

RNN

LSTM → Long / short term Memory

is a special type of RNN, designed to handle long term dependencies and solve vanishing gradient problem. It handles both long term & short term dependencies.

Jmp

- Standard RNNs forget earlier context in long sequences.
- LSTM solves this by introducing a **cell state** — a “conveyor belt” that runs through the sequence, carrying information with minimal changes.
- It uses **gates** to control what to keep, forget, and output.

Think of LSTM as a **smart memory manager** that knows exactly when to remember and when to forget.

⇒

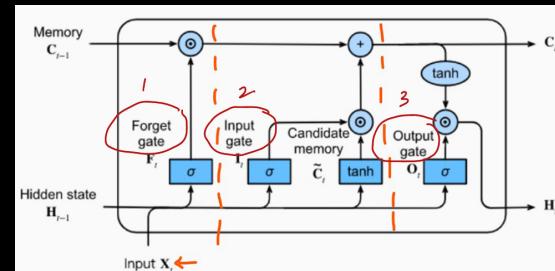
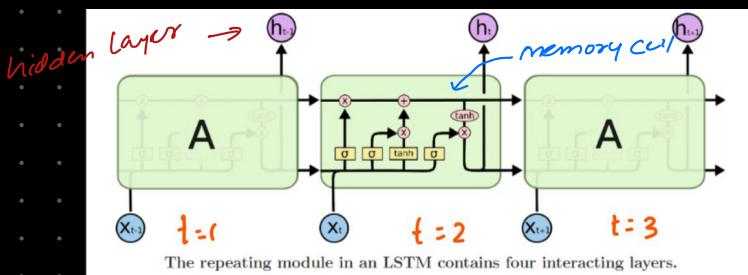
2 LSTM Structure

LSTM has three gates:

1. Forget Gate (f_t) – Decides what past information to erase.
2. Input Gate (i_t) – Decides what new information to store.
3. Output Gate (o_t) – Decides what part of the memory to output as the hidden state.

It maintains two states:

- Cell state (C_t) → Long-term memory. → *memory cell (conveyor belt)*
- Hidden state (h_t) → Short-term/contextual memory.



Jmp

LSTM computation;

Memory cell = conveyor belt

Example Sentence:

"The dog chased the cat"

At the word "cat", here's what happens:

1. Forget Gate

"What can I forget from earlier?"
→ Maybe "the dog" isn't needed now.

2. Input Gate

"Is this new word ('cat') important?"
→ Yes — likely the object of the chase.

3. Candidate Memory

"Based on 'cat' and previous context, what should I suggest adding to memory?"
→ Suggest: "Cat is the thing being chased."

4. Update Memory

Forget old parts (like 'dog'), add accepted parts of the new idea ('cat being chased').

5. Output Gate

"What should I use from this new memory to predict the next word?"
→ Maybe something like "into" (e.g., "into the yard").

Candidate memory is just a suggestion of new info based on the input and context. The input gate decides whether to use it.

Step	Role
Input	The current word ("cat")
Forget gate	Decides which old ideas to discard ("the dog" isn't needed anymore)
Candidate memory	Suggests new meaning from current input ("cat" is object of action)
Input gate	Decides how much of that suggestion to accept (yes, cat is important)
Memory update	Combines past memory with accepted new info
Output gate	Chooses what parts of memory help predict the next word (e.g., "into")

VVimp

An LSTM (Long Short-Term Memory) is a special type of RNN designed to remember information for longer and avoid the vanishing gradient problem that simple RNNs suffer from.

Why it exists

- Simple RNNs can only remember information for a few timesteps before older info fades away.
- LSTMs add a memory cell and gates that control what to keep, what to forget, and what to output.

VVimp

In simple words:

- Forget gate:** "What should I erase from memory?"
- Input gate:** "What new info should I add?"
- Cell state:** The memory highway that carries important info across many timesteps.
- Output gate:** "What part of memory do I share as output now?"

Memory Gate → maintains memory (items)



LSTM problems = Complex / time consuming

LSTM Variants

1. Classical LSTM
2. GRU LSTM

3. Bidirectional LSTM
4. LSTM with Attention
5. Convolutional LSTM

- Classic LSTM:** The original architecture with input, forget, and output gates, plus a memory cell to retain information across long sequences. 2 3 4
- Bidirectional LSTM (BiLSTM):** Processes input sequences in both forward and backward directions, helping capture both past and future context in sequence data. This is especially beneficial in tasks like language modeling and translation. 3 2
- Gated Recurrent Unit (GRU):** A streamlined variant that merges the forget and input gates into a single "update" gate, simplifying the overall architecture and computation while often performing comparably to LSTM on many tasks. 1 4 2 3
- ConvLSTM (Convolutional LSTM):** Integrates convolution operations into the LSTM architecture, making it suitable for spatial-temporal data such as video or images where both time and space information are important. 2
- LSTMs with Attention:** Enhance the LSTM model by allowing the network to focus on relevant parts of the input sequence selectively at each time step, which improves performance in complex tasks like translation or summarization. 2

Forget input = update

2 Ganks

GRU Gated Recurrent Unit

It is a simpler and faster version of LSTM

- In a standard RNN, gradients vanish or explode when sequences are long, so the network forgets earlier context.
- GRU fixes this using **gates** (like LSTM), but with **fewer parameters** and **no separate cell state**.

Think of GRU as **LSTM Lite** — keeps the essence, removes the extras.

2 GRU Structure

GRU has **two gates** instead of LSTM's three:

→ GRU →

1. Update Gate (z_t) – Decides **how much of the past information** to keep vs. update with new info.
2. Reset Gate (r_t) – Decides **how much of the past information to forget** before combining with new input.

Key difference from LSTM:

- **No separate cell state** → GRU keeps everything in a single **hidden state** (h_t).
- This means **less memory** and **faster training**.

Imp: 2 Gates, No separate hidden state

3 Step-by-Step Working

At time step t :

1. Update gate

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

→ If z_t is close to 1 → keep most of the old memory.
If z_t is close to 0 → replace it with new info.

2. Reset gate

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

→ If r_t is close to 0 → forget the past completely.
If r_t is close to 1 → use full past info.

3. Candidate hidden state

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

→ New information considering the reset gate.

4. Final hidden state

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

→ A blend of old state and candidate state, controlled by update gate.

4 LSTM vs GRU — Interview-Ready Table

Feature	LSTM	GRU
Gates	3 (input, forget, output)	2 (update, reset)
Memory storage	Separate cell state and hidden state	Single hidden state only
Parameters	More (slower training)	Fewer (faster training)
Long-term dependencies	Better at capturing	Good but sometimes less effective
Complexity	Higher	Lower

How GRU Maintains Sequence Without a Separate Memory State

1 LSTM's Two-Memory Approach

- Cell state (C_t) → Long-term memory (like a storage hard drive).
- Hidden state (h_t) → Short-term, immediate output memory (like RAM).

This separation allows LSTM to explicitly keep long-term information mostly untouched in the cell state.

2 GRU's Single-Memory Approach

- GRU merges cell state and hidden state into one hidden state (h_t).
- This single state stores both long-term and short-term dependencies in a blended form.

So, instead of having:

$$(C_t, h_t)$$

like LSTM, GRU just has:

$$h_t$$

3 How It Works

- Update gate (z_t) controls how much of the old h_{t-1} to keep.
If z_t is high → past info is preserved (similar to cell state persistence in LSTM). → to keep
- Reset gate (r_t) decides how much past info to ignore when computing the new candidate state. → to forget

Key point: The update gate's mechanism effectively lets old information flow unchanged through many time steps, which is how GRU keeps memory without explicitly having a cell state.

4 Why This Still Works for Sequences

- Even though GRU has only one vector h_t , it:
 - Retains sequence context via gating (update gate works like a "forget + keep" controller).
 - Learns to encode both long- and short-term information in different parts of the same vector.
 - Passes information through time steps selectively, avoiding the vanishing gradient problem.

5 Analogy

Think of LSTM as having two separate notebooks — one for the full story (cell state) and one for today's notes (hidden state).

GRU says: "Why not just have one notebook but write carefully so old important notes remain while adding new ones?"



📌 Practical Steps to Build an LSTM Model for Text

1. Preprocess text → clean, tokenize, convert to sequences of integers.
2. Word embeddings → use Word2Vec, GloVe, or an embedding layer.
3. Model architecture → Sequential → Embedding → LSTM layer (instead of simple RNN) → Dense (softmax).
4. Training → cross-entropy loss + Adam optimizer.
5. Evaluation → accuracy, perplexity, or F1 score depending on task.
6. Prediction → next-word generation or text classification.

⌚ Interview-friendly one-liner

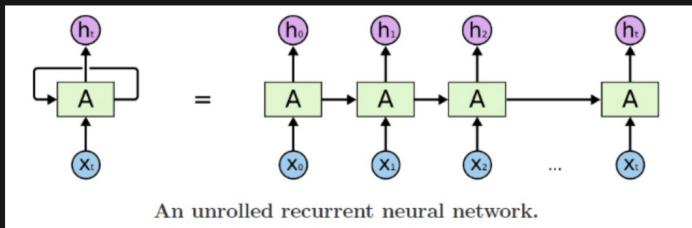
"I would preprocess the text, use embeddings, then build a Sequential model with an Embedding layer, an LSTM layer to capture long-term dependencies, and a Dense softmax output. I'd train with cross-entropy and Adam, evaluate, and finally use it for prediction."

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
 → Long Term Memory ✓ ↗
 → Short Term Memory ✓ ↘
- ④ LSTM Architecture
- ⑤ Working of LSTM RNN

Problems With RNN → Long Term Dependency



For small sentences
Simple RNN will just work fine

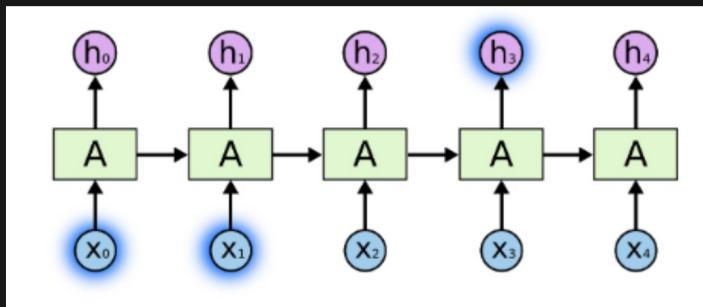
Vanishing Gradient Problem

Task:

Next Word In a

Sentence

The color of the sky is blue
— further context



Gap is 1cm

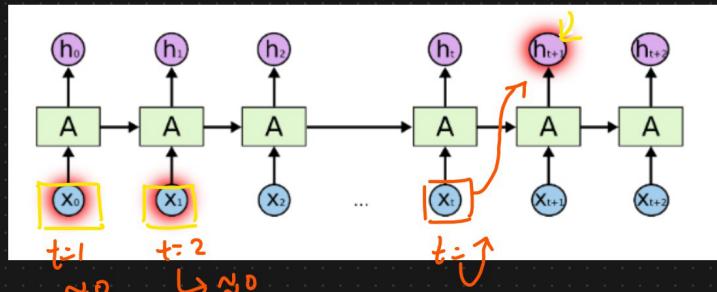
Huge gap O/P ← Context

I grew up in India... I speak fluent $\downarrow \leftarrow \uparrow$ L_{Context}

Name of language



further context



Huge Gap → Long Term Dependency

10-0.25
0-1

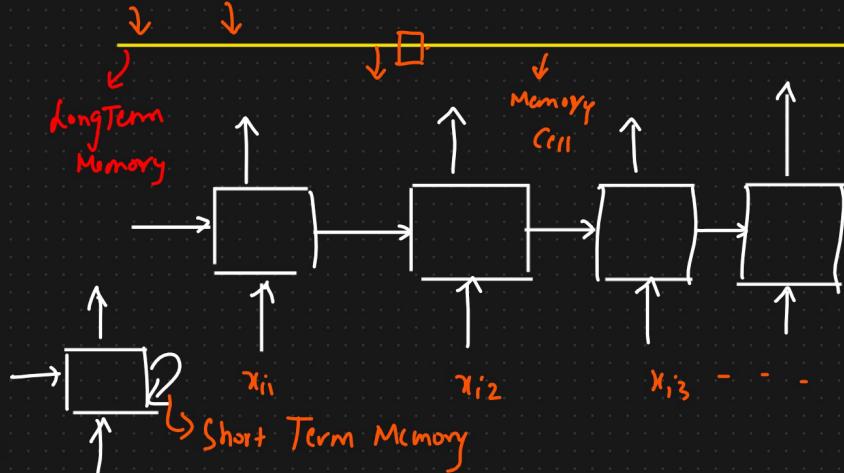
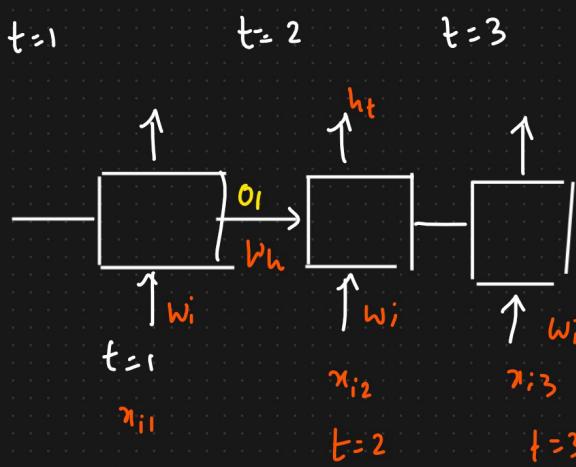
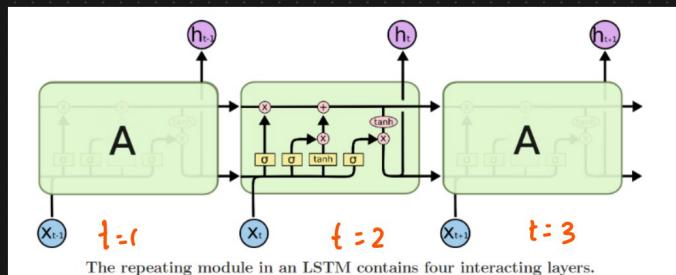
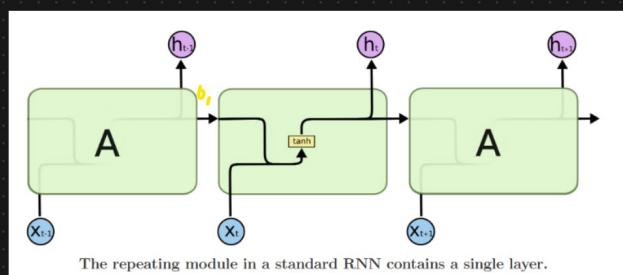
RNN → Long Term Dependency → Vanishing Gradient Problem

Chain Rule $\rightarrow \approx 0$.



Basic Representation of RNN And LSTM RNN

LSTM RNN



LSTM RNN → Long Term Memory (Memory cell)
 LSTM RNN → Short Term Memory (Hidden State)

Convoyana But : huggages

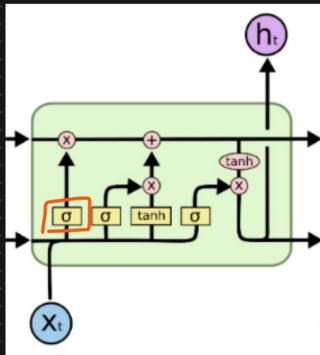


Imp. why is LSTM ; Long- Short Term ?

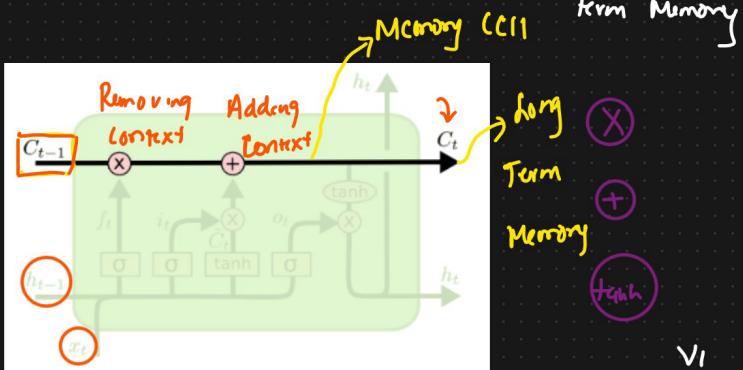
Yes, traditional RNNs function as "short-term memory" models—they maintain a single hidden state that passes through each timestep, capturing recent information. However, they struggle to retain information over long sequences due to the vanishing gradient problem, where signals fade as they propagate backward during training [Medium](#) [Wikipedia](#).

LSTMs, on the other hand, extend this capability by introducing a memory cell along with gating mechanisms (*input, forget, and output gates*) that control what to keep, forget, or pass on. These gates enable the model to preserve and manage long-term dependencies more effectively, earning the name "Long Short-Term Memory".

LSTM Architecture



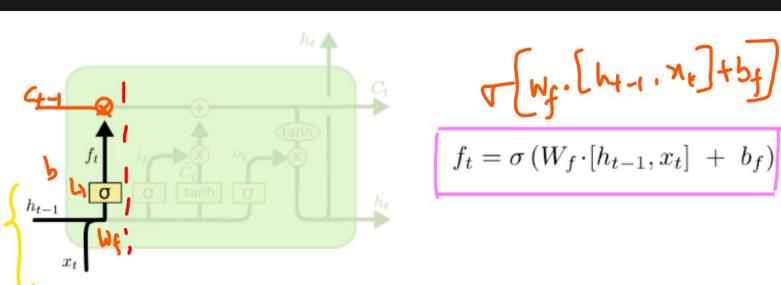
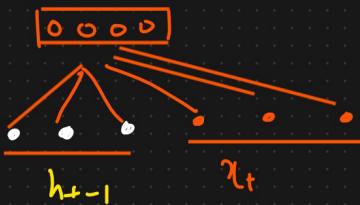
Basic Architecture



Combining 2 vectors

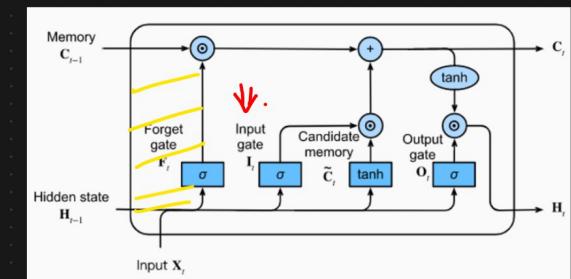
$$h_{t-1} = [1 \ 2 \ 3]$$

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



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

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



Forget Gate

Text Next Word
 x_1, x_2, x_3, x_4 y_5



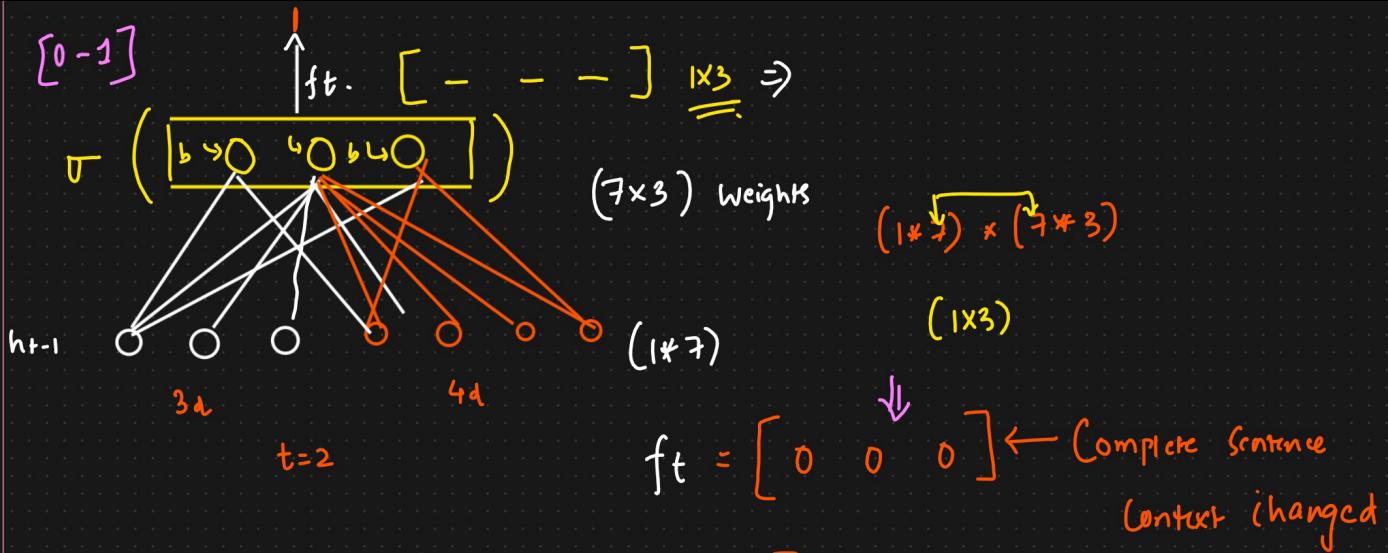
h_{t-1} = Hidden state of previous time stamp
 x_t = Word passed as i/p in the current time stamp

$$\begin{bmatrix} 0 & 2 & 4 & 1 \end{bmatrix} \quad \begin{bmatrix} 4 & 5 & 12 \end{bmatrix} \quad \dots$$

x_t x_{t+1}

$$h_{t-1} = \begin{bmatrix} 1 & 2 & 4 \end{bmatrix} \quad 3d$$

$$C_{t-1} = \begin{bmatrix} 4 & 2 & 1 \end{bmatrix} \Leftarrow 3d$$



$$\textcircled{1} \quad c_{t-1} = [6 \ 8 \ 9] \otimes [0 \ 0 \ 0]$$

$$= [0 \ 0 \ 0] \leftarrow \text{Removing all the previous context}$$

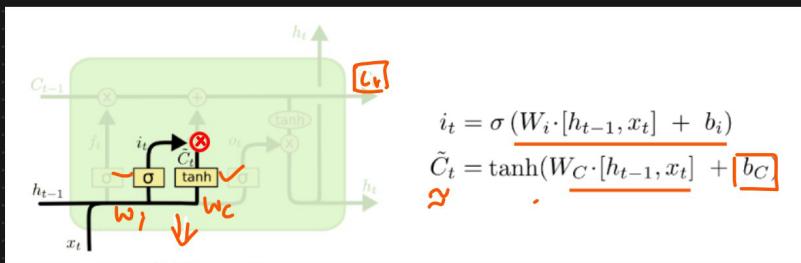
$$ft = [1 \ 1 \ 1]$$

$$\textcircled{2} \quad c_{t-1} = [6 \ 8 \ 9] \otimes [1 \ 1 \ 1] = [6 \ 8 \ 9]$$

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

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



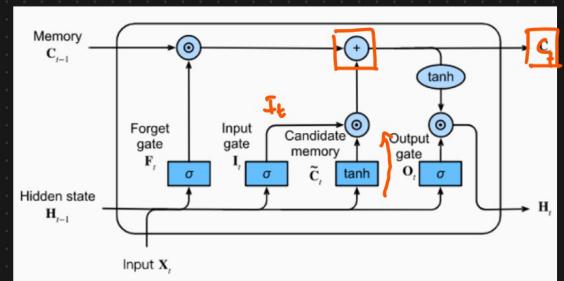
Adding Info

$$I_t = [2 \quad 4 \quad 1] \xrightarrow{\sigma} I_t = [0 \quad 8 \quad 0] \Rightarrow \text{Up Gate}$$

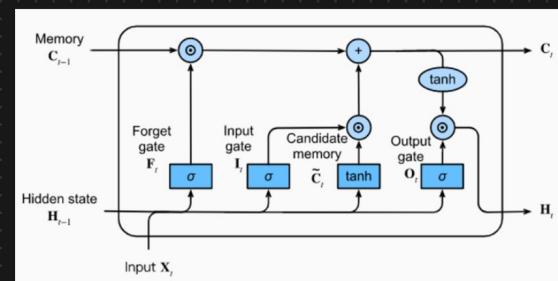
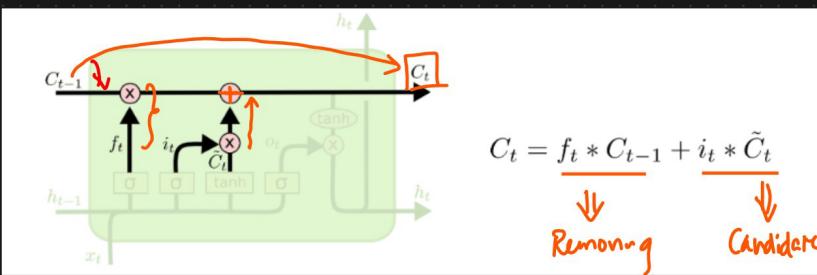
$$b \rightarrow [0 \quad 0 \quad 0]$$

$h_{t-1} \leftarrow \underbrace{0 \quad 0 \quad 0}_{w_I} \quad \underbrace{0 \quad 0 \quad 0}_{x_t}$

$b \rightarrow \begin{matrix} 0 & 0 & 0 \end{matrix} \xrightarrow{\text{tanh}} \begin{matrix} 0 & 0 & 0 \end{matrix} \xrightarrow{(1 \times 3)} \begin{matrix} 0 & 0 & 0 \end{matrix} \xrightarrow{w_C} \begin{matrix} 0 & 0 & 0 \end{matrix} \xrightarrow{(1 \times 7)}$



Context = If any information needed to be added to the memory
 $c_{t-1} \rightarrow$ The information will be added



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

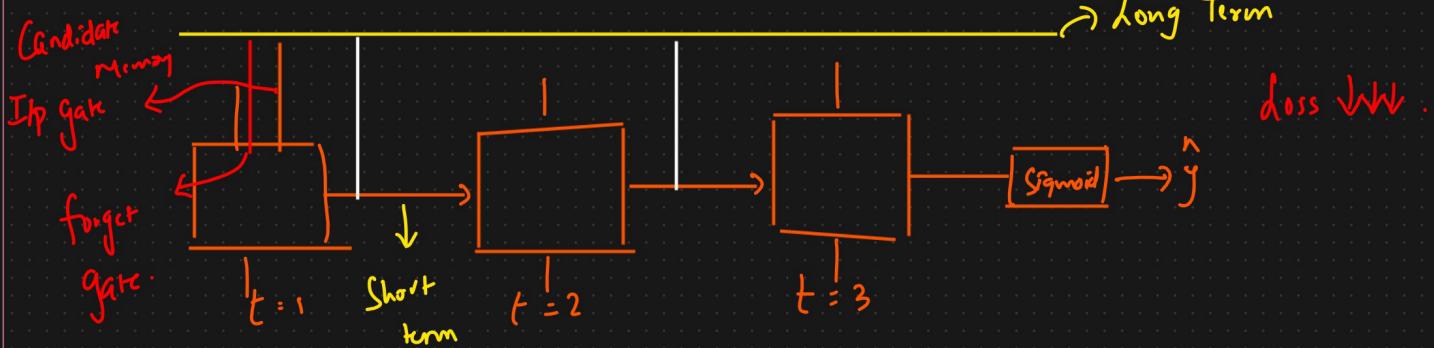
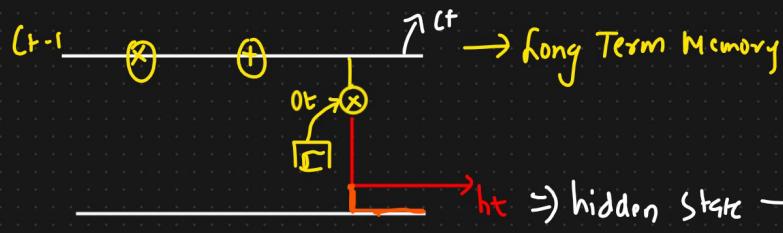
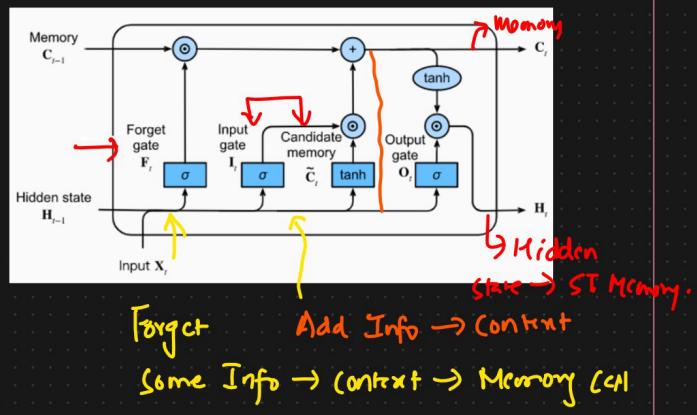
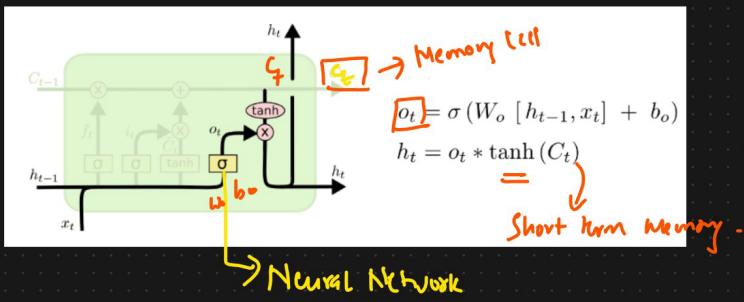
or
 Forgetting
 Some info

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
 \downarrow \downarrow \downarrow

GRU RNN ⇒ dSTM Variant

Training Data With LSTM RNN

{Training} Text Paragraph

I Went to Restaurant and order burger

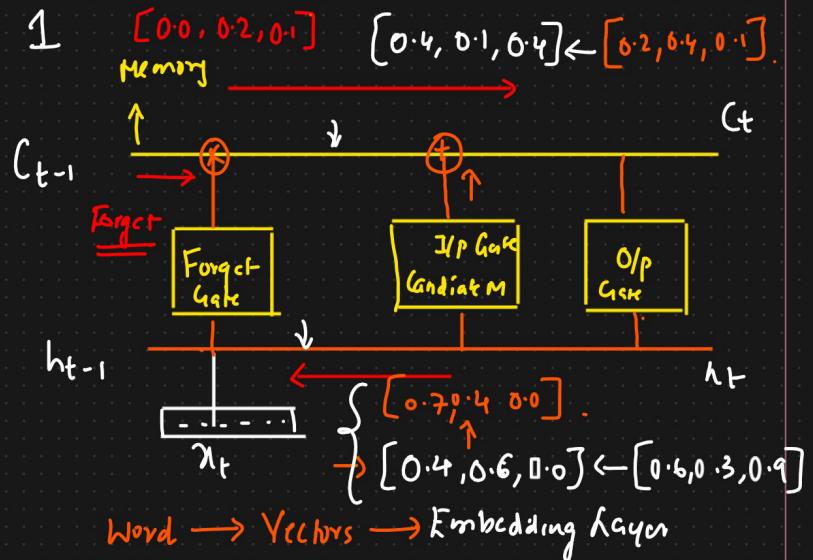
The burger looked tasty and crispy

→ 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.

O/P [good/bad] gate Forget gate I/p gate O/p gate



Word2Vec [3 dimension - vecnr]

→ [Good] [Bad] [Healthy] ← Black Box

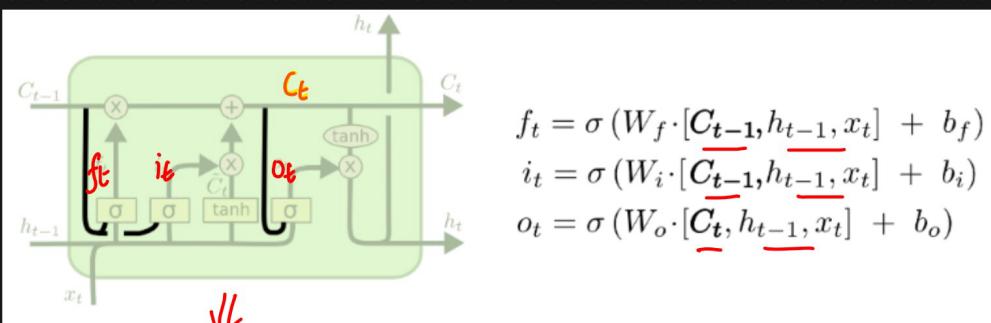
Tasty [0.9 0.0 0.1] ← 3 d.

Variants of LSTM RNN

LSTM Variants Introduced By Gers & Schmidhuber [2000]

LSTM RNN [1970-80]

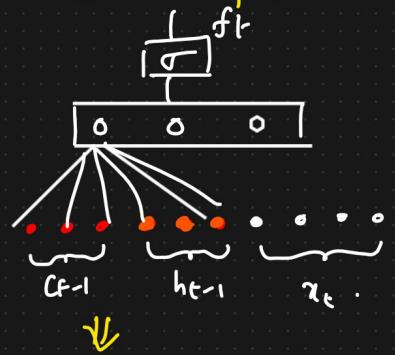
↳ Research paper



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

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

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

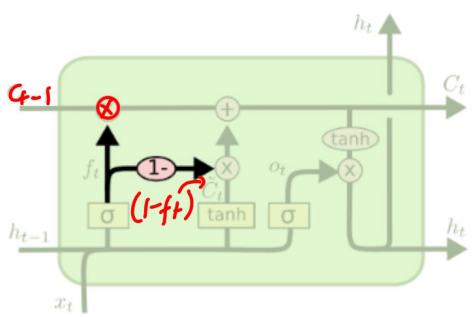


↓
focus

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



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

We add new value only when we forget \rightarrow gav.

Goal: We only forget

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



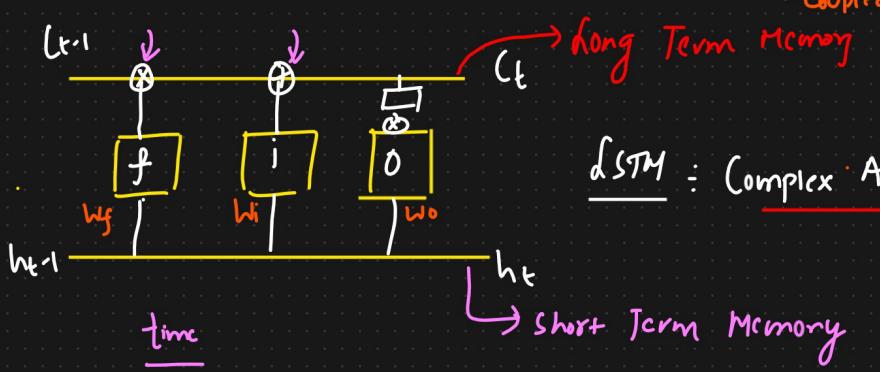
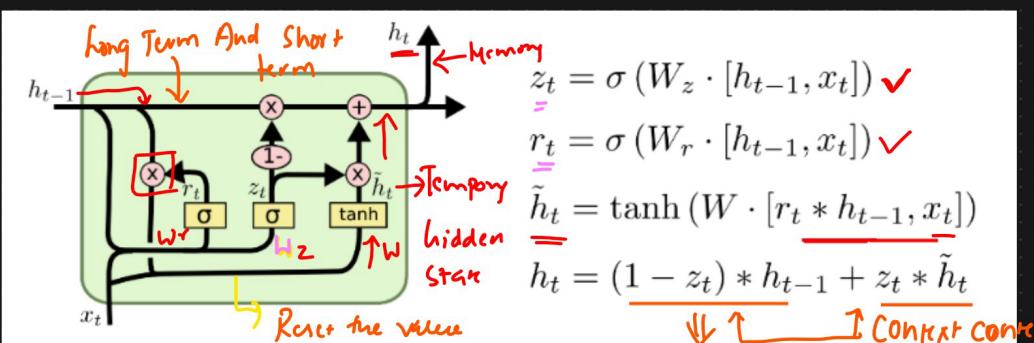
We only i/p new values to the state when we forget something older.

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

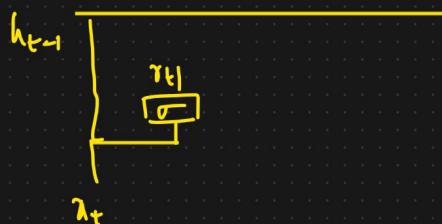
GRU \rightarrow Gated Recurrent Unit [Cho, et al [2014]]

1980 \rightarrow LSTM
2000 - variants
2014 \rightarrow GRU

$z_t \Rightarrow$ Update Gate \leftarrow
 $r_t \Rightarrow$ Reset Gate \leftarrow
 $\tilde{h}_t \Rightarrow$ Temporary hidden state.



Training Time $\uparrow\uparrow$



z_t ,

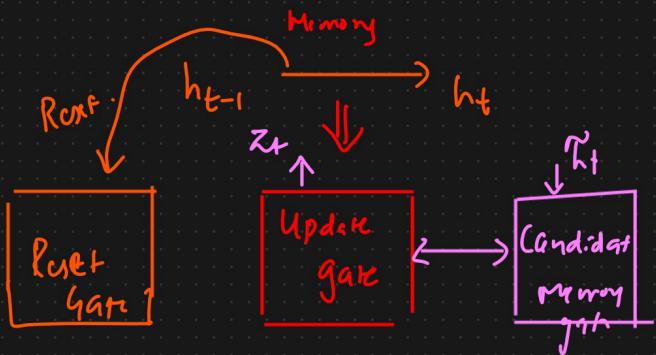
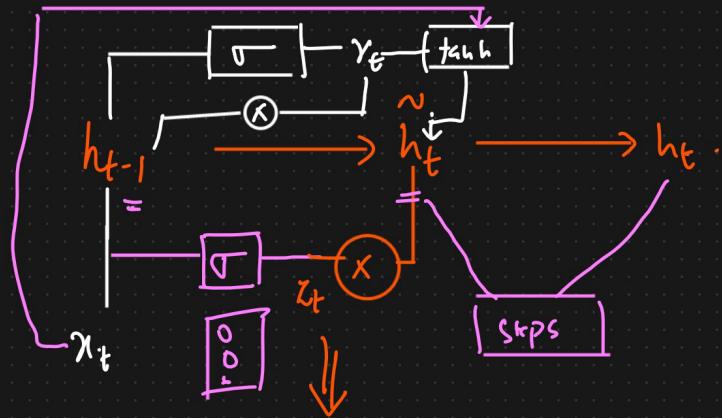
Reset gate $= r_t =$

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

LSTM: Complex Architecture
forget
i/p + candidate memory
O/P

Trainable parameters
[w_f, w_i, w_o]

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



What Context Info needs to be Added



Candidate hidden state · [Current Context]

\Downarrow
Imp \rightarrow Add Info

$L_t \Rightarrow$ fun Info

==