi-directional LSTM / Bi-directional GRU



Bi-direc



Bi-directional RNN/LSTM/GRU is applicable in the case where Past Output is affected by Future Input or vice versa. Or in other words we need to traverse in both forward and backward directions in order to correctly determine the context of the word/words in the text or any sequence based data.

Note that RNN/LSTM/GRU structures that we have seen until now are simple or stacked RNN which are uni-directional where the information is going to flow in only one direction. Whereas, in Bi-directional architecture information is going to flow in both forward and backward direction.

The main reason to basically use mentioned deep architectures is to increase the number of hidden layers. As we increase the number of hidden layers we are able to capture or retrieve much more complex patterns from the input feed.

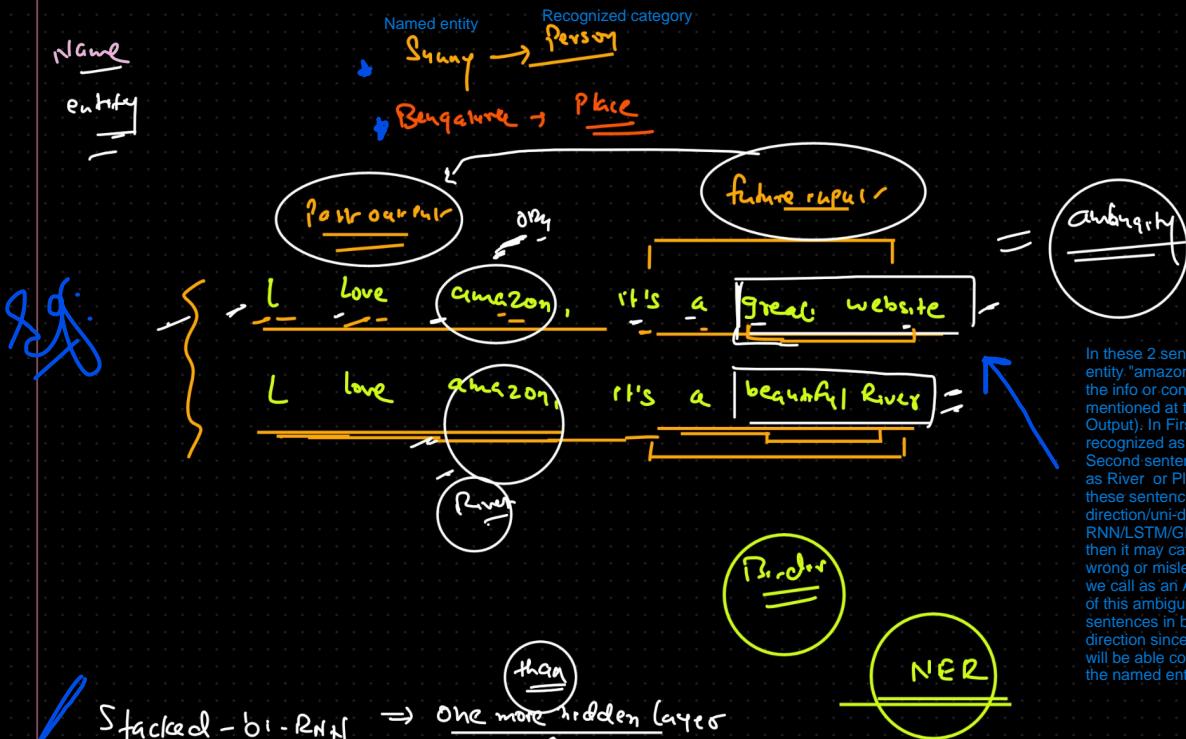
One such use case is when we are working with sequential data having some hierarchy like hierarchical text classification which focuses on classifying one text into multiple labels, which are organized as a hierarchical taxonomy.

Another such use case is the Speech recognition.

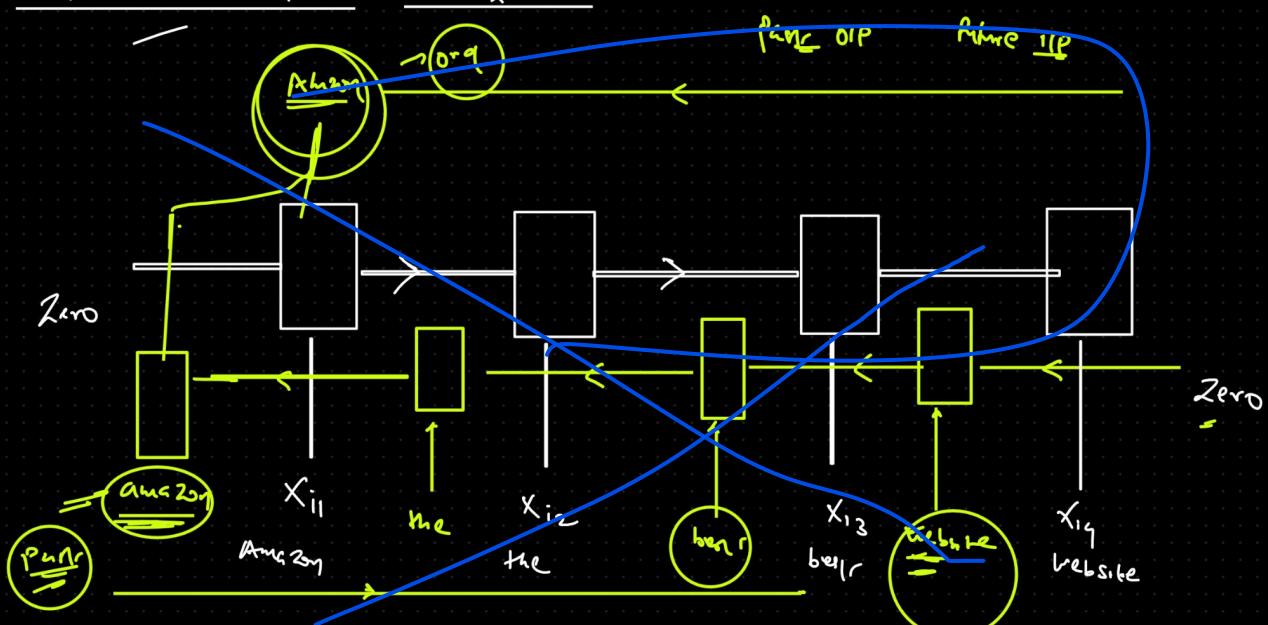
Let's Understand this Bidirectional concept using NER use case which stands for the named entity recognition

NER:

Named-entity recognition is a subtask of information extraction that seeks to locate and classify named entities mentioned in unstructured text into pre-defined categories such as person names, organizations, locations, medical codes, time expressions, quantities, monetary values, percentages, etc.



Stacked - bi - RNN \Rightarrow One more hidden layers



Understanding the Bi-directional architecture using an example:
 Sentence: Amazon is the best website.
 Problem statement is to correctly classify the category of named entity "Amazon"

