

(RNN)

- ④ Mini batch training: Updating the weights once every N steps:

$$\delta = \frac{1}{N} \sum_{k=1}^N \delta_{ijk}$$

→ It helps reducing the complexity of training

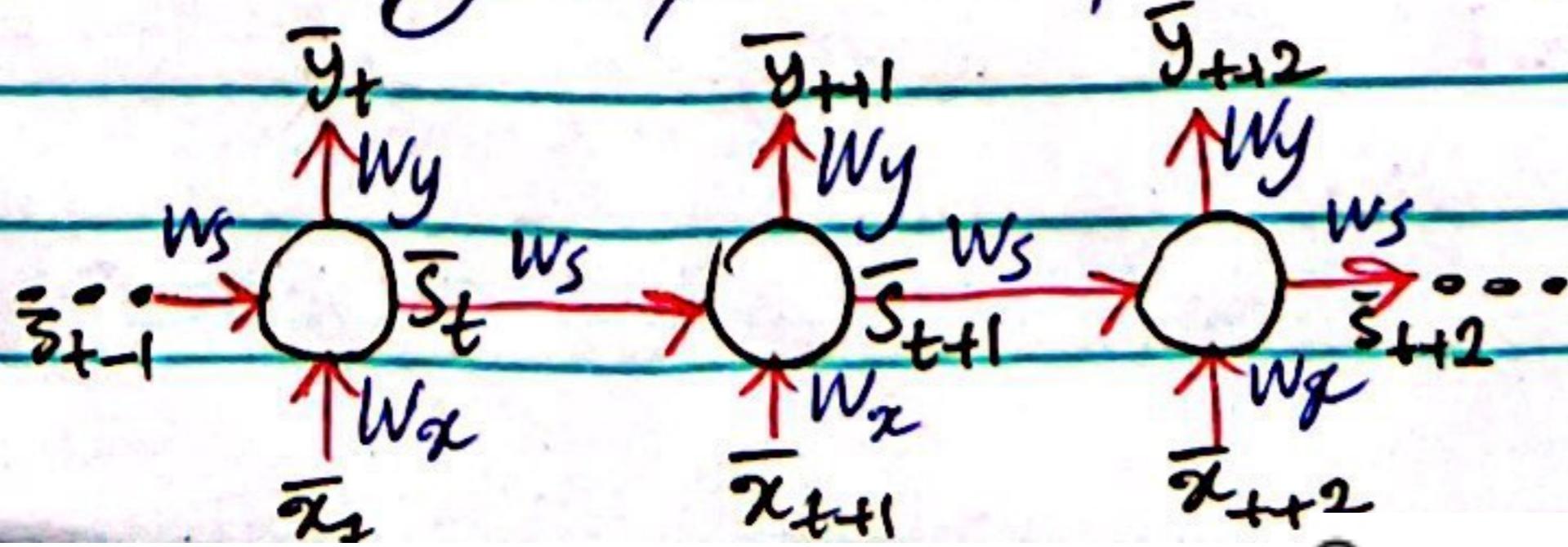
→ Good for noise reduction in training stage

- ④ Recurrent means occurring often or frequently.

→ In RNN we perform the same task for inputs in a sequence.

- ④ RNN Captures the information from the previous input by maintaining an internal memory elements called states.

- ④ It is useful for applications with temporal dependency
 ≡ the current output depends not only on the current input but also on a memory element which takes into account past inputs. E.g., Predicting the next word in a sentence: Sentiment analysis, Speech recognition, time-series analysis, NLP, and gesture recognition.



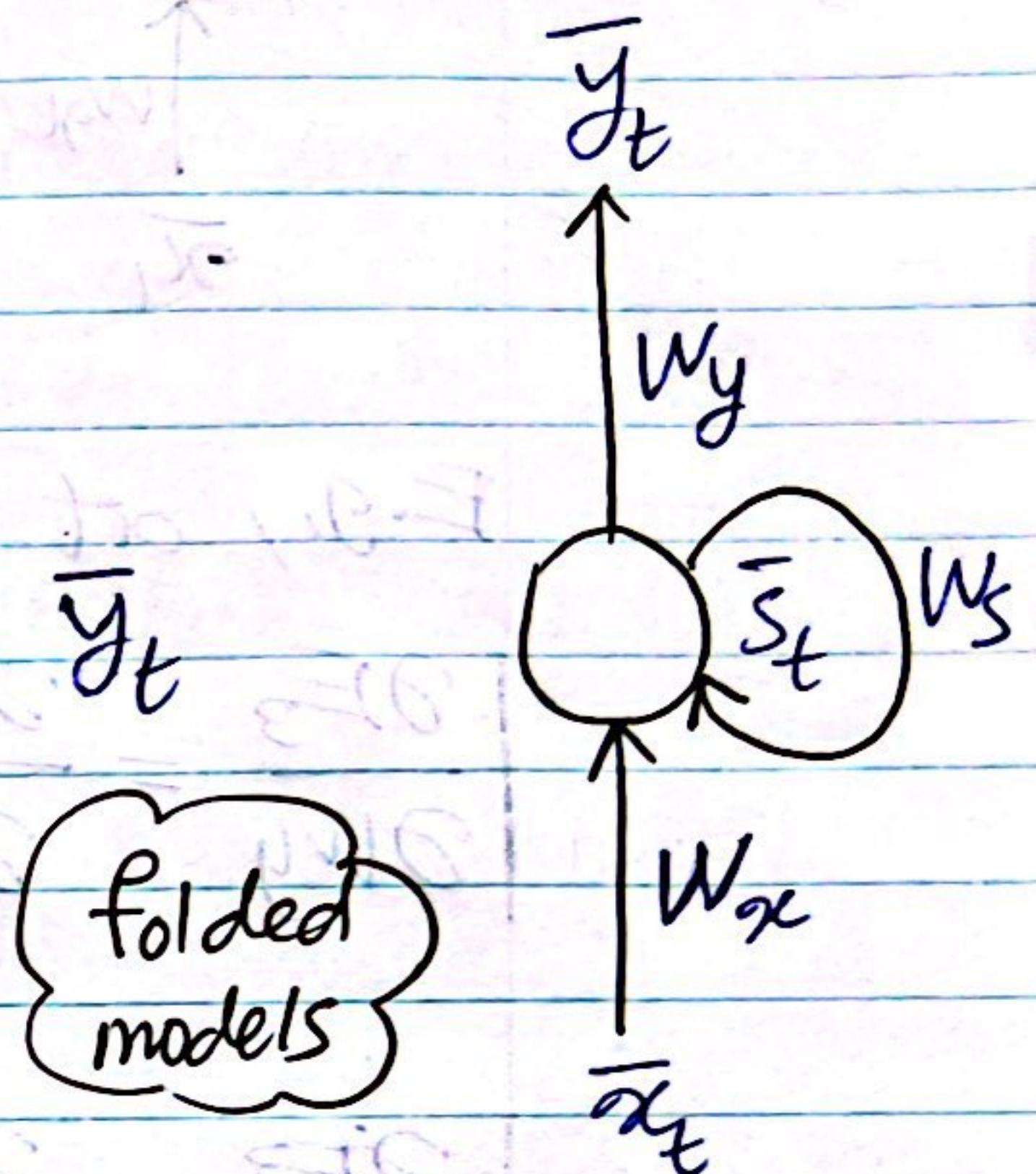
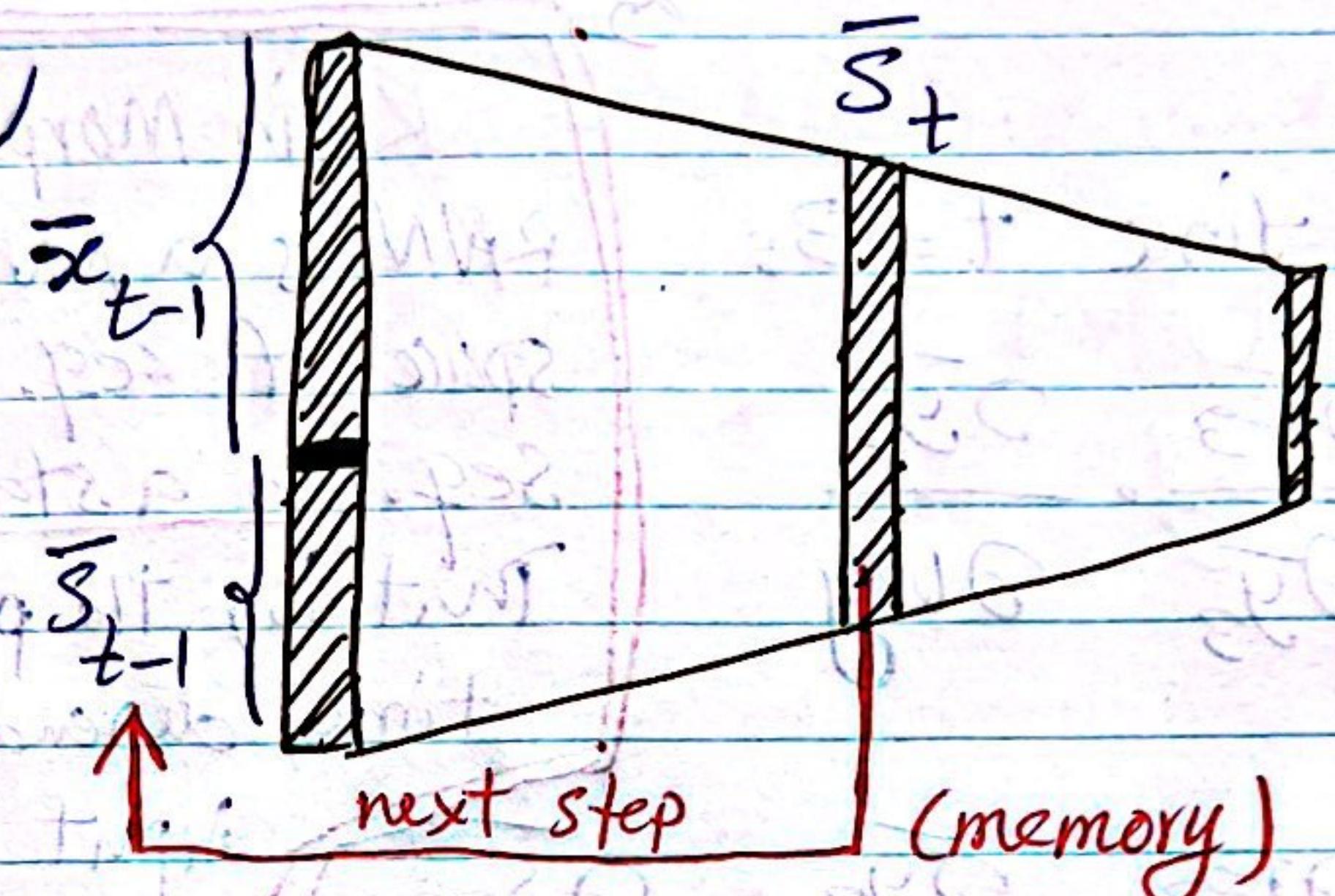
* The core differences between RNN and FFNN are:

b) the manner we define our outputs:

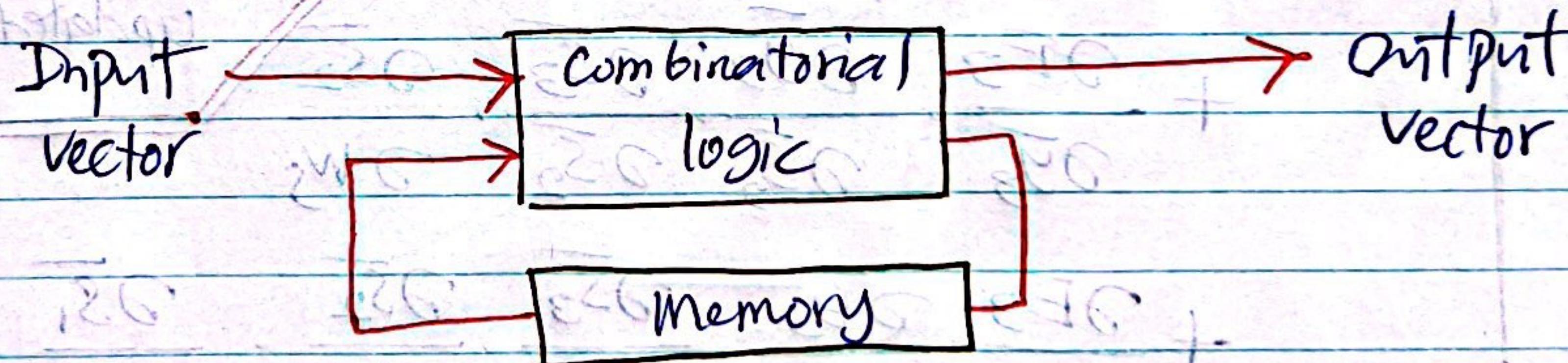
2) memory elements or states.

* Simple RNN

(Elman
Network)



It is equivalent to a state machine:

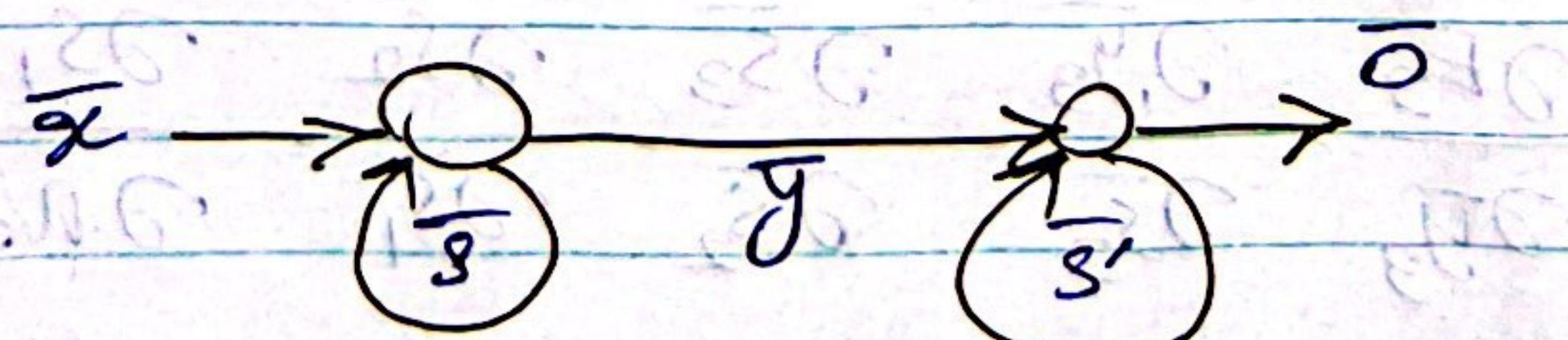


$$\bar{s}_t = \phi(\bar{x}_t \cdot W_x + \bar{s}_{t-1} \cdot W_s)$$

$$\bar{y}_t = \bar{s}_t \cdot W_y \text{ or } \sigma(\bar{s}_t \cdot W_y)$$

weights W remain the same for a given sequence of input!

* We can stack RNNs:

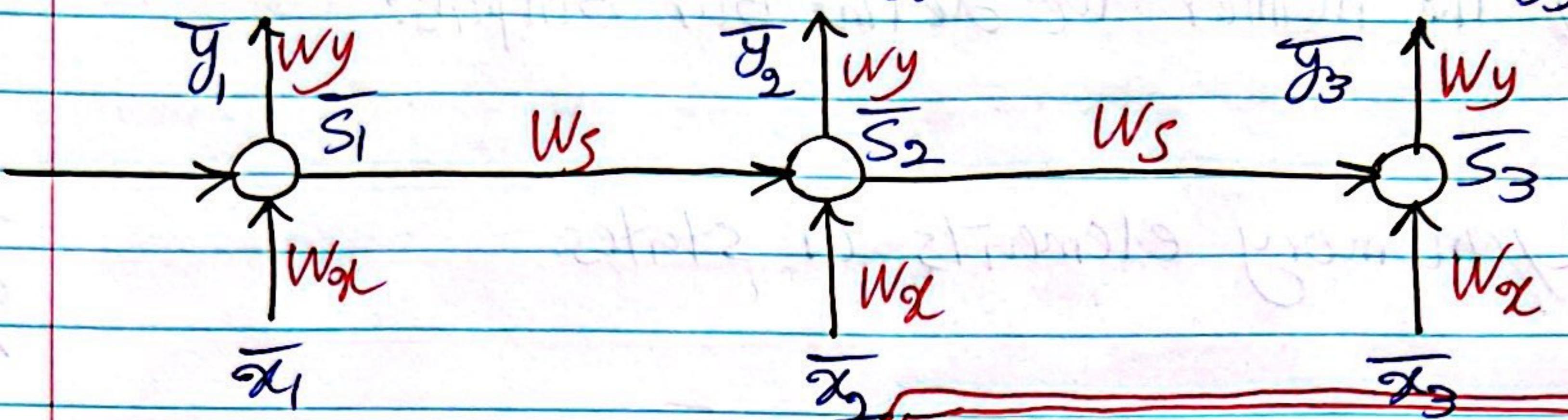


* Backpropagation Through Time (BPTT)

$$E_1 = (\bar{d}_1 - \bar{y}_1)^2$$

$$E_2 = (\bar{d}_2 - \bar{y}_2)^2$$

$$E_3 = (\bar{d}_3 - \bar{y}_3)^2$$



E.g., at time $t=3$:

$$\frac{\partial E_3}{\partial W_y} = \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial W_y}$$

$$\begin{aligned} \frac{\partial E_3}{\partial W_s} &= \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial W_s} \\ &+ \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \cdot \frac{\partial \bar{s}_2}{\partial W_s} \end{aligned}$$

$$+ \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \cdot \frac{\partial \bar{s}_2}{\partial \bar{s}_1} \cdot \frac{\partial \bar{s}_1}{\partial W_s}$$

$$\begin{aligned} \frac{\partial E_3}{\partial W_x} &= \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial W_x} \\ &+ \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \cdot \frac{\partial \bar{s}_2}{\partial W_x} \end{aligned}$$

$$+ \frac{\partial E_3}{\partial \bar{y}_3} \cdot \frac{\partial \bar{y}_3}{\partial \bar{s}_3} \cdot \frac{\partial \bar{s}_3}{\partial \bar{s}_2} \cdot \frac{\partial \bar{s}_2}{\partial \bar{s}_1} \cdot \frac{\partial \bar{s}_1}{\partial W_x}$$

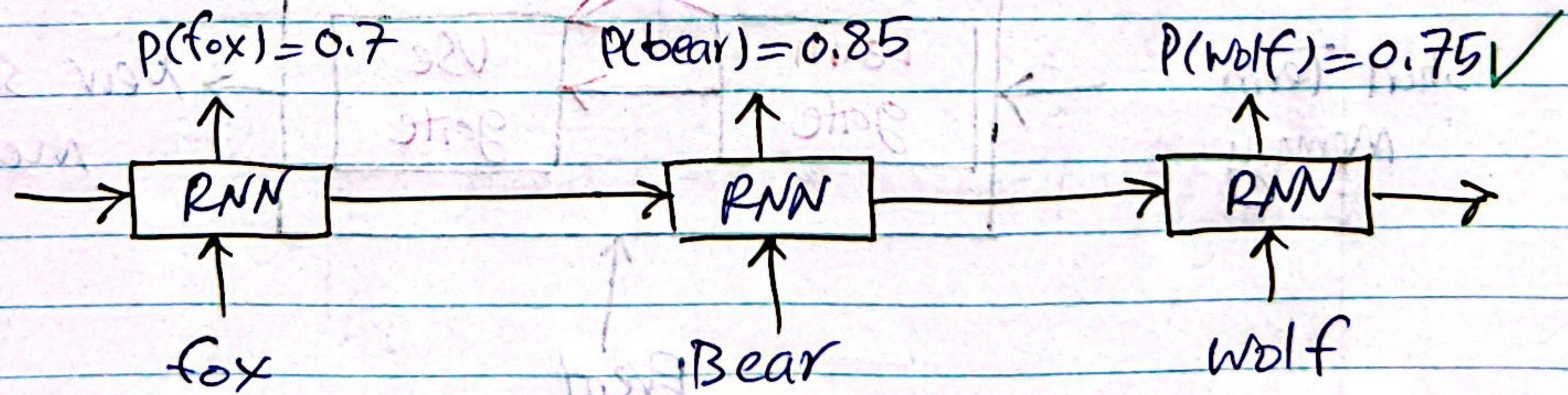
Kevin Murphy:
 RANN is a MN that maps the input space of seq. to output space of seq. in a stateful way!
 That is, the prediction of y_t over time, depends not only on the input x_t , but also on the hidden state of the sys, s_t , which gets updated over time!

* Problems: 1/ Vanishing Gradient remedy LSTM
 with RNN 2/ Exploding Gradient remedy Gradient clipping

LSTM: Long Short Term Memory Networks

These nets are useful when we need to switch between remembering recent things, and things from long time ago.

* Let's say, we have an image classifier and it has a picture of wolf \rightarrow There is a high risk of missclassification of being dog! If we feed the net with a video which previous images are fox and bear, it is more probable that net classifies the image as wolf!



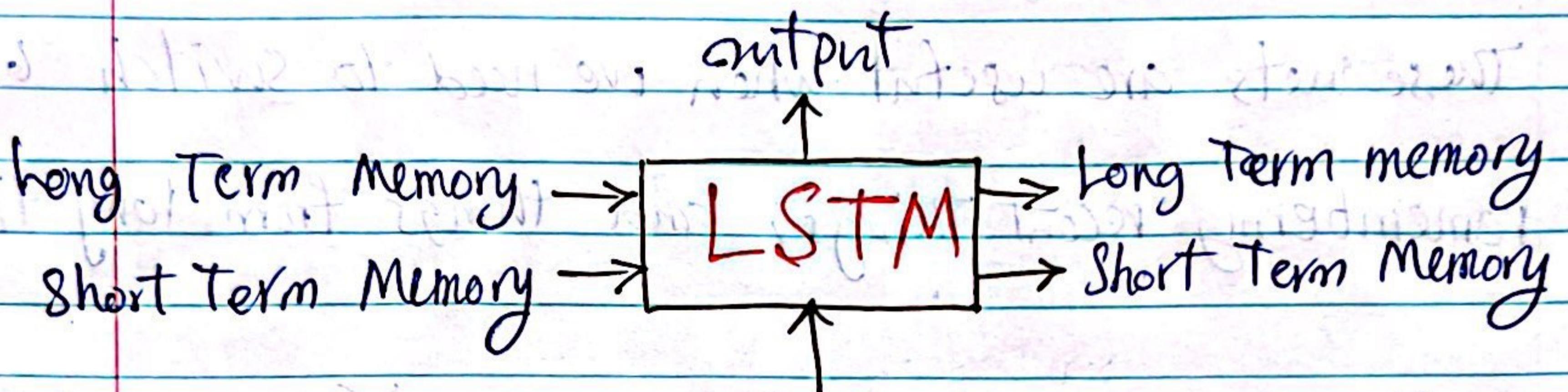
* This is because the network knows they are wild!

But, if we have these frames:

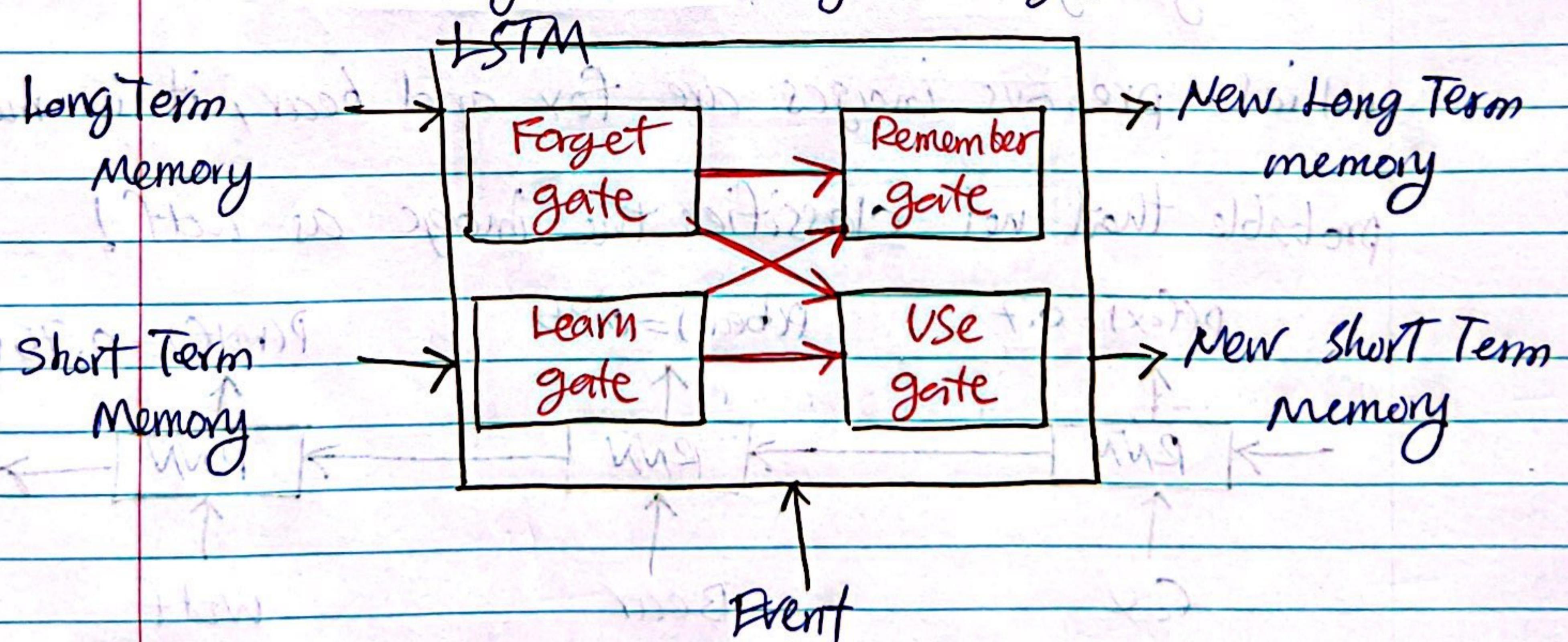
Bear \rightarrow tree \rightarrow squirrel \rightarrow wolf

(23)

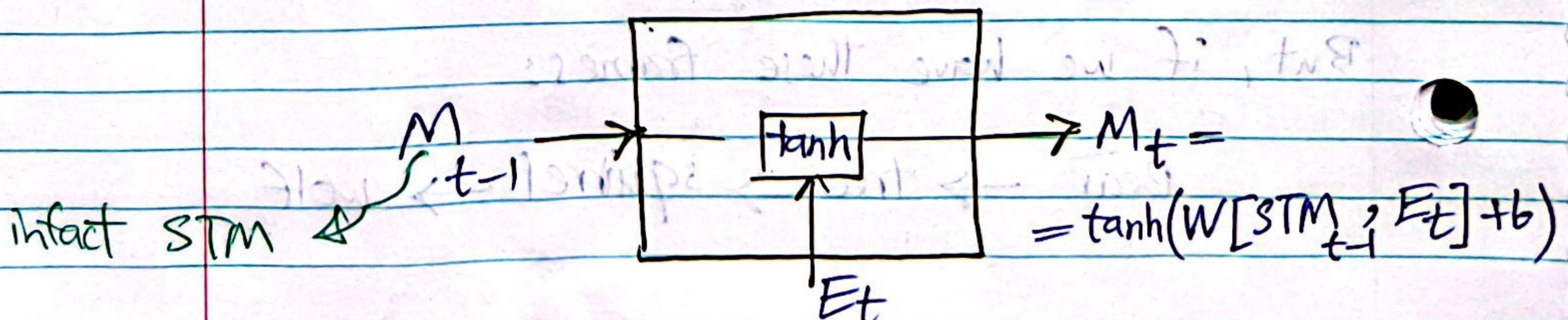
in this case, the net may forget about wild animals and missclassify wolf as dog! This is the problem with RNN! Due to vanishing gradient. This is where LSTM helps us! MT2



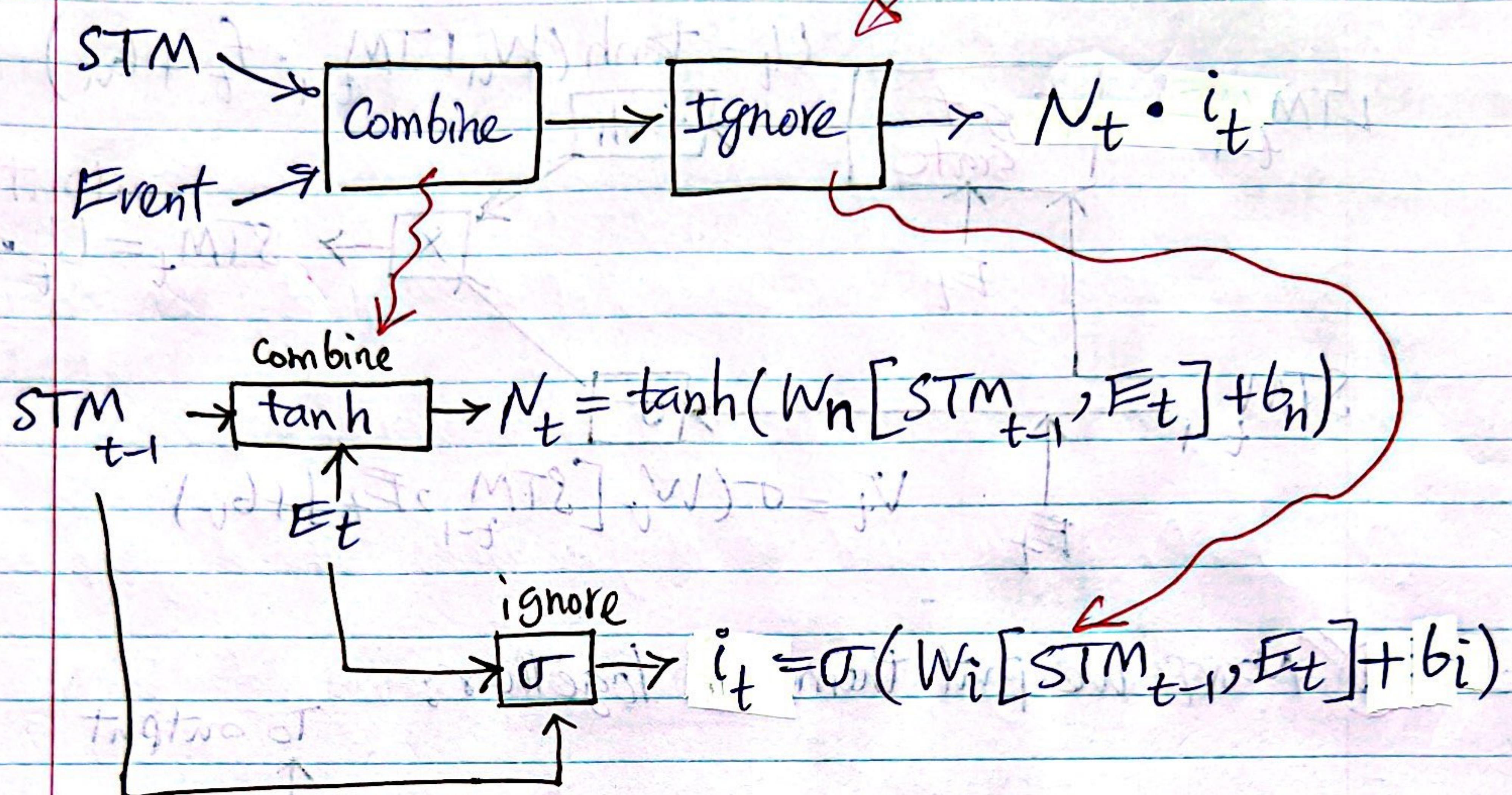
In LSTM, if we deem it is necessary, the network will remember thing from a long time ago!



RNN & revisiting:

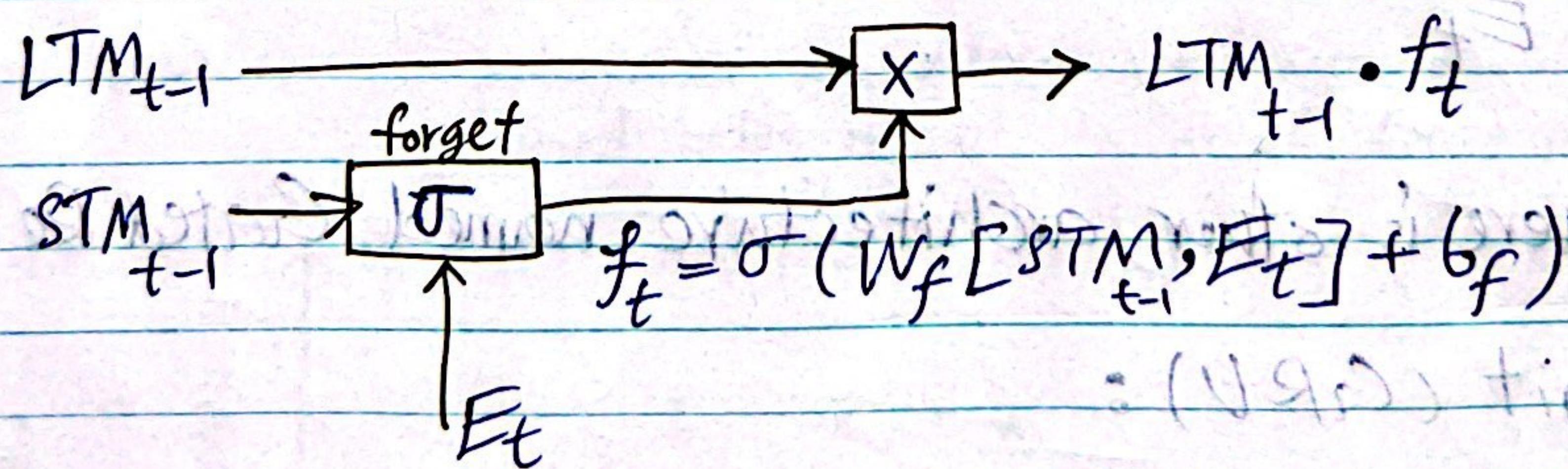


① We will start from the learn gate:

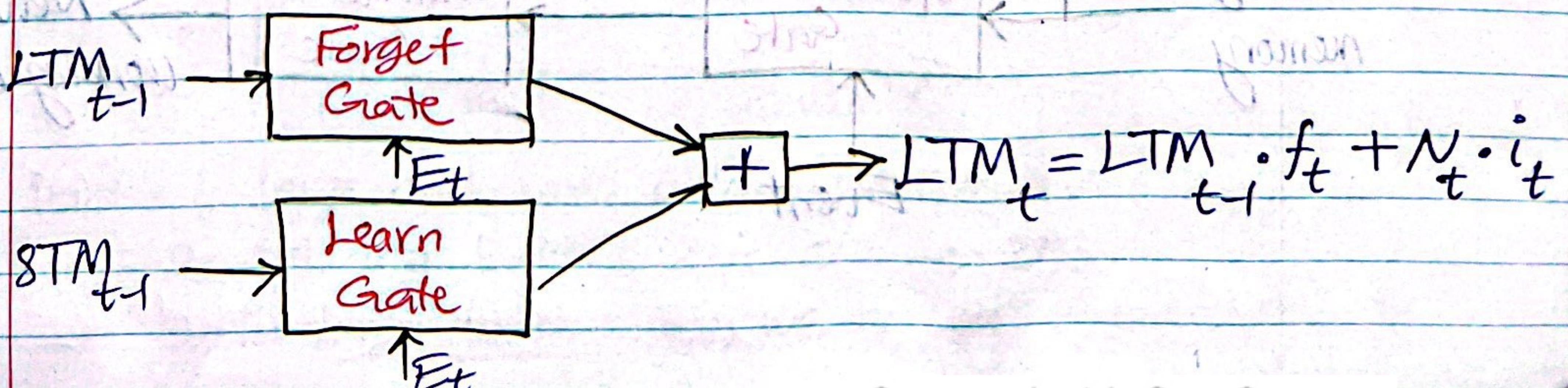


② The result of learn gate is element-wise product of N_t and i_t !

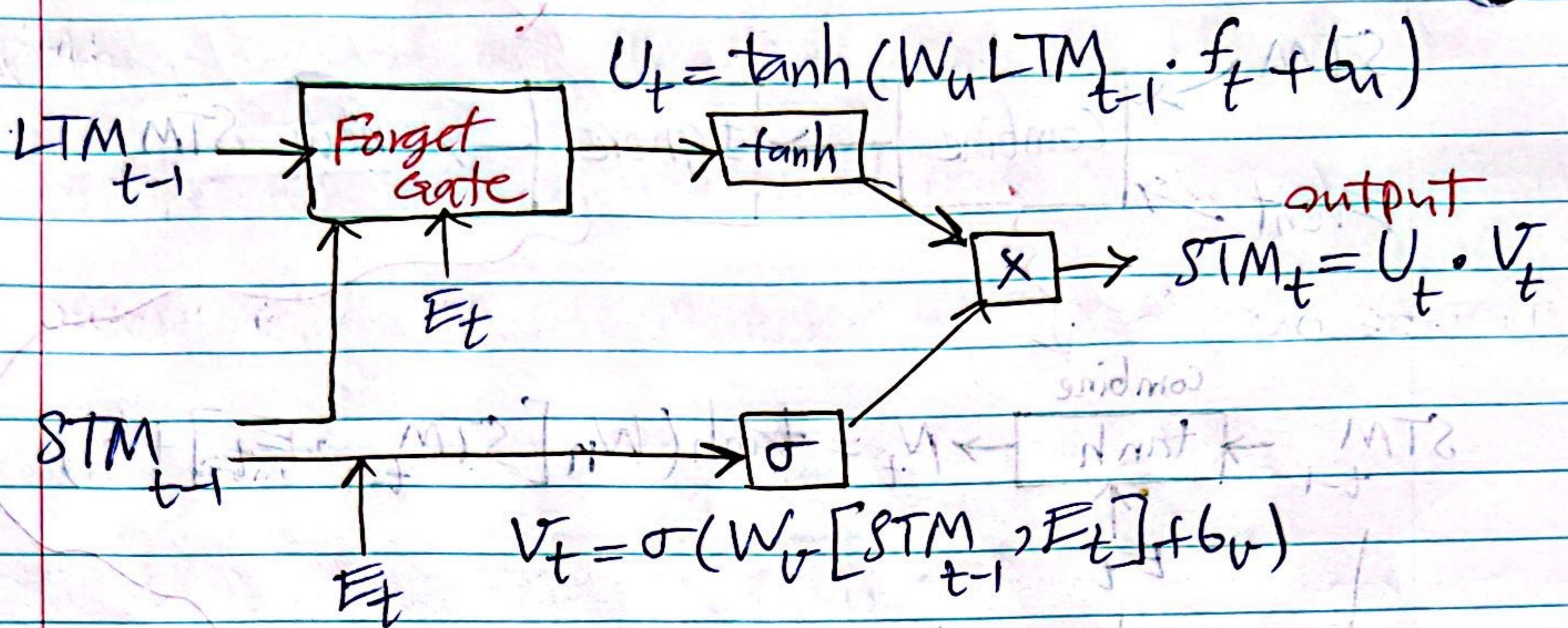
③ Now, we will explain the Forget Gate:



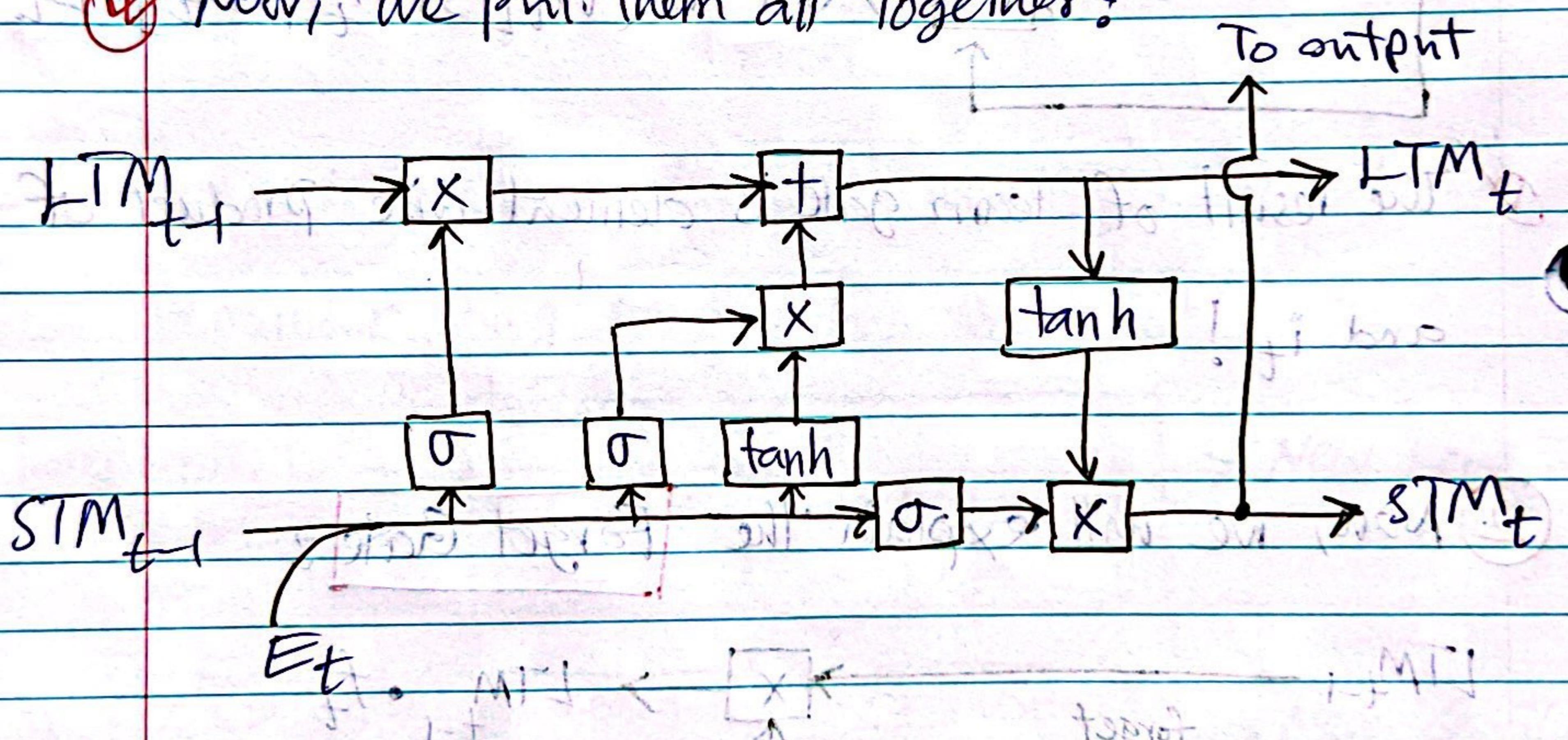
④ Remember Gate is relatively straight-forward



(4) Use Gate or Output gate:

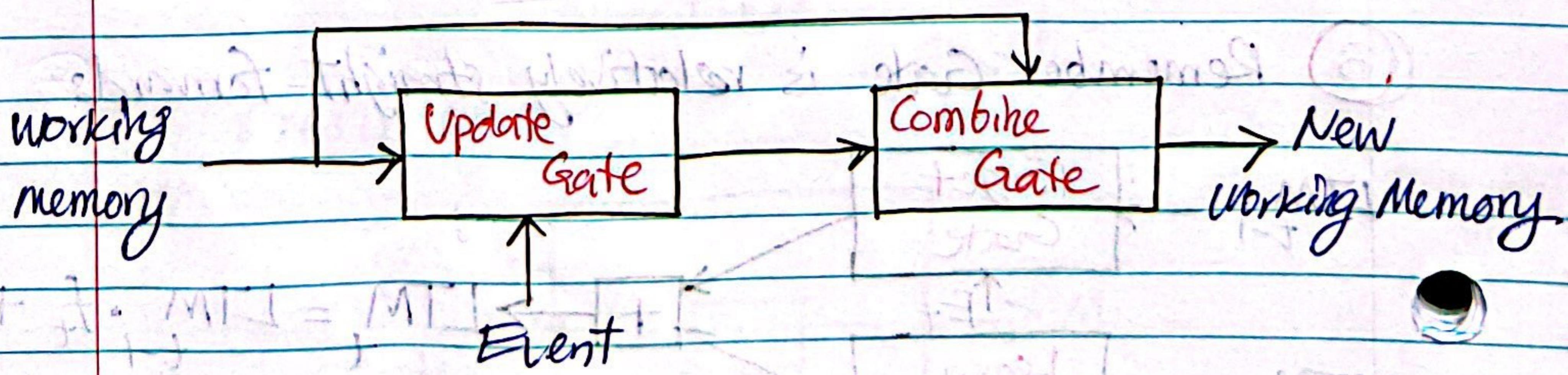


(4) Now, we put them all together:

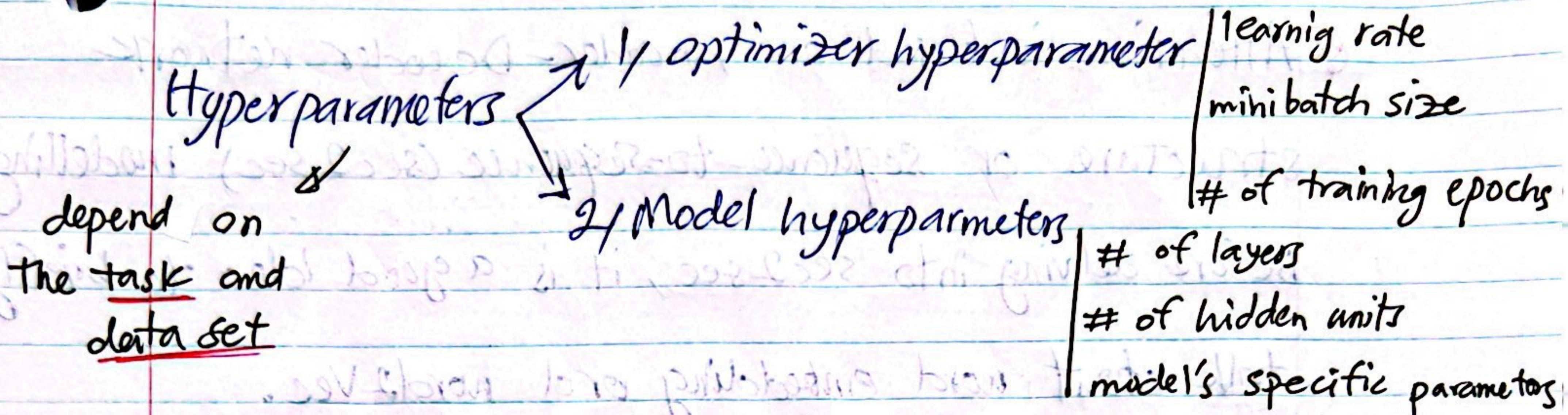


* There is other architecture named Gated Recurrent

Unit (GRU):



④ An important consideration in DL is hyperparameter tuning.



Minibatch sizes: 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, ... good whole data

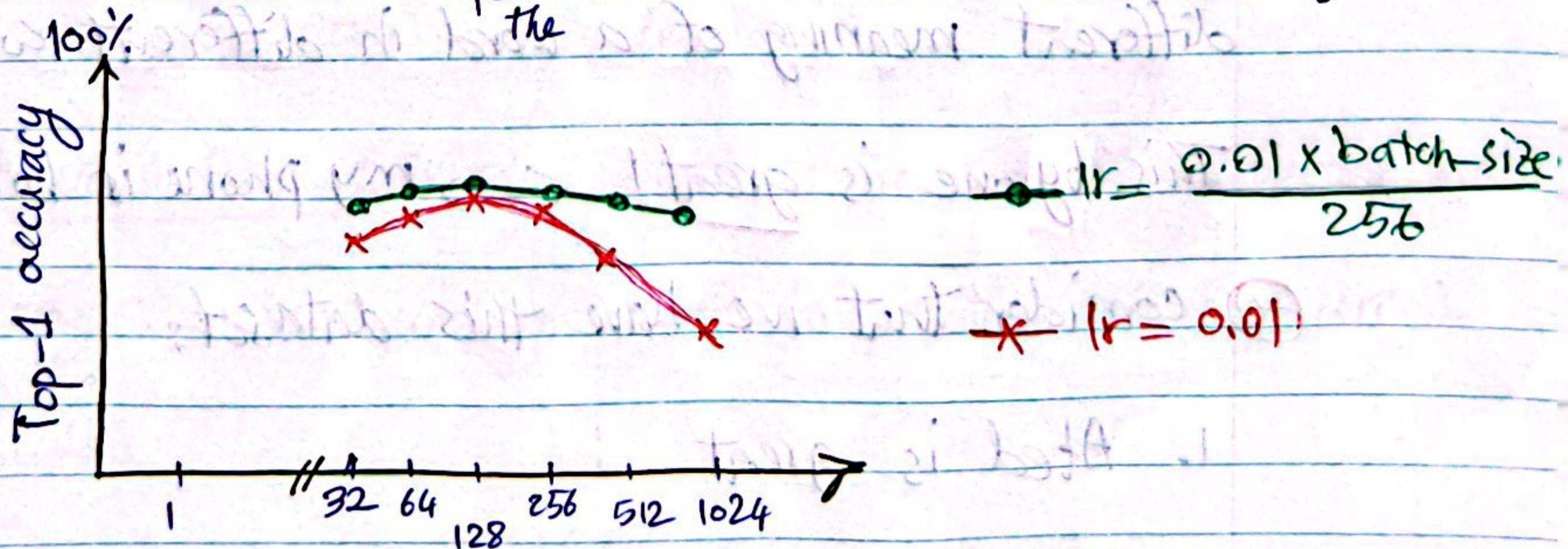
- less memory demand
- more noisy learning
- long learning process

- allows computational boost
- need large memories
⇒ more computational resources

~~the computational boost comes at a cost: decreasing model accuracy~~

This is because noise is smaller

~~batches introduce noise into the learning and helps model to evade local minimas!~~



Points 1, large and small batch sizes are bad!

2, optimum batch size is $32 \rightarrow 256$

3, for larger batch sizes, we have to increase the lr!