

⇒ LSTM RNN (Long Short Term Memory RNN):

- Problem with simple RNN:

RNN → can't handle Long term dependencies

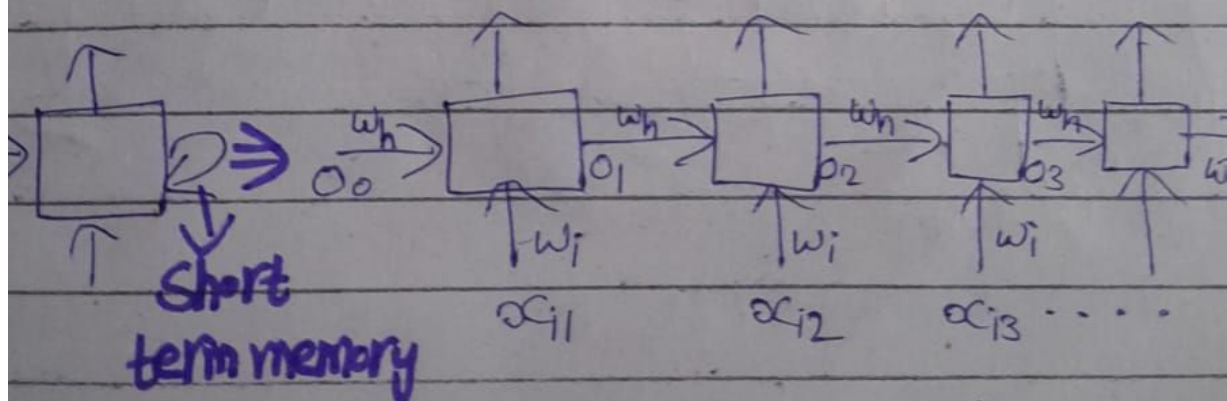
→ Vanishing Gradient Problem

• Basic Representation of LSTM

RNN:

Long Term memory

which remember the context till when it is required after that remove it

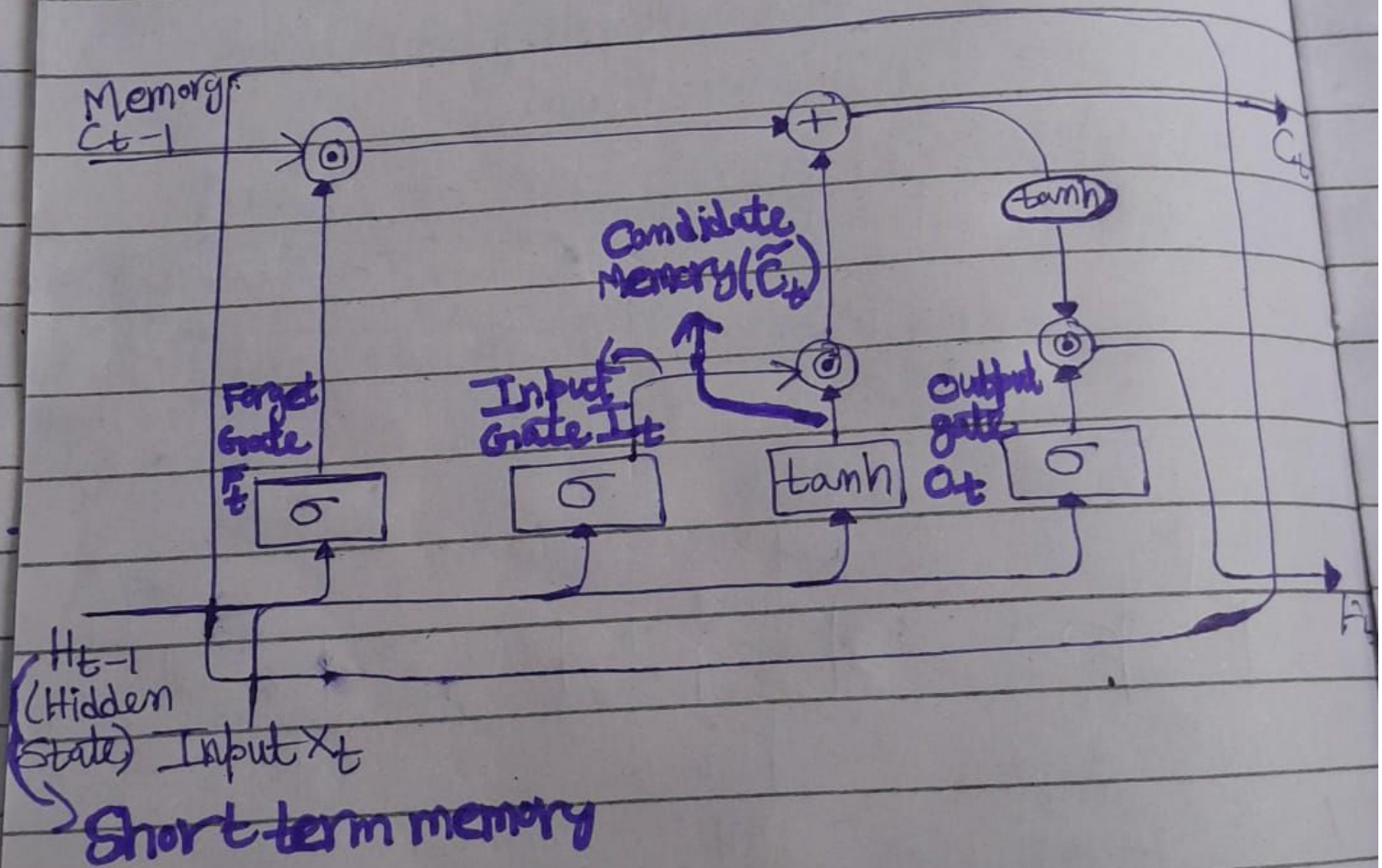


In LSTM RNN, we have short term memory like Simple RNN but we also have Long Term memory which remember the context till it

is required then drop it

LSTM RNN $\begin{cases} \text{Long Term Memory} \\ \text{Short Term Memory} \end{cases}$

\Rightarrow LSTM Architecture:



$\square \Rightarrow$ Neural Network Layer

$\rightarrow \Rightarrow$ Vector transfer

$\nabla \Rightarrow$ Copy

➔ ⇒ Concatenate

e.g. Combining two vectors

$$h_{t-1} = [1 \ 2 \ 3]$$

$$x_t = [4 \ 5 \ 6]$$

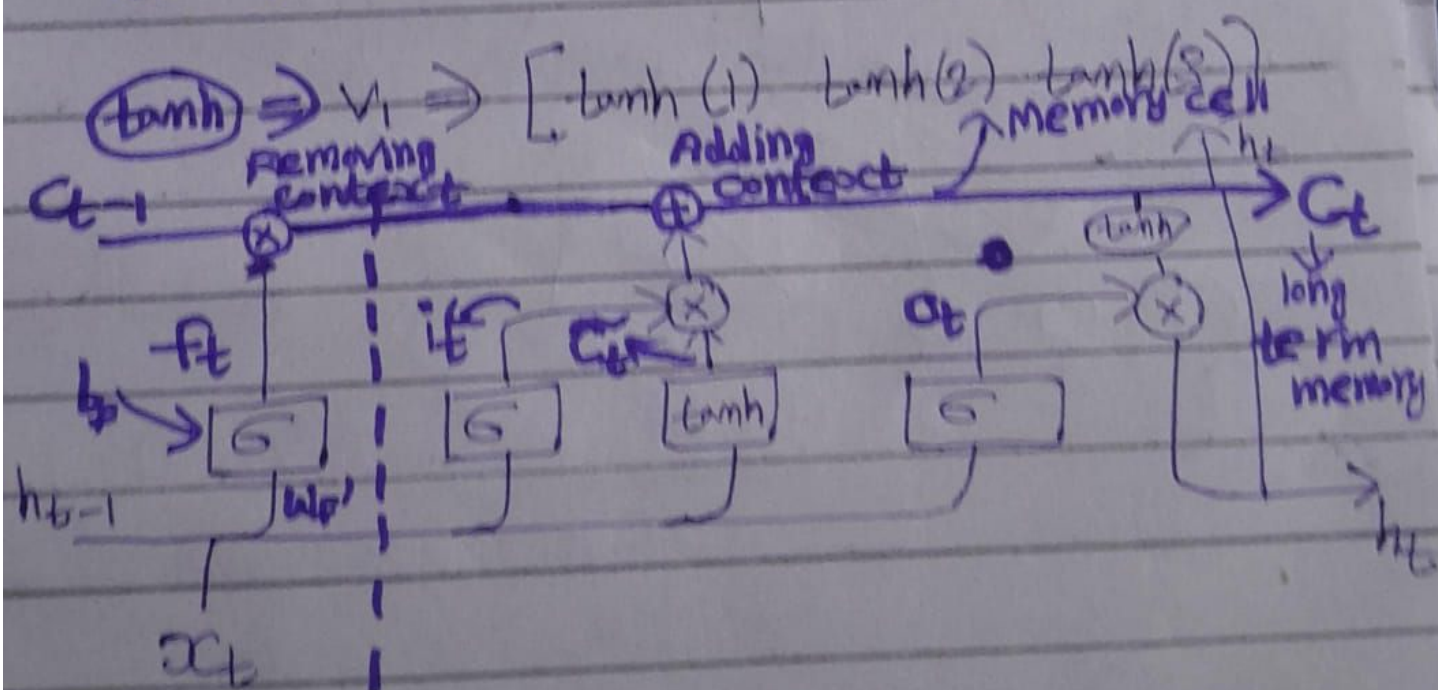
So concatenate $[1 \ 2 \ 3 \ 4 \ 5 \ 6]$

⊗ ⇒ Pointwise operation

$$\text{Let } v_1 = [1 \ 2 \ 3], v_2 = [4 \ 5 \ 6]$$

$$\otimes \Rightarrow [1 \times 4 \ 2 \times 5 \ 3 \times 6] \Rightarrow [4 \ 10 \ 18]$$

$$\oplus \Rightarrow [1+4 \ 2+5 \ 3+6] \Rightarrow [5 \ 7 \ 9]$$



Forget
gate
 f_t

So,

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

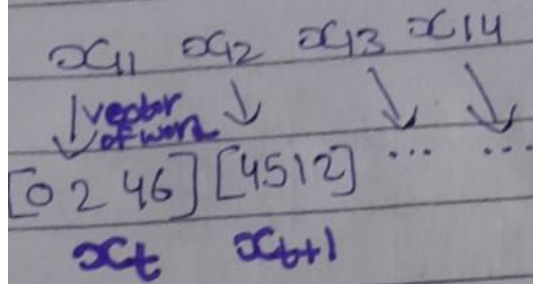
Forget Gate:

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

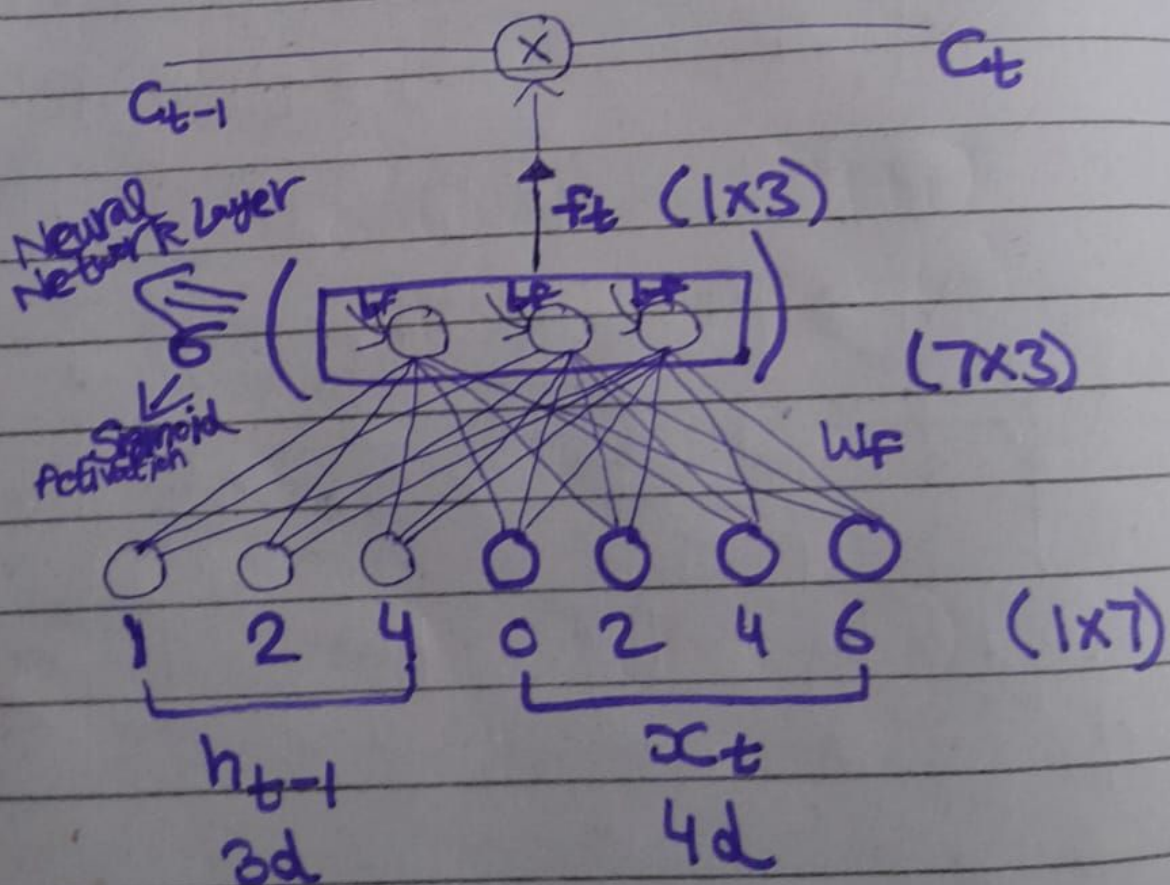
Next word to predict

Text

y_{15}



Let $h_{t-1} = [1 \ 2 \ 4]$ (3 dimensional)



Now Let

$$C_{t-1} = [6 \ 8 \ 9]$$

① Let $f_t = [0 \ 0 \ 0]$

$$C_t = C_{t-1} \otimes f_t = [6 \ 8 \ 9] \otimes [0 \ 0 \ 0]$$

$$C_t = [0 \ 0 \ 0] \Rightarrow \text{Removed all the previous context}$$

② let $f_t = [1 \ 1 \ 1]$

$$C_t = [6 \ 8 \ 9] \otimes [1 \ 1 \ 1]$$

$$C_t = [6 \ 8 \ 9] \Rightarrow \text{No context is removed}$$

③ let $f_t = [0.5, 1, 0.5]$

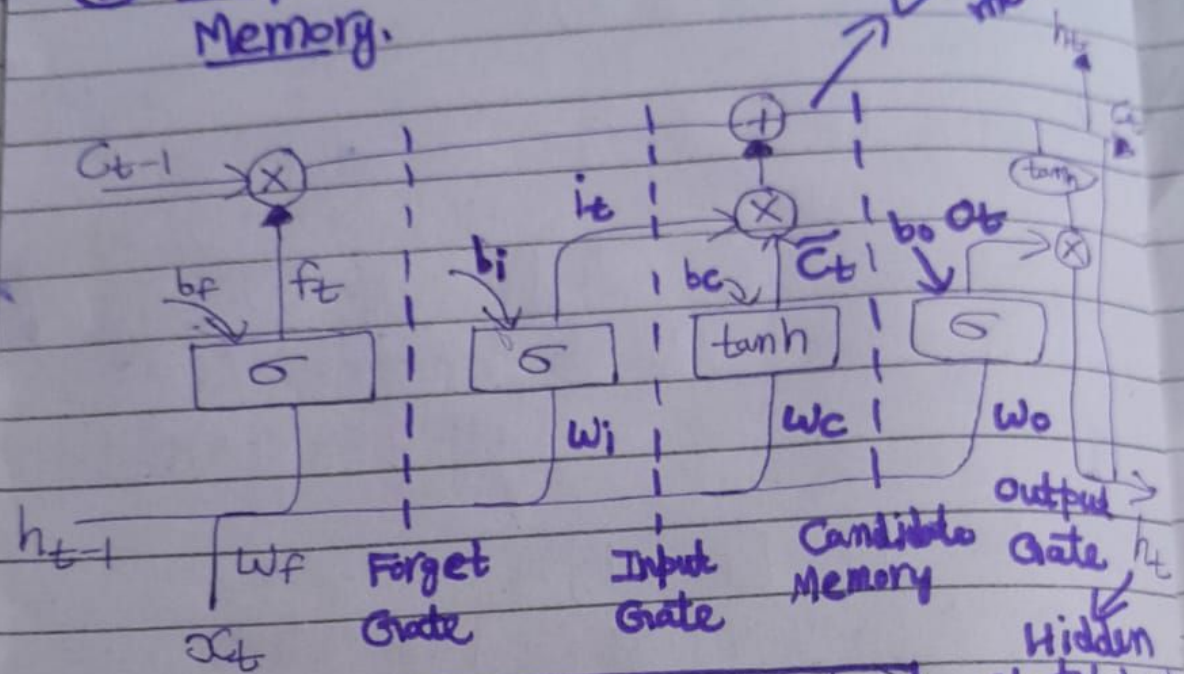
$$C_t = [6 \ 8 \ 9] \otimes [0.5, 1, 0.5]$$

$$C_t = [3 \ 8 \ 4.5] \Rightarrow \text{Some previous context is removed}$$

Conclusion:

“Based on the context Forget gate will forget some information or not forget some information.”

② Input Gate and Candidate Memory.



$$i_t = \sigma(w_i[h_{t-1} \parallel x_t] + b_i)$$

$$\tilde{C}_t = \tanh(w_c[h_{t-1} \parallel x_t] + b_c)$$

$$\text{Final Input} = i_t \otimes \tilde{C}_t$$

$$C_t = C_{t-1} \oplus (i_t \otimes \tilde{C}_t)$$

→ Final Input

Conclusion:

If any information needed to add in memory it will get added.

So, till now

$$C_t = \frac{C_{t-1} \otimes f_t}{\otimes} \oplus \frac{i_t \otimes \tilde{C}_t}{\otimes}$$

⇓
Removing
Context
(Forget gate)

⇓
Adding
Context
(Input gate
and candidate
memory)

3) Output gate:

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

$$h_t = O_t \otimes \tanh(C_t)$$

Training Data with LSTM RNN.

Test Paragraph

Output (good/bad)

I went to Restaurant

1

and order burger

The burger looked

tasty and crispy

But burger is not good
for health.

It has lot of fats,
cholesterol

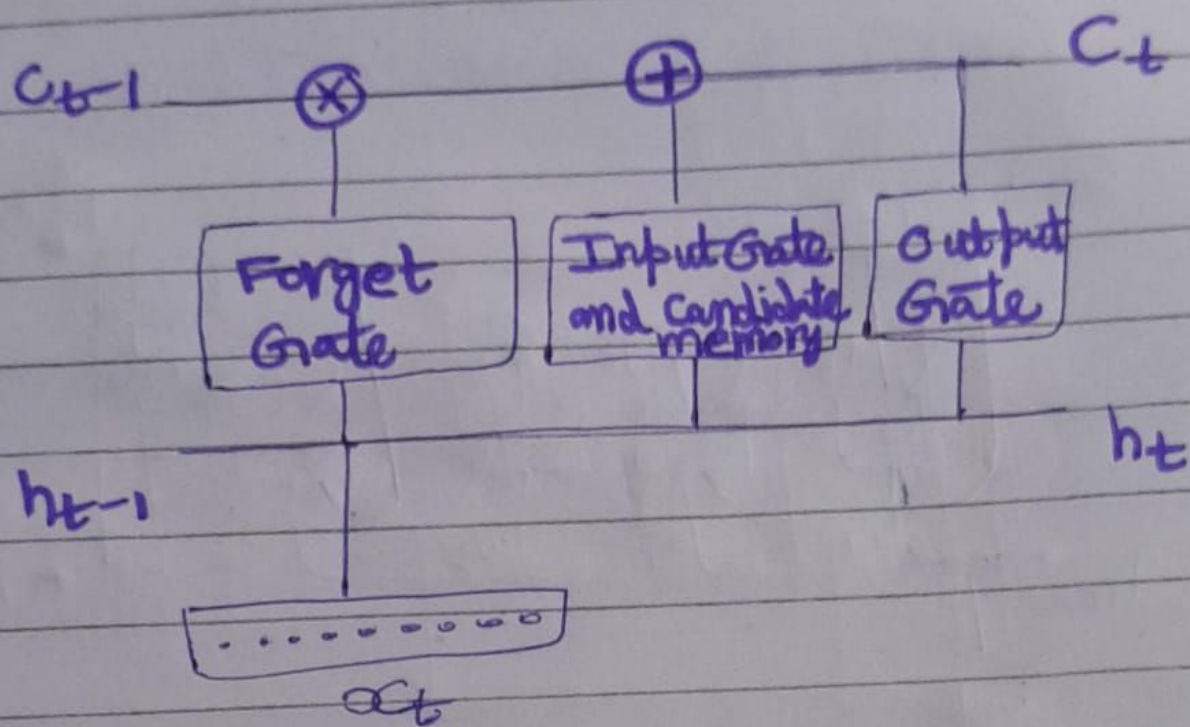
But this burger was

made with whey protein

and only vegetables were

used so it was good

Training:



Step 1: Words \rightarrow Vectors \rightarrow Embedding Layer

- Word2vec can be used

Let Word2Vec (3 dimensions)

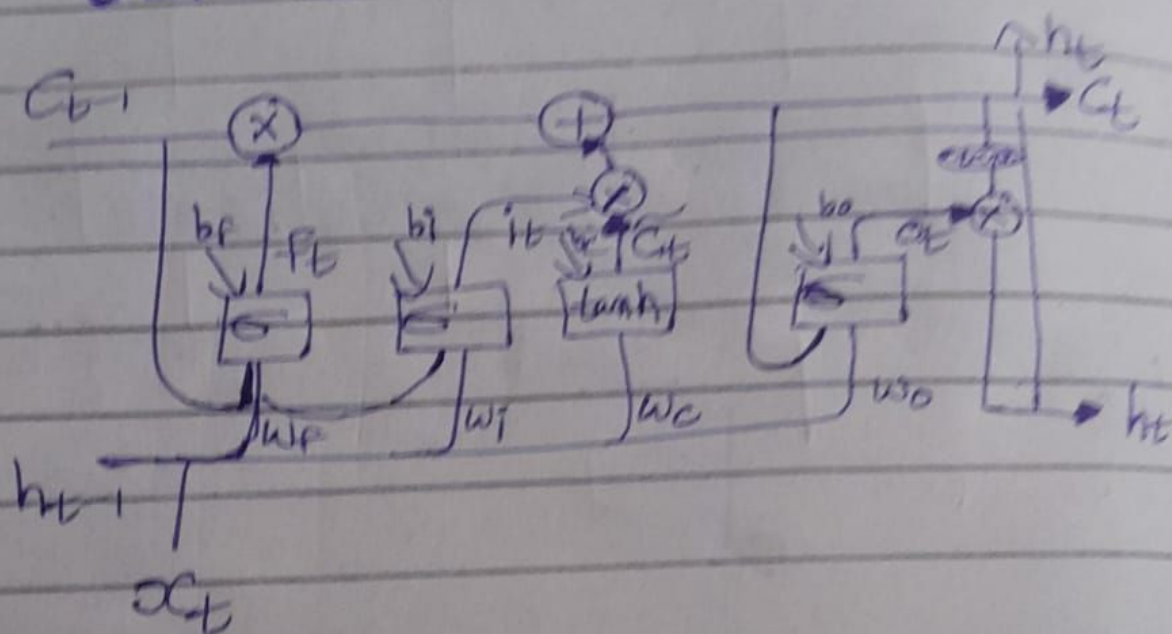
Good bad Healthy \leftarrow Black Box

Tasty $\begin{bmatrix} 0.9 & 0.0 & 0.1 \end{bmatrix} \leftarrow 3d$

- We will forget the non-important context through forget gate and add further context through Input gate and candidate memory.

→ Variants of LSTM RNN:

LSTM variants introduced by Gers & Schmidhuber (2000):



Eqn

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

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

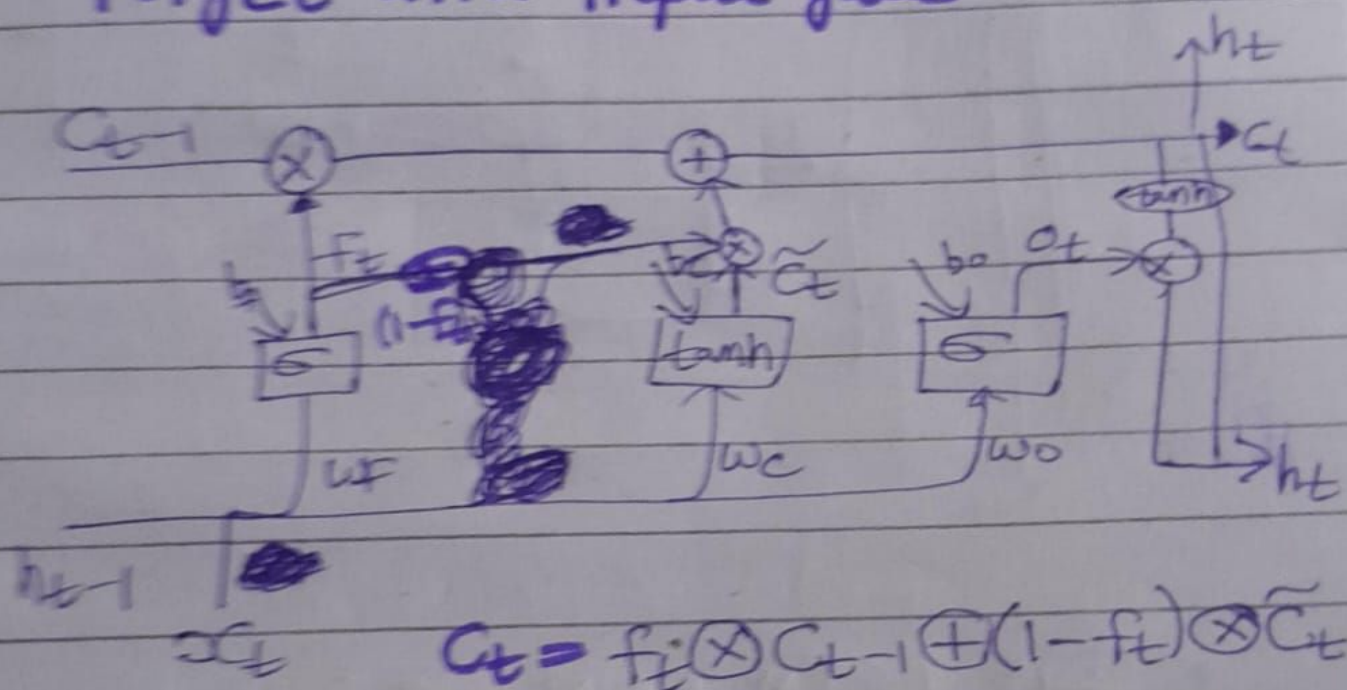
$$o_t = \sigma(W_o [C_t, h_{t-1}, x_t] + b_o)$$

Connections → From memory cell to Forget gate, input gate and output gate ⇒ Peephole Connections

Peephole Connections:

We let the gate layers look at the cell state

Another variant \rightarrow Coupling
Forget and input gates:



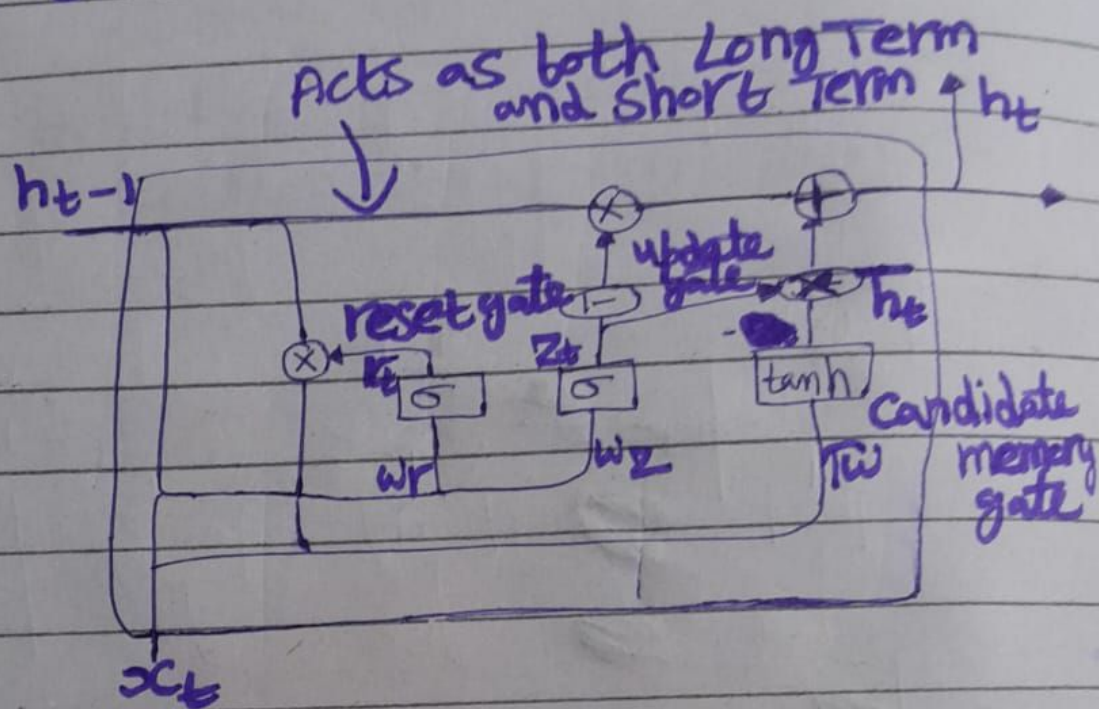
\Rightarrow Instead of deciding separately what to forget and what to add new info we make this decision together.

Conclusion:

We only forget when we are going to input something in its place or we only add/input new

values only when we forget.

⇒ GRU (Gated Recurrent unit)



- In LSTM RNN we have too many number of gates, weights and bias and complex architecture due to which training time increase because there are many trainable parameters.

Update gate

$$z_t = \sigma(W_z \cdot [h_{t-1} \parallel x_t])$$

Reset gate

$$r_t = \sigma(W_r \cdot [h_{t-1} \parallel x_t])$$

Temporary Hidden State

$$\tilde{h}_t = \tanh(W \cdot [r_t \otimes h_{t-1} \parallel x_t])$$

$$h_t = (1 - z_t) \otimes h_{t-1} \oplus z_t \otimes \tilde{h}_t$$

Reset gate (r_t):

Responsible for resetting/forgetting
some information from h_{t-1}

Update Gate (z_t):

"It defines what context
info needs to be added"

⇓ Depends on

Candidate Hidden State [Current context]