

Recurrent Neural Networks

A RNN Language Model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(U h^{(t)} + b_2) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$$

$h^{(0)}$ is the initial hidden state

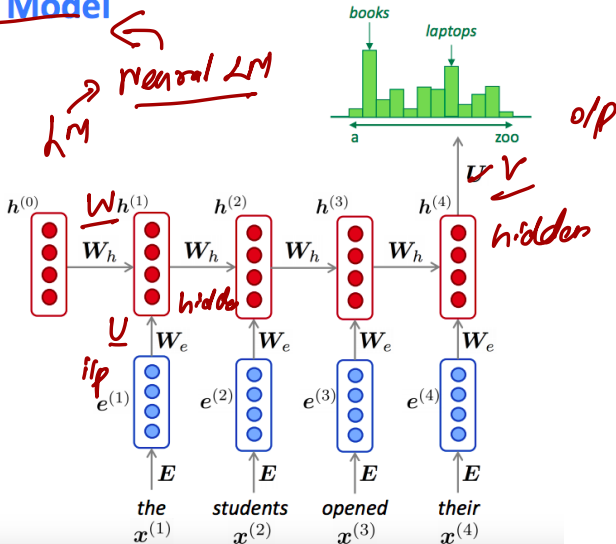
word embeddings

$$e^{(t)} = E x^{(t)}$$

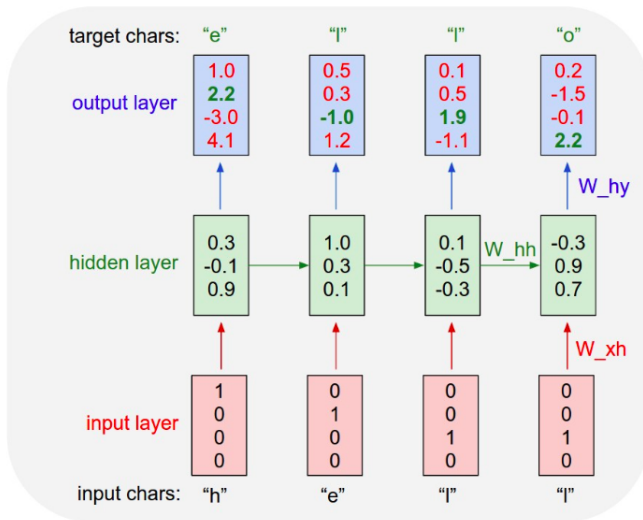
words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

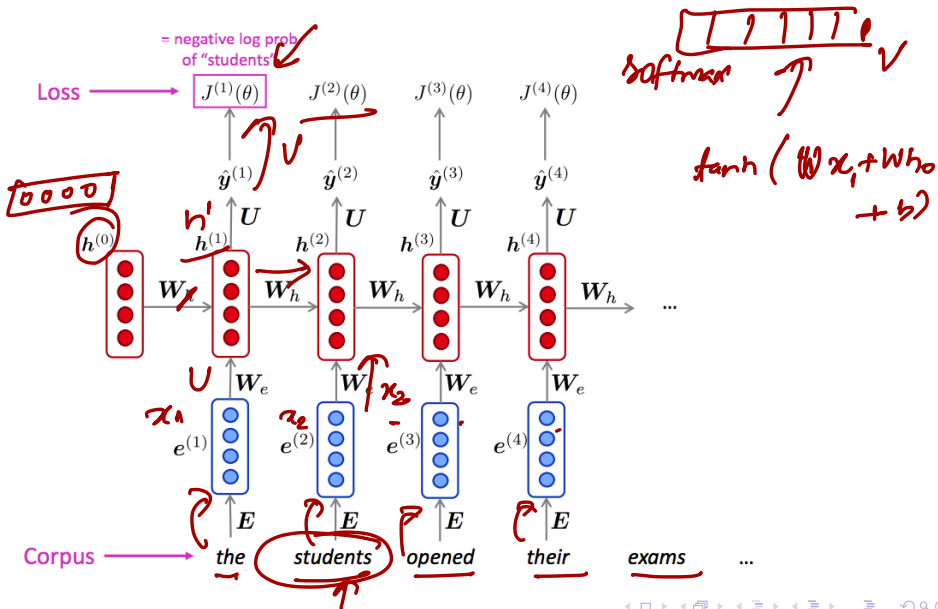
$$\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$$



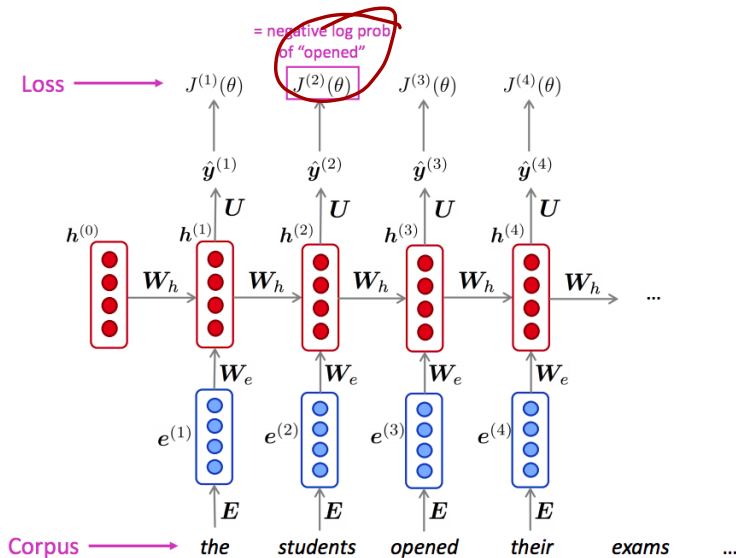
Example RNN



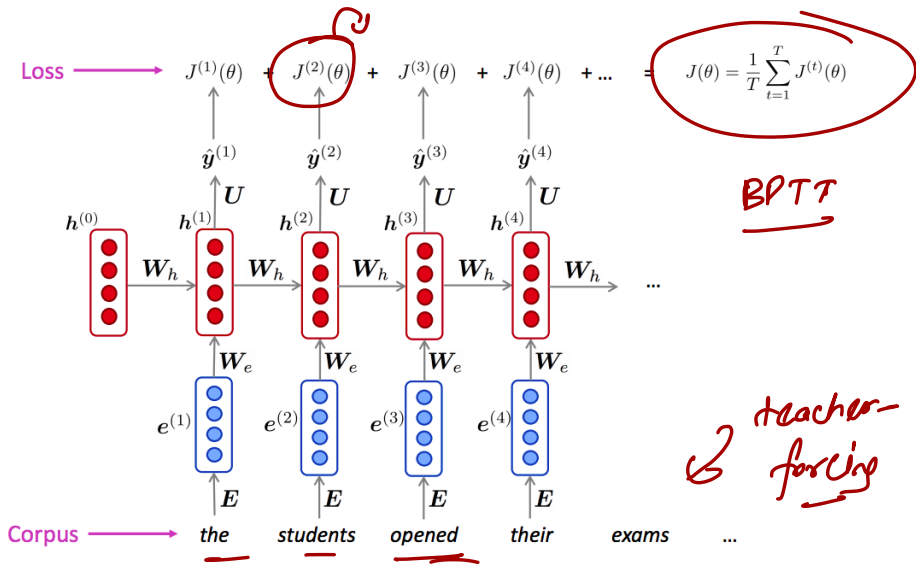
Training a RNN language model

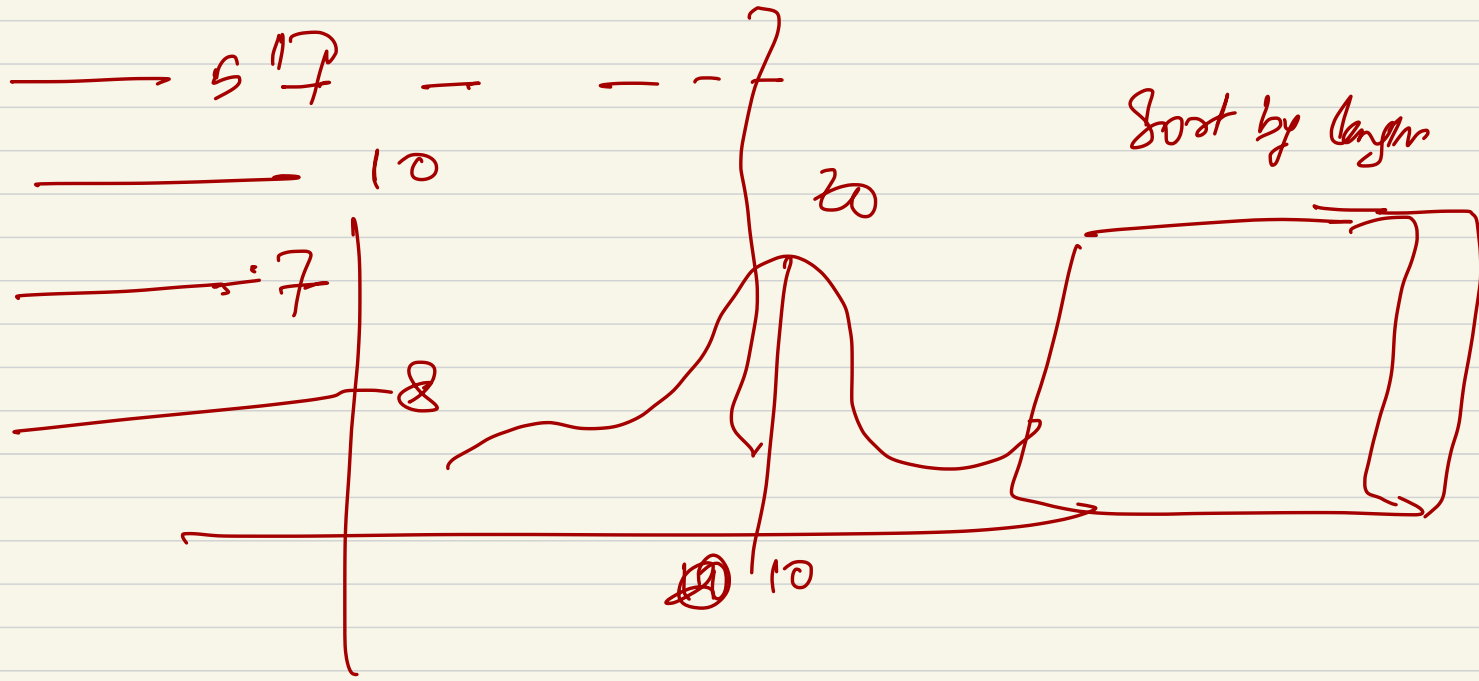


Training a RNN language model

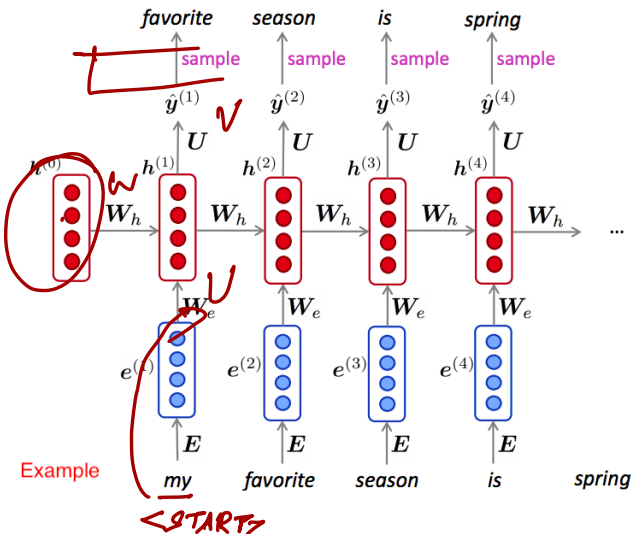


Training a RNN language model

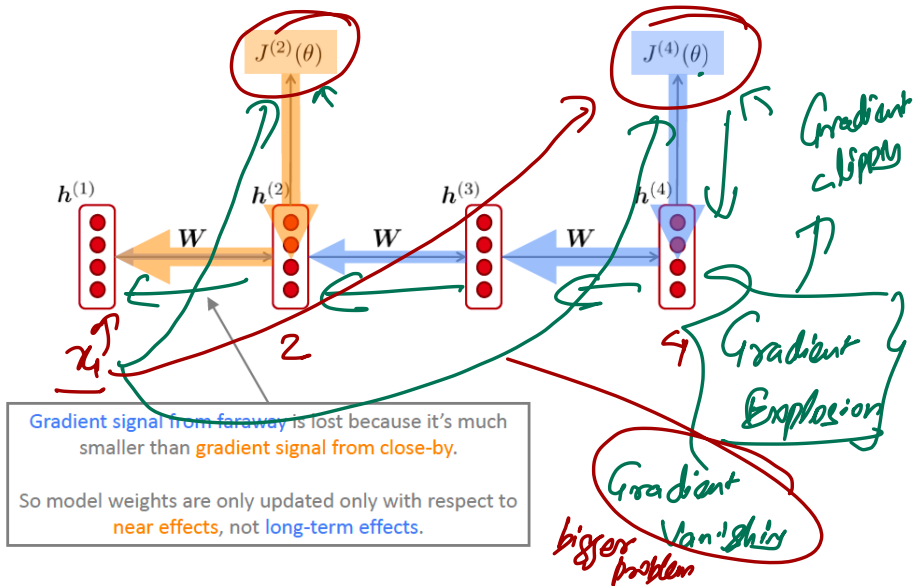




Generating text with a RNN Language Model



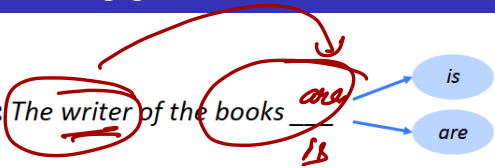


Vanishing Gradient Problem with RNNs



Effect of vanishing gradient on RNN LM

- **LM task:** *When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her _____*
- To learn from this training example, the RNN-LM needs to **model the dependency** between “tickets” on the 7th step and the target word “tickets” at the end.
- But if gradient is small, the model **can't learn this dependency**
 - So the model is **unable to predict similar long-distance dependencies** at test time

Effect of vanishing gradient on RNN LM

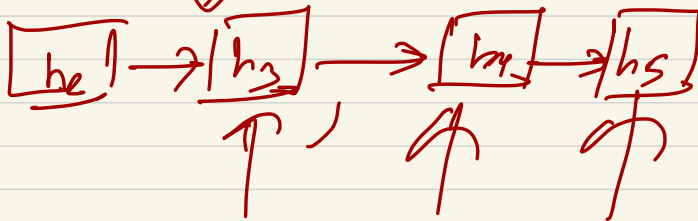
- **LM task:** *The writer of the books are* 
- **Correct answer:** *The writer of the books is planning a sequel*
- **Syntactic recency:** *The writer of the books is* (correct) 
- **Sequential recency:** *The writer of the books are* (incorrect) 
- Due to vanishing gradient, RNN-LMs are better at learning from sequential recency than syntactic recency, so they make this type of error more often than we'd like [Linzen et al 2016]



sub. singular

sub. singular

h_1



x_1

x_2

x_3

x_4

x_5

the

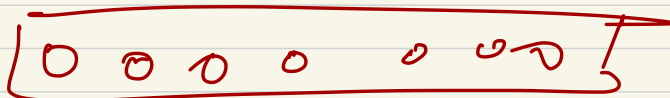
write

of

the

books

$$h_3 \sim \text{rank} (Ux + W \boxed{h_2} \rightarrow b)$$

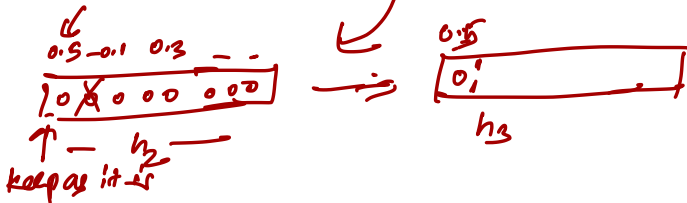


How to fix vanishing gradient problem?

- The main problem is that it is too difficult for the RNN to learn to preserve information over many timesteps.
- In a vanilla RNN, the hidden state is constantly being rewritten

$$h^{(t)} = \tanh(W h^{(t-1)} + U x^{(t)} + b)$$

- How about better RNN units?



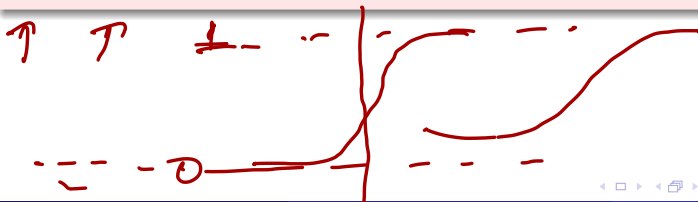
Using Gates for better RNN units

τ

- The gates are also vectors
- On each timestep, each element of the gates can be open (1), close (0) or somewhere in-between.
- The gates are dynamic: their value is computed based on the current context.

Two famous architectures

GRUs, LSTMs



RNN

$$\tanh(Ux_t + Wh_{t-1} + b)$$

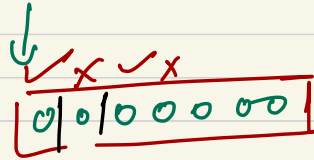
GRU

(2014)

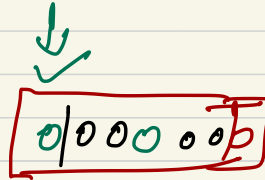
r_t Reset Gate

u_t Update Gate

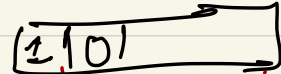
LSTM (1997)



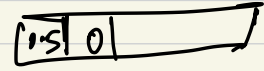
h_{t-1}



h_t



u_t



x_t

$$u_t = 1$$

$r_t = \text{anything}$

GRU



$$r_t = \sigma(U_r x_t + W_r h_{t-1} + b_r)$$

$$u_t = \sigma(U_u x_t + W_u h_{t-1} + b_u)$$

element wise

$$\tilde{h}_t = \tanh(U_h x_t + W_h (r_t \odot h_{t-1}) + b_h)$$

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

Gated Recurrent Units (GRU)

$[u_t = 1, r_t = 0]$ forget

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep t we have input $x^{(t)}$ and hidden state $h^{(t)}$ (no cell state).

$u_t = 0$
remember

Update gate: controls what parts of hidden state are updated vs preserved

Reset gate: controls what parts of previous hidden state are used to compute new content

New hidden state content: reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

Hidden state: update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

$$u^{(t)} = \sigma(W_u h^{(t-1)} + U_u x^{(t)} + b_u)$$

$$r^{(t)} = \sigma(W_r h^{(t-1)} + U_r x^{(t)} + b_r)$$

gates

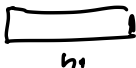
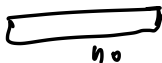
$$\tilde{h}^{(t)} = \tanh(W_h(r^{(t)} \circ h^{(t-1)}) + U_h x^{(t)} + b_h)$$

$$h^{(t)} = (1 - u^{(t)}) \circ h^{(t-1)} + u^{(t)} \circ \tilde{h}^{(t)}$$

$u_t = 0$

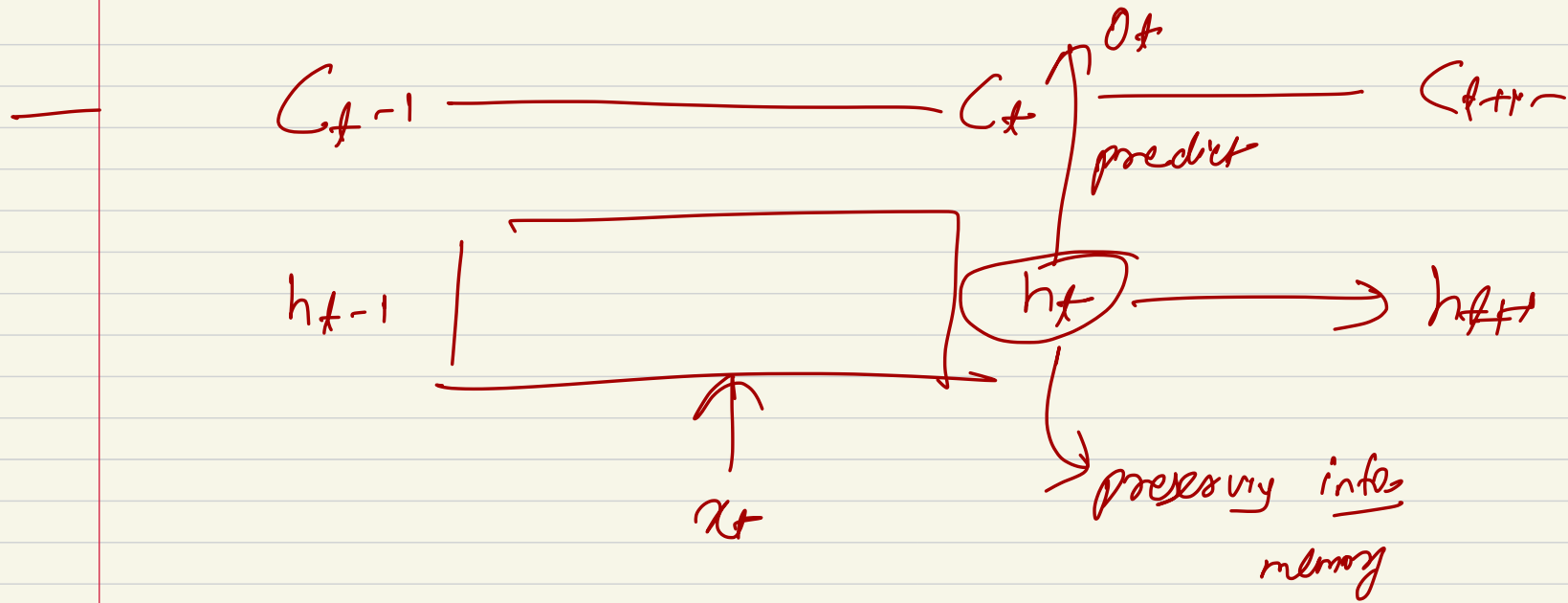
How does this solve vanishing gradient?

Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)



writer





LSTM

cell state

LSTM

three gates

$$\frac{C_{t-1}}{h_{t-1}}$$

$$\frac{C_t}{h_t}$$

i/p gate

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

o/p gate

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

forget gate

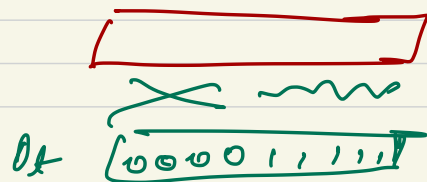
$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

1. Update cell content

$$\tilde{C}_t = \tanh(W_c h_{t-1} + U_c x_t + b_c)$$

$$\rightarrow \underline{C_t} = f_t \circ C_{t-1} + i_t \circ \underline{\tilde{C}_t}$$

$$\rightarrow \underline{h_t} = \underline{o_t \circ \tanh(C_t)}$$



Long Short Term Memory (LSTM)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep t :

Forget gate: controls what is kept vs forgotten, from previous cell state

Input gate: controls what parts of the new cell content are written to cell

Output gate: controls what parts of cell are output to hidden state

New cell content: this is the new content to be written to the cell

Cell state: erase ("forget") some content from last cell state, and write ("input") some new cell content

Hidden state: read ("output") some content from the cell

Sigmoid function: all gate values are between 0 and 1

$$\begin{aligned}f^{(t)} &= \sigma(W_f h^{(t-1)} + U_f x^{(t)} + b_f) \\i^{(t)} &= \sigma(W_i h^{(t-1)} + U_i x^{(t)} + b_i) \\o^{(t)} &= \sigma(W_o h^{(t-1)} + U_o x^{(t)} + b_o)\end{aligned}$$

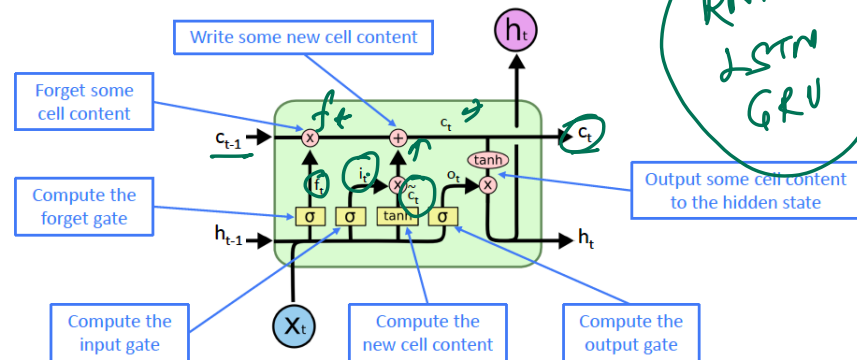
$$\begin{aligned}\tilde{c}^{(t)} &= \tanh(W_c h^{(t-1)} + U_c x^{(t)} + b_c) \\c^{(t)} &= f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)} \\h^{(t)} &= o^{(t)} \circ \tanh c^{(t)}\end{aligned}$$

Gates are applied using element-wise product

All these are vectors of same length n

Long Short Term Memory (LSTM)

You can think of the LSTM equations visually like this:



RNN
LSTM
GRU

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

