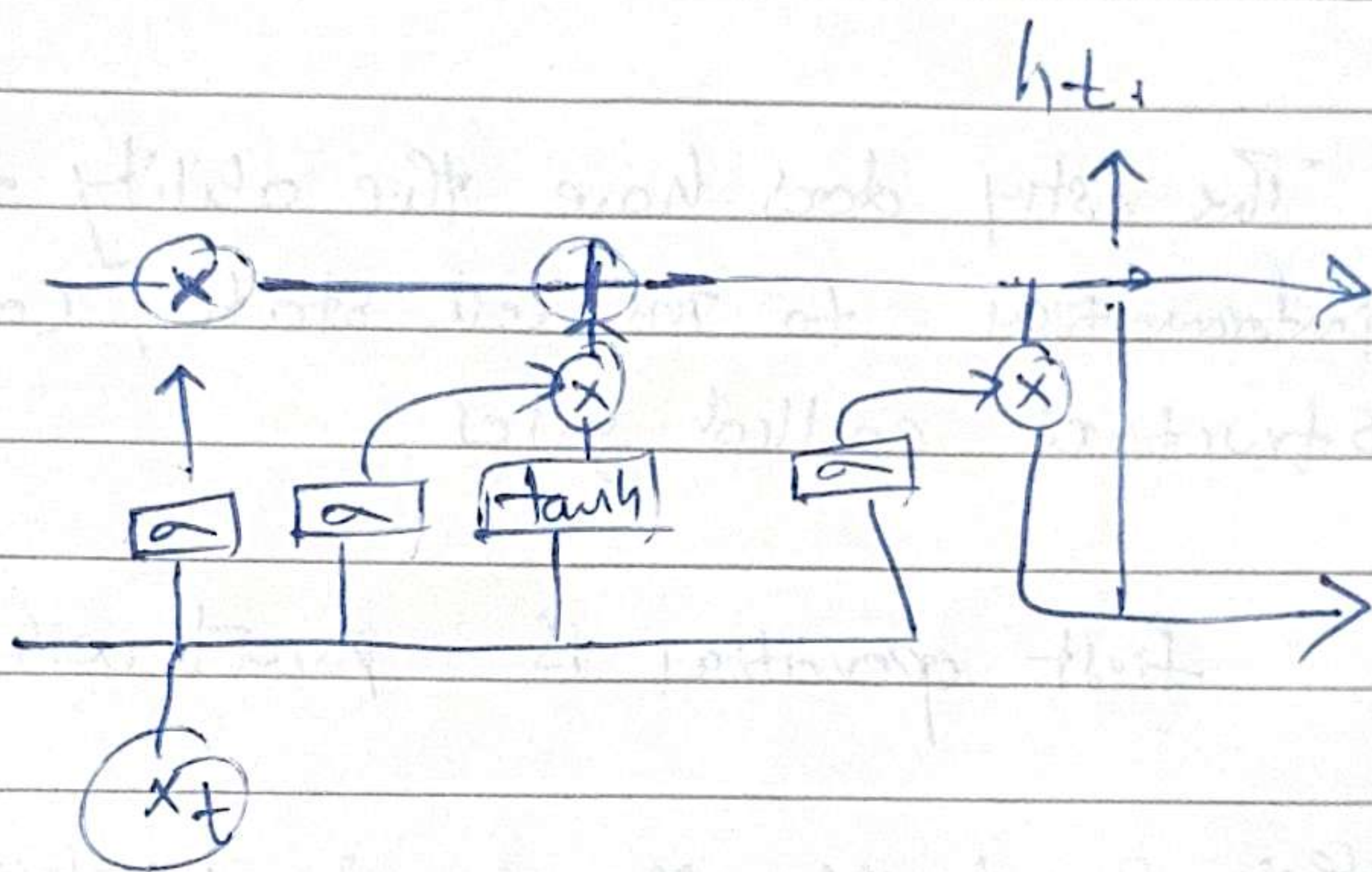


# LSTM:

LSTMs are explicitly designed to avoid the long term dependency problem. Remembering information for long periods of time is practically their default behavior not something they struggle to learn.

In LSTMs have this chain like structure, but repeating module has a different structure instead of having the single neural network layer, there are four interacting layers.



Memory cell

$\square \Rightarrow$  Neural network layer

forget gate

$\circ \Rightarrow$  pointwise operation

input gate

$\rightarrow \Rightarrow$  vector transfer

output gate

$\Rightarrow \Rightarrow$  concatenate

$\rightleftarrows \Rightarrow$  copy

Memory cell is used to remember or forget the memory based on the context of the input.

ex: text generator the nouns which are not been used are forgotten.



The key to LSTMs is the cell state, the horizontal line running through the top of the diagram.

The cell state is kind of like a conveyor belt. It runs straight down the entire chain with only some minor interaction. It's very easy for information to just flow along it unchanged.



The LSTM does have the ability to remove (or add) information to the cell state, carefully regulated by structures called gates.

first operation is pointwise operation.

this pointwise operation is a matrix multiplication to the  $\begin{bmatrix} 1 & 0 \end{bmatrix}$   $0 \rightarrow$  to forget the data  $1 \rightarrow$  retain the data.

$\oplus$  addition operation in last step we are forgetting some information so additional information need to be added.

Gates are a way to optimally let information through. They are composed out of a sigmoid neural net layer and a pointwise multiplication operation.

The Sigmoid layer outputs numbers between 0 to 1. how much content should be let through. A value of zero means let nothing through while a value of one means let everything through.



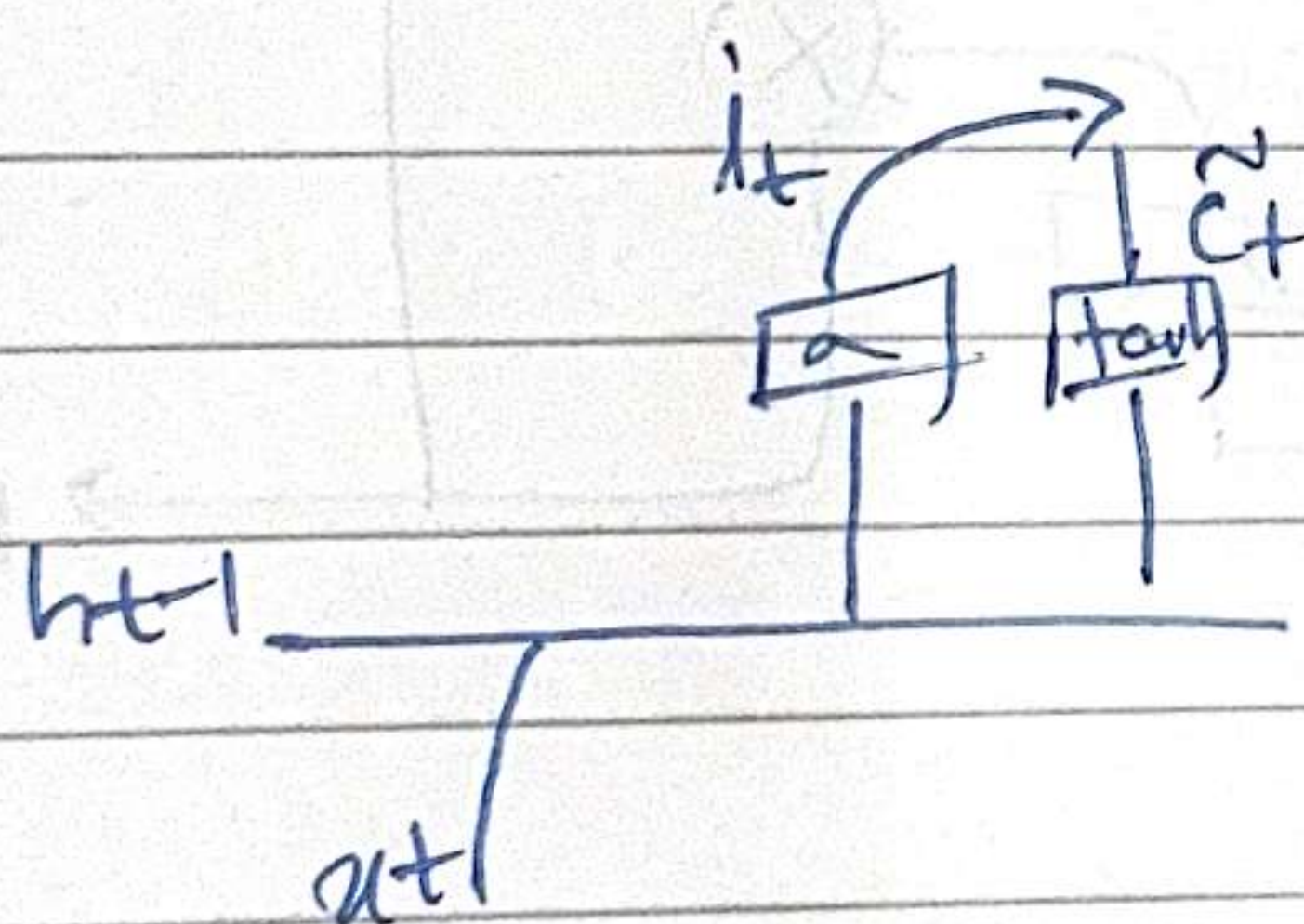
① The first Step is to decide what information we are going to throw away from the cell state. And the decision is made by a Sigmoid layer and the gate is called forget gate.

$\Rightarrow$  It looks at  $h_{t-1}$  and  $x_t$  and outputs a number between  $0 \rightarrow 1$  for each number in the cell state  $c_t$ .

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

② The next Step is to decide what information we are going to store in the cell state. This has two parts.

first Sigmoid layer called the input gate layer decides which values we will update next a tanh layer creates a vector of new candidate values,  $\tilde{c}_t$



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

$$\tilde{c}_t = \tanh(\omega_c \cdot [h_{t-1}, x_t] + b_c)$$



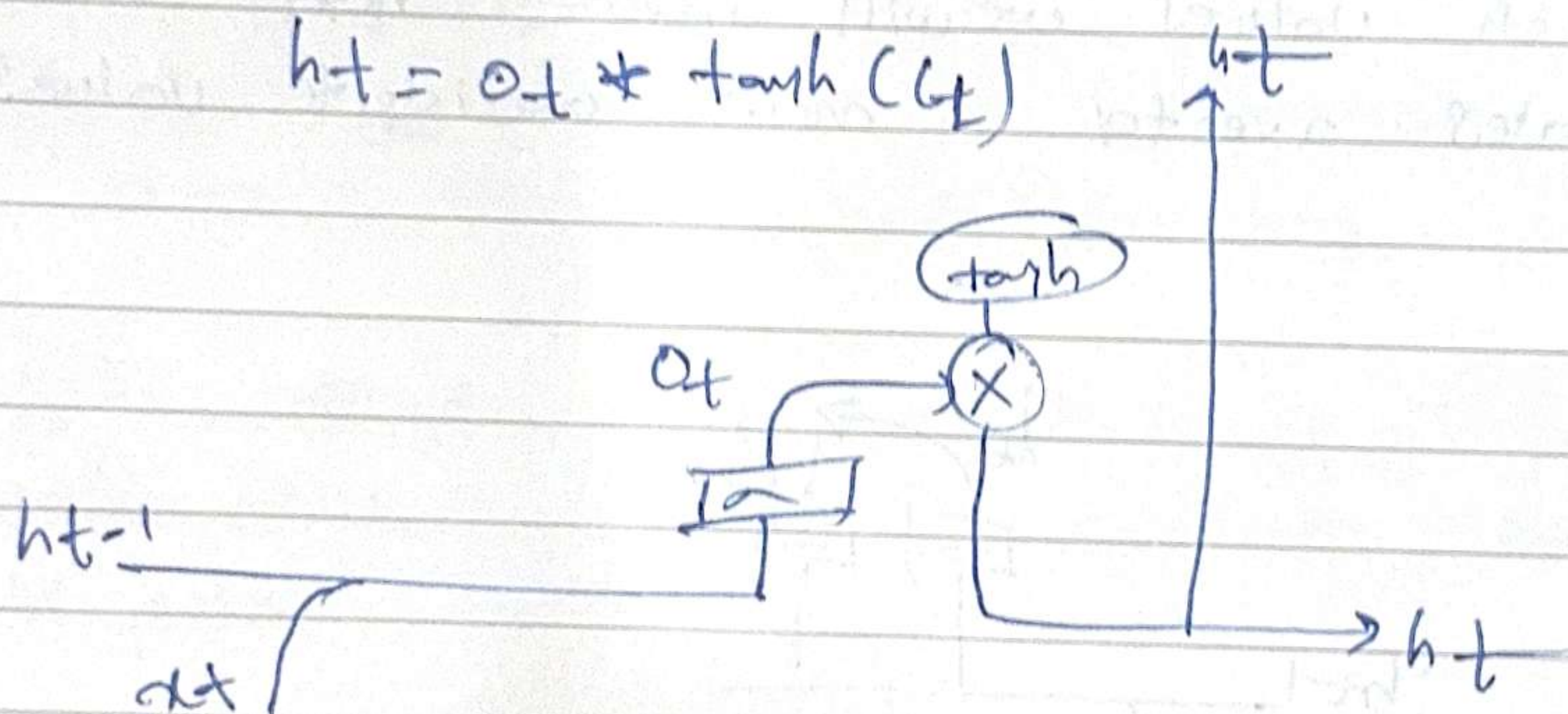
Now we will update the old cell state  $C_{t-1} \rightarrow C_t$

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

This output will be based on our cell state first we will run the Sigmoid layer which decides what parts of the cell state we are going to output then we put the cell state through tanh and multiply it by the output of the sigmoid gate

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

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



And the hidden state  $h_t$  short memory will update and then goes to next  $h_{t+1}$  hidden state.



Gers & Schmidhuber: adding peephole connections

Here we let the gate to look at the cell state.

$$f_t = \alpha(w_f \cdot [c_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \alpha(w_i \cdot [c_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \alpha(w_o \cdot [c_t, h_{t-1}, x_t] + b_o)$$

Another variation is to use the coupled forget and input gate. Instead of separately deciding what to forget and what to add we make these two decisions together.

$$c_t = f_t * c_{t-1} + (1 - f_t) * \tilde{c}_t$$

1. Hidden state  $\rightarrow h_t$

Short term memory

Directly used for predictions

Shape (hidden-unit, 1)

2. cell state  $\rightarrow c_t$  (Long term)

Shape (hidden-unit, 1)