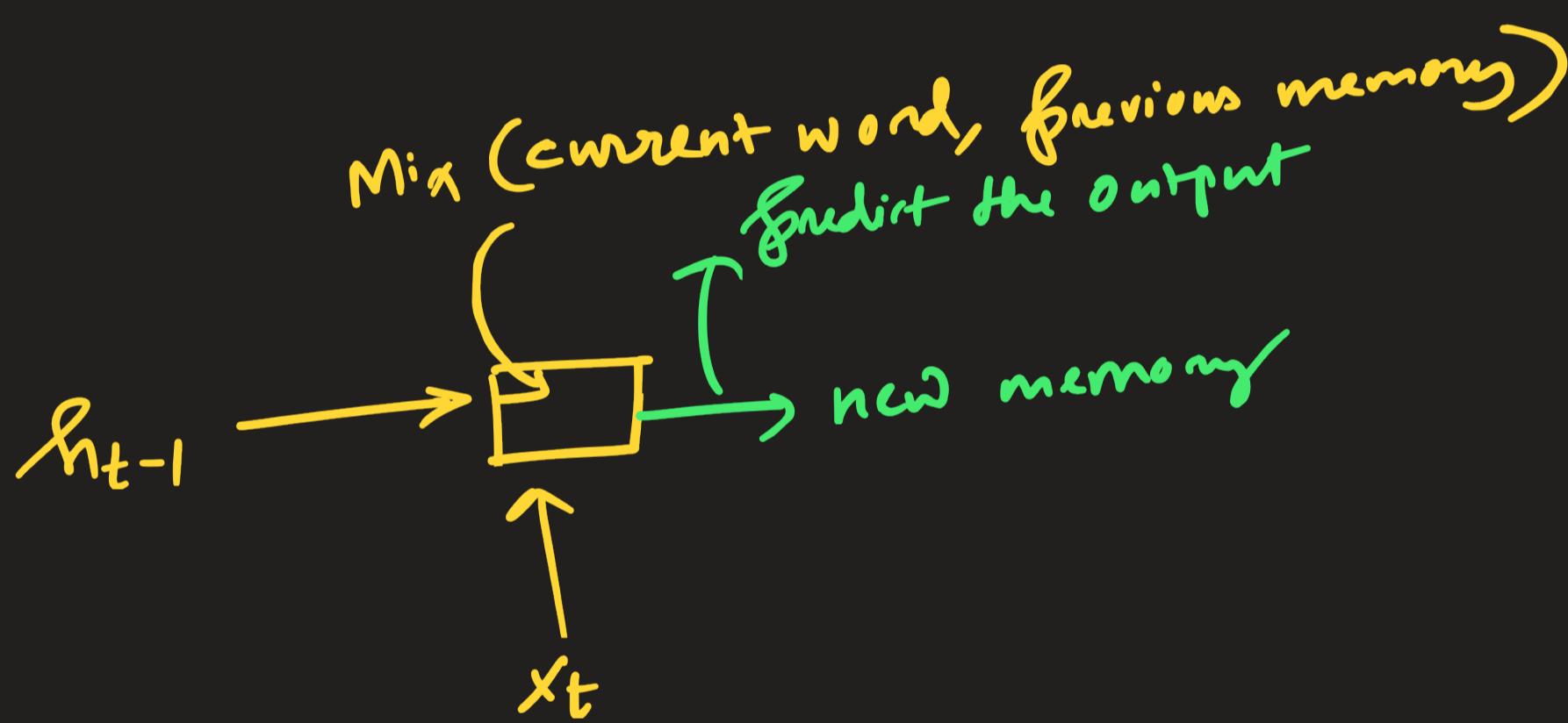
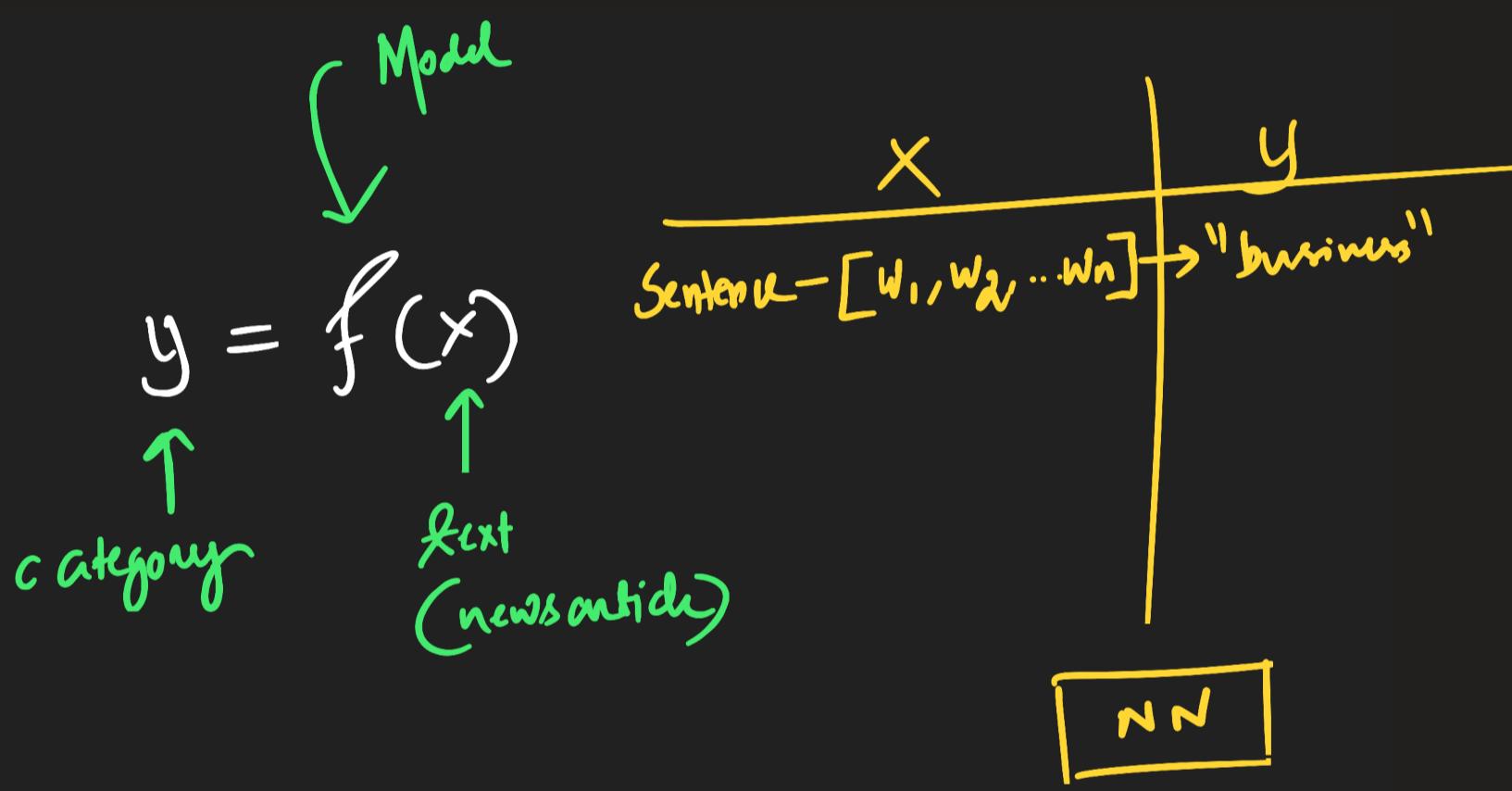


# *RNNs*

*Satya Pattnaik*

RNN - thinks one word at a time



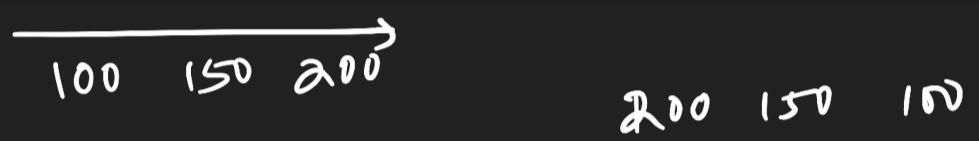


## *Why MCP is not meant for this problem?*

Word → [Emb] [ ]

S1: dog bites man      |      |      |      }  
S2: man bites dog      |      |      |      }

dog  
bites  
man



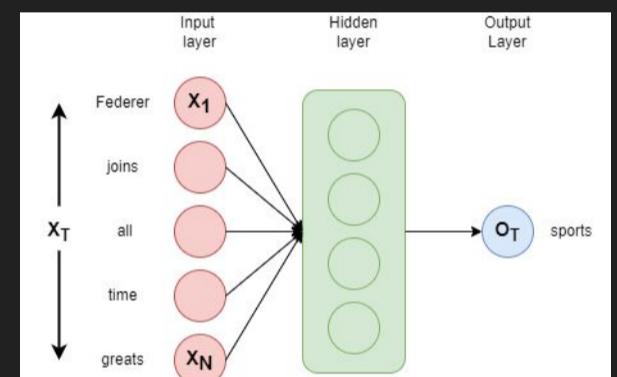
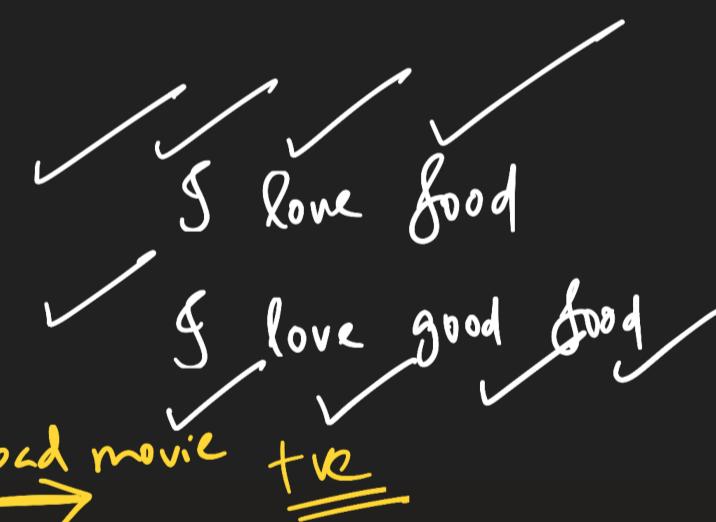
What if we use MLP?

All inputs at once

not → -ve

I am not a bad singer

This is not a bad movie



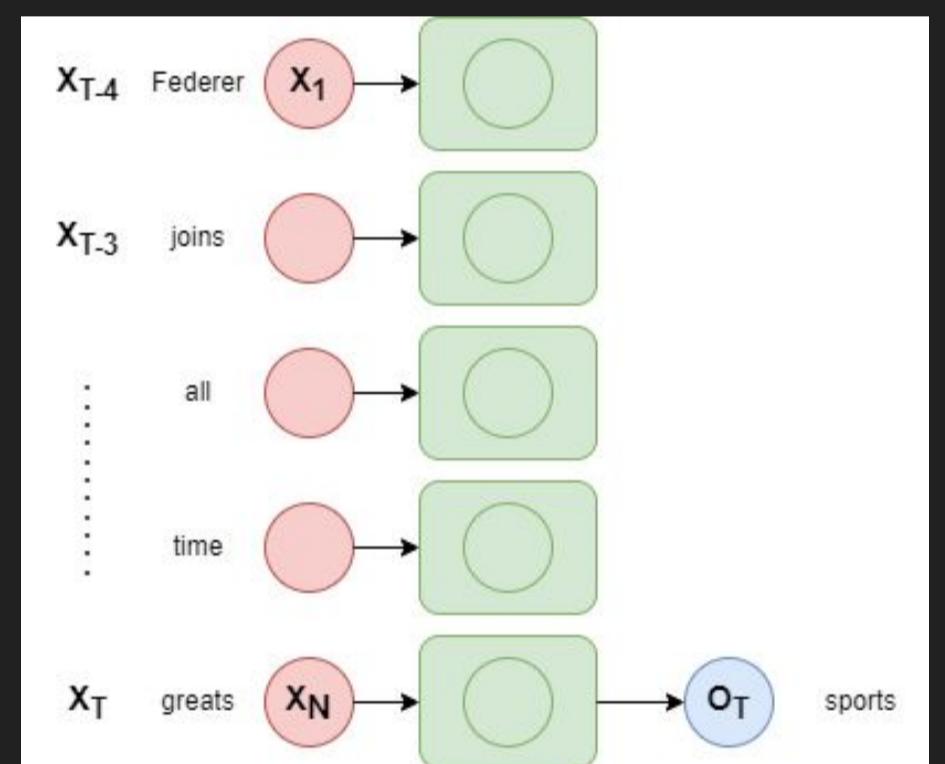
100 200 300 → +vecT

What is missing in this Architecture?

- no connection with previous input words
- each word is isolated from each other
- varying lengths of sentences also will have issues

200  
300  
100

*Can we modify this architecture?*

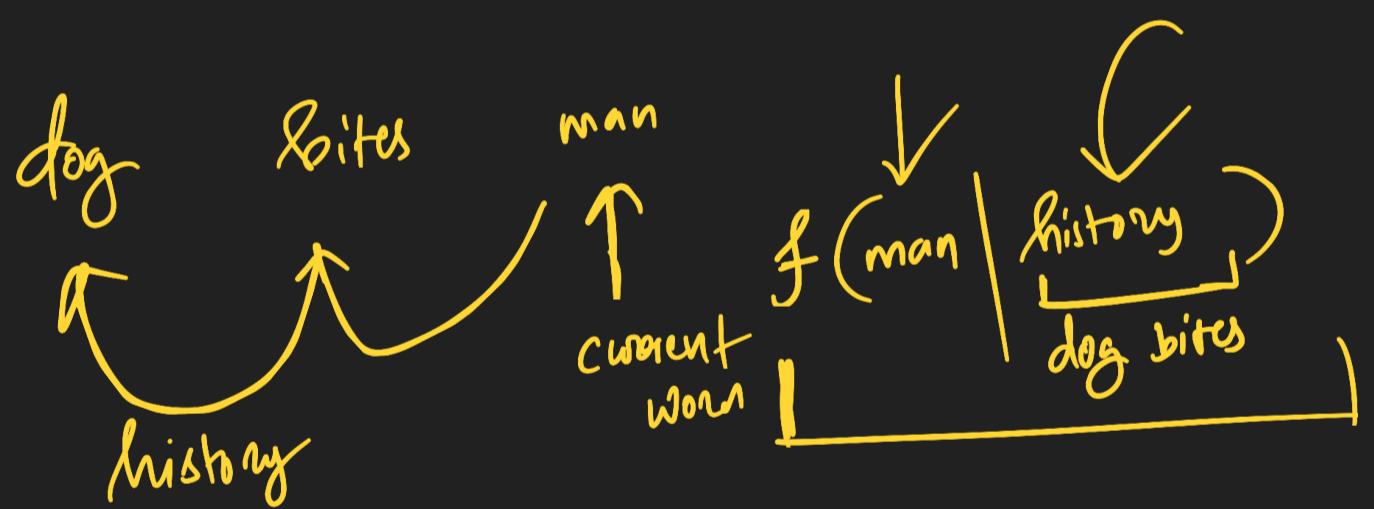


*What is missing?*

Summarize the issues

dog bites man  
man bites dog

$$NN_{ot} = f(w_0)$$



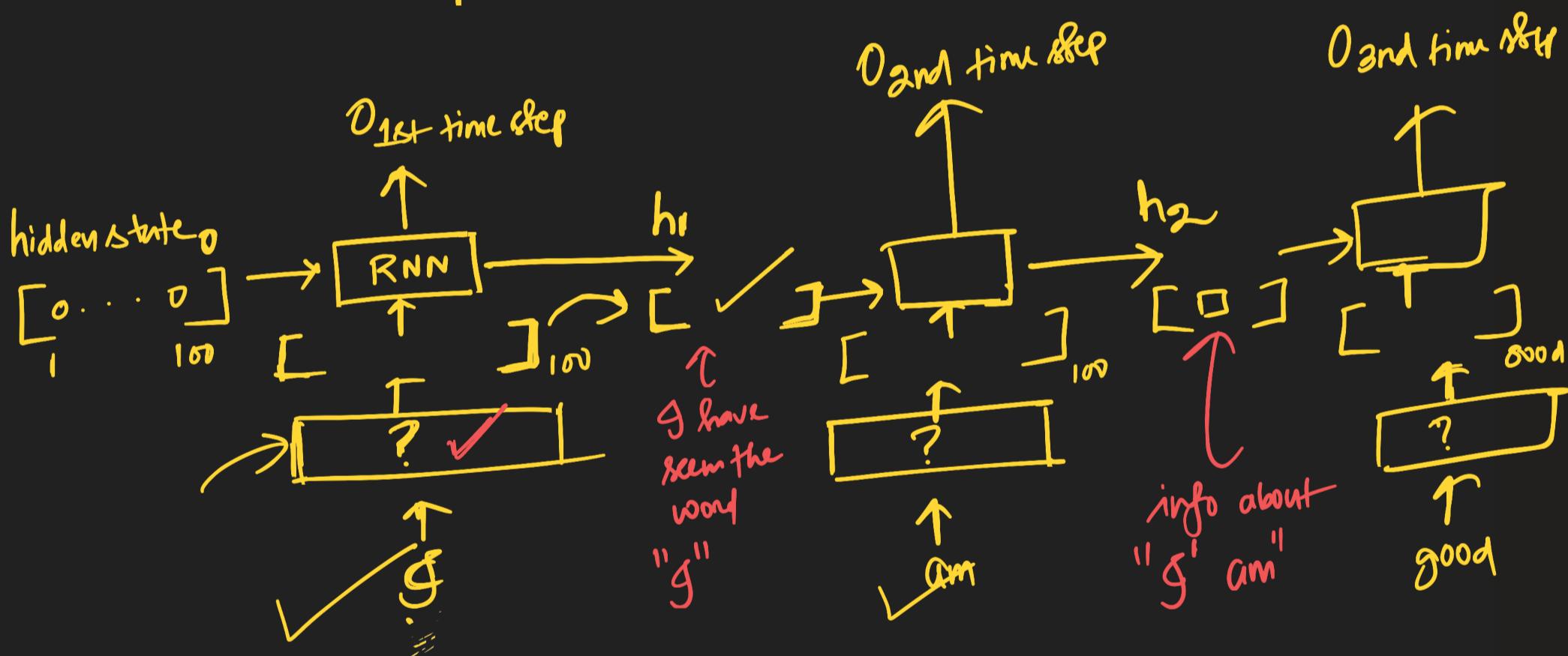
RNNs -> Intro

g  
↑

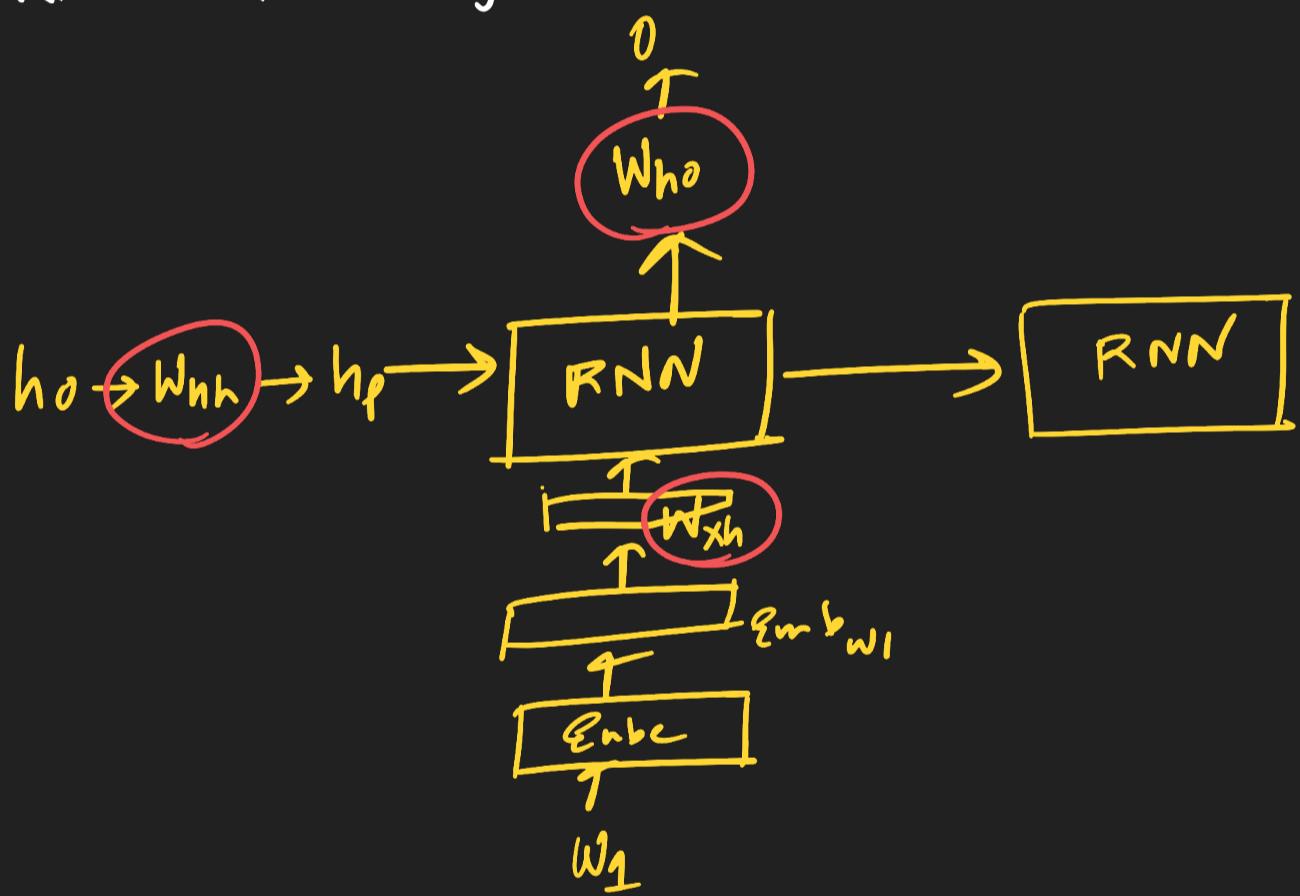
am

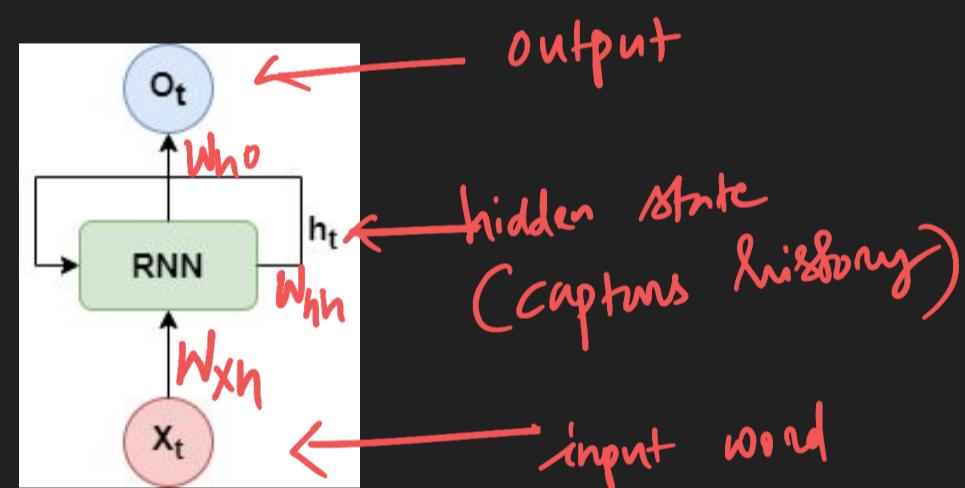
good

10 → 20 ↑  
 $\stackrel{\text{Stock}}{\text{↑}}$   
 $\stackrel{\text{↑ up}}{\text{↑}}$

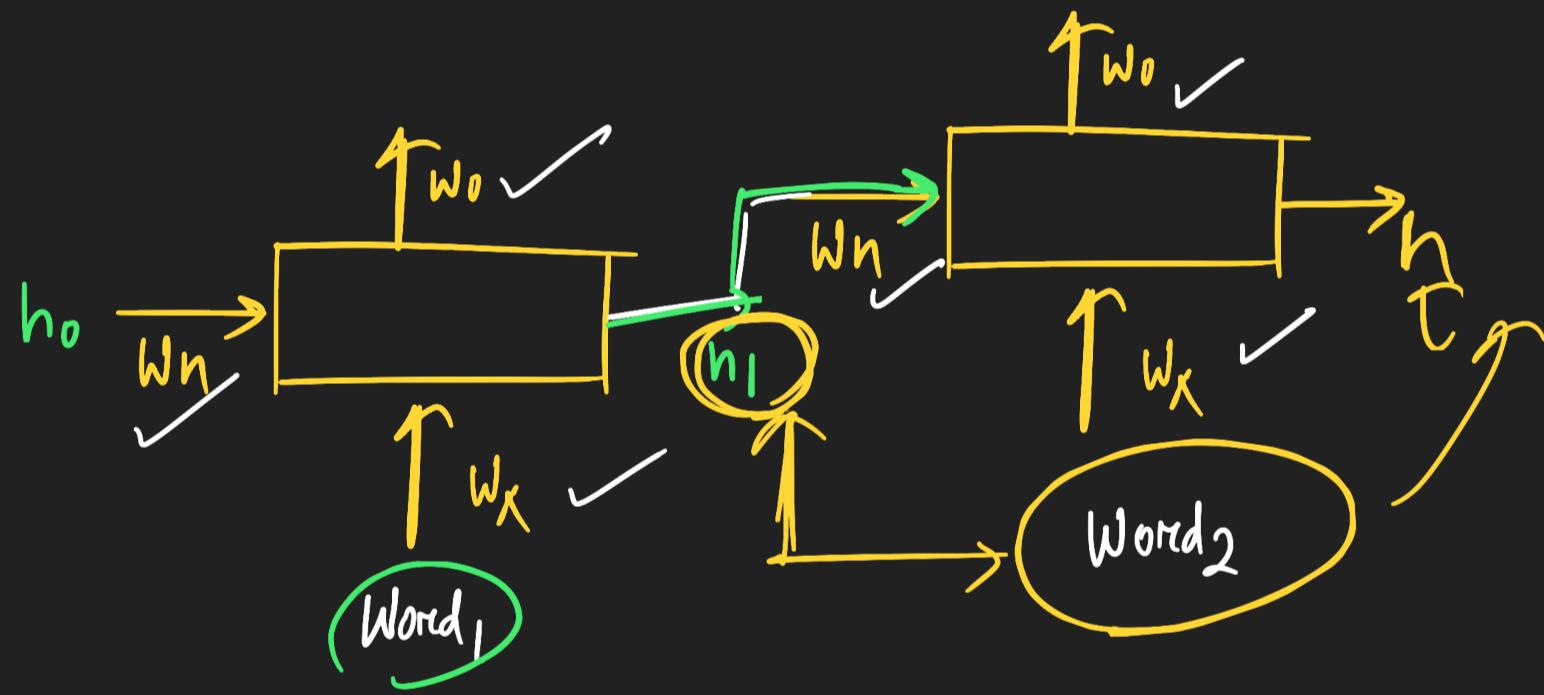


RNNs -> Advantages

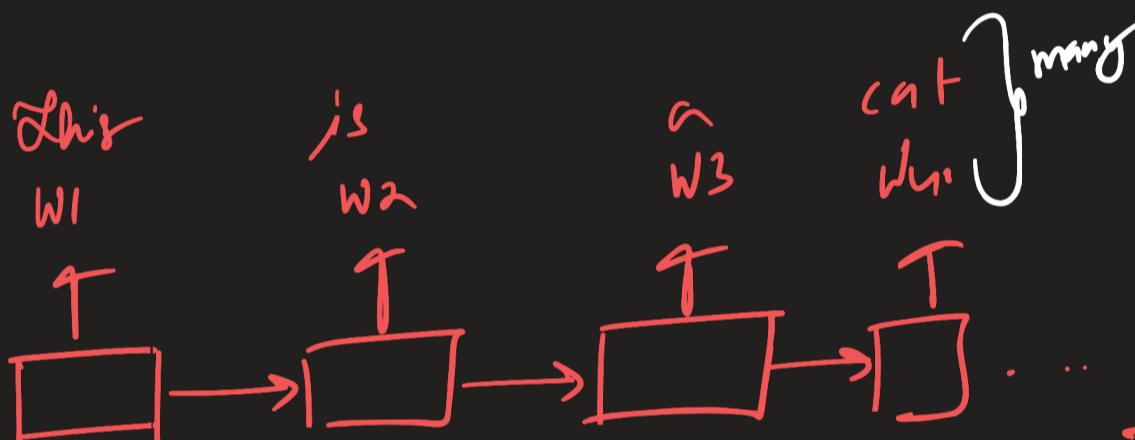




How is it unravelled over time?



## Image captioning



Rep

CNN  
image

Embedding of  
picture of  
a cat

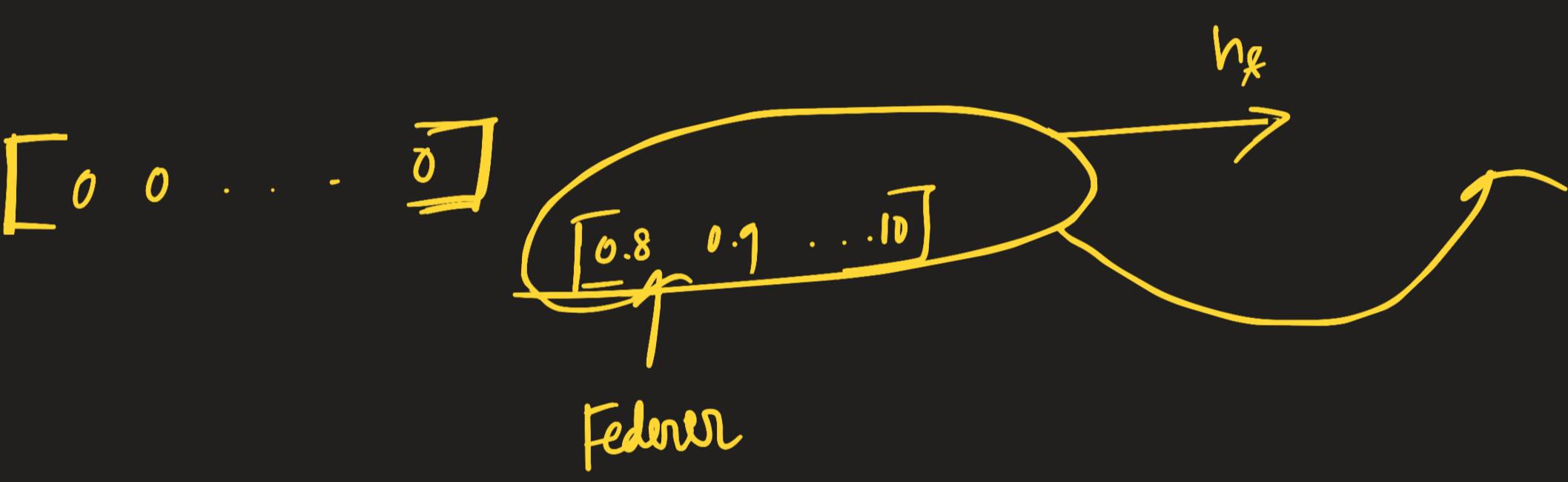
Sentiment classification

op<sup>tre</sup> } one  
T  
This song is good } many  
T

more

~~Song is a good one~~

$h_t [ ]$  tve

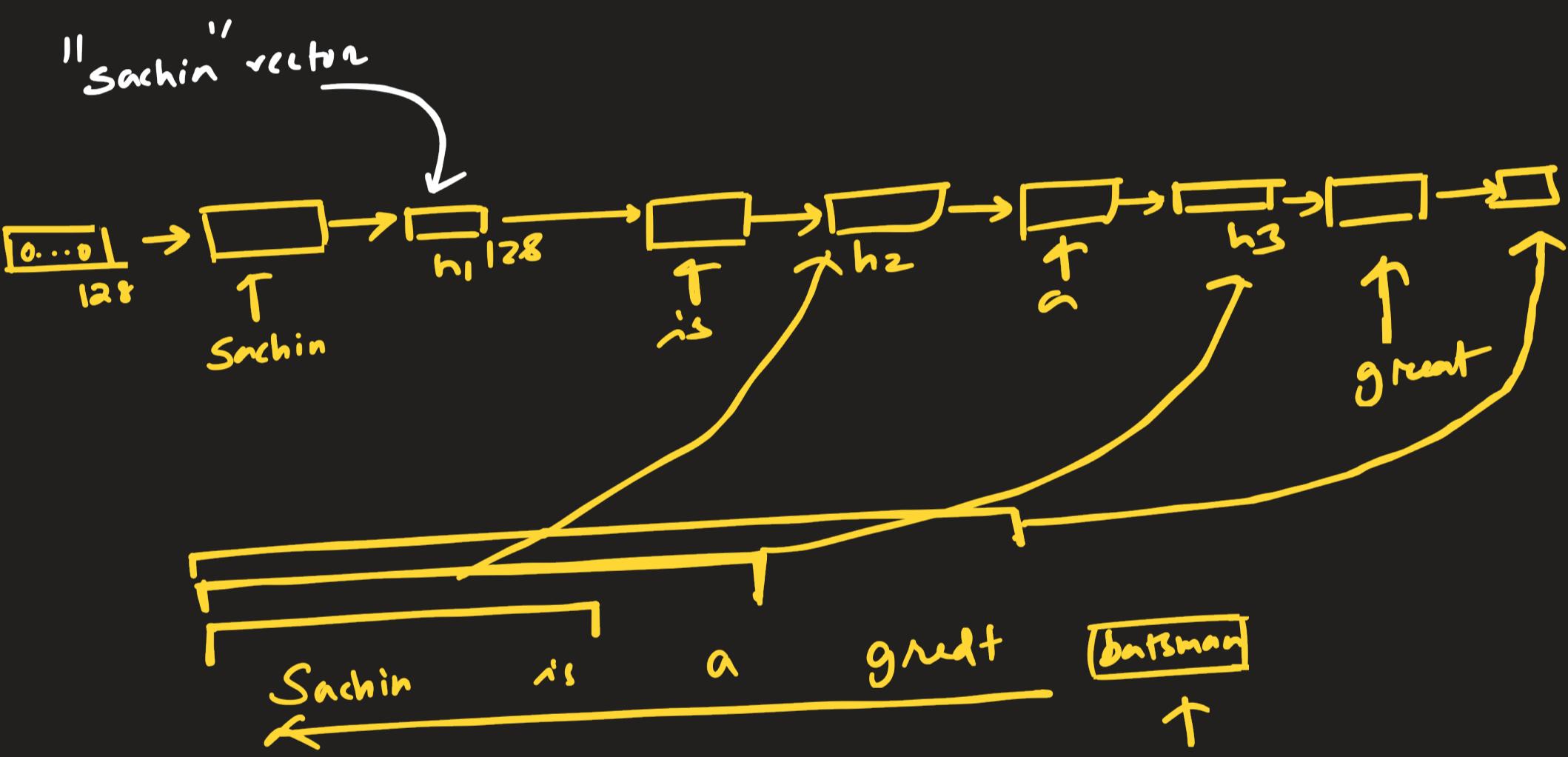


$$w_1 \quad w_2 \quad \rightarrow h_2 = f(e_{mb_1}, e_{mb_2})$$

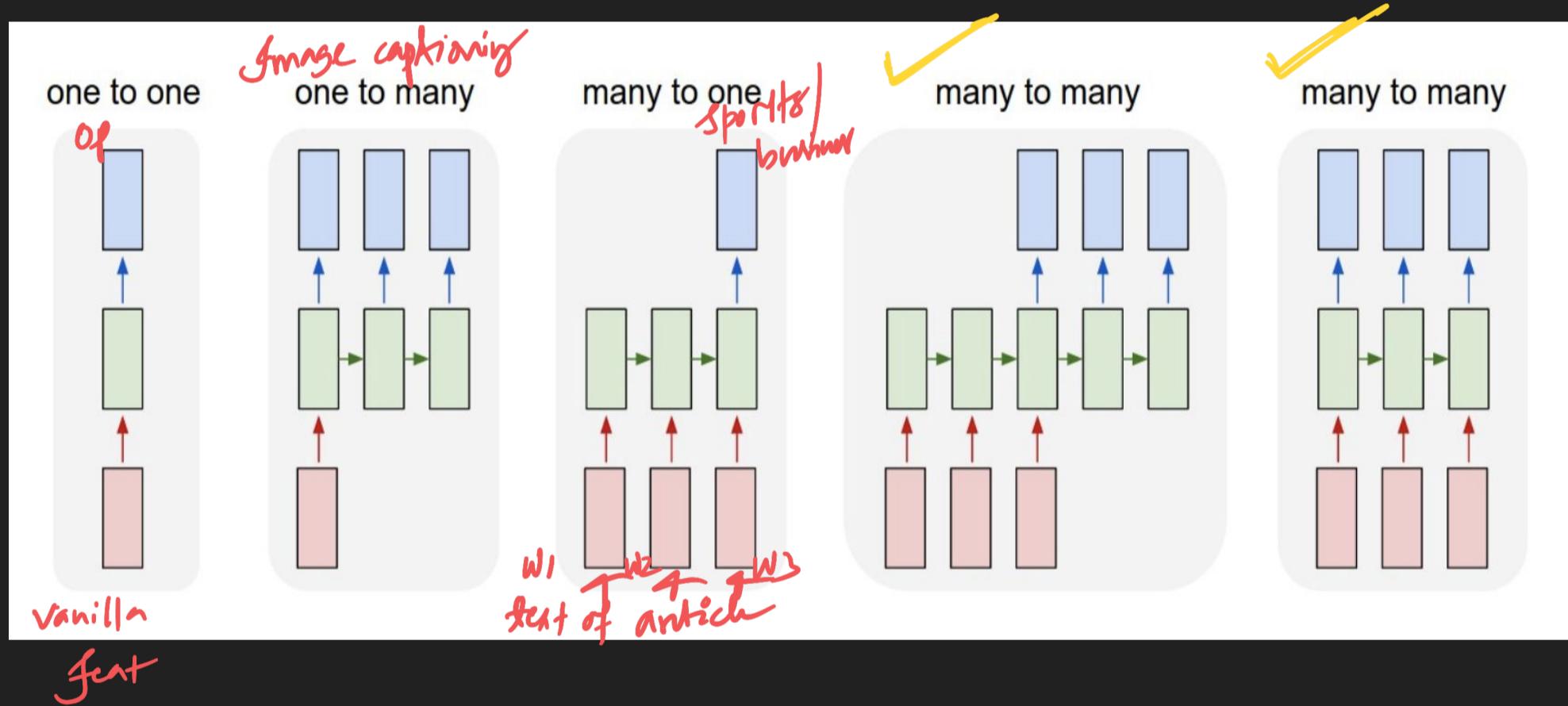
great

301

great game



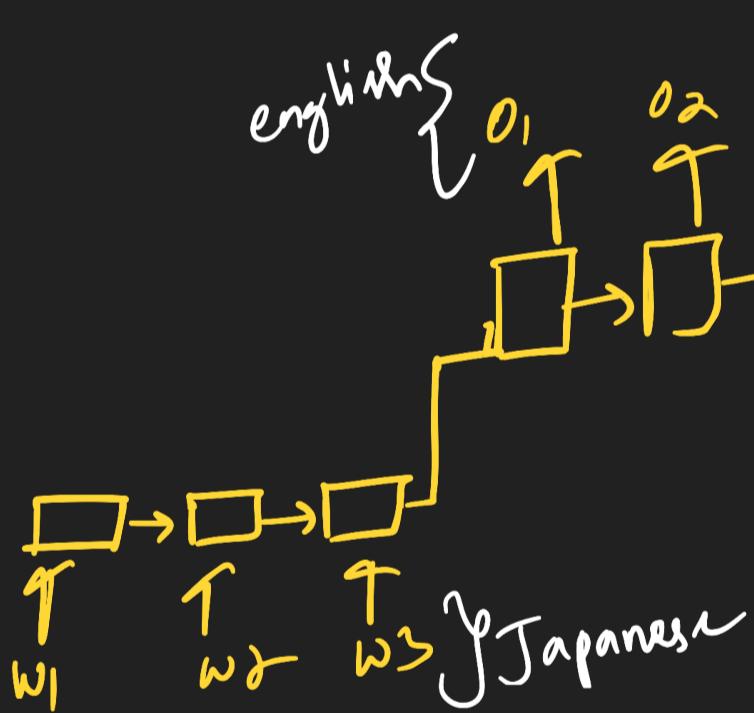
## Types of RNN



## Types of RNN

Many to Many

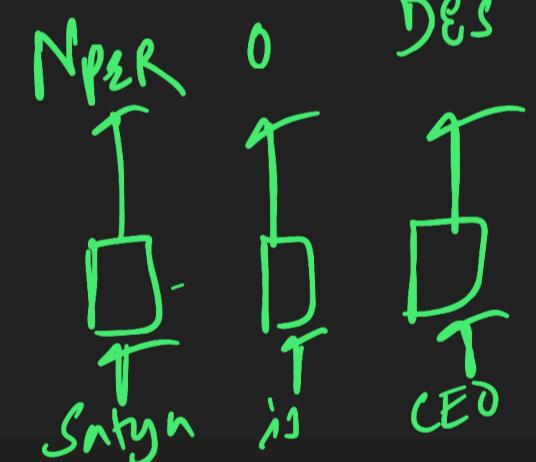
no. of ip  $\neq$  no. of op



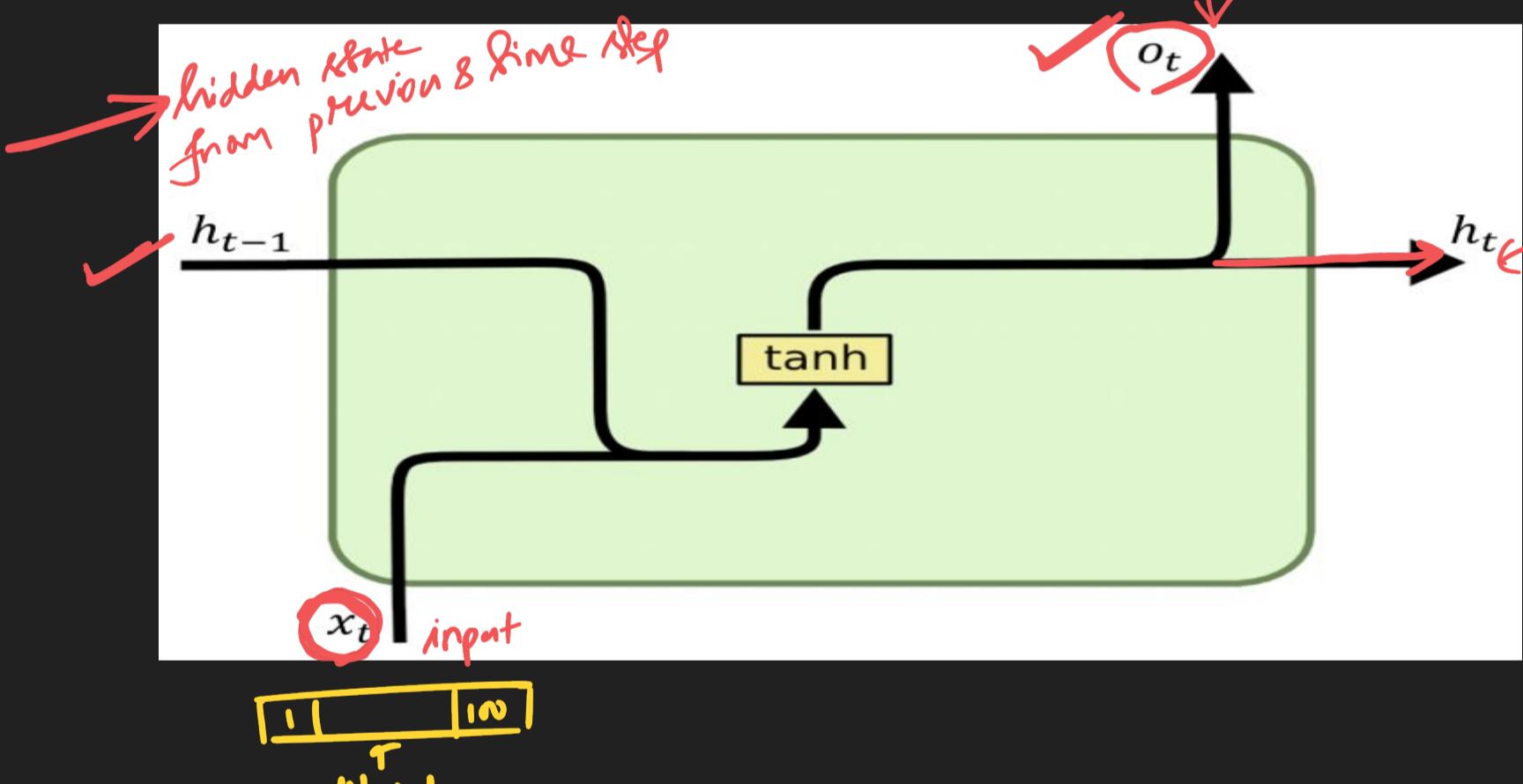
Many to Many

no. of ip = no. of op

NOON                          NOUN



Different Components of RNN



$$\text{output} = Z(h_t)$$

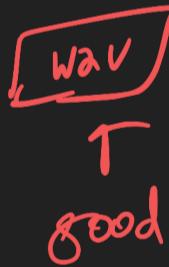
$$h_t = F(x_t, h_{t-1})$$

But do we pass the words as it is to RNN cell?



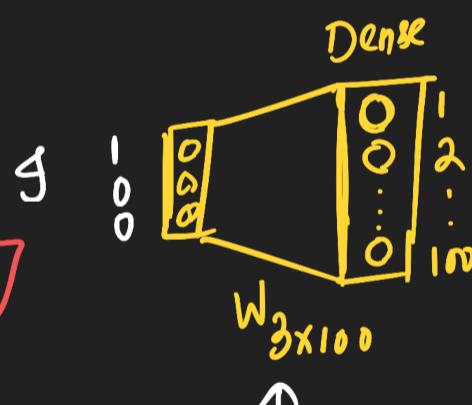
But do we pass the words as it is to RNN cell?

- 1] Pretrained embedding
- 2] Trainable embedding layer



vocab

$$\begin{array}{rcl} \text{I} & \rightarrow & 1 \\ \text{am} & \rightarrow & 2 \\ \text{good} & \rightarrow & 3 \end{array}$$

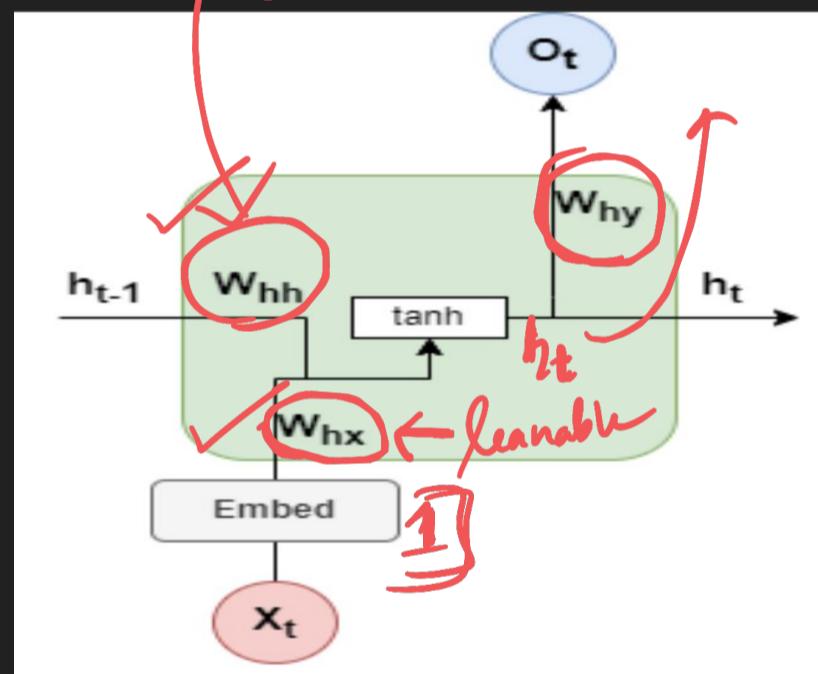


$$\begin{matrix} \text{I} & 1 & \text{am} & 0 & \text{good} & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \end{matrix}$$

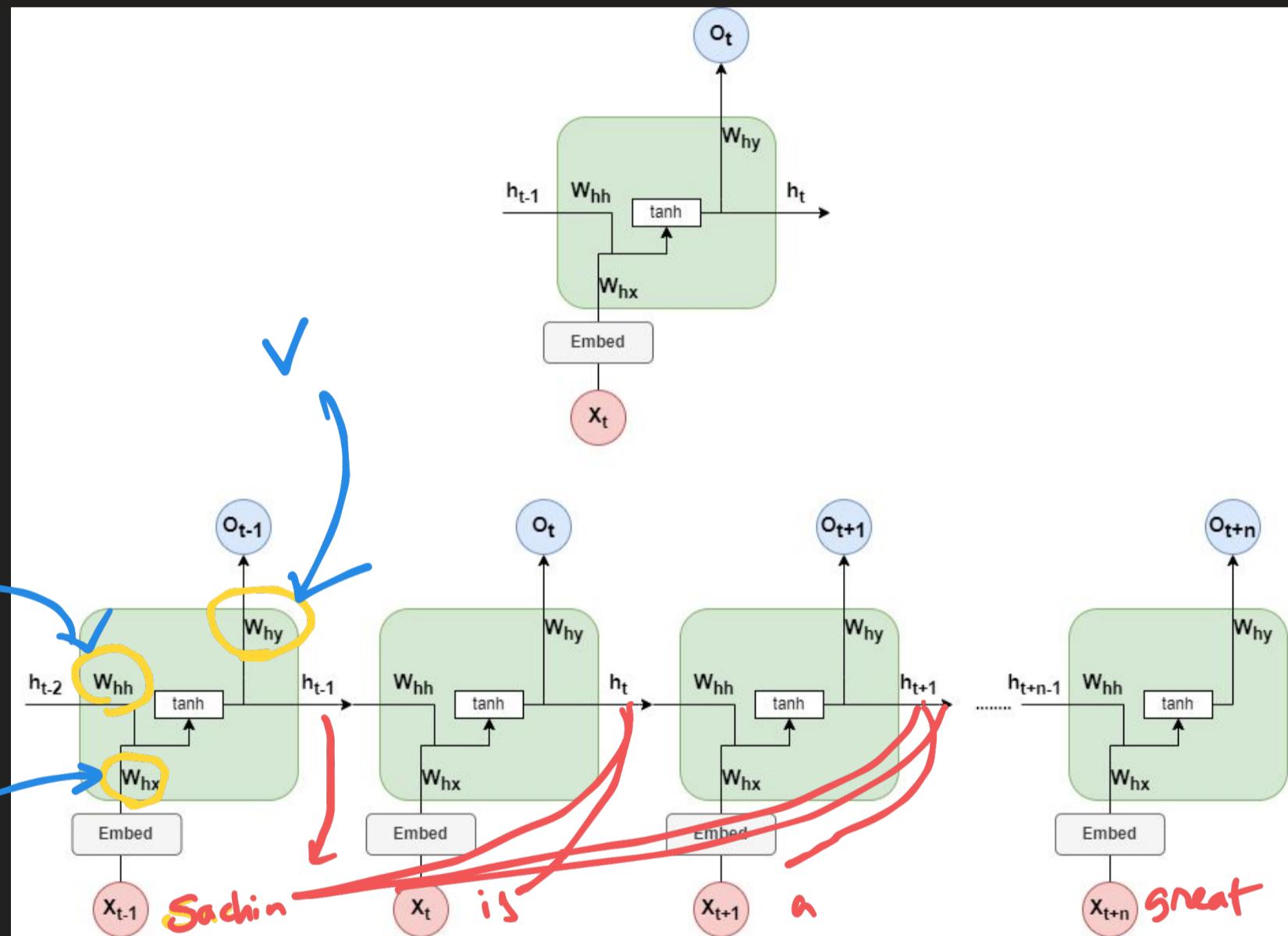
↑ Weight is trainable

What is being learnt?

hidden state  
over timestep  $\rightarrow [$

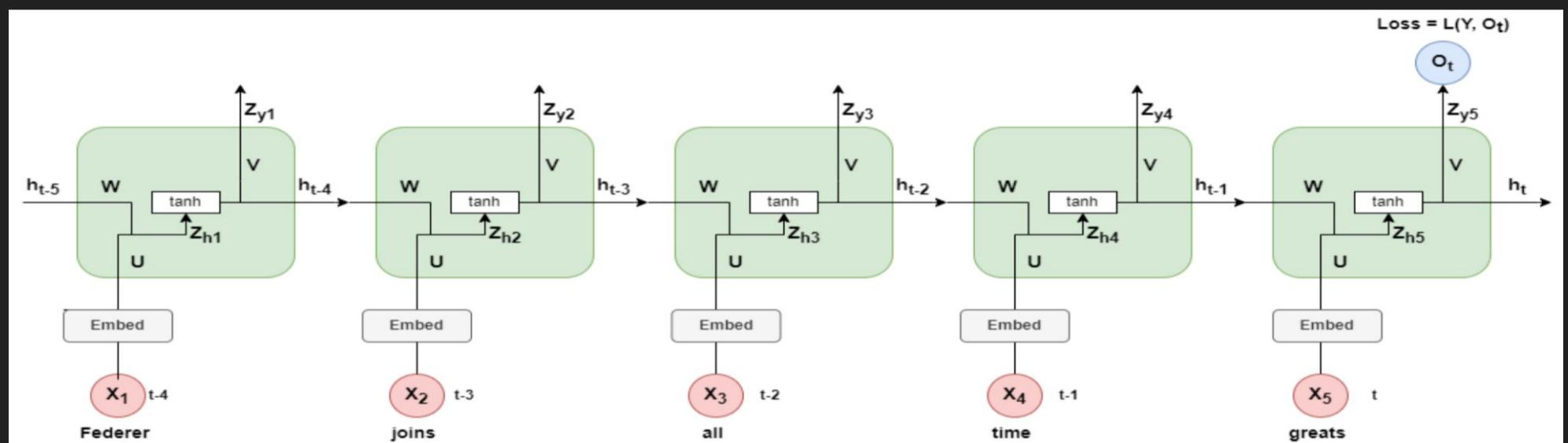


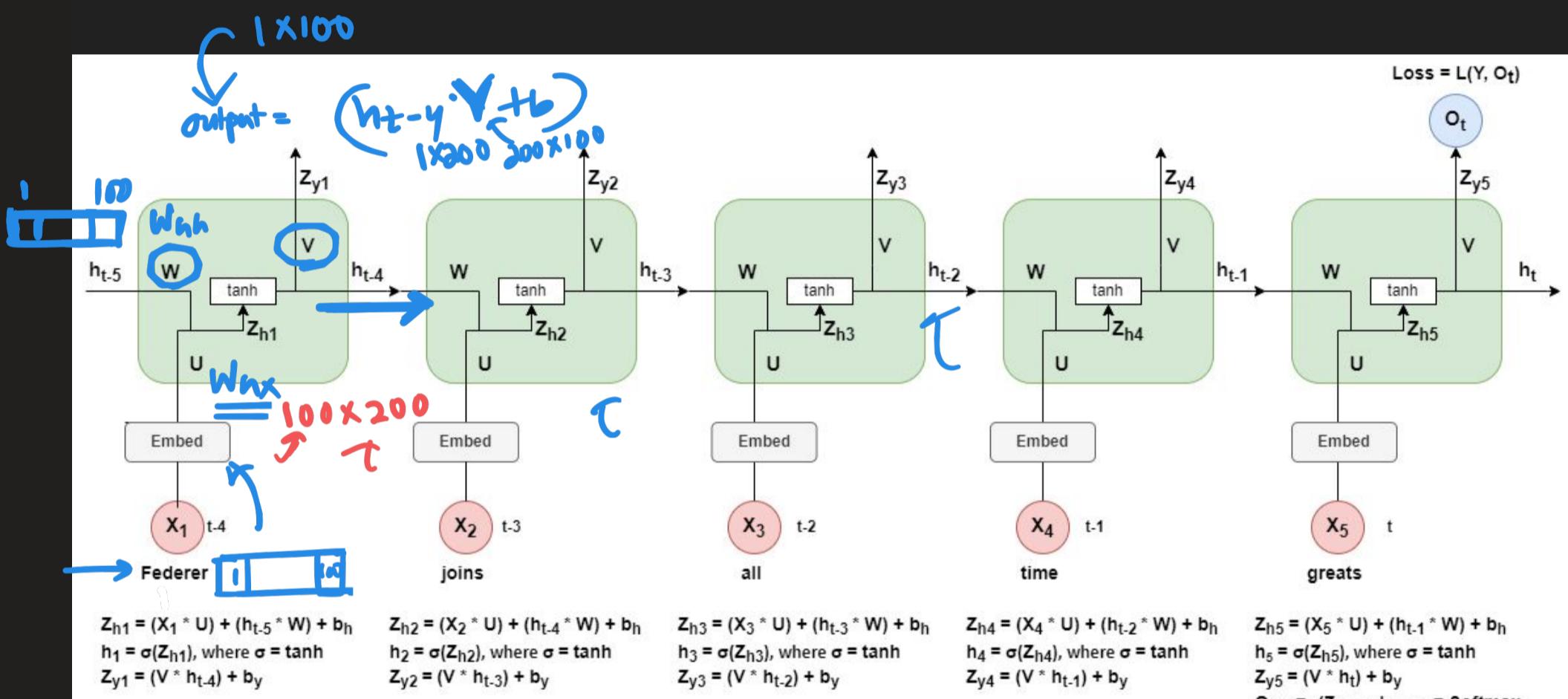
input word (embedding)



*RNN Forward Pass*

## RNN Forward Pass





Handwritten derivation of the hidden state update:

$$Z_{h1} = \underbrace{\sum_{\text{Word}} U}_{1 \times 100} \underbrace{100 \times 200}_{1 \times 200 \text{ currently}} + \underbrace{h_{t-5} \cdot W}_{1 \times 100 \quad 100 \times 200} + \underbrace{b_h}_{1 \times 200}$$

$$h_{t-4} = \tanh(Z_{h1}) \quad 1 \times 200$$

independent  $\rightarrow s_1$   
 $\rightarrow s_2$   
 $\vdots$   
 $\rightarrow s_{10}$

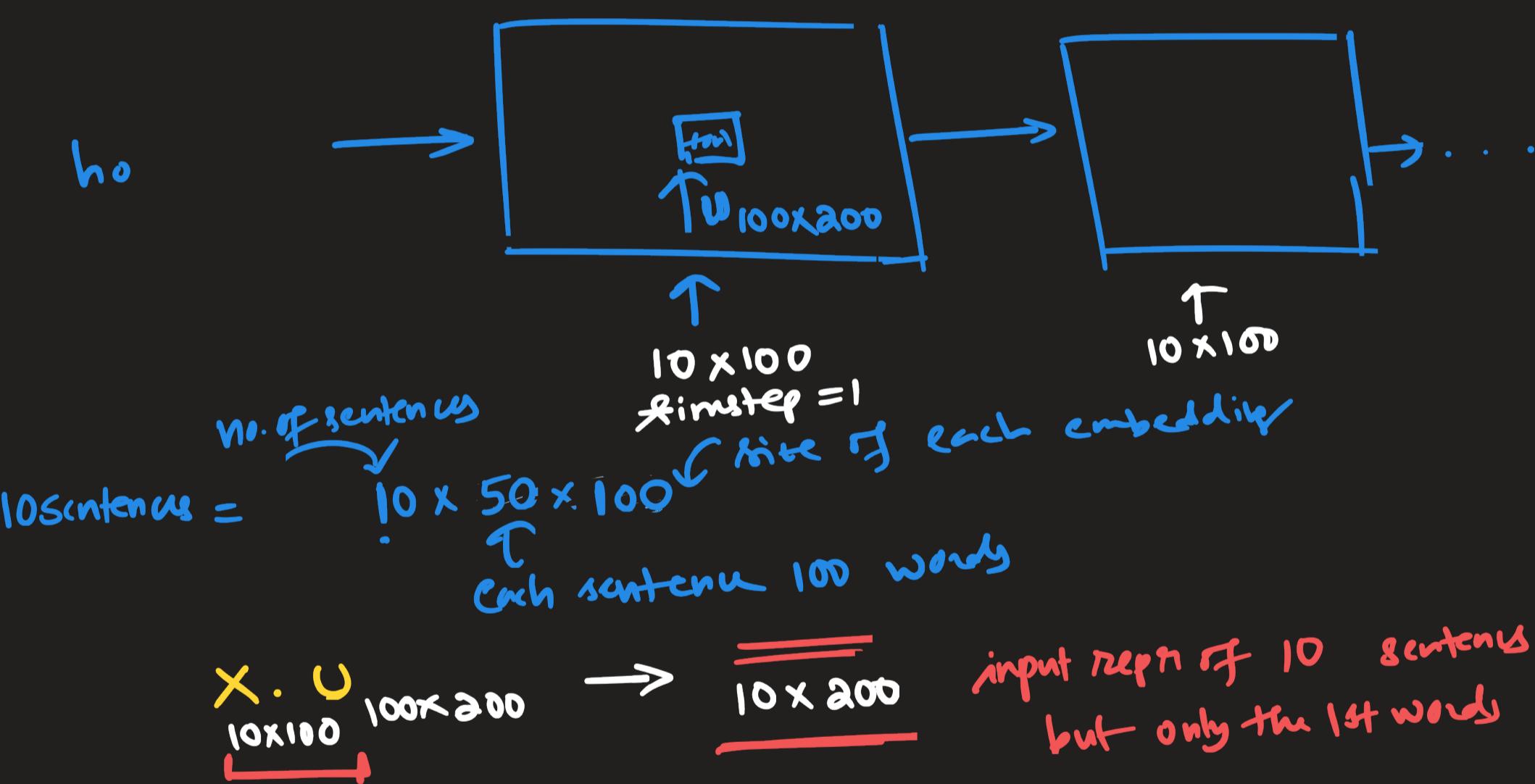
$w_1$   $\left[ \begin{array}{c} \\ \\ \end{array} \right] \cdot U_{100 \times 200}$   
 $w_2$   $\left[ \begin{array}{c} \\ \\ \end{array} \right] \cdot U_{100 \times 100}$

$w \times 100$

~~10~~ ✓  
 Sentence  $\times$  words  $\times$  size of embedding



How many such times step = 5



## Steps

U  $\times$  current word

W  $\times$  previous memory

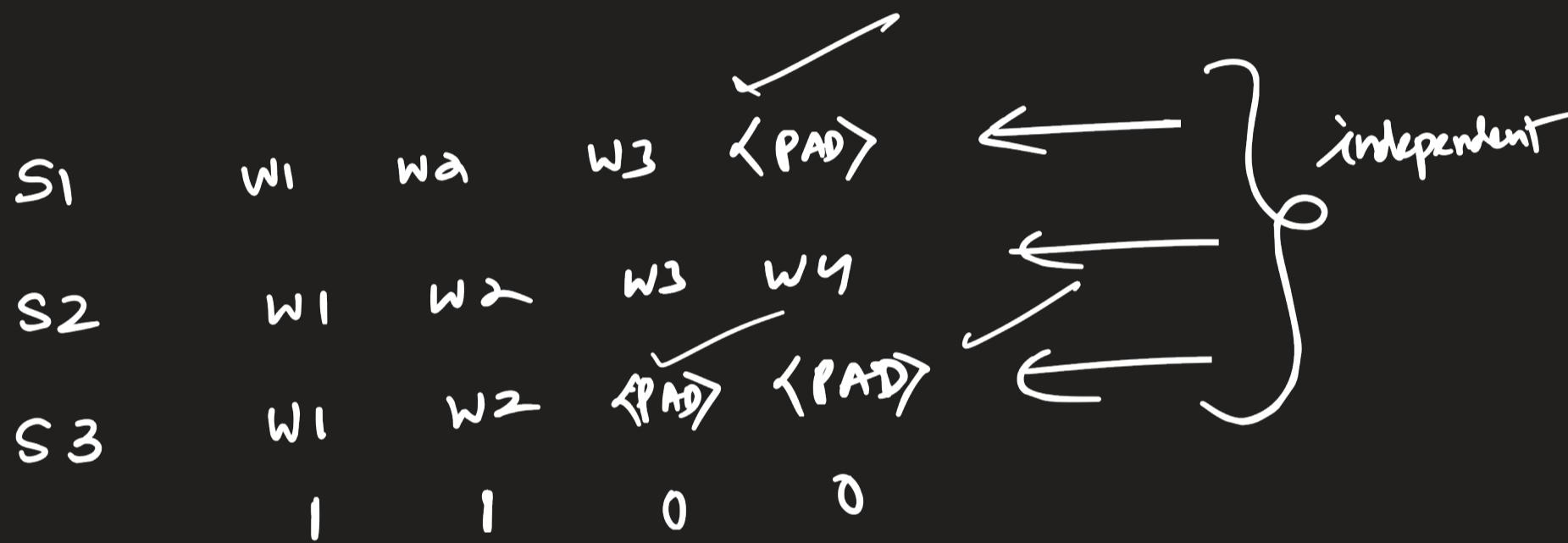
V  $\times$  memory

What does this word mean

What S was thinking so far

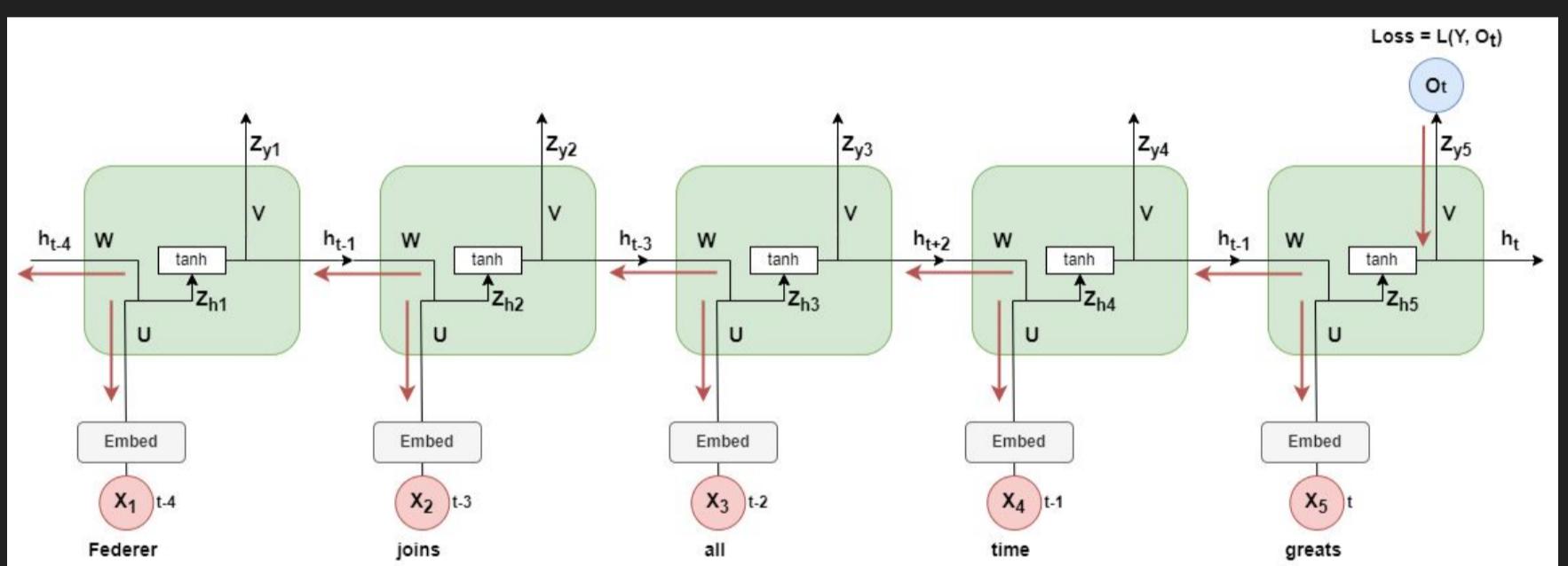
Based on my understanding  
so far what should I  
answer

QPM



$h_{0,000}$   
 $S_1 \rightarrow \square + \square \rightarrow \square \rightarrow \boxed{\square}$

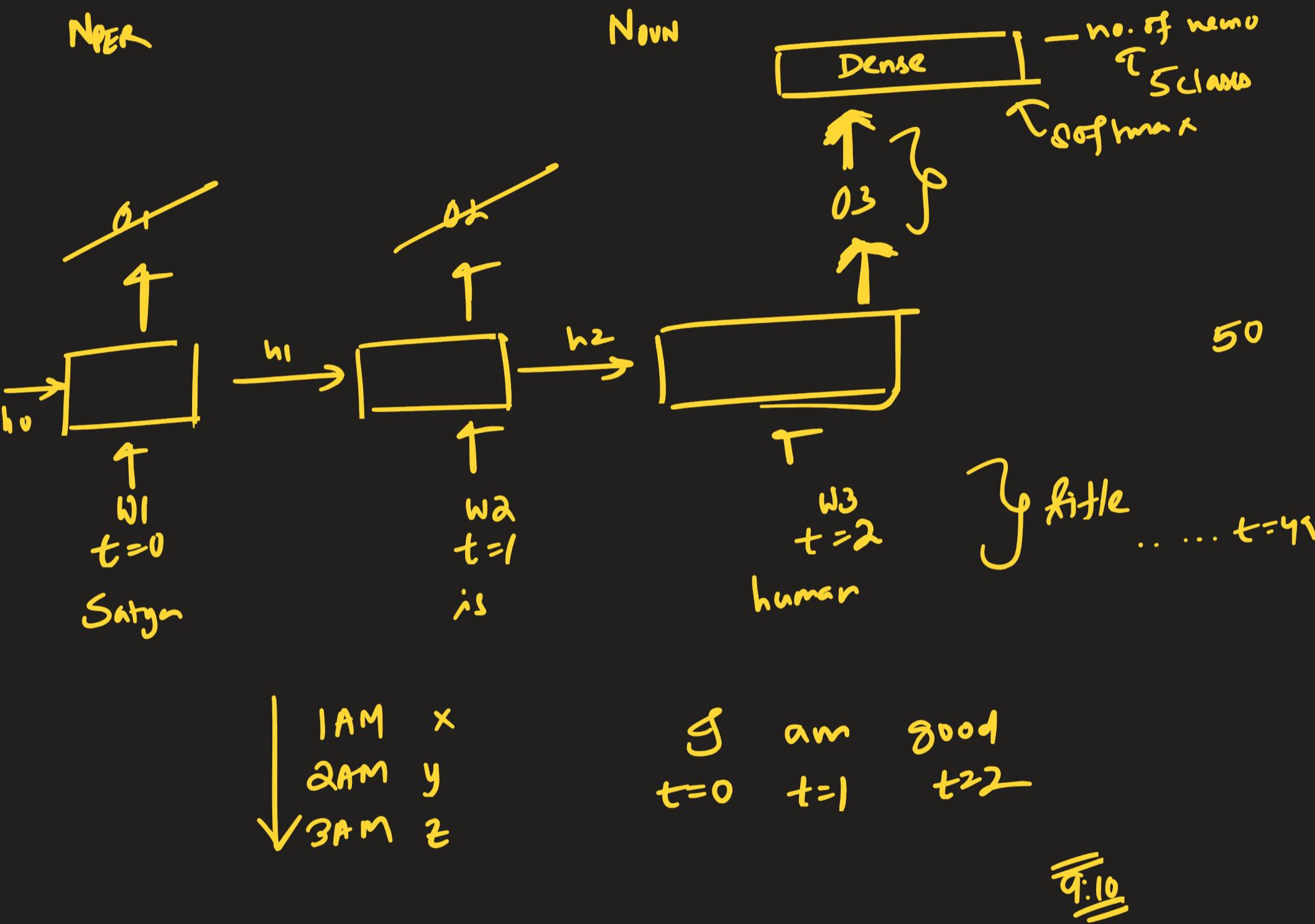
$\rightarrow S_2 \rightarrow \square$

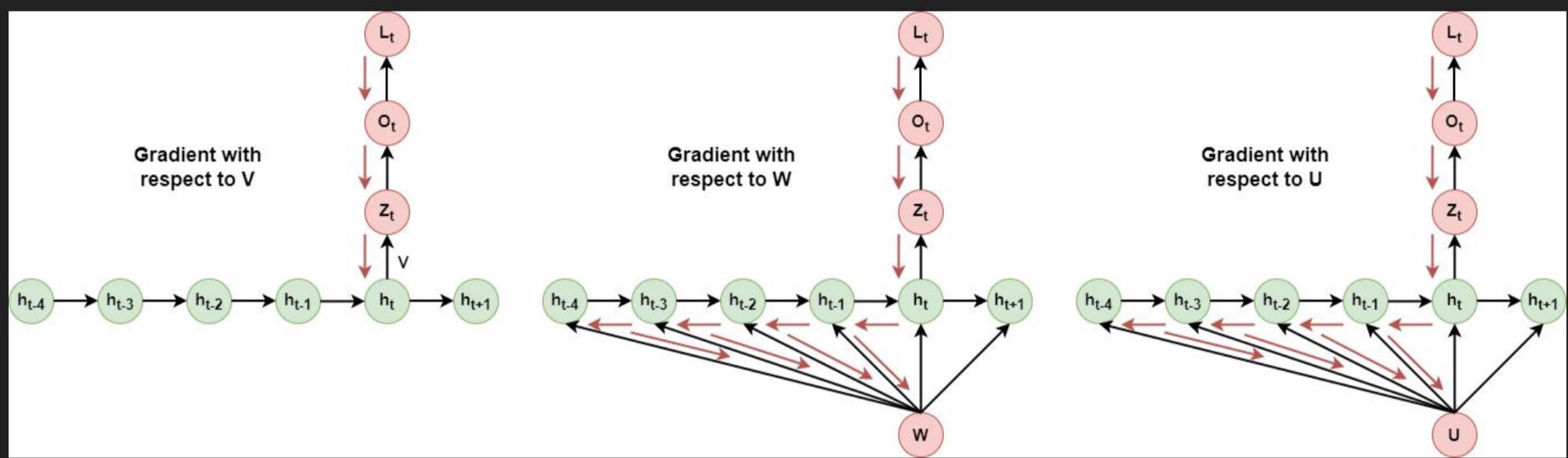


Great god is

{ god is great } ✓  
→ → ↑

*Weight Updates*





*Loss*





Let's generalize this to a single formula

$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial O_t} \frac{\partial O_t}{\partial Z_t} \frac{\partial Z_t}{\partial h_t} \sum_{r=1}^t \frac{\partial h_t}{\partial h_r} \frac{\partial h_r}{\partial W}$$

Update the weight matrix  $W$

$$W = W - \alpha \frac{\partial L_t}{\partial W}$$

Sachin Tendulkar is a cricketer

Sachin Tendulkar is a MP of Rijya Sabha

$$\vec{h}_t = f(Vx_t + Wh_{t-1})$$

$$\overleftarrow{h}_t = f'(W'x_t + W'\overleftarrow{h}_{t+1})$$

$$o_t = \text{concat}_{\text{avg}} \left[ \vec{h}_t, \overleftarrow{h}_t \right]$$

tree - vr  
I  
[ I am good ]

PER ✓ DES ORG  
T T T T T  
Sachin MP from Rajyasabha-  
batsman  
↓  
PLAYER

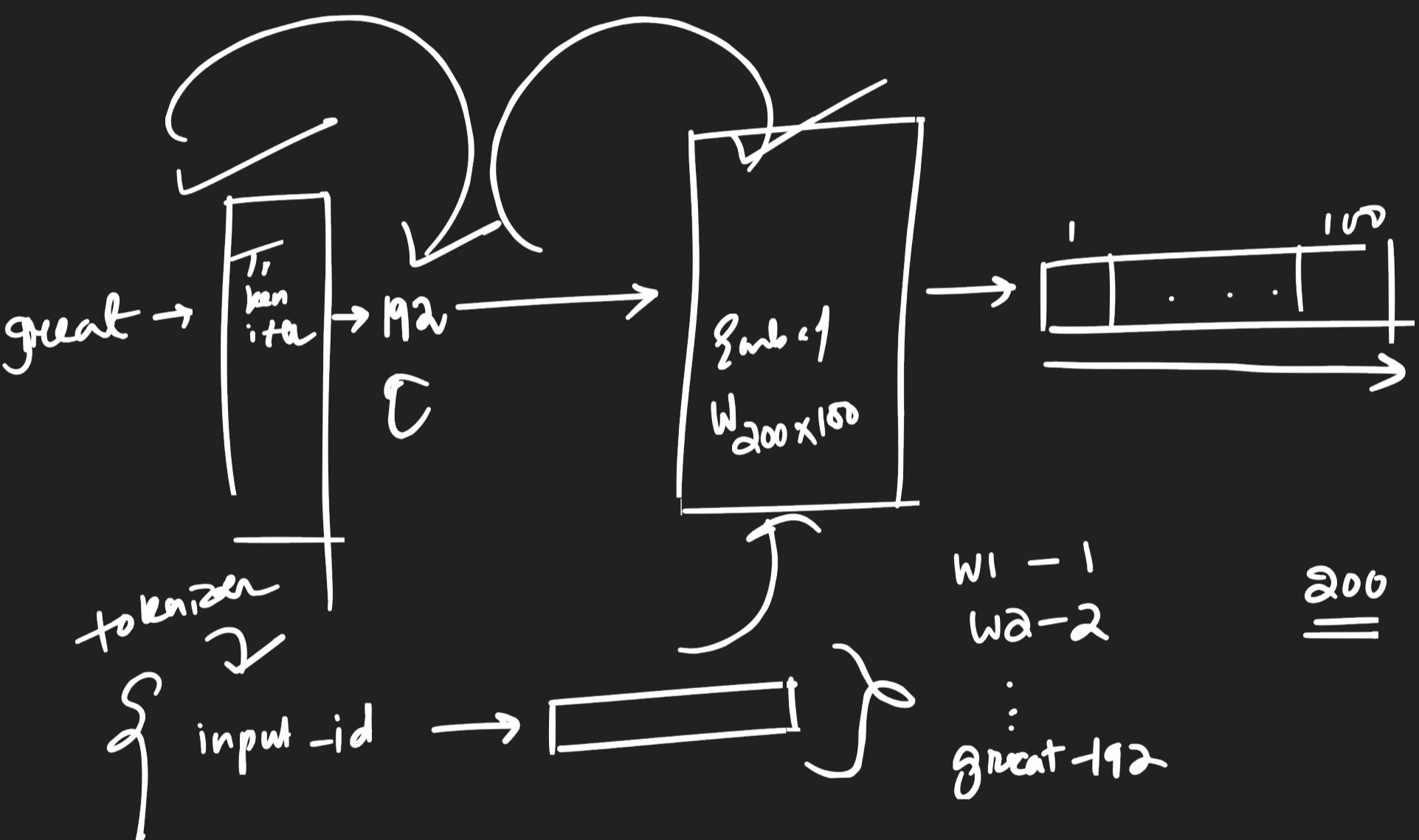
$$\text{Ran}' h = 1$$



---

$$(0.25 * 1.25) \times 100 \cong 2$$

Sequence words



$\text{10}$   
 $\text{500}$   
 (None, None, 100)

$\checkmark$  (None, 50)



each word



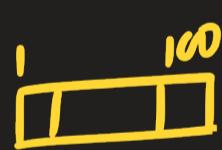
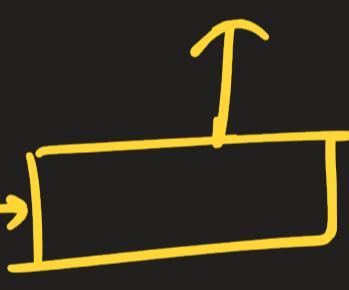
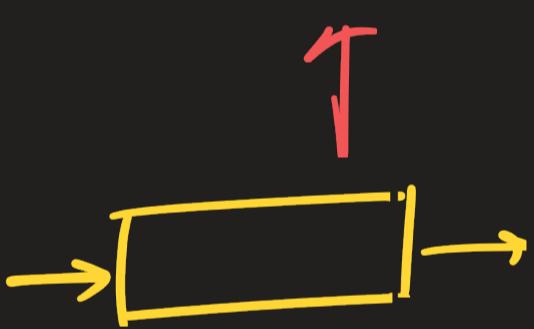
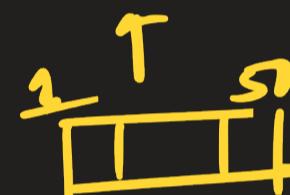
$t=1$   
 $w_1$   
 $t=2$   
 $w_2$   
 $\dots$   
 $t=500$   
 $w_{500}$

getrus segfum

return sig = False



$S_N$   
segfum



w<sub>1</sub>

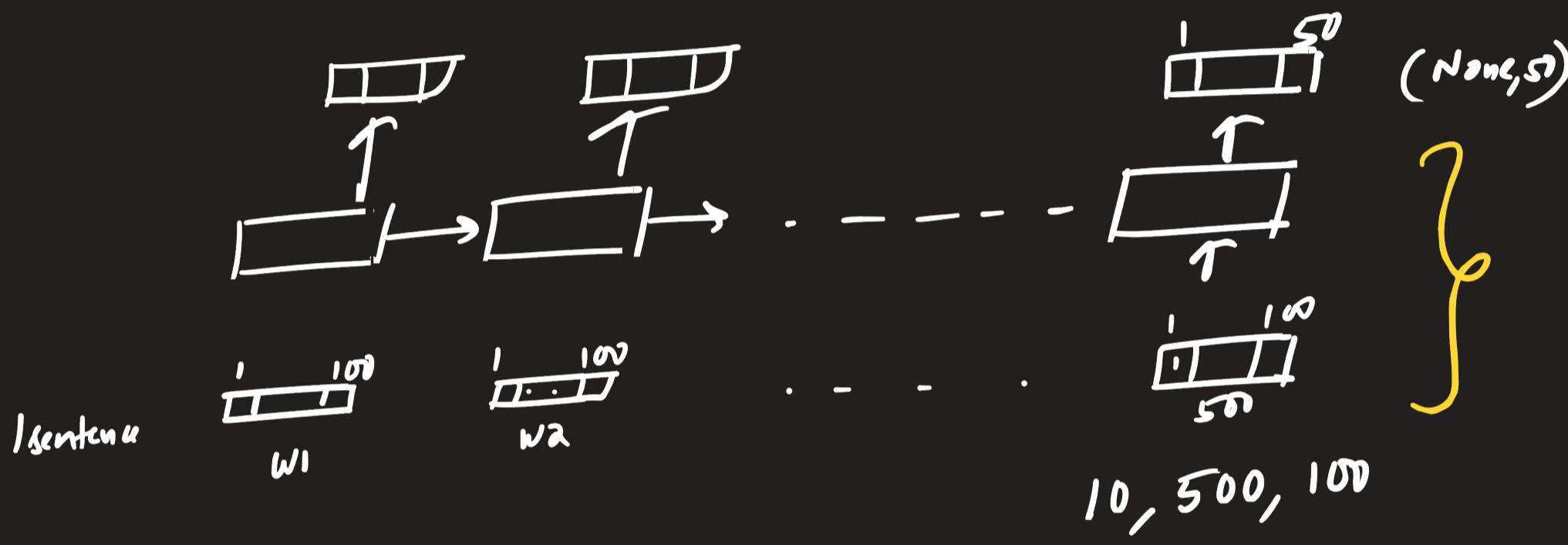


w<sub>2</sub>



w<sub>600</sub>

10  
=



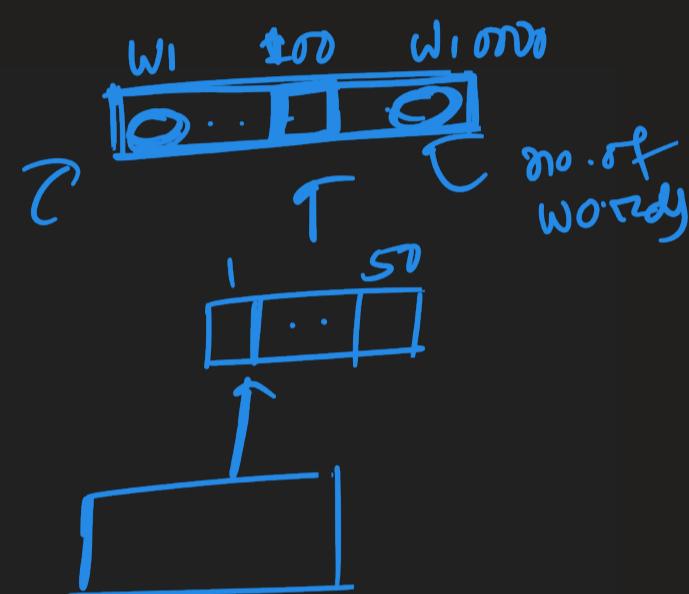
no. of sentence, each sentence of 500  
 None, None, 100 ← word dim = 100  
 None, 1, 50:

$100 \rightarrow \text{great}$

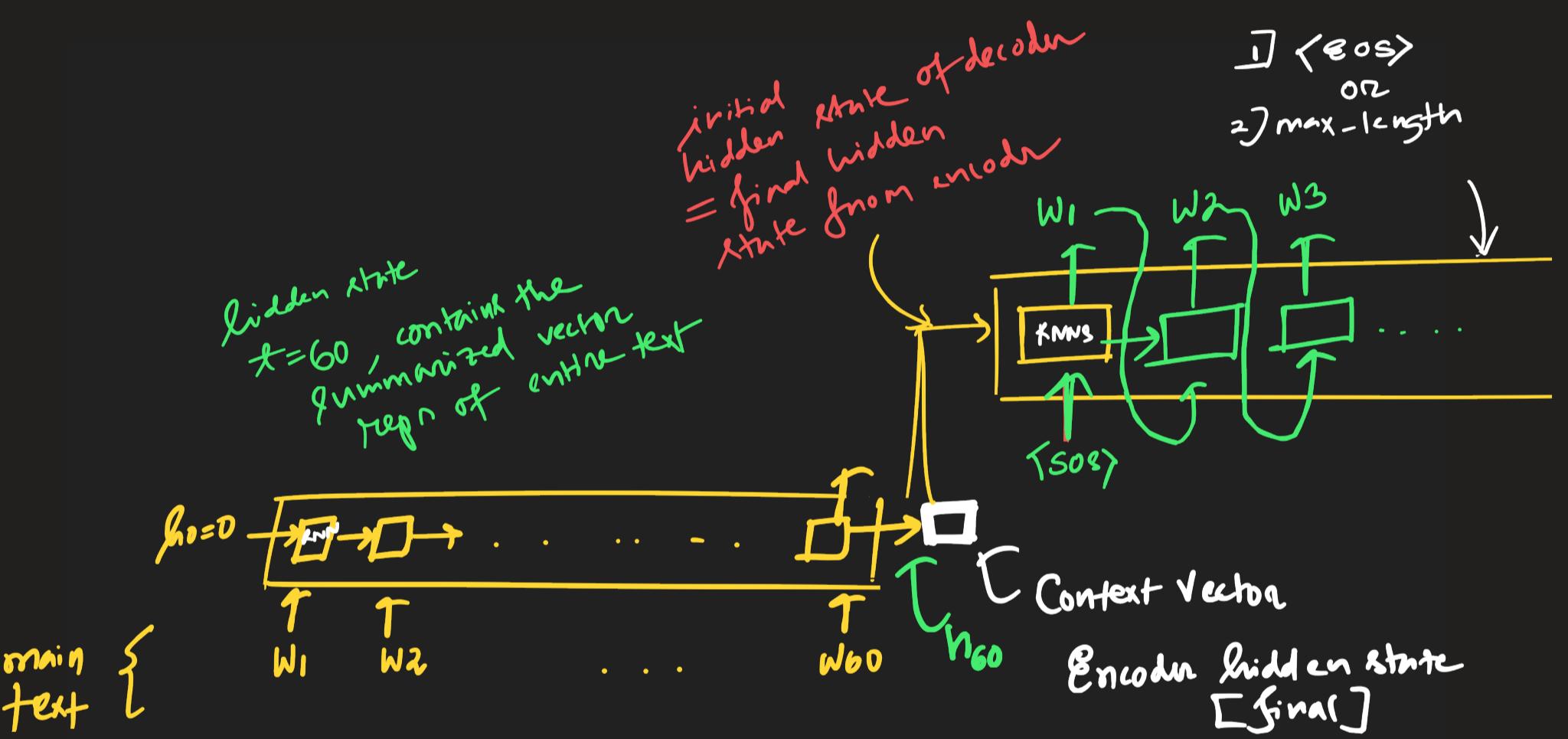
*LSTM*

Dense

softmax



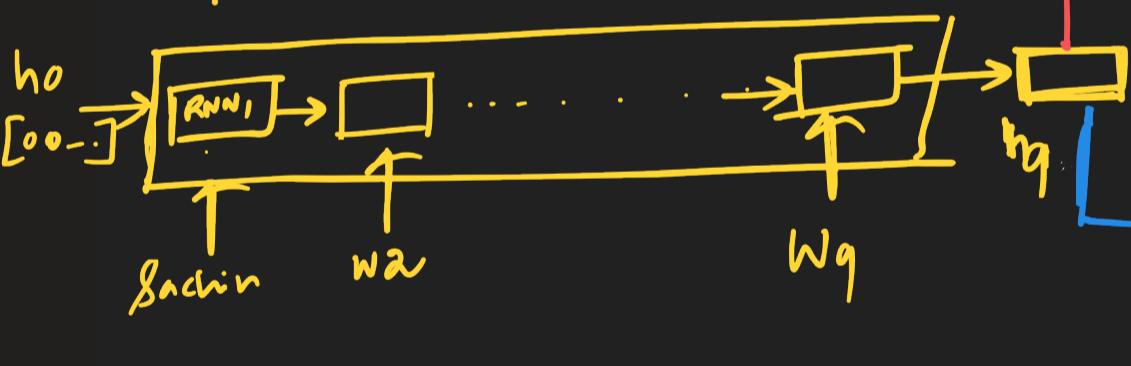
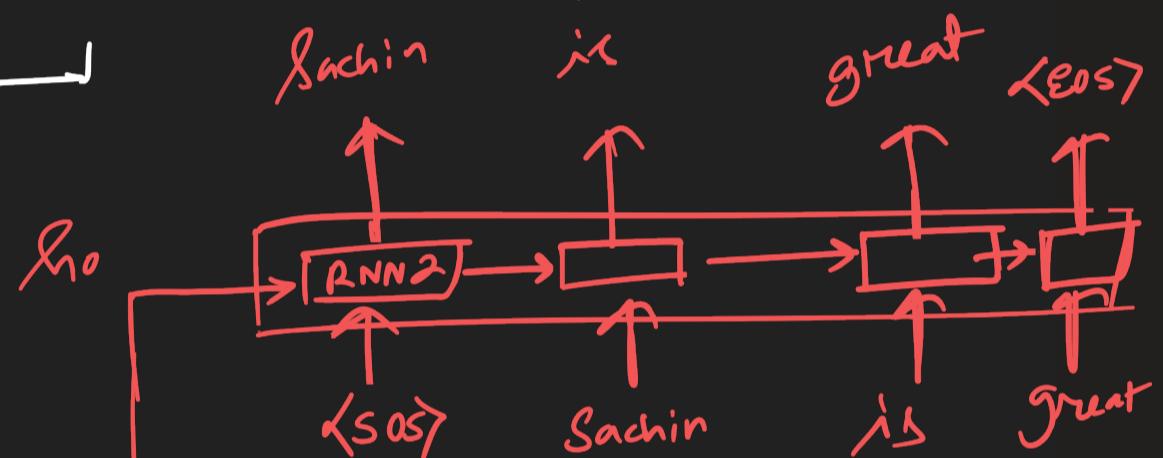
Satya Pattnaik



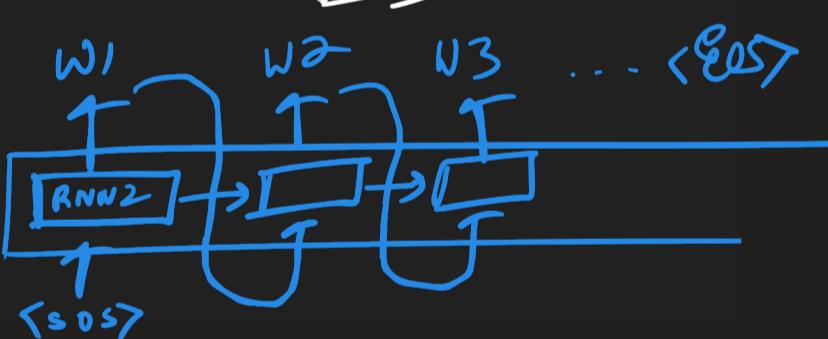
$X$   
 Sachin is a good  
 cricketer, he is from  
 Maharashtra  
 $\hookrightarrow q$

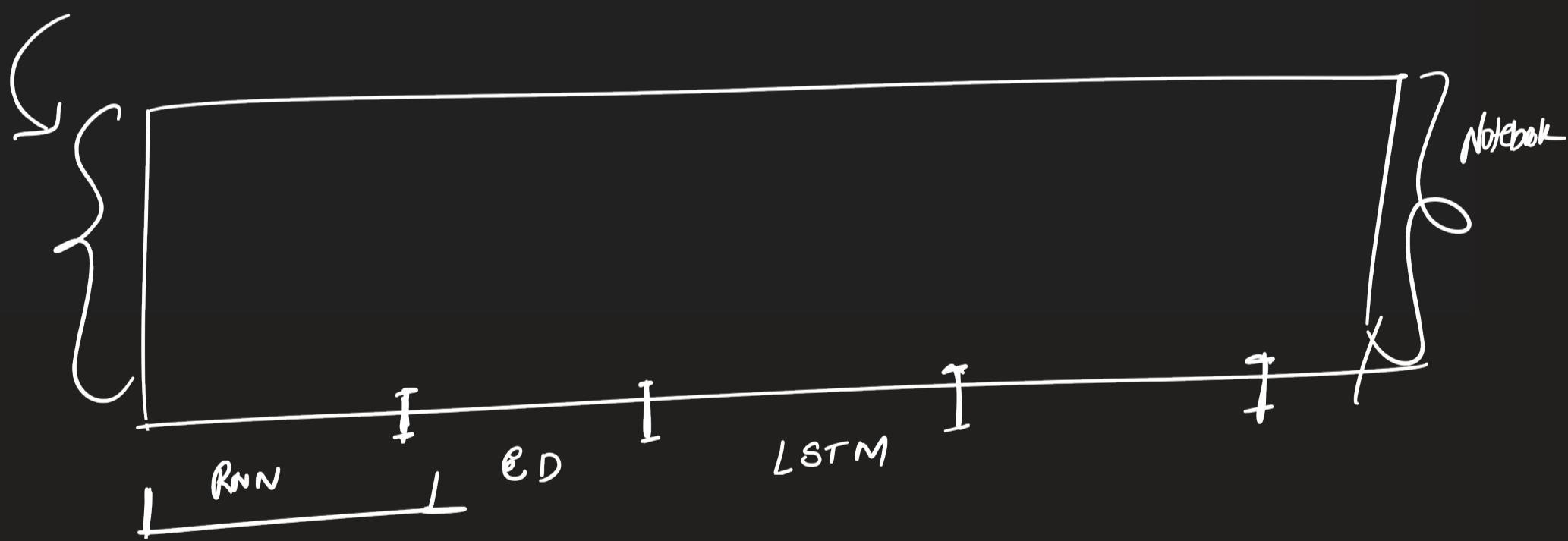
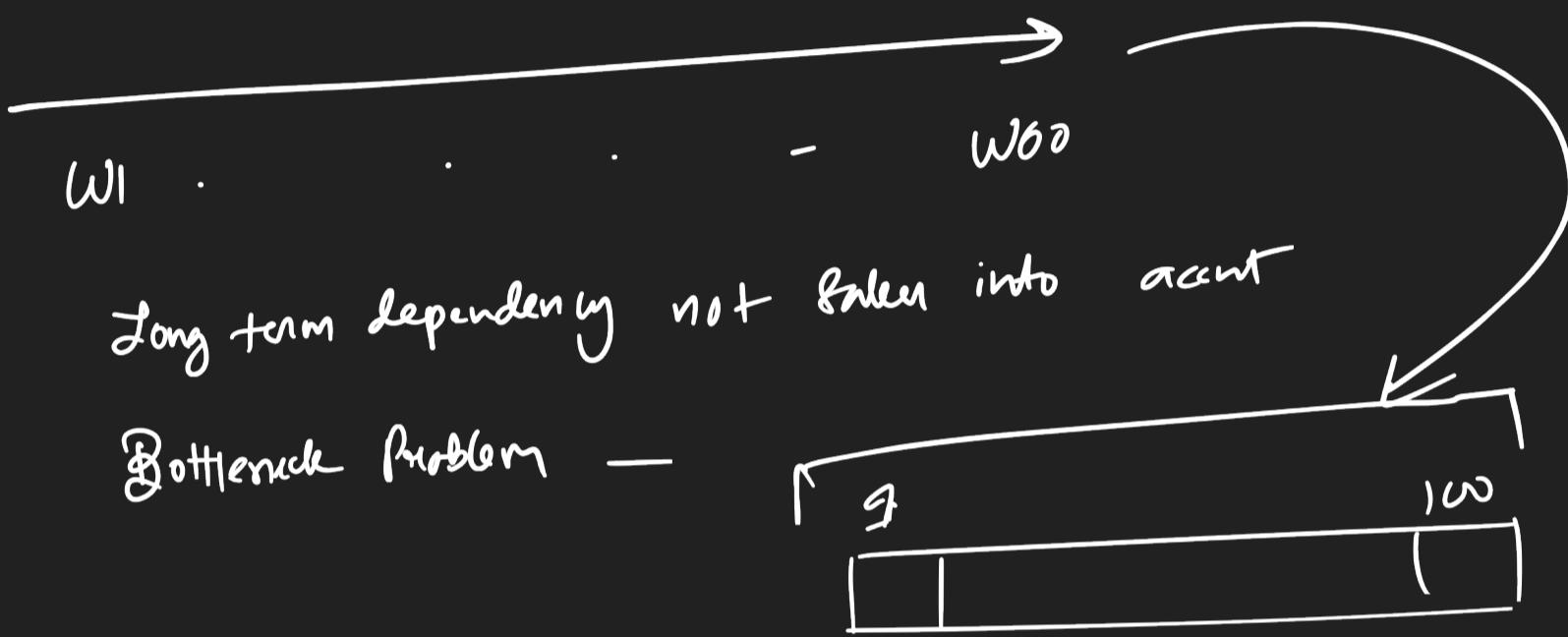
$y$   
 Sachin is great

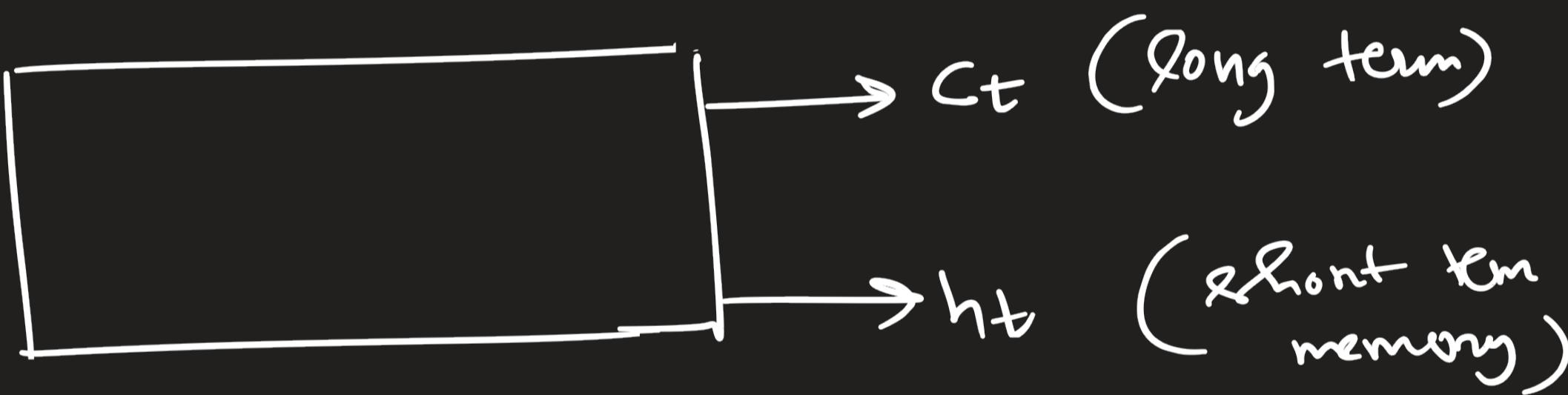
Training

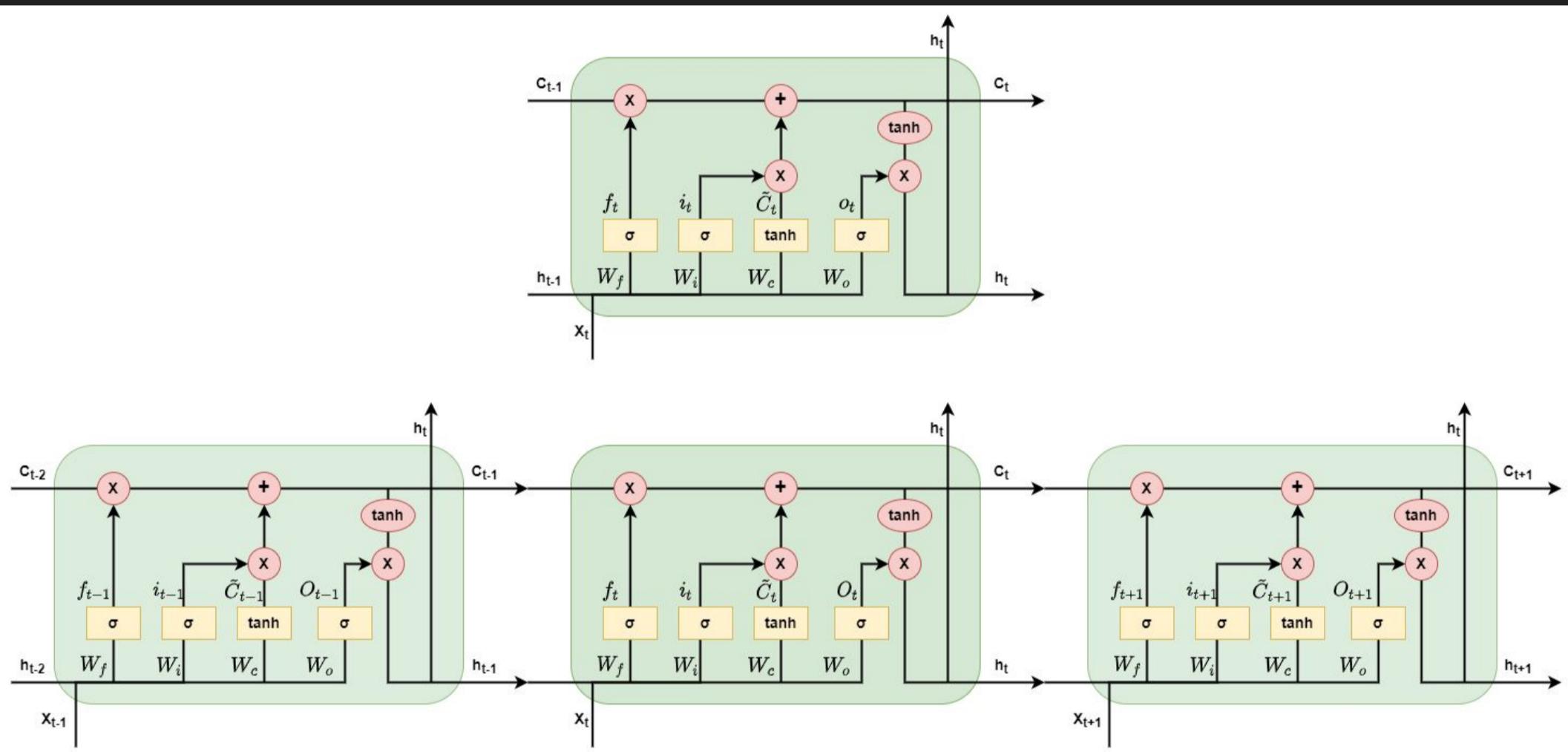


Inference

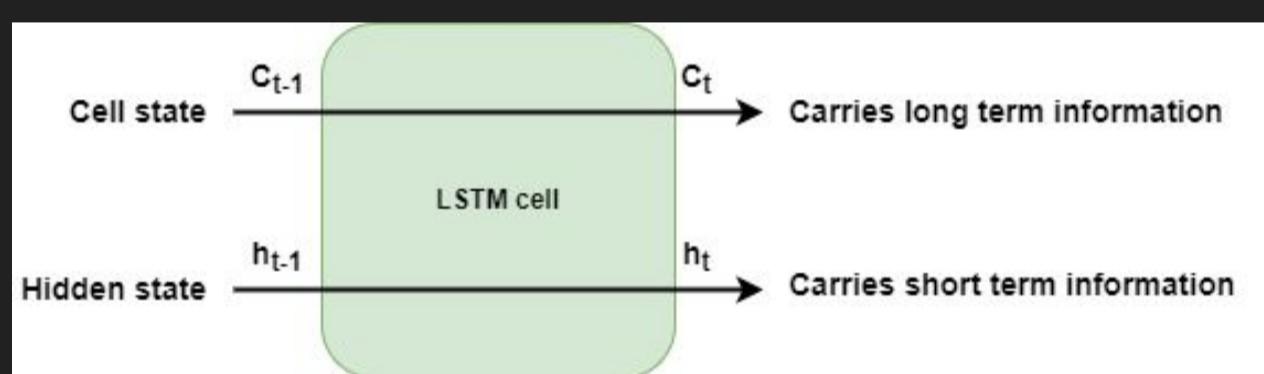








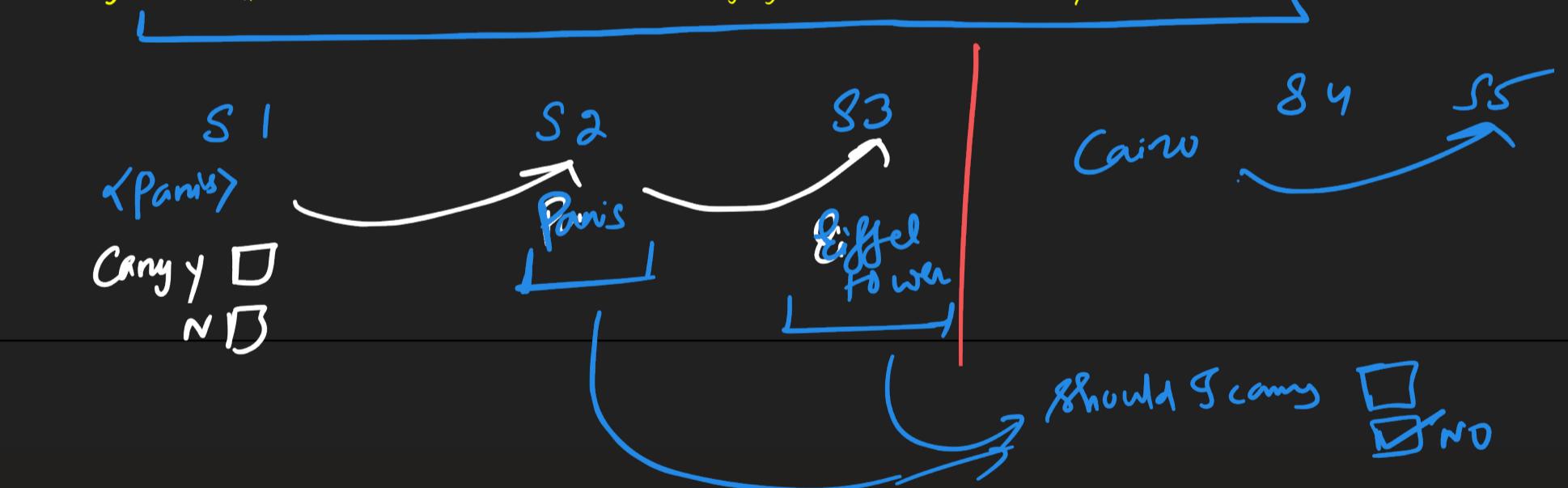
*How does LSTM overcome short term memory problem*

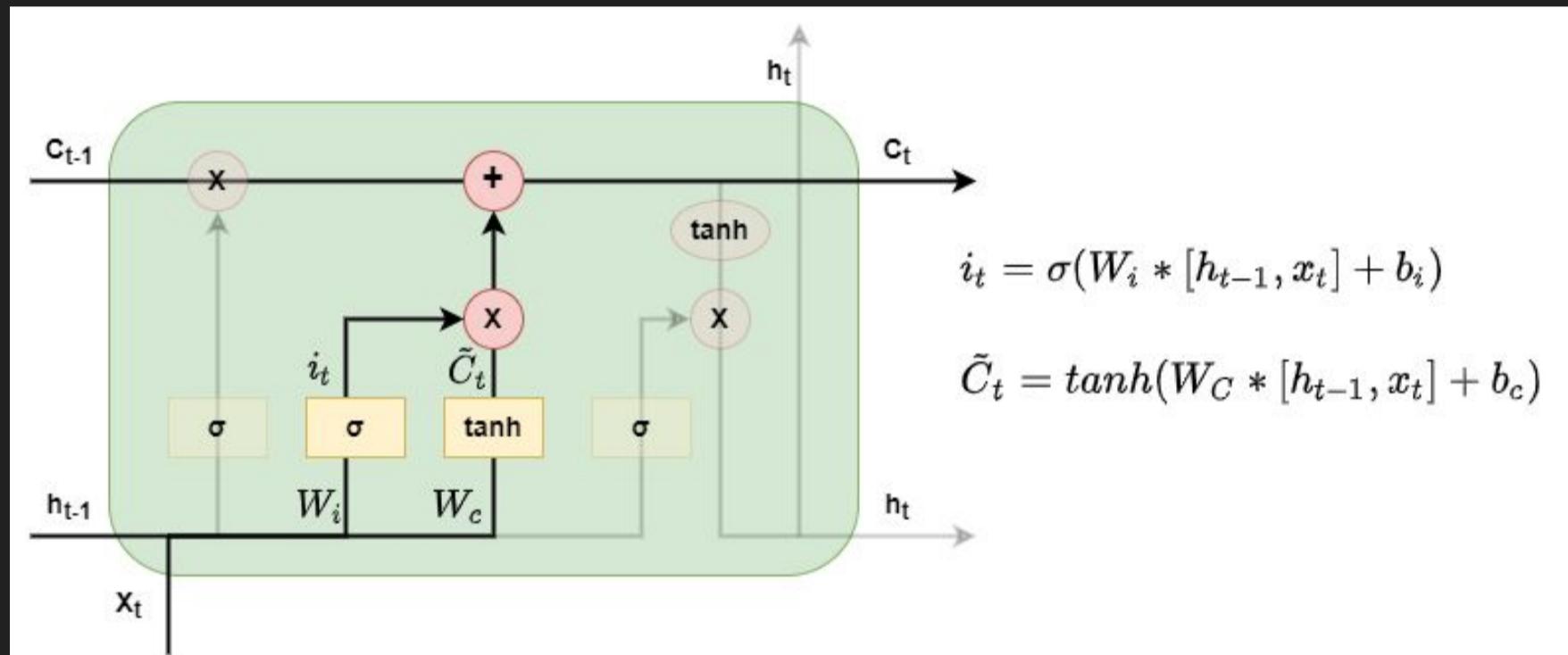


Let's understand it with an example sentence

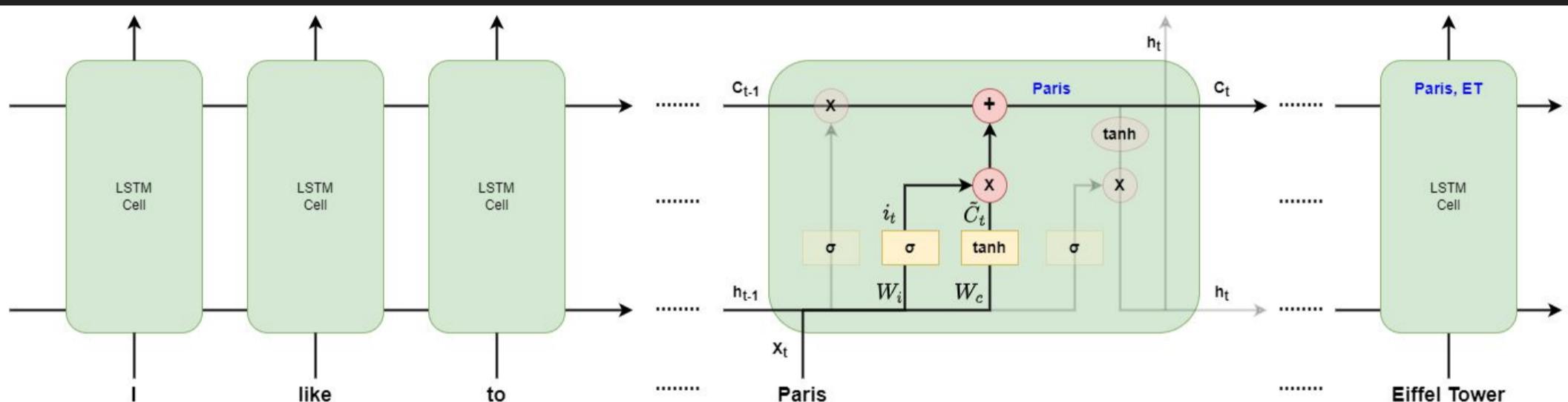
Sentence: I like to travel and I visited Paris last year as a part of my business trip.] It was winter season and too cold for someone from the sub-continent.] I stayed at Hilton hotel and I was fortunate enough to get a room where from my window I could see the Eiffel tower] I took my family on vacation to Cairo, Egypt last month.] They were quite a lot of market area for shopping to buy souvenirs and my kids enjoyed seeing the great pyramids for the first time.] 55

Summary: I saw Eiffel tower when I was in Paris and during my vacation I visited the Pyramids in Cairo.

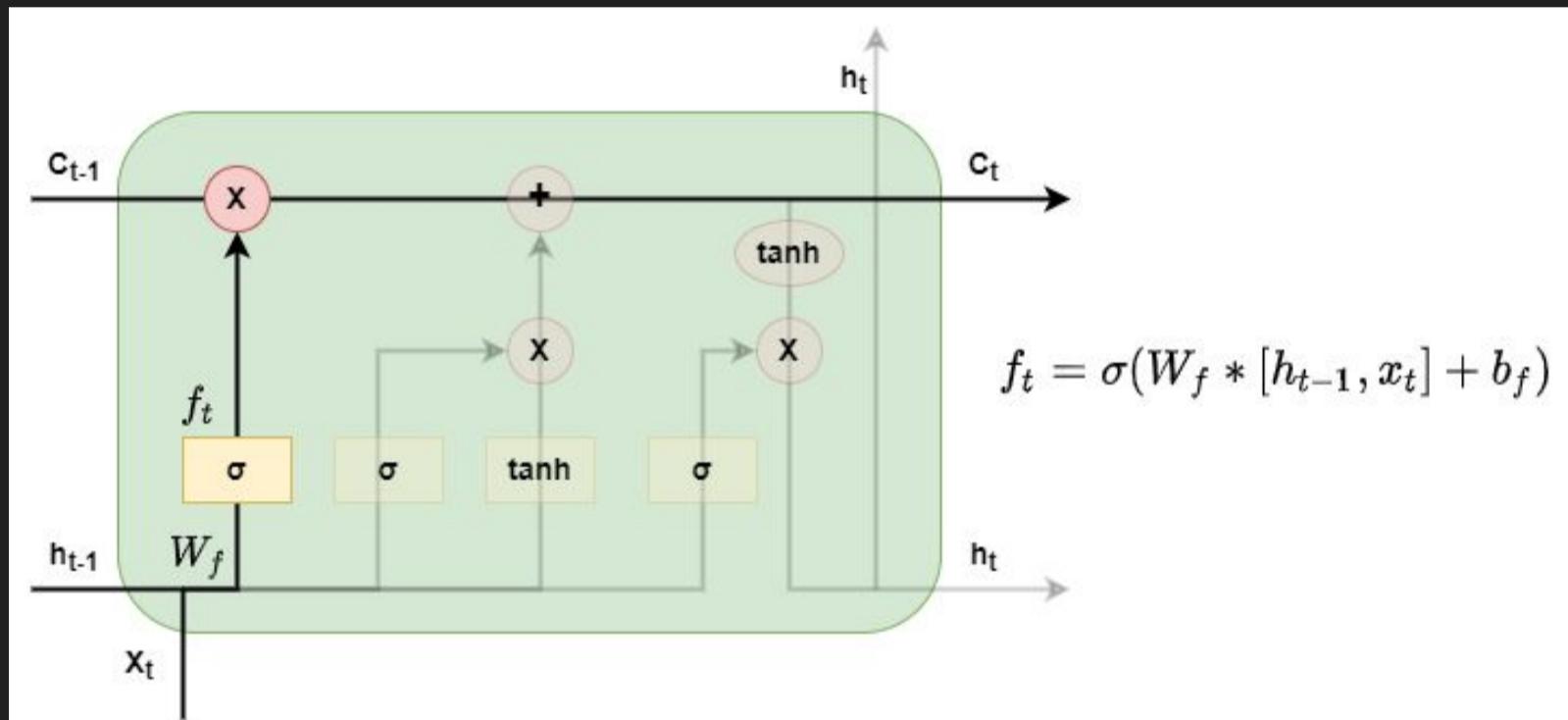




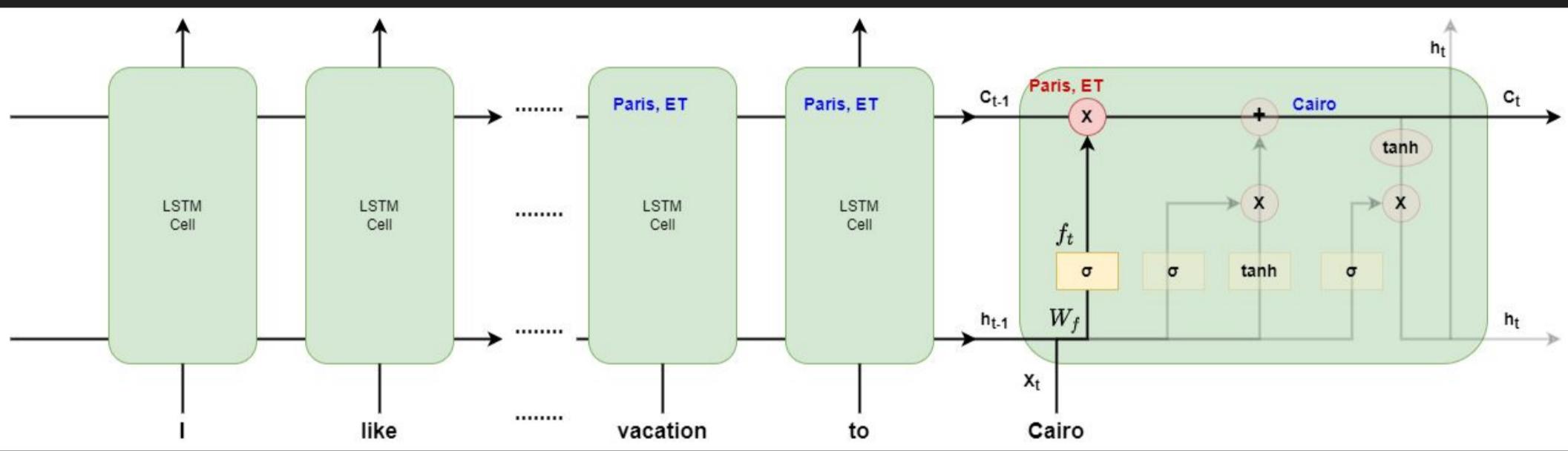




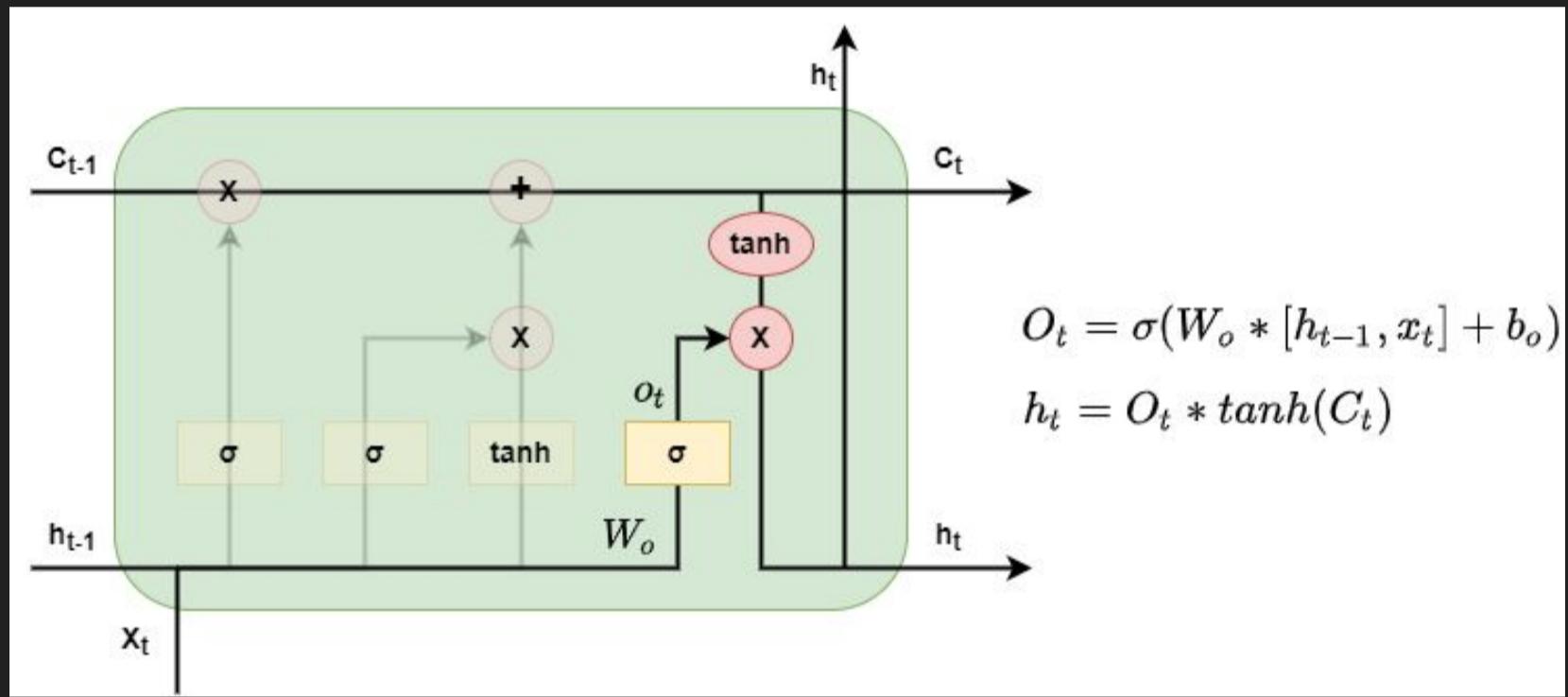


















*LSTM(Code)*

$$C_t = (i_t * \tilde{C_t}) + (f_t * C_{t-1})$$

$$\frac{\partial \mathbf{C_t}}{\partial \mathbf{C_{t-1}}}~=~\frac{\partial \mathbf{C_t}}{\partial \mathbf{f_t}}~\frac{\partial \mathbf{f_t}}{\partial \mathbf{h_{t-1}}}~\frac{\partial \mathbf{h_{t-1}}}{\partial \mathbf{C_{t-1}}}~+~\frac{\partial \mathbf{C_t}}{\partial \mathbf{i_t}}~\frac{\partial \mathbf{i_t}}{\partial \mathbf{h_{t-1}}}~\frac{\partial \mathbf{h_{t-1}}}{\partial \mathbf{C_{t-1}}}~+~\frac{\partial \mathbf{C_t}}{\partial \tilde{\mathbf{C_t}}}~\frac{\partial \tilde{\mathbf{C_t}}}{\partial \mathbf{h_{t-1}}}~\frac{\partial \mathbf{h_{t-1}}}{\partial \mathbf{C_{t-1}}}~+~\frac{\partial \mathbf{C_t}}{\partial \mathbf{C_{t-1}}}$$

$$C_t = (i_t * \tilde{C_t}) + (f_t * C_{t-1})$$

$$\frac{\partial \mathbf{C_t}}{\partial \mathbf{C_{t-1}}}~=~\frac{\partial \mathbf{C_t}}{\partial \mathbf{f_t}}~\frac{\partial \mathbf{f_t}}{\partial \mathbf{h_{t-1}}}~\frac{\partial \mathbf{h_{t-1}}}{\partial \mathbf{C_{t-1}}}~+~\frac{\partial \mathbf{C_t}}{\partial \mathbf{i_t}}~\frac{\partial \mathbf{i_t}}{\partial \mathbf{h_{t-1}}}~\frac{\partial \mathbf{h_{t-1}}}{\partial \mathbf{C_{t-1}}}~+~\frac{\partial \mathbf{C_t}}{\partial \tilde{\mathbf{C_t}}}~\frac{\partial \tilde{\mathbf{C_t}}}{\partial \mathbf{h_{t-1}}}~\frac{\partial \mathbf{h_{t-1}}}{\partial \mathbf{C_{t-1}}}~+~\frac{\partial \mathbf{C_t}}{\partial \mathbf{C_{t-1}}}$$

## What is learnt

You can calculate the **ROUGE-1** score like this:

$$\text{ROUGE-1}_{\text{recall}} = \frac{\text{unigram cand.} \cap \text{unigram ref.}}{|\text{unigram ref.}|}$$

$$\text{ROUGE-1}_{\text{precision}} = \frac{\text{unigram cand.} \cap \text{unigram ref.}}{|\text{unigram cand.}|}$$

$$\text{ROUGE-1}_{\text{F1}} = 2 \cdot \frac{\text{recall} \cdot \text{precision}}{\text{recall} + \text{precision}}$$

$$\text{ROUGE-L}_{\text{recall}} = \frac{\text{LCS}(\text{cand.}, \text{ref.})}{\#\text{words in ref.}}$$

$$\text{ROUGE-L}_{\text{precision}} = \frac{\text{LCS}(\text{cand.}, \text{ref.})}{\#\text{ words in cand.}}$$

$$\text{ROUGE-L}_{\text{F1}} = 2 \cdot \frac{\text{recall} \cdot \text{precision}}{\text{recall} + \text{precision}}$$

	Candidate	Reference
Text	I really loved reading the Hunger Games.	I loved reading the Hunger Games.
Unigrams	[('I',), ('really',), ('loved',), ('reading',), ('the',), ('Hunger',), ('Games.',)]	[('I',), ('loved',), ('reading',), ('the',), ('Hunger',), ('Games.',)]
Bigrams	[('I', 'really'), ('really', 'loved'), ('loved', 'reading'), ('reading', 'the'), ('the', 'Hunger'), ('Hunger', 'Games.')]	[('I', 'loved'), ('loved', 'reading'), ('reading', 'the'), ('the', 'Hunger'), ('Hunger', 'Games.')]
LCS	I loved reading the Hunger Games.	

The ROUGE-1 calculation is performed as

$$\text{ROUGE-1}_{\text{recall}} = \frac{\text{unigram cand.} \cap \text{unigram ref.}}{|\text{unigram ref.}|} = \frac{6}{6} = 1$$

$$\text{ROUGE-1}_{\text{precision}} = \frac{\text{unigram cand.} \cap \text{unigram ref.}}{|\text{unigram cand.}|} = \frac{6}{7}$$

$$\text{ROUGE-1}_{\text{F1}} = 2 \cdot \frac{\text{recall} \cdot \text{precision}}{\text{recall} + \text{precision}} = \frac{12}{13} \approx 0.923$$

Similarly, the ROUGE-L is calculated as

$$\text{ROUGE-L}_{\text{recall}} = \frac{\text{LCS(cand., ref.)}}{\#\text{words in ref.}} = \frac{6}{6} = 1$$

$$\text{ROUGE-L}_{\text{precision}} = \frac{\text{LCS(cand., ref.)}}{\#\text{ words in cand.}} = \frac{6}{7}$$

$$\text{ROUGE-L}_{\text{F1}} = 2 \cdot \frac{\text{recall} \cdot \text{precision}}{\text{recall} + \text{precision}} = \frac{12}{13} \approx 0.923$$





