

LSTM Part 1: The What?

16 October 2024 19:10

LSTM : The Why?

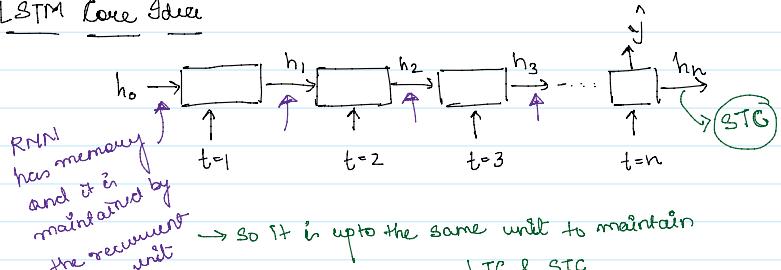


Long Short Term Memory



They were made because RNNs couldn't handle long term dependency

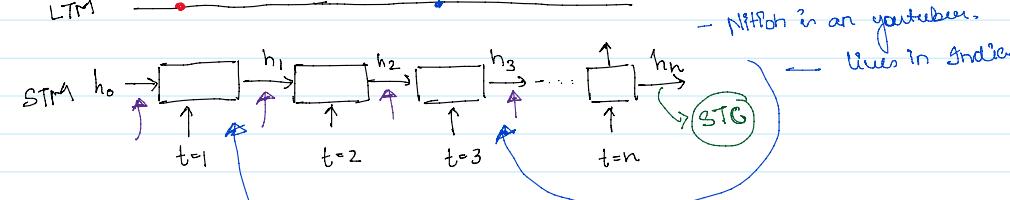
LSTM Core Idea



mathematically we saw in BPTT in RNN due to vanishing/exploding gradient only STC dominates

This is when scientists thought that why not maintain a connection to get LTC

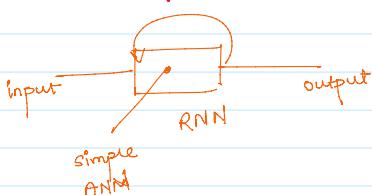
We can use this info in next sentence
Ankitा girl → Remove previous LTM & add new Nithi in youtubeer



Job current input
ke lie lagta hai ki wo
Long term more important
hair, hair use long term
memory me deal denge

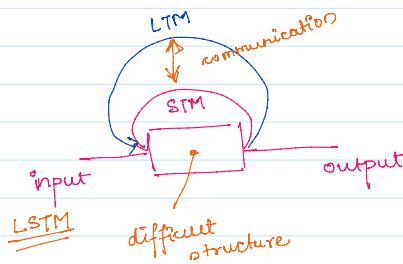
Example
- Ankitा in a good girl - ?
- Nithi in an youtubeer
predict

LTM remain until it is removed
and it is propagated till end

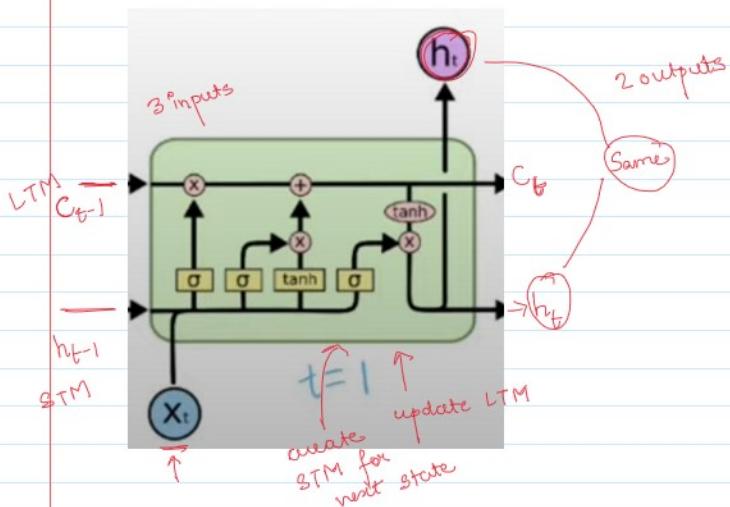
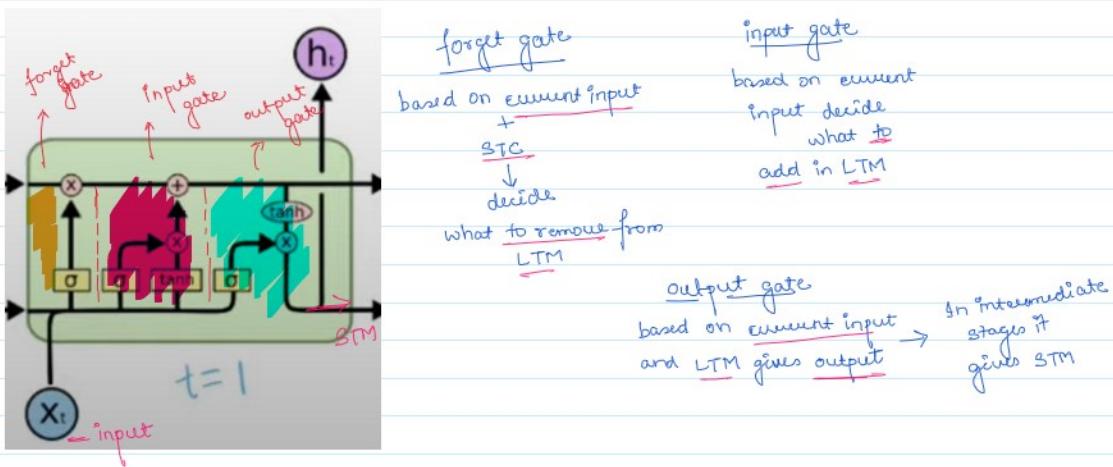
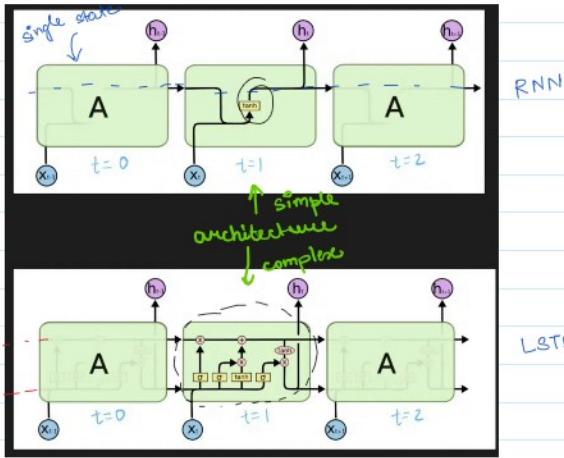


2 differences

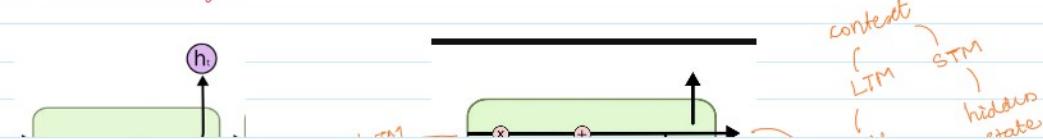
- 2 states (LTM + STM)
- architectural differences due to communication

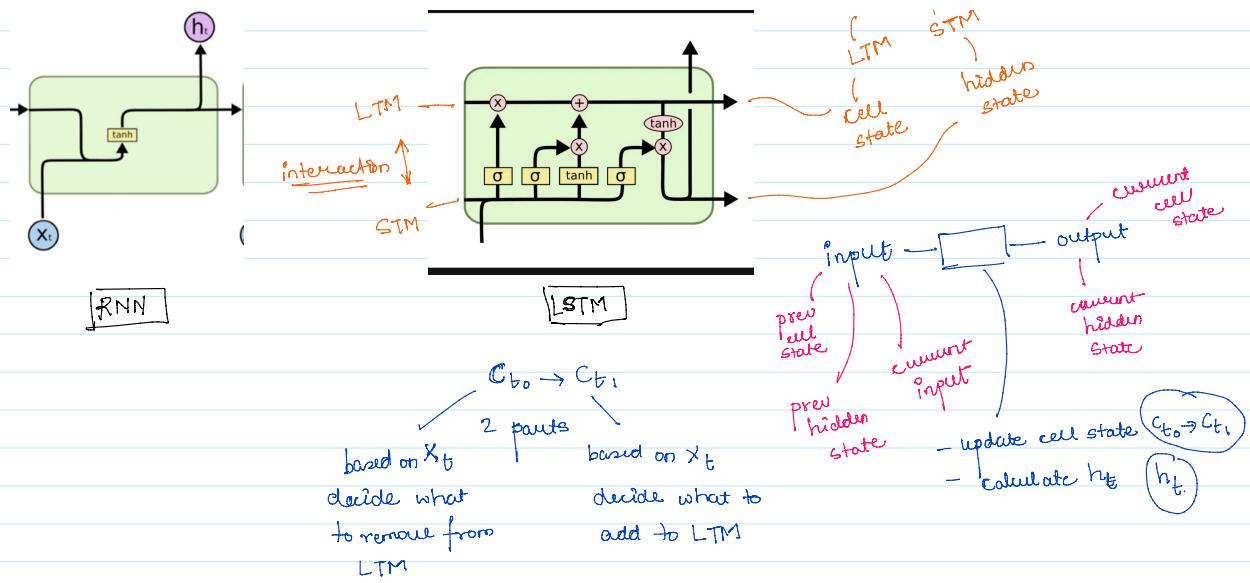


different structure
(because STM & LTM communicate to decide what's important)

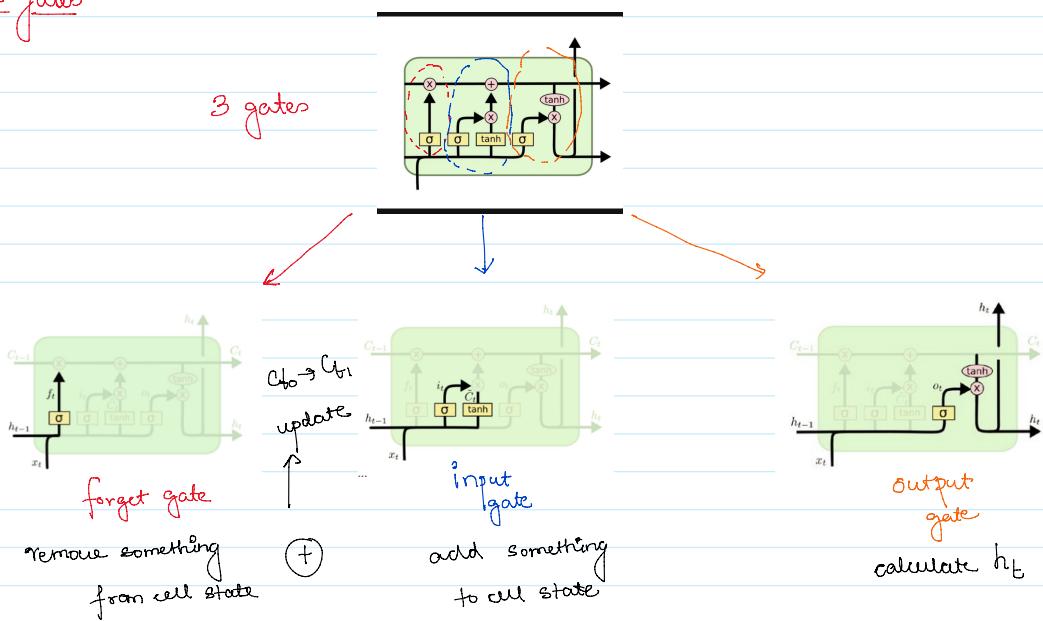


Architecture of LSTMs

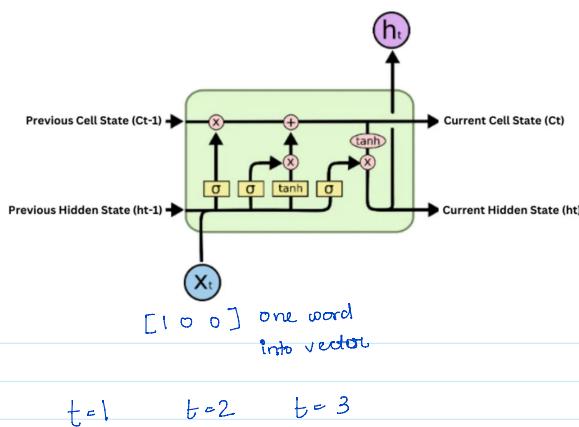




The Gates



What are h_t and C_t ?

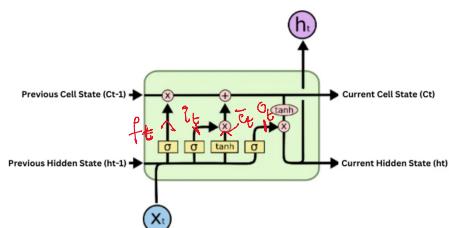


- They are vectors
- They have same dimension

$$\begin{array}{l|l} \text{cat mat rat} & 0 \\ \text{cat rat rad} & 0 \\ \text{mat mat cat} & 1 \end{array}$$

vectorize one

$$\begin{array}{l|l} [1\ 0\ 0] [0\ 1\ 0] [0\ 0\ 1] & 0 \\ [1\ 0\ 0] [0\ 0\ 1] [0\ 0\ 1] & 0 \\ [0\ 1\ 0] [0\ 1\ 0] [1\ 0\ 0] & 1 \end{array}$$

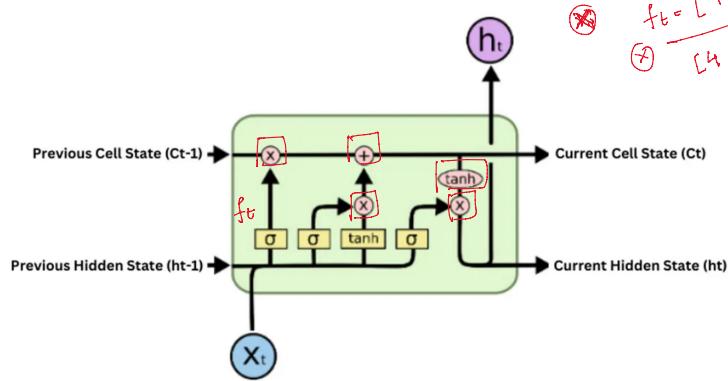


What are $f_t, i_t, o_t, \tilde{c}_t$?

f_t forget gate
 i_t input gate
 \tilde{c}_t candidate cell state
 o_t output gate

They have same dimension as C_t, h_t

Pointwise Operations



$$C_{t-1} = \begin{bmatrix} 4 & 5 & 6 \\ 1 & 1 & 1 \\ 2 & 2 & 2 \end{bmatrix}$$

$$f_t = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

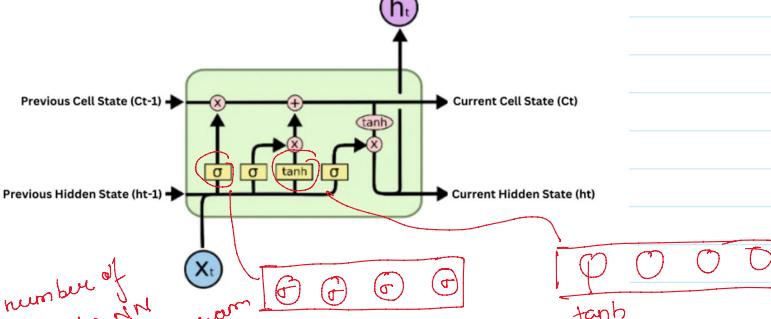
$$\oplus \quad \begin{bmatrix} 4 & 10 & 18 \end{bmatrix}$$

$$C_t = \begin{bmatrix} 4 & 5 & 6 \end{bmatrix}$$

$$f_t = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

$$\begin{aligned}
 C_t &= \begin{bmatrix} 4 & 5 & 6 \end{bmatrix} \\
 f_t &= \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \\
 &\quad \begin{bmatrix} 5 & 7 & 9 \end{bmatrix} \\
 &\quad \text{(tanh)} \quad C_t = \begin{bmatrix} 4 & 5 & 6 \end{bmatrix} \quad \text{(tanh}(6)\text{)} \\
 &\quad \text{(tanh}(4)\text{)} \quad \text{(tanh}(5)\text{)} \\
 &\quad \text{---} \quad \text{---} \quad \text{---} \quad \text{---} \quad \text{---} \quad \text{---}
 \end{aligned}$$

Neural Network Layers



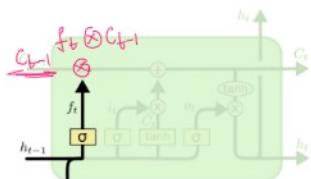
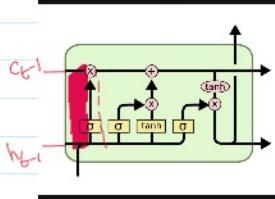
- The number of nodes in the NN layer are hyperparam
- The nodes will NN layer
- The nodes will be same for all the 4 NNs

The Forget Gate



input
output
from

The Forget Gate



input
 h_{t-1}, C_{t-1}
 x_t

output
remove from
all states

- 1) calculate f_t
- 2) $C_{t-1} \otimes f_t$

$$[f_1 \ f_2 \ f_3] = f_t \quad (1 \times 3)$$

$$W_{3 \times 7} = W_f \quad 21 \text{ weights} + 3 \text{ biases}$$

$$(3 \times 7)(7 \times 1) = (3 \times 1) + (3 \times 1) \\ = (3 \times 1)^T = f_t$$

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

$$\begin{matrix} 3 \times 7 & \otimes & 7 \times 1 \\ & & 3 \times 1 \end{matrix}$$

$$+ (3 \times 1) \rightarrow (3 \times 1) f_t$$

$$C_{t-1} (3 \times 1) \otimes f_t (3 \times 1)$$

$$\downarrow (3 \times 1) \quad \begin{matrix} 4 \cdot 5 \cdot 6 \\ \dots \end{matrix} \quad C_{t-1} \dots \otimes \rightarrow [2 \ 2 \cdot 5 \ 3]$$

f_t basically this means we are halving the values in LTM and forgetting half LTM

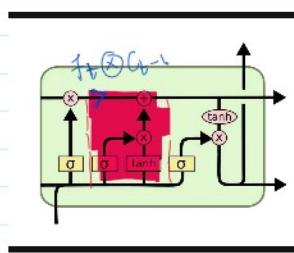
$$\text{If } f_t = [0 \ 0 \ 1] \rightarrow \text{forget all LTM}$$

$f_t = [1 \ 1 \ 1] \rightarrow \text{remember all forget nothing}$

$\therefore f_t$ decides how much information to let flow

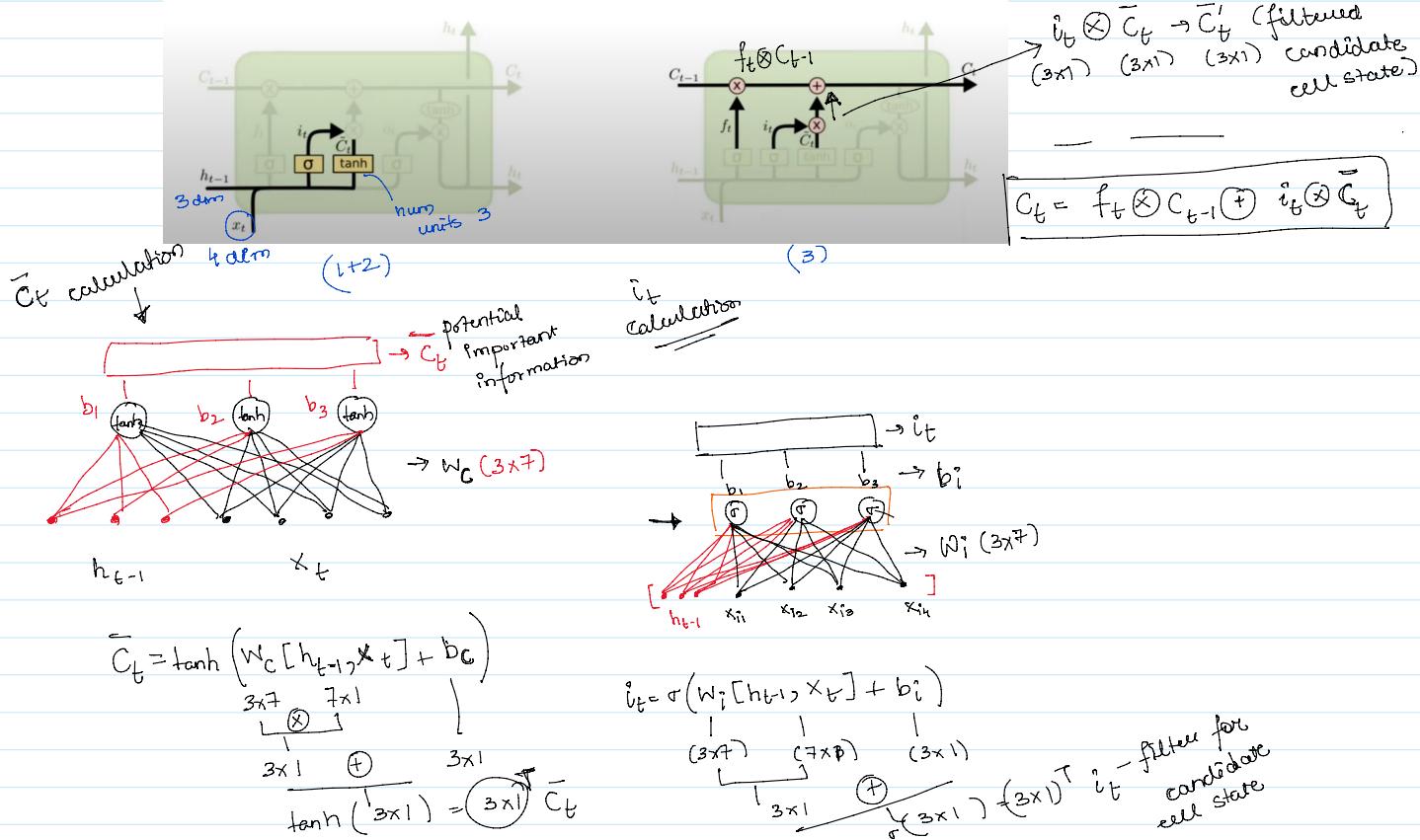
Input Gate

Add some new important information to C_t



Stages

- 1) C_t candidate cell state
- 2) i_t decides what values from C_t to add in C_t
- 3) C_t final cell state



Output gate

