

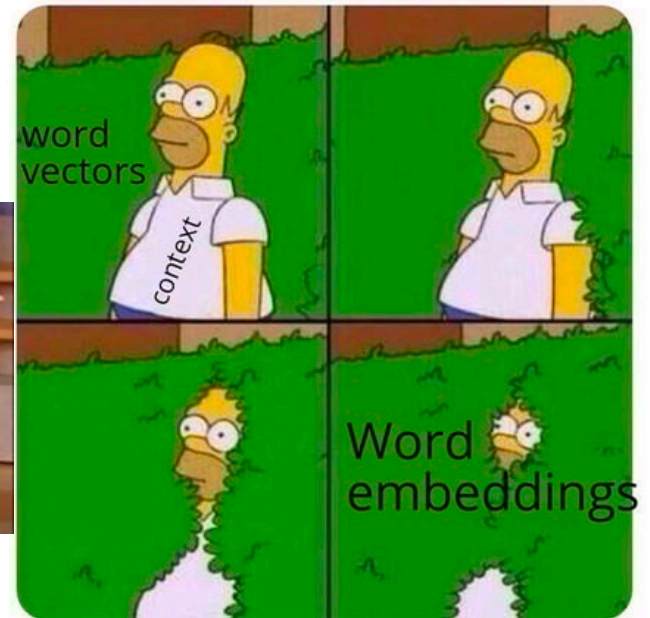
May 5, 2023.

DSML: NLP module.

Word embeddings
in a nutshell

Long Short-Term Memory
(LSTM)

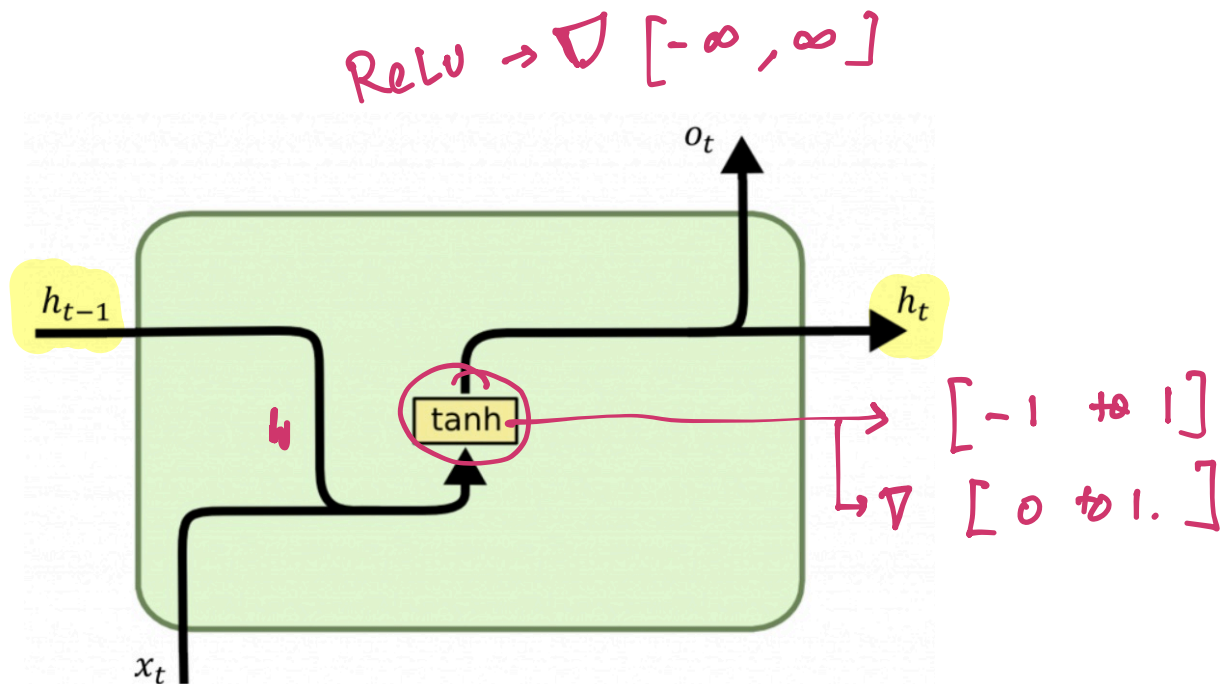
Class starts
@ 9:05



When you penalize your Natural
Language Generation model for
large sentence lengths



Recap:



* Recurrent Neural Networks (RNNs).

→ What problems do they solve?

* Variable input/output lengths. * Explicit encoding of context - hidden state.

→ Types?

* one-to-one * one-to-many * Many-to-one * Many-to-many

→ Training?

Backpropagation through time (BPTT).

Agenda:

* A new type of Architecture:

long Short-term Memory (LSTM).

→ Why? RNN training is difficult because of vanishing/exploding gradients.

* Application: Text Summarization.

→ Encoder - Decoder Architecture.

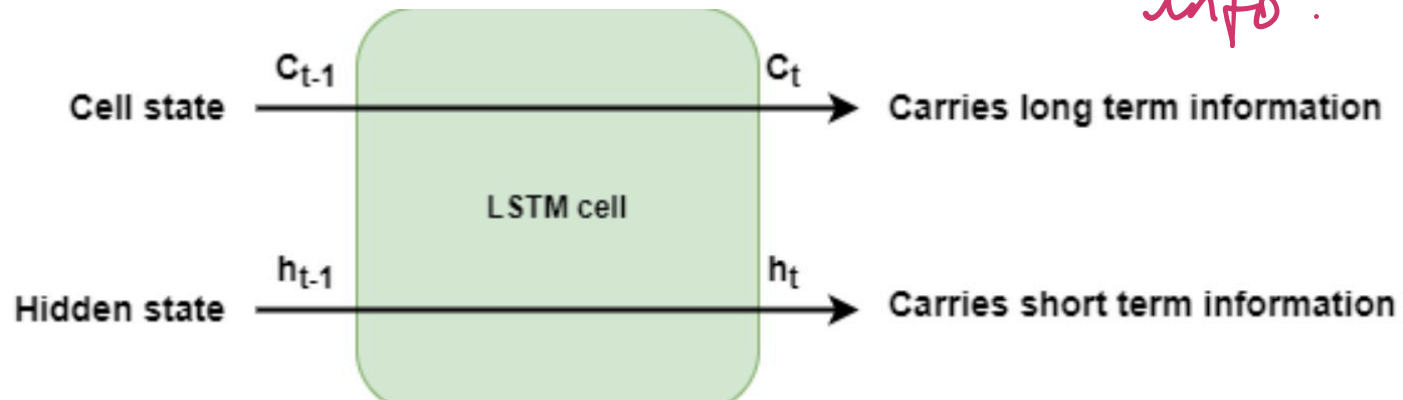
→ Input, Forget and Output gates.

→ Gated Recurrent Units.

→ ROUGE score.

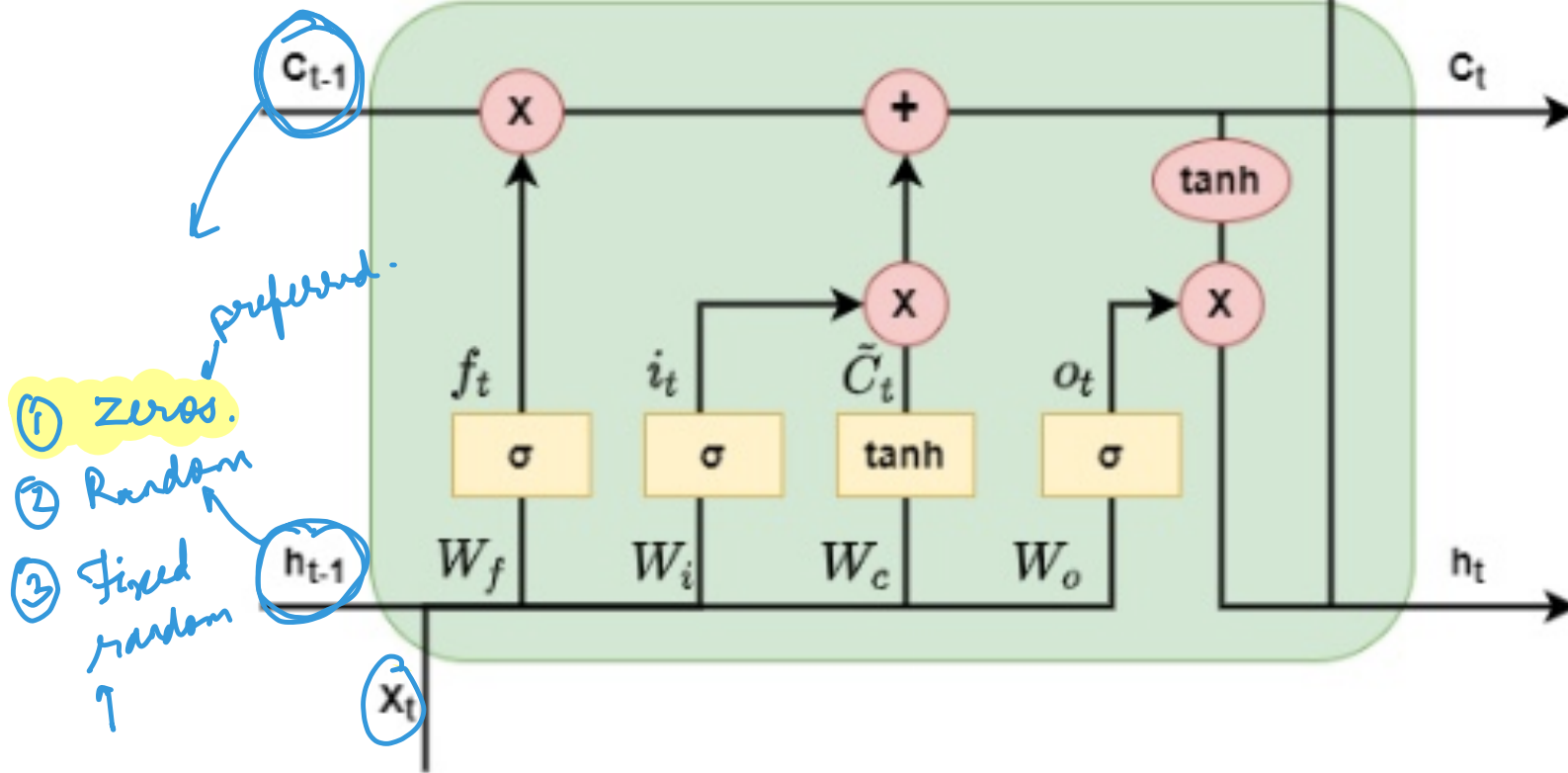
Let's understand the Industry Requirement !!

Conceptual change: Add memory to store long-term info.

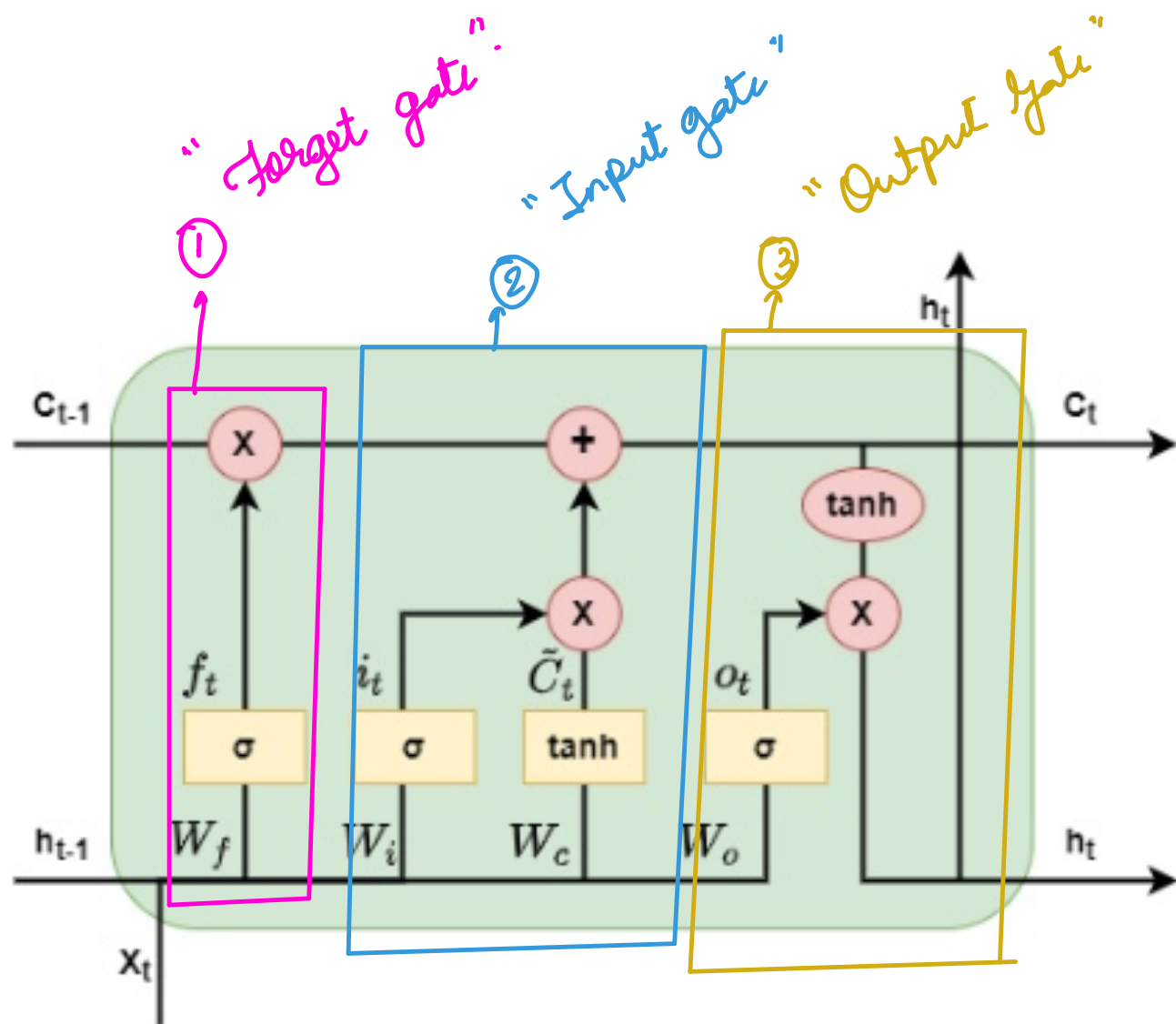


Practical Implementation.

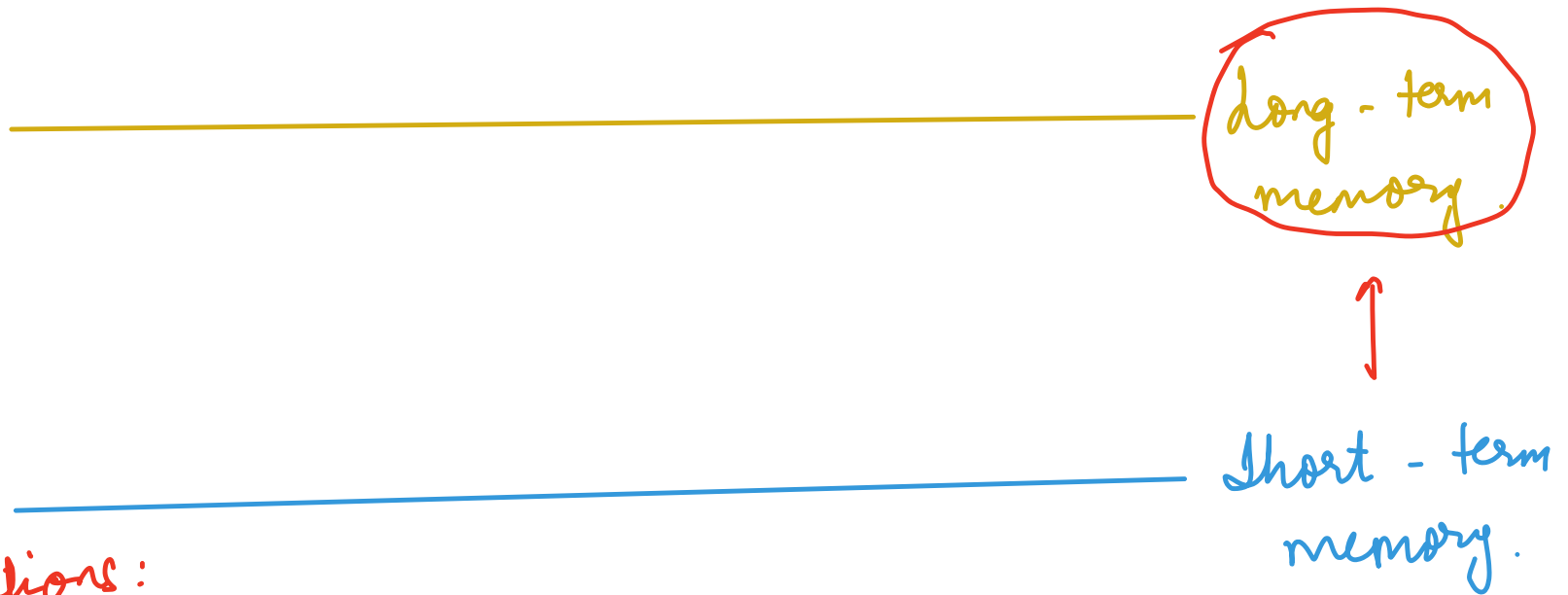
g have a fat cat.
1 2 3 4 5



Input word at
time step t .



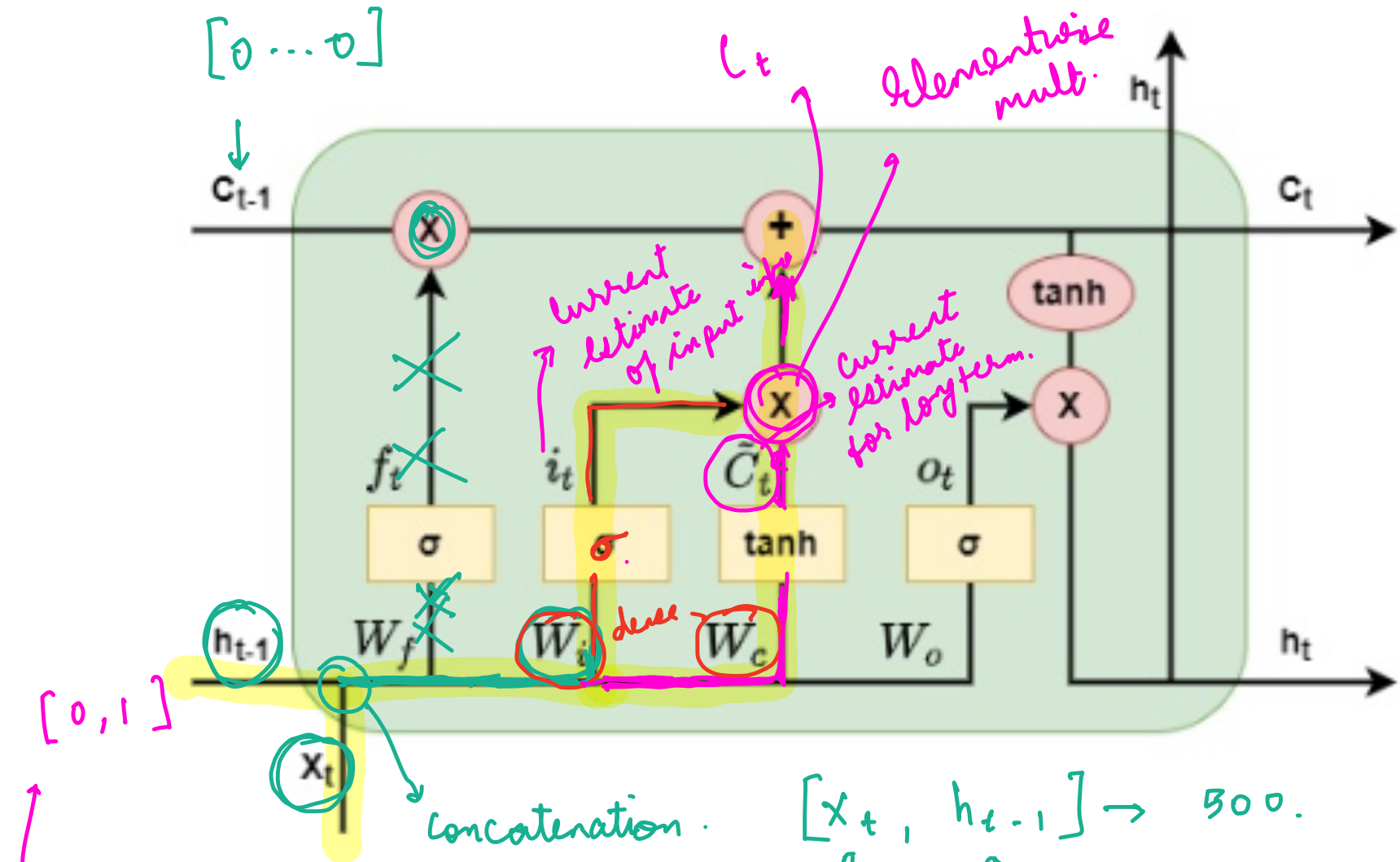
Conceptual Explanation.



Operations:

- (i) Input operation: Allow info to flow from short-term to long-term.
- (ii) Forget operation: Allow unnecessary information to be removed from the long-term.
- (iii) Output operation: Retrieve info from long term.

Input gate: how information flows into long-term memory.



$$\begin{aligned}
 i_t &= \sigma(W_i [x_t, h_{t-1}] + b_i) \\
 \tilde{c}_t &= \tanh(W_c [x_t, h_{t-1}] + b_c) \rightarrow (-1, 1)
 \end{aligned}$$

Dimensional analysis: $[x_t, h_{t-1}] \rightarrow 500$.
 x_t is 200, h_{t-1} is 300.

what to put in long term.

C_t

=

i_t

Relative weights
wrt. current input.

*

↑
elementwise

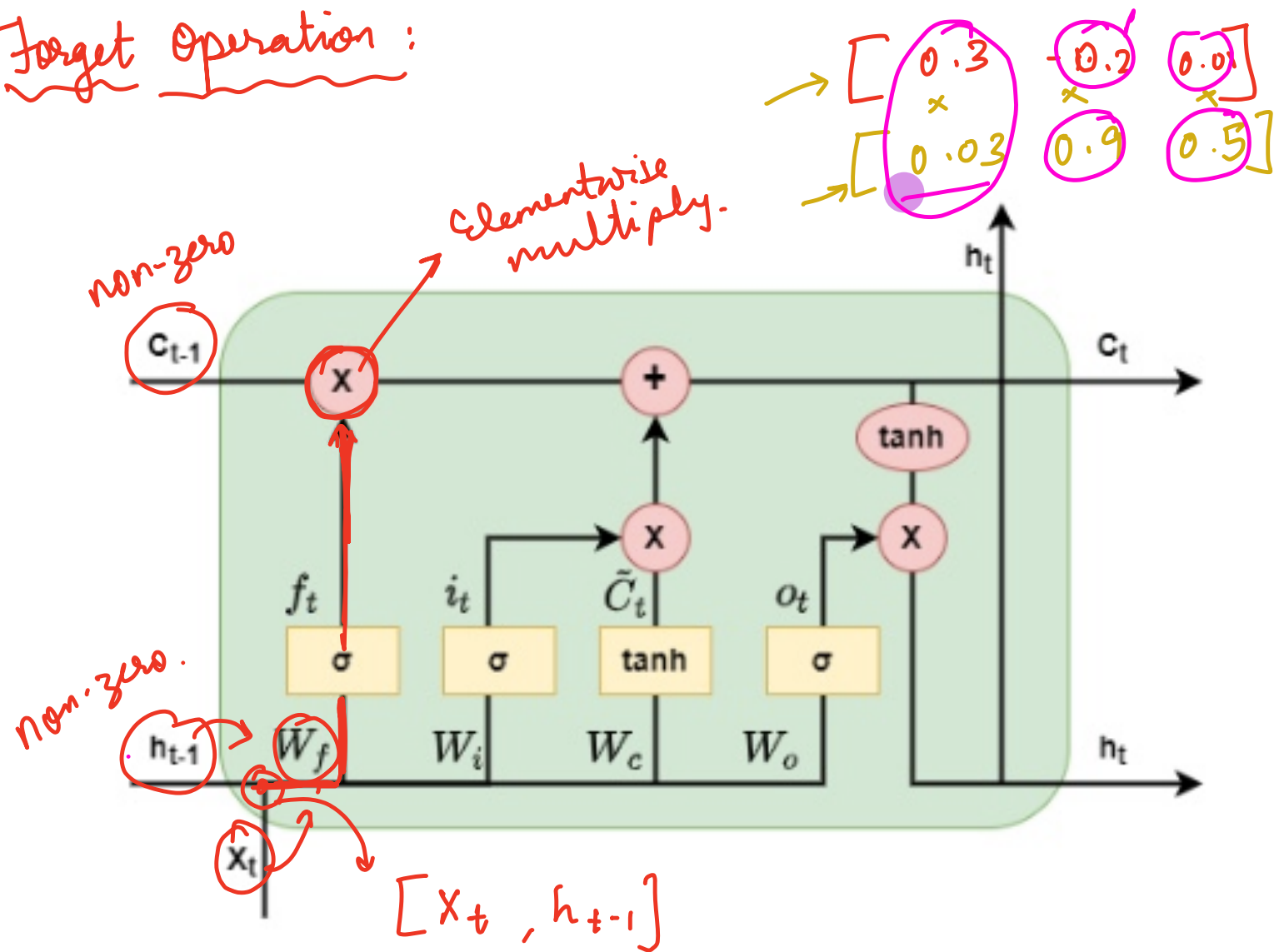
\tilde{C}_t

Context info from
current word.

$\begin{bmatrix} 0.1 & 0.5 & 0.3 \end{bmatrix}$
↑

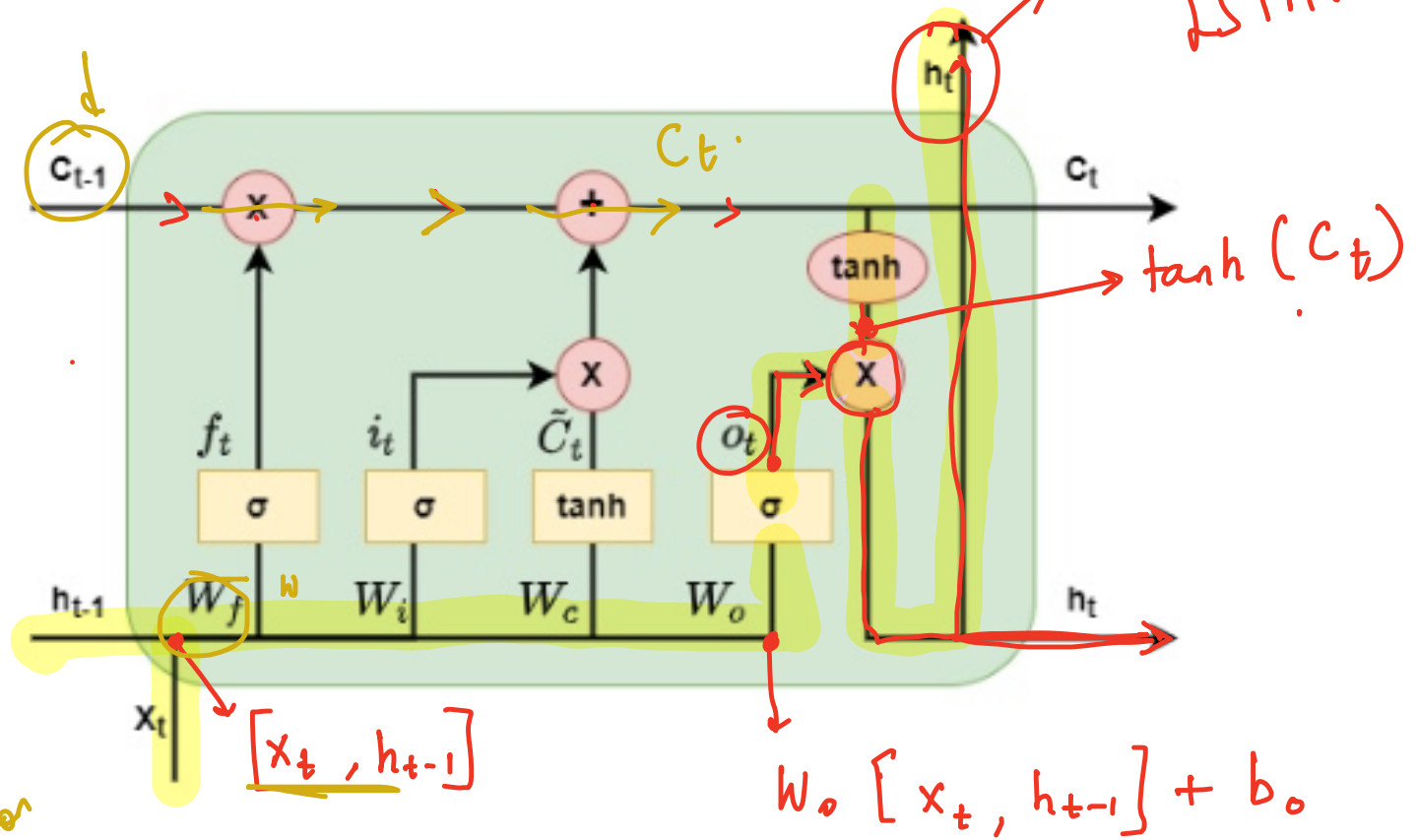
$\begin{bmatrix} \downarrow & \downarrow & \\ -0.8 & -0.3 & 0.7 \end{bmatrix}$
↑

Forget Operation:



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

Output gate:



contribution from long-term

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

$$h_t = \tanh(C_t) * o_t$$

↑
elementwise

contribution from short term

Hidden State

Which of the following statements is correct about Hidden State ?

 HINTS 

 [Complete Solution](#)

You will get full points if and only if you give CORRECT ANSWER in first attempt. All later attempts will get you ZERO score.

- ☐ It is present in both MLP and RNN
- ☐ When we are passing the first input the Hidden states are all generally initialized according to HE initialization.
- ☒ When we are passing the first input the Hidden states are all typically initialized to zeros
- ☒ Activation Function typically used in Vanilla RNN is tanh

Feedback:
Include more detailed explanation.