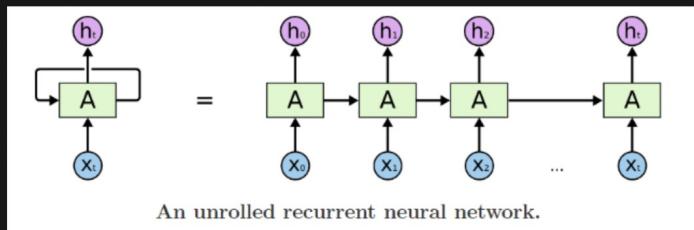


LSTM RNN [Long Short Term Memory RNN]

RNN → Long Term Dependencies → Vanishing Gradient Problem

- ① RNN → Problem? ✓
- ② Why LSTM RNN? ✓ Basic Representation.
- ③ How LSTM RNN works
 
- ④ LSTM Architecture
- ⑤ Working of LSTM RNN

Problems With RNN → Long Term Dependency

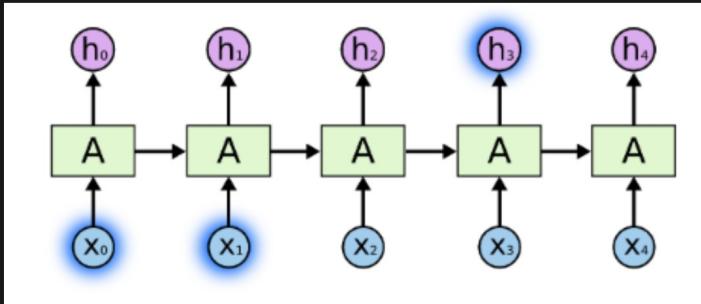


Vanishing Gradient Problem

Task:

Next Word In a Sentence

The color of the Sky is blue
— further context



Gap is 1cm

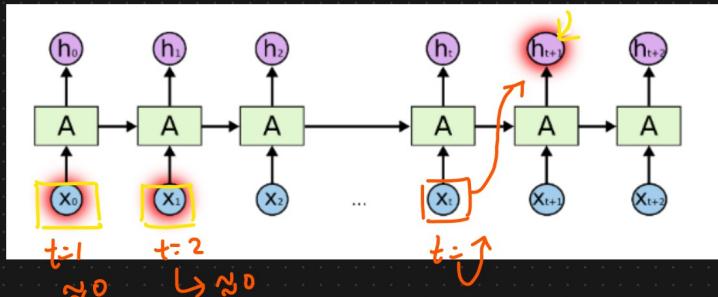
Huge gap O/P ← Context

I grew up in India ... I speak
fluent English ← Context

Name of language



further context



10-0.2r

0-1

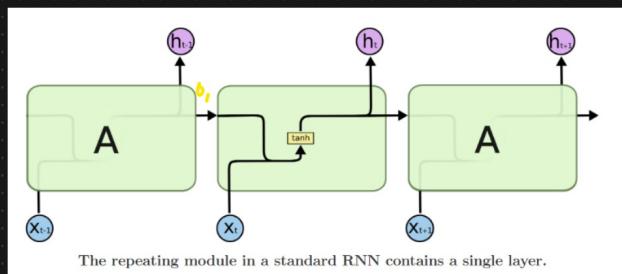
Huge Gap → Long Term Dependency

RNN → Long Term Dependency → Vanishing Gradient Problem

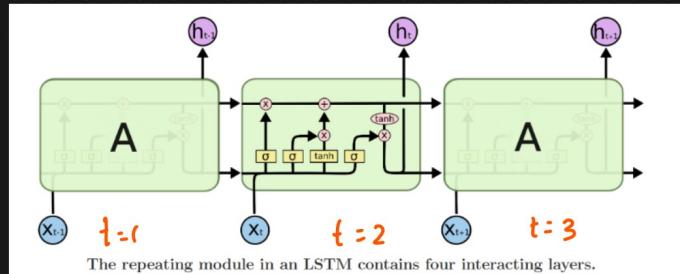
Chain Rule → ≈ 0 .

Basic Representation of RNN And LSTM RNN

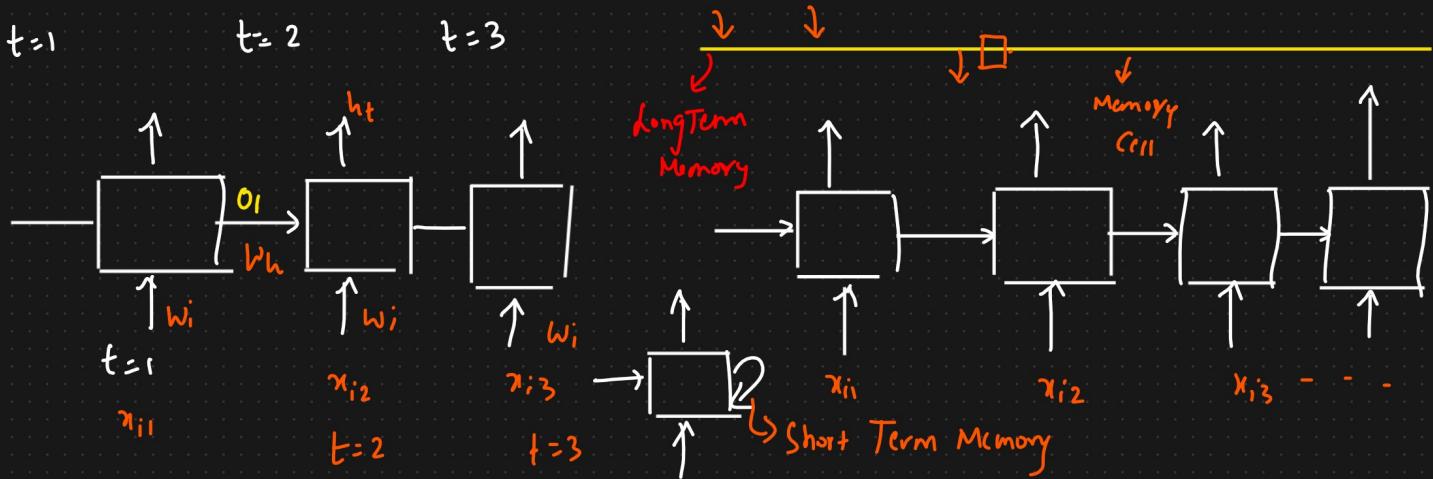
LSTM RNN



The repeating module in a standard RNN contains a single layer.



The repeating module in an LSTM contains four interacting layers.

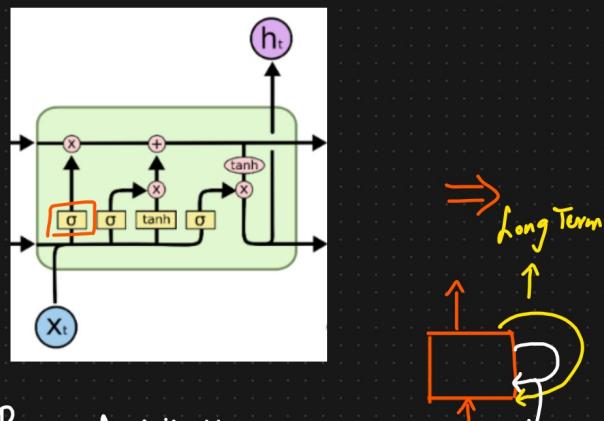


LSTM RNN → Long Term Memory
 → Short Term Memory

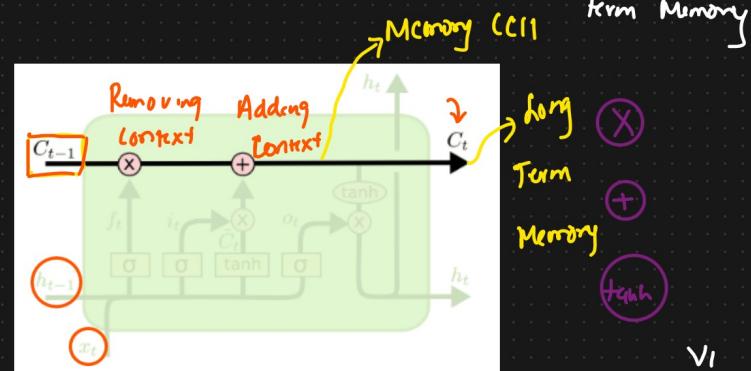
Convoyana But : fuggages



LSTM Architecture



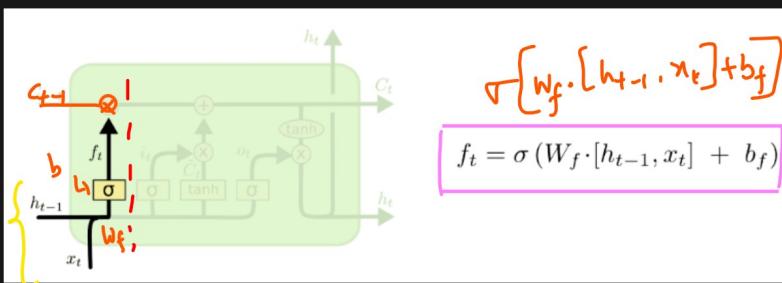
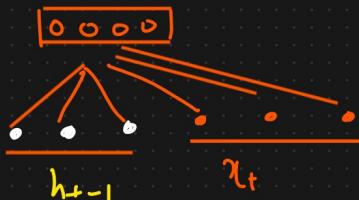
Basic Architecture



Combining 2 vectors

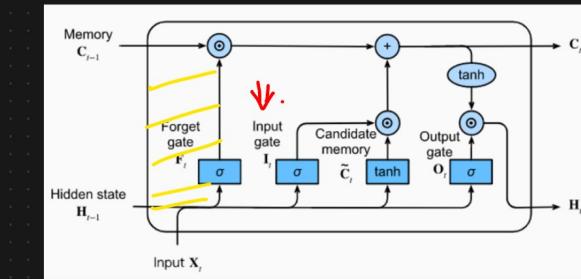
$$h_{t-1} = [1 \ 2 \ 3]$$

$$x_t = [2 \ 3 \ 4]$$



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

$$\sqrt{w_f \cdot [h_{t-1}, x_t] + b_f}$$



Forget Gate

Text Next Word
 $x_1 \ x_2 \ x_3 \ x_4 \quad y_5$



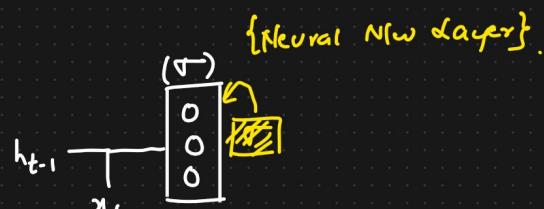
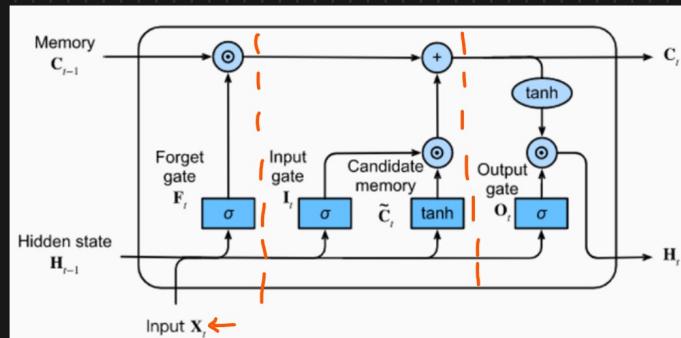
h_{t-1} = Hidden state of previous time stamp
 x_t = Word passed as i/p in the current time stamp

$$x_t \quad [0 \ 2 \ 4 \ 1] \quad h_{t-1} \quad [1 \ 2 \ 4] \quad 3d$$

$$x_{t+1} \quad [4 \ 5 \ 1 \ 2] \quad - \ - \ -$$

$$c_{t-1} \quad [4 \ 2 \ 1] \leftarrow 3d$$

LSTM RNN

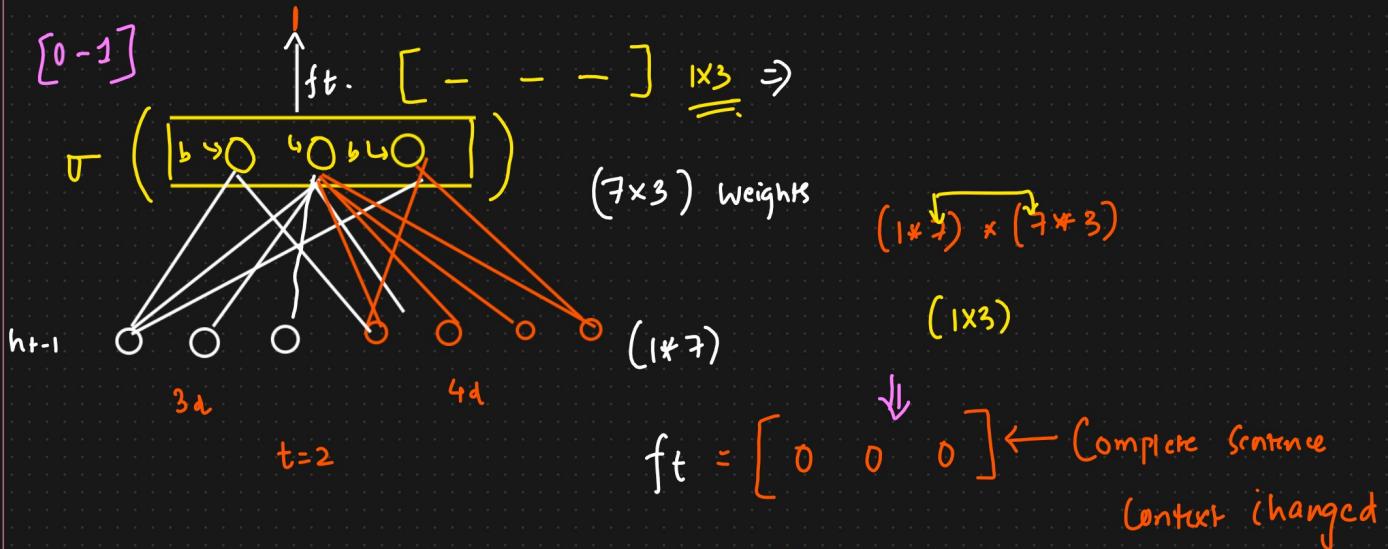


$$v_1 = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \xrightarrow{\text{tanh}}$$

$$\times = [4 \ 10 \ 18]$$

$$+ = [5 \ 7 \ 9]$$

$$\tanh \Rightarrow [\tanh(1) \ \tanh(2) \ \tanh(3)]$$



$$\textcircled{1} C_{t-1} = \begin{bmatrix} 6 & 8 & 9 \end{bmatrix} \otimes \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \leftarrow \text{Removing all the previous context}$$

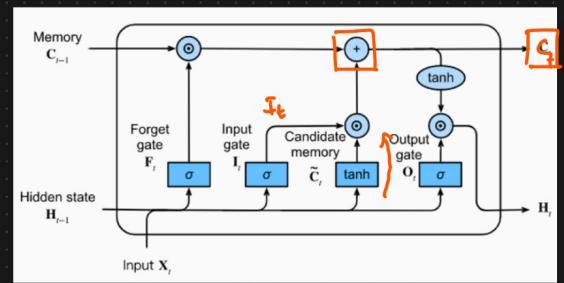
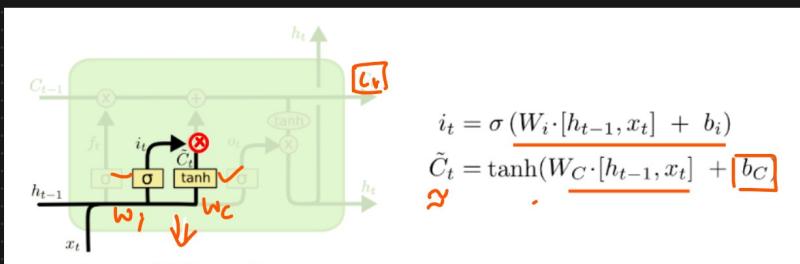
$$ft = [1 \ 1 \ 1]$$

$$\textcircled{2} C_{t-1} = \begin{bmatrix} 6 & 8 & 9 \end{bmatrix} \otimes \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 6 & 8 & 9 \end{bmatrix}$$

$$\textcircled{3} C_{t-1} = \begin{bmatrix} \underline{6} & \underline{8} & \underline{9} \end{bmatrix} \otimes \begin{bmatrix} 0.5 & 1 & 0.5 \end{bmatrix} = \begin{bmatrix} 3 & 8 & 4.5 \end{bmatrix}$$

Conclusion : Based on the context \rightarrow Forget gate will let go some information or will not let go some info {Forgetting}.

② Input Gate And Candidate Memory



Adding Info

$$I_t = \begin{bmatrix} 2 & 4 & 1 \end{bmatrix} \xrightarrow{\sigma} \begin{bmatrix} 0.8 & 0 \end{bmatrix} \xrightarrow{\oplus} \begin{bmatrix} 0 & 2 & 0 \end{bmatrix} \Rightarrow \text{Input Gate}$$

$b \rightarrow \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$

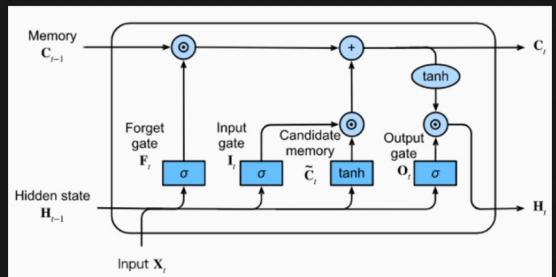
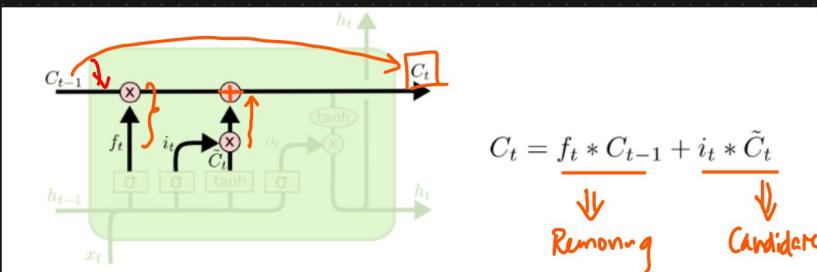
$h_{t-1} \leftarrow \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \xrightarrow{W_i} \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$

$x_t \xrightarrow{W_c} \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$

$b \rightarrow \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \xrightarrow{(1 \times 3)} \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \xrightarrow{(1 \times 3)} \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \xrightarrow{(1 \times 3)} \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$

Context = If any information needed to be added it the memory

$C_{t-1} \rightarrow$ The information will be added



I stay in India - - - - -
 and I speak English Hindi

or
 Forgetting
 Some info

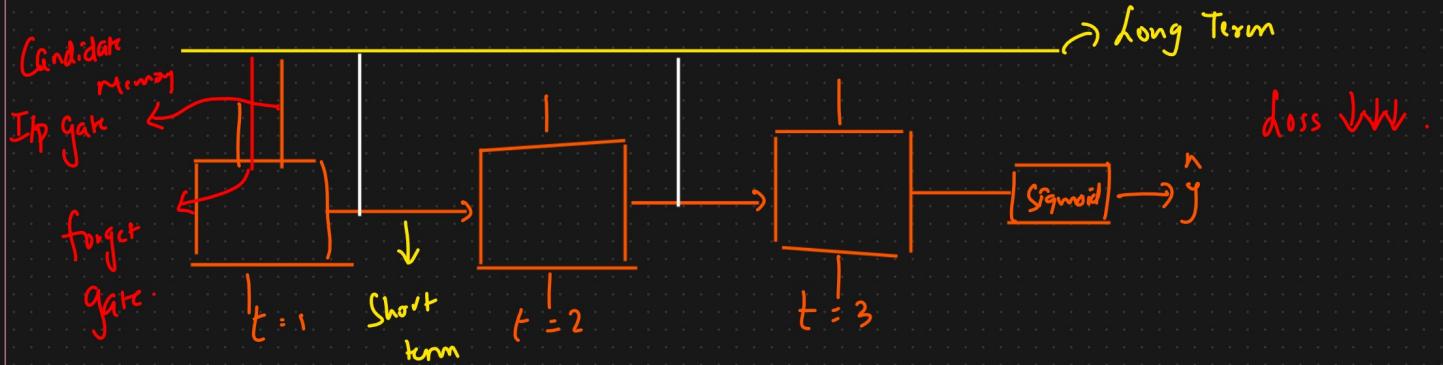
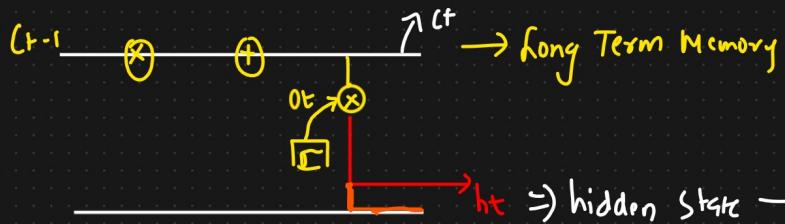
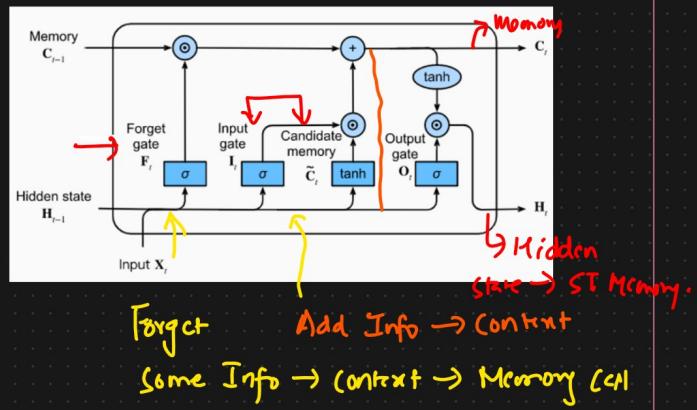
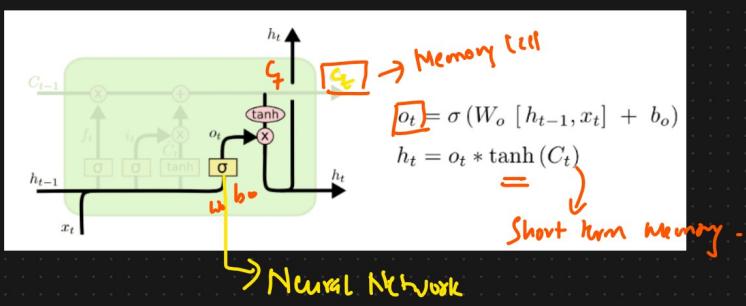
Memory -

Forget Gate

I/P gate \otimes Candidate memory

+
 $C_{t-1} \Rightarrow C_t$

Output gate LSTM RNN



$[W_i, W_c, W_o]$ → Updating ← Back Propagation

GRU RNN ⇒ LSTM Variant

Training Data With LSTM RNN

{Training} Text Paragraph

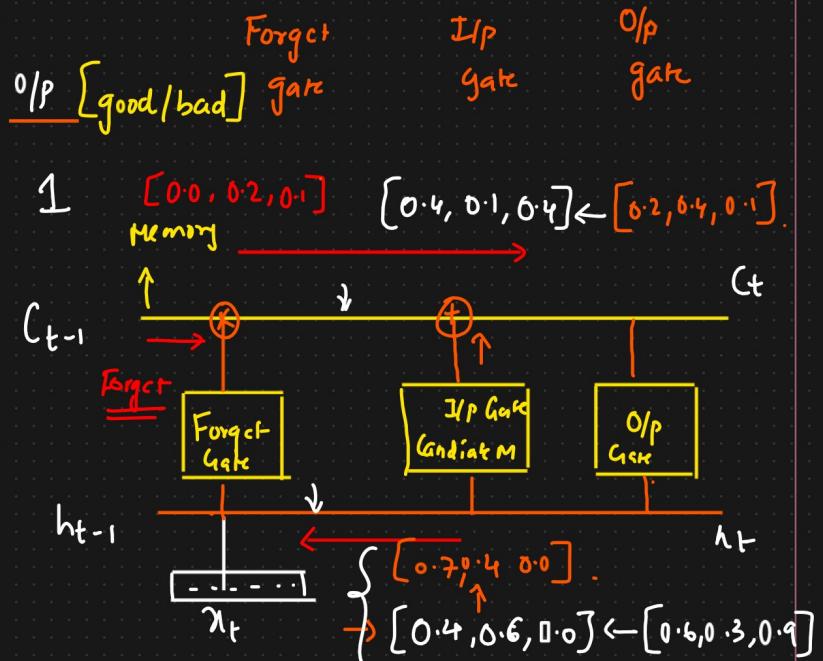
I Went to Restaurant and order burger

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



Word → Vectors → Embedding Layer

Word2Vec [3 dimension - vector]

→ $\begin{bmatrix} \text{Good} \\ \text{Bad} \\ \text{Healthy} \end{bmatrix}$ ← Black Box

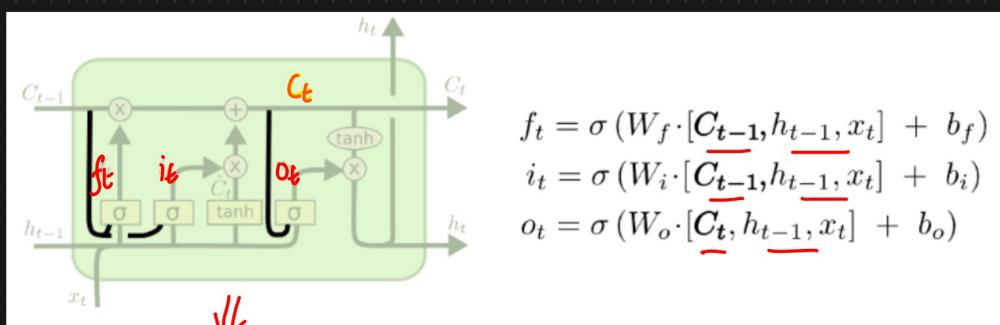
Tasty $[0.9 \quad 0.0 \quad 0.1]$ ← 3 d.

Variants of LSTM RNN

LSTM Variants Introduced By Gers & Schmidhuber [2000]

LSTM RNN [1970-80]

↳ Research paper



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

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

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

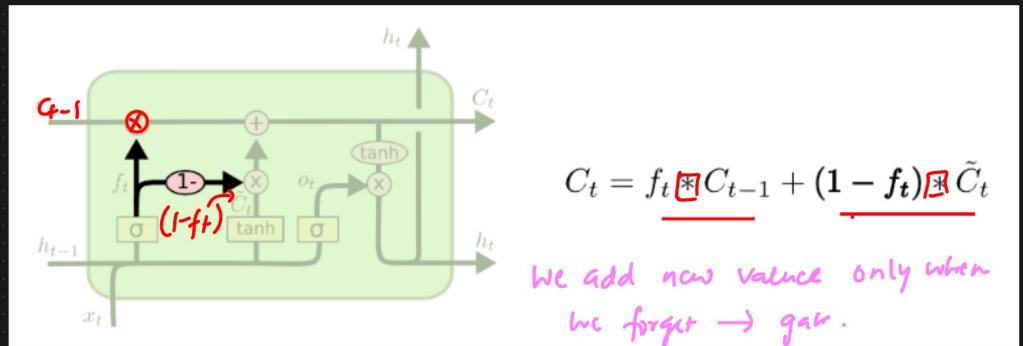
↓
focus
=====

Connections → From memory cell to

forget gate
i/p gate ⇒ Peephole
connections
o/p gate

Peephole Connections: We let the gate layers look at the cell state

Another variation \rightarrow Coupling Forget And I/p Gates



Goal: We only forget

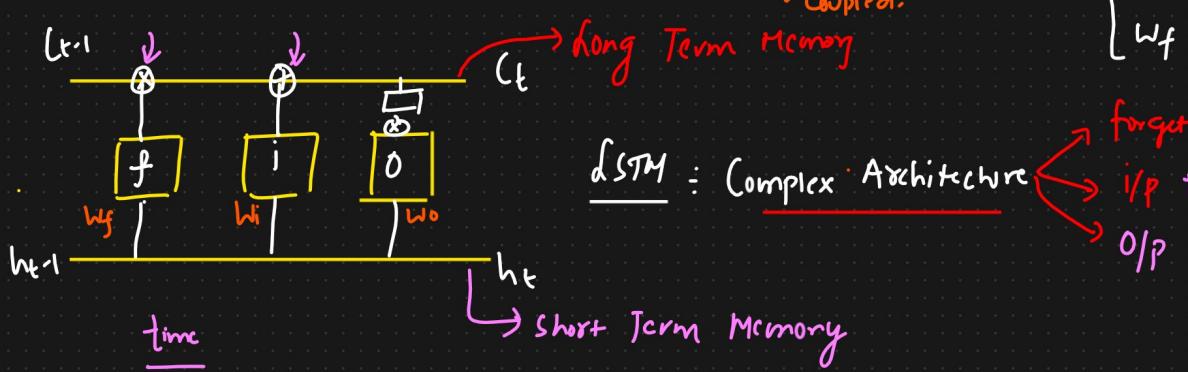
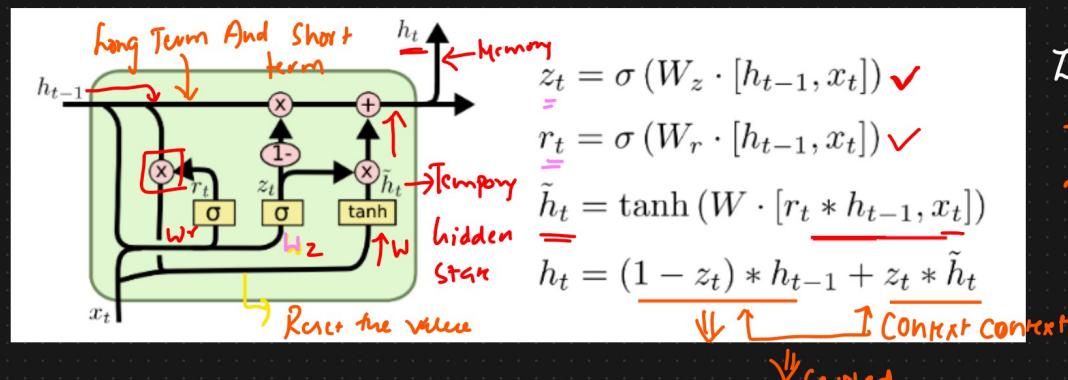
when we're going to i/p something in its place.

We only i/p new values to the state when we forget something older.

Instead of separately deciding what to forget and what we should add new Info, we make this decision together.

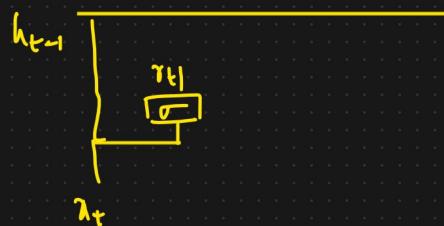
Gated Recurrent Unit $\left[\text{Cho, et al [2014]} \right]$

1980 \rightarrow LSTM
2000 - variants
2014 \rightarrow GRU



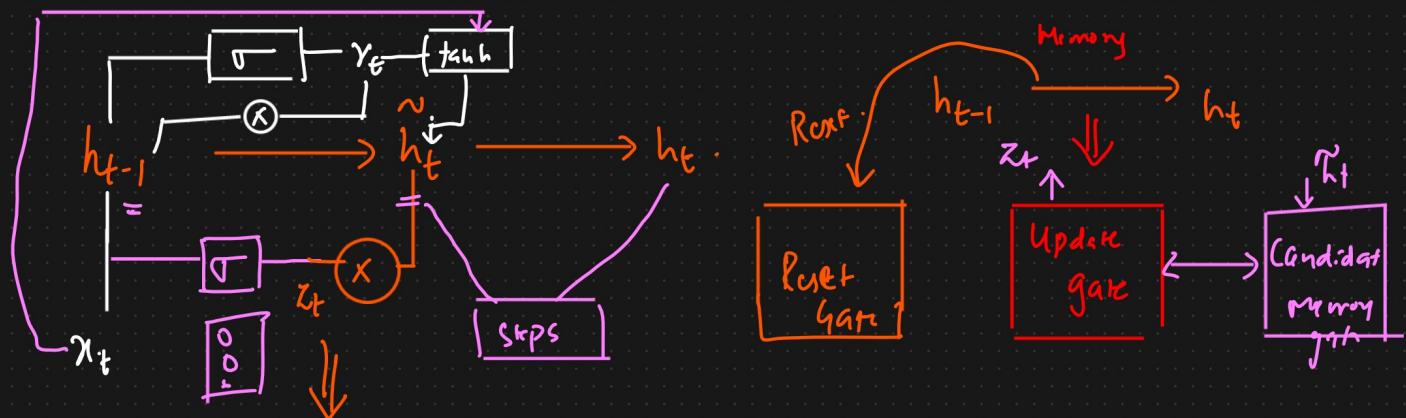
Training Time $\uparrow \uparrow$

$$\text{Reset gate } = r_t =$$



\rightarrow Resetting some info from $h_{t-1} \Rightarrow$ Memory \rightarrow LTM + STM

$$\begin{aligned}
 h_{t-1} &= [0.6 \quad 0.5 \quad 0.3 \quad 0.9] \\
 r_t &\leftarrow [0.2 \quad 0.4 \quad 0.8 \quad 0.2] \\
 \downarrow & \quad \downarrow \quad \downarrow \quad \downarrow \\
 x_t &\rightarrow [0.12 \quad 0.20 \quad 0.24 \quad 0.18] \leftarrow \text{Rescaling} \rightarrow \text{Context}
 \end{aligned}$$



What Context Info needs to be Added



Candidate hidden state · [Current Context]

↓
Imp → Add Info

$\tilde{h}_t \Rightarrow$ New Info

—