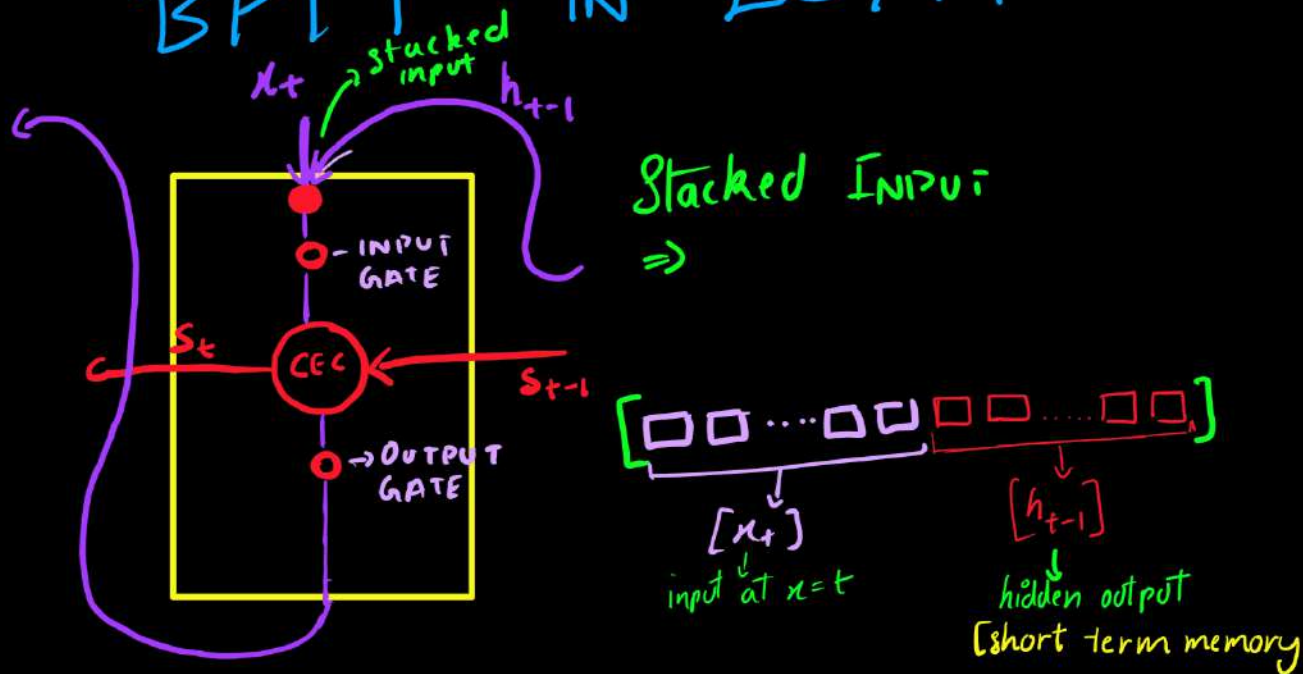
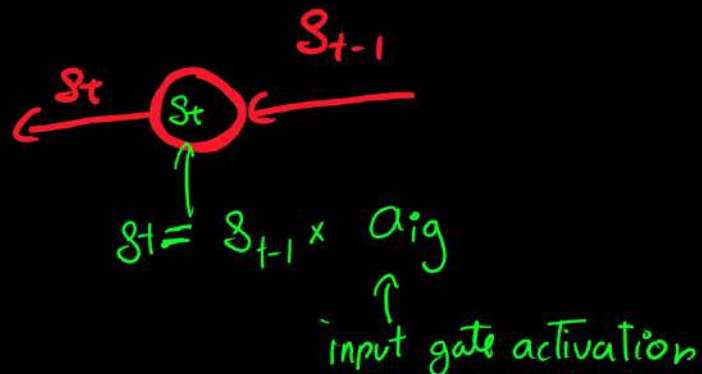


# BPTT IN LSTM

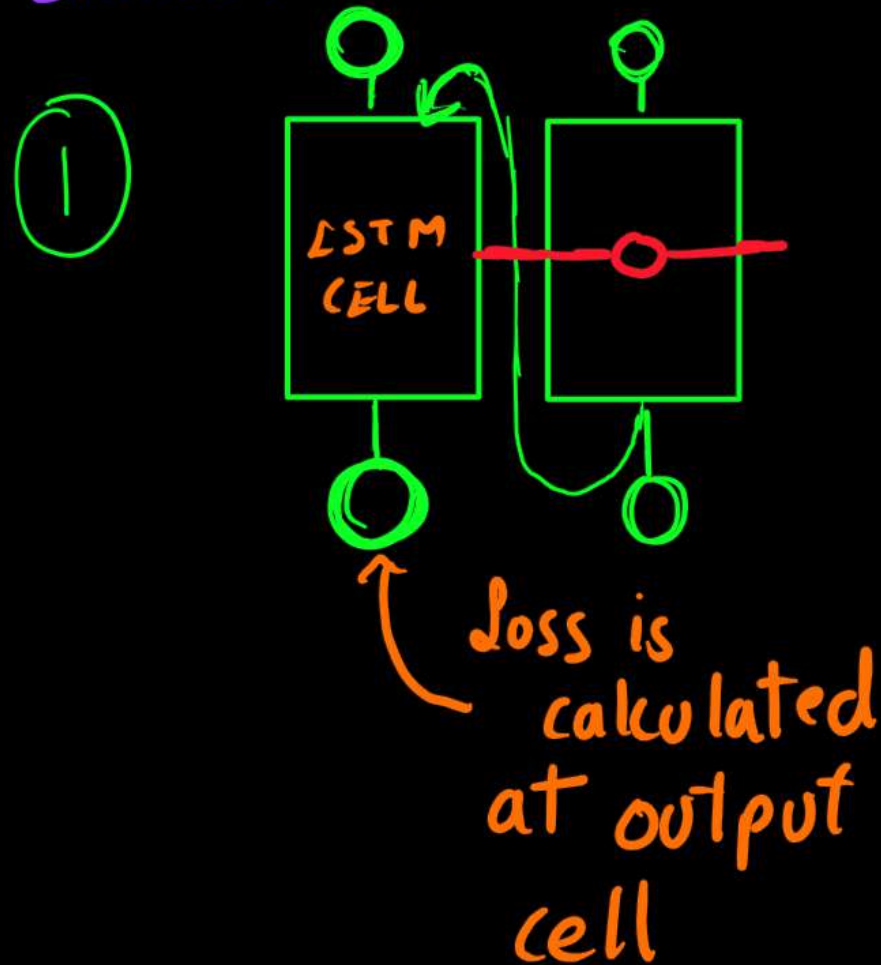


CEC Cell [long term memory]

At every Time step



# BACKPROPAGATION



\* FOR MSE

$$\text{Loss} = (\underbrace{y}_{\text{True } y} - \underbrace{\hat{y}}_{\text{predicted } y})^2$$

$$\frac{dL}{d\hat{y}} = 2(y - \hat{y})$$

Now we find out how the loss affects the activation and the weights of the output cell

$$\textcircled{2} \quad \hat{y} = a(w_y h_x + b_x)$$

$\downarrow$  activation function       $\downarrow$  output weights       $\downarrow$  hidden output at layer  $x$        $\downarrow$  bias at  $x$

NOTE:- All the derivatives are partial derivative here

We need to find  $\frac{dL}{dw_y}$  Slope of the loss function wrt  $w_y$

$$\frac{dL}{dw_y} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_y}$$

$$\frac{d\hat{y}}{dw_y} = \text{activation derivative} * h_x \quad \left[ \begin{array}{l} \text{this is the gradient} \\ \text{change at time step } t \end{array} \right]$$

$$\frac{d\hat{y}}{db_y} = \text{activation derivative} * 1$$

$\downarrow$  Bias gradient

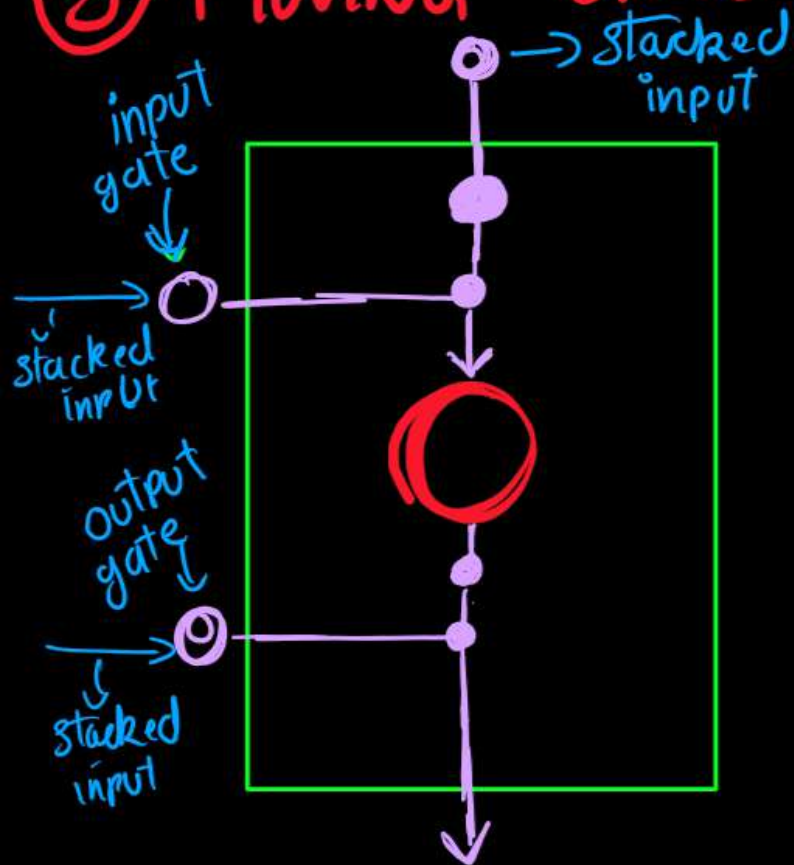
★ A Gradient accumulator sums up these gradients (changes) over all the time steps

We also find  $\frac{d\hat{y}}{dh_x}$   $\left[ \begin{array}{l} \text{Slope of the loss function} \\ \text{wrt activation of the hidden layer} \end{array} \right]$

$$\frac{d\hat{y}}{dh_x} = \text{activation derivative} * w_y$$



## ③ MOVING ONTO



★ Taking the stacked input as  $i_s_t$  at time  $t$

## THE LSTM MEMORY CELL

Right now we have

- Updated the output weights & bias to the gradient accumulator
- We have the gradient of how loss function affects the output of the hidden layer
- In an LSTM cell we have to update the weights and biases of 3-thing

- (i) Input weights & biases
- (ii) Input gate weights & biases
- (iii) Output gate weights & biases

Backpropagation starts from the output gate

$$a_{og} = \text{activation}(i s_t \cdot w_{og} + b_{og})$$

activation  
of output gate

stacked  
input

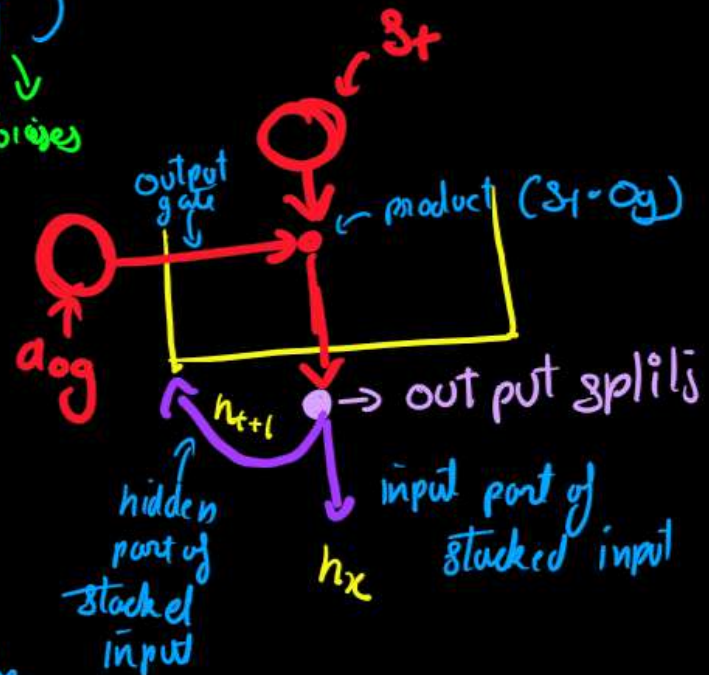
weights  
of  
output  
gate

biases

hidden part  
which goes to next  
time step

$$[h_{next}, h_x] = \text{activation}(a_{og} \cdot s_t)$$

output  
of hidden  
layer



For now let's assume we get the gradient of the hidden part of the stacked output from the next time step

We also have  $dh_x$  from the output





Now that we have the gradient of stacked output we need

$$\frac{dL}{dw_{og}} = \underbrace{\frac{dL}{do_s}}_{\text{gradient stacked output}} \cdot \underbrace{\frac{do_s}{da_{og}}}_{\text{chain rule}} \cdot \underbrace{\frac{da_{og}}{dw_{og}}}_{\text{chain rule}}$$

this will get us the gradient to update output gate

gradient stacked output

chain rule  
Till we get  $dw_{og}$

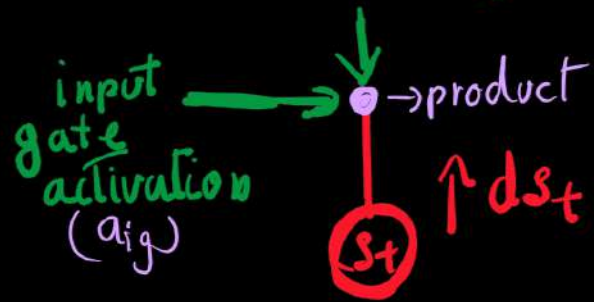
So calculate  $\frac{dL}{ds_t} = \frac{dL}{do_s} \cdot \frac{do_s}{ds_t}$

\* I don't want to calculate everything again so I will just write them in derivative form



done

Now we need to update the input gate, and cell input weights



$$a_{ig} = \text{act}(i_s \cdot w_{ig} + b)$$
$$a_i = \text{act}(i_s \cdot w_i + b)$$

$$s_t = s_{t-1} * \text{product}$$

$$s_t = s_{t-1} * a_{ig} * a_i$$

same  $\frac{dL}{ds_t}$  is used to calculate.

$$\frac{dL}{dw_{ig}} = \frac{dL}{ds_t} \cdot \frac{ds_t}{da_{ig}} \cdot \frac{da_{ig}}{dw_{ig}}$$

$$\frac{dL}{dw_i} = \frac{dL}{ds_t} \cdot \frac{ds_t}{da_i} \cdot \frac{da_i}{dw_i}$$

calculate this chain rule

All these calculated gradient are added to their respective gradient accumulators

$\therefore$  Gradients from all the three will constitute the gradient of the input vector

$$\frac{dL}{dis} \text{ (from output gate)} = \frac{dL}{da_{og}} \cdot \frac{da_{og}}{dis} = w_{og} + \text{activation derivative (output gate)} \quad \text{--- (1)}$$

$$\frac{dL}{dis} \text{ (from input gate)} = \frac{dL}{da_{ig}} \cdot \frac{da_{ig}}{dis} = w_{ig} + \text{activation derivative (input gate)} \quad \text{--- (2)}$$

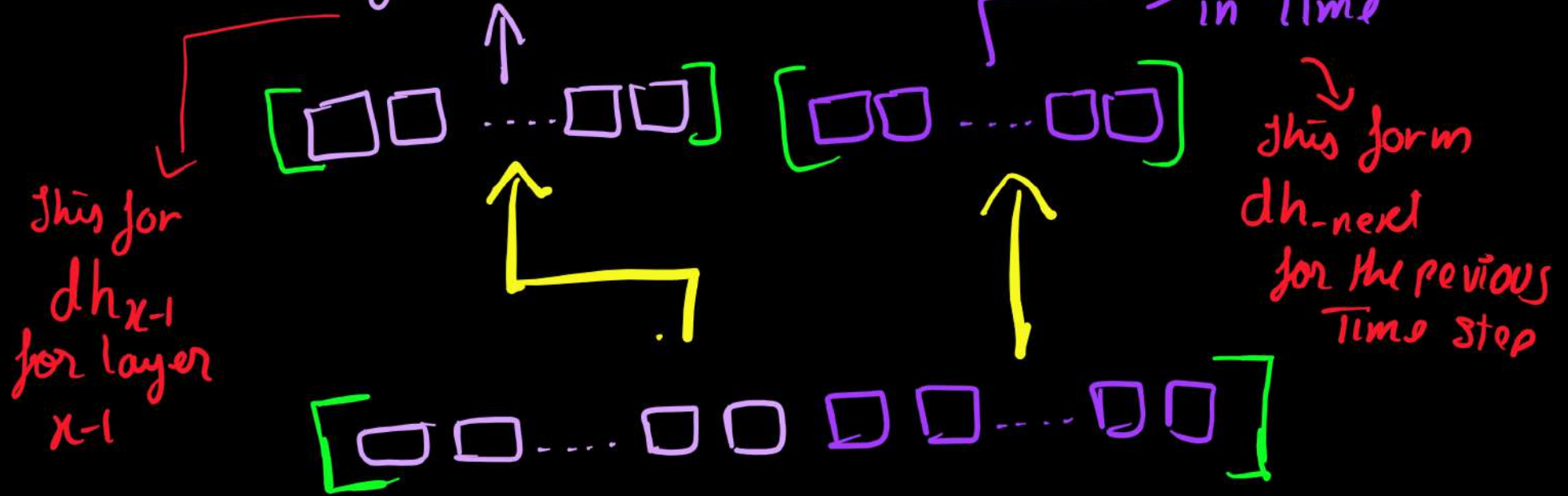
$$\frac{dL}{dis} \text{ (from cell input)} = \frac{dL}{da_i} \cdot \frac{da_i}{dis} = w_i + \text{activation derivative (cell input)} \quad \text{--- (3)}$$

$$\frac{dL}{dis} = \text{(1)} + \text{(2)} + \text{(3)}$$



FINALLY  
split into

this stacked input gradient is then



---

With This BPTT through an LSTM is  
complete