

## Sequential Data

22 October 2022 13:09

ANN → tabular data      } CNN → images

{ RNN → Recurrent NN  
is type of sequential model  
to work on sequential data }

iq	marks	gender	placement
19	0	No	does not matter
marks	0	0	
gender	0	0	

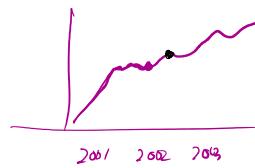
RNN → NLP → ML

CNN → images → computer vision

eg → text → sequential data

Hi my name is Nitish

Time series



Speech

sequential

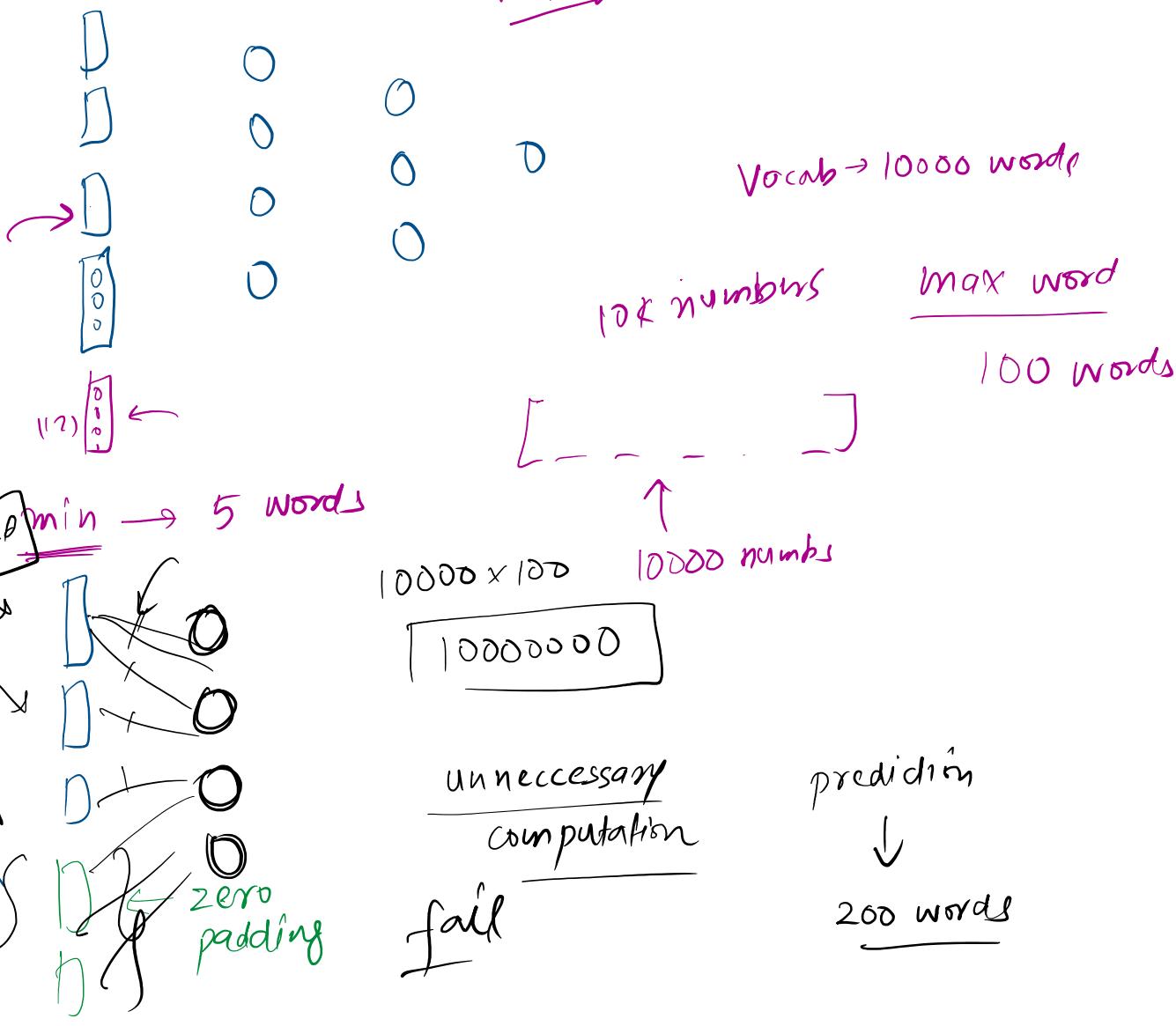
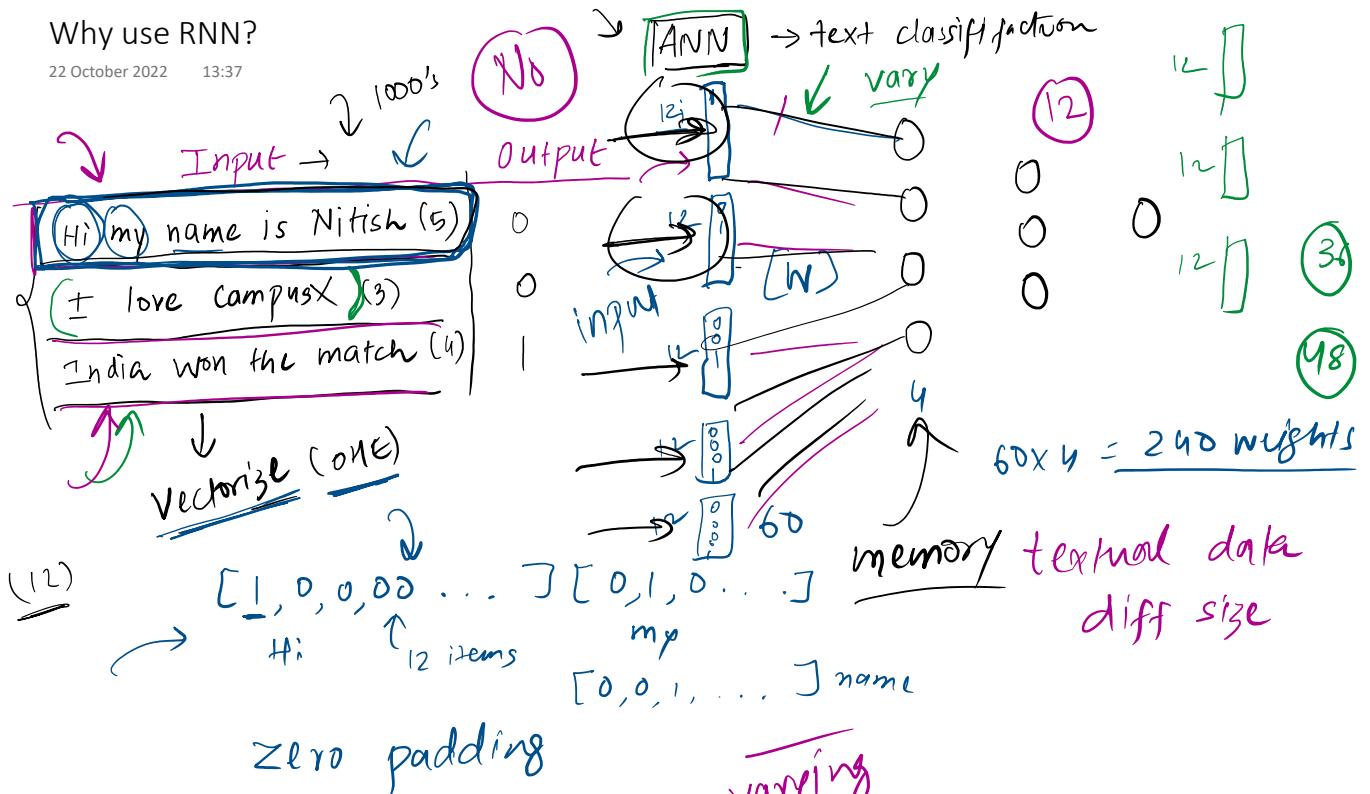
sequential

DNA sequence

RNN  
→ RNN Why?  
→ Application → RNN  
→ Roadmap ↴

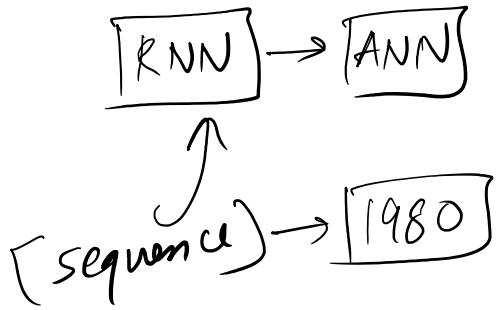
## Why use RNN?

22 October 2022 13:37



n in text input → varying size

- 1) text input → varying size
- 2) zero padding → unnecessary computation
- 3) Prediction problem
- 4) Totally disregarding the sequence info

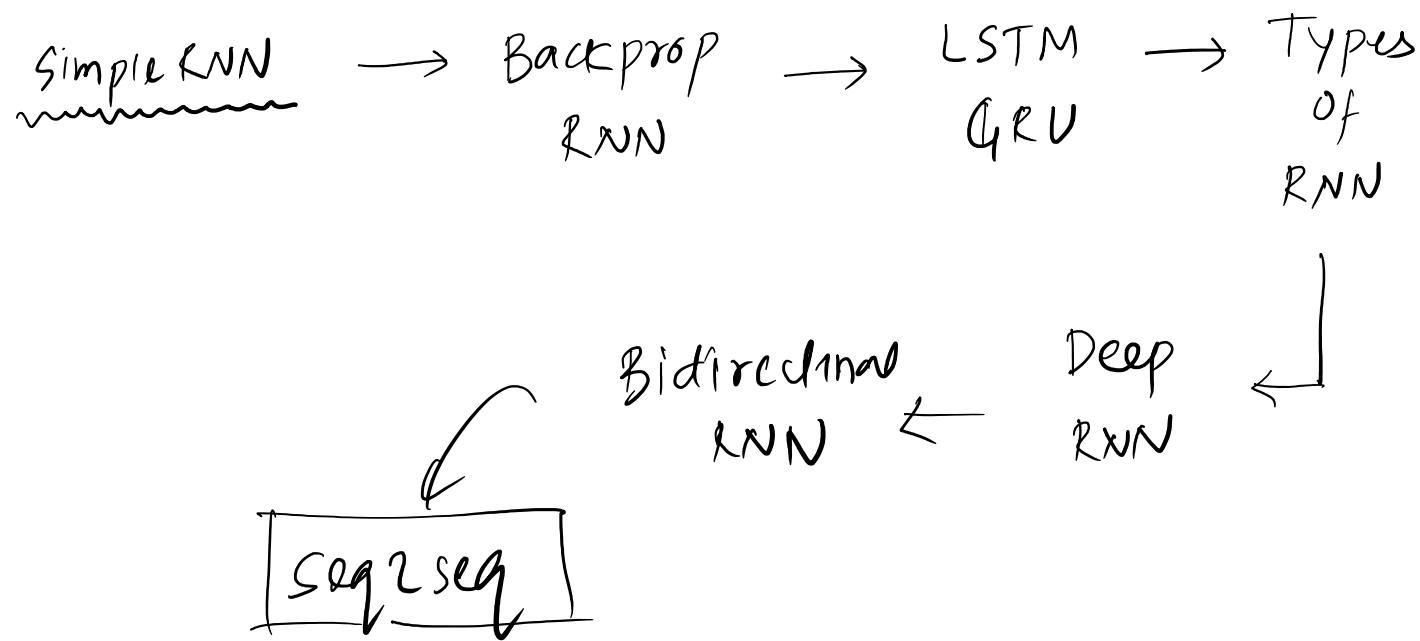


# RNN Applications

22 October 2022 13:37

# Roadmap

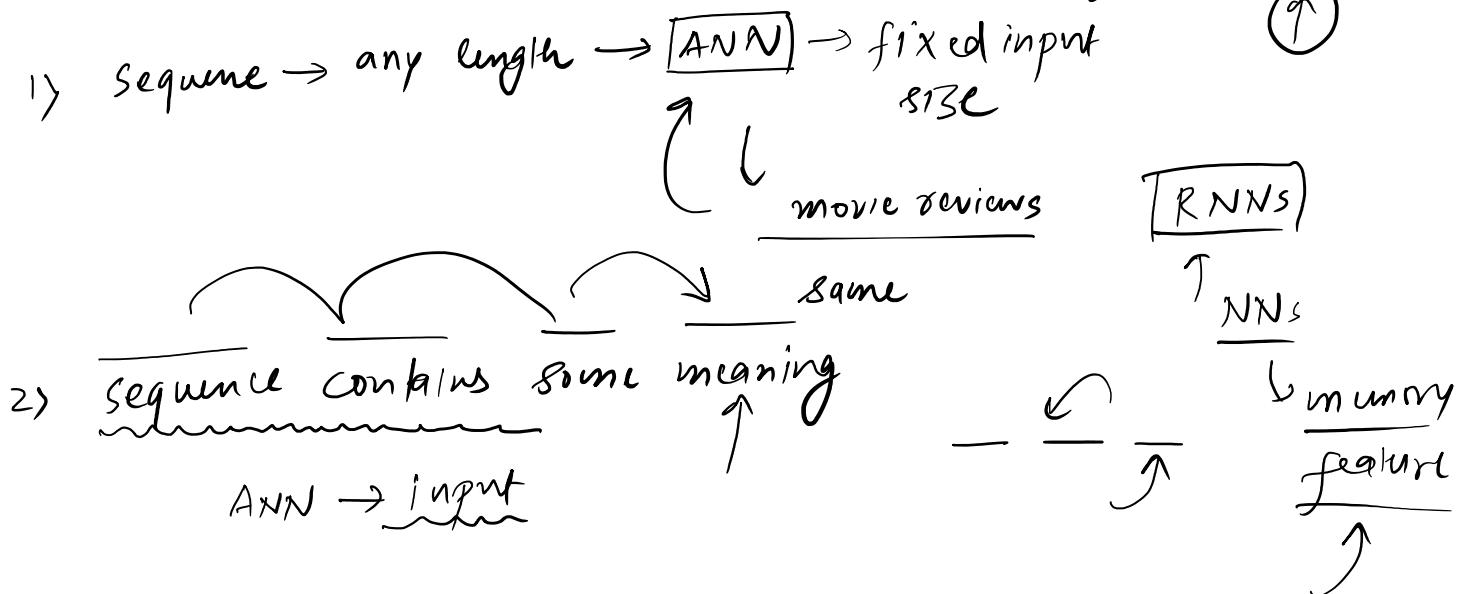
22 October 2022 13:37



## Why RNNs?

29 October 2022 13:30

zero padding  $\rightarrow$  cost of computation



### RNN architecture

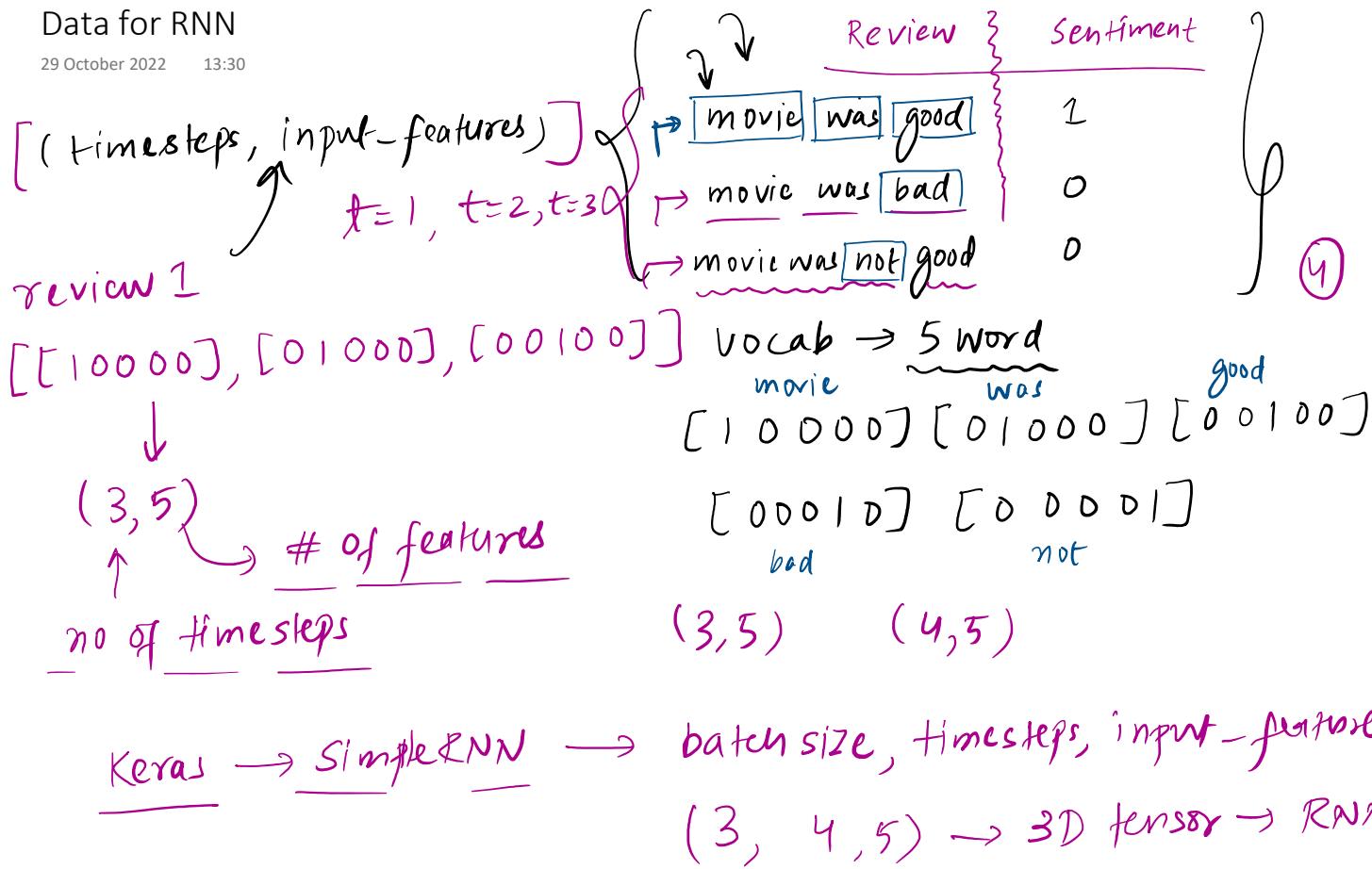
RNN forward prop  $\rightarrow$  prediction

input  $\rightarrow$  output

Codes  $\rightarrow$  Solidify

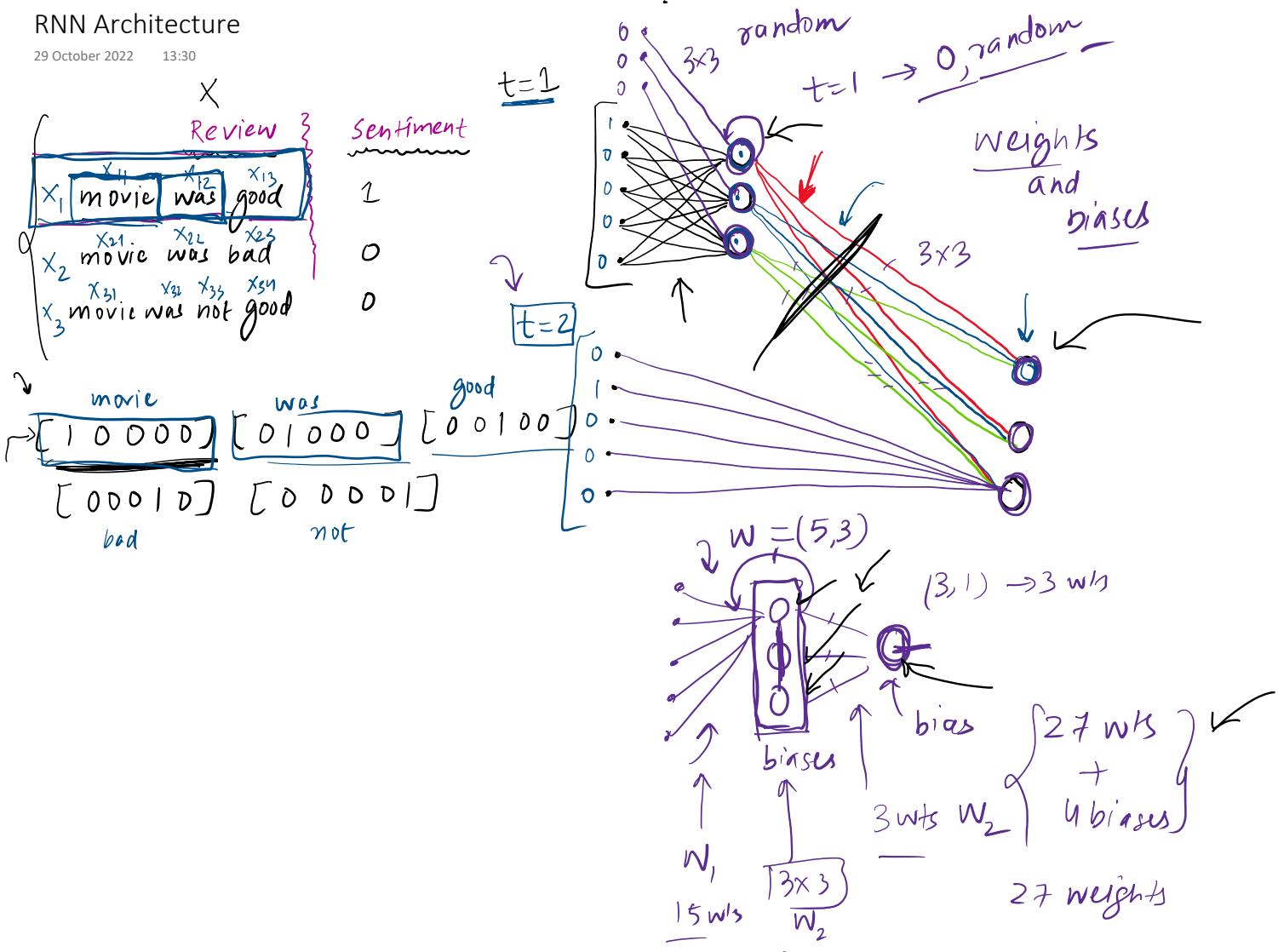
## Data for RNN

29 October 2022 13:30



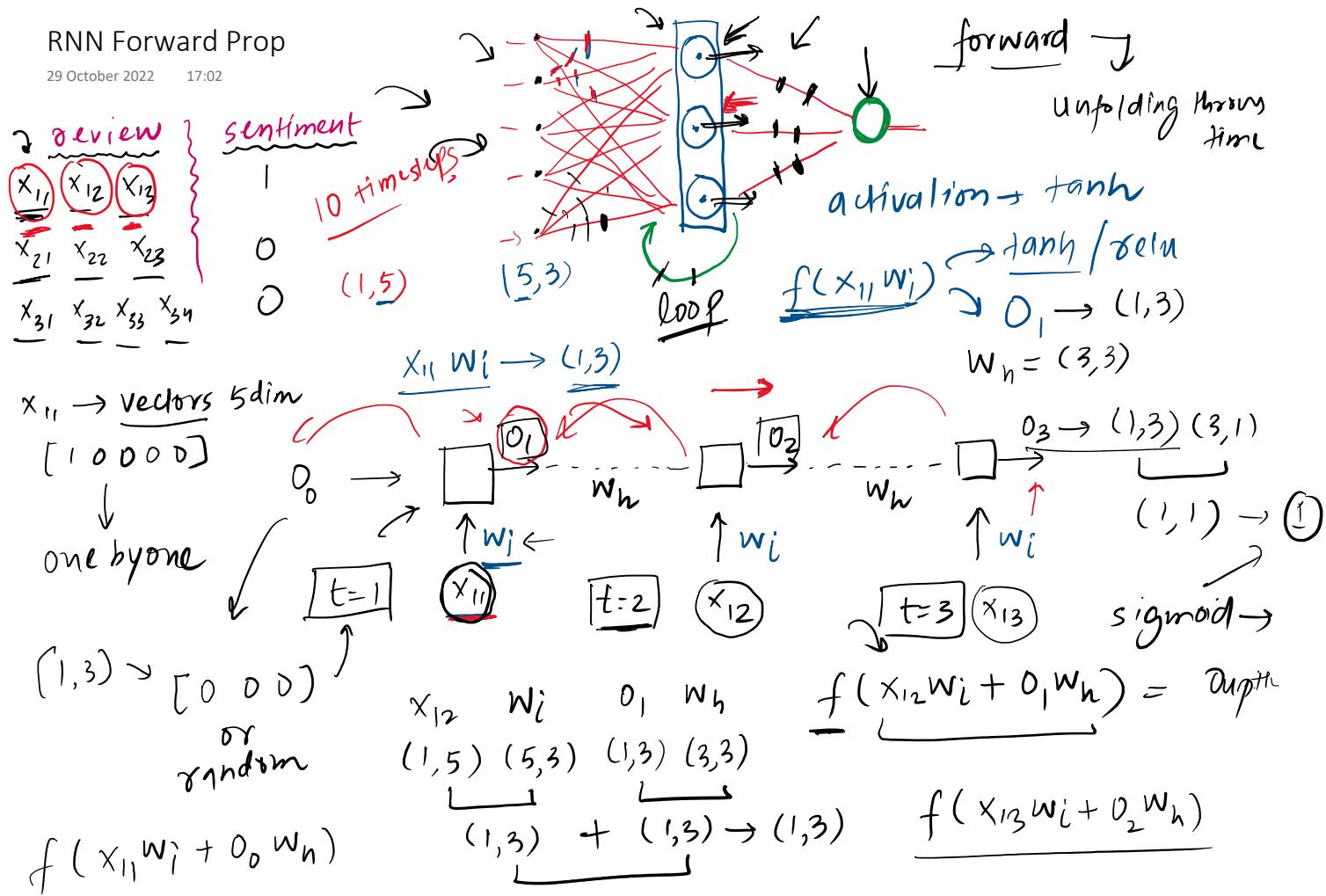
## RNN Architecture

29 October 2022 13:30



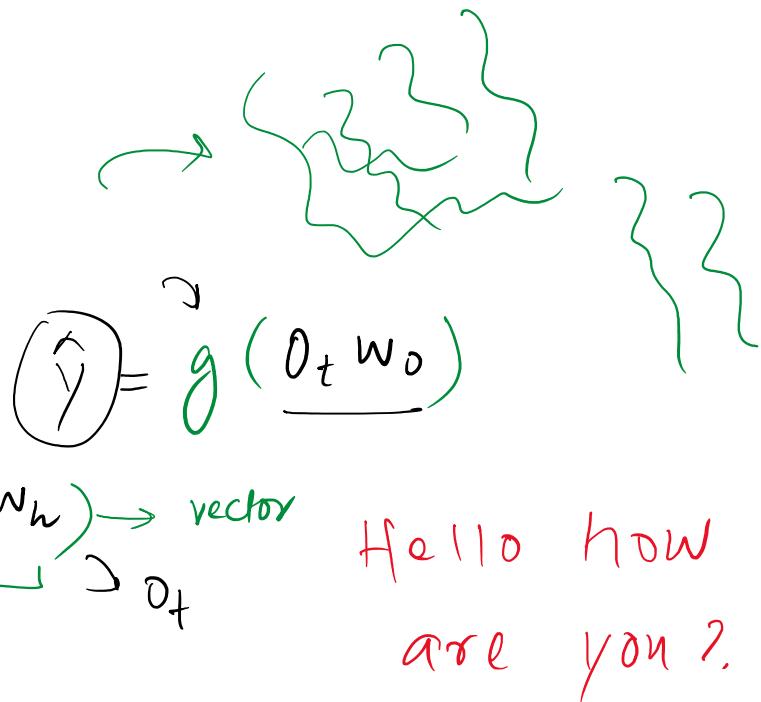
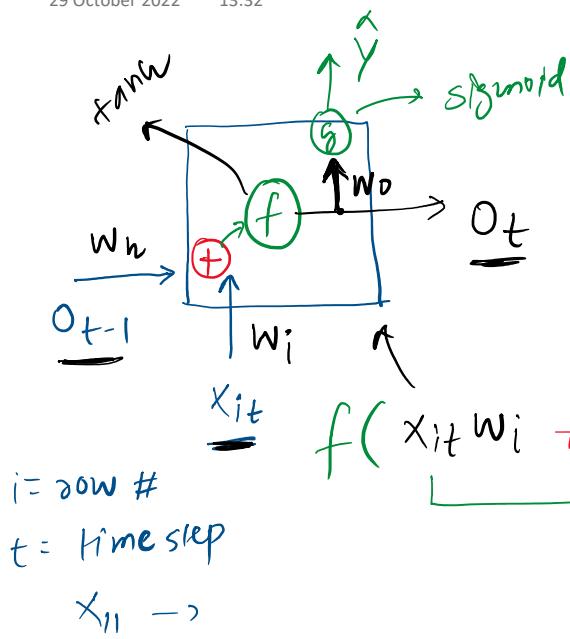
## RNN Forward Prop

29 October 2022 17:02



## Simplified Representation

29 October 2022 13:32

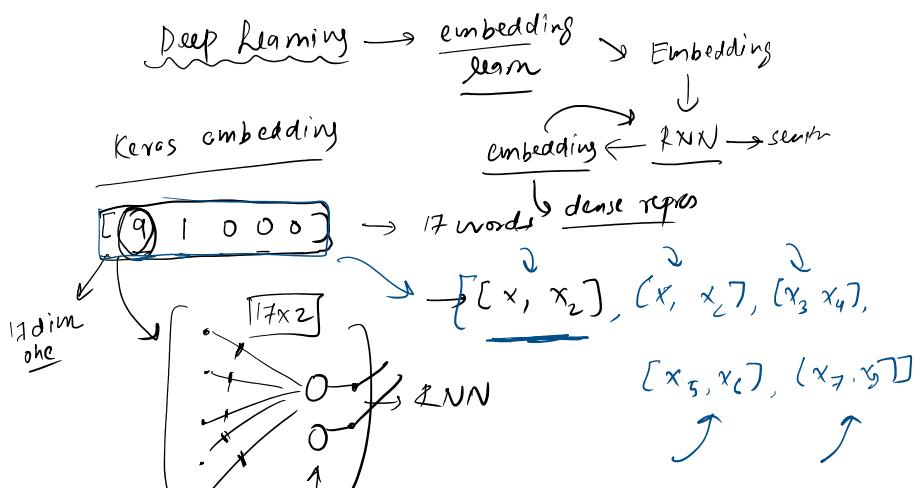
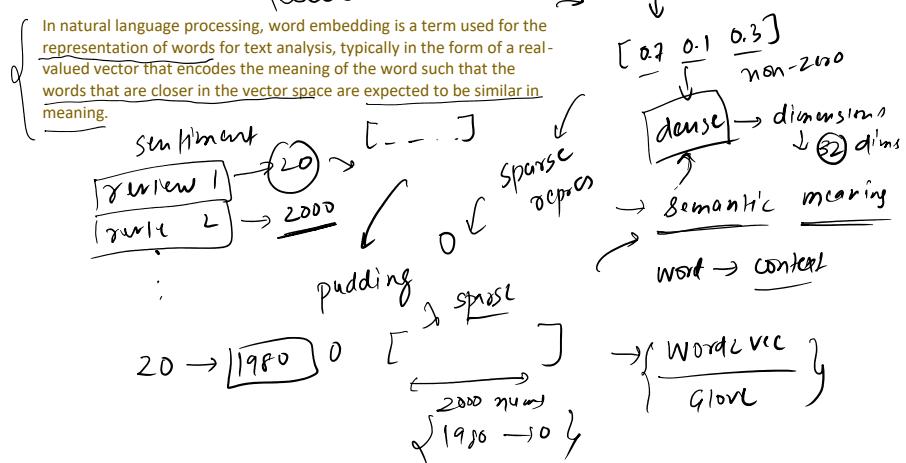
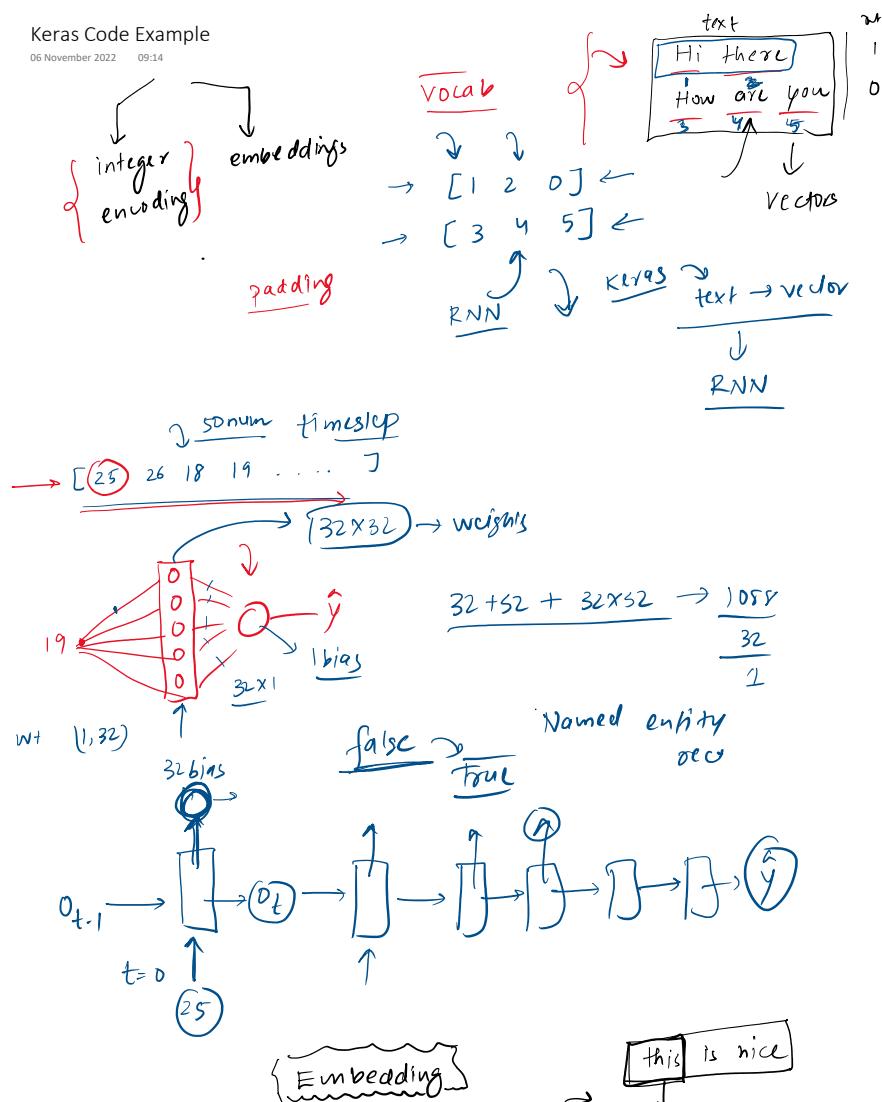


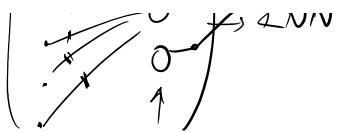
# Code

29 October 2022 13:32

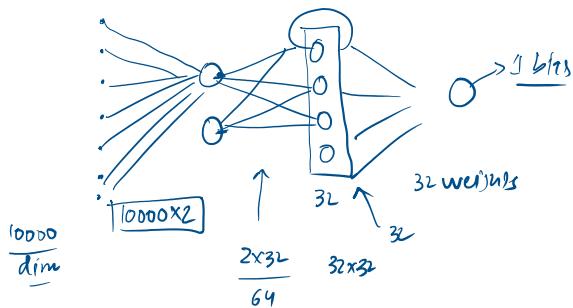
# State and Memory

29 October 2022 13:33





17 nodes



population in  $n$ th year  $\rightarrow x$

$$x + \frac{10\% \text{ of } x}{= 10000} = \underline{\underline{10000}} \quad ((n-1))$$

$$\rightarrow x + 0.1x = 10000$$

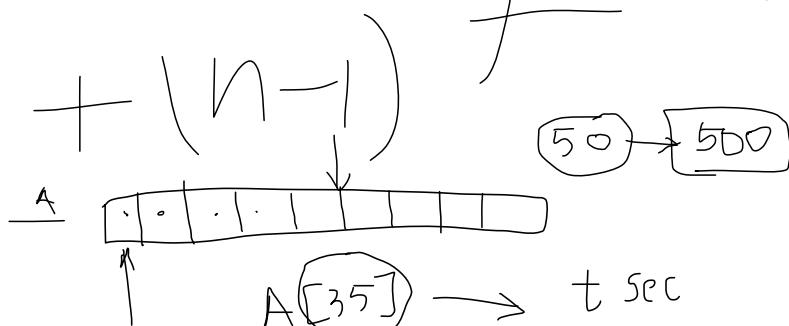
$$\frac{1.1x}{= 10000}$$

$$x = \boxed{10000}$$

$$1.1$$

$$\frac{x-1}{x} + \frac{1}{2} \left( \frac{x-1}{x} \right)^2 + \frac{1}{2} \left( \frac{x-1}{x} \right)^3 + \frac{1}{2} \left( \frac{x-1}{x} \right)^4 + \dots$$

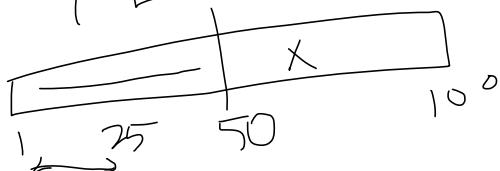
2+



$A[35] \rightarrow t \text{ sec}$

$A[35] \rightarrow$

$\boxed{35} \rightarrow \boxed{1 \times 4 \times 35}$



$O(n)$

$O(n^2) \rightarrow \text{nested loops}$

input  $\rightarrow 10 \text{ loops} \times 10 \text{ loops}$

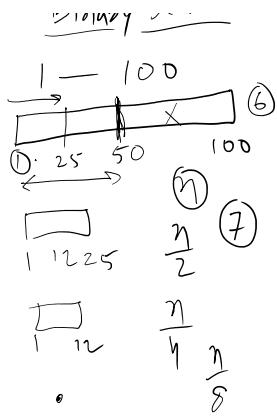
time  $\rightarrow (n)^2 \rightarrow n^2$

$O(n^2)$

$\sqrt{(x)}$

Binary Search

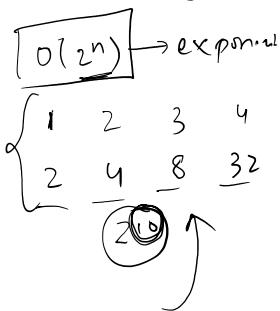
$1 \rightarrow 100$



$1 \quad 100$

$\overline{)1225}$

          
Sorting



{  
 for i in range  
 \_\_\_\_\_ O(n)  
 for j in range  
 \_\_\_\_\_ O(n)

$O(n+n) \rightarrow$

$O(2n)$

$\rightarrow O(n)$

$O(n + n^2) \quad O(n^2)$

for i in range

    for j in

    0 1 2 3 4 5 6 7 8 9 10

    25 → '25'

str()

    n = 345 % 10

digits[5]

    5 + ''  
 \_\_\_\_\_ '5'

    345 // 10 → 34

34 // 10 → 3

    1 1 1 5 1 .

$$\underline{34 \cdot 110} \rightarrow (4)$$

$$4 + \underline{15} \rightarrow \\ 45$$

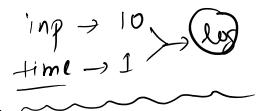
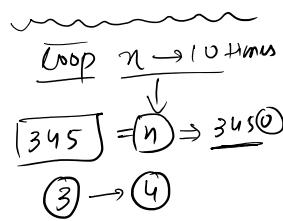
$$\underline{34 \cdot 110} \rightarrow (3)$$

$$3^1 = 0$$

$$3^1 \cdot 10 \Rightarrow 3$$

`digits[3]`

$$3 + '45' \\ = '345'$$



$$\begin{matrix} O(n) \\ O(n) \end{matrix} \rightarrow O(n+n) \\ O(n) \quad O(2n) \rightarrow O(n)$$

$$\begin{matrix} O(n) \\ O(n) \end{matrix} \rightarrow O(n \times n) \\ O(n) \quad O(n^2)$$

$$O(n)$$

$$O(n)$$

$$O(1000000)$$

$$n^2 \boxed{1000000}$$

$\times$

$$O(n^2)$$

$$1 \rightarrow \left(\frac{n}{2}\right)^{\frac{O(n \times L)}{2}}$$

$$4 \frac{n}{2} \quad 2n$$

$$\overbrace{O(n)}^{n=100}$$

$$150, 100 \boxed{\frac{n}{2}}$$

$$2 \rightarrow \boxed{100=n} \boxed{\frac{n}{2}}$$

$$j=1 \rightarrow 2$$

$$j=2 \rightarrow 4 \boxed{2-100}$$

$$j=3 \rightarrow 8$$

$$j=4 \rightarrow 32$$

$$\frac{n}{2} \times \log n \\ \boxed{O(n \log n)}$$

$O(1)$  → constant

$$n = \boxed{345}$$

$$3+4+5 \rightarrow 12$$

(5)

$$\begin{array}{r} 345 \\ 3 \quad 3 \\ \hline 10 \end{array}$$

$$3450 \rightarrow 4$$

log

$$\begin{array}{l} \text{inp} \rightarrow 10 \quad 100 \\ \text{out} \rightarrow 1 \quad 2 \quad 3 \end{array}$$

fibonacci

↓ function  
↳ recursion

$$\begin{array}{ccccccc} 0 & 1 & 1 & 2 & (3) & (5) & (8) \\ \downarrow & \uparrow & \downarrow & \uparrow & \downarrow & \uparrow & \uparrow \\ n=1 & n=0 & \rightarrow (1) & & & & \end{array}$$

fib(n)

↓  
function calls

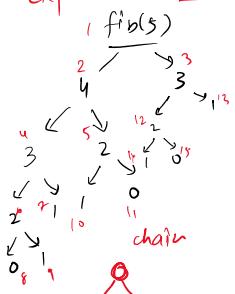
$$\text{input} \rightarrow (2) \rightarrow 10$$

#fcalls →

$$\text{fib}((3)) \xrightarrow{n=3} (1)$$

$$\begin{array}{c} \swarrow \quad \searrow \\ \text{fib}(2) \quad \text{fib}(1) \xrightarrow{n=1} 1 \\ \swarrow \quad \searrow \\ \text{fib}(1) \quad \text{fib}(0) \\ \hline 1 \quad 0 \end{array}$$

exponential → bad



$$\begin{array}{l} O(2^n) \quad 15 \quad 20 \\ \text{input} \quad 1 \quad 2 \quad 3 \quad 4 \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ 1 \quad 2 \quad 4 \quad 8 \quad 16 \end{array}$$

$$O(2^n)$$

input	1	2	3	4
$+ 1$	2	4	8	16

$$\boxed{n=50, 100, 500}$$

$\uparrow$  weeks  
exponential  
 $\hookrightarrow$  days/weeks

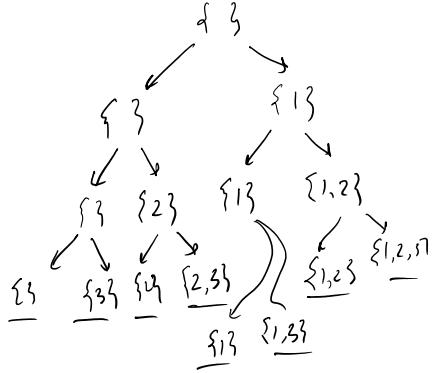
subset  $\rightarrow O(?)$   
power set

$$\{1, 2, 3 \rightarrow \{\{3, \{1\}, \{2\}, \{1, 2\}\}$$

$$\{1, 2, 3\} \rightarrow \{3, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{2, 3\}, \{1, 2, 3\}$$

$$\{3, \{1\}, \{2\}, \{1, 2\}, \{2, 3\}, \{1, 2, 3\}$$

$$\{1, 2, 3\}$$



reduce  $\rightarrow$  divide  $\rightarrow$  log

increase  $\rightarrow$  multi  $\rightarrow$  exp

$\rightarrow$  exponentiation

$$\{1, 2\} \rightarrow \textcircled{4} \quad 2^2 = \underline{4}$$

$$\{1, 2, 3\} \rightarrow \textcircled{8} \quad 2^3 = \underline{8}$$

$$\{1, 2, 3\} \rightarrow \textcircled{8} \quad 2^3 = \underline{8}$$

$$O(2^n)$$

$$O(?)$$

$$T(n) = \begin{cases} 3T(n-1) & \text{if } n > 0 \\ 1, & \text{otherwise} \end{cases}$$

$$n > 0$$

$$T(n) = \underline{3T(n-1)}$$

$$= 3[\underline{3T(n-2)}]$$

$$= 3^2 \underline{T(n-2)}$$

$$= 3^2 [3\underline{T(n-3)}]$$

$$= 3^3 T(n-3)$$

$$= 3^n T(n-n)$$

$$= 3^n \underline{T(0)}$$

$$T(n) = \boxed{3^n} \rightarrow O(3^n)$$

$$T(n) = \begin{cases} 2T(n-1)-1 & \text{if } n>0 \\ 1, \text{ otherwise} & \rightarrow \text{constant} \end{cases}$$

$$T(n) = \underline{2T(n-1)-1}$$

$$= 2[2T(n-2)-1]-1$$

$$= 2^2 \underline{T(n-2)-2}-1$$

$$= 2^2 [2T(n-3)-1]-2-1$$

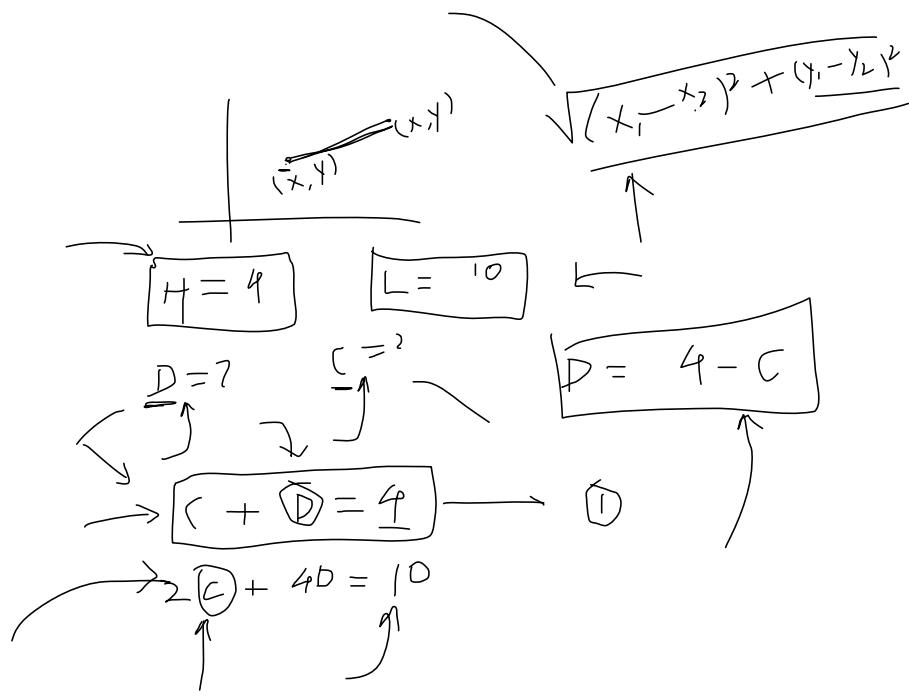
$$= 2^3 T(n-3) - 2^2 - 2^1 - 2^0$$

$$= \underline{2^n T(n-n)} - 2^{n-1} - 2^{n-2} - \dots - 2^1 - 2^0$$

$$= 2^n - [2^{n-1} + 2^{n-2} + \dots + 2^1 + 2^0]$$

$$= 2^n - [2^n - 1] = 2^n - 2^n + 1$$

$O(1) \rightarrow \text{constant}$



15    5     $15^2 + 5^2$

$\leq$

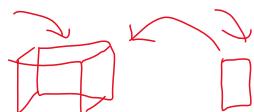
$$17 \times 15 = 2^2 + 3^2 + 4^2 \times 5^n \quad n=5$$

$$26 \quad [3, 6] \rightarrow 5^{+n}$$

$$a=3 \quad n=5$$

$$d=6-3$$

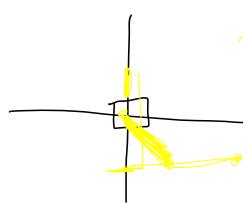
$$\begin{array}{r} 1000 \\ , 060 \\ \hline 2 \end{array} \frac{2}{3} \cancel{+} \frac{4}{5} = \boxed{\frac{10+12}{15}} \rightarrow \frac{22}{15}$$



$$0 \quad 1 \quad 1 \quad 2 \quad 2 \quad 5$$

$$\boxed{1000}$$

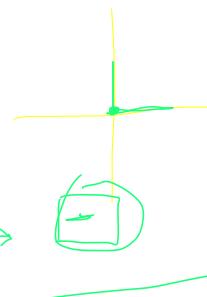
$$\# \quad 1002 \quad \boxed{2222} \quad 51 = \overbrace{5}^{5 \times 10 \times 1 \times 2 \times 1}$$



$$4P \rightarrow 1 \\ D - 3 \\ R - 4 \\ L - 3$$

$$a=[1, 2] \\ a=b \\ a=b[;]$$

$$q \rightarrow$$

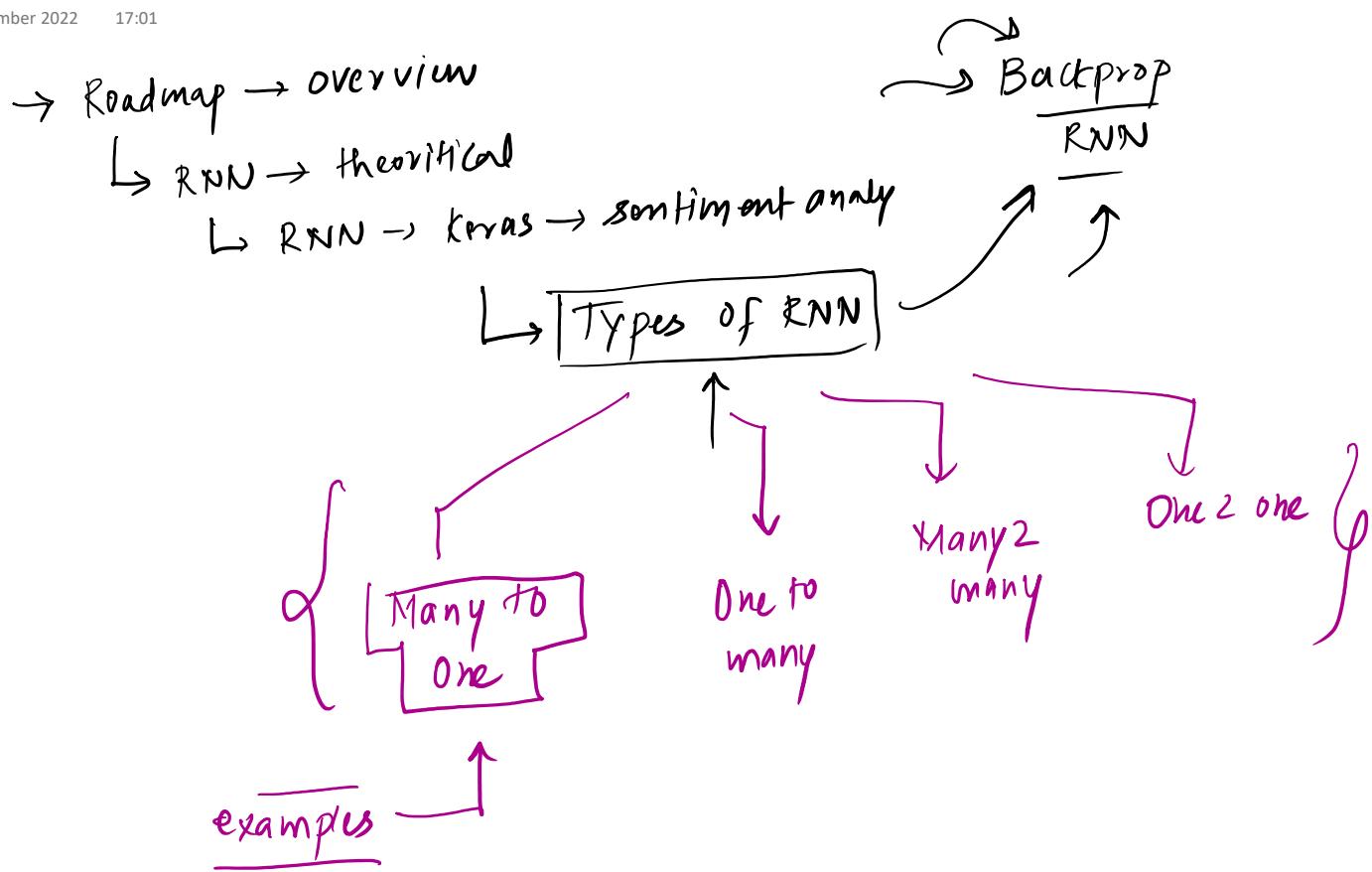


$$b \rightarrow$$

$$\boxed{0}$$

## Till Now

17 November 2022 17:01

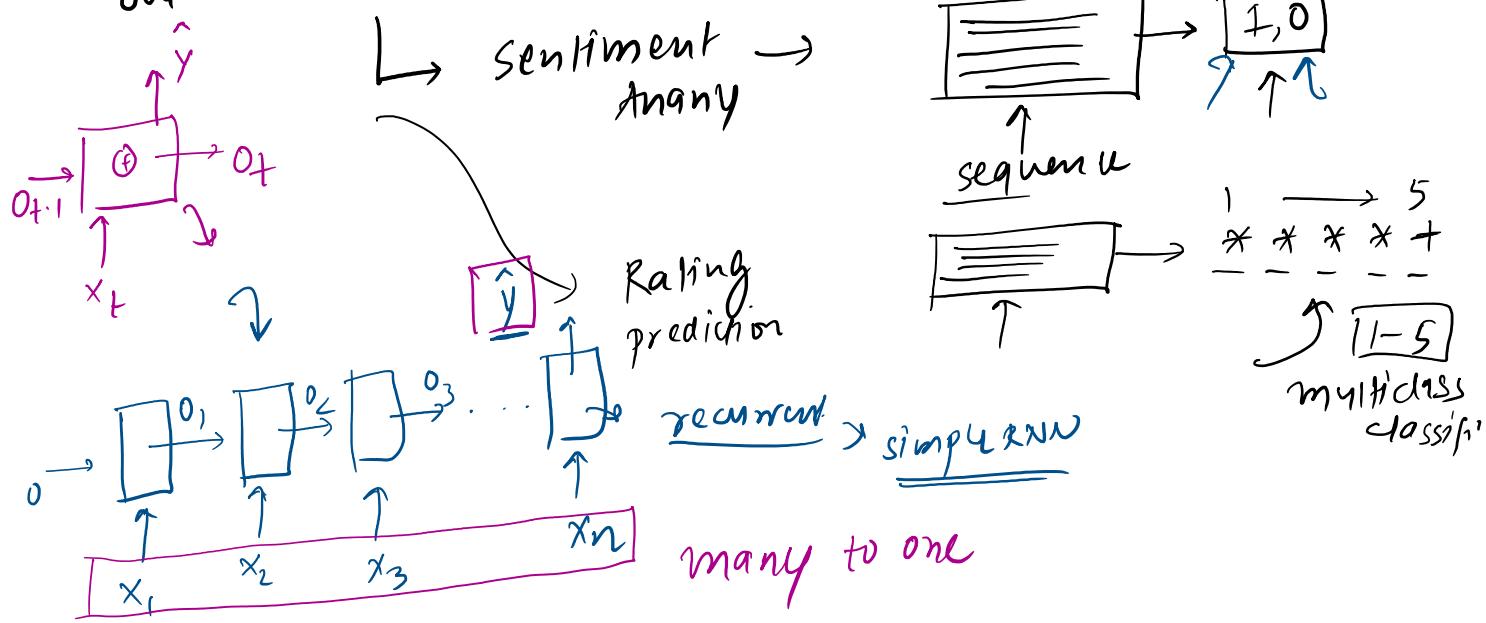


## Many to One

17 November 2022 17:02

ember 2022 17:02 → scenario → sent, change timelines

input → sequence →  $\text{seq}_1, \dots, \text{seq}_n$   
out → non-seq → int/num → scalars - (1, 0)

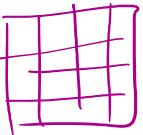


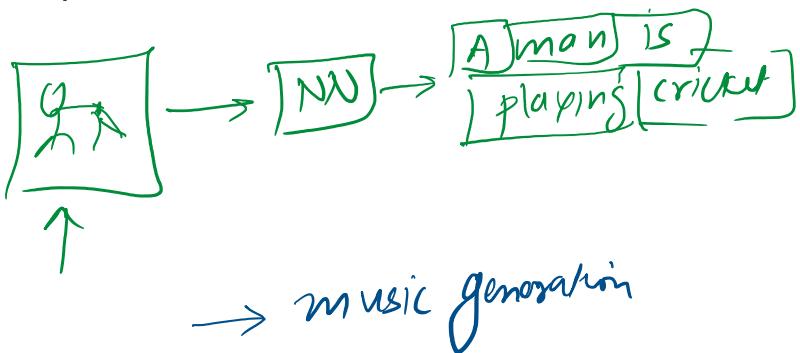
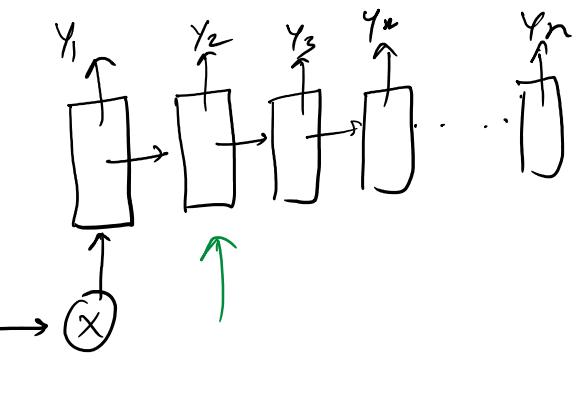
## One to Many

17 November 2022 17:02

→ normal non sequential  
↳ 

→ Output → sequences  
image captioning

 → textual  
depres.



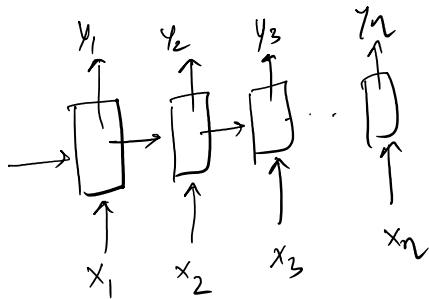
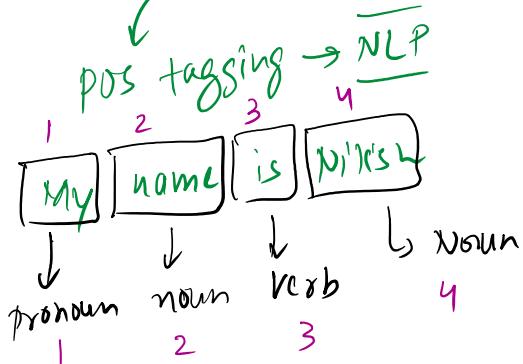
## Many to Many

17 November 2022 17:02

input  $\rightarrow$  segment  $\rightarrow$  seq2seq  
 out  $\rightarrow$  sequence

Variable length

input seq == output seq



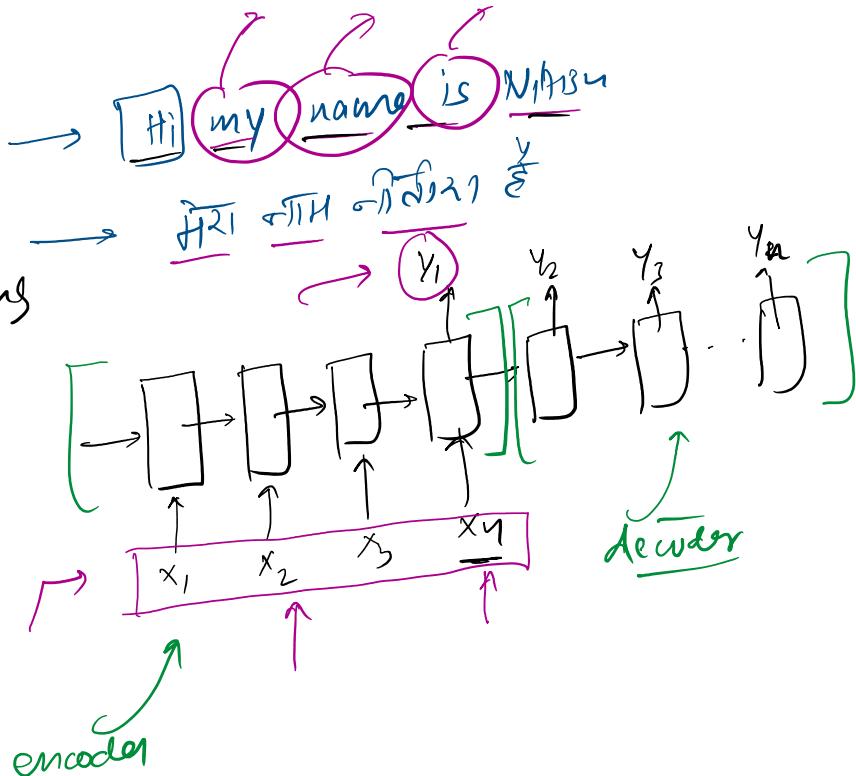
same length  
 many  $\rightarrow$  many  
 RNN

NER  
 Lets meet at 7pm at the airport

Variable length  
 machine trans  
 $\hookrightarrow$  1 lang  $\rightarrow$  2 lang

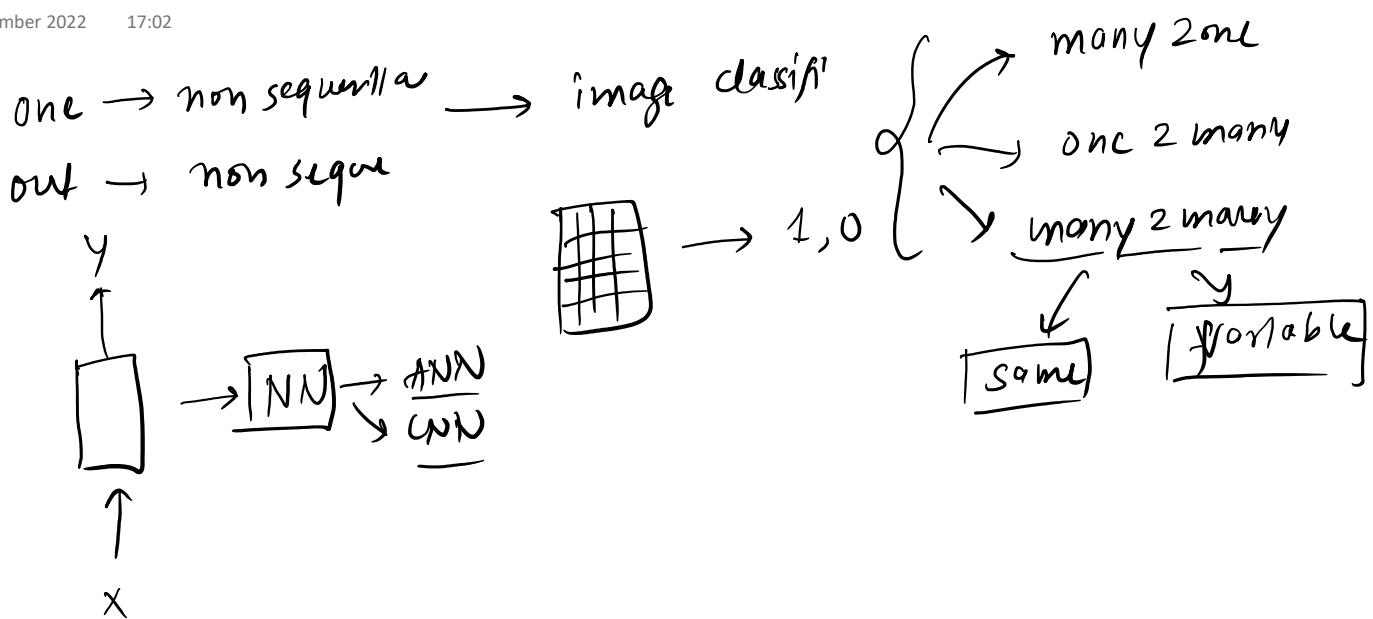
google translate

encoder  
 decoder



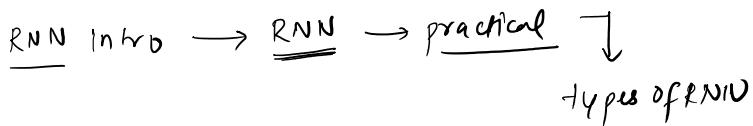
## One to One

17 November 2022 17:02



## Backpropagation in RNN

01 December 2022 16:43

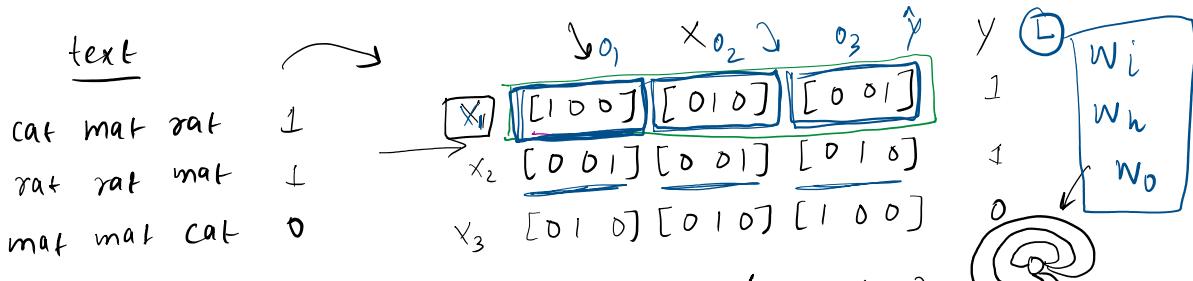


Many to One RNN

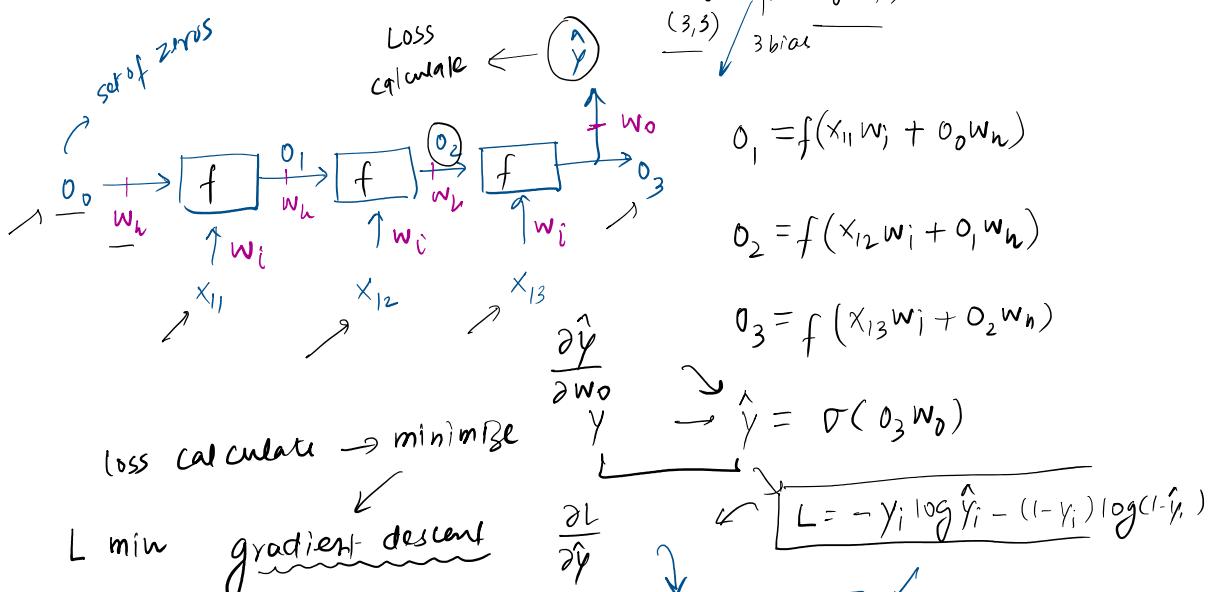
Sentiment Analysis

text → I/O

$$\hat{y} \rightarrow \mathbb{C}$$



forward prop



$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_0}$

$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_i}$

$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial w_i}$

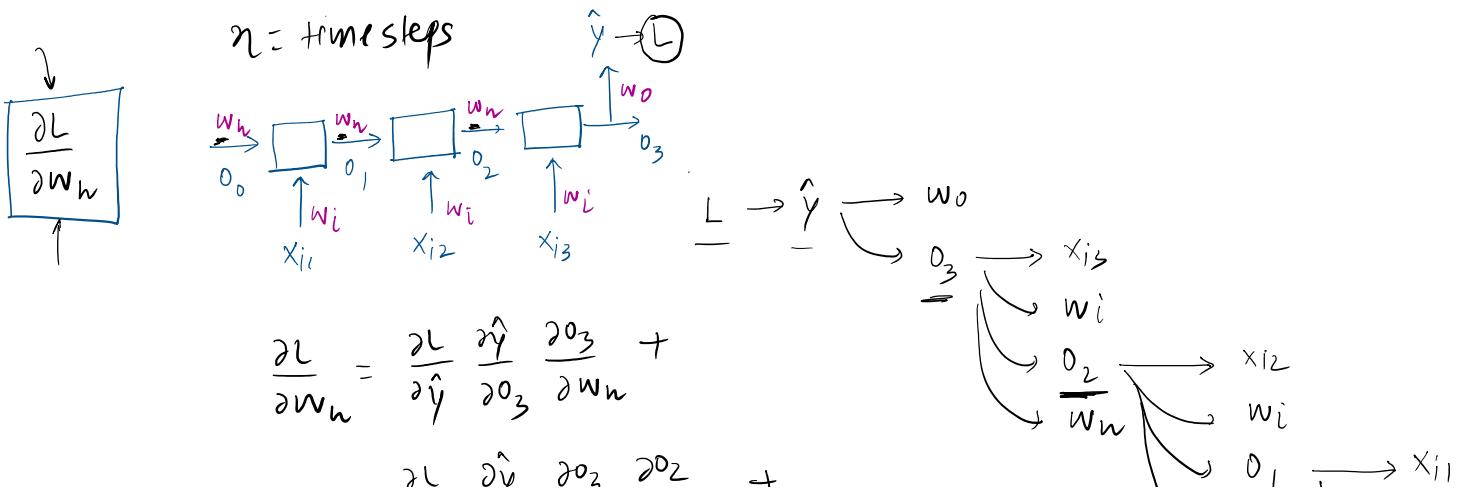
$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial w_i}$

$\frac{\partial L}{\partial w_i} = \sum_{j=1}^3 \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_j} \frac{\partial o_j}{\partial w_i}$

$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \left[ \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial w_i} \right] + \left[ \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_i} \right]$

$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \left[ \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial w_i} \right]$

$\frac{\partial L}{\partial w_i} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_j} \frac{\partial o_j}{\partial w_i}$

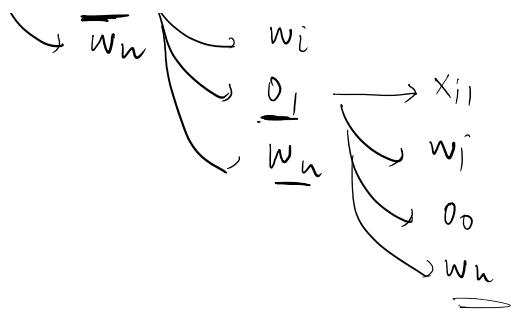


$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial w_h} +$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial w_h}$$

$$\boxed{\frac{\partial L}{\partial w_h} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_j} \frac{\partial o_j}{\partial w_h}}$$

$\eta = \text{timesteps}$



for  $j=3$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_h} \rightarrow \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial o_1} \frac{\partial o_1}{\partial w_h}$$

for  $j=10$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_{10}} \frac{\partial o_{10}}{\partial w_h} \frac{\partial o_t}{\partial o_{t-1}}$$

$j$	$\frac{\partial o_t}{\partial o_{t-1}}$
$t=2$	$\vdots$

$$\frac{\partial o_t}{\partial o_{t-1}} = \frac{\partial o_2}{\partial o_1} \frac{\partial o_3}{\partial o_2}$$

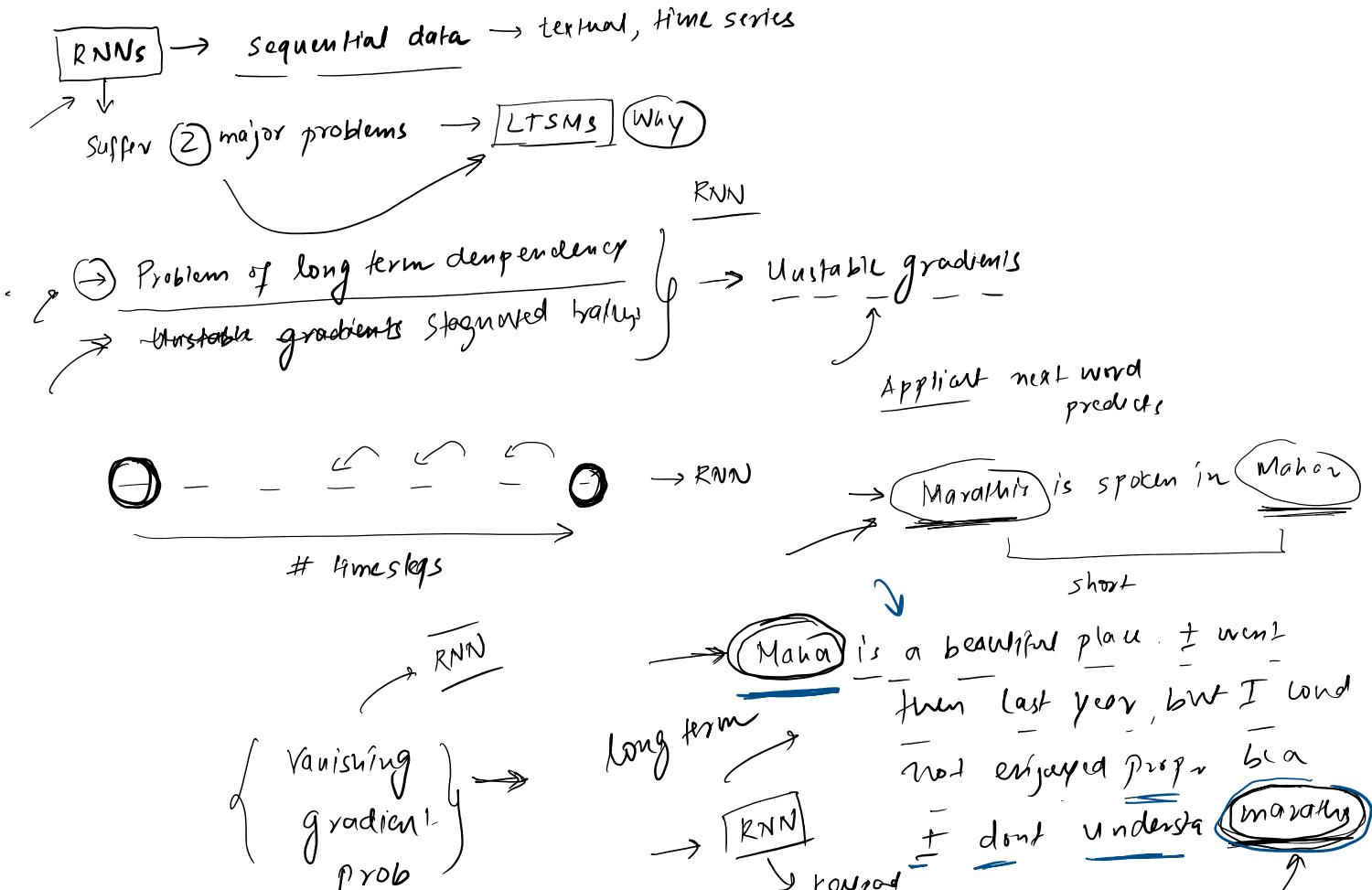
$$o_t = f(x_{it} w_{inp} + o_{t-1} w_h)$$

$$\frac{\partial o_t}{\partial o_{t-1}} = \frac{\partial o_t}{\partial f'(x_{it} w_{inp} + o_{t-1} w_h) w_h}$$

$\uparrow \quad \downarrow$   
[0-1]

## Problem with RNN

19 December 2022 16:33



### Problem #1 → Problem of long term dependency → Vanishing

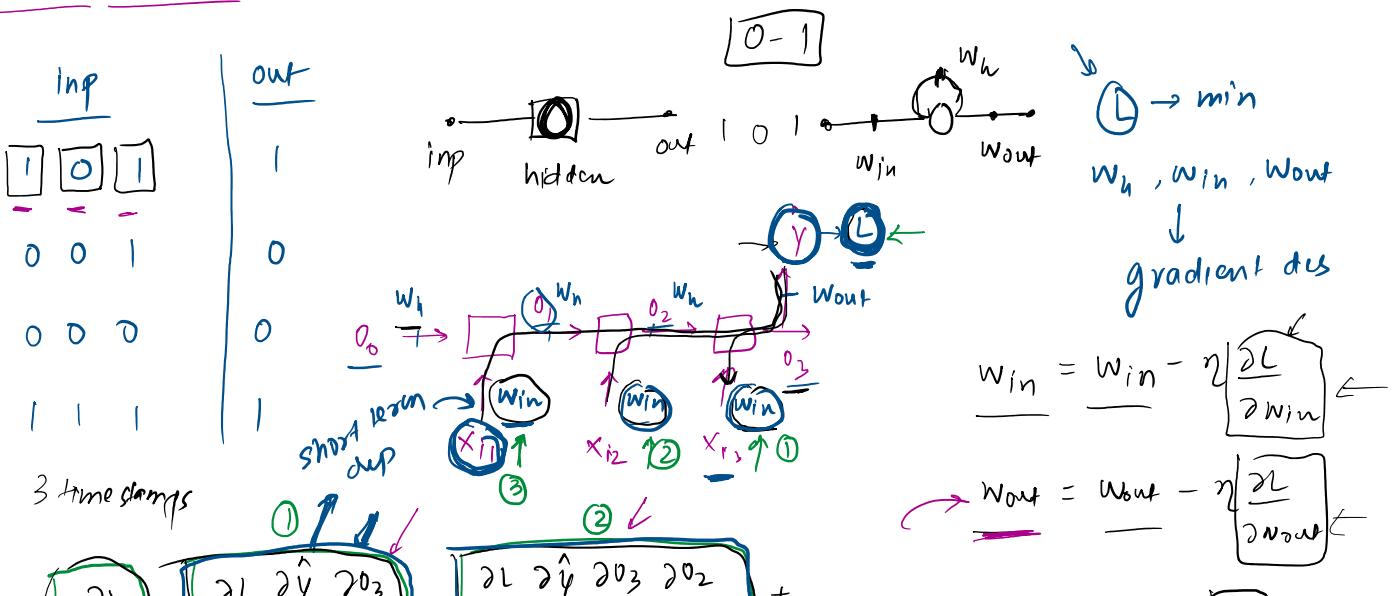


Diagram illustrating the backpropagation through time (BPTT) for an LSTM cell. The diagram shows the computation of gradients for hidden states  $h_t$  and cell states  $c_t$  over multiple time steps.

**Top Left:** A diagram showing the gradient flow from the loss function  $L$  through the hidden states  $h_1, h_2, \dots, h_T$  and cell states  $c_1, c_2, \dots, c_T$ . The diagram highlights the forget gate's role in determining the cell state update. Labels include  $\frac{\partial L}{\partial h_i}$ ,  $\frac{\partial h_i}{\partial c_j}$ , and  $\frac{\partial c_i}{\partial w_{in}}$ .

**Top Right:** A diagram showing the update rule for the hidden state weight  $w_h$  at time step  $t$ :

$$w_h = w_h - \eta \frac{\partial L}{\partial h_t}$$

**Middle Left:** A diagram showing the gradient flow for the hidden state  $h_t$  over time steps  $t=2$  to  $T$ . It highlights the "long term dep" (dependency) of the gradient flow.

**Middle Right:** A diagram showing the gradient flow for the hidden state  $h_t$  over time steps  $t=2$  to  $T$  using the chain rule:

$$\frac{\partial L}{\partial h_t} = \prod_{t=2}^{100} \left( \frac{\partial h_t}{\partial h_{t-1}} \right) \frac{\partial h_1}{\partial w_{in}}$$

**Bottom Left:** A diagram showing the gradient flow for the hidden state  $h_t$  over time steps  $t=2$  to  $T$  using the chain rule and highlighting the vanishing gradient problem:

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(x_{it} w_{in} + o_{t-1} w_h) w_h$$

**Bottom Right:** A diagram showing the gradient flow for the hidden state  $h_t$  over time steps  $t=2$  to  $T$  using the chain rule and highlighting the vanishing gradient problem:

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(x_{it} w_{in} + o_{t-1} w_h) w_h$$

**Bottom Center:** A diagram showing the gradient flow for the hidden state  $h_t$  over time steps  $t=2$  to  $T$  using the chain rule and highlighting the vanishing gradient problem:

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(x_{it} w_{in} + o_{t-1} w_h) w_h$$

Sol ④

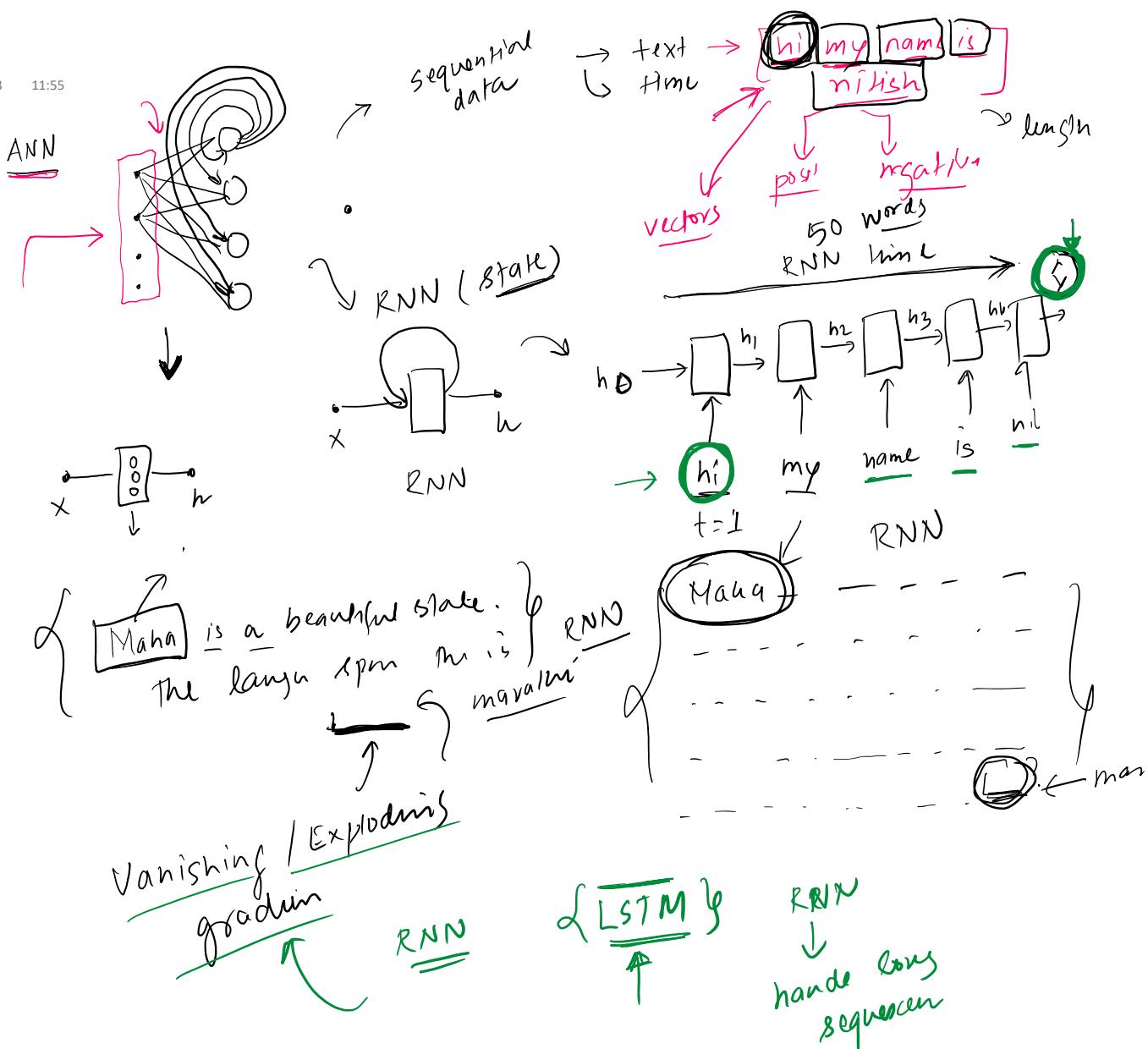
- 1) Diff activation  $\rightarrow$  relu / leaky relu
- 2) Better weight init
- 3) Skip conn
- 4) LSTM

Problem #2  $\rightarrow$  Unstable Training (Exploding gradients)

- 1) Gradient Clipping
- 2) Controlled learning rate
- 3) LSTM

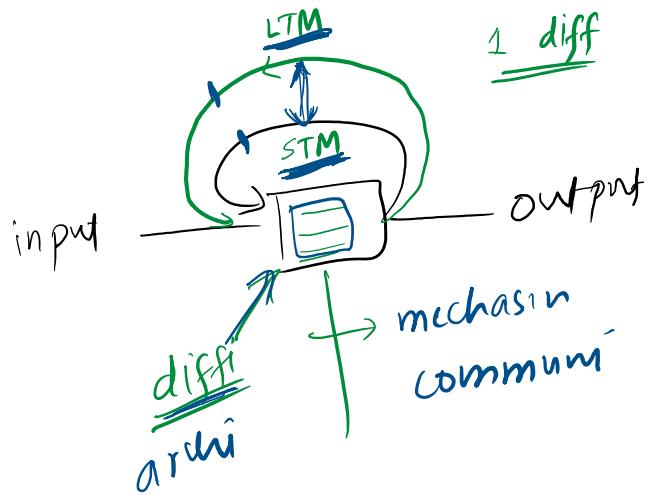
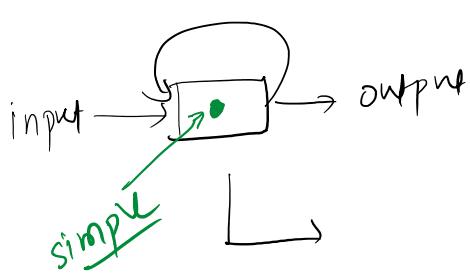
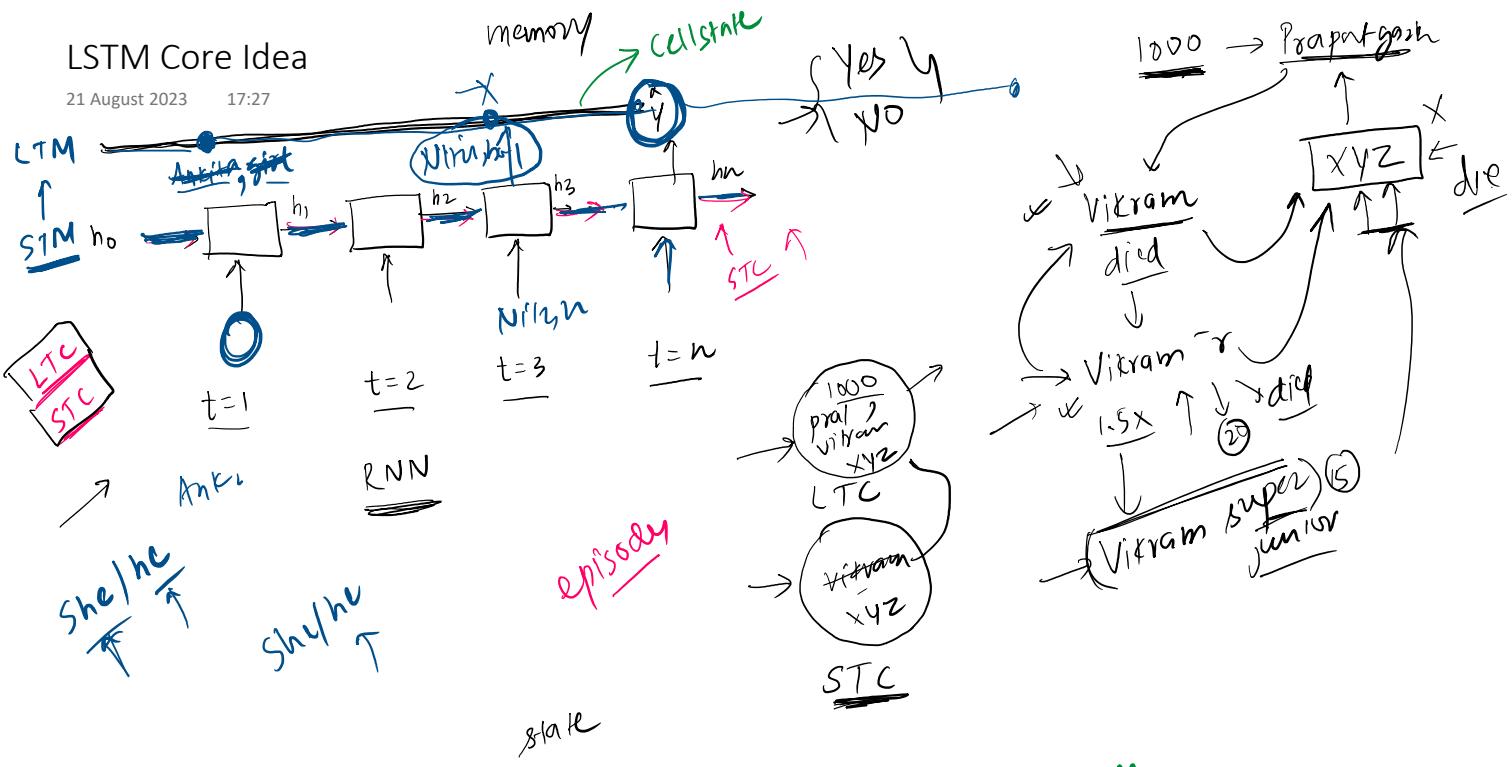
## Recap

21 August 2023 11:55



# LSTM Core Idea

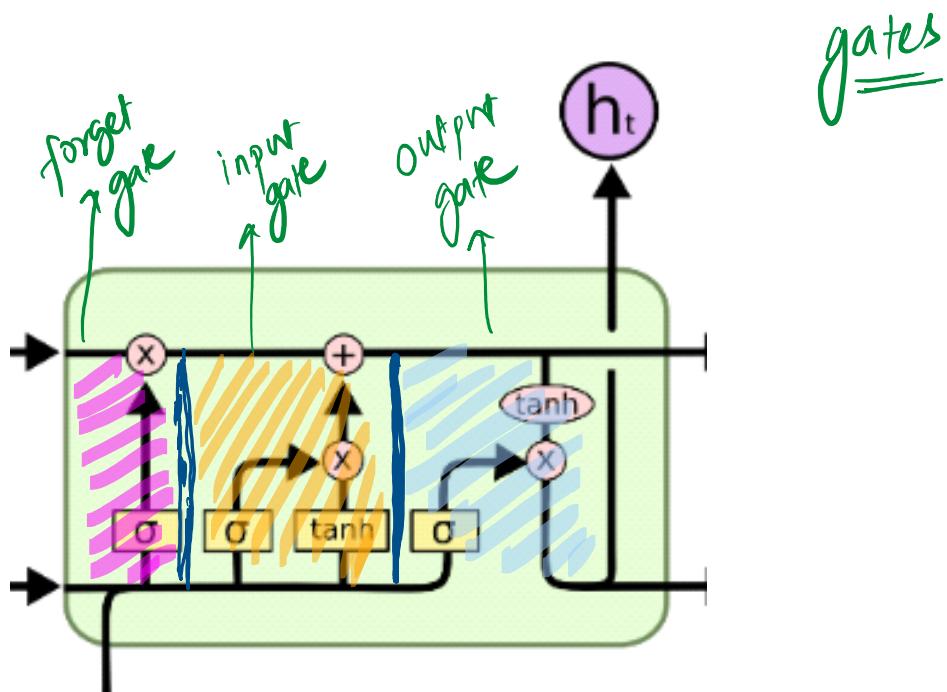
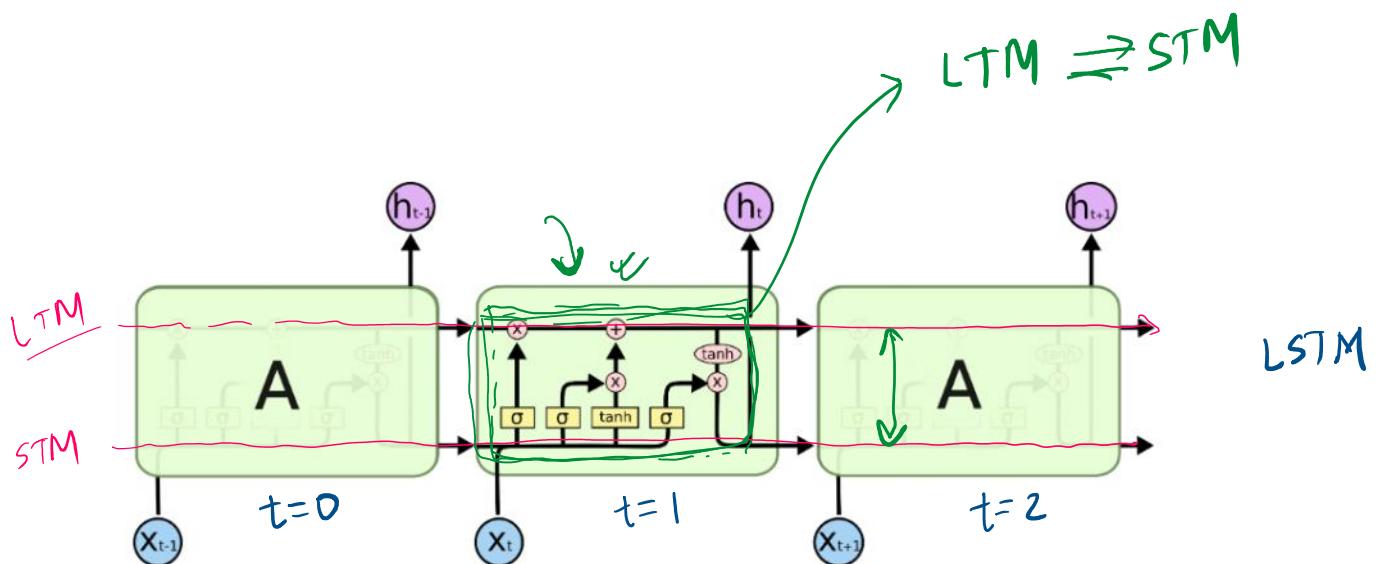
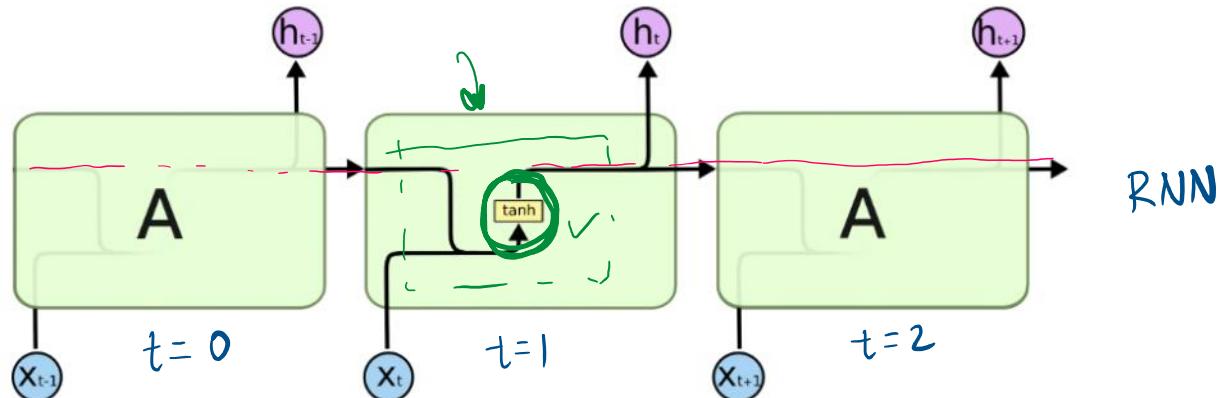
21 August 2023 17:27

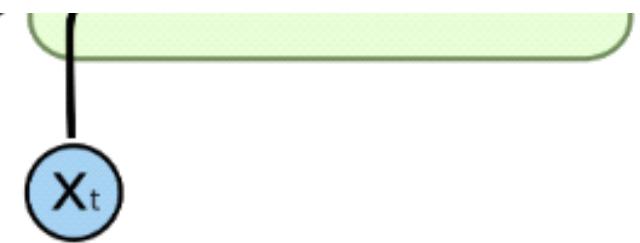


# LSTM Architecture

21 August 2023

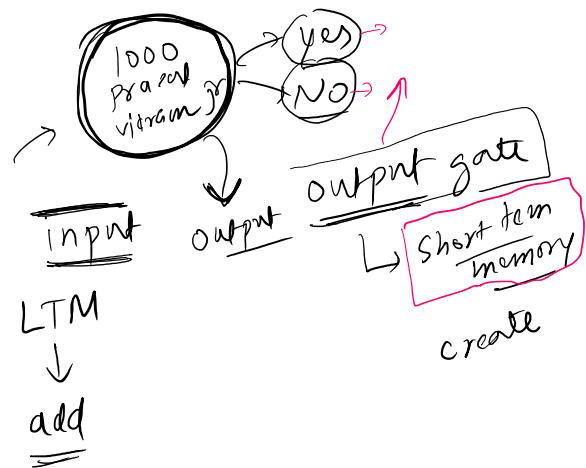
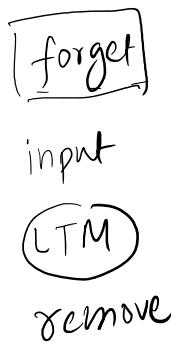
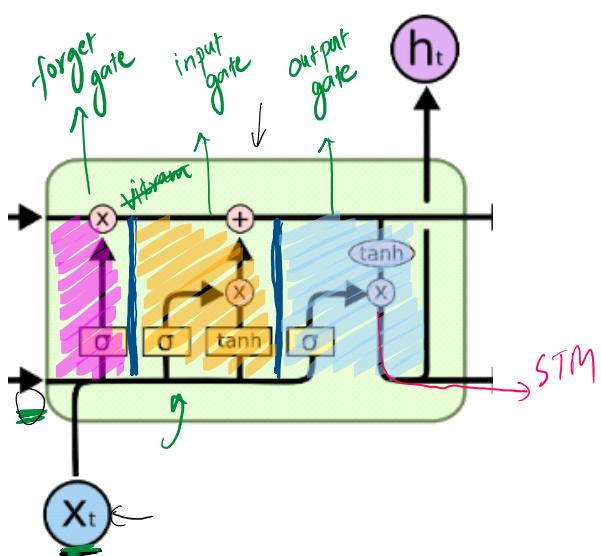
18:41





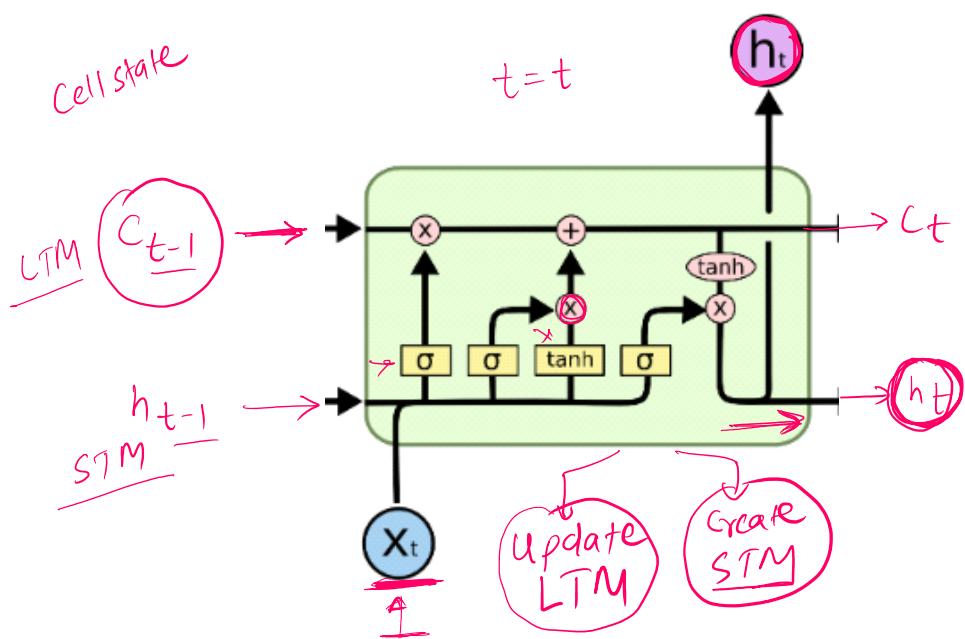
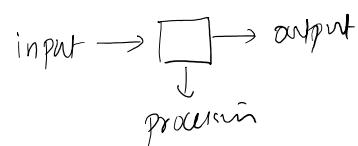
# LSTM Gates

21 August 2023 19:06



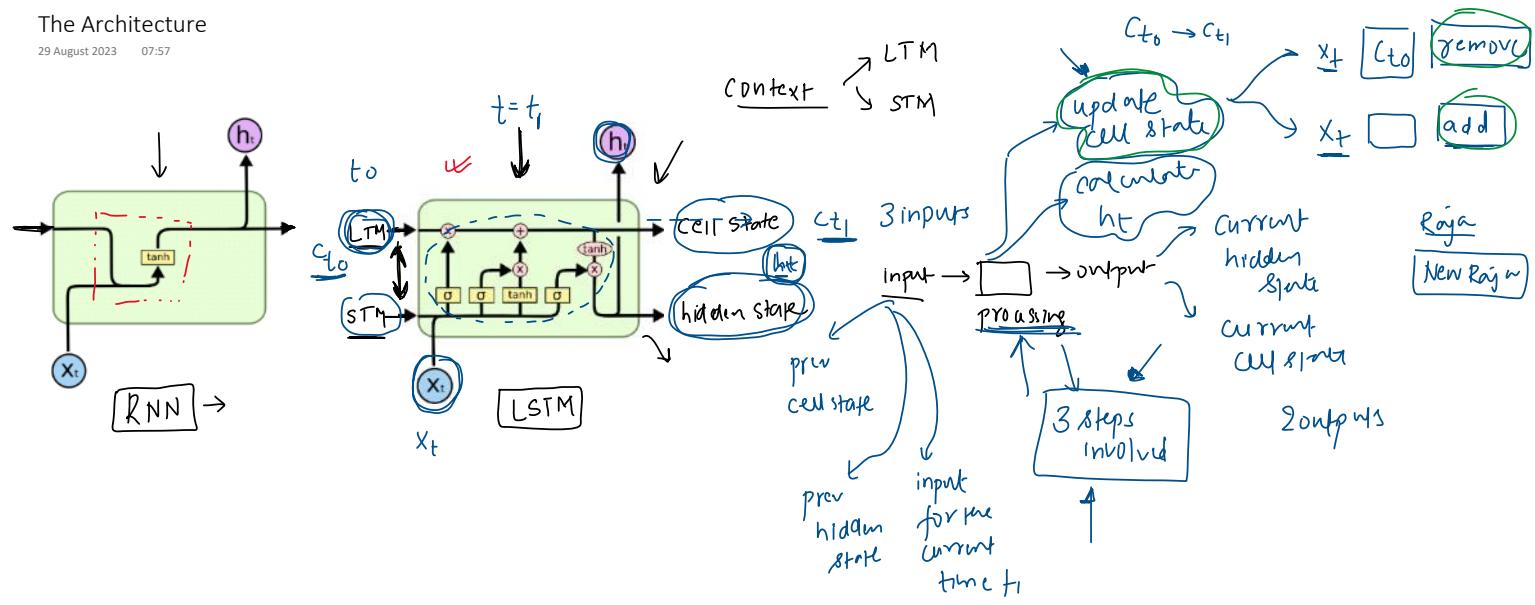
# Summary

21 August 2023 19:29



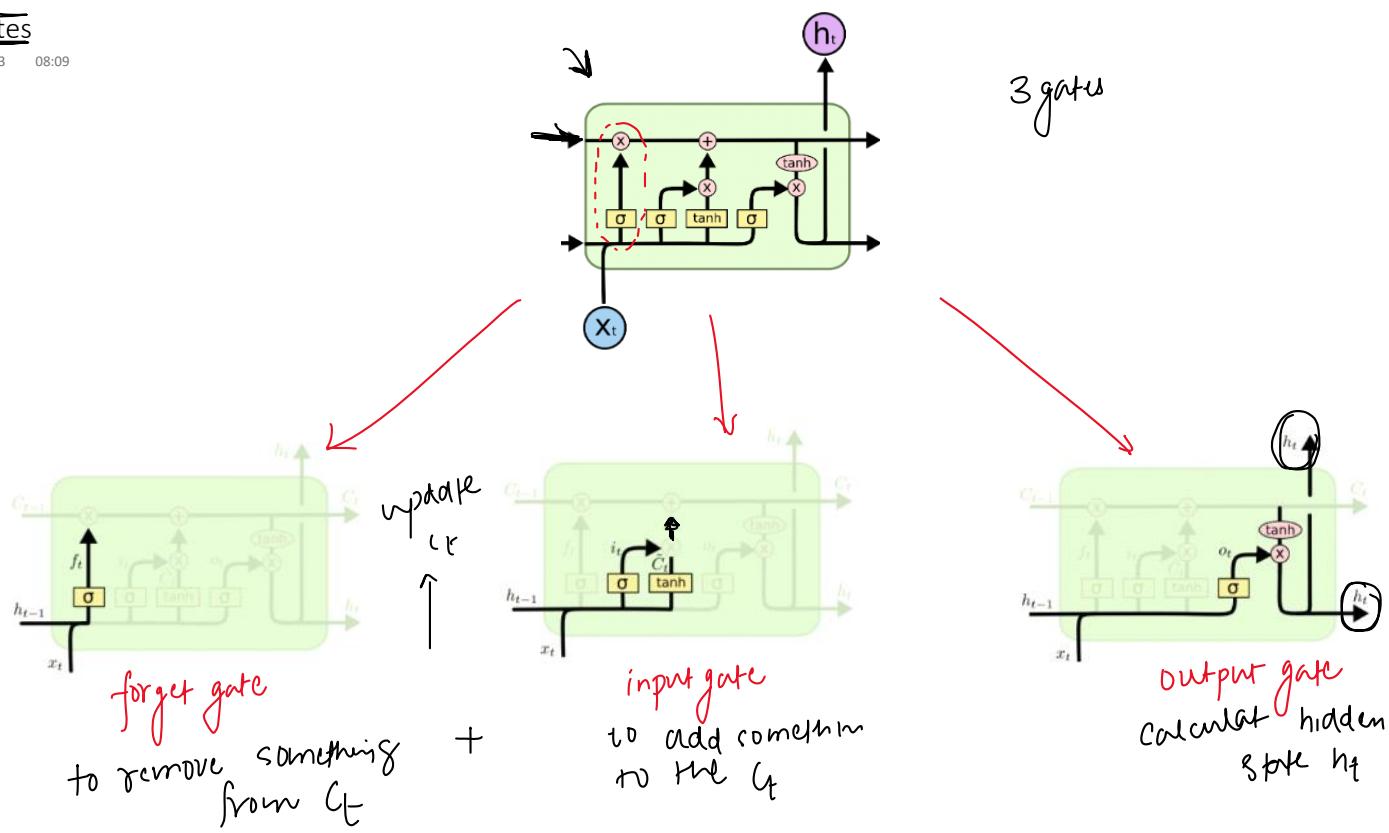
## The Architecture

29 August 2023 07:57



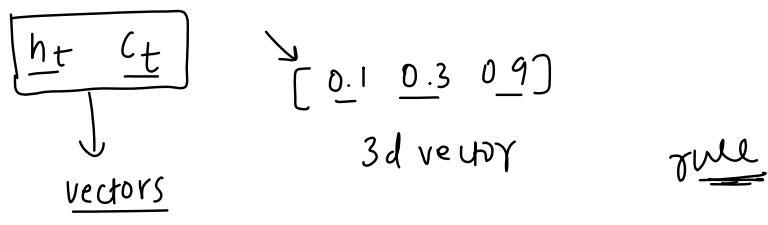
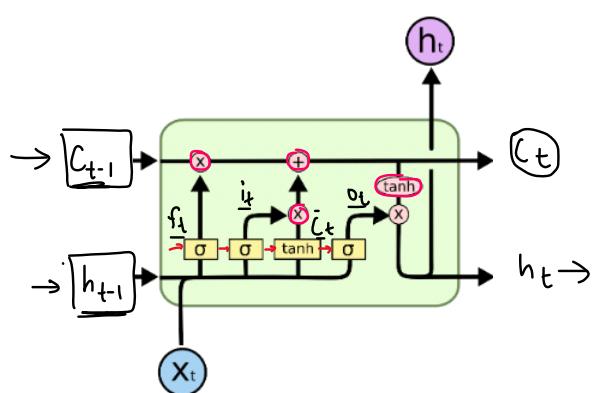
# The Gates

29 August 2023 08:09



## What are $C_t$ and $h_t$

29 August 2023 08:08



$h_t \quad C_t$  dim equal

$h_t [0.1 \quad 0.45 \quad 0.6]$

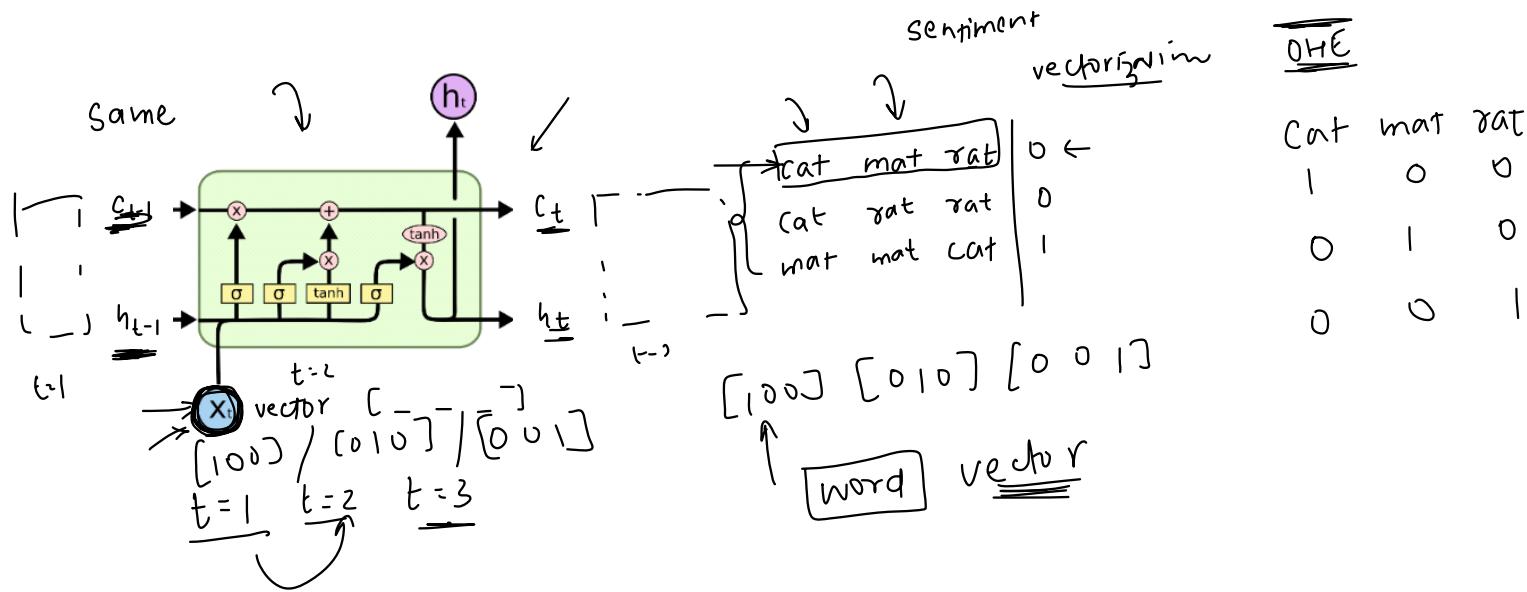
$C_t [0.55 \quad 0.6 \quad 0.0]$

same

# What is $X_t$

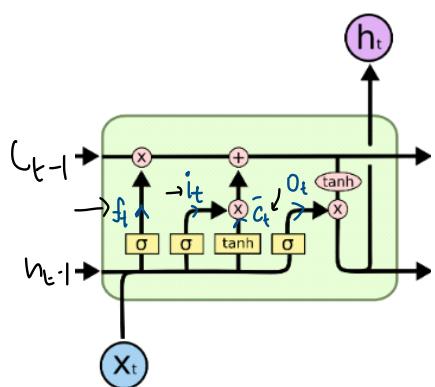
29 August 2023 17:40

RNN



What are  $f_t$ ,  $i_t$ ,  $o_t$  and  $\bar{C}_t$

29 August 2023 08:09



$f_t$  forget gate  
 $i_t$  Input gate  
 $\bar{C}_t$  candidate cell state  
 $o_t$  output gate

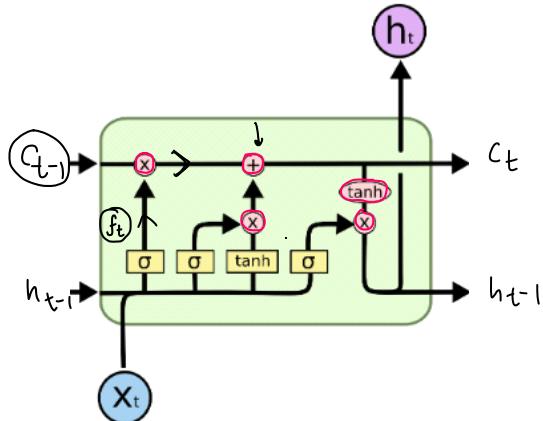
vectors

$$\begin{matrix} C_t & h_t \end{matrix}$$
$$f_t \quad i_t \quad \bar{C}_t \quad o_t$$

[ x 4 2 ]  
[ ] 7

## Pointwise Operations

29 August 2023 18:26



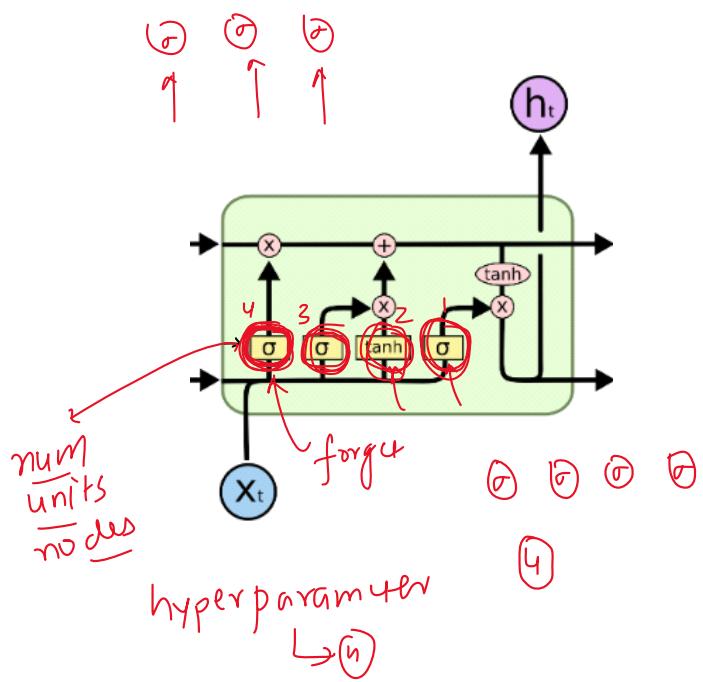
$$\begin{aligned}
 & \rightarrow \otimes \\
 & \rightarrow + \\
 & \rightarrow \text{tanh}
 \end{aligned}$$

$c_{t-1} = \begin{bmatrix} 4 & 5 & 6 \\ 1 & 2 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 0.26 & 0.34 & 0.53 \end{bmatrix}$   
 $\text{tanh}(4) = [5 \ 7 \ 9]$

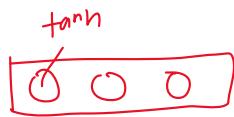
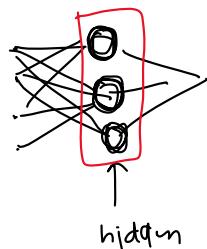
$f_t = \underline{\text{shape(dim)}} \quad \underline{\text{vector}}$   
 $c_{t-1} \otimes f_t \rightarrow \text{vector} \rightarrow [n \ 10 \ 18]$

## → Neural Network Layers

29 August 2023 18:34

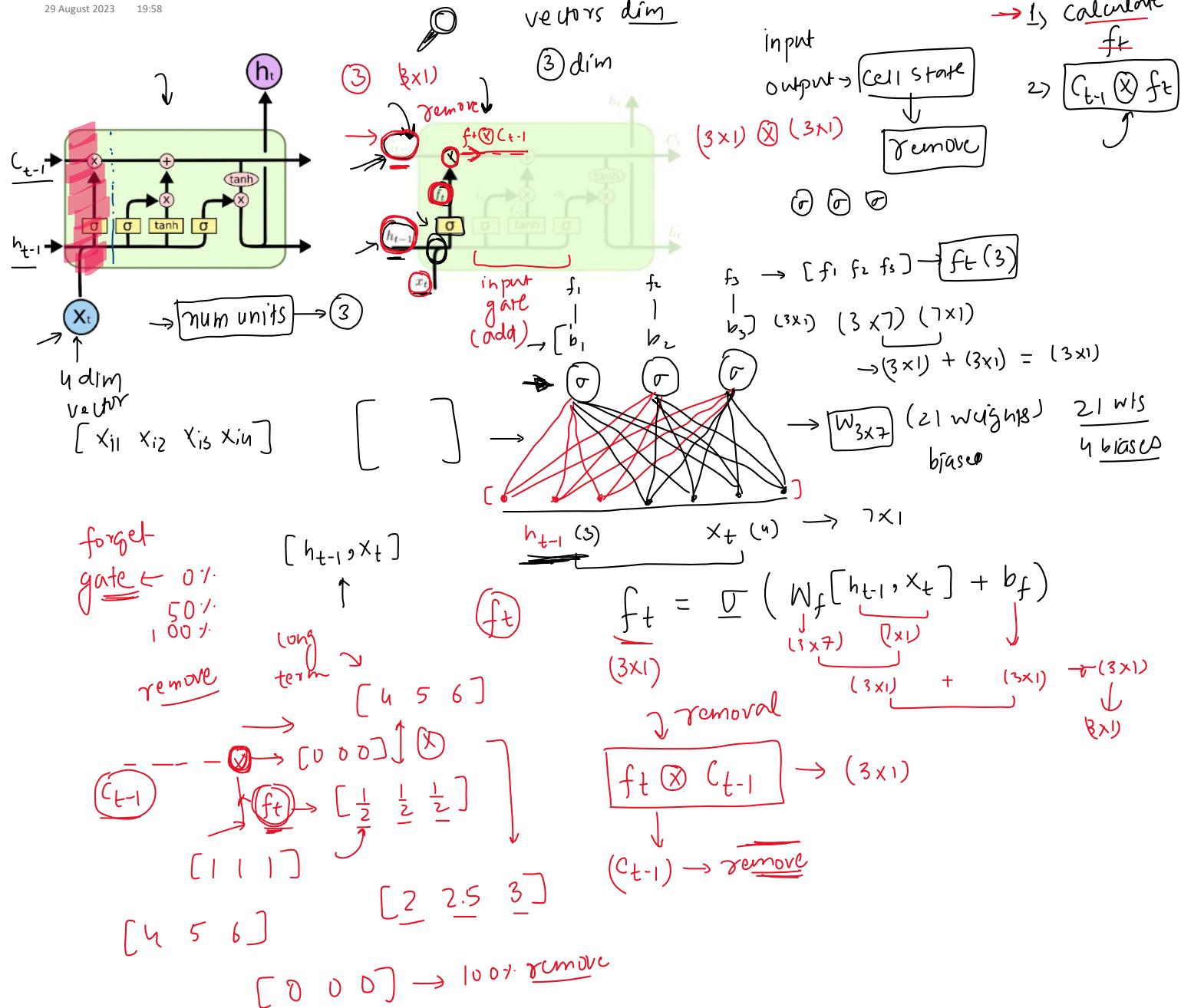


ANN



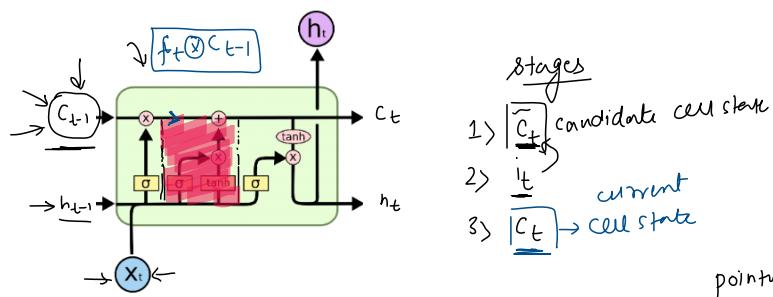
## The Forget Gate

29 August 2023 19:58

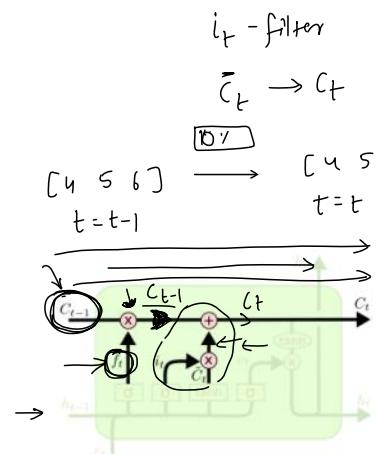
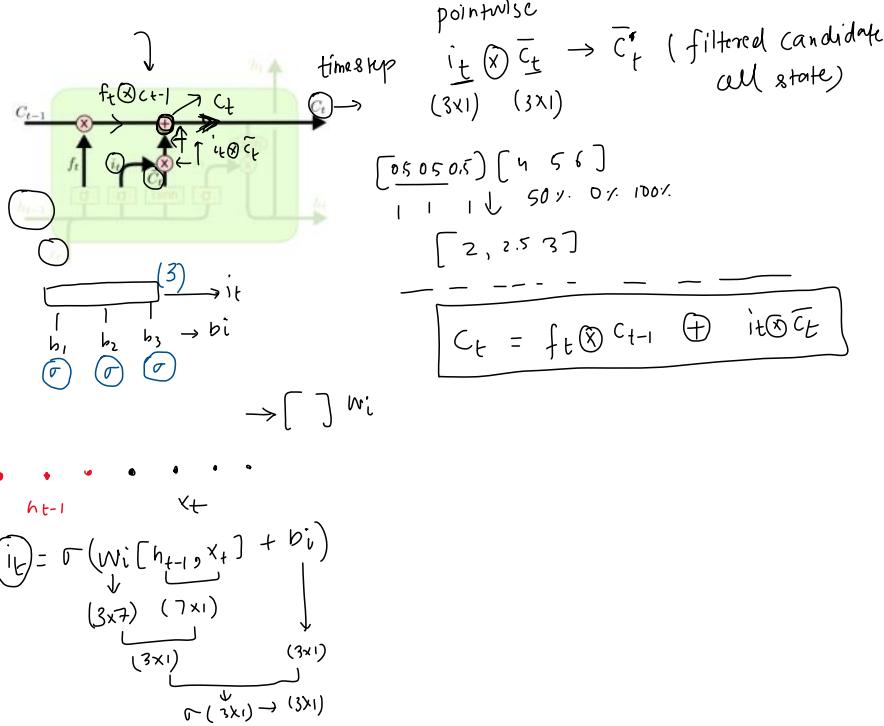
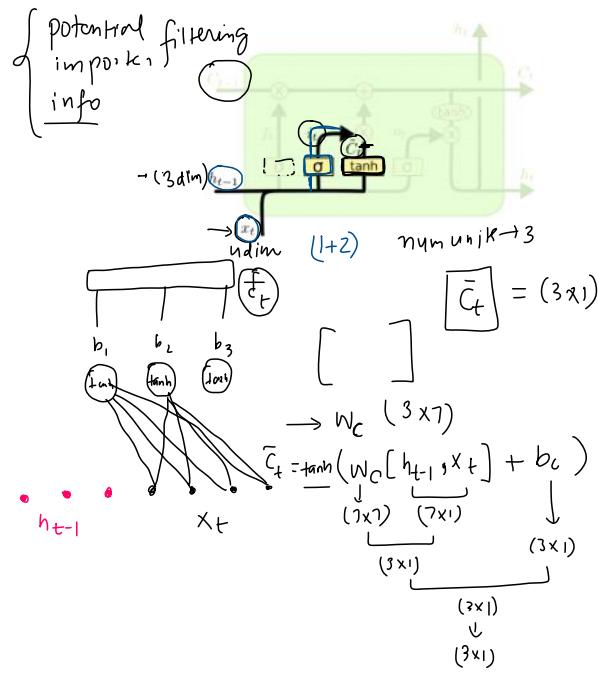


The Input Gate  
30 August 2023 04:38

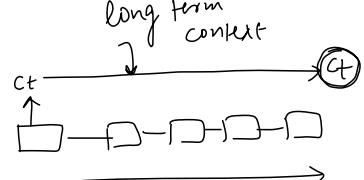
add some new imp info to  $c_t$



- stages
- 1)  $\underline{c_t}$  candidate cell state
  - 2)  $\underline{i_t}$  current
  - 3)  $\underline{c_t}$  cell state

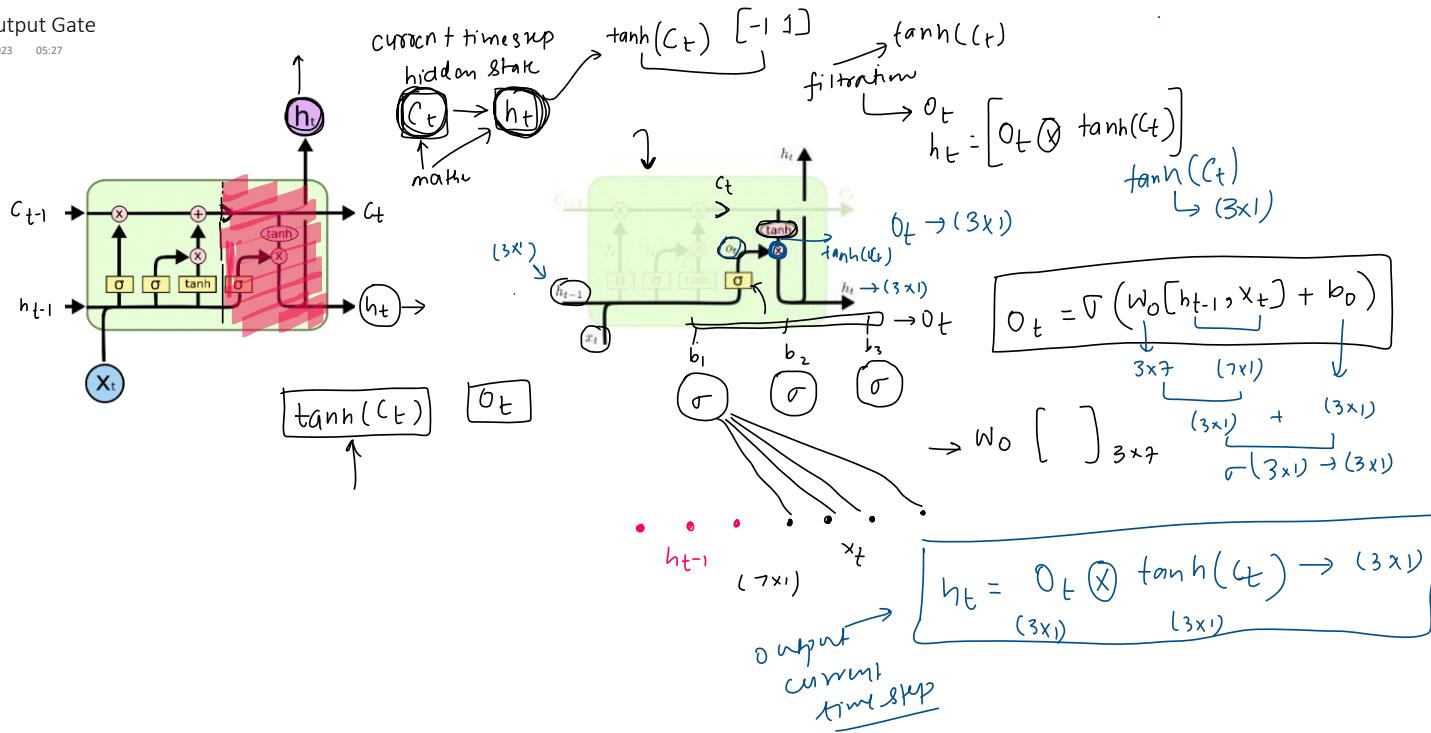


$$\begin{aligned} & f_{t-1} = [1 \ 1 \ 1] \\ & i_t \otimes \bar{c}_t = [0 \ 0 \ 0] \\ & C_t = [4 \ 5 \ 6] \end{aligned}$$



## The Output Gate

30 August 2023 05:27



# What is a Next Word Predictor

08 September 2023 08:50



code

Eran Brauer  
Mobile • 1h ago

Guy Katabi • 8:47 AM  
Hi Eran

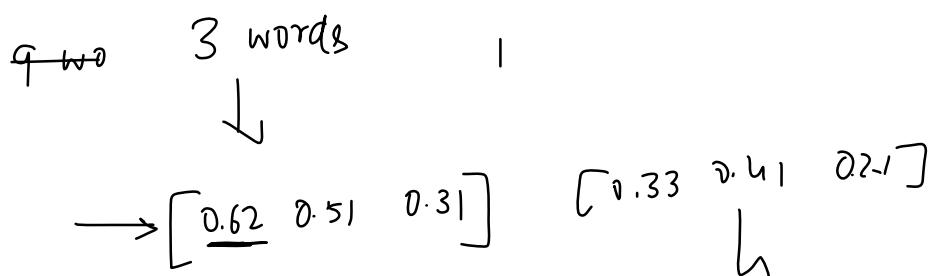
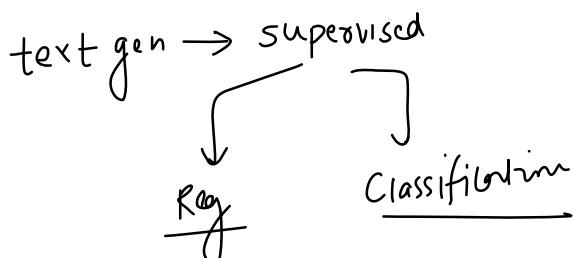
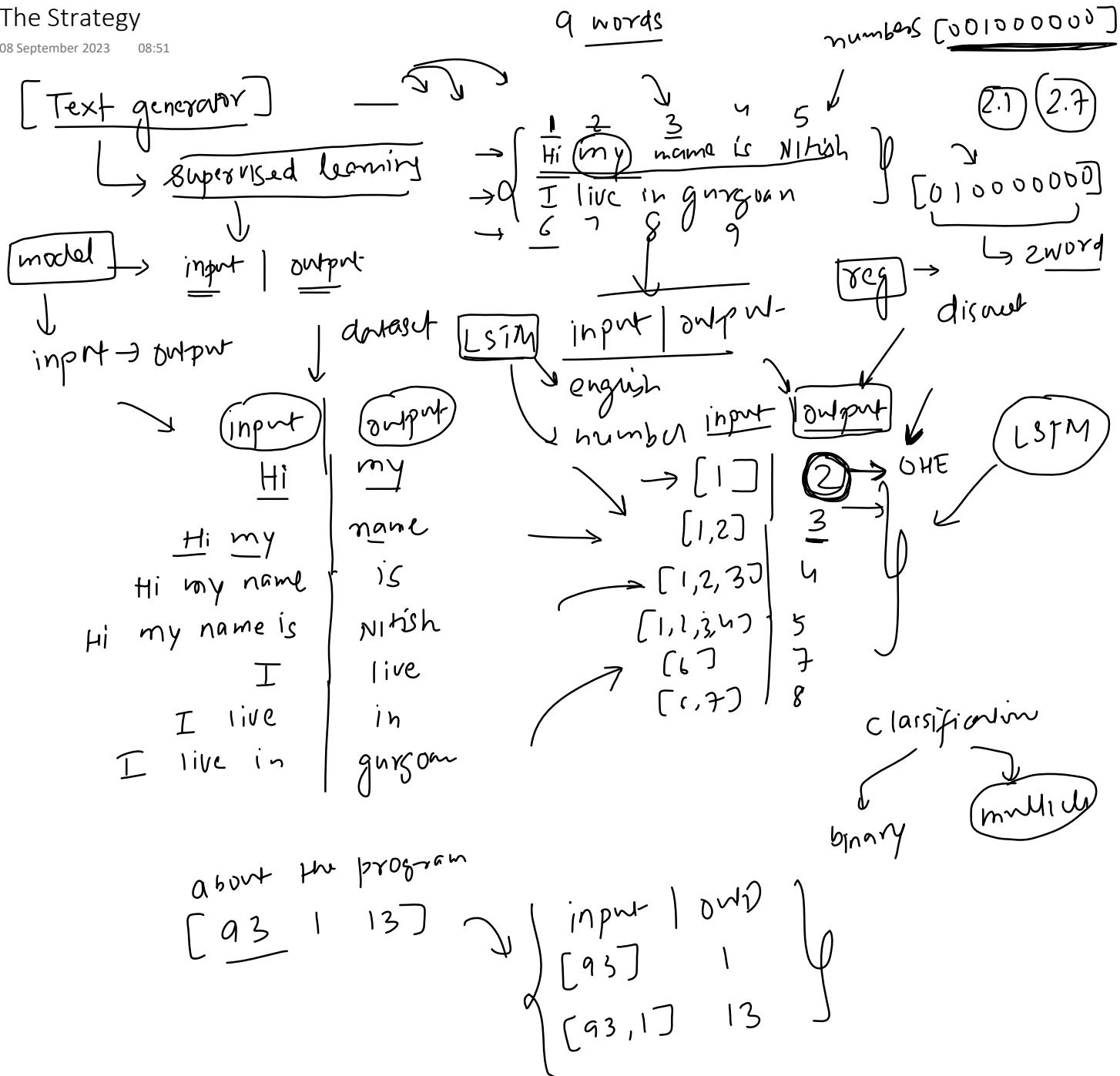
Thanks for reaching out and glad to be in your network.

Image, Video, GIF, Smileys

Send

# The Strategy

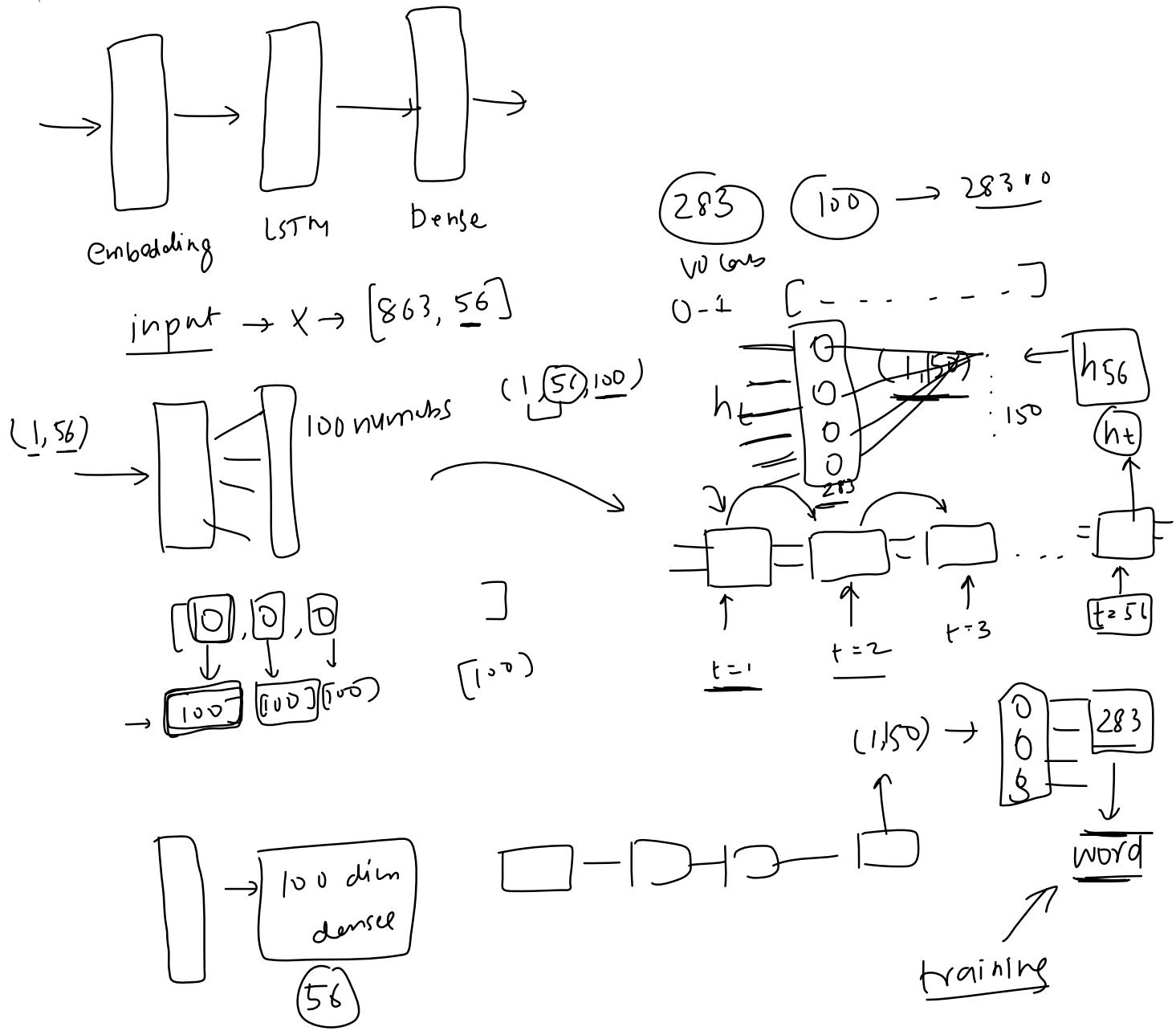
08 September 2023 08:51



$[ \text{I} \text{ D O}]$        $[ \text{O} \text{ I D}]$   
↳ first word      ↳ second word

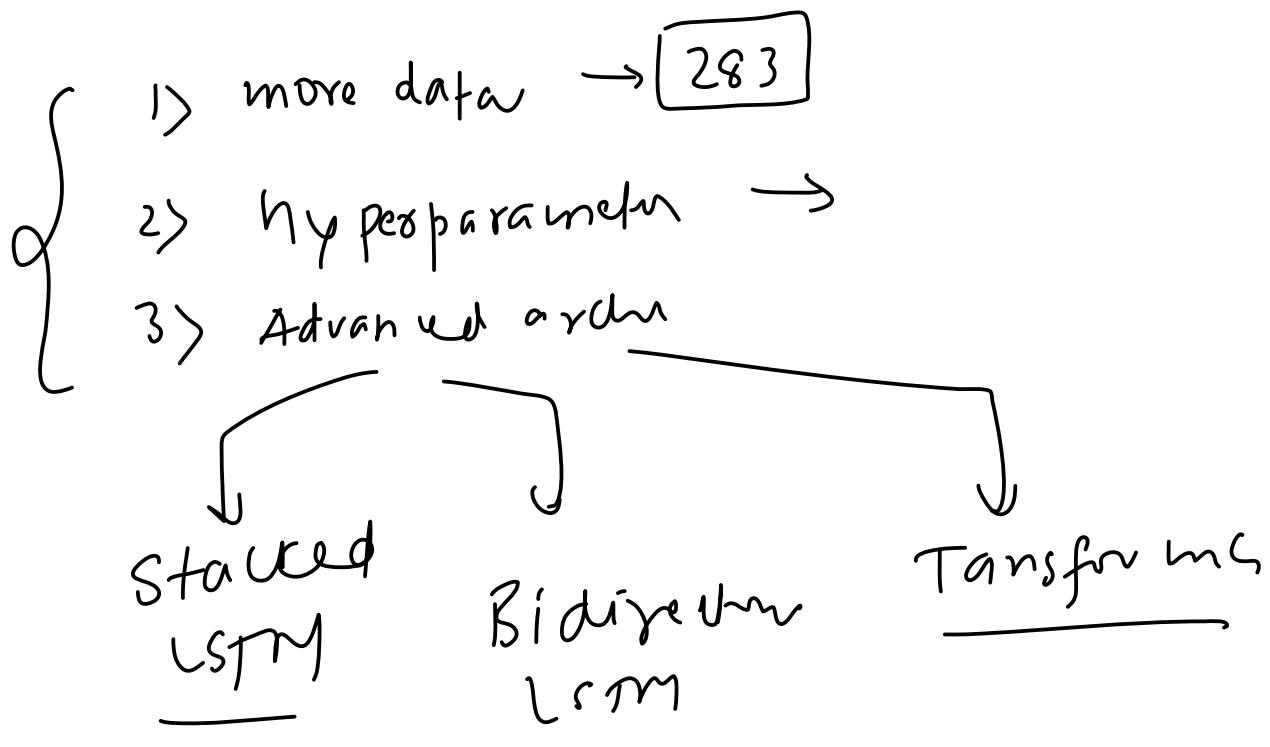
# The Architecture

08 September 2023 08:55



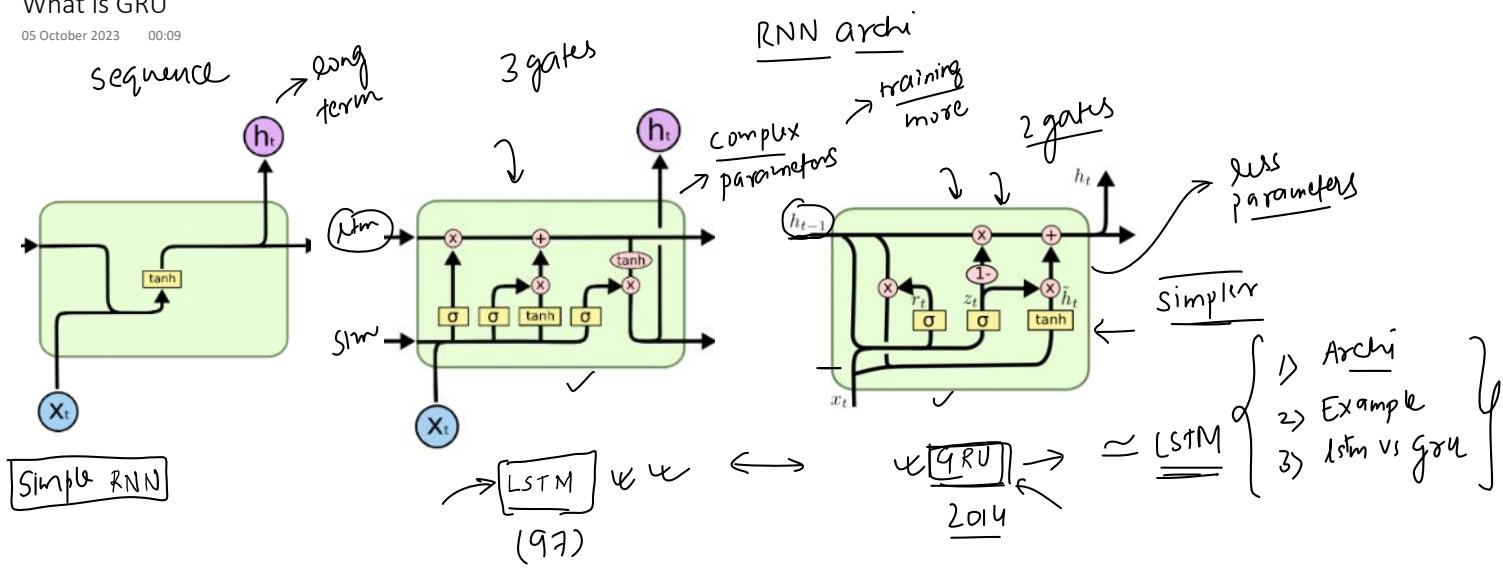
# How to improve performance?

08 September 2023 08:51



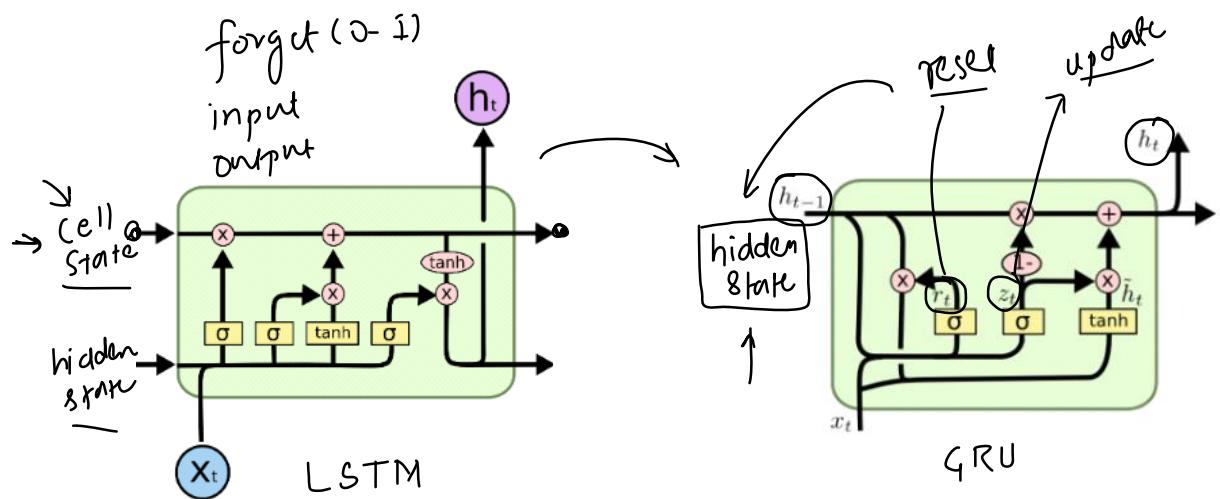
# What is GRU

05 October 2023 00:09



# The Big Idea Behind GRU

05 October 2023 00:47

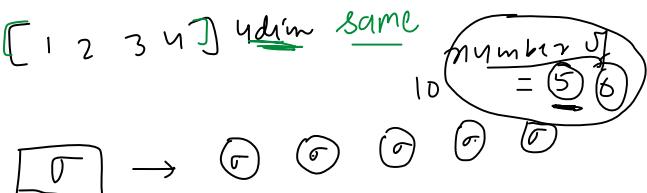
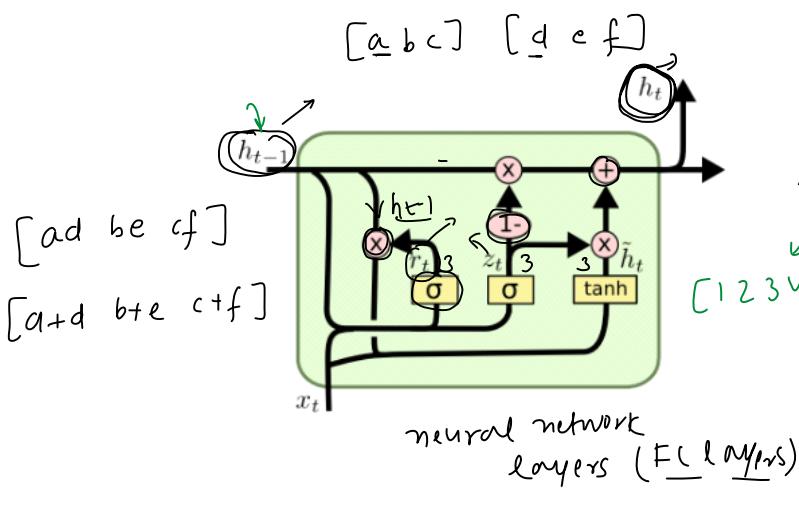
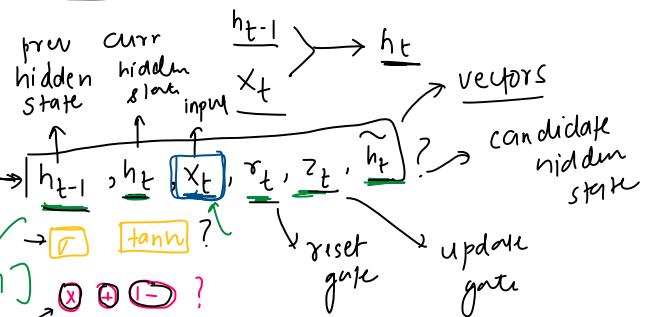


## The Setup

05 October 2023 01:07

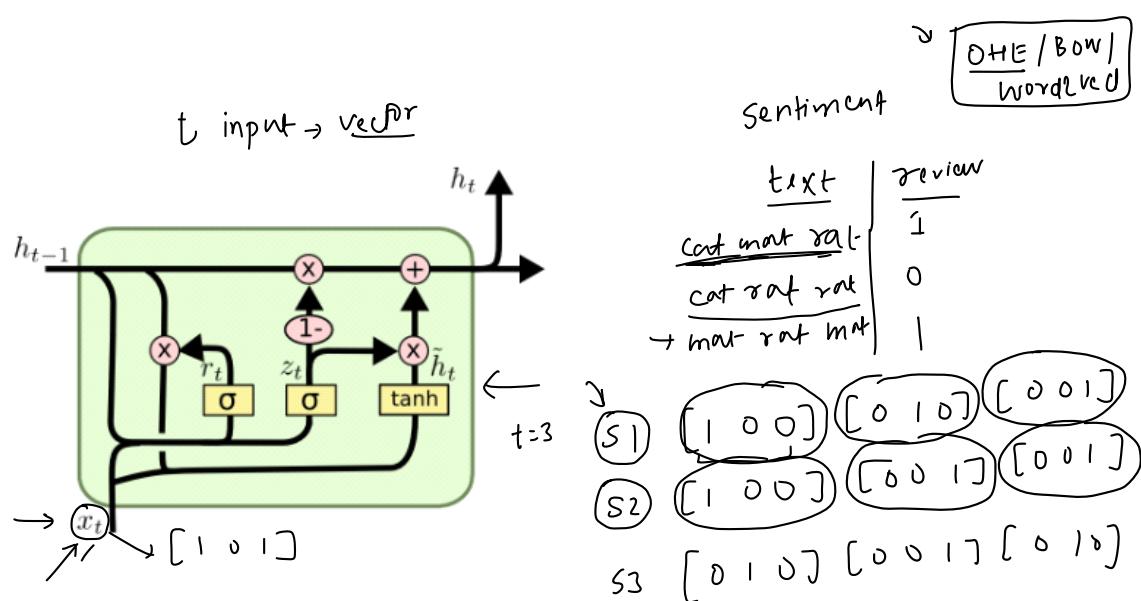
→ Advise → LSTM / GRU → confusing

goal →  $\boxed{t}$



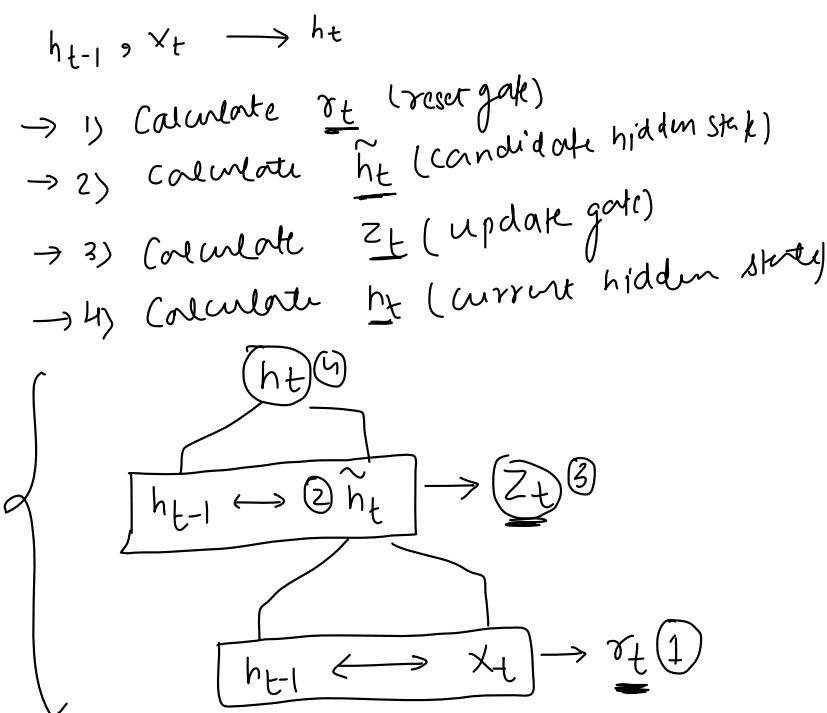
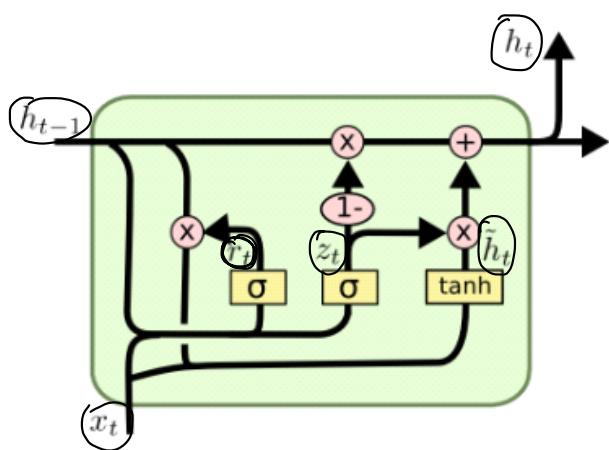
# The Input Xt

05 October 2023 01:52



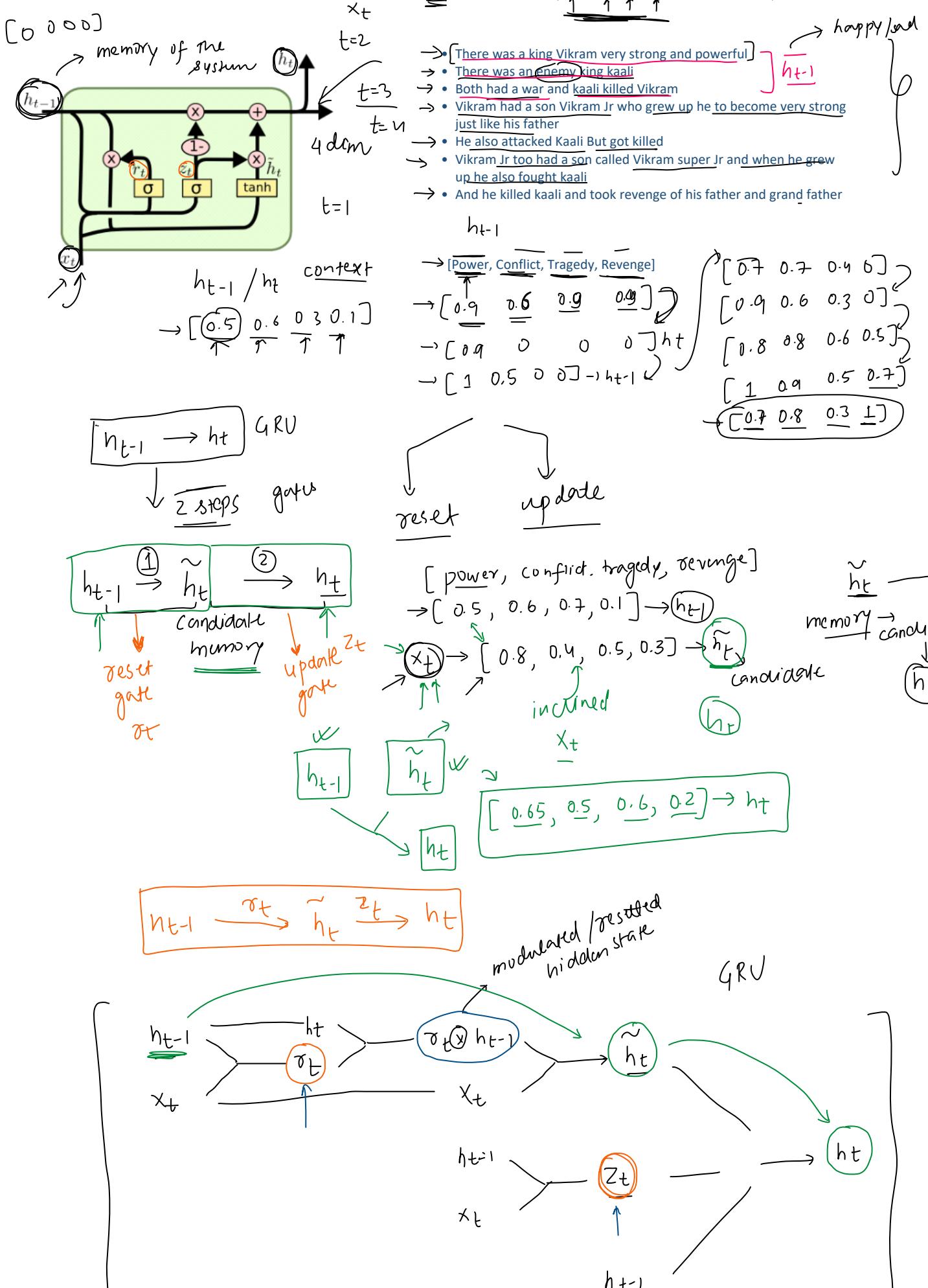
# Architecture

05 October 2023 02:10



What exactly is hidden state?

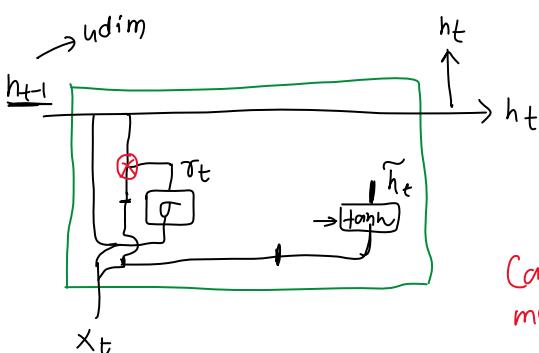
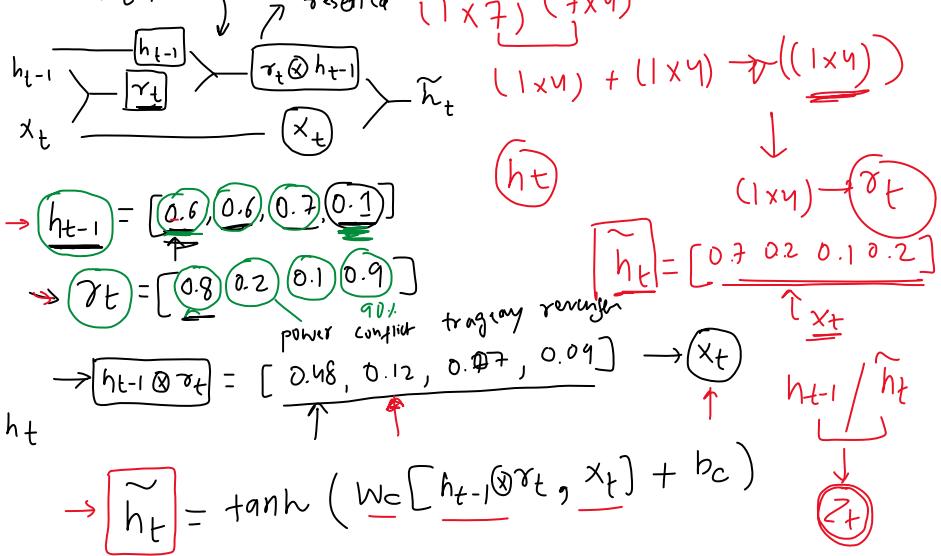
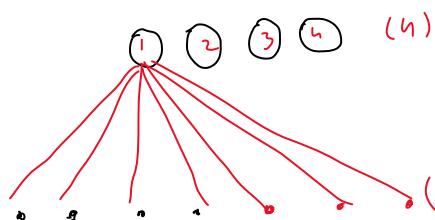
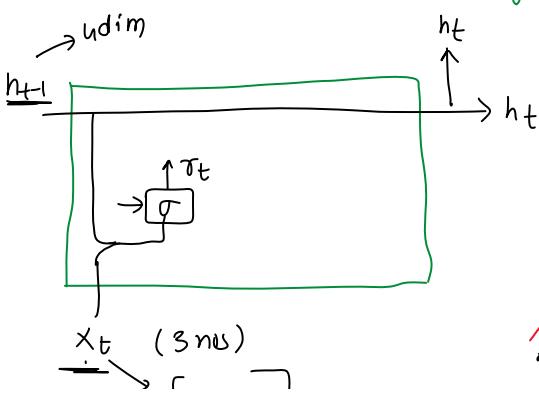
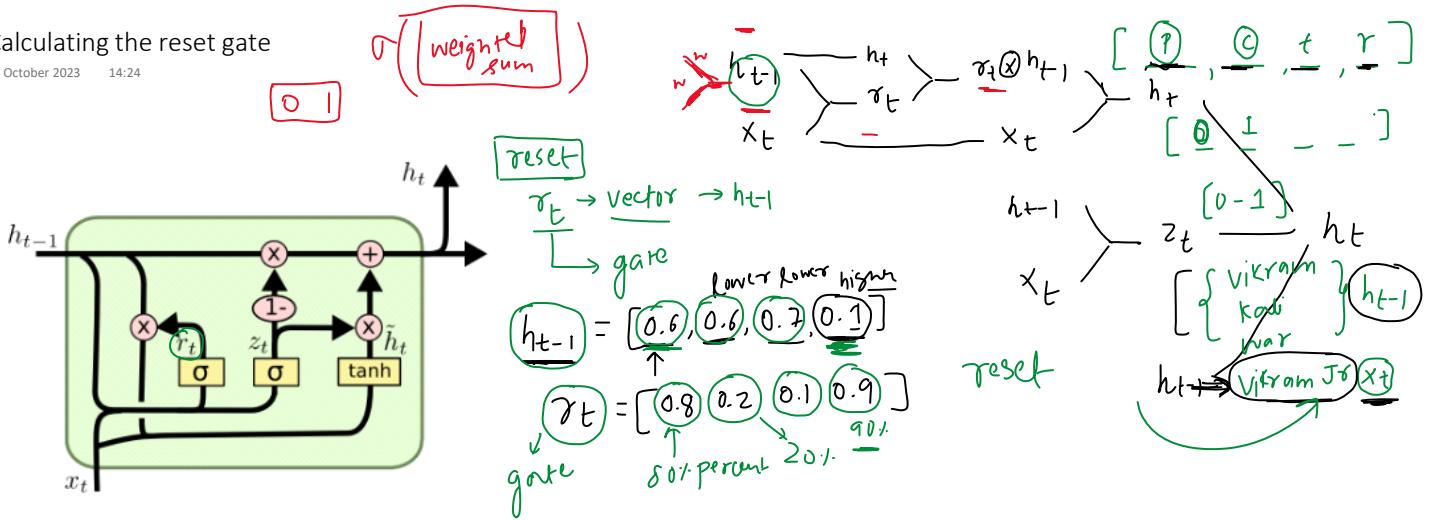
05 October 2023 02:19





## Calculating the reset gate

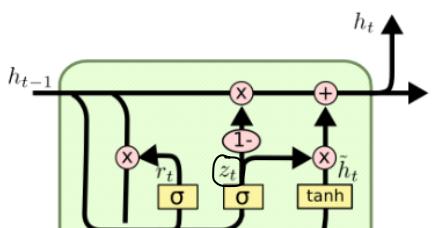
05 October 2023 14:24



Candidate memory  $\rightarrow b_C$

$\rightarrow W_C (7 \times 4)$

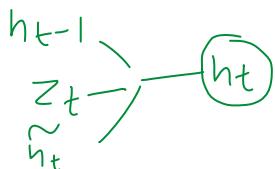
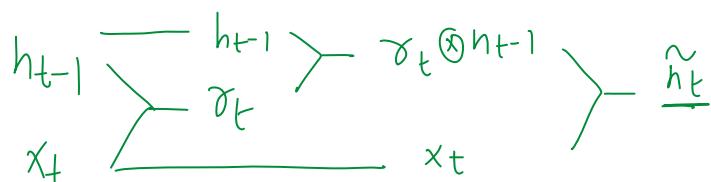
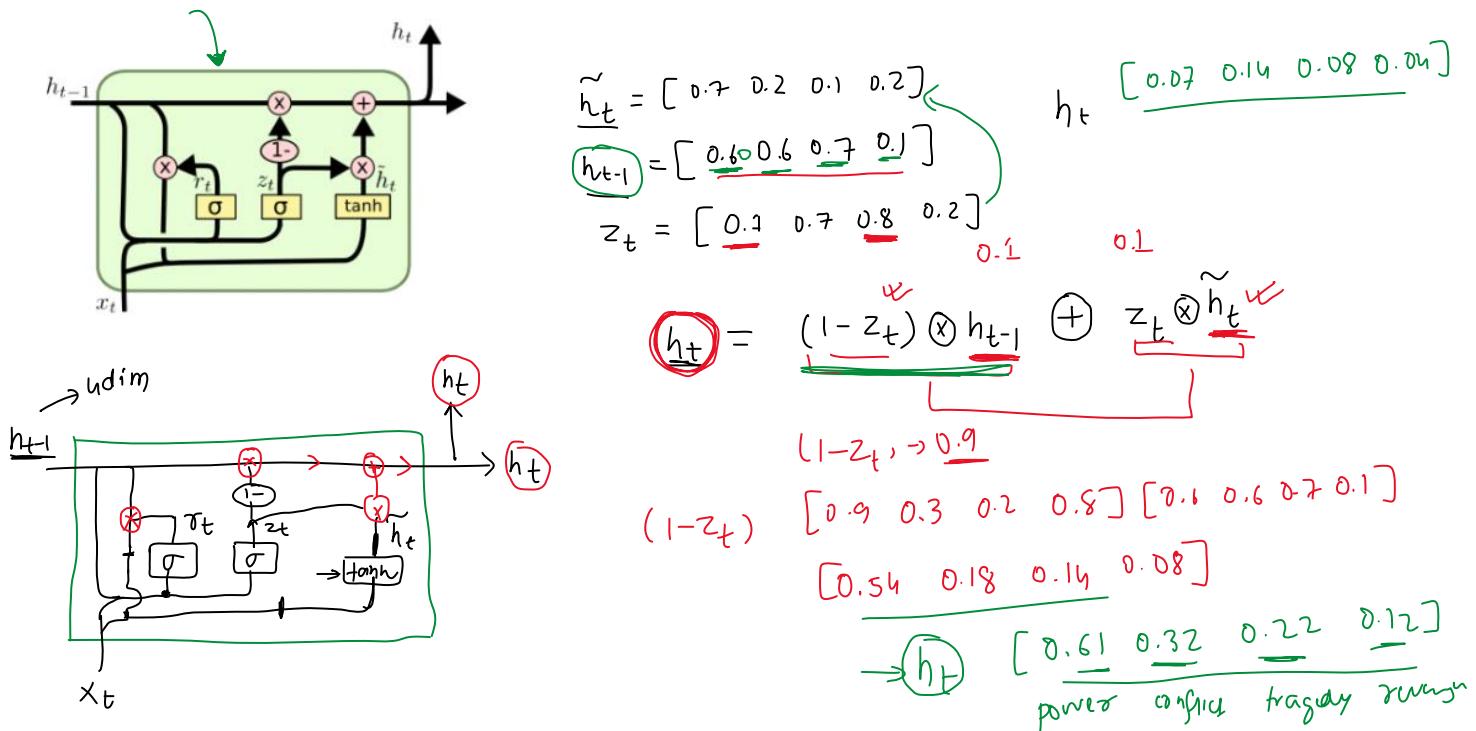
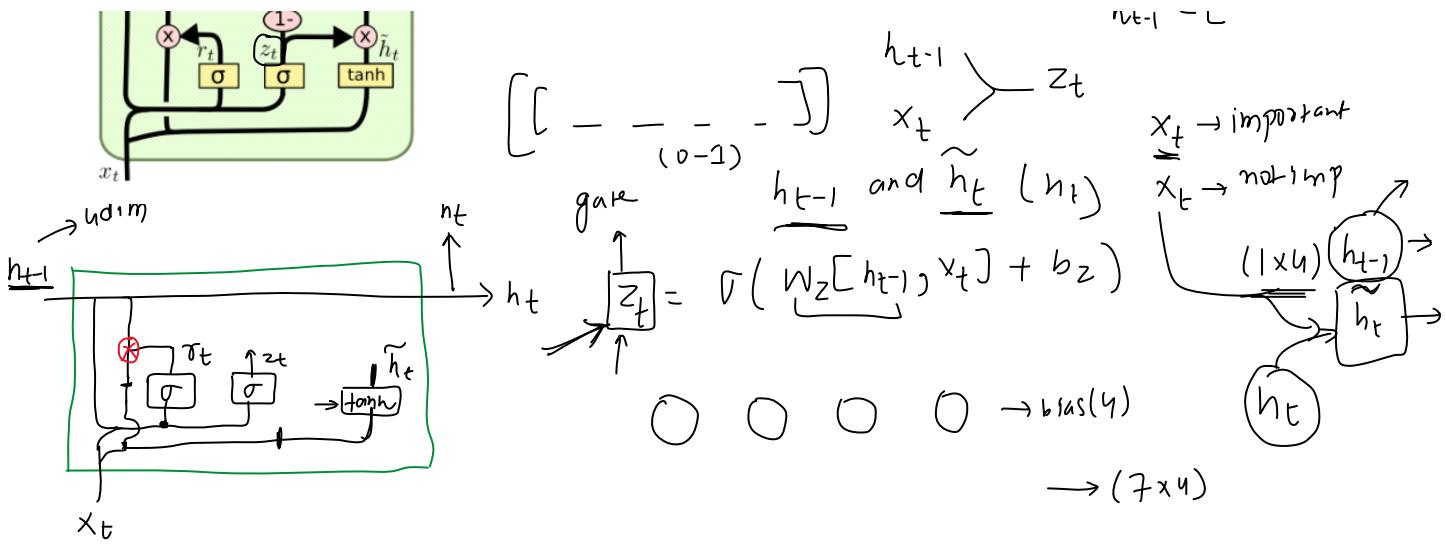
$(h_{t-1} \otimes r_t)$



$$h_{t-1} \rightarrow r_t \otimes h_{t-1} \rightarrow \tilde{h}_t = [0.7, 0.2, 0.1, 0.2]$$

$$h_{t-1} = [0.6, 0.6, 0.7, 0.1]$$

$$\tilde{h}_t = r_t \otimes h_{t-1} \rightarrow z_t$$



# LSTM vs GRU

05 October 2023 16:45 ✓

Here are the main differences between LSTM and GRU:

## 1. Number of Gates:

- LSTM: Has three gates — input (or update) gate, forget gate, and output gate.
- GRU: Has two gates — reset gate and update gate.

## 2. Memory Units:

- LSTM: Uses two separate states - the cell state ( $c_t$ ) and the hidden state ( $h_t$ ). The cell state acts as an "internal memory" and is crucial for carrying long-term dependencies.
- GRU: Simplifies this by using a single hidden state ( $h_t$ ) to both capture and output the memory.

## 3. Parameter Count:

- LSTM: Generally has more parameters than a GRU because of its additional gate and separate cell state. For an input size of  $d$  and a hidden size of  $h$ , the LSTM has  $4 \times ((d \times h) + (h \times h) + h)$  parameters.
- GRU: Has fewer parameters. For the same sizes, the GRU has  $3 \times ((d \times h) + (h \times h) + h)$  parameters.

## 4. Computational Complexity:

- LSTM: Due to the extra gate and cell state, LSTMs are typically more computationally intensive than GRUs.
- GRU: Is simpler and can be faster to compute, especially on smaller datasets or when computational resources are limited.

## 5. Empirical Performance:

- LSTM: In many tasks, especially more complex ones, LSTMs have been observed to perform slightly better than GRUs.
- GRU: Can perform comparably to LSTMs on certain tasks, especially when data is limited or tasks are simpler. They can also train faster due to fewer parameters.

## 6. Choice in Practice:

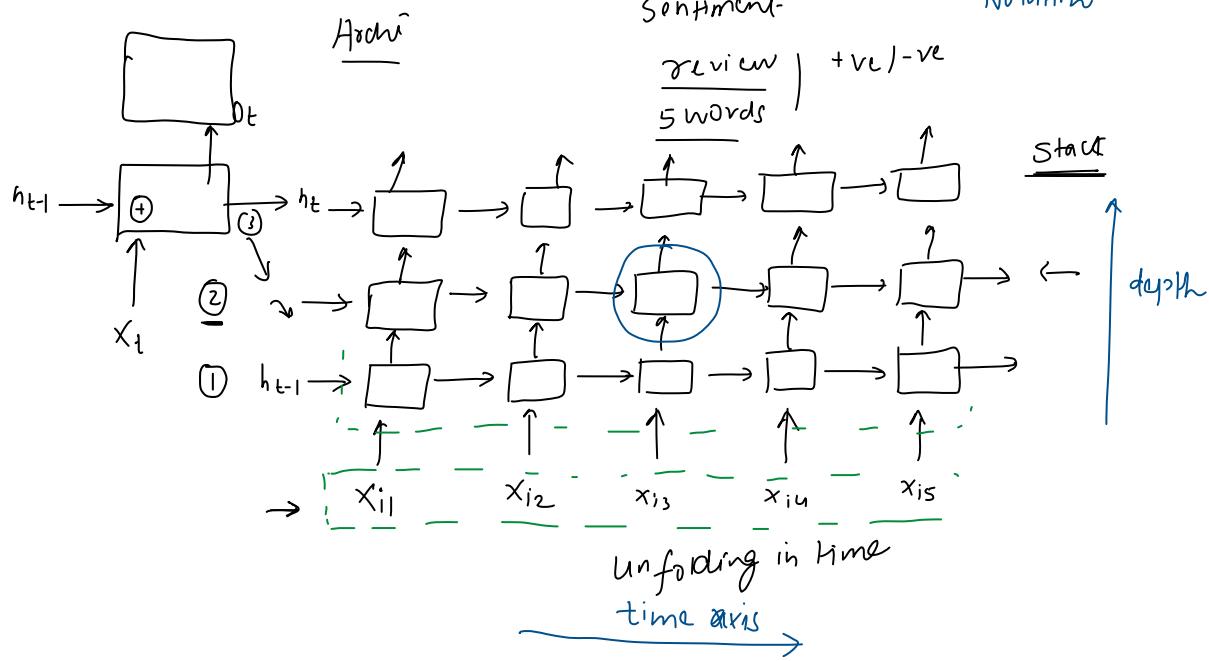
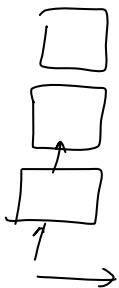
- The choice between LSTM and GRU often comes down to empirical testing. Depending on the dataset and task, one might outperform the other. However, GRUs, due to their simplicity, are often the first choice when starting out.

# What is Deep RNN →

17 October 2023

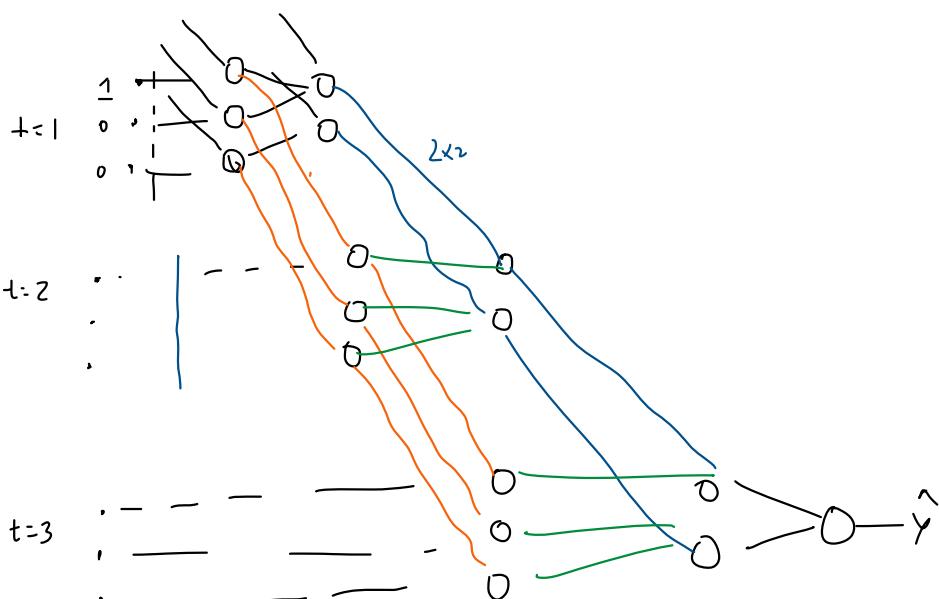
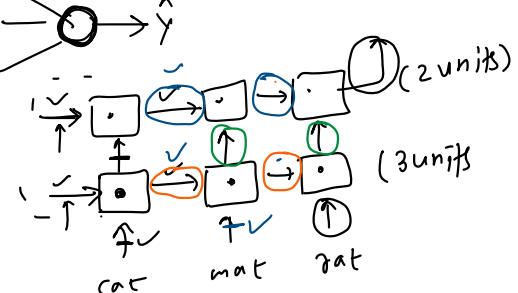
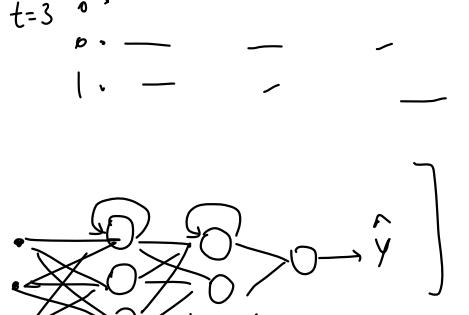
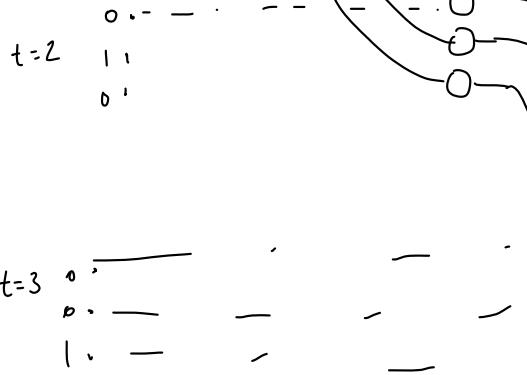
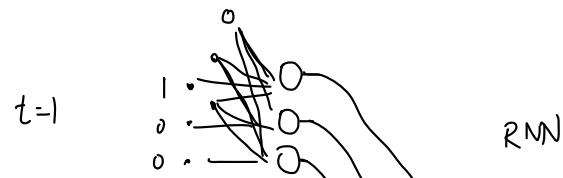
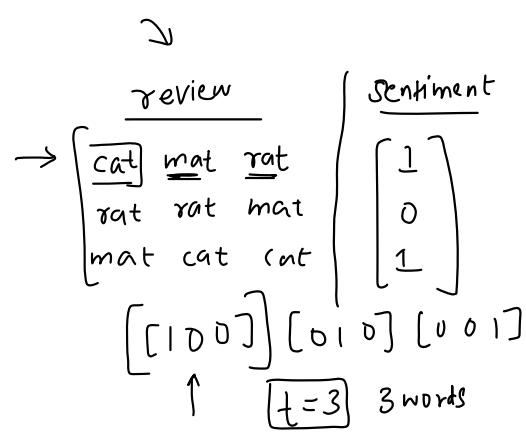
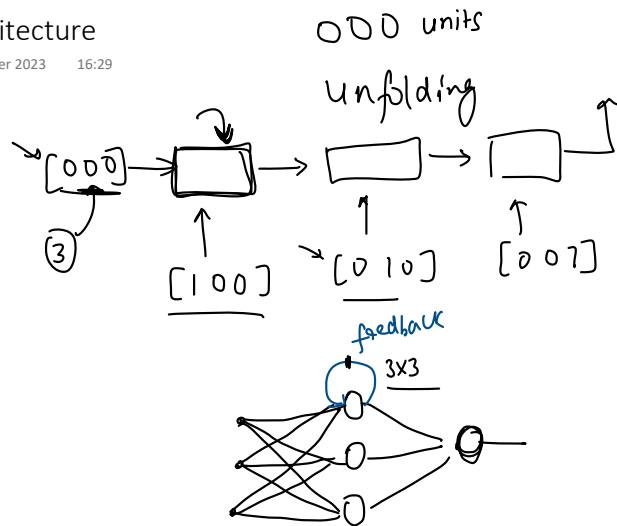
ANN

J

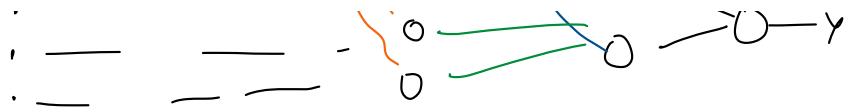


## Architecture

17 October 2023 16:29



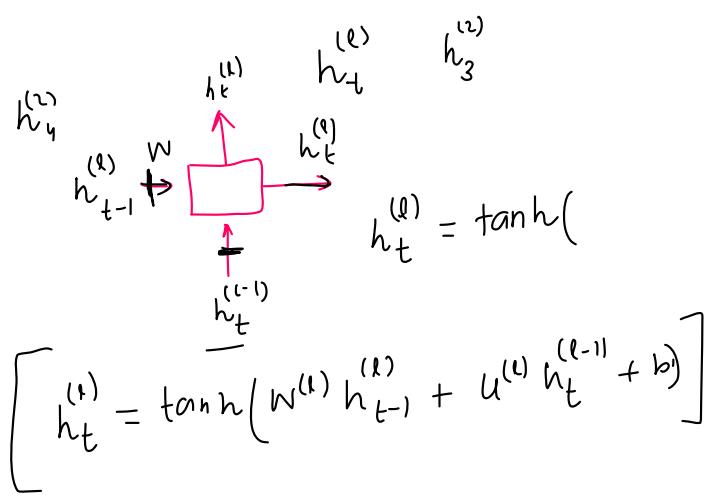
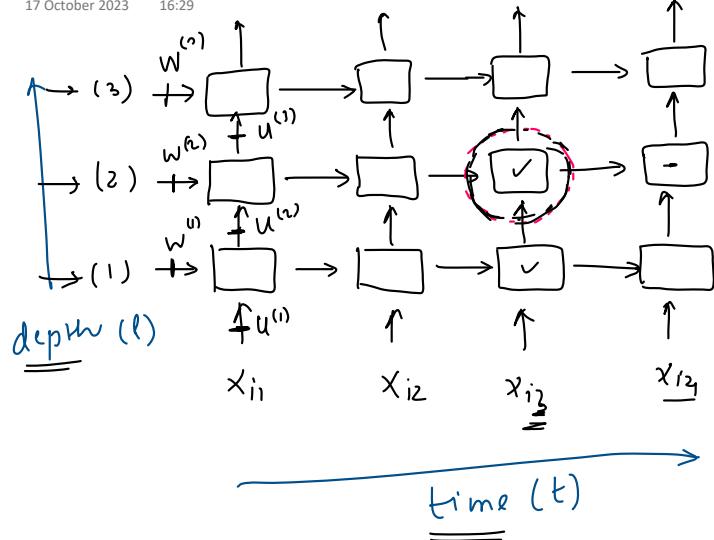
$t=3$



## Notation

17 October 2023

16:29



## Why and When to use?

17 October 2023 16:29

- {
- 1. Hierarchical Representation ✓
- 2. Customization for Advanced Tasks
- }

### deep KNN

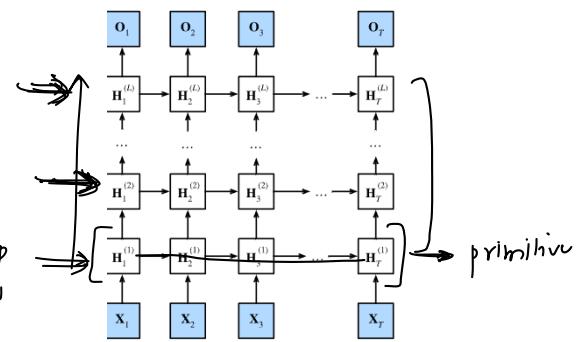
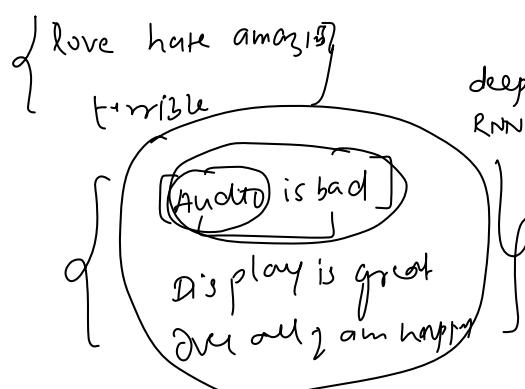
product

### stack

sentence

encoder-decoder  
↓  
machine

{   
 deep  
 KNNs  
 } ↗



→ sentence



### When to use Deep RNNs?

Complex  
tasks

{ speech recg  
Machine translation }

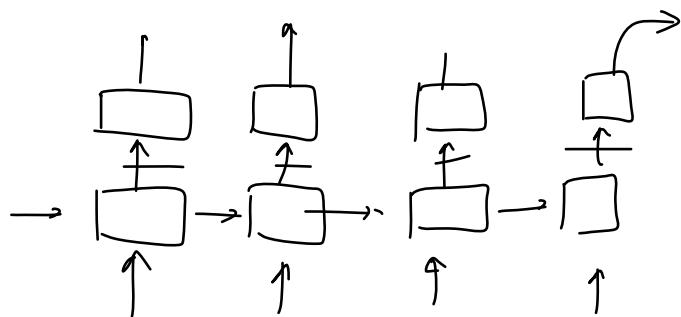
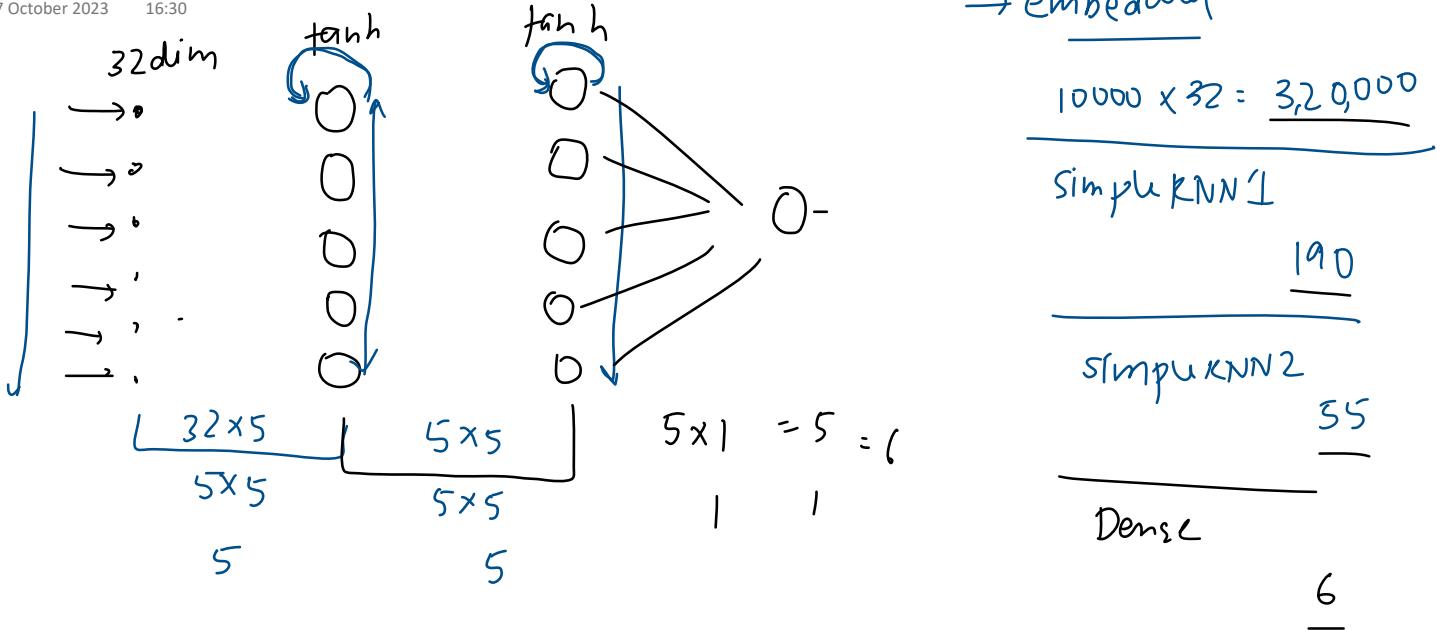
Large datasets  
Overfitting

Computational

Simpler Models  
↓  
Deep RNN

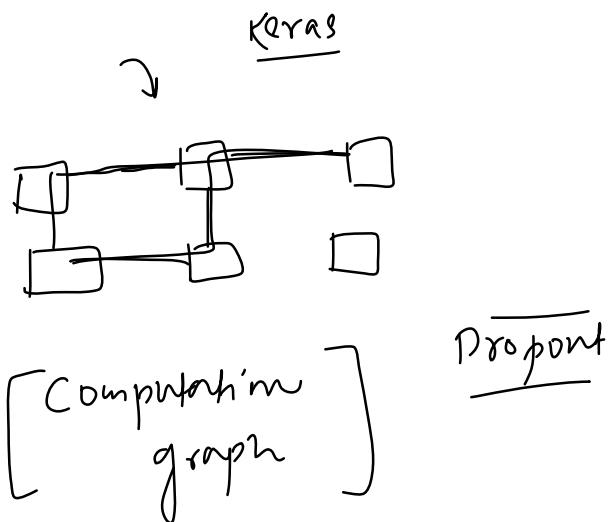
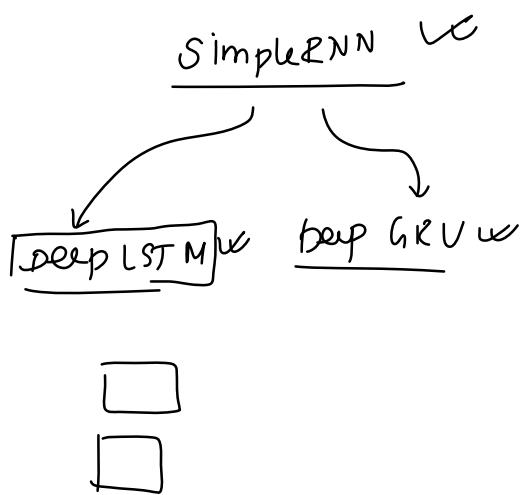
## Code Example

17 October 2023 16:30



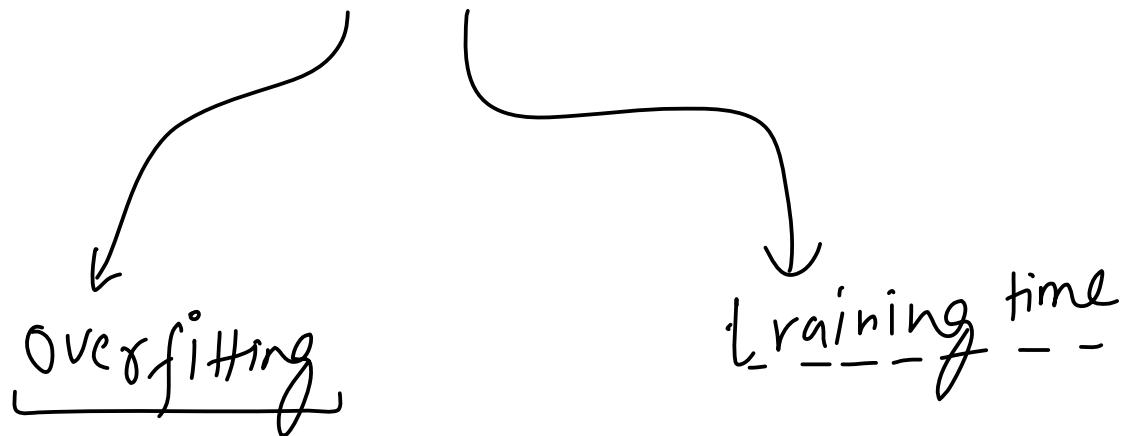
## Variants

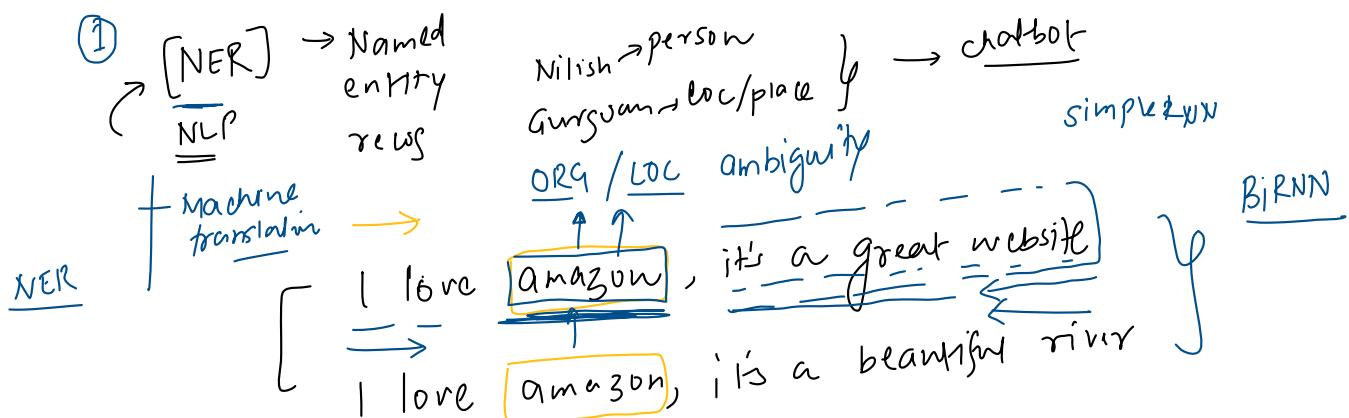
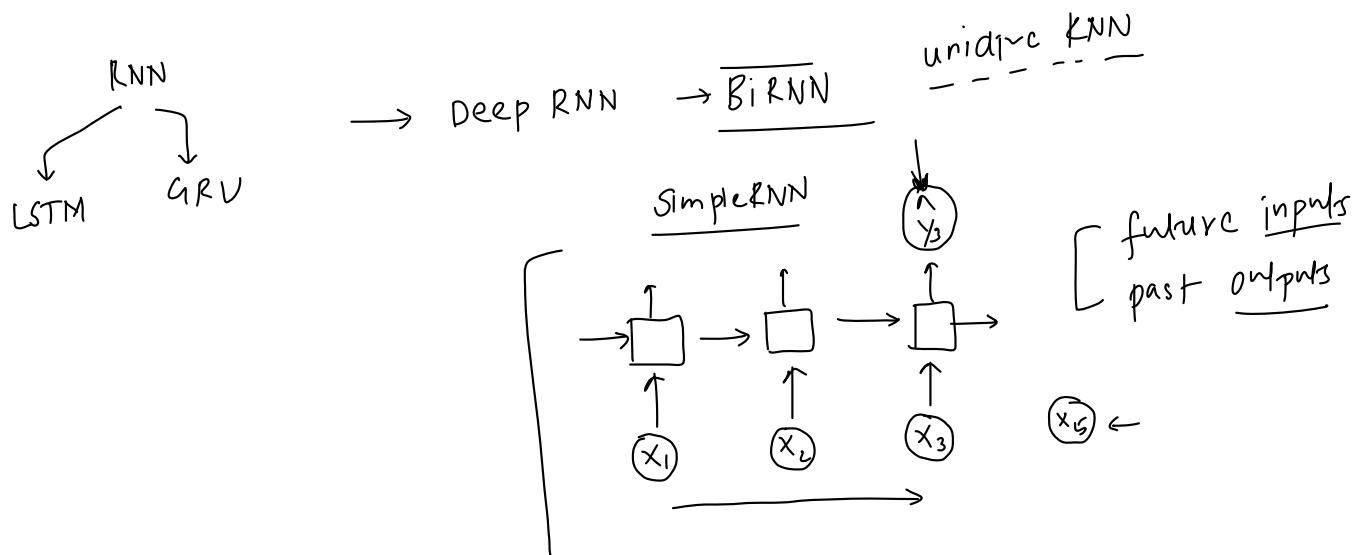
17 October 2023 16:30



# Disadvantages

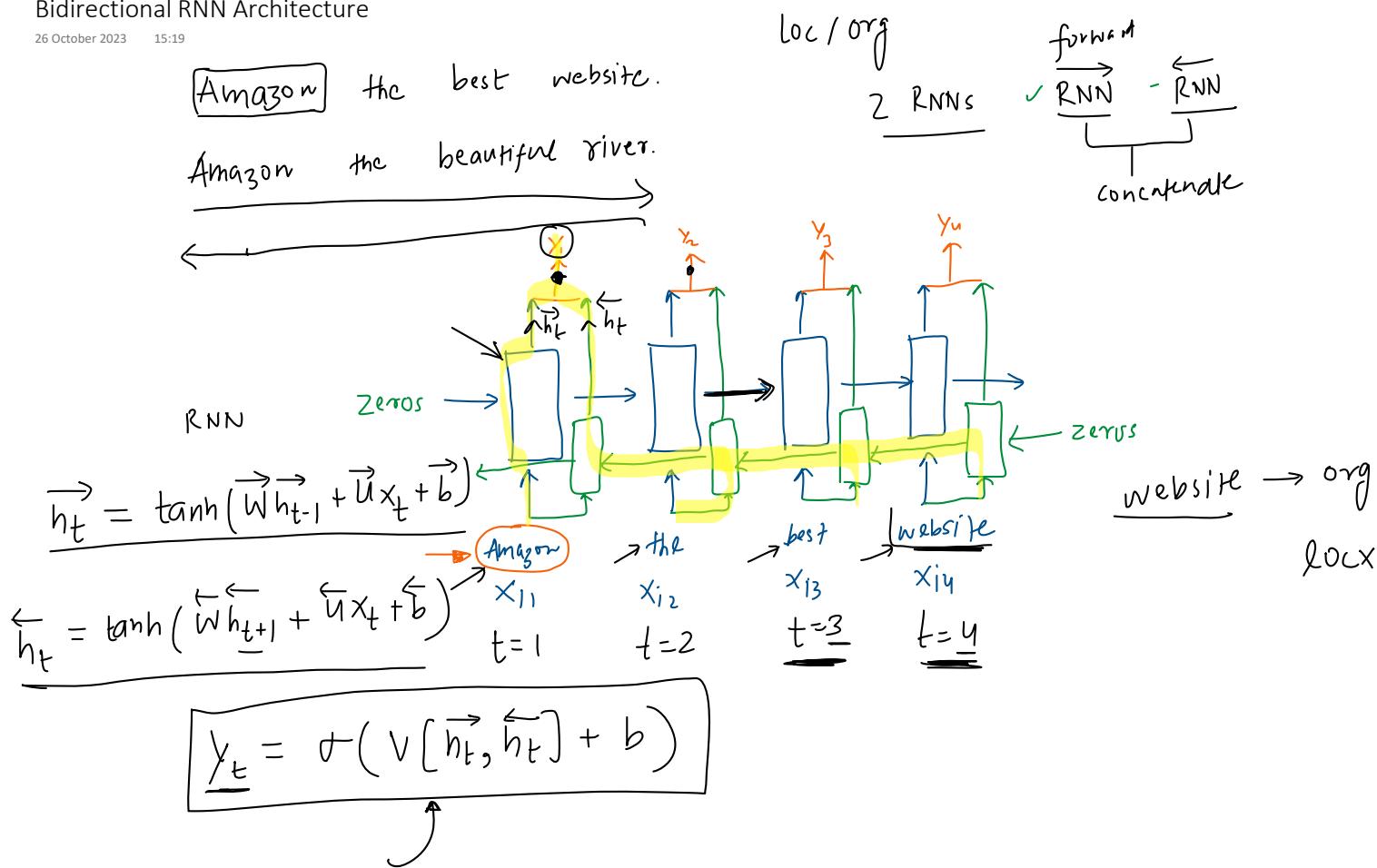
17 October 2023 16:30





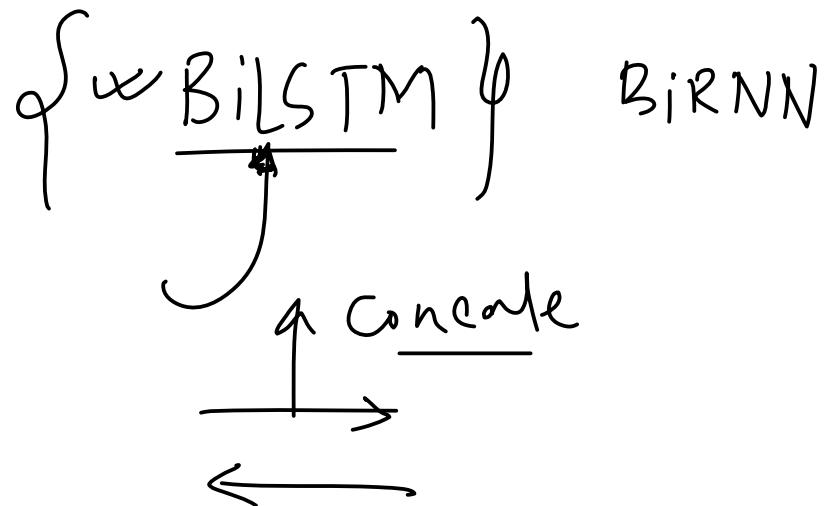
## Bidirectional RNN Architecture

26 October 2023 15:19



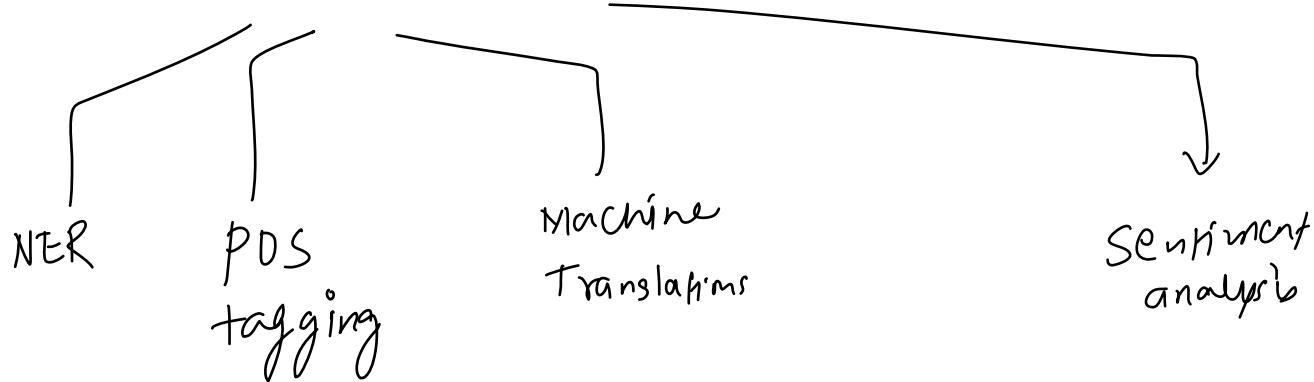
# Code

26 October 2023 15:21



## Applications and Drawbacks

26 October 2023 15:21



[Time series forecasting]

→ ←

→ Complexity → 190 → 380  
↓ ↓ ↓  
downdown

→ training → overfitting

→ → ← [Speech recog.] → birnn  
↓ ↓ ↓  
latency → slow