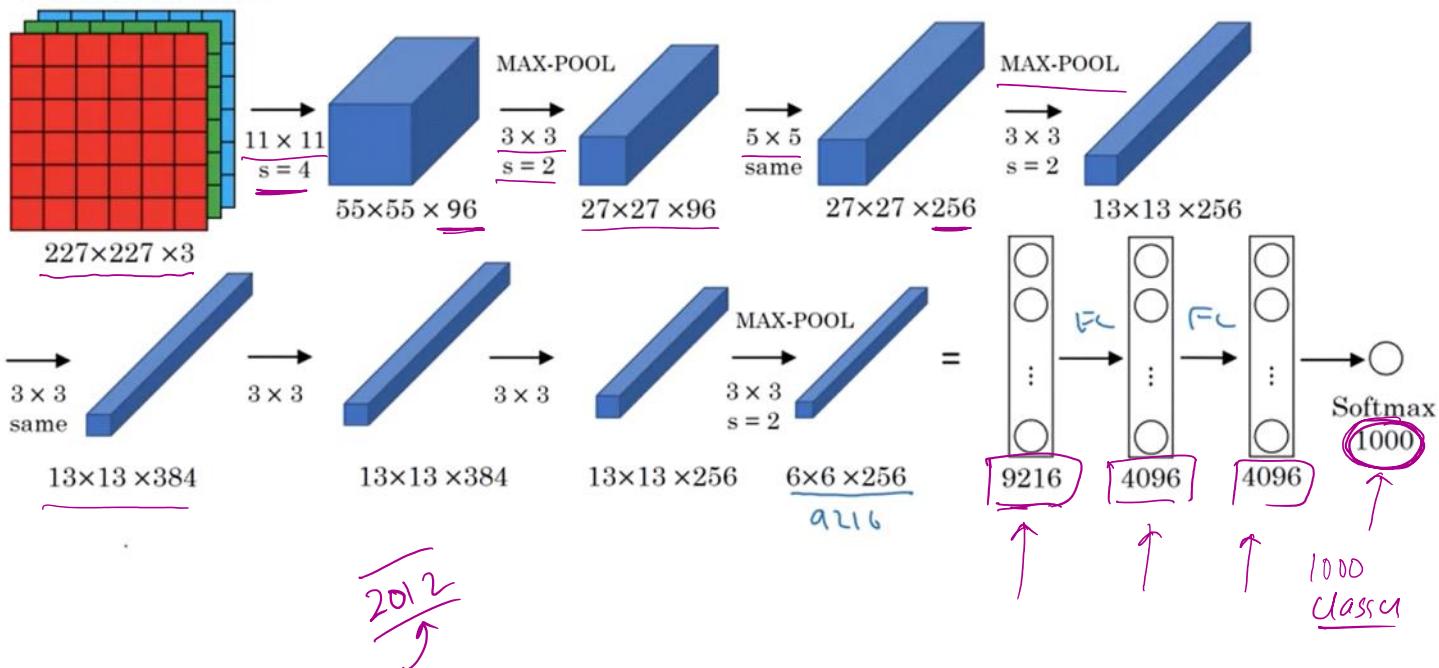


AlexNet

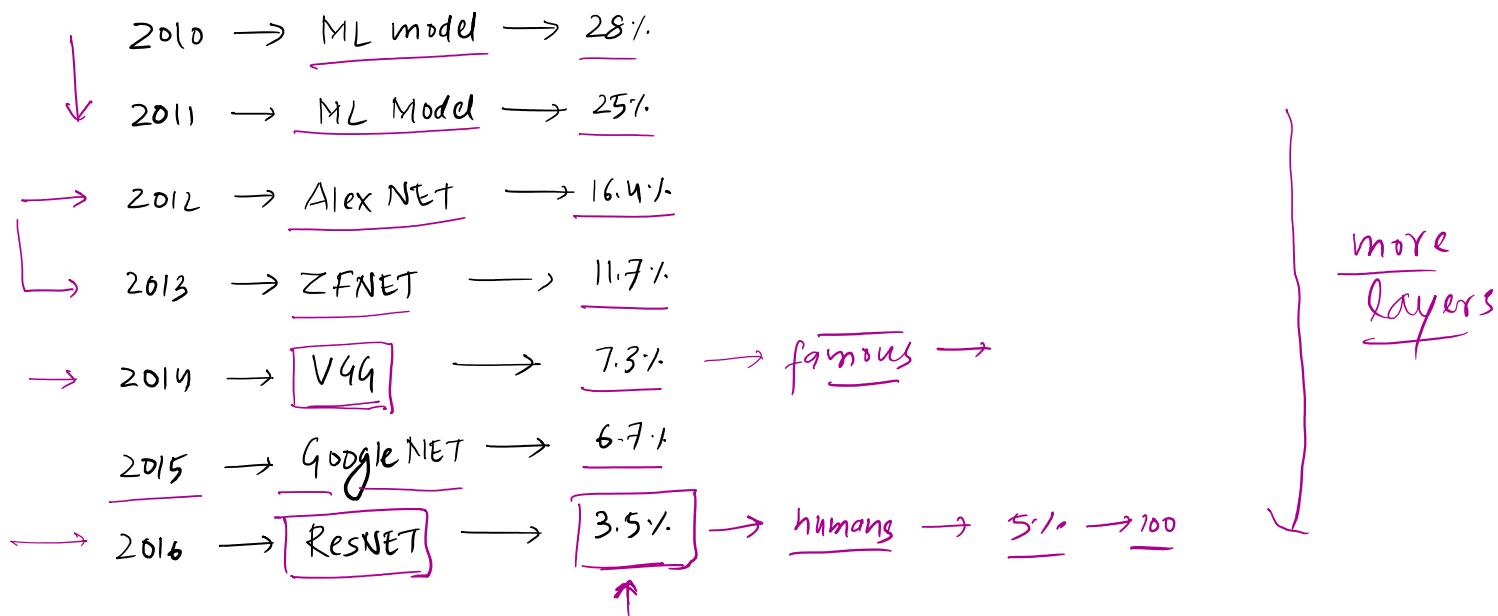


[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]

Andrew Ng

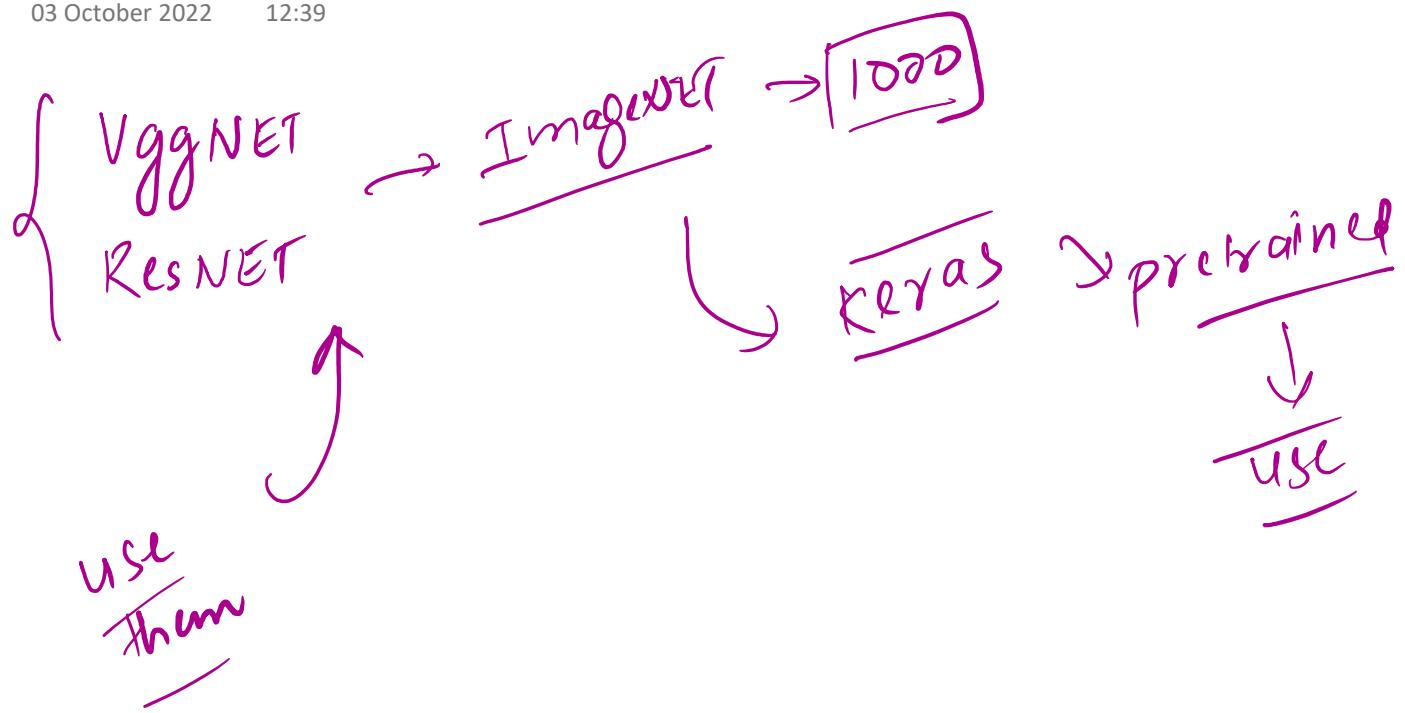
Famous Architectures

03 October 2022 12:39



Idea of Pretrained Models

03 October 2022 12:39



Keras Demo

03 October 2022 12:40

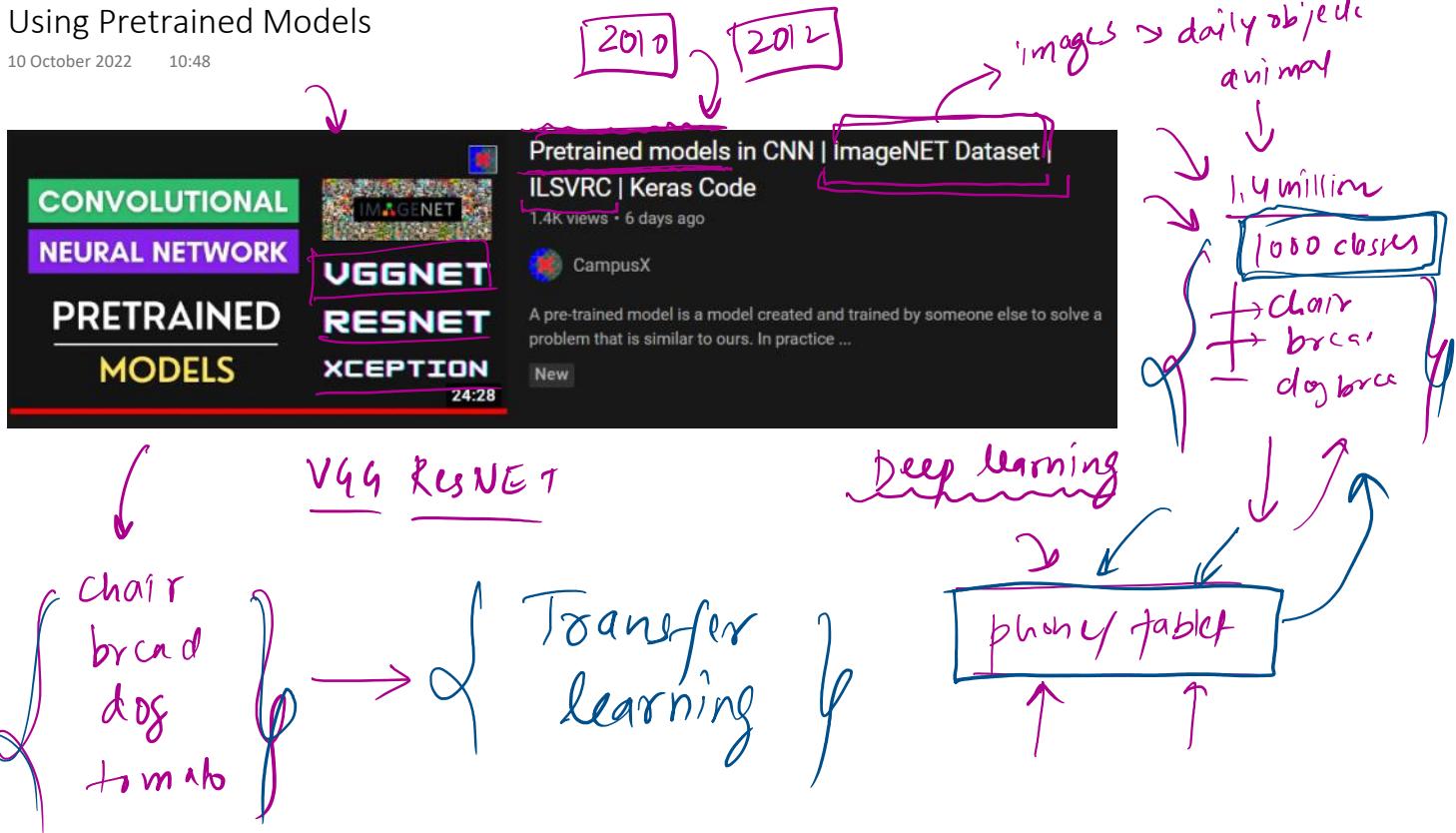
Problem with training your own model

10 October 2022 10:48

- 1) Data hungry → labelled → 10,000 → google
↓
cat/dog → manan
labor
- 2) lot of time → days/weeks)

Using Pretrained Models

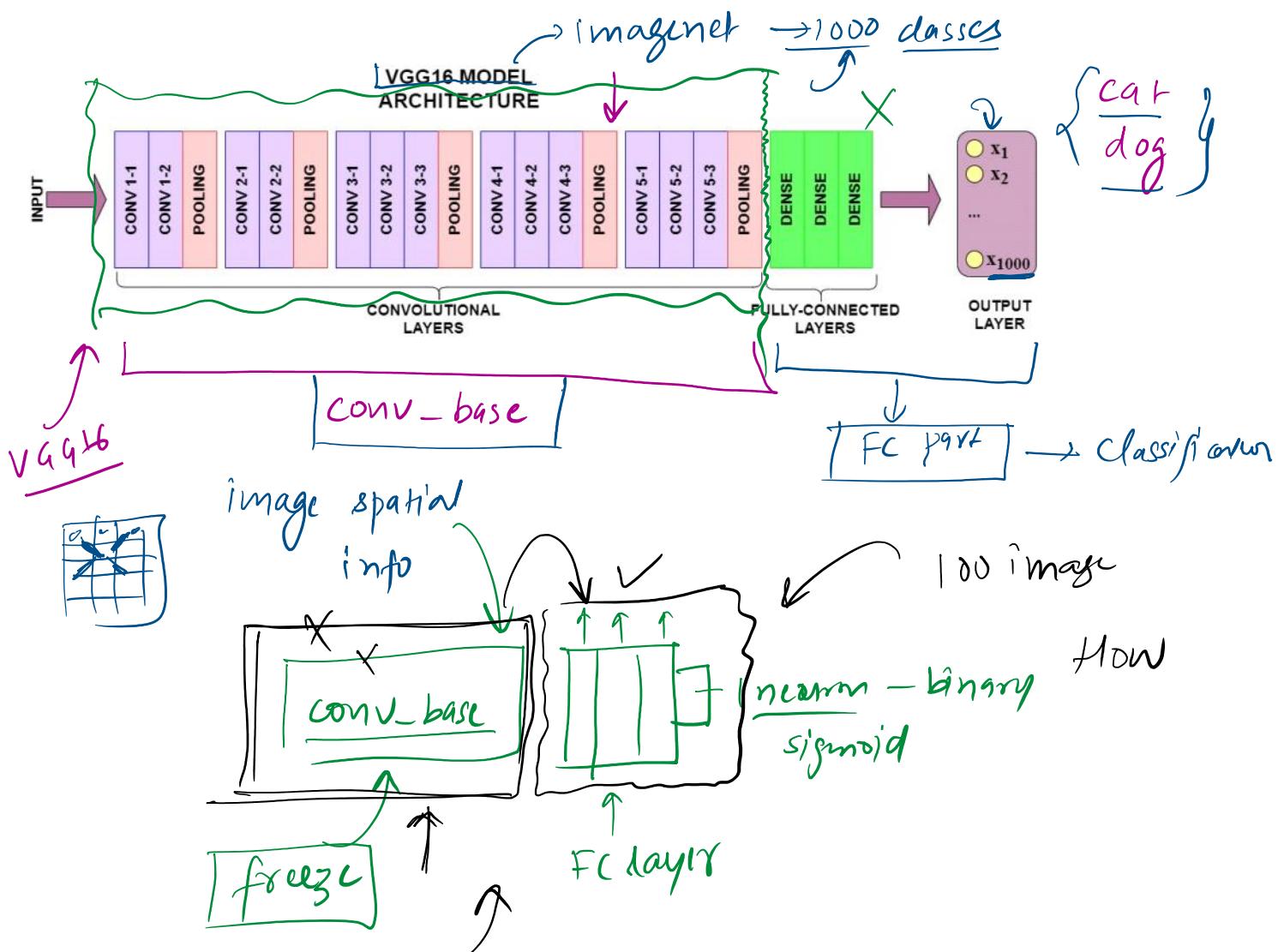
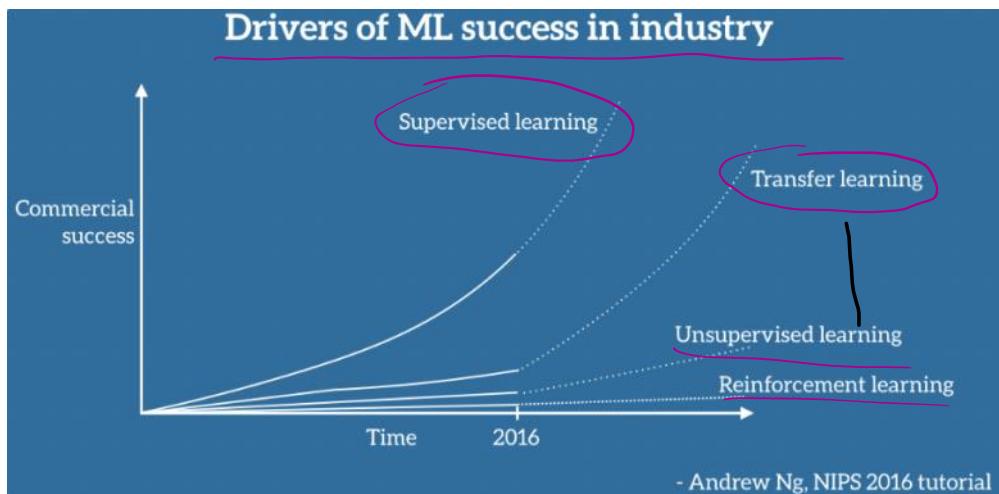
10 October 2022 10:48



Transfer Learning

10 October 2022 10:49

Transfer learning is a research problem in machine learning that focuses on storing knowledge gained while solving one problem and applying it to a different but related problem.

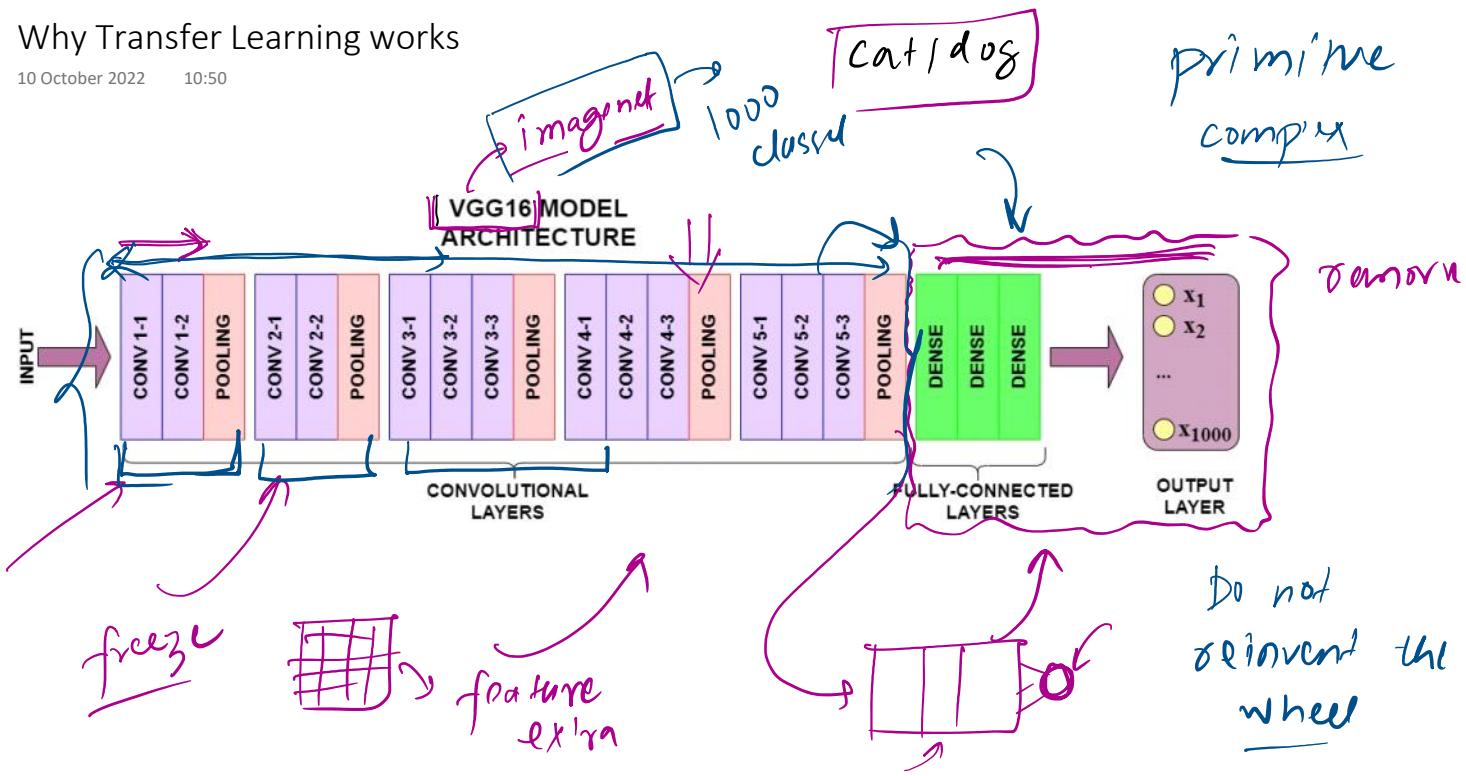


IT^os



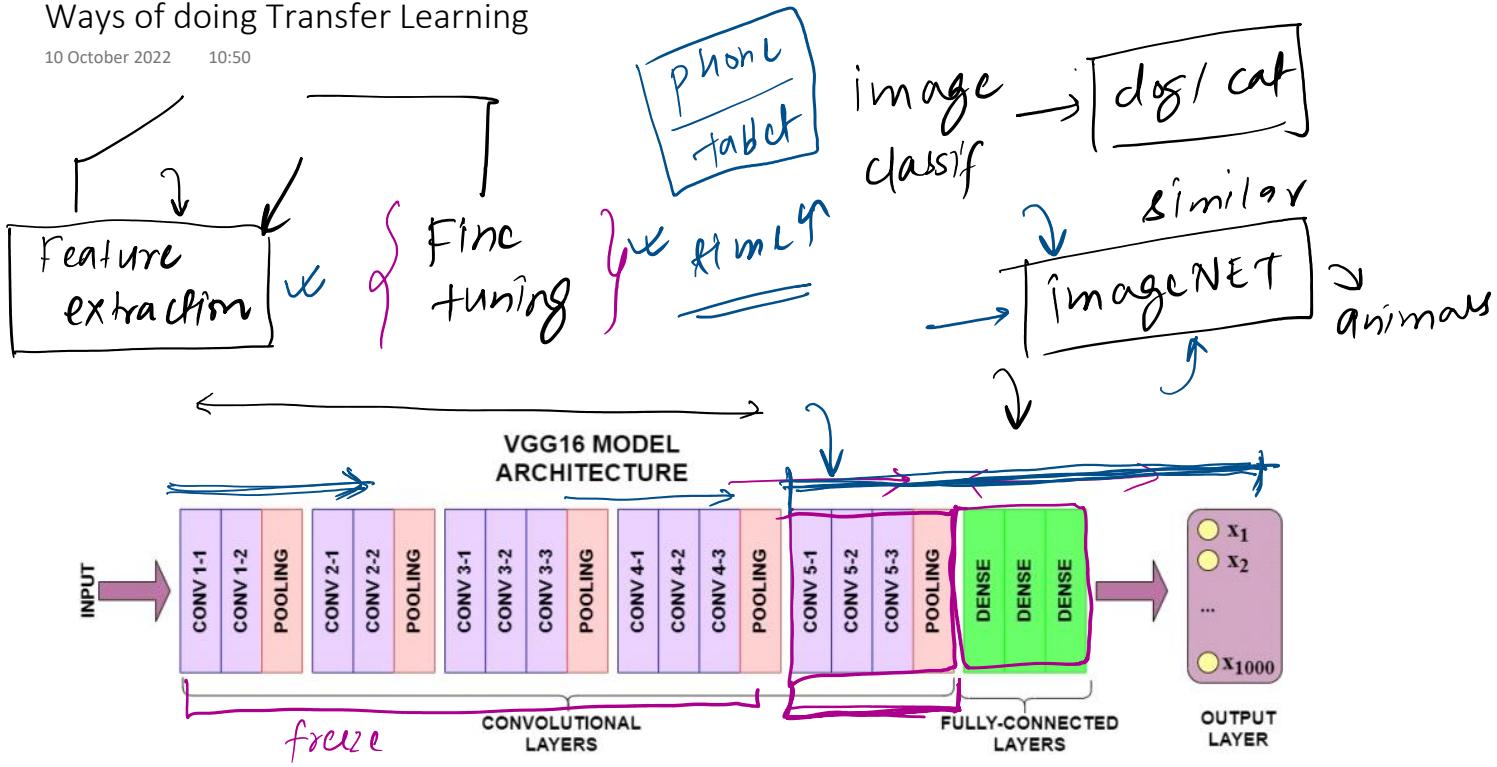
Why Transfer Learning works

10 October 2022 10:50



Ways of doing Transfer Learning

10 October 2022 10:50



Code

10 October 2022 10:50

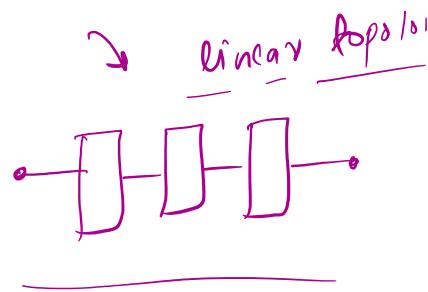
Problem with Sequential Model

14 October 2022 16:00

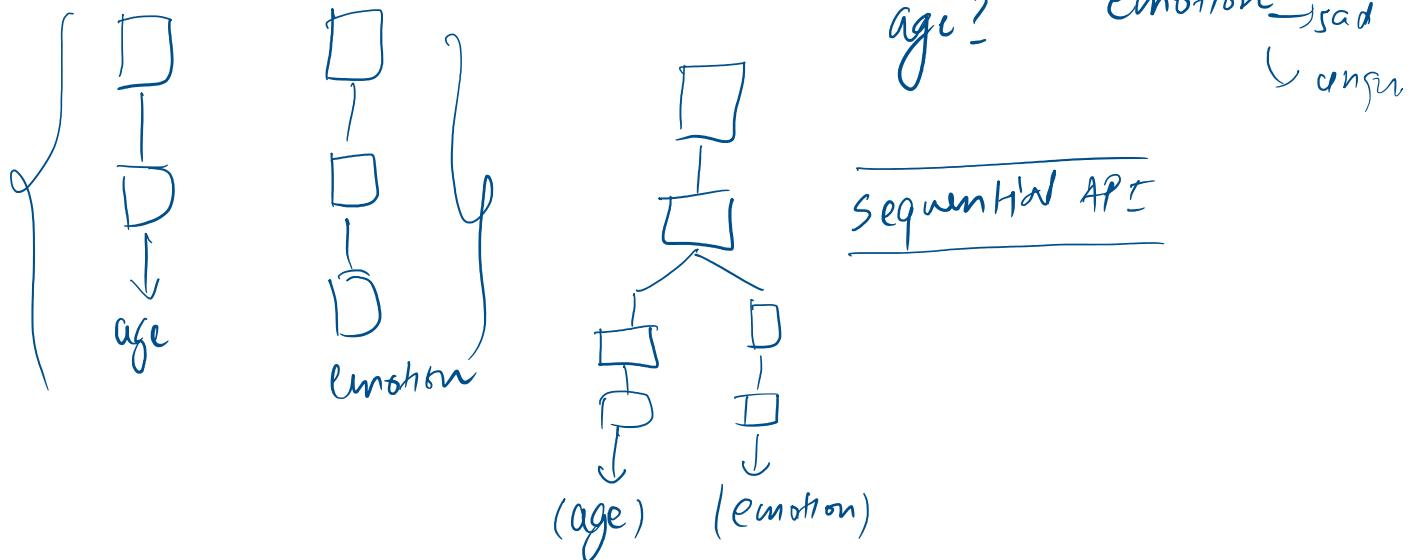
ANN \leftrightarrow CNN

↳ sequential - Keras

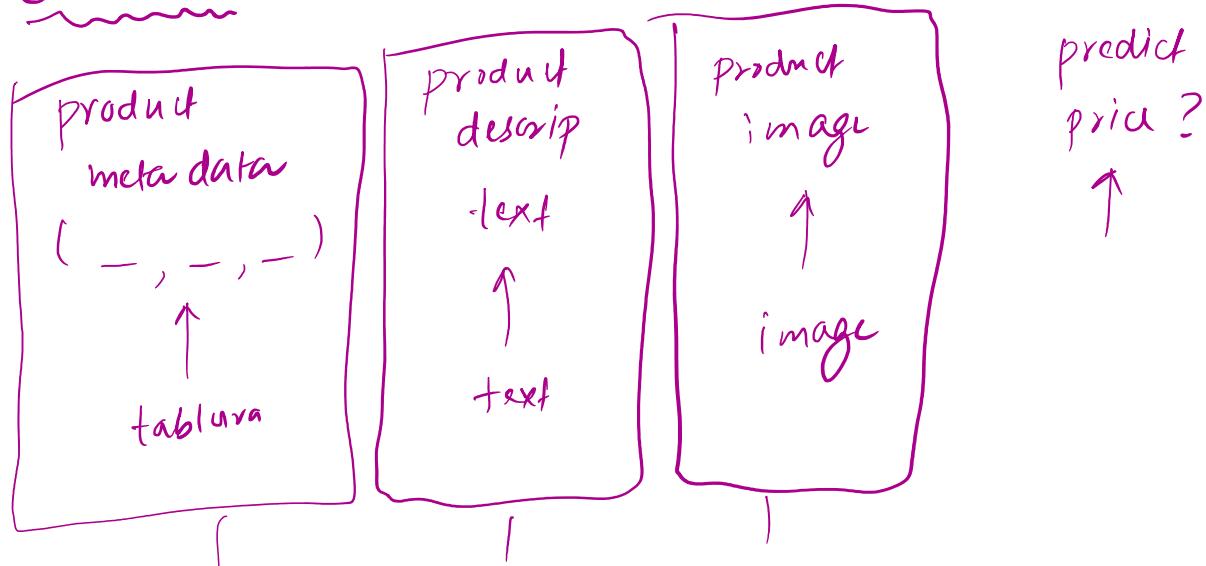
- 1 input
- 1 output
- linear

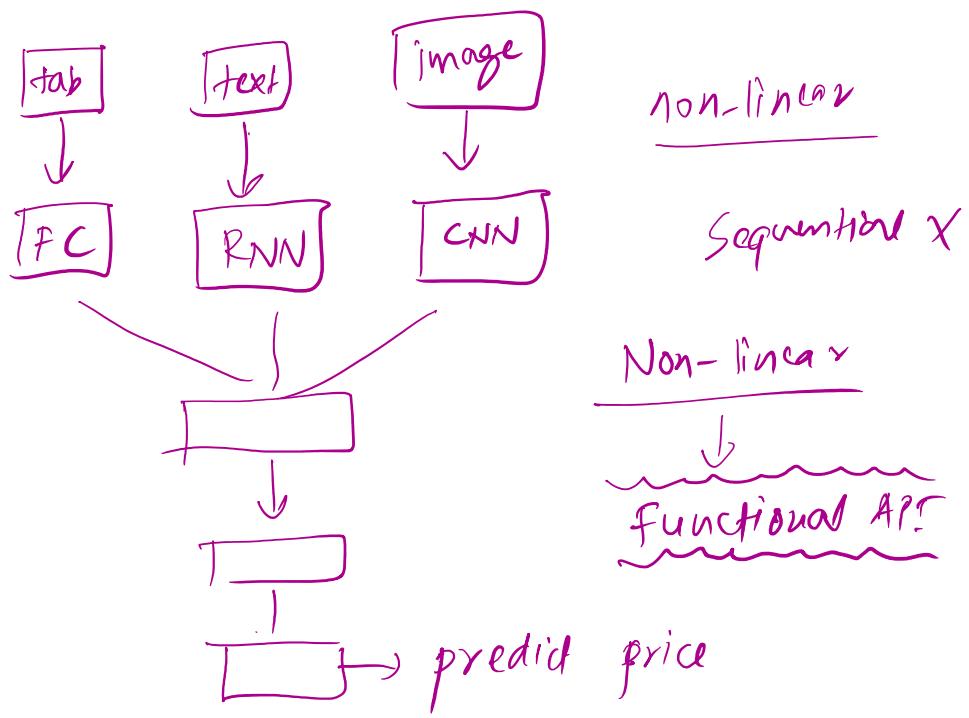


Example →  → image dataset → human faces → 25000 images



E-commerce



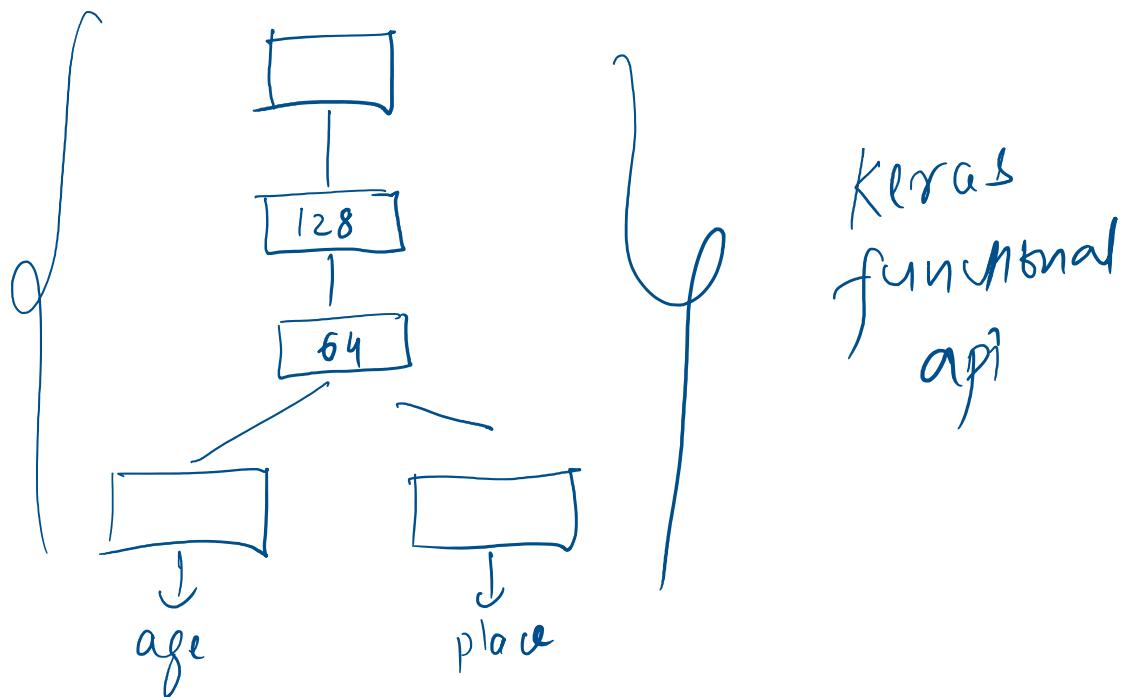


A Simple Example

14 October 2022 16:01

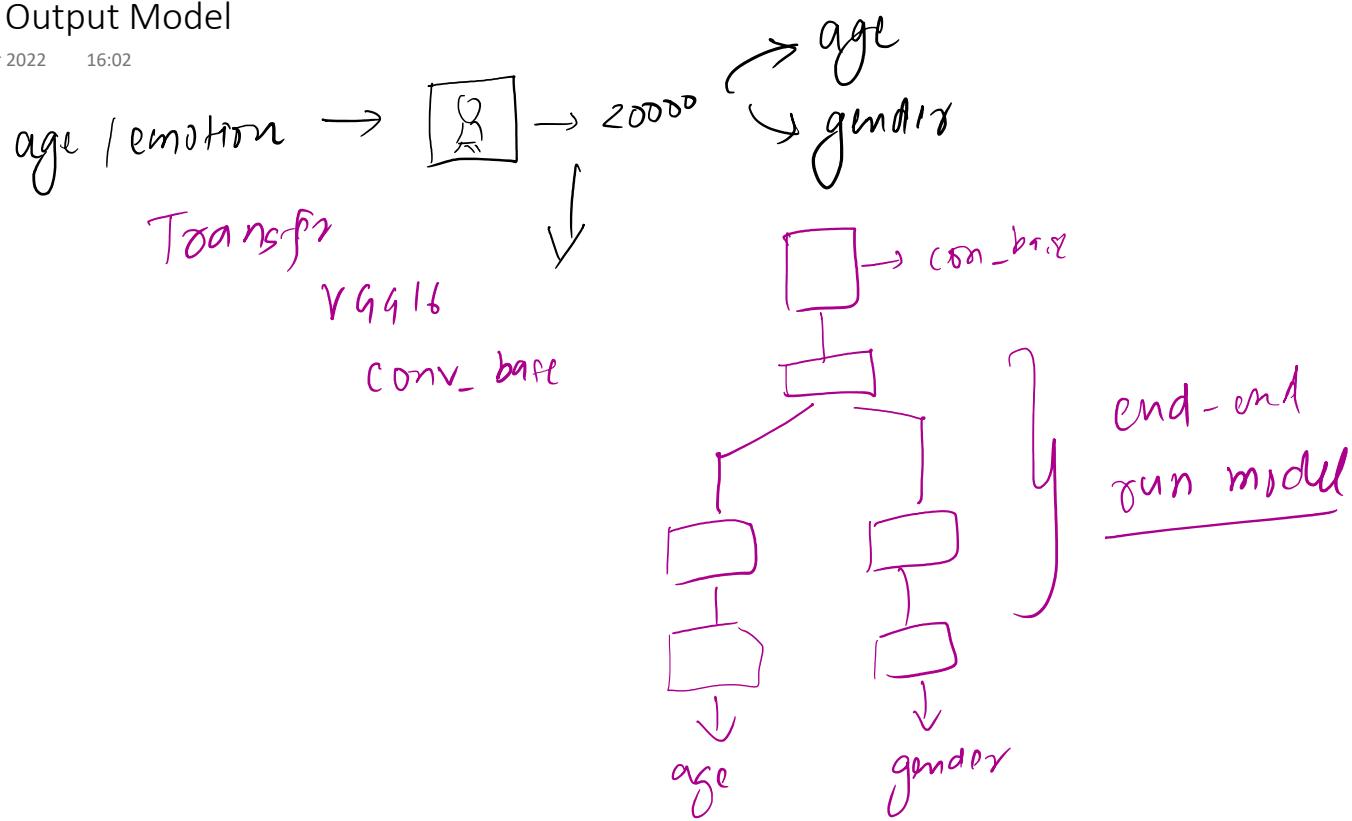
yearly salary	height	marital status
---------------	--------	----------------

3 cols
age
delhi/mumbai



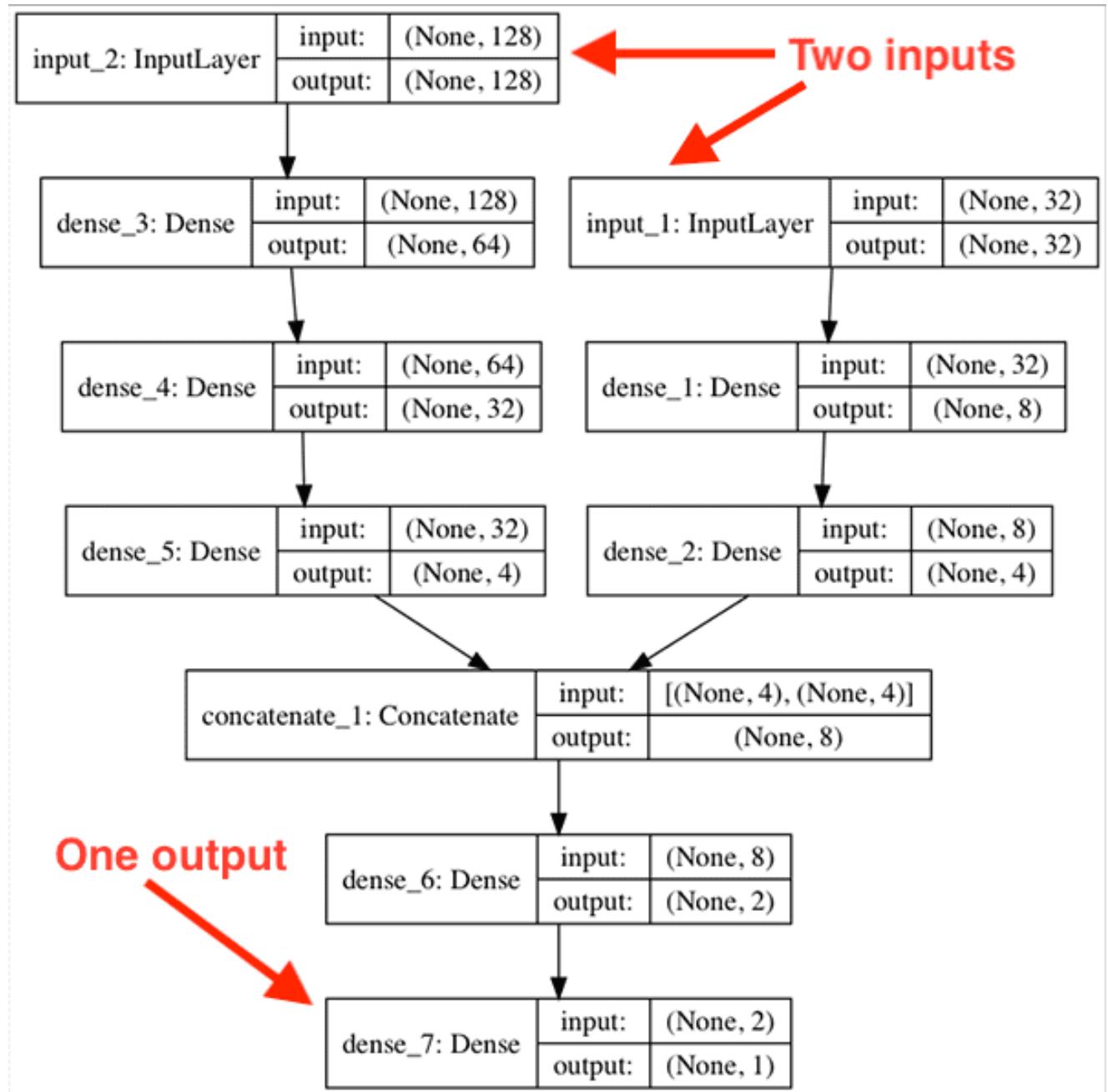
Multi Output Model

14 October 2022 16:02



Multi Input Model

14 October 2022 16:02



Shared Layers Model

14 October 2022 16:02

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	

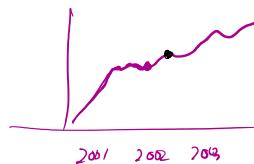
RNN → NLP → ML

CNN → images → computer vision

eg → text → sequential data

Hi my name is Nitish

Time series



Speech

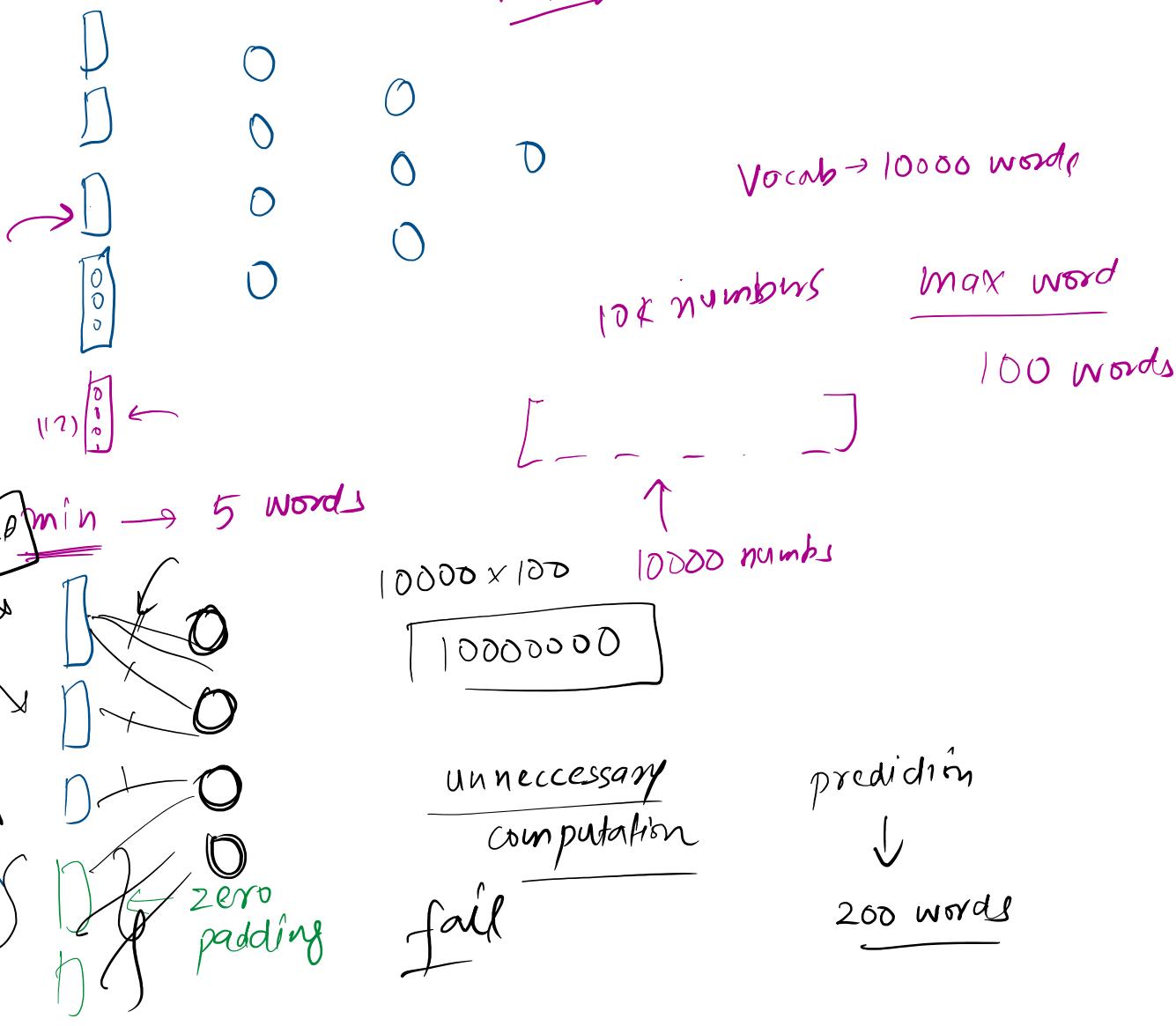
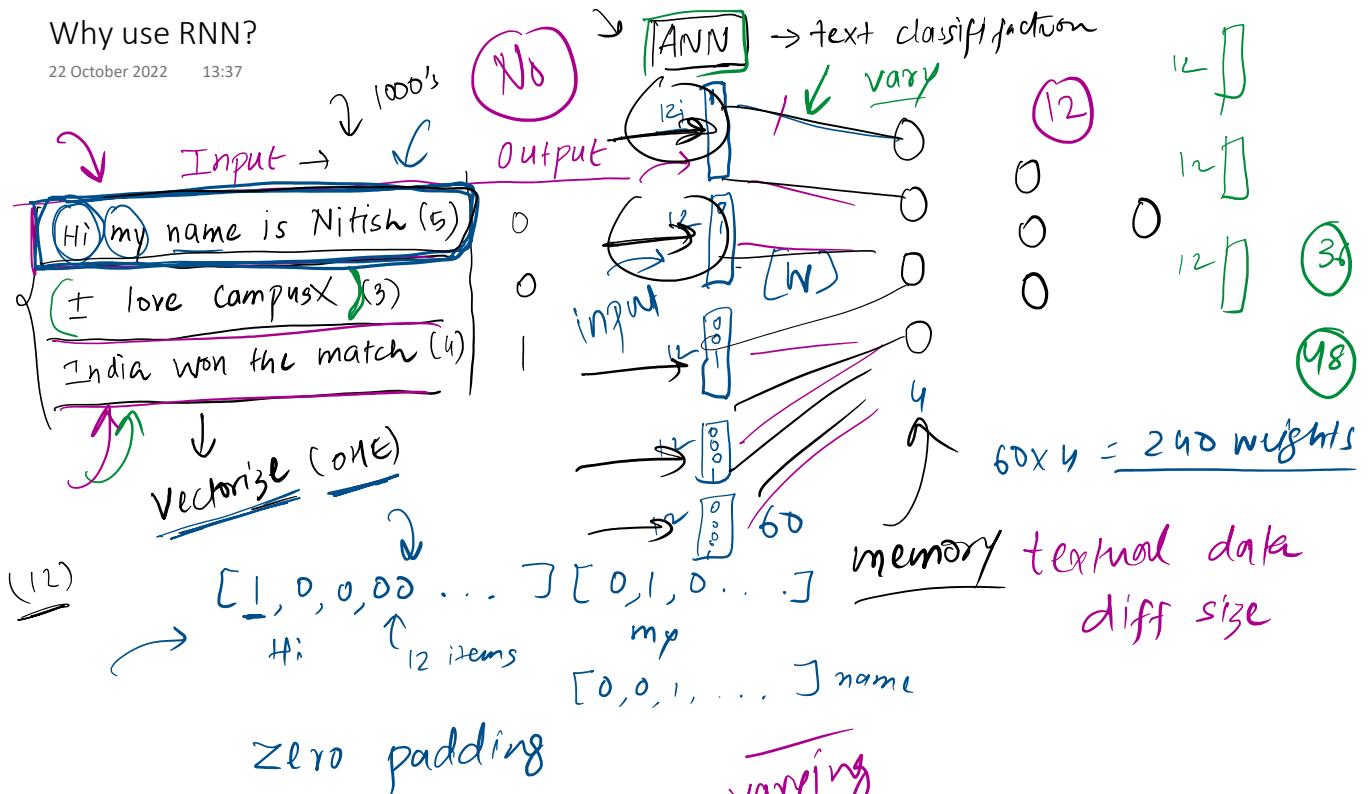
sequential

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

DNA sequence

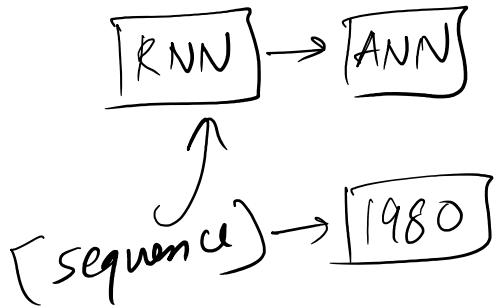
Why use RNN?

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

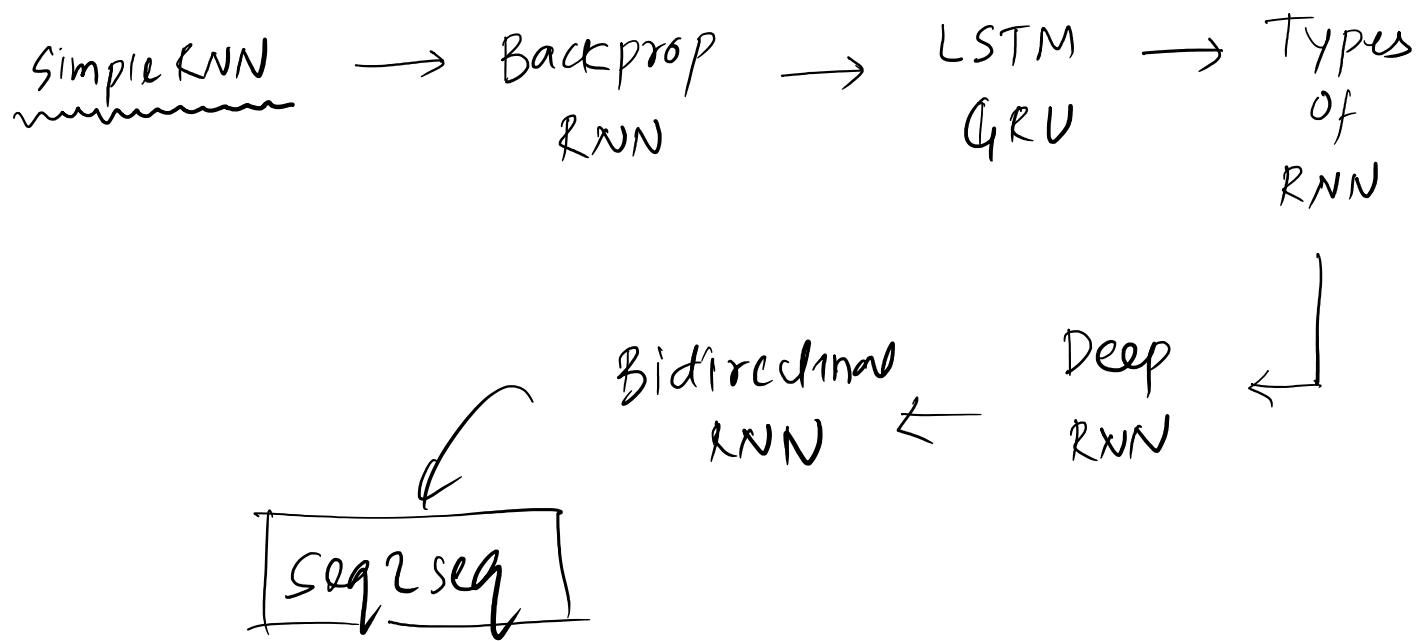


RNN Applications

22 October 2022 13:37

Roadmap

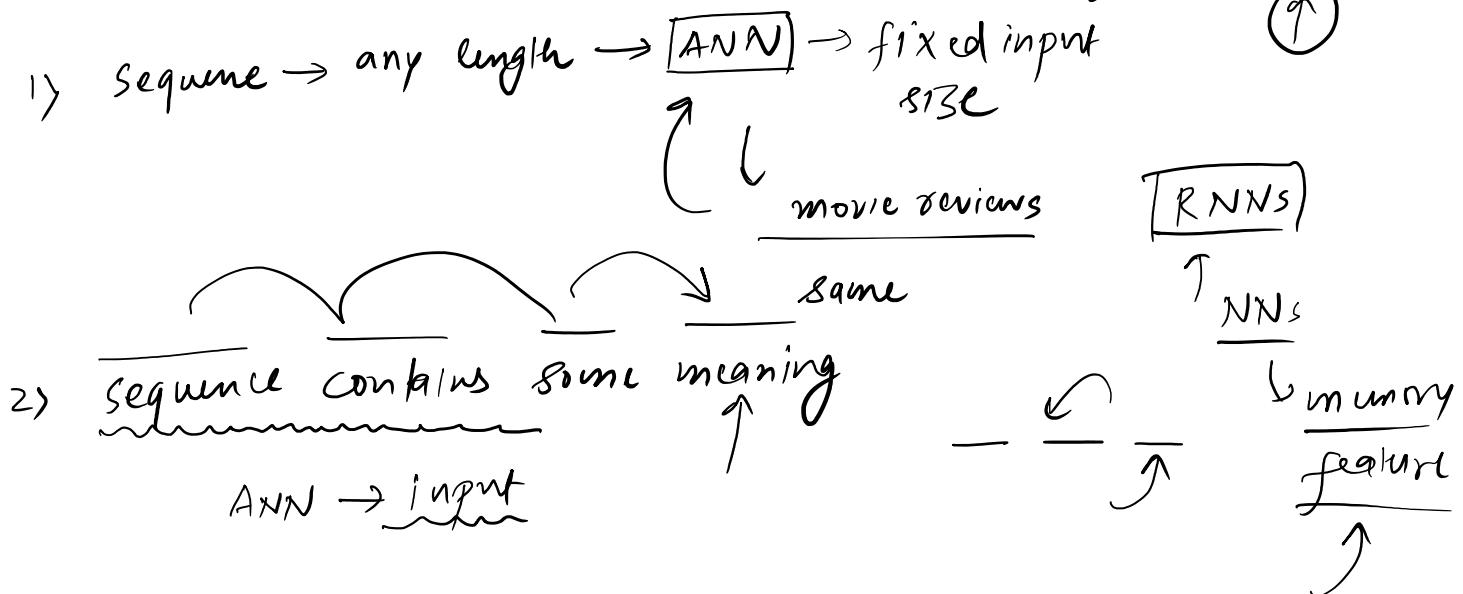
22 October 2022 13:37



Why RNNs?

29 October 2022 13:30

zero padding \rightarrow cost of computation



RNN architecture

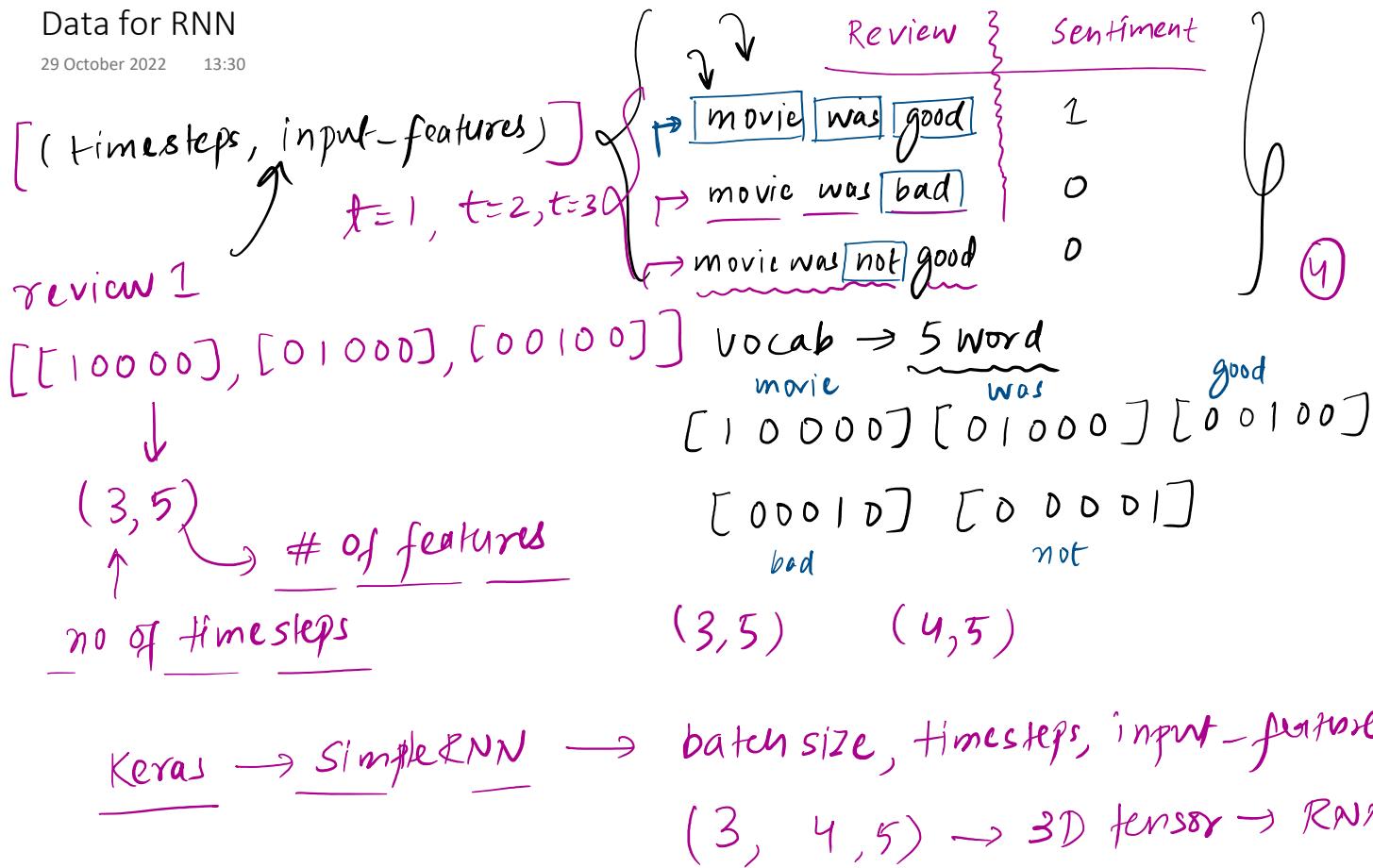
RNN forward prop \rightarrow prediction

input \rightarrow output

Codes \rightarrow Solidify

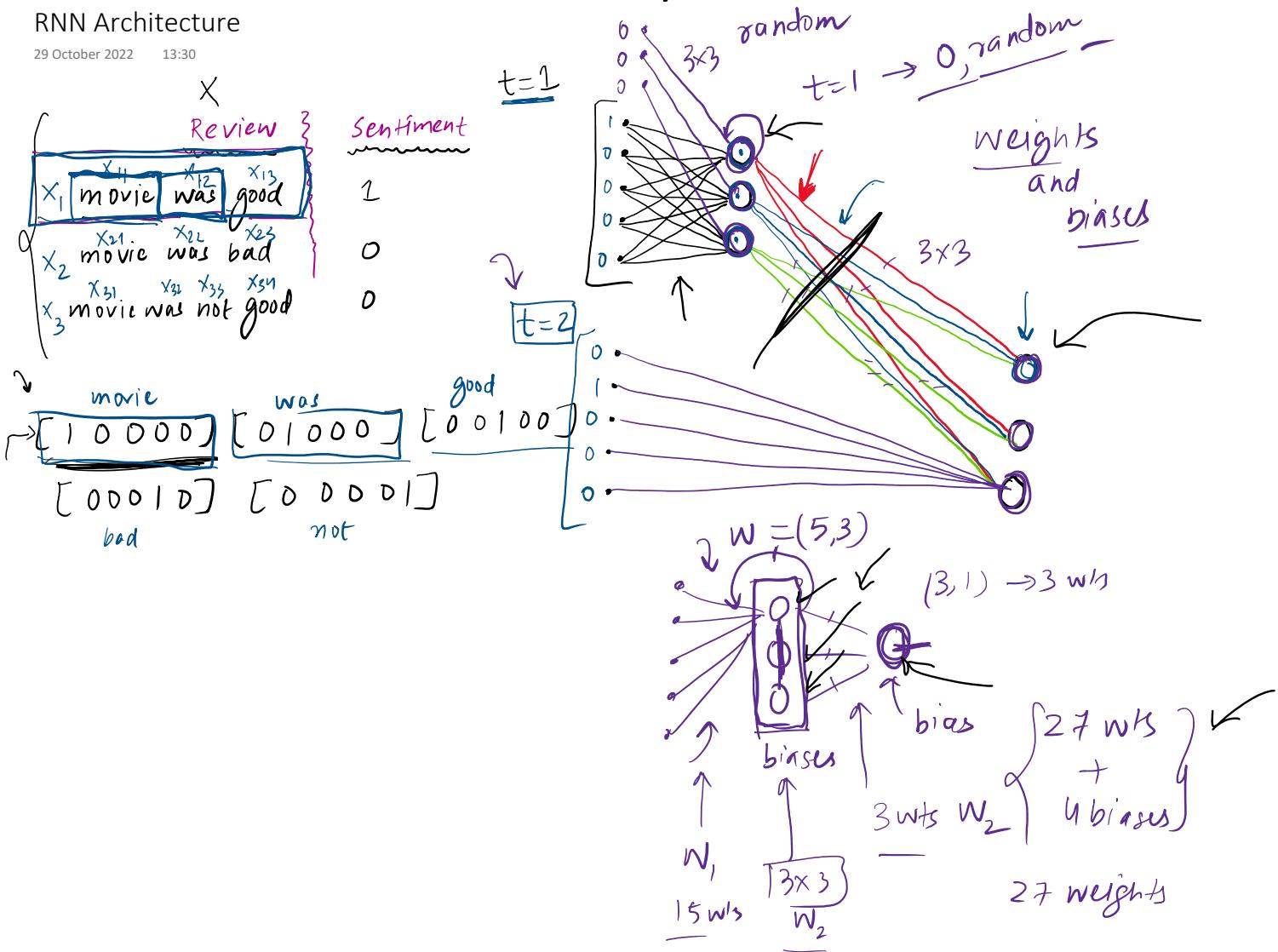
Data for RNN

29 October 2022 13:30



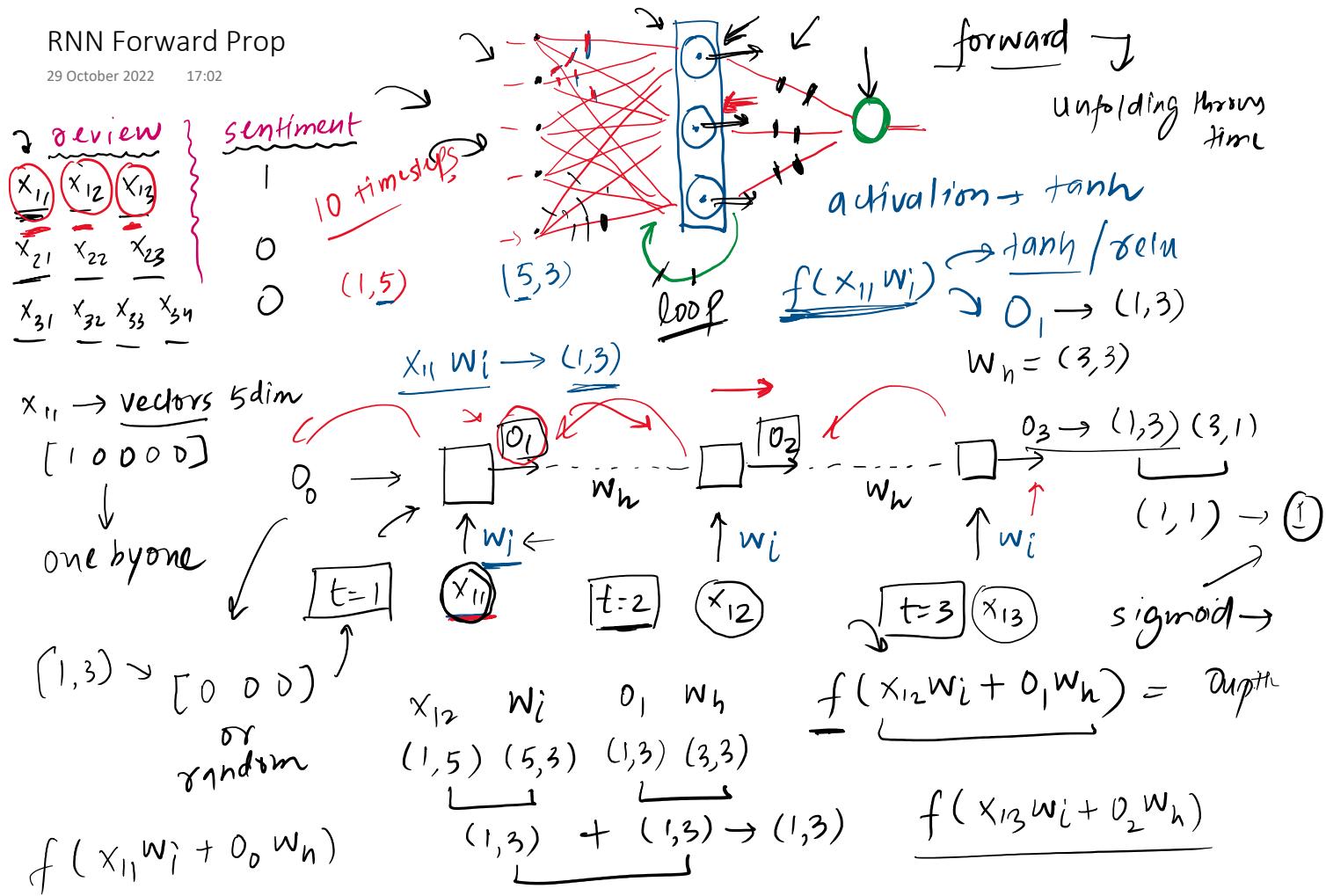
RNN Architecture

29 October 2022 13:30



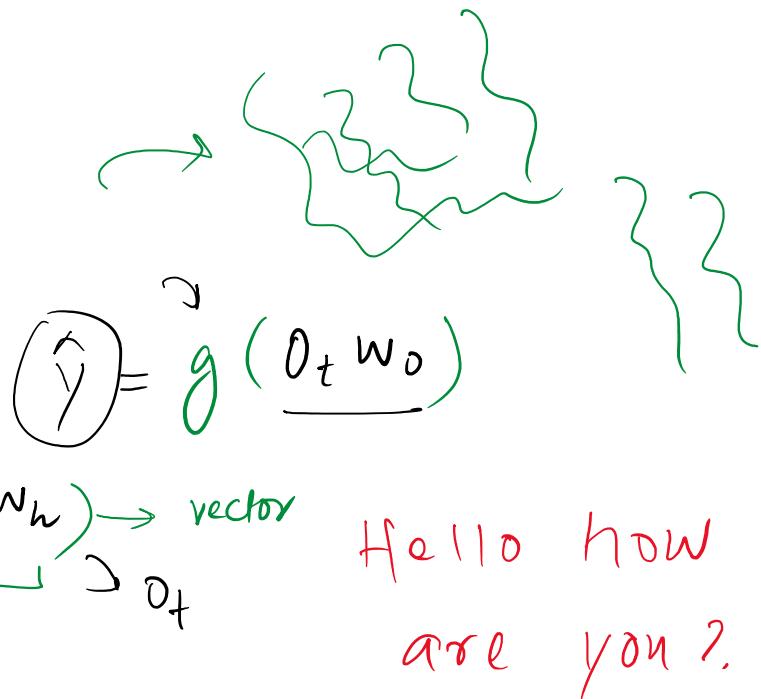
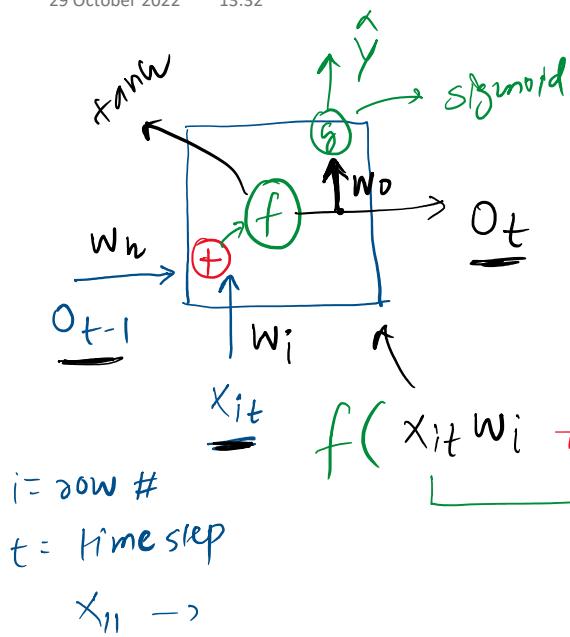
RNN Forward Prop

29 October 2022 17:02



Simplified Representation

29 October 2022 13:32

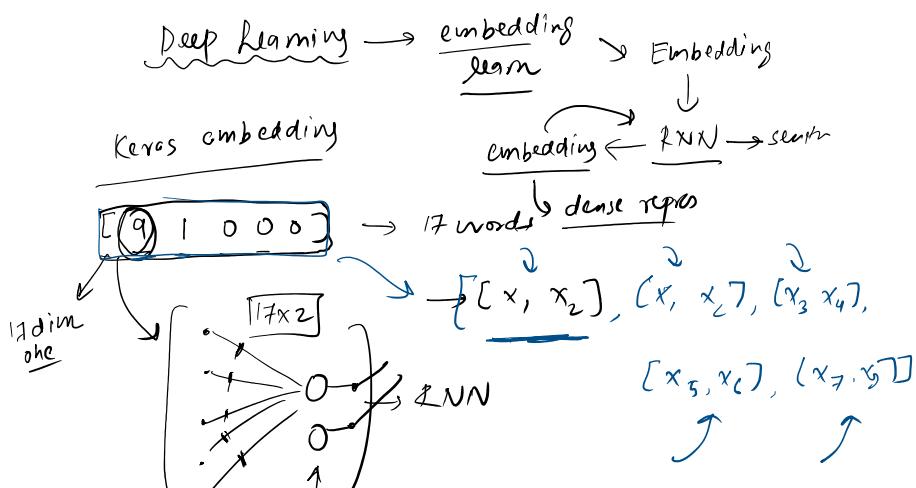
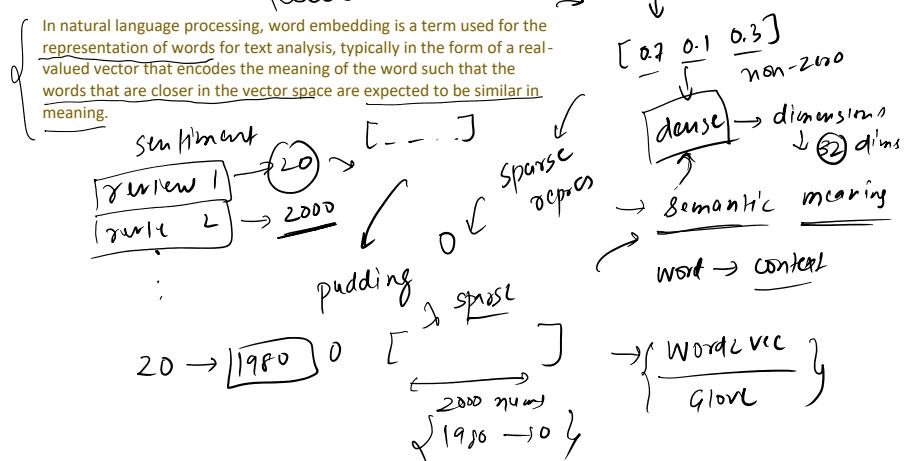
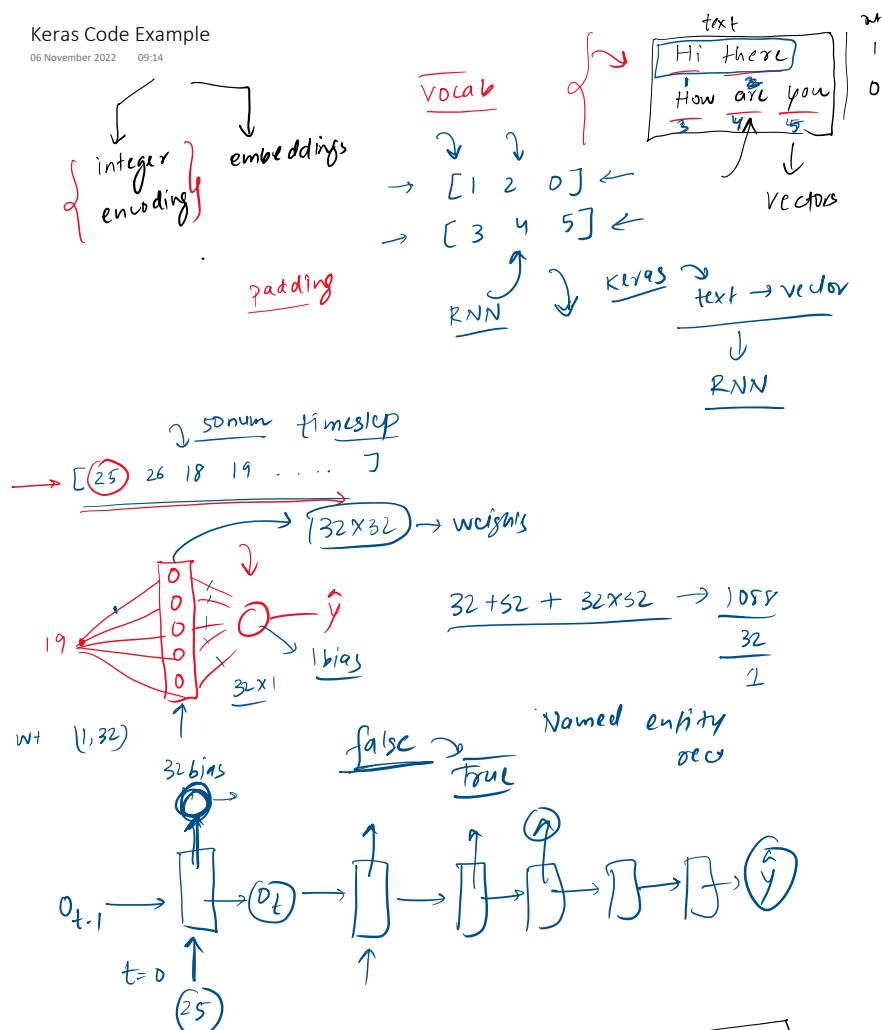


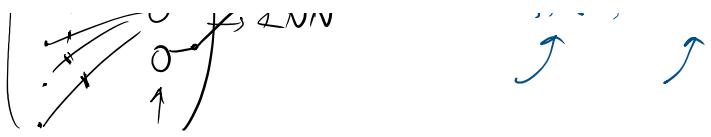
Code

29 October 2022 13:32

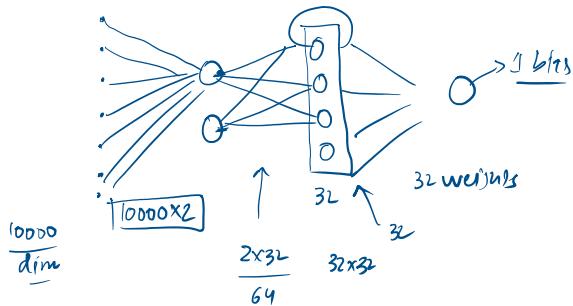
State and Memory

29 October 2022 13:33





17 nodes



population in n th year $\rightarrow x$

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

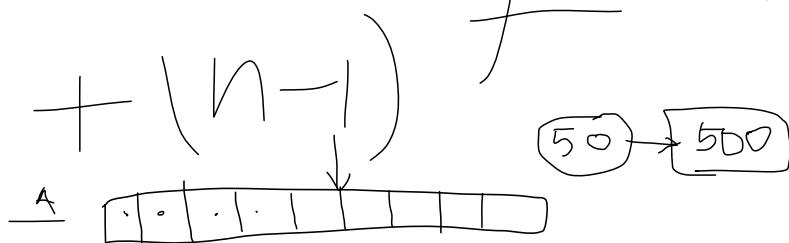
$$\frac{x + 0.1x}{= 10000}$$

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

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

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

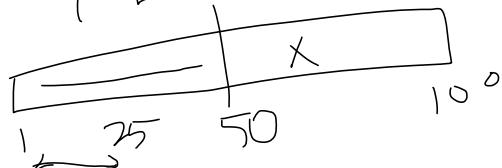
2+



A[35] → t sec

A[35] →

35 → $1 \times 4 \times 35$



O(n)

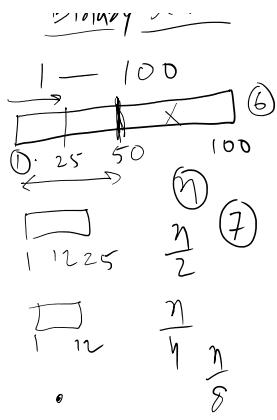
O(n²) → nested loops

input → 10 (100x100)
time $(10)^2 = 100$
 $O(\log(n))$

$\sqrt{(x)}$

Binary Search

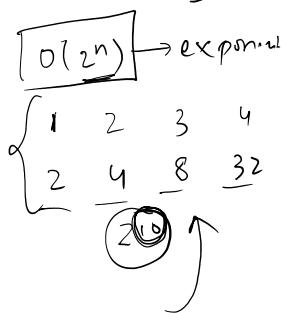
1 → 100 → 1(6)



1000
→

$O(n \log n)$ ↓

Sorting



{
 for i in range
 _____ $O(n)$
 for j in range
 _____ $O(n^2)$
 $O(n)$

$O(n+n) \rightarrow$

$O(2n)$

$\rightarrow O(n)$

$O(\cancel{n} + \underline{n^2}) \quad O(n^2)$

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]

$$5 + 1 \\ = '5'$$

$$345 // 10 \rightarrow 34$$

$$\underline{34 // 10} \rightarrow 3$$

1 1 1 5 1 .

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

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

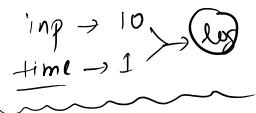
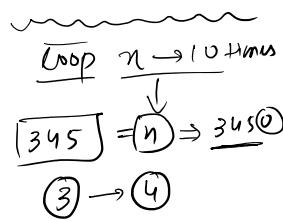
$$\underline{34 \cdot 110} \rightarrow (3)$$

$$3^1 = 0$$

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

`digits[3]`

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



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

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

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

$$O(n)$$

$$O(n)$$

$$O(1000000)$$

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

$$\times$$

$$O(n^2)$$

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

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

$$\begin{matrix} O(n) \\ O(n) \end{matrix} \quad n=100$$

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

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

$$j=1 \rightarrow (2)$$

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

$$j=3 \rightarrow 8$$

$$j=4 \rightarrow 32$$

$$\frac{n}{2} \times \log n$$

$$\boxed{n \log n}$$

O(1)) → constant

$$n = \textcircled{345}$$

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

5

30

4

3 3

10

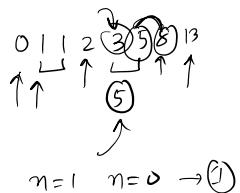
3450 → ④

los

inp → 10 100
out → 1 2 3

fibonacci

↳ function
↳ recursion



$\text{fib}(n)$

↓
function calls

input → ⑦ → 10

#fcalls →

$$\text{fib}(\overset{n=3}{(3)})_0$$

$$\checkmark \quad \underline{\text{fib}(2)}^{\circledcirc} \quad \underline{\text{fib}(1)}^{\circledcirc} \rightarrow 1$$

fib(1) fib(0)

1) 5
exponentials -

$$\frac{f'(b)}{1} > 0$$

2
4

$\begin{array}{r} & \swarrow & \searrow \\ & r_1 & \\ \downarrow & & \downarrow \\ 2 & r_1 & 1 & 0 & 0 \\ & & & & \text{shallow} \end{array}$

char

1

15

$O(2^n)$ (5)

input 1 2

15 2 4

$$O(2^n)$$

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

$$n = 50, 100, 500$$

↑ weeks

exponential
↳ days/weeks

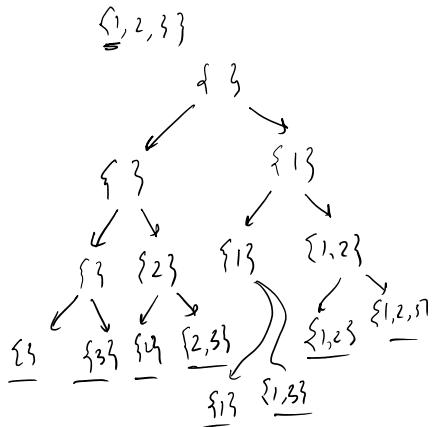
subset
power set $\rightarrow O(?)$

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

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

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

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



reduce → divide → log

increase → multi → exp

→ exponentiation

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

$$2^3 = 8$$

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

$$2^4 = 16$$

$$O(2^n)$$

$$O(?)$$

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

$$n > 0$$

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

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

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

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

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

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

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

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

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

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

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

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

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

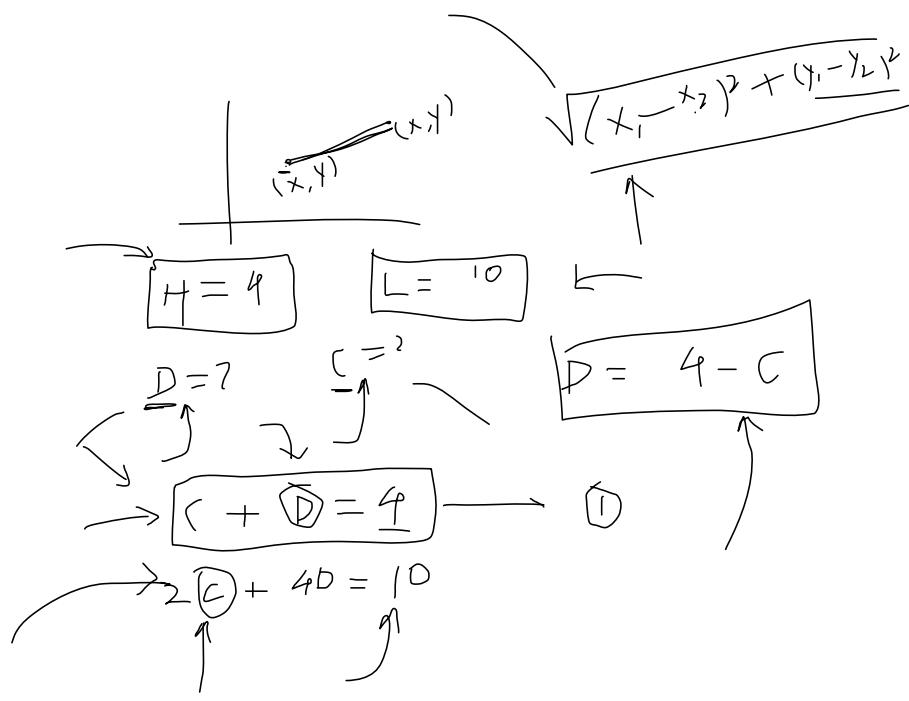
$$= 2^3 T(n-3) - 2^2 - 2^1 - 2^0$$

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

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

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

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



15 5 $15^2 + 5^2$

\leq

$$1 \quad -$$

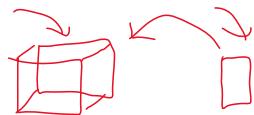
$$\begin{matrix} & 15 \\ a & 12 & 5 \\ & 1^2 + 2^2 + 3^2 + 4^2 \times 5^2 \end{matrix} \quad n = 5$$

$$2 \quad 6 \quad \begin{bmatrix} 3, 6 \end{bmatrix} \rightarrow 5^{+n}$$

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

$$d = 6 - 3$$

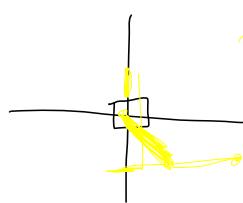
$$\begin{matrix} 1 & 0 & 0 & 0 \\ , & 0 & 6 & 0 \end{matrix} \quad \frac{2}{3} \cancel{+} \frac{4}{5} = \boxed{\frac{10 + 12}{15}} \rightarrow \frac{22}{15}$$



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

$$\boxed{1 \quad 0 \quad 0 \quad 0}$$

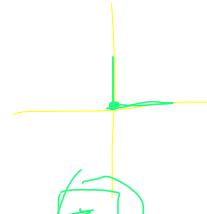
$$\begin{matrix} 1 & 0 & 0 & 2 \\ , & 2 & 2 & 2 \end{matrix} \quad 5^1 = \boxed{5} \xrightarrow{5 \times 0 \times 1 \times 2 \times 1}$$



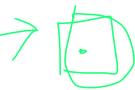
$$\boxed{4P \rightarrow 1 \\ D - 3 \\ R - 4 \\ L - 3}$$

$$\begin{matrix} a = [1, 2] \\ a = b \\ \hline a = b [;] \end{matrix}$$

$$q \rightarrow$$

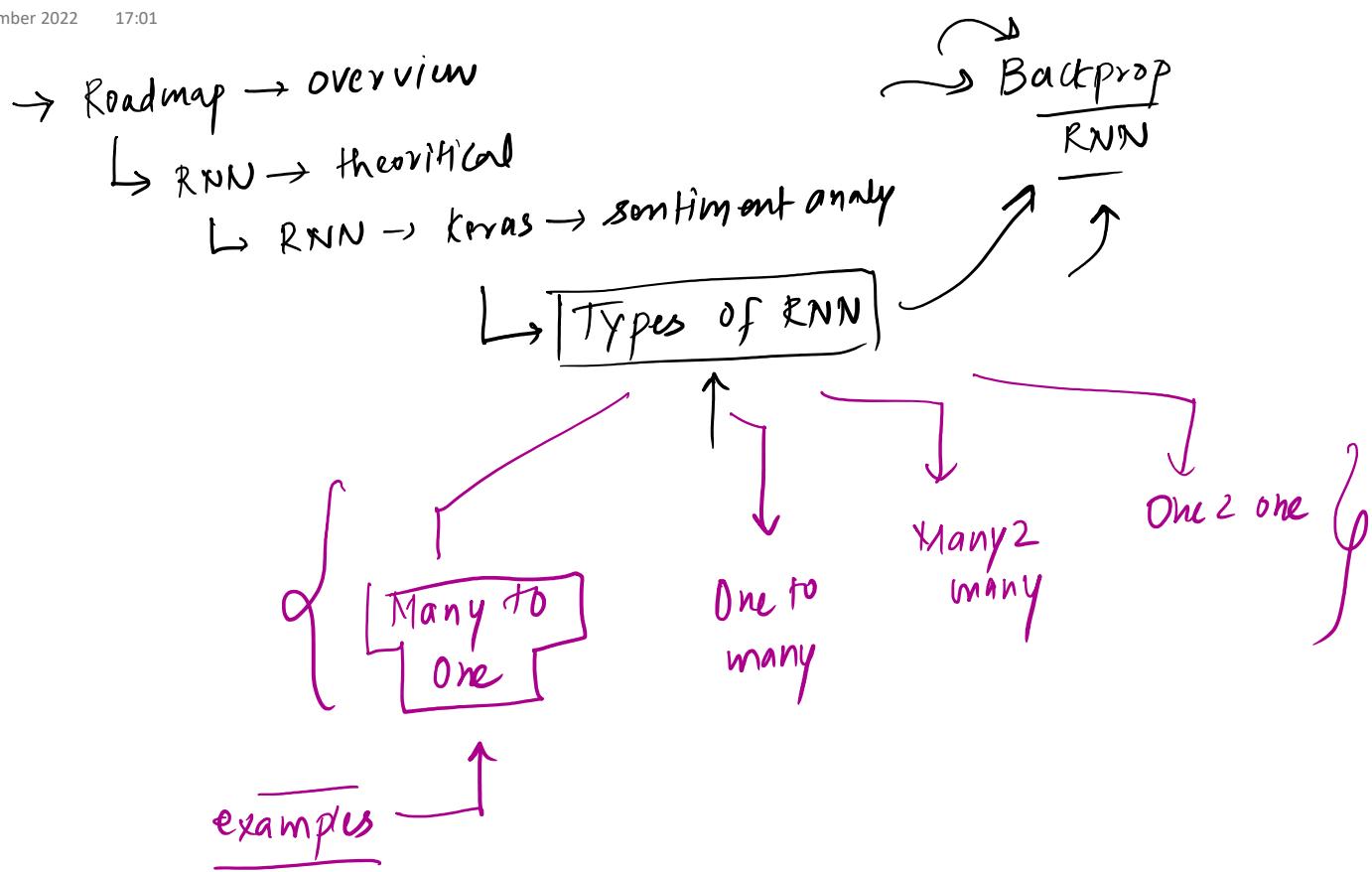


$$b \rightarrow$$



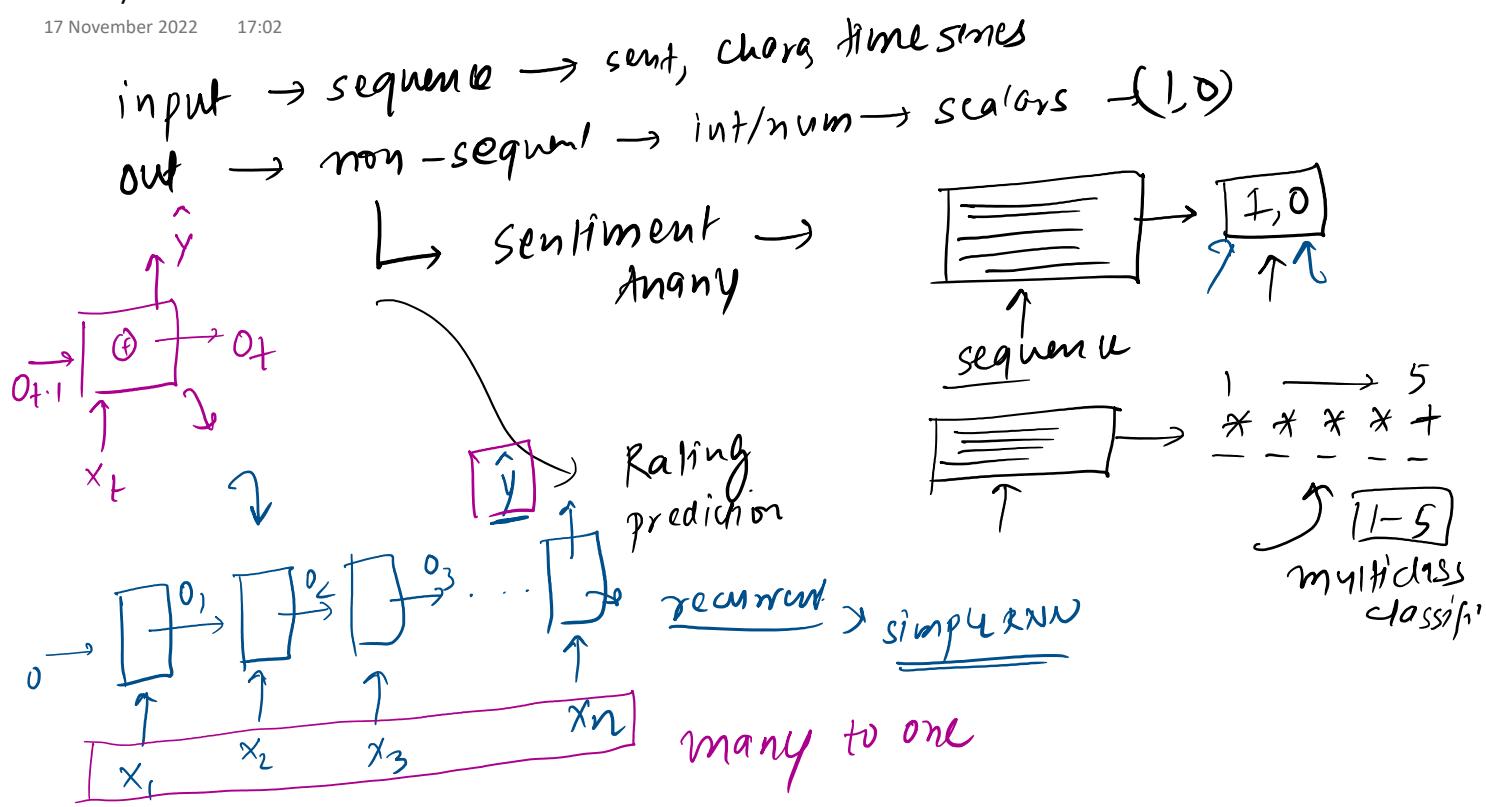
Till Now

17 November 2022 17:01



Many to One

17 November 2022 17:02

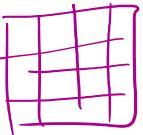


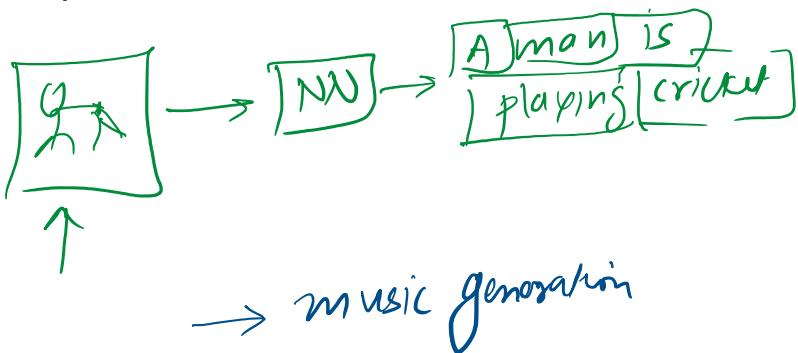
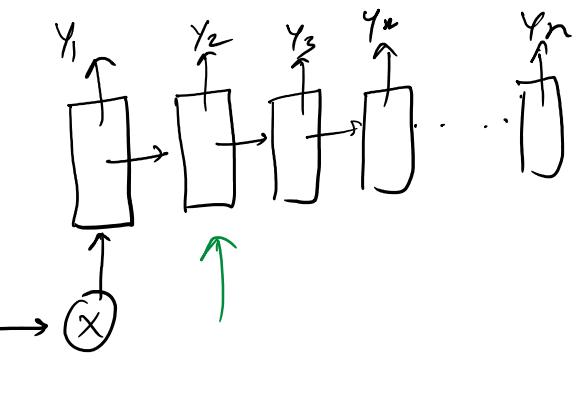
One to Many

17 November 2022 17:02

→ normal non sequential
↳ 

→ Output → sequences
image captioning

 → textual
depres.



Many to Many

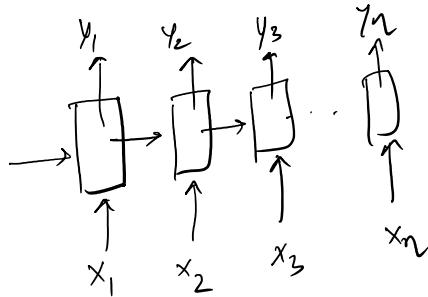
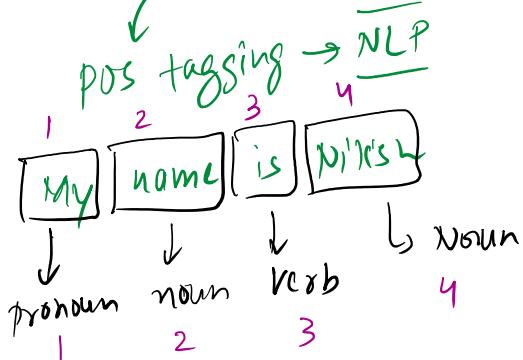
17 November 2022 17:02

input \rightarrow segment \rightarrow seq2seq
 out \rightarrow sequence

Same length

Variable length

input seq == output seq

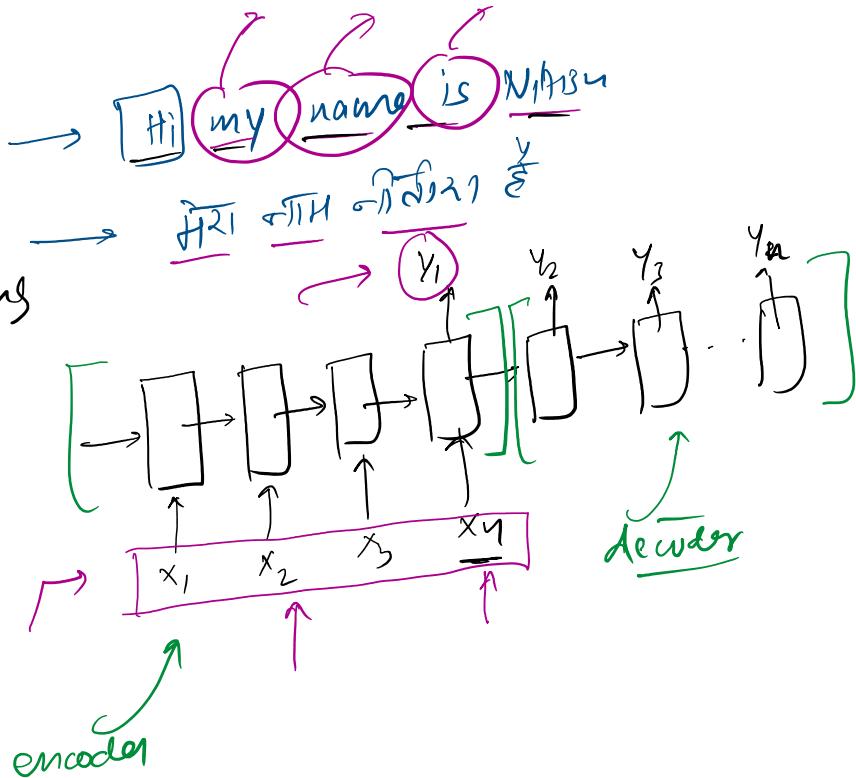


same length
 many \leftrightarrow many
 RNN

NER
 Lets meet at 7pm at the airport

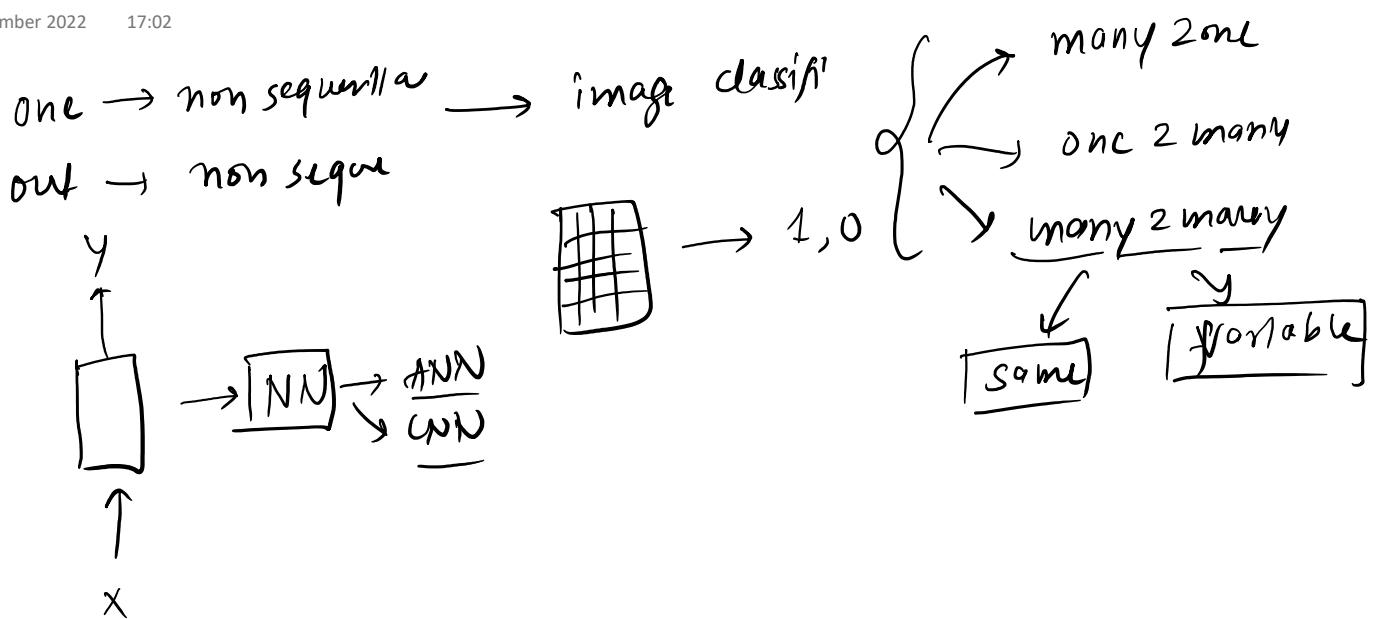
Variable length
 machine trans
 \hookrightarrow 1 lang \rightarrow 2 lang
 google translate

encoder
 decoder



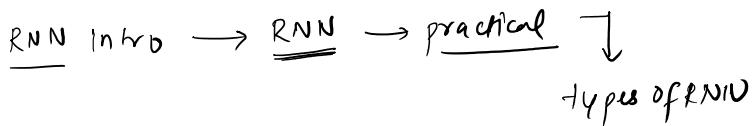
One to One

17 November 2022 17:02



Backpropagation in RNN

01 December 2022 16:43

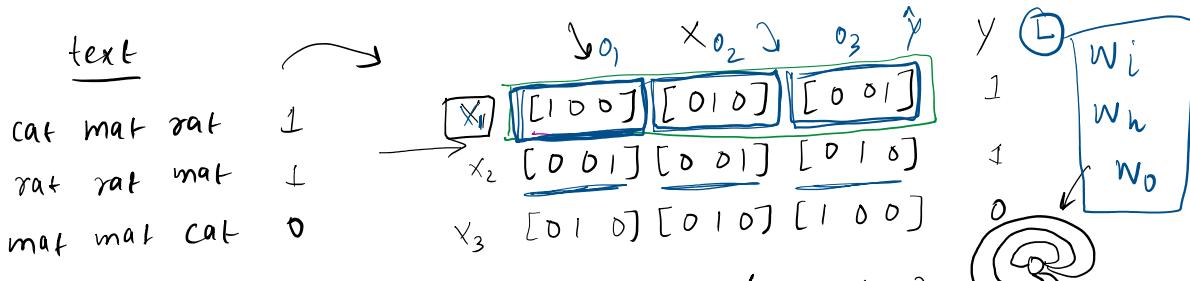


Many to One RNN

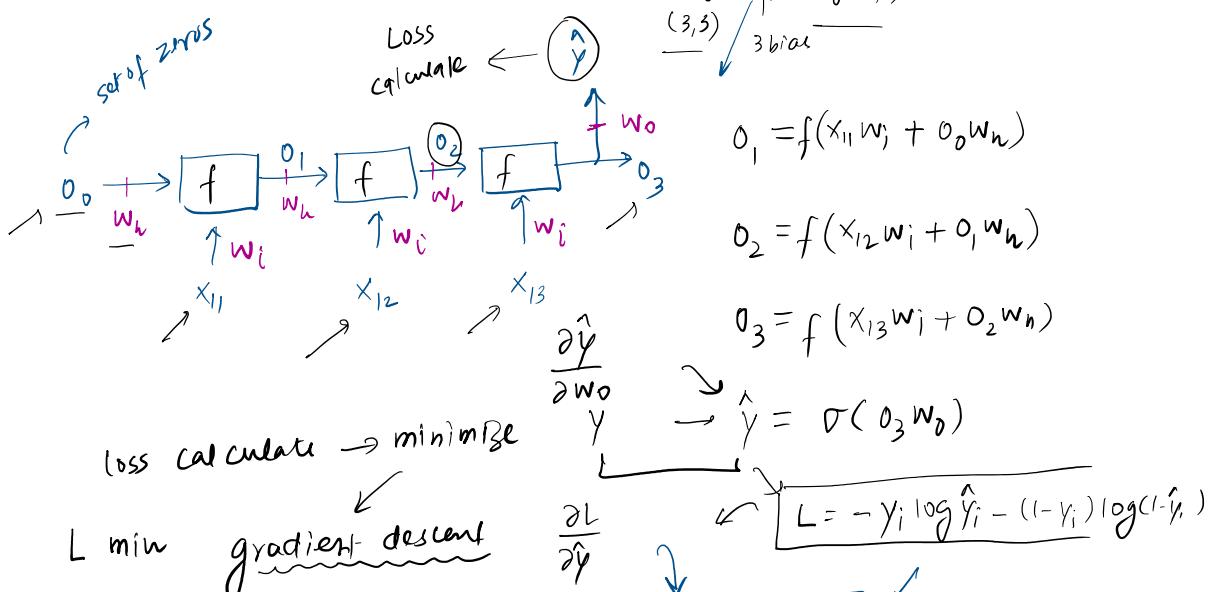
Sentiment Analysis

text → I/O

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



forward prop



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

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

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

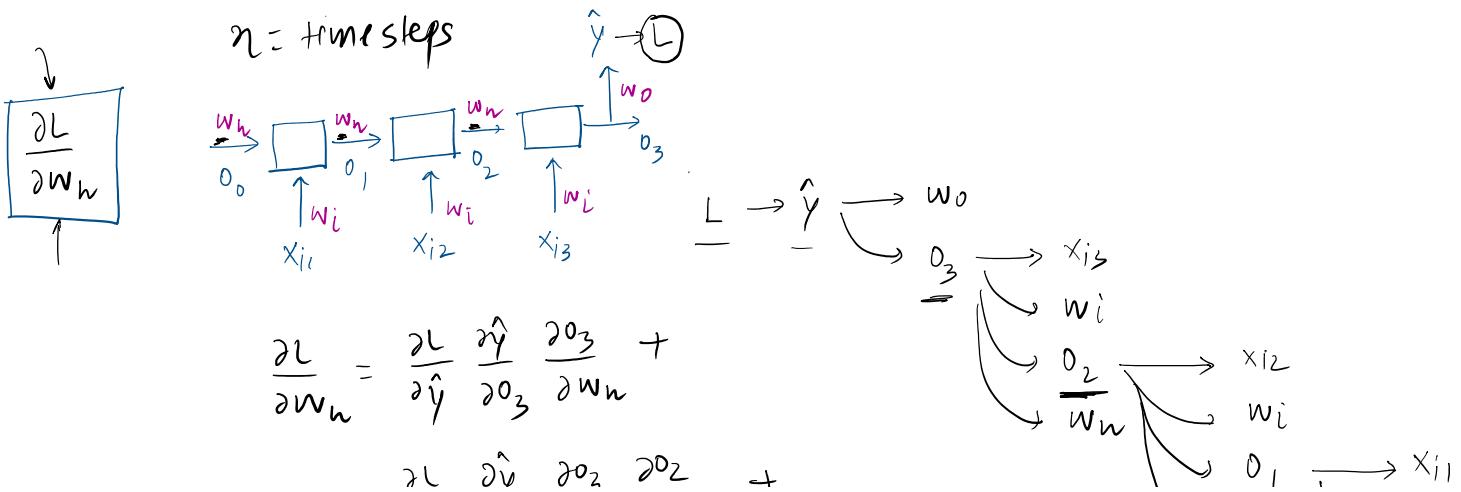
$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_3} \frac{\partial o_3}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial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

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

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

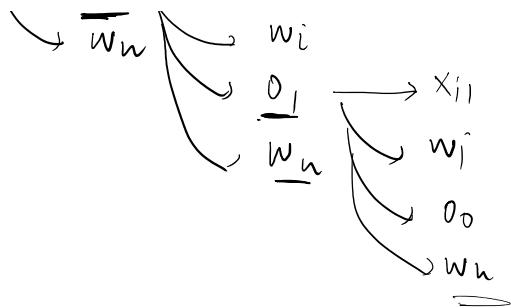
$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \left[\frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial w_i} \right]$

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



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

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



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

$n = \text{timesteps}$

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

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

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

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

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

Problem with RNN

19 December 2022 16:33



Problem #1 → Problem of long term dependency → Vanishing

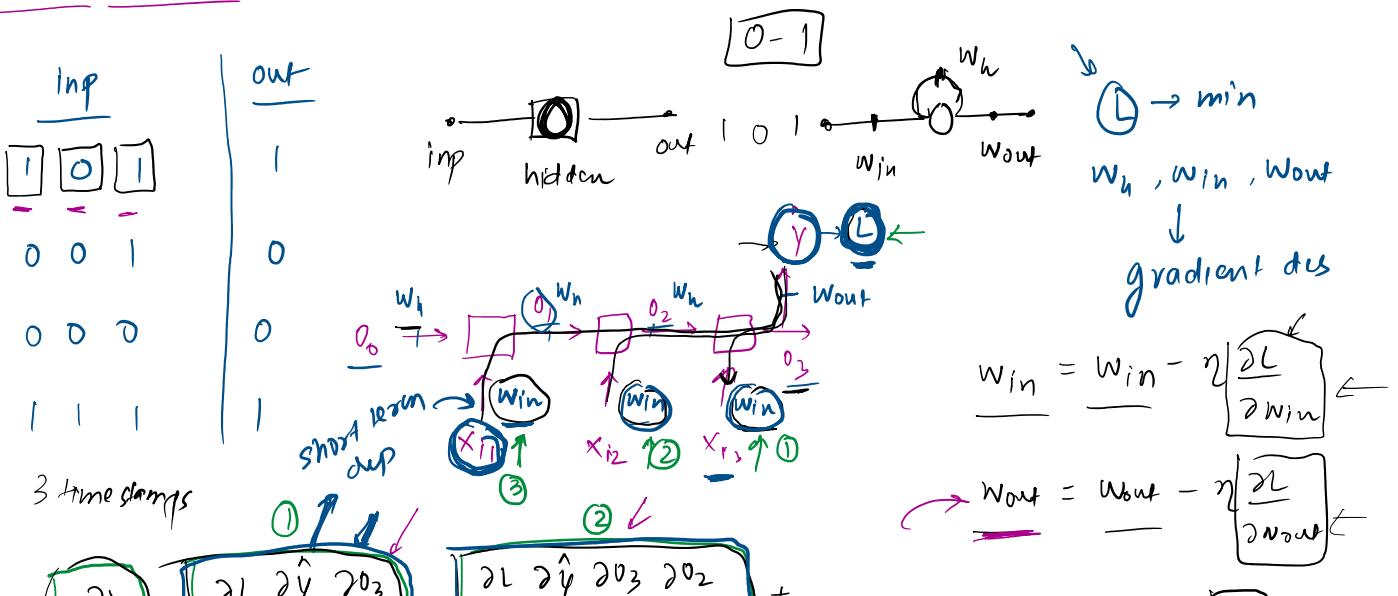


Diagram illustrating the backpropagation through time (BPTT) for an LSTM cell. The diagram shows the computation of gradients for hidden states h_t and cell states c_t over multiple time steps.

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. A red arrow labeled "inf" points to the forget gate's output, indicating it can become infinite.

Top Right: A diagram showing the update rule for the hidden state h_t at time step t , where the gradient is scaled by the forget gate's value (f_t) to prevent vanishing gradients.

Middle Left: A diagram showing the gradient flow for the hidden state h_t at time step t through the forget gate f_t and the hidden-to-hidden weight W_{hh} .

Middle Right: A diagram showing the gradient flow for the hidden state h_t at time step t through the hidden-to-hidden weight W_{hh} and the forget gate f_t .

Bottom Left: A diagram showing the gradient flow for the hidden state h_t at time step t through the hidden-to-hidden weight W_{hh} and the forget gate f_t .

Bottom Right: A diagram showing the gradient flow for the hidden state h_t at time step t through the hidden-to-hidden weight W_{hh} and the forget gate f_t .

Sol ④

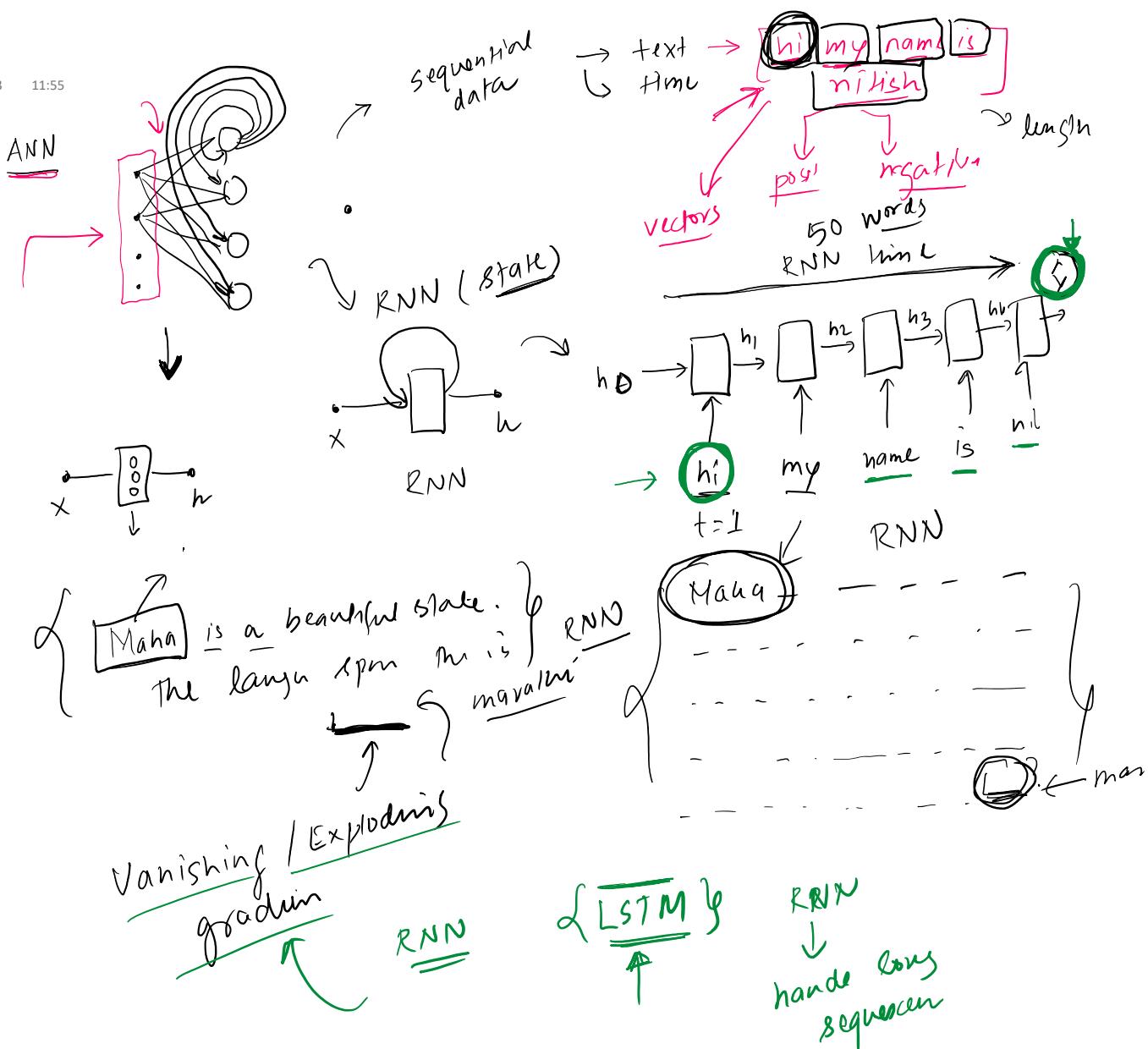
- 1) Diff activation \rightarrow relu / leaky relu
- 2) Better weight init
- 3) Skip conn
- 4) LSTM

Problem #2 \rightarrow Unstable Training (Exploding gradients)

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

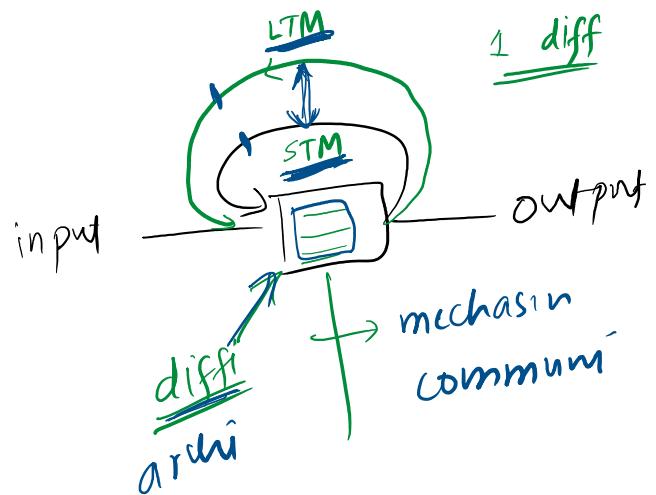
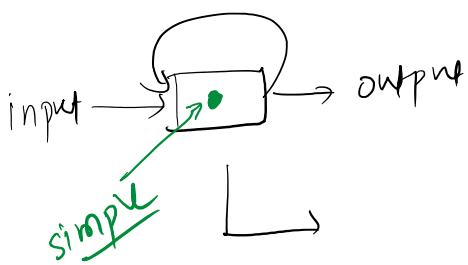
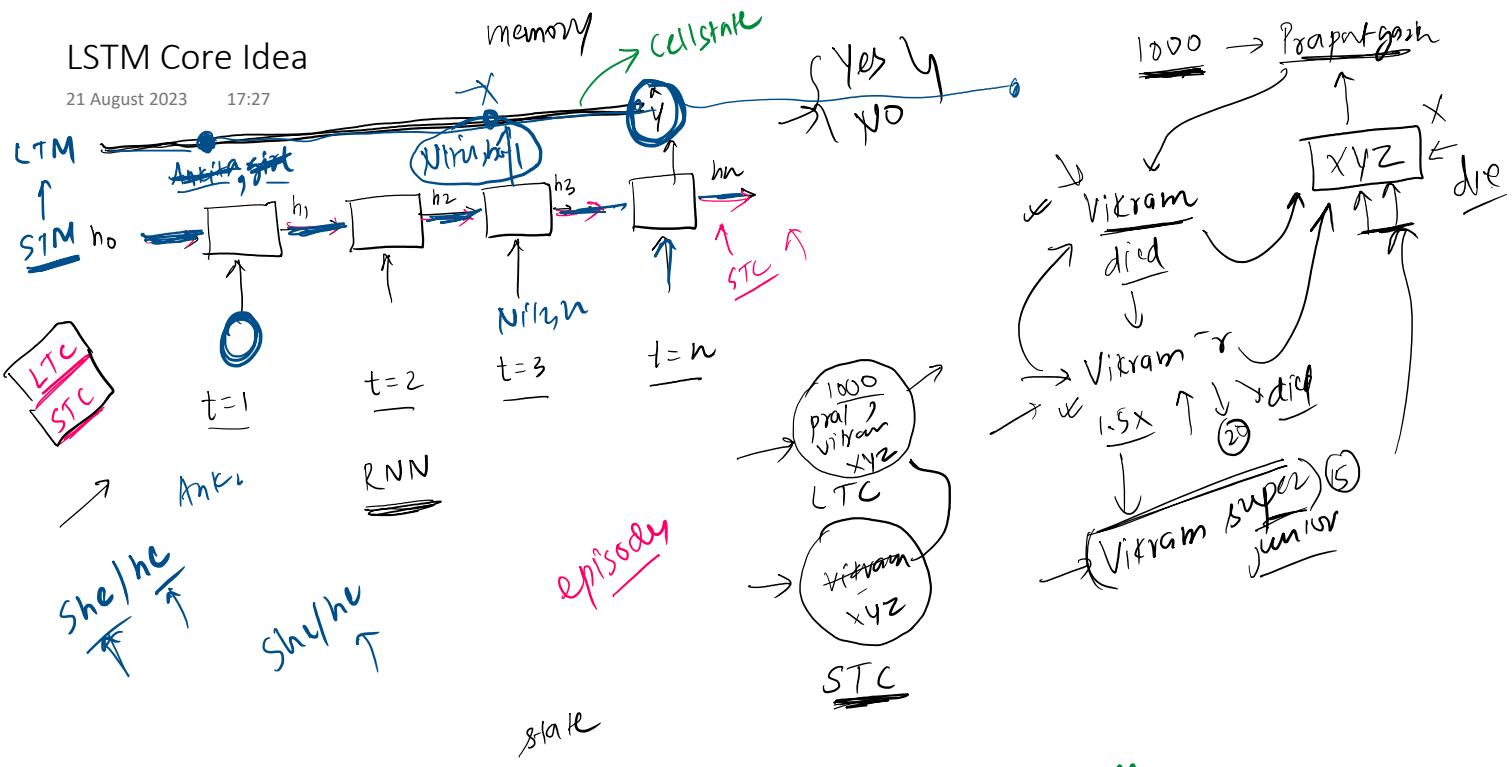
Recap

21 August 2023 11:55



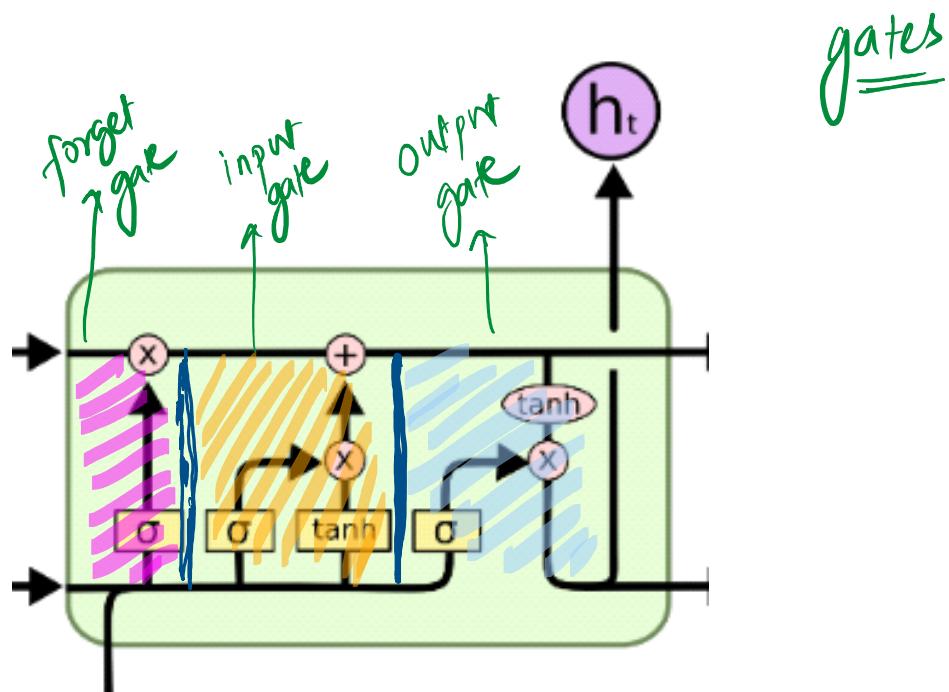
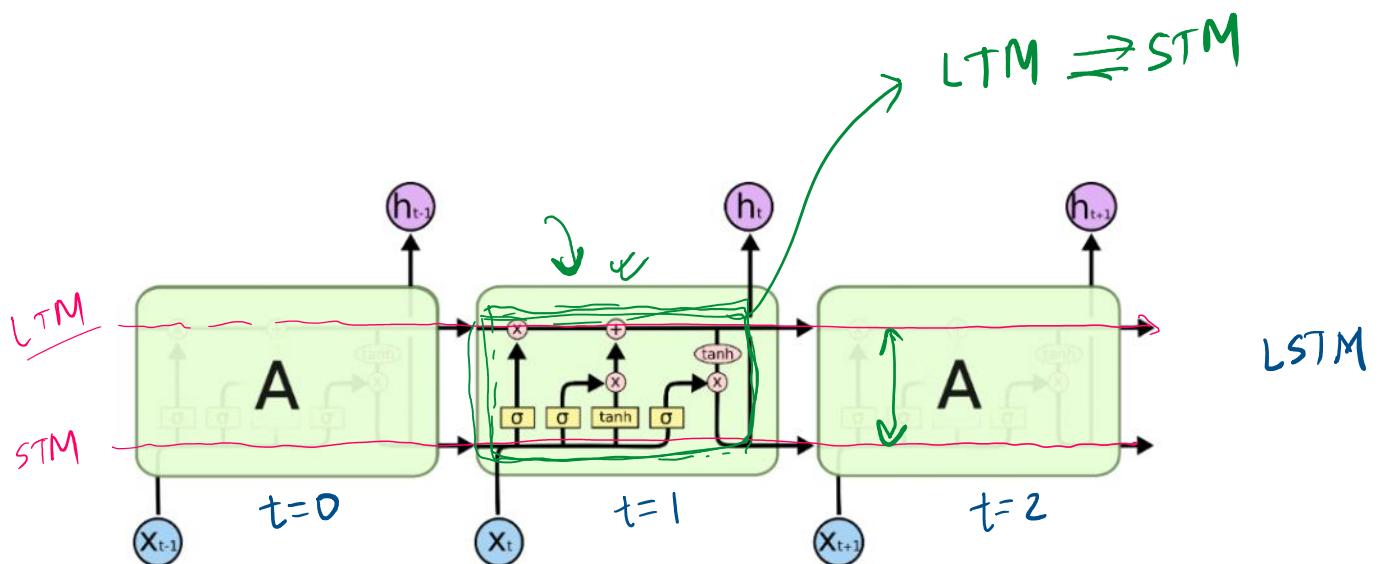
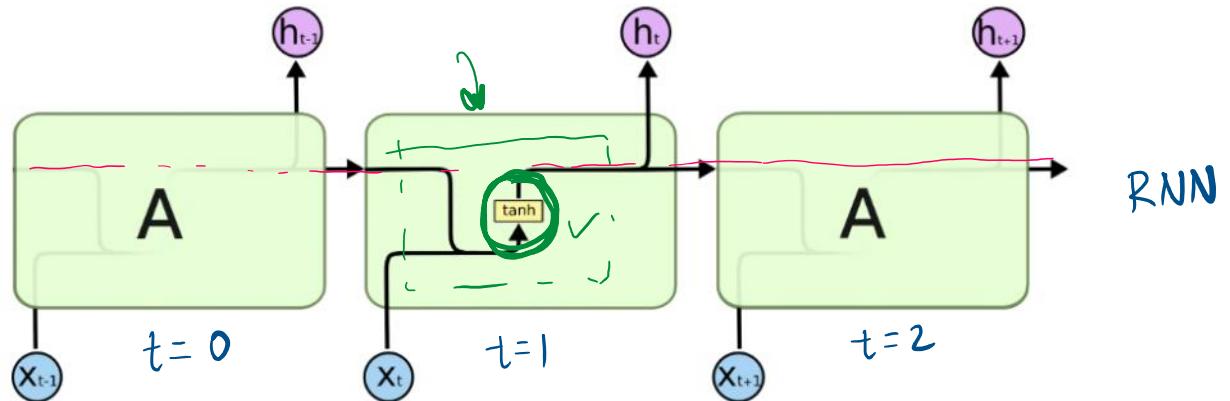
LSTM Core Idea

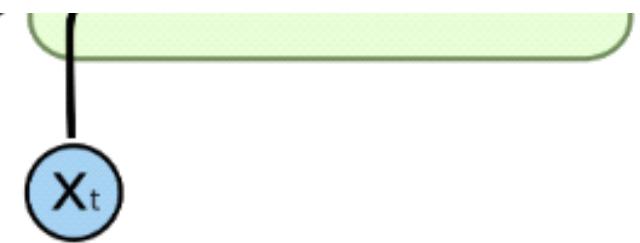
21 August 2023 17:27



LSTM Architecture

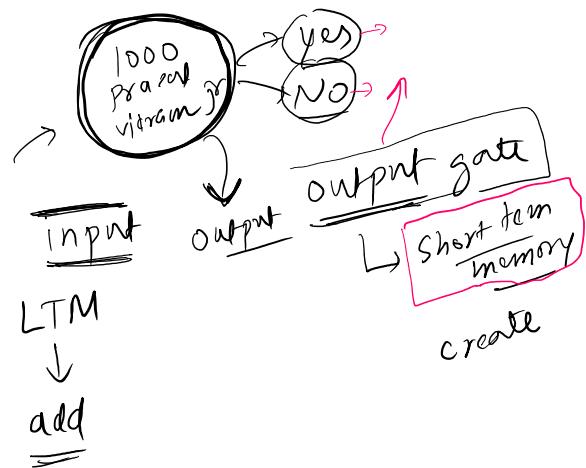
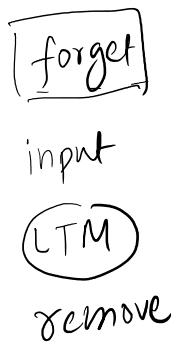
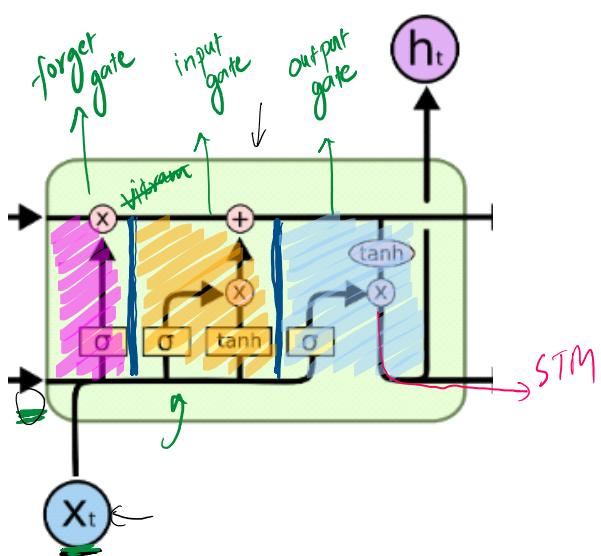
21 August 2023 18:41





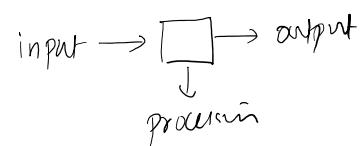
LSTM Gates

21 August 2023 19:06



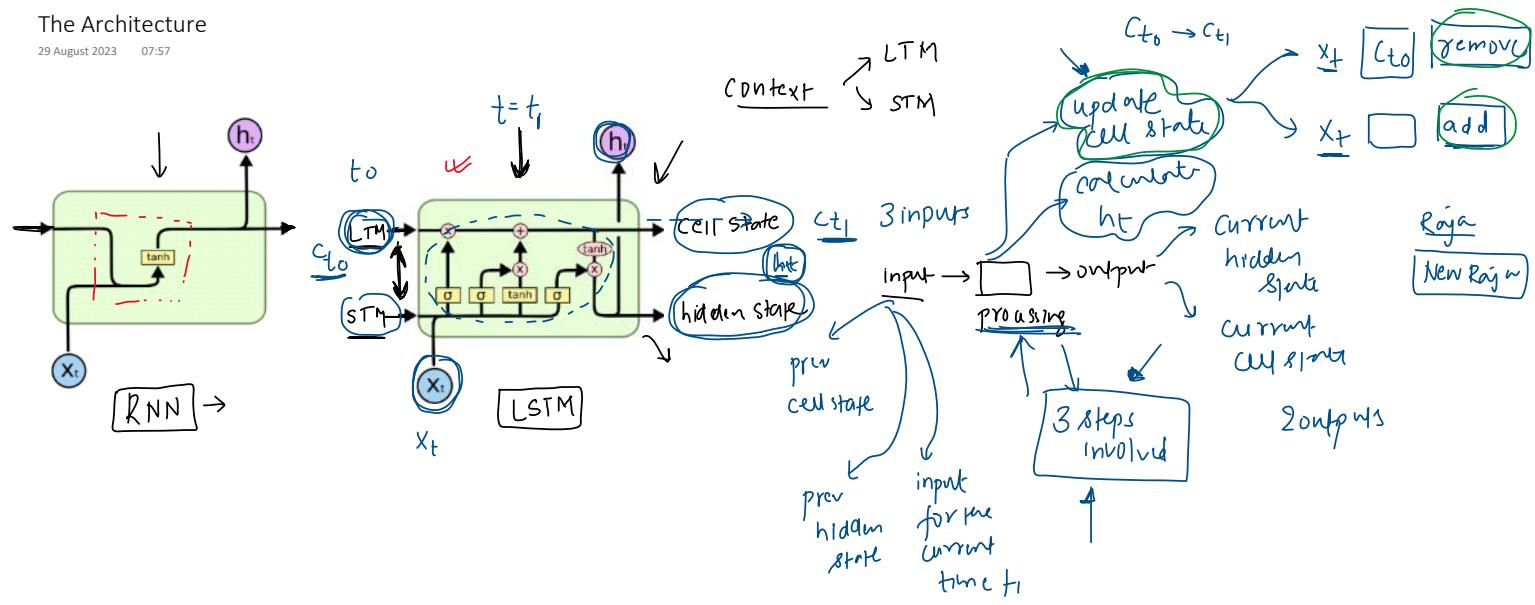
Summary

21 August 2023 19:29



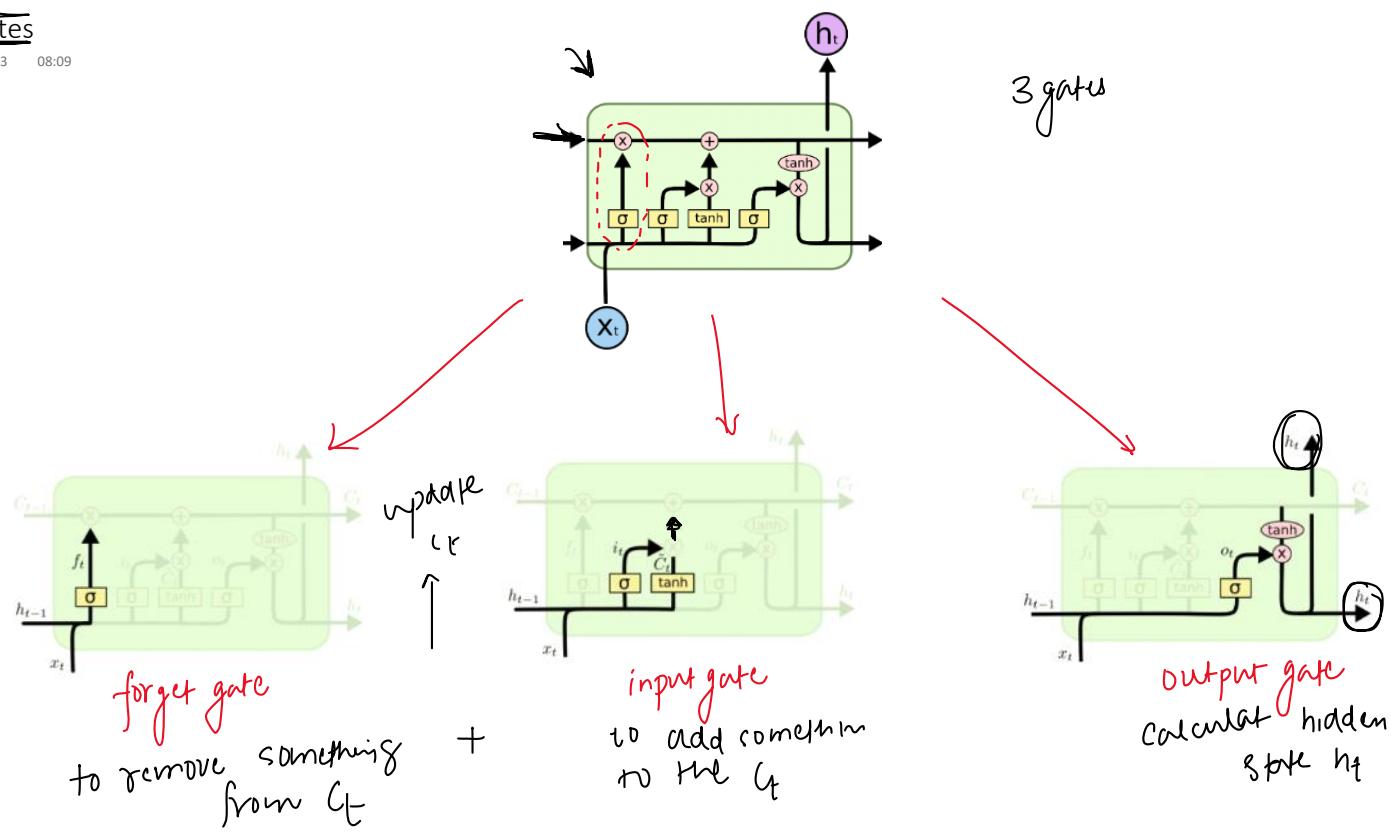
The Architecture

29 August 2023 07:57



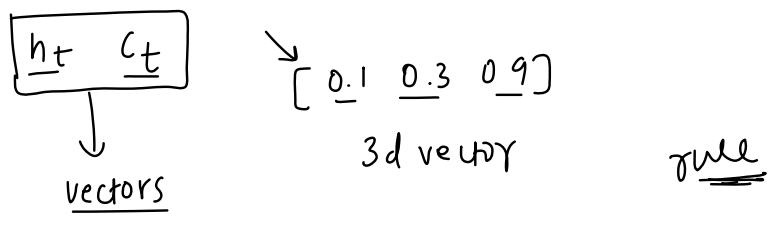
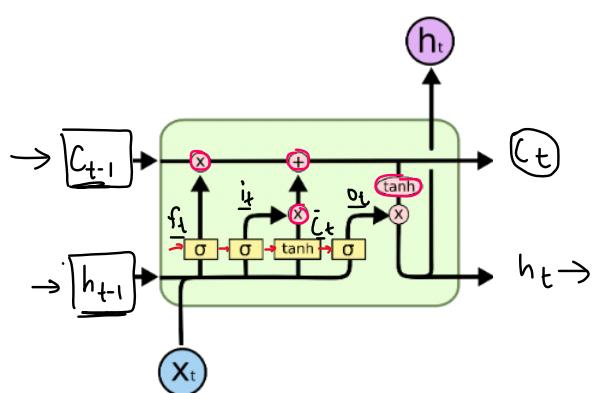
The Gates

29 August 2023 08:09



What are C_t and h_t

29 August 2023 08:08



$h_t \quad C_t$ dim equal

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

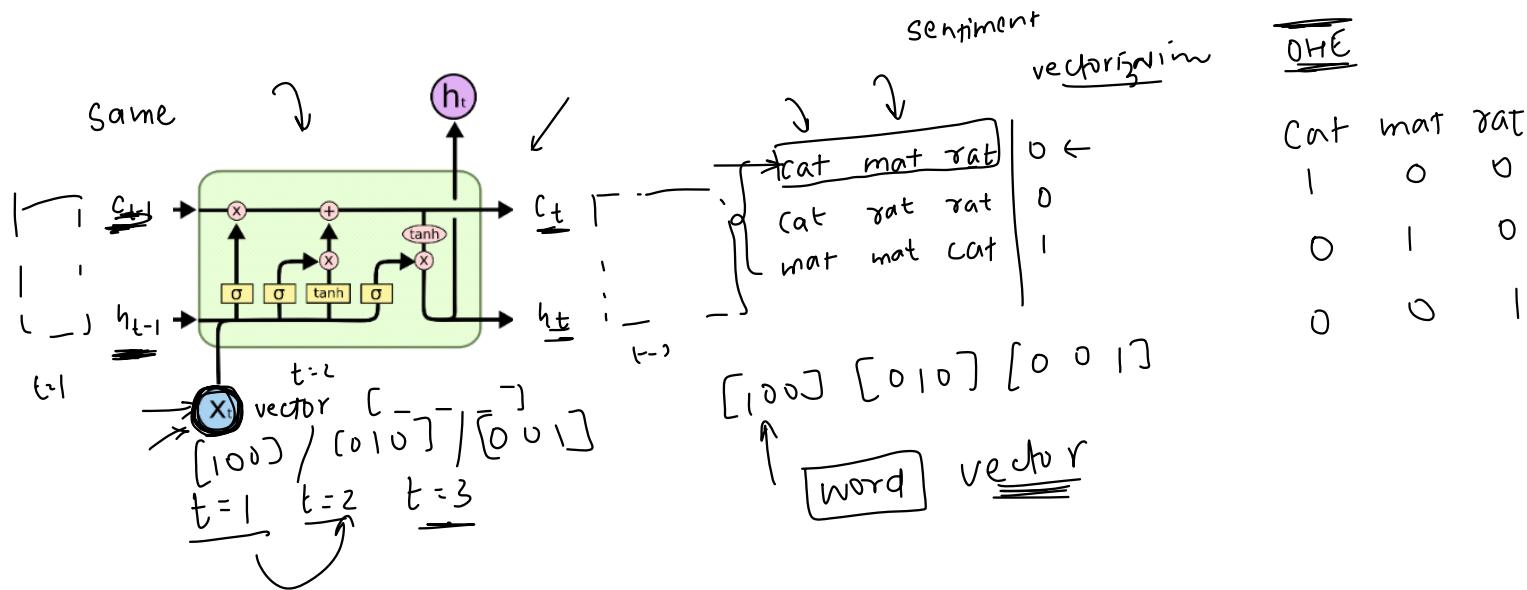
$C_t [0.55 \quad 0.6 \quad 0.0]$

same

What is X_t

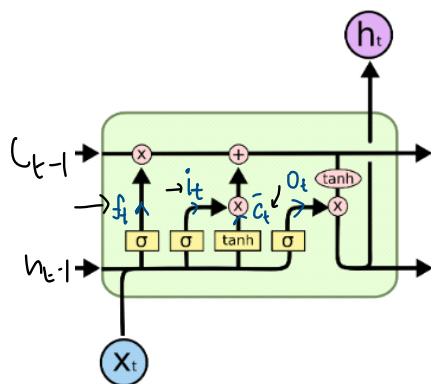
29 August 2023 17:40

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

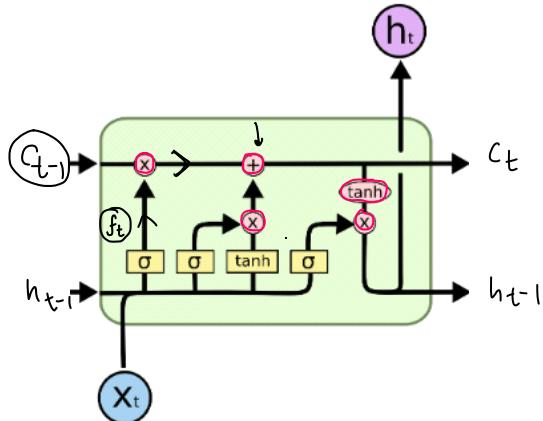
vectors

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

[x 4 2]
[] 7

Pointwise Operations

29 August 2023 18:26

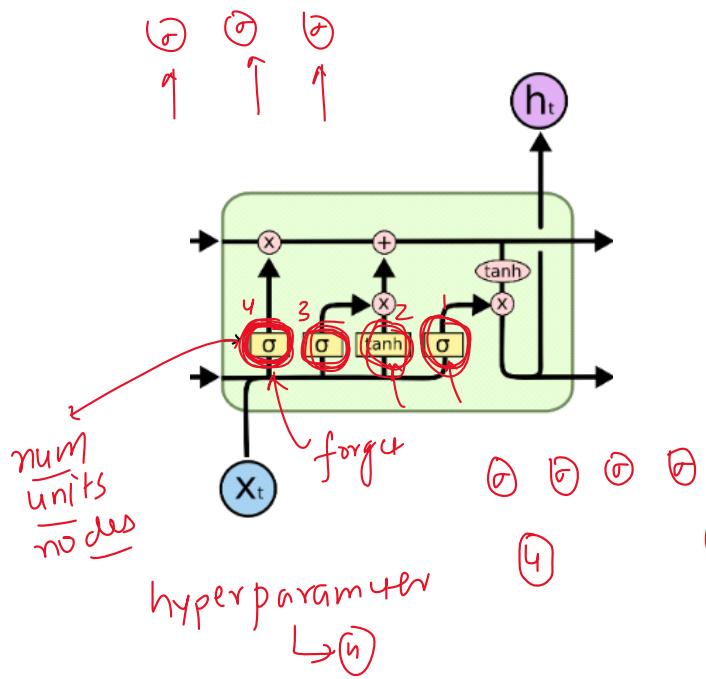


$$\begin{aligned}
 & \rightarrow \otimes \\
 & \rightarrow + \\
 & \rightarrow \tanh
 \end{aligned}$$

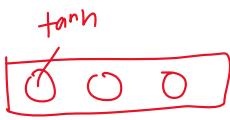
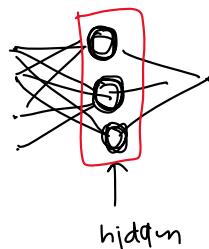
$c_{t-1} = \begin{bmatrix} 4 & 5 & 6 \\ 1 & 2 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 0.26 & 0.34 & 0.53 \end{bmatrix}$
 $\tanh(u)$
 $f_t = \underline{\text{shape(dim)}} \quad \downarrow \text{vector}$
 $c_{t-1} \otimes f_t \rightarrow \text{vector} \rightarrow [5 \ 7 \ 9]$
 $\rightarrow [n \ 10 \ 18]$

→ Neural Network Layers

29 August 2023 18:34

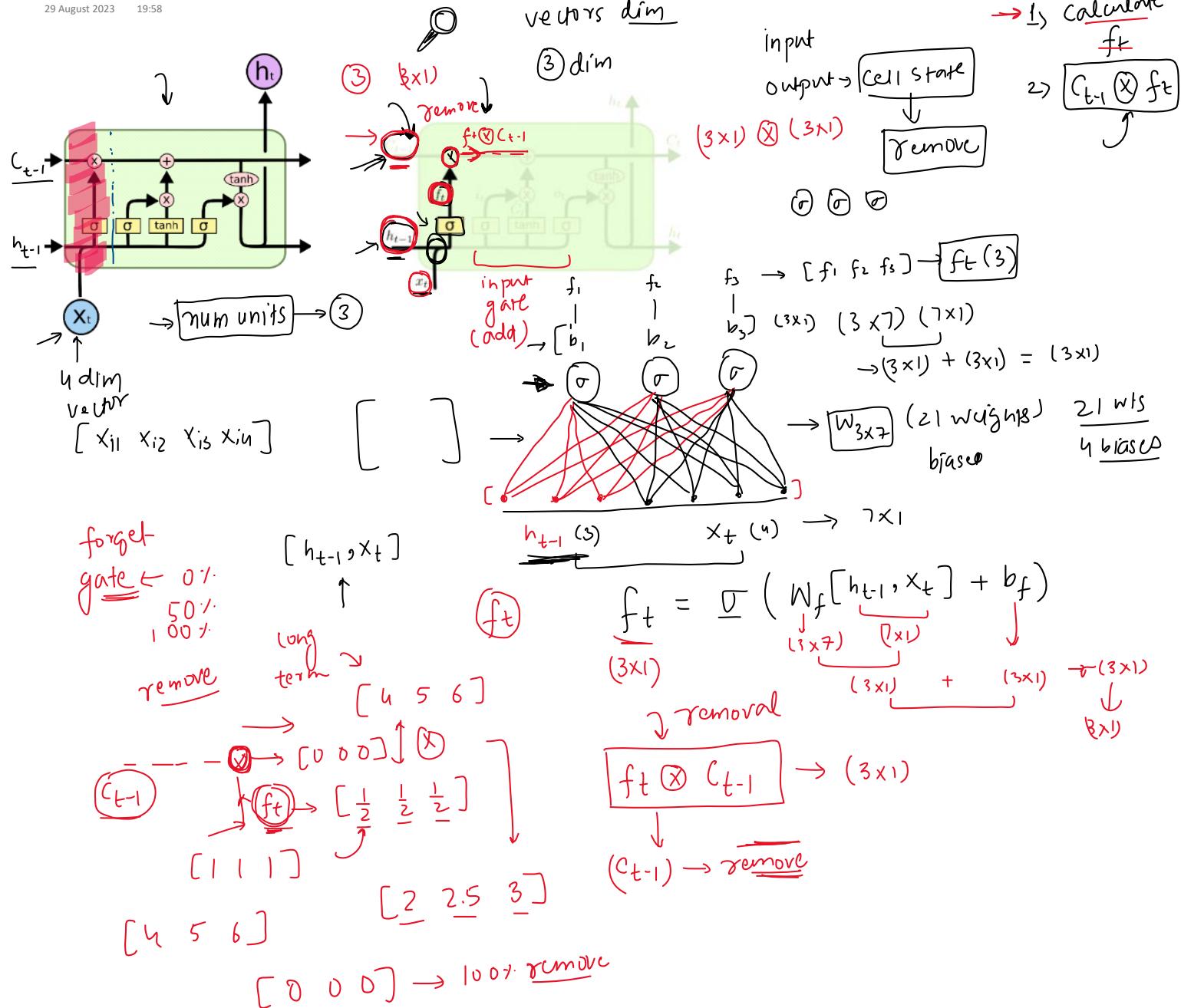


ANN



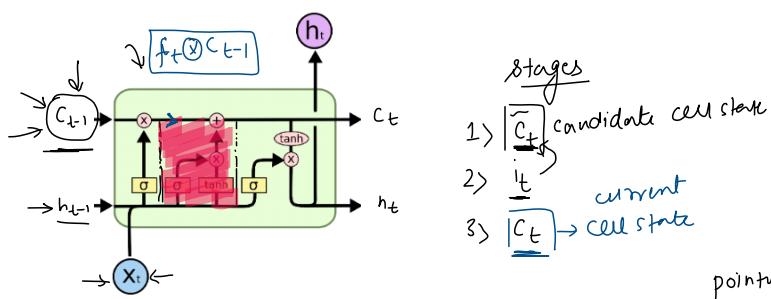
The Forget Gate

29 August 2023 19:58

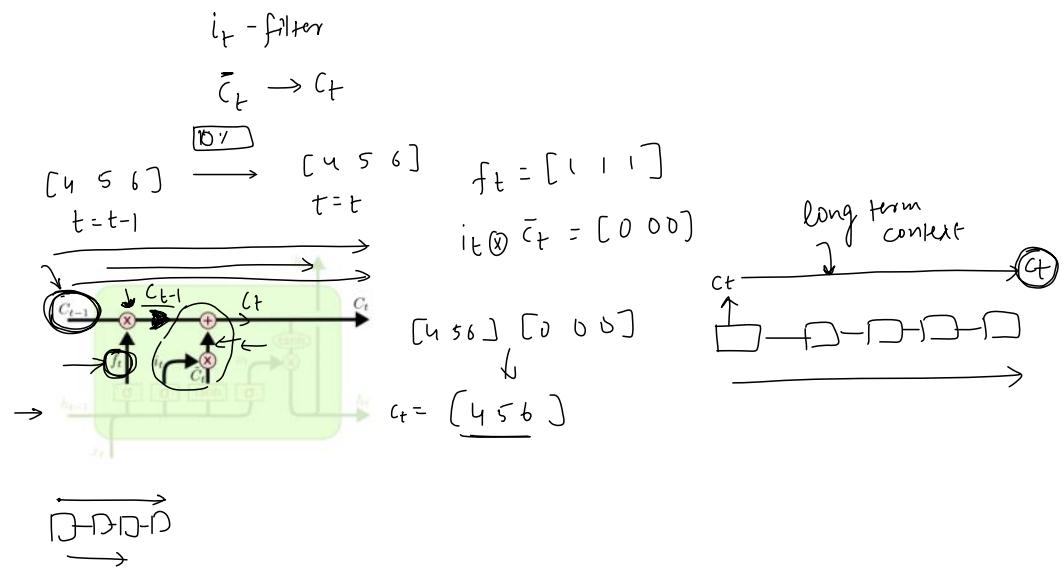
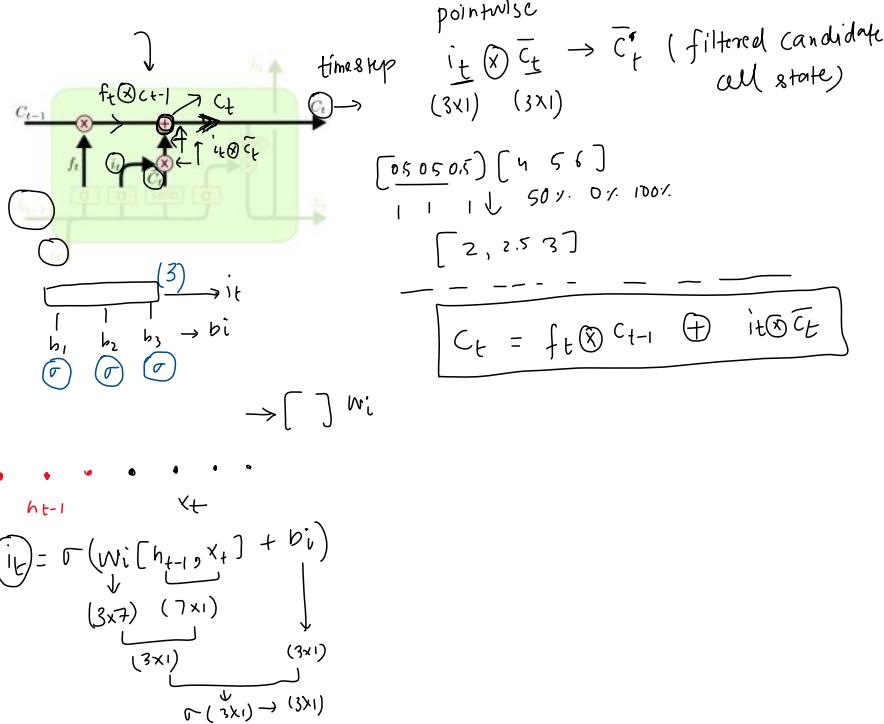
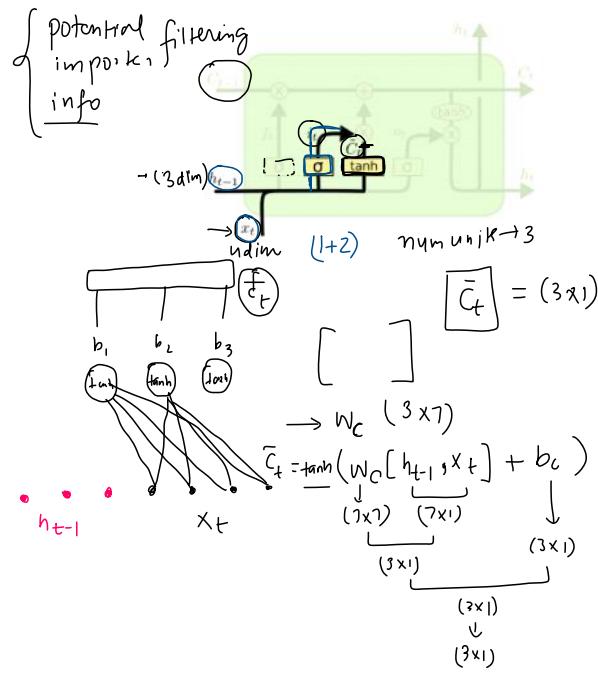


The Input Gate
30 August 2023 04:38

add some new imp info to c_t

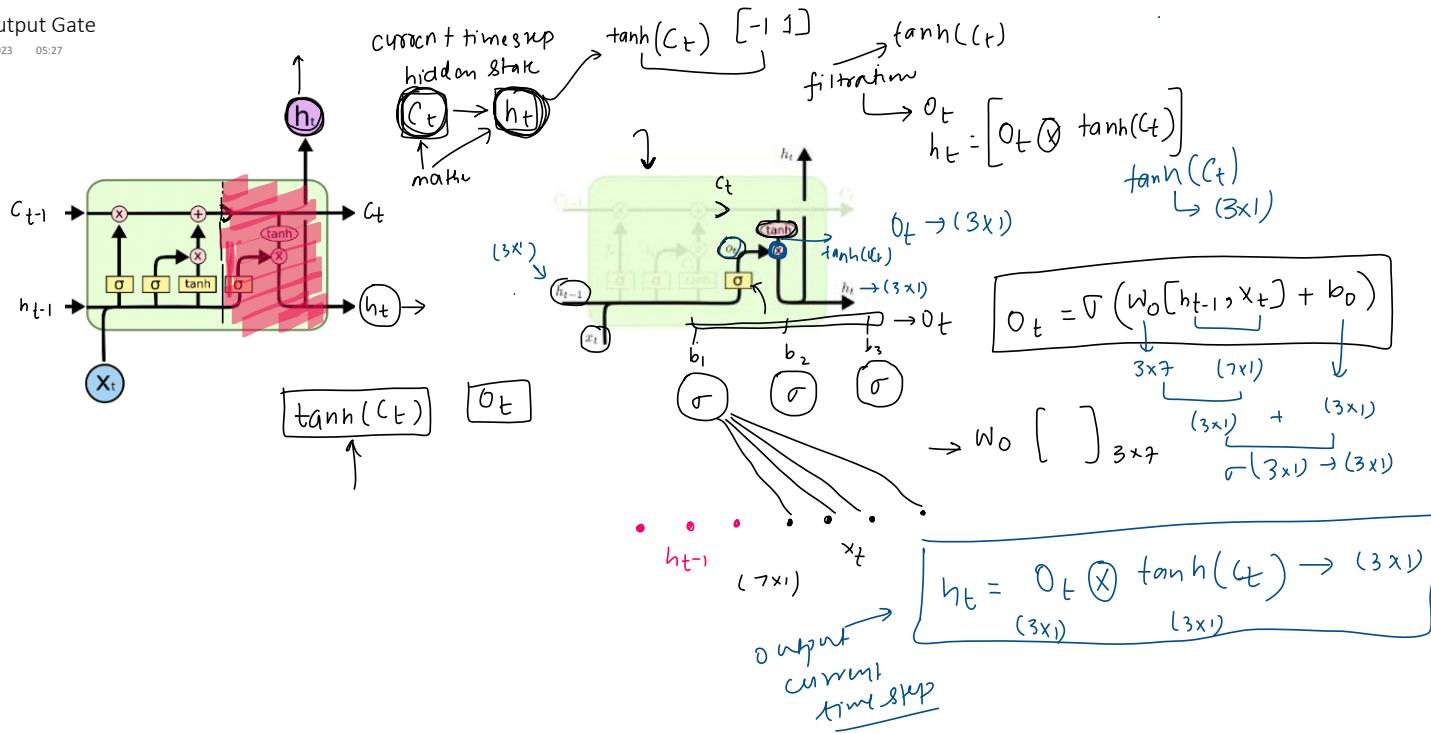


- stages
- 1) $\underline{c_t}$ candidate cell state
 - 2) $\underline{i_t}$ current
 - 3) $\underline{c_t}$ cell state



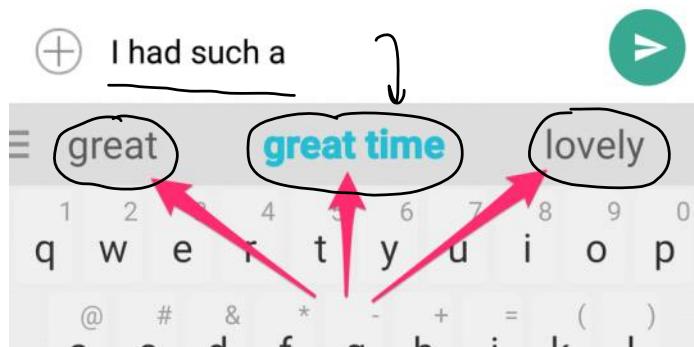
The Output Gate

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

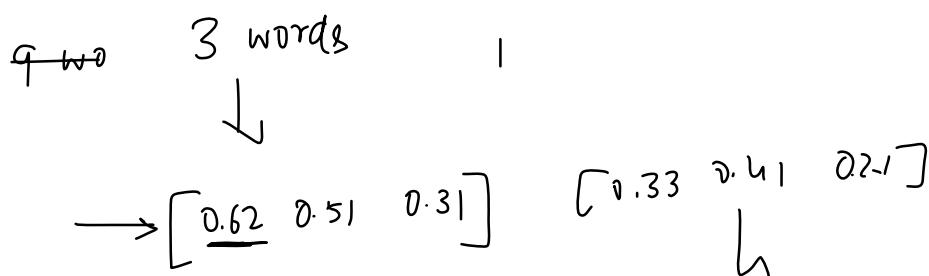
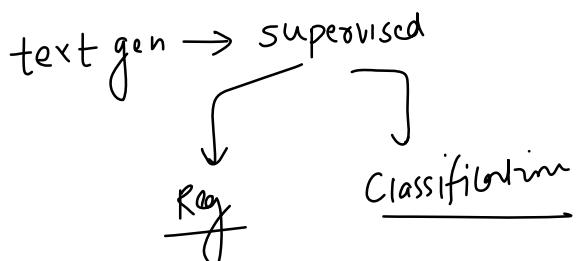
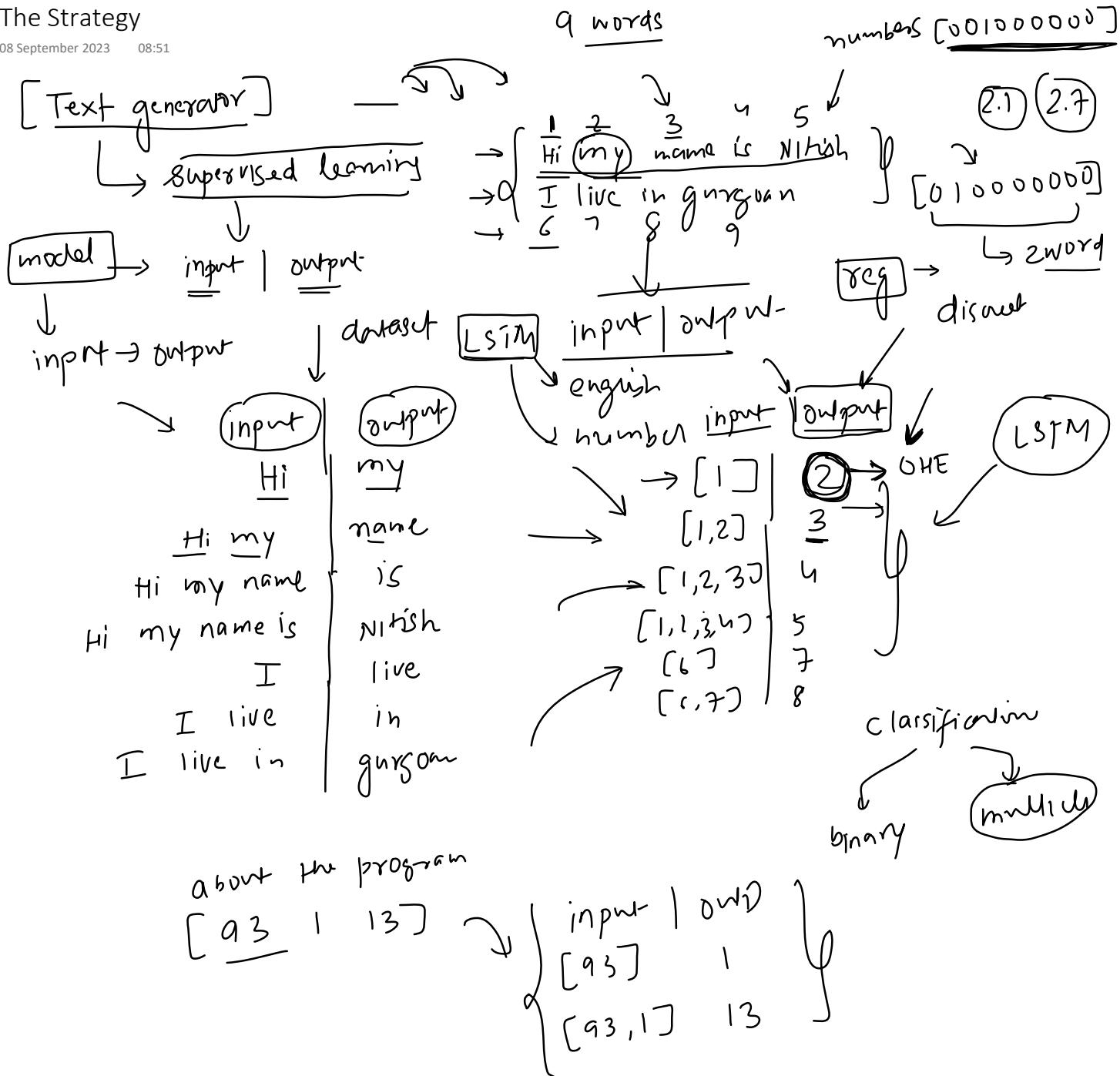
Thanks for reaching out and glad to be in your network.

Image, Video, GIF, Smileys

Send

The Strategy

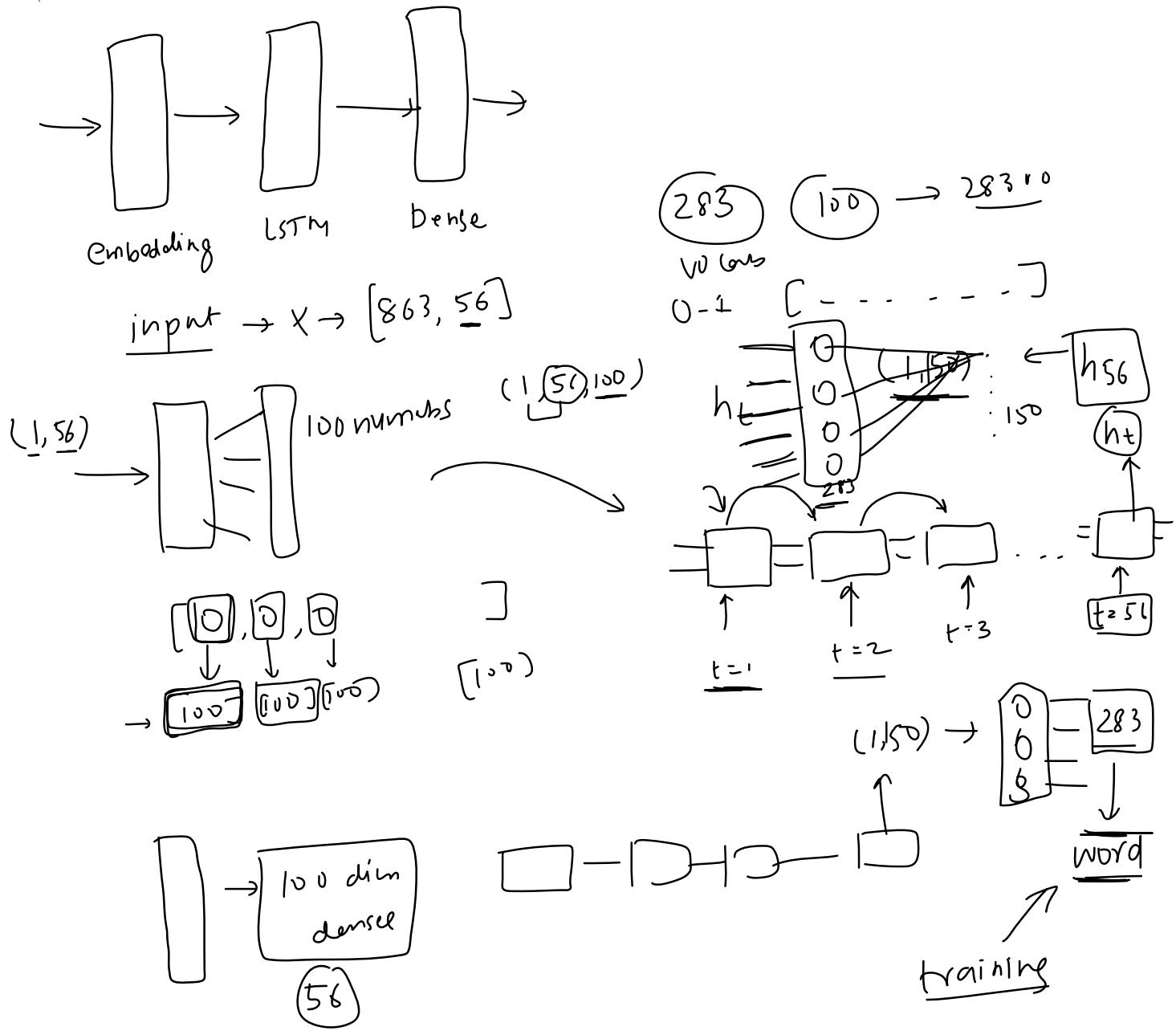
08 September 2023 08:51



[| 0 0] [0 | 0]
↳ first word ↳ second word

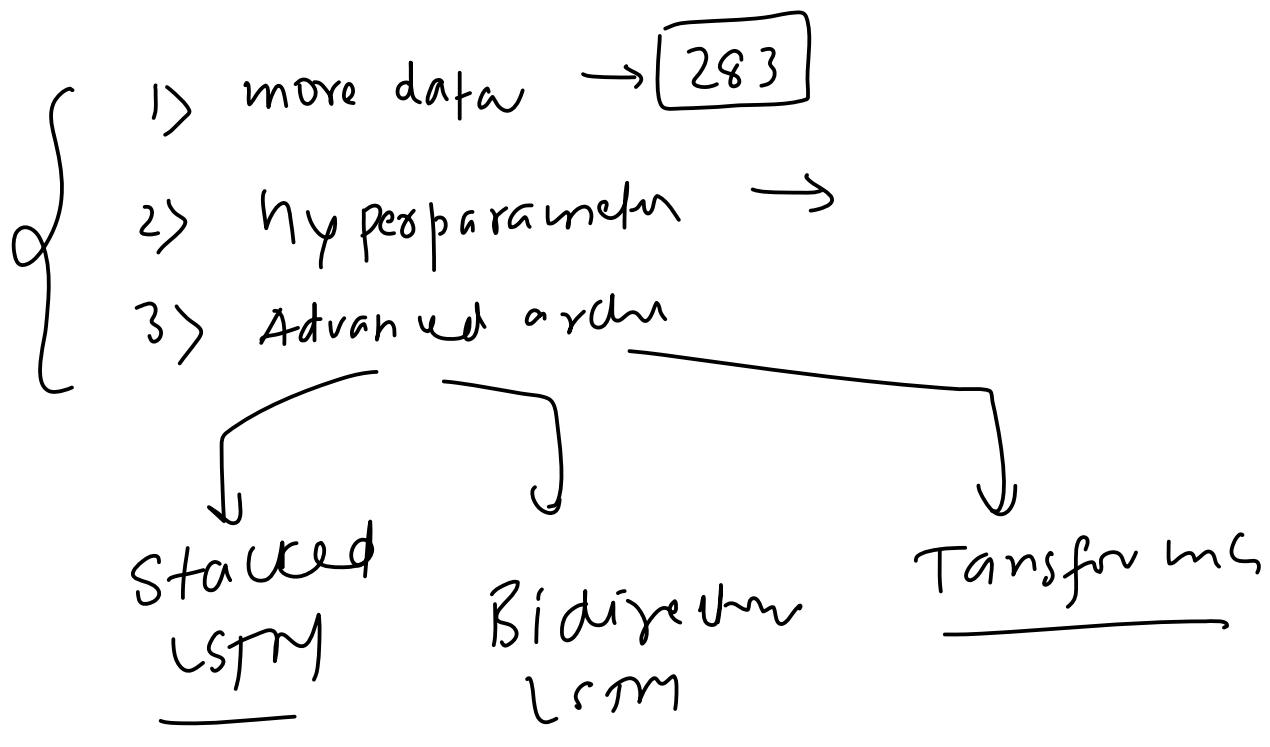
The Architecture

08 September 2023 08:55



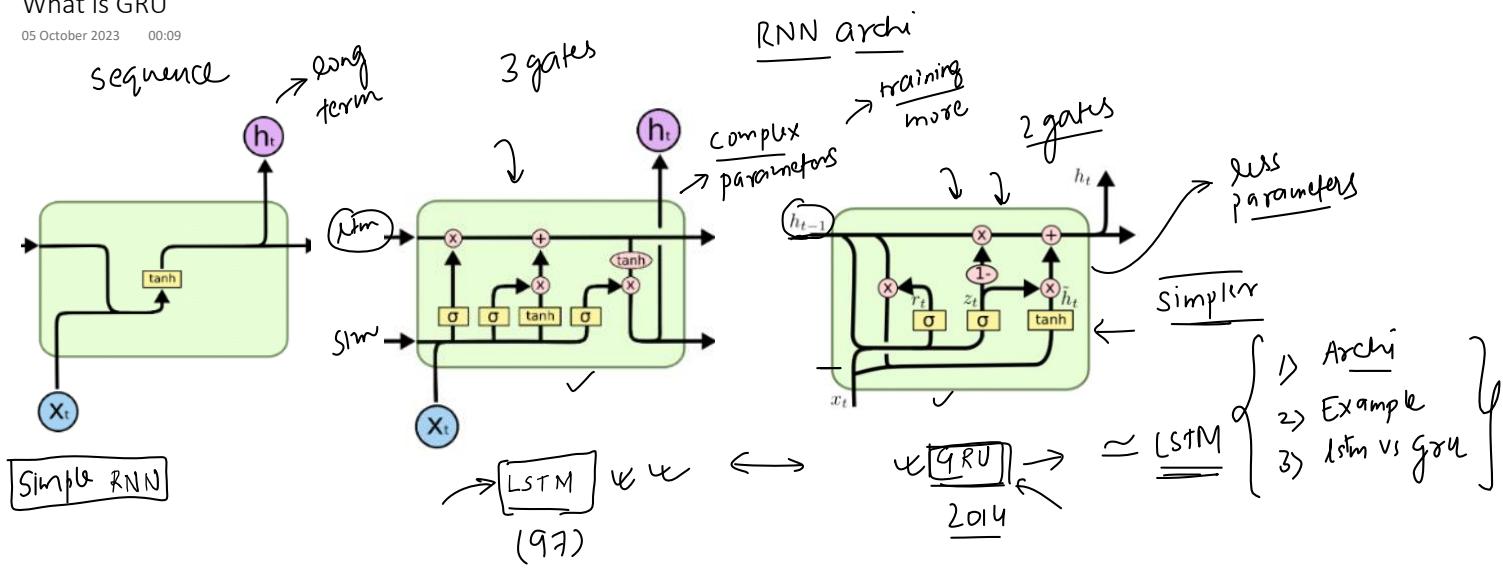
How to improve performance?

08 September 2023 08:51



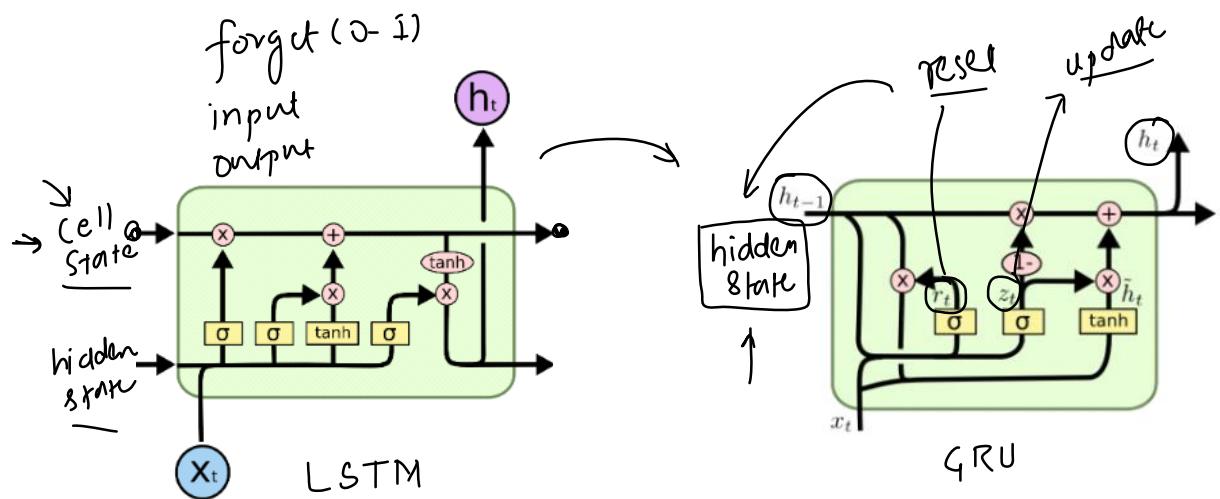
What is GRU

05 October 2023 00:09



The Big Idea Behind GRU

05 October 2023 00:47

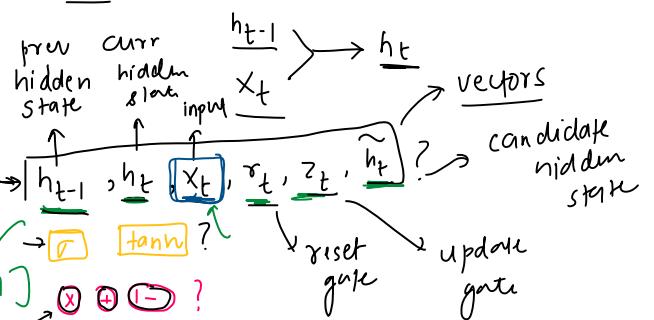


The Setup

05 October 2023 01:07

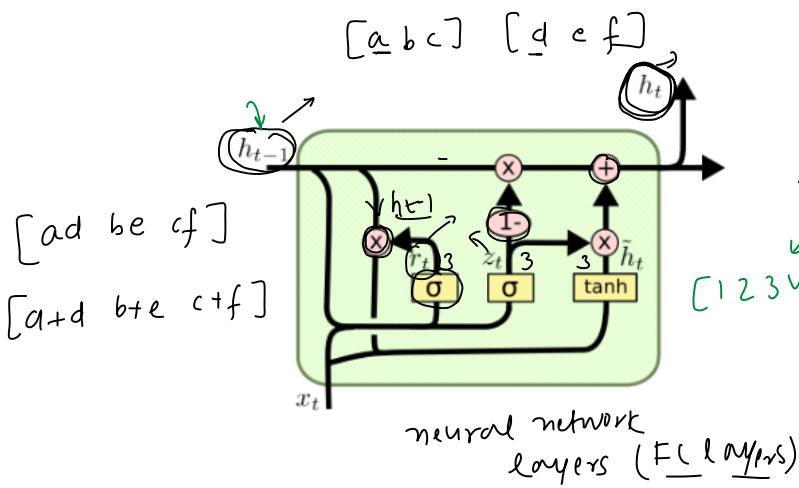
→ Advise → LSTM / GRU → confusing

goal → \boxed{t}



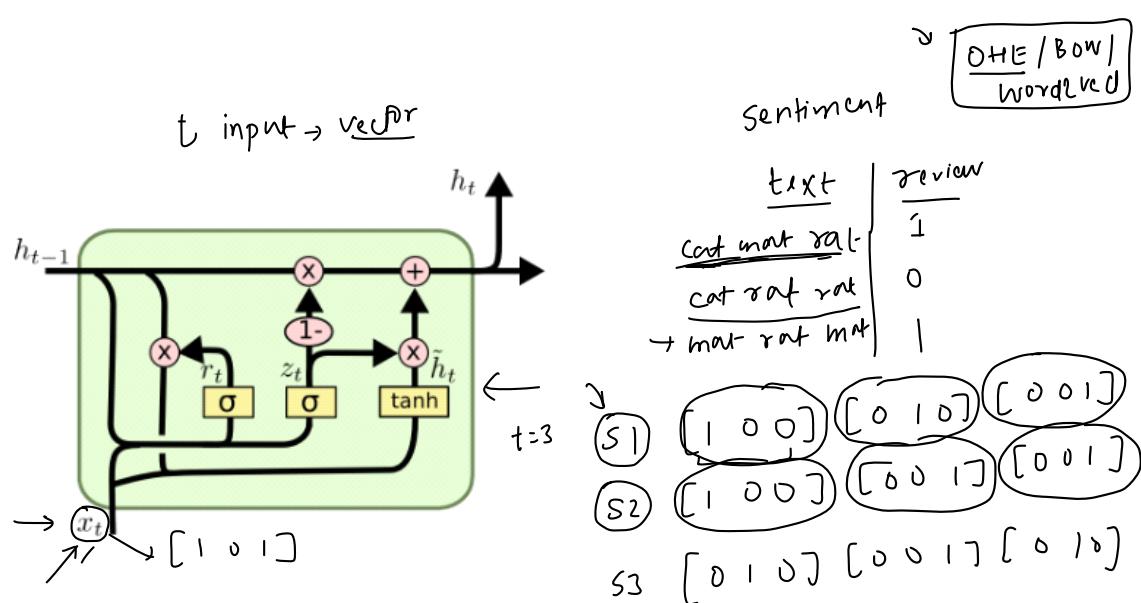
$\boxed{1234}$ hidden same

number of
10 = $\boxed{5} \quad \boxed{6}$



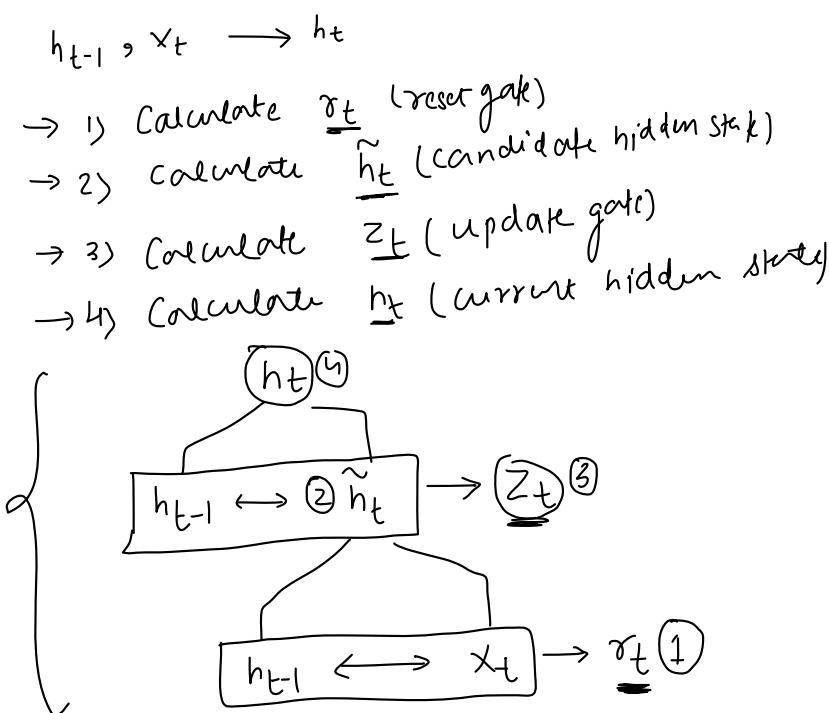
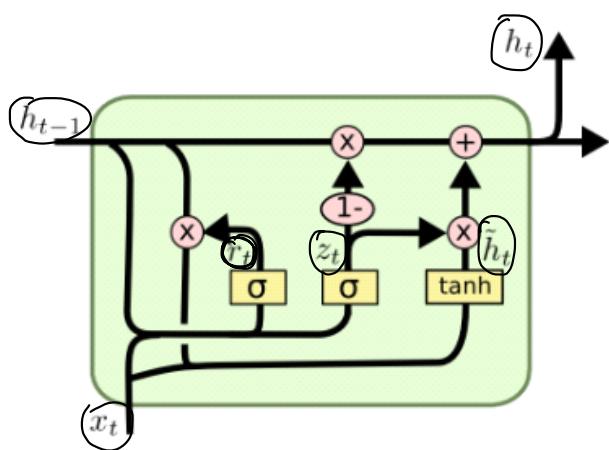
The Input Xt

05 October 2023 01:52



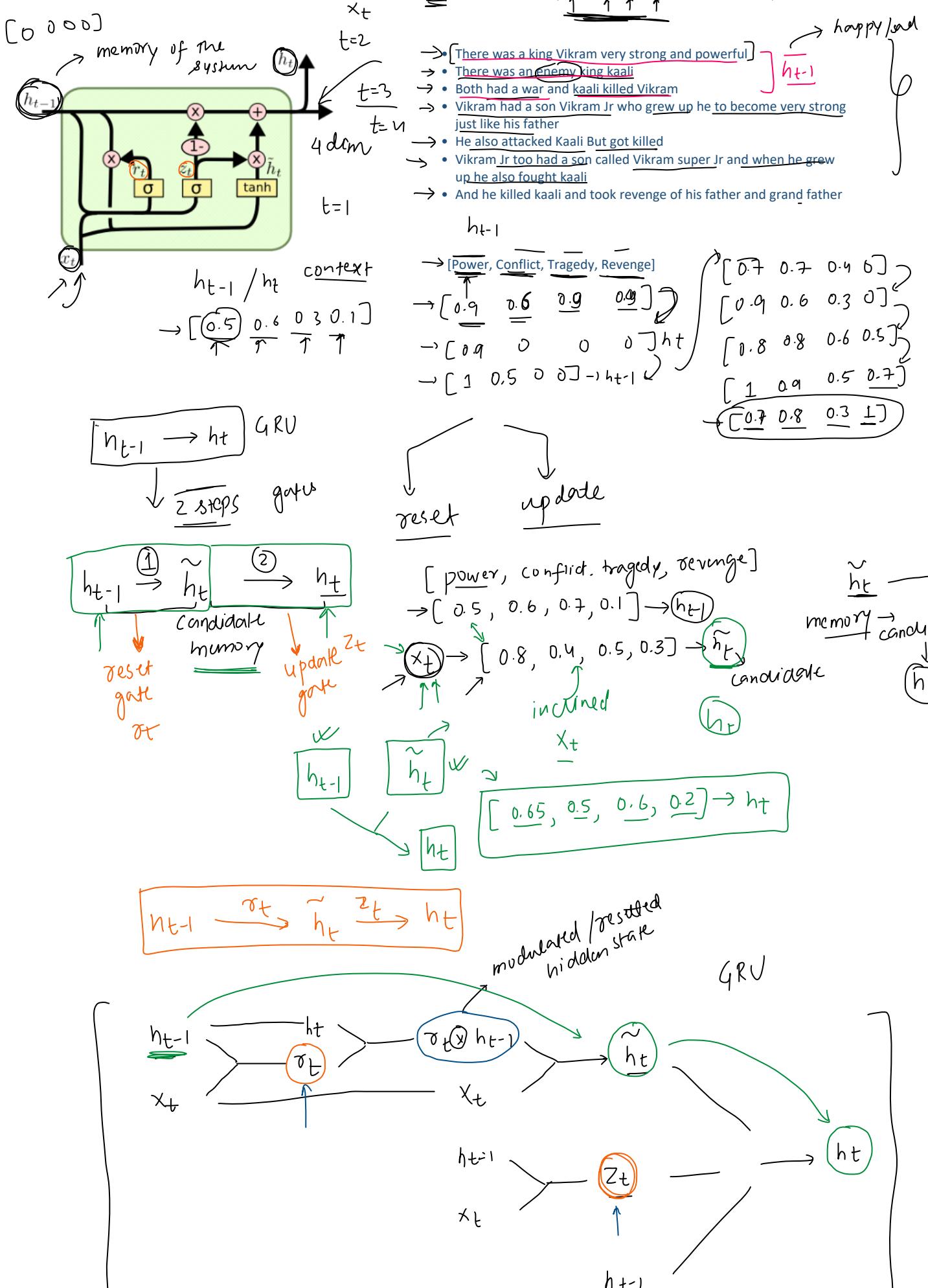
Architecture

05 October 2023 02:10



What exactly is hidden state?

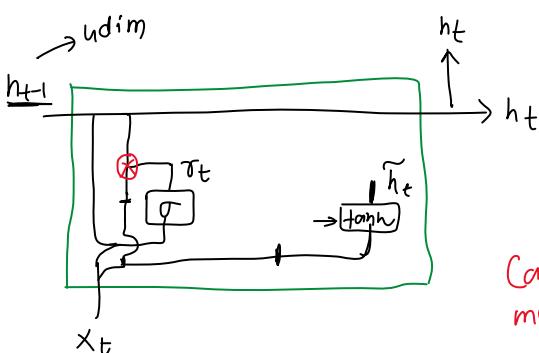
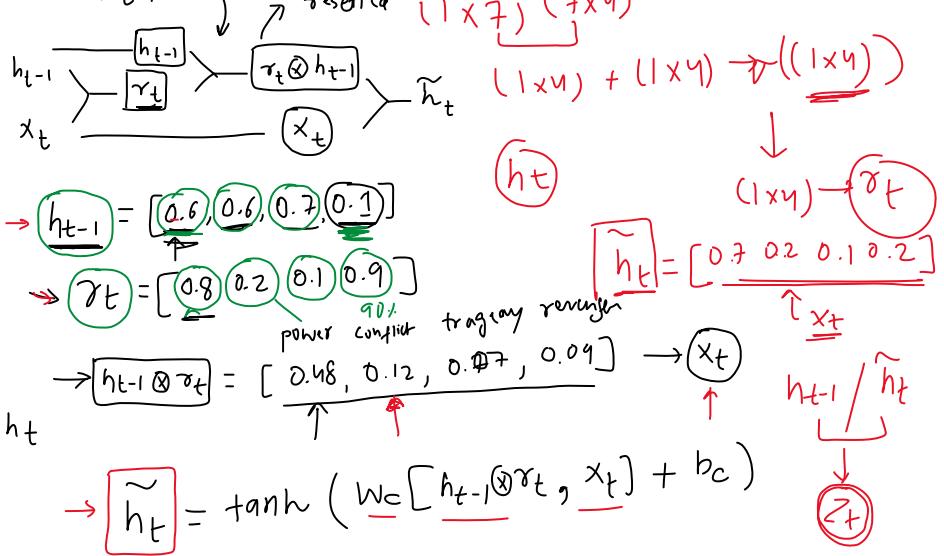
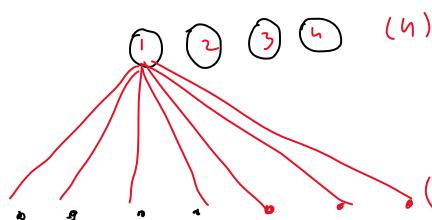
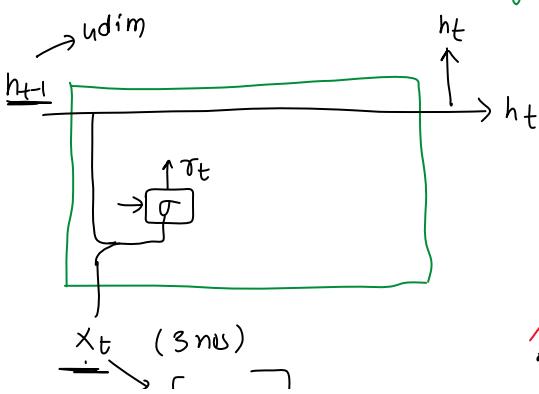
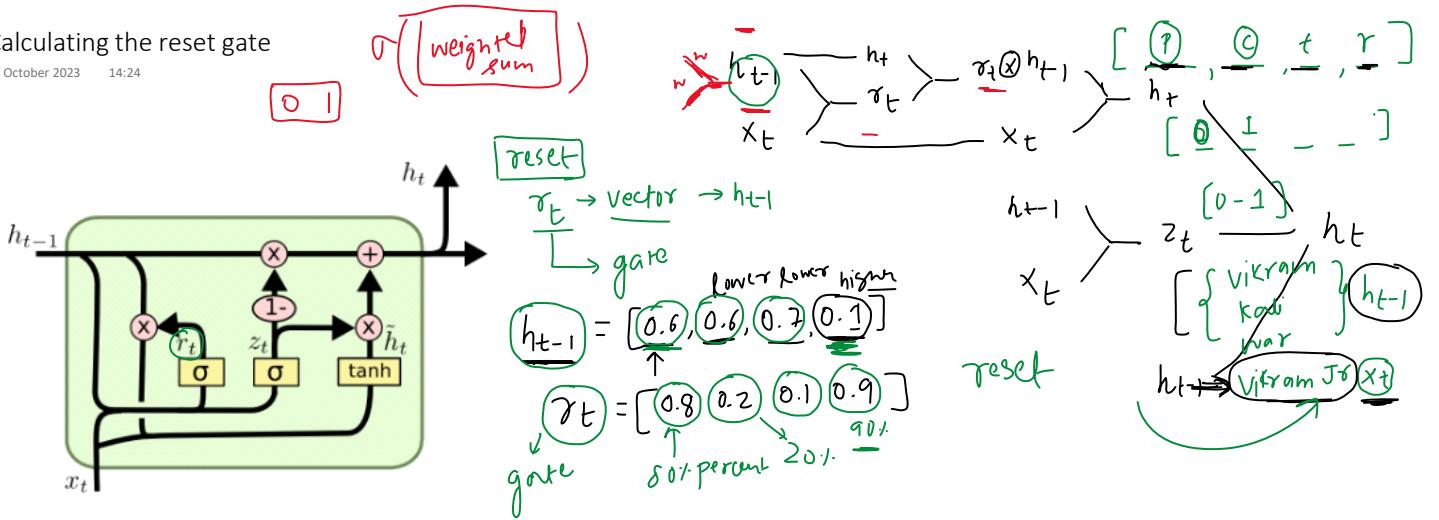
05 October 2023 02:19





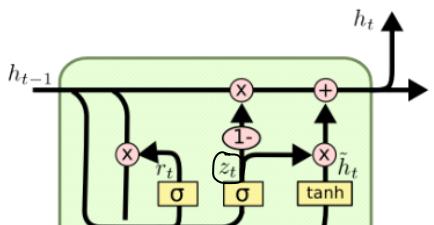
Calculating the reset gate

05 October 2023 14:24



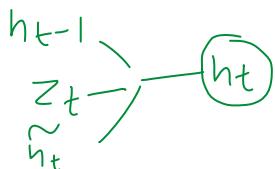
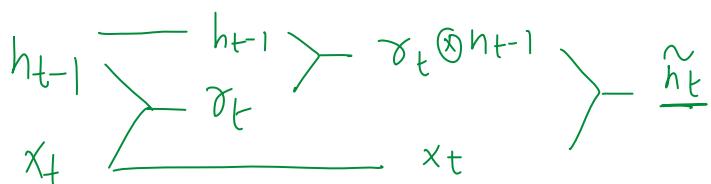
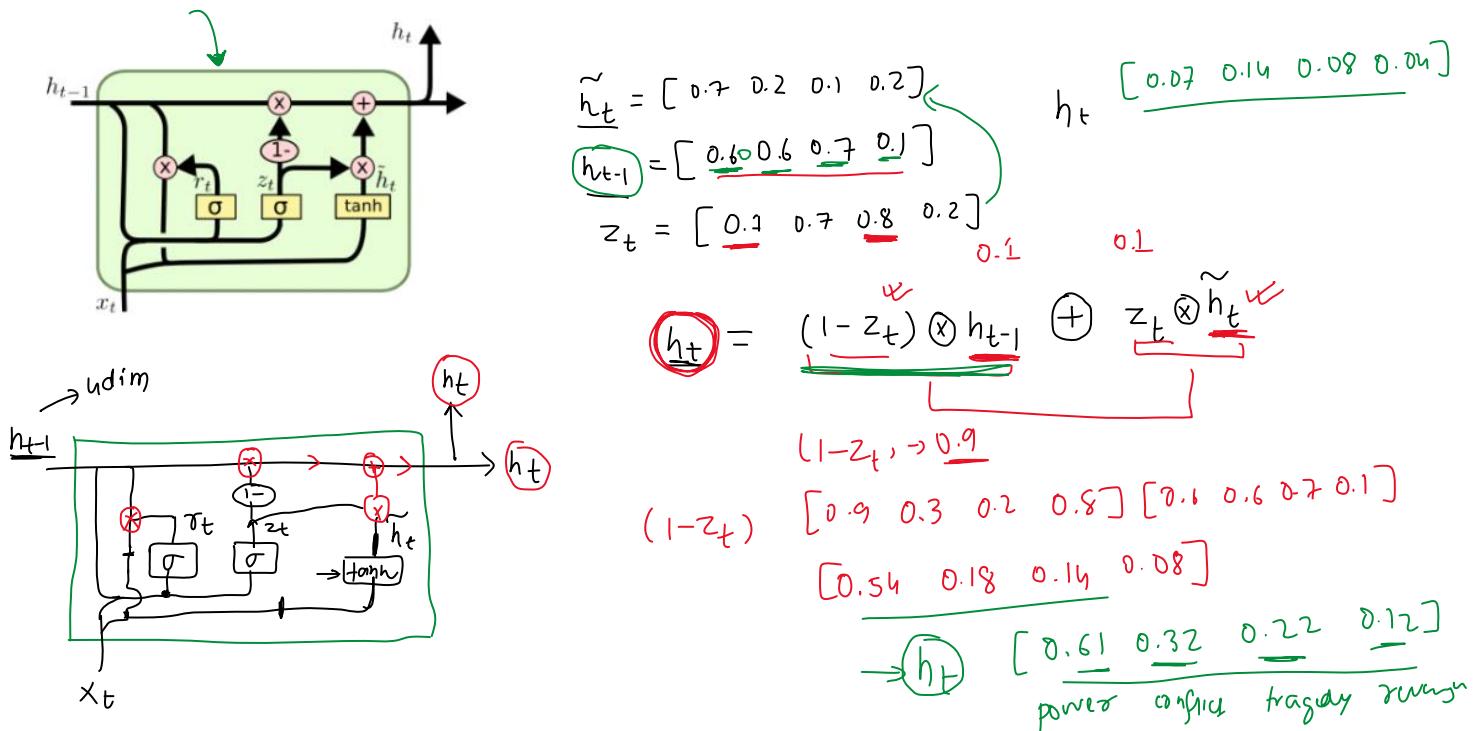
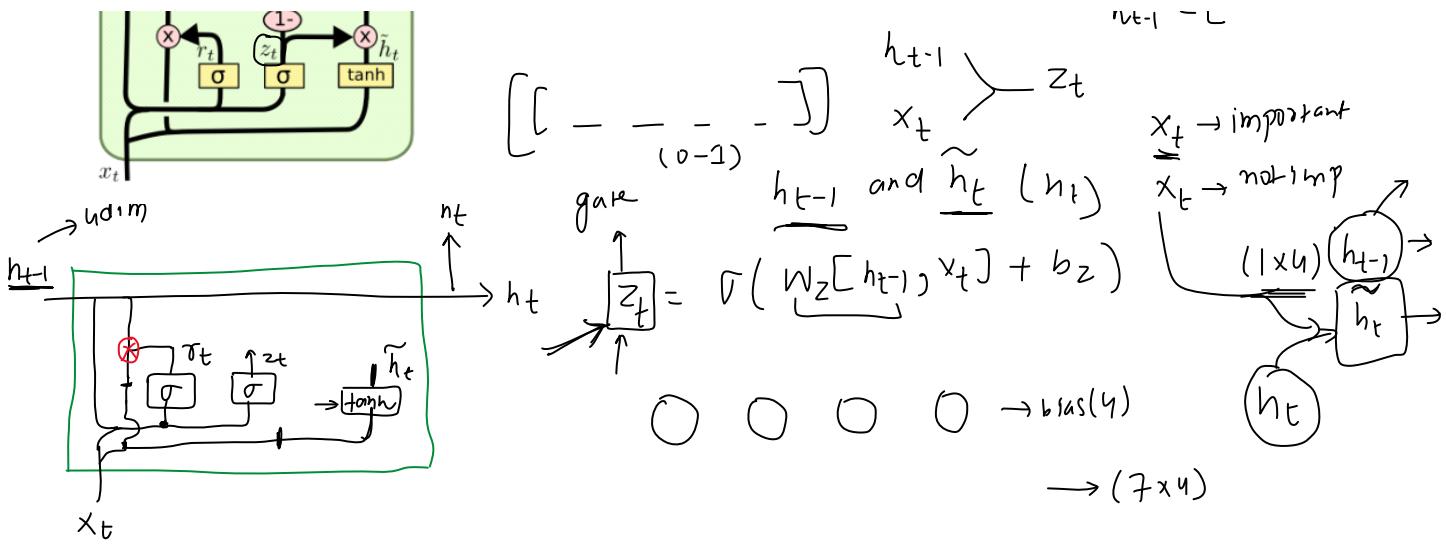
Candidate memory $\rightarrow b_c$

$$\rightarrow W_c \quad (7 \times 4)$$



$$h_{t-1} \rightarrow r_t \rightarrow h_t$$

$$h_{t-1} \rightarrow z_t$$



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:

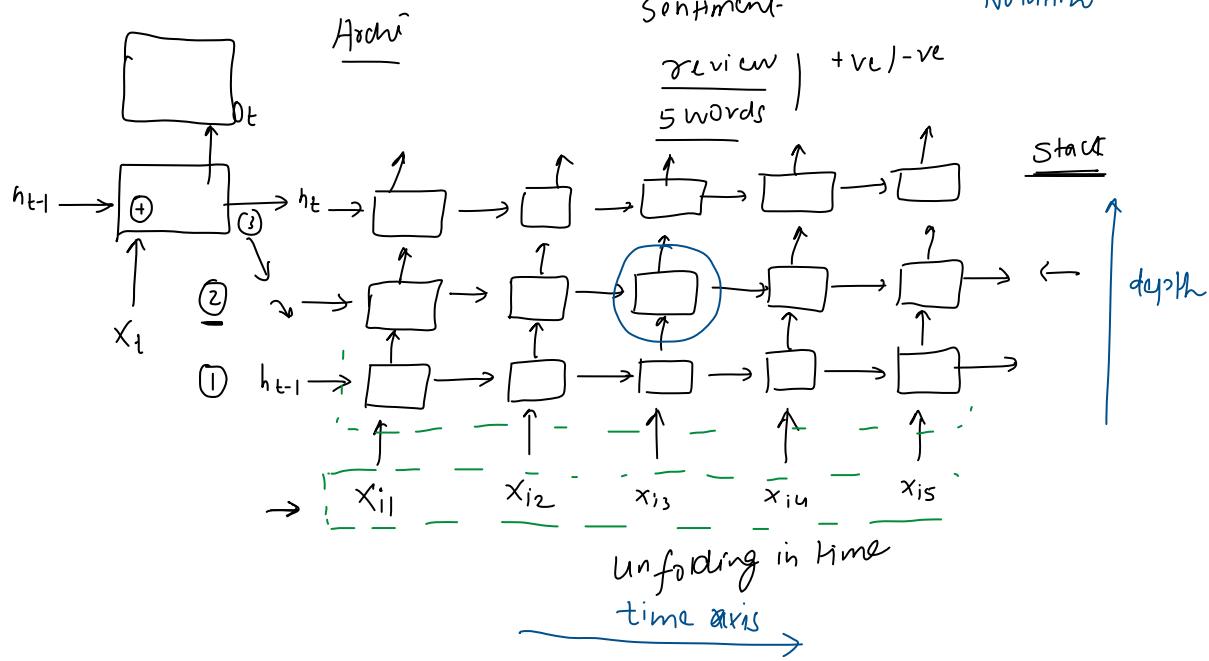
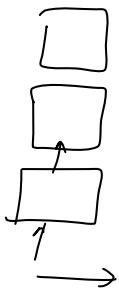
- The choice between LSTM and GRU often comes down to empirical testing. Depending on the dataset and task, one might outperform the other. However, GRUs, due to their simplicity, are often the first choice when starting out.

What is Deep RNN →

17 October 2023

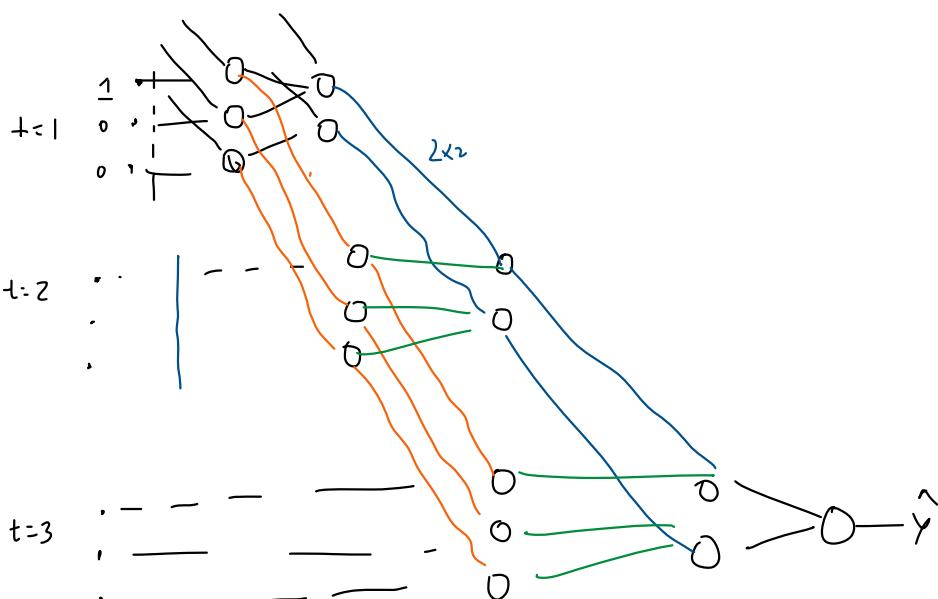
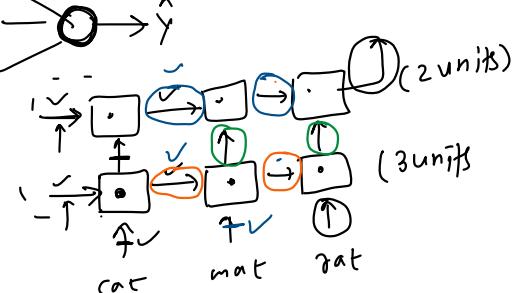
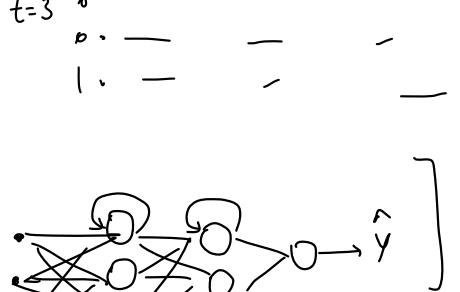
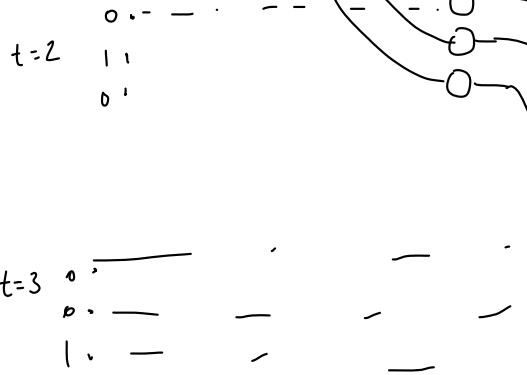
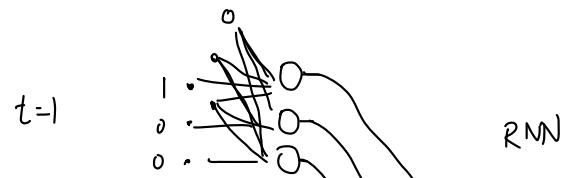
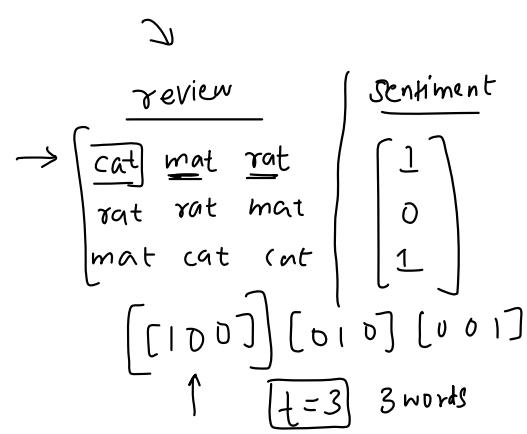
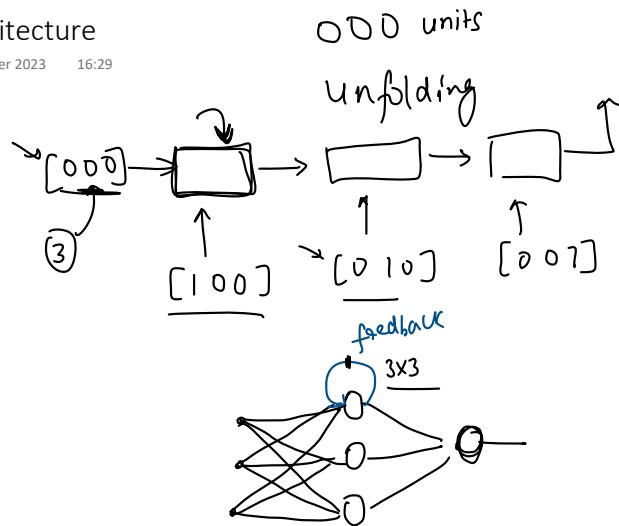
ANN

J

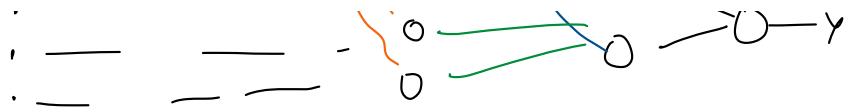


Architecture

17 October 2023 16:29



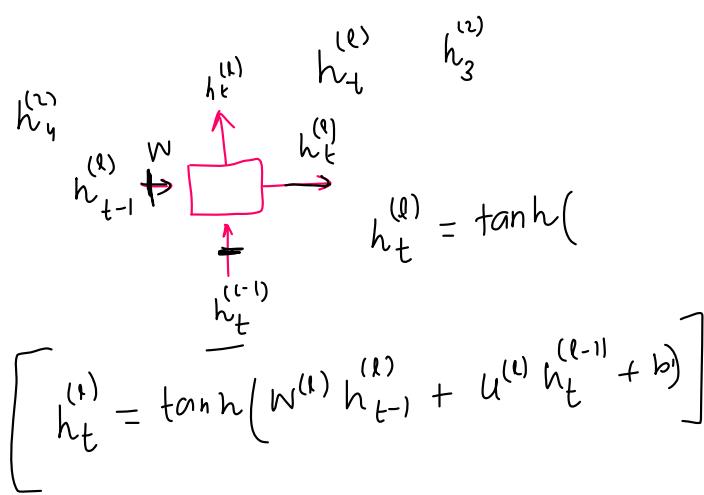
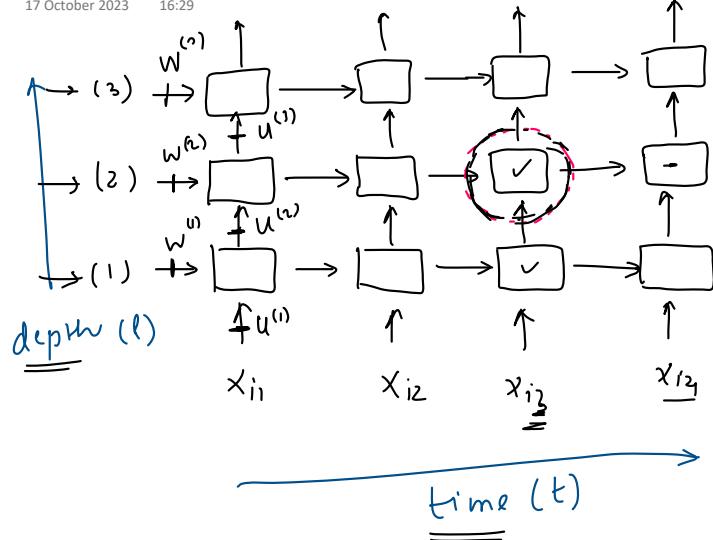
$t=3$



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

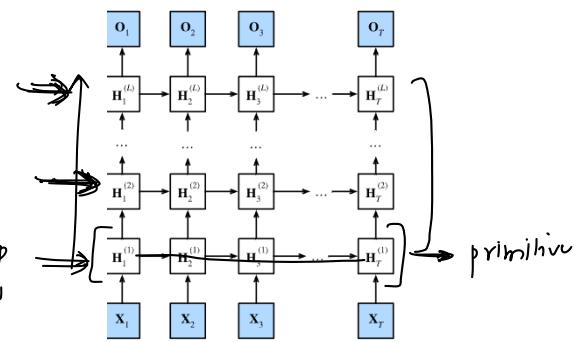
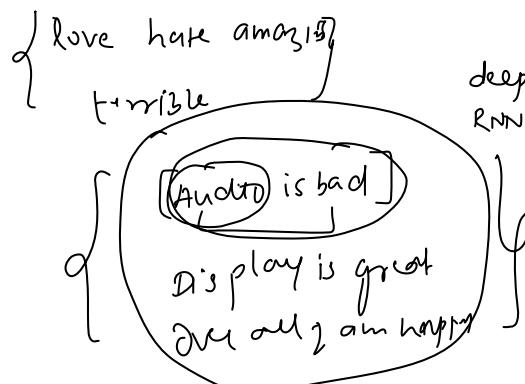
product

stack

sentence

encoder-decoder
↓
machine

{
 deep
 KNNs
 } ↗



→ sentence



When to use Deep RNNs?

Complex
tasks

{ speech recg
Machine translation }

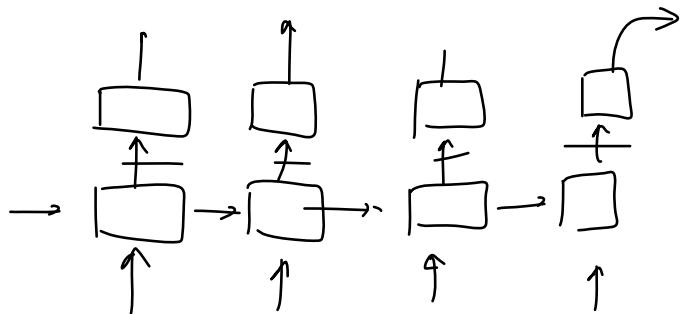
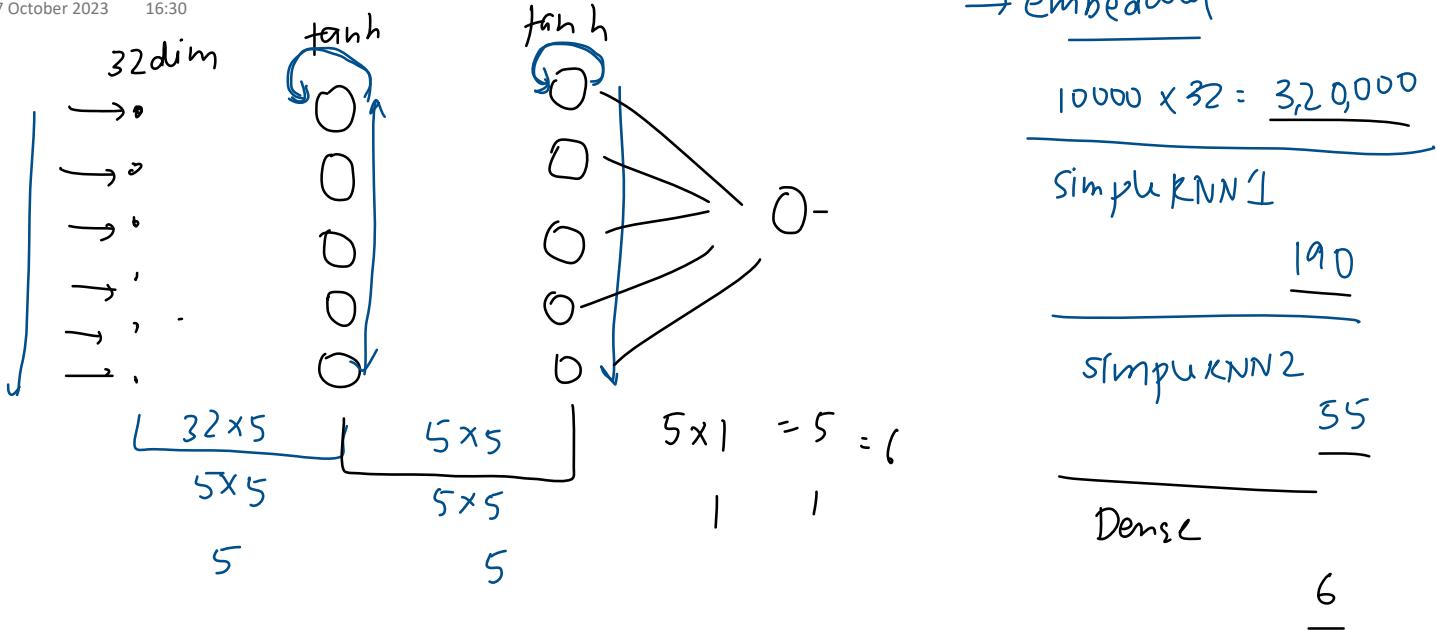
Large datasets
Overfitting

Computational

Simpler Models
↓
Deep RNN

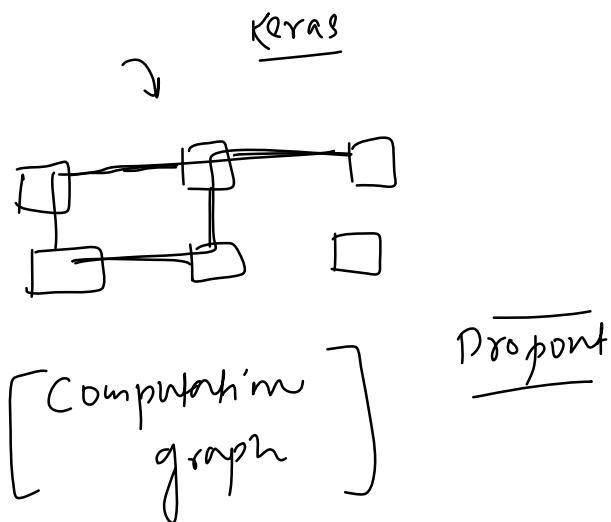
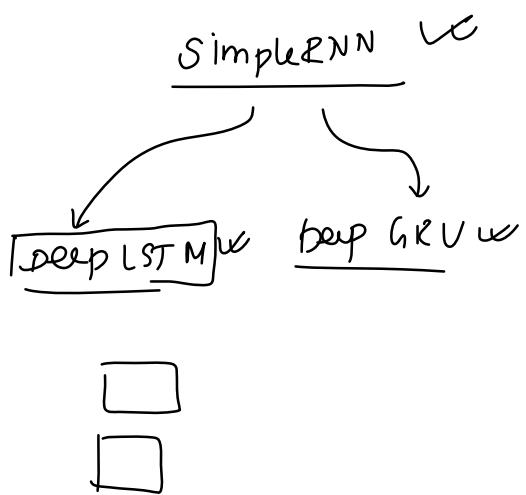
Code Example

17 October 2023 16:30



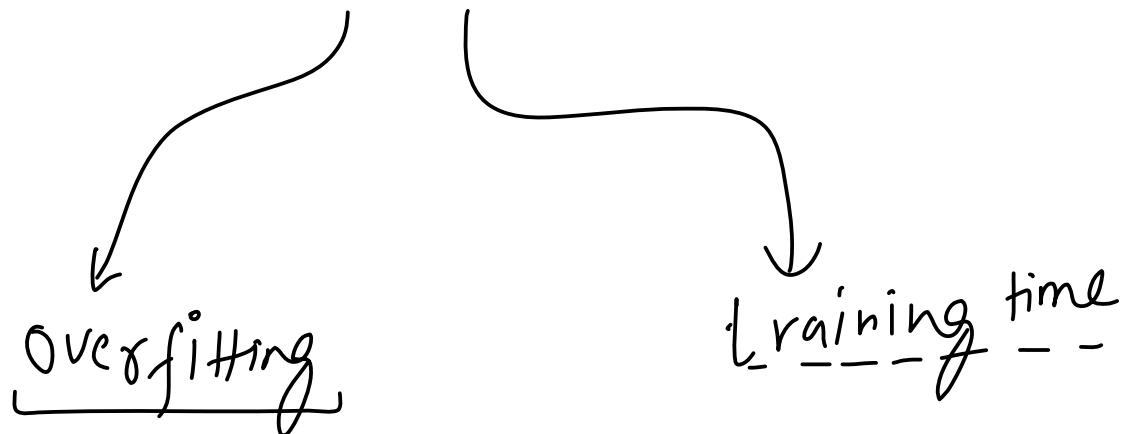
Variants

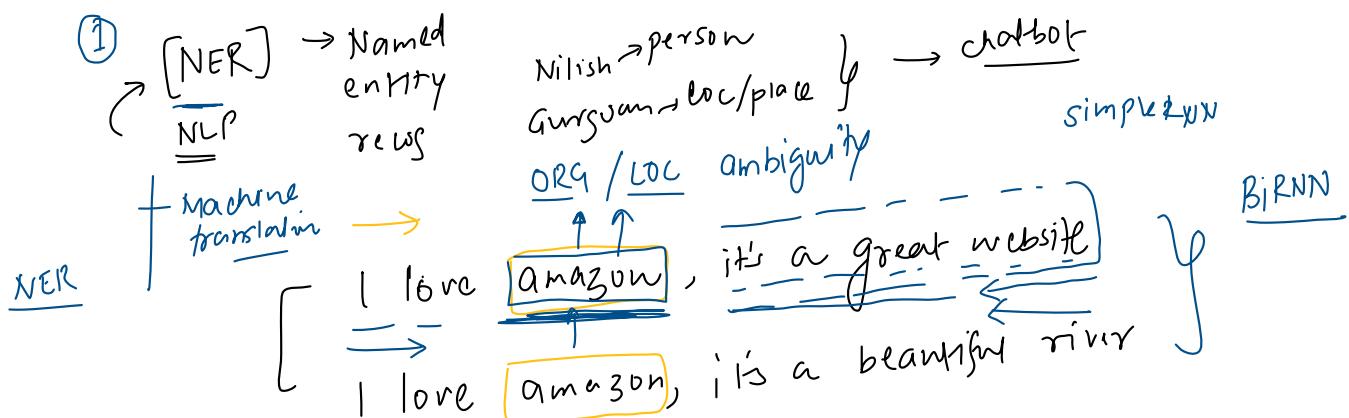
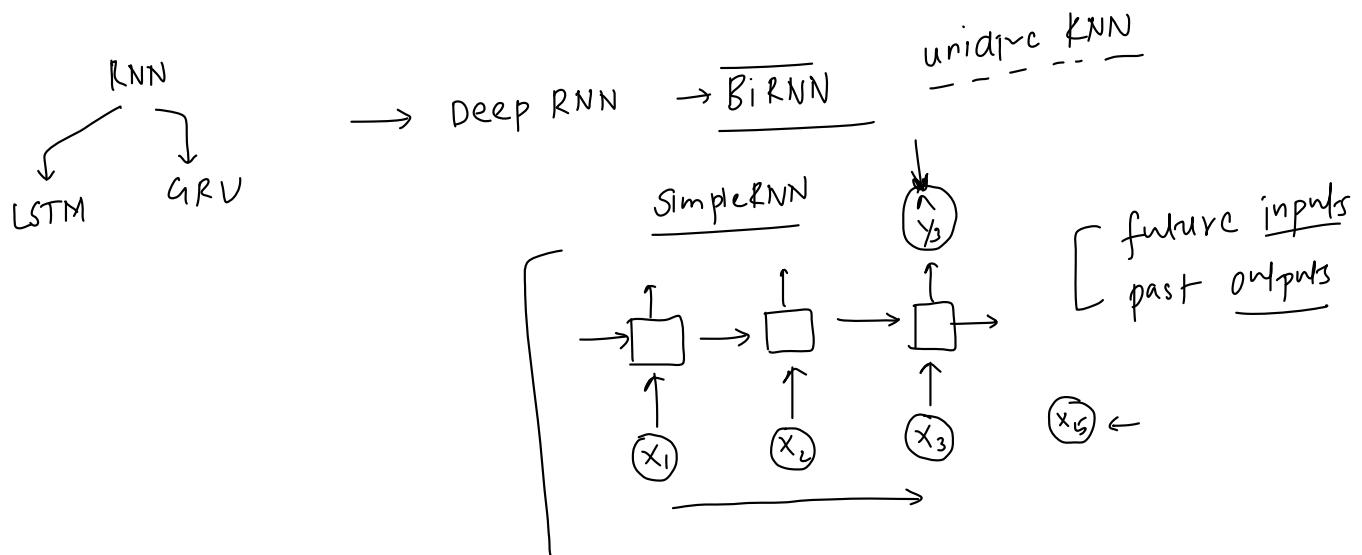
17 October 2023 16:30



Disadvantages

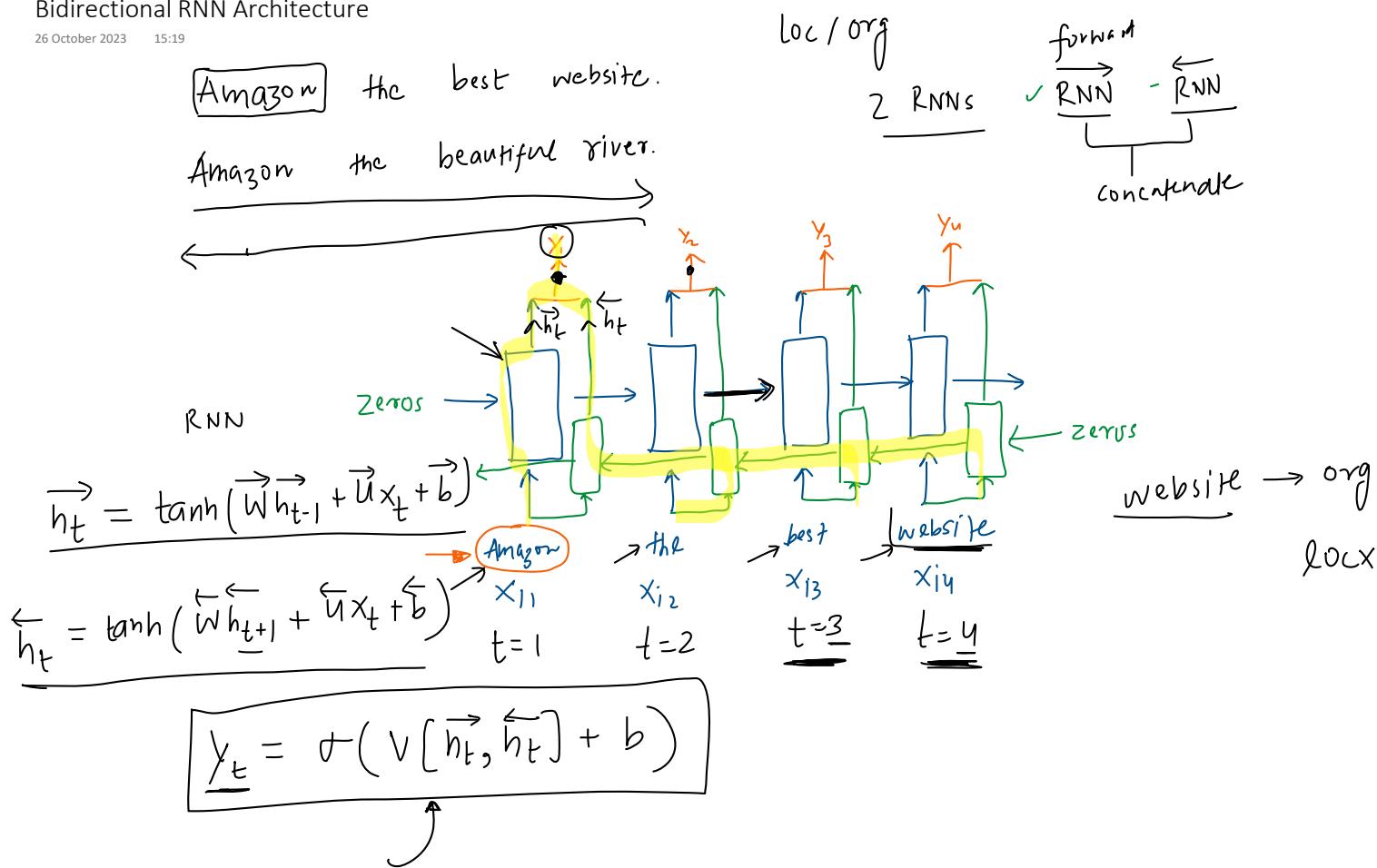
17 October 2023 16:30





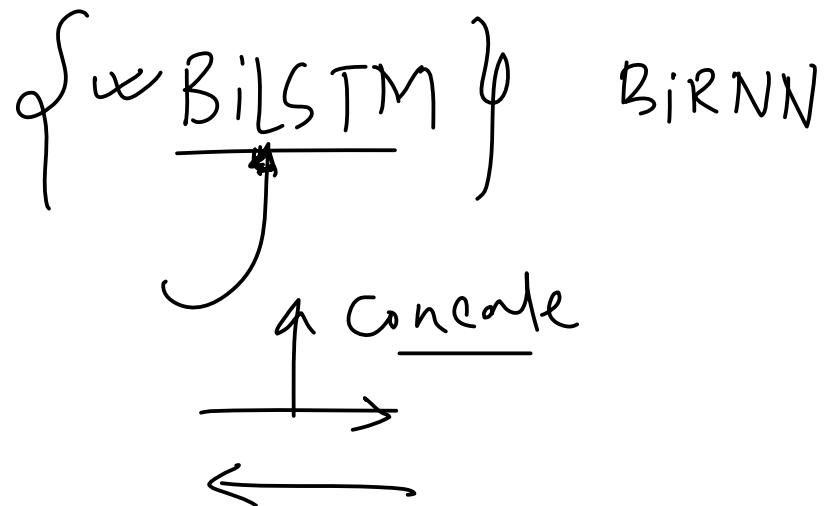
Bidirectional RNN Architecture

26 October 2023 15:19



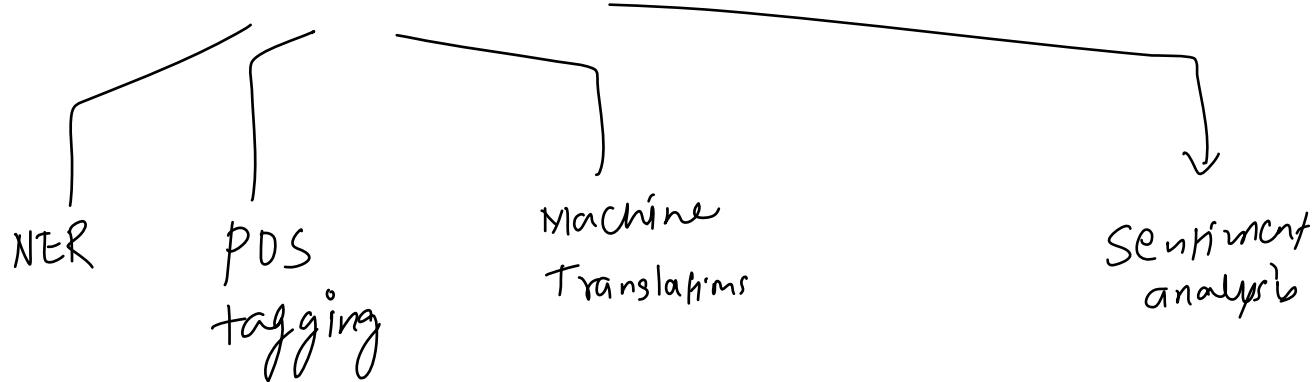
Code

26 October 2023 15:21



Applications and Drawbacks

26 October 2023 15:21



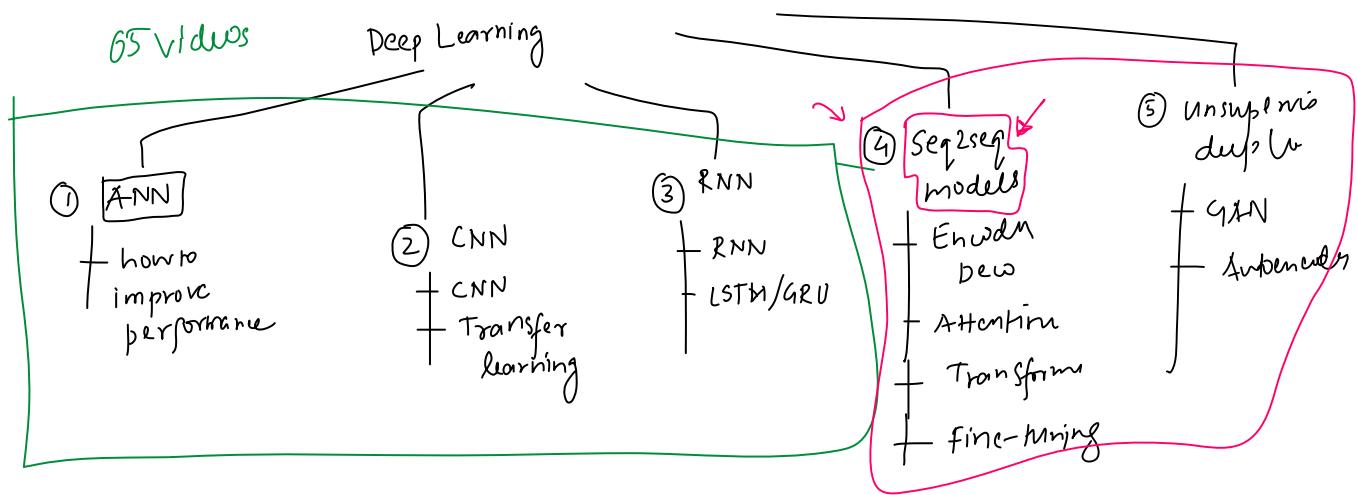
[Time series forecasting]

→ ←

→ Complexity → 190 → 380
↓ ↓
downdown

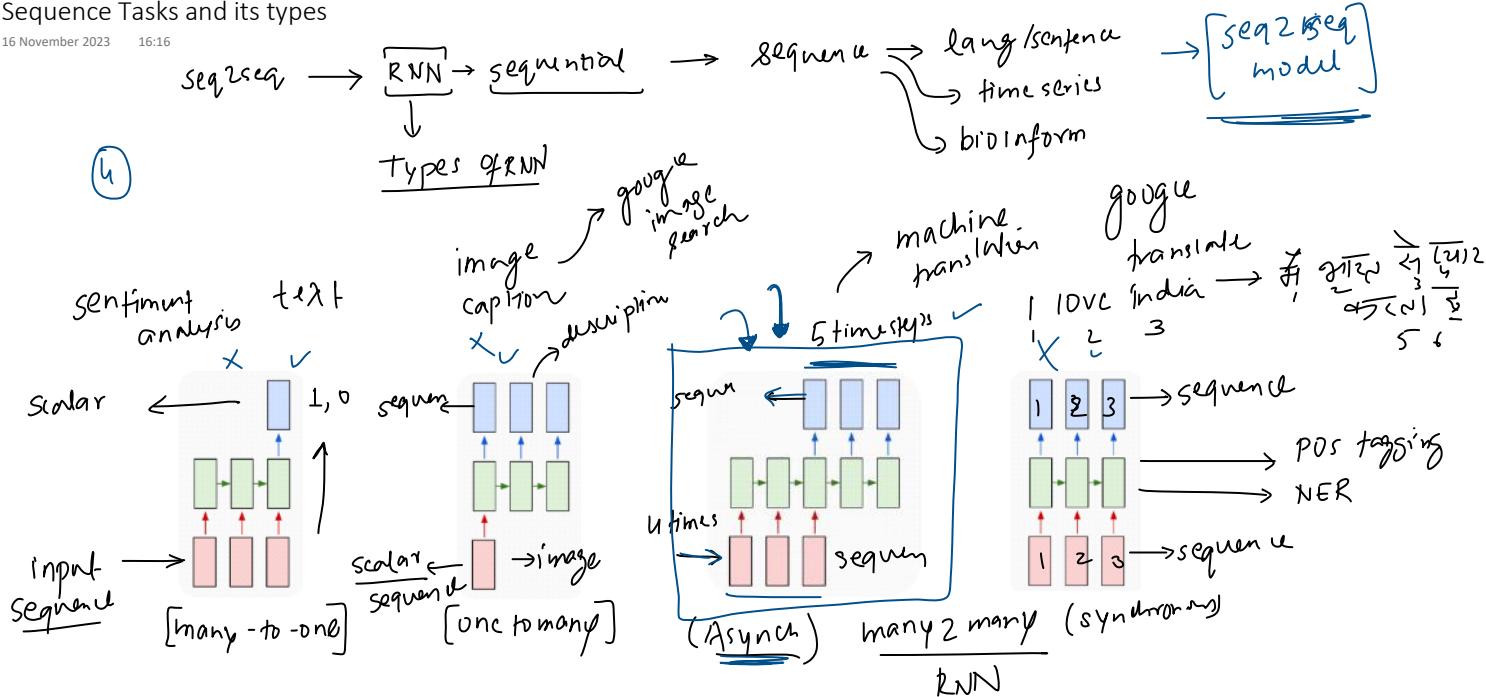
→ training → overfitting

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



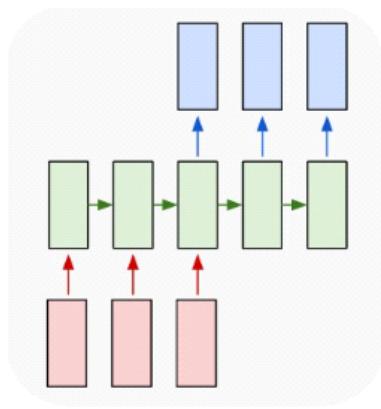
Sequence Tasks and its types

16 November 2023 16:16



Seq2Seq tasks

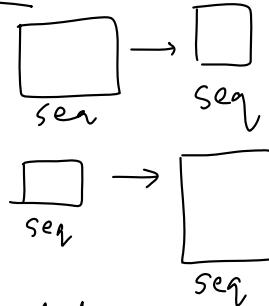
16 November 2023 16:16



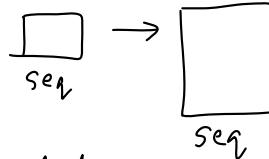
NLP

Seq2seq → machine trans

1) text summariz →



2) Question answer →



3) chatbot → input (text) → output (*ex)

4) speech-to-text →

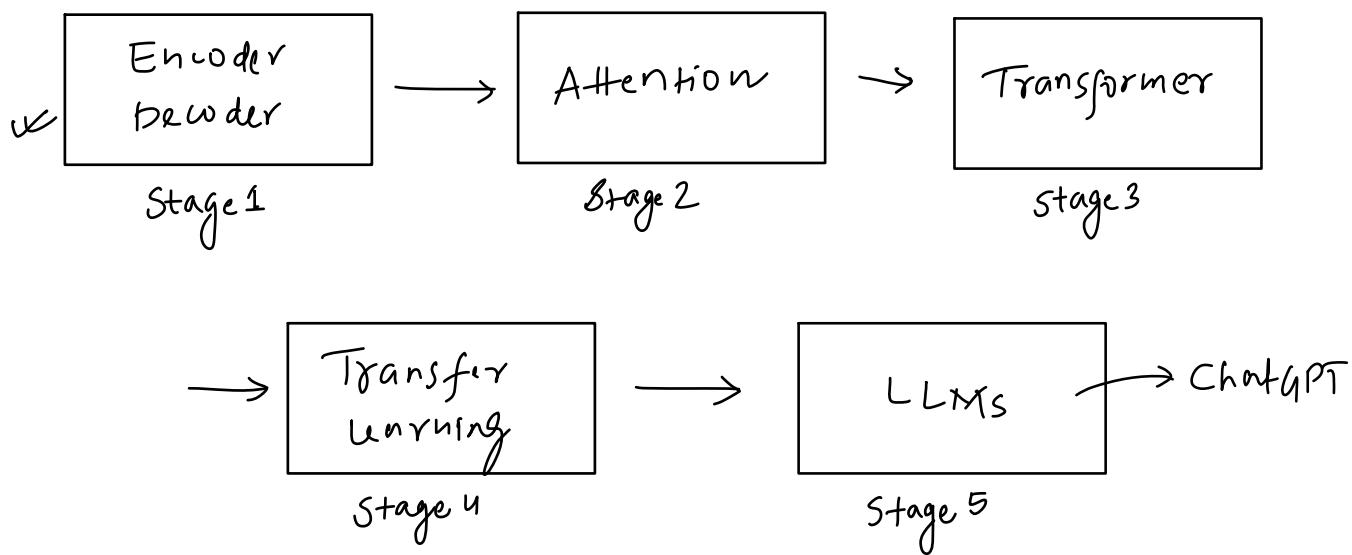
seq seq

knowledge base

History of Seq2Seq Models

16 November 2023 16:16

ChatGPT



2014 Seminal

seq2seq
 ↓
 diff
 ↳ encoder
 decoder

Sequence to Sequence Learning with Neural Networks

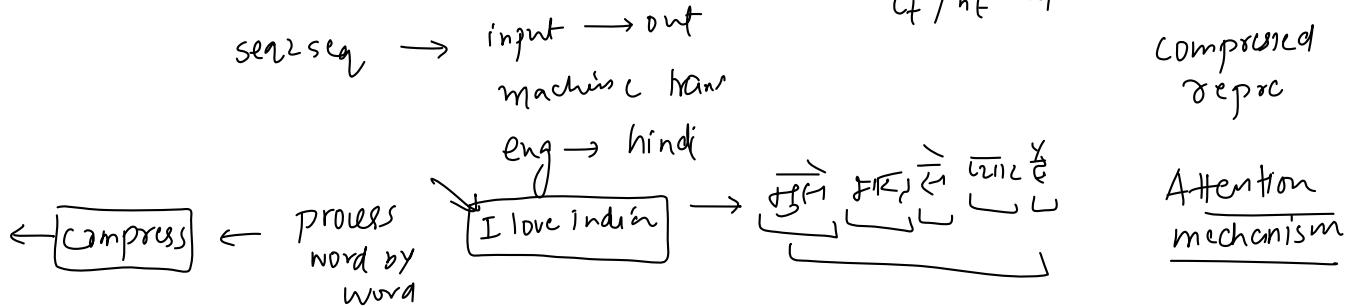
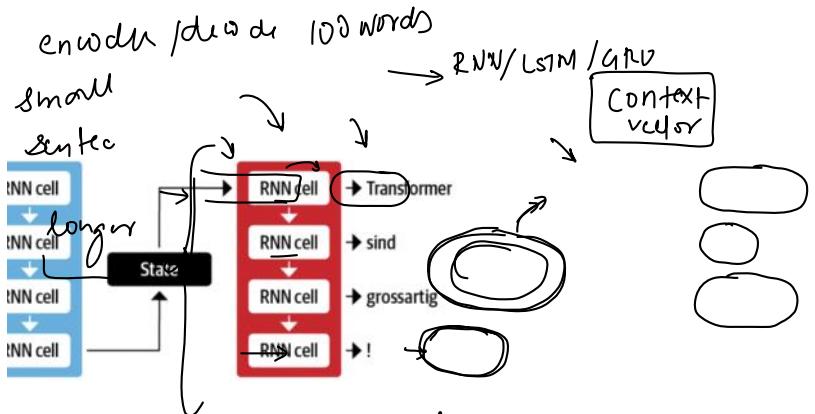
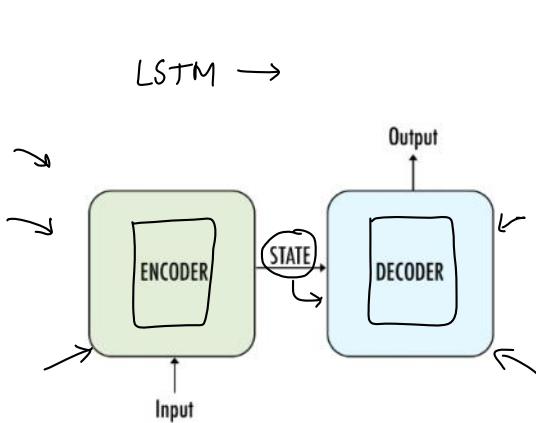
→ Ilya Sutskever
 Google
 ilyasu@google.com

[Oriol Vinyals]
 Google
 vinyals@google.com

[Quoc V. Le]
 Google
 qvl@google.com

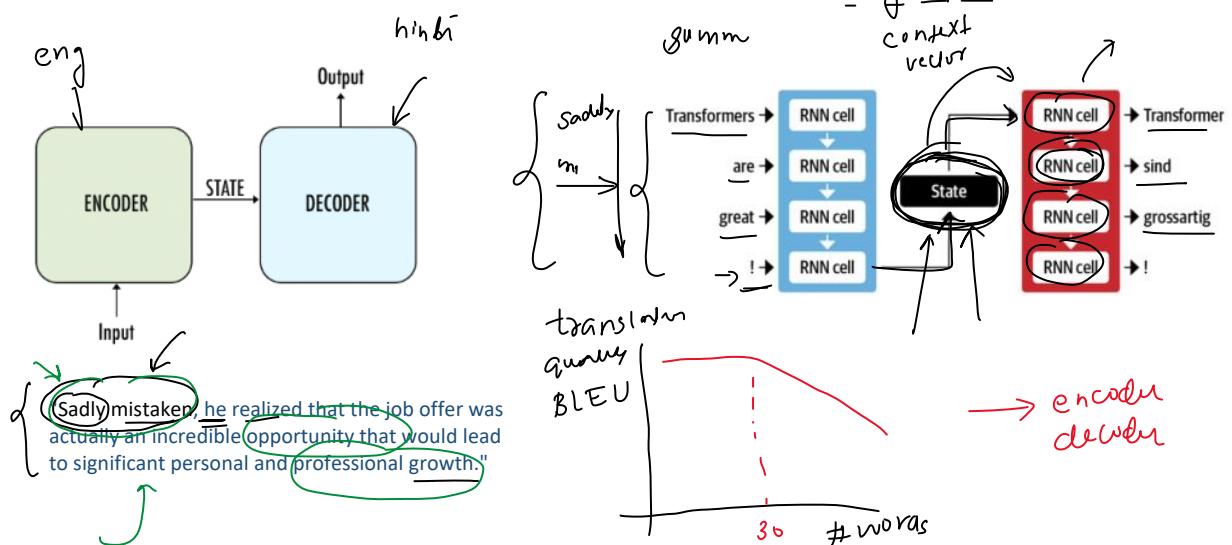
Abstract

Deep Neural Networks (DNNs) are powerful models that have achieved excellent performance on difficult learning tasks. Although DNNs work well whenever large labeled training sets are available, they cannot be used to map sequences to sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a multilayered Long Short-Term Memory (LSTM) to map the input sequence to a vector of a fixed dimensionality, and then another deep LSTM to decode the target sequence from the vector. Our main result is that on an English to French translation task from the WMT'14 dataset, the translations produced by the LSTM achieve a BLEU score of 34.8 on the entire test set, where the LSTM's BLEU score was penalized on out-of-vocabulary words. Additionally, the LSTM did not have difficulty on long sentences. For comparison, a phrase-based SMT system achieves a BLEU score of 33.3 on the same dataset. When we used the LSTM to rerank the 1000 hypotheses produced by the aforementioned SMT system, its BLEU score increases to 36.5, which is close to the previous best result on this task. The LSTM also learned sensible phrase and sentence representations that are sensitive to word order and are relatively invariant to the active and the passive voice. Finally, we found that reversing the order of the words in all source sentences (but not target sentences) improved the LSTM's performance markedly, because doing so introduced many short term dependencies between the source and the target sentence which made the optimization problem easier.



Stage 2 - Attention Mechanism

20 November 2023 10:59



2015 → A Henkim

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau
Jacobs University Bremen, Germany
KyungHyun Cho [Yoshua Bengio]
Université de Montréal

ABSTRACT

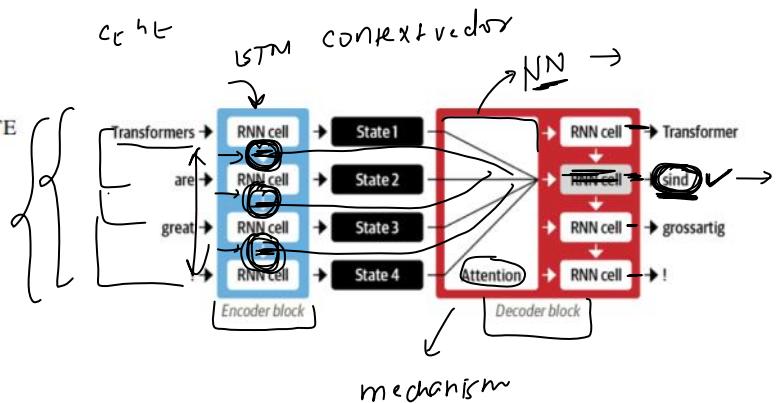
Neural machine translation is a recently proposed approach to machine translation. Unlike the traditional statistical machine translation, the neural machine translation aims at building a single neural network that can be jointly tuned to maximize the translation performance. The models proposed recently for neural machine translation often belong to a family of encoder-decoders and encode a source sentence into a fixed-length vector from which a decoder generates a translation. In this paper, we conjecture that the use of a fixed-length vector is a bottleneck in improving the performance of this basic encoder-decoder architecture, and propose to extend this by allowing a model to automatically (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly. With this new approach, we achieve a translation performance comparable to the existing state-of-the-art phrase-based system on the task of English-to-French translation. Furthermore, qualitative analysis reveals that the (soft-)alignments found by the model agree well with our intuition.

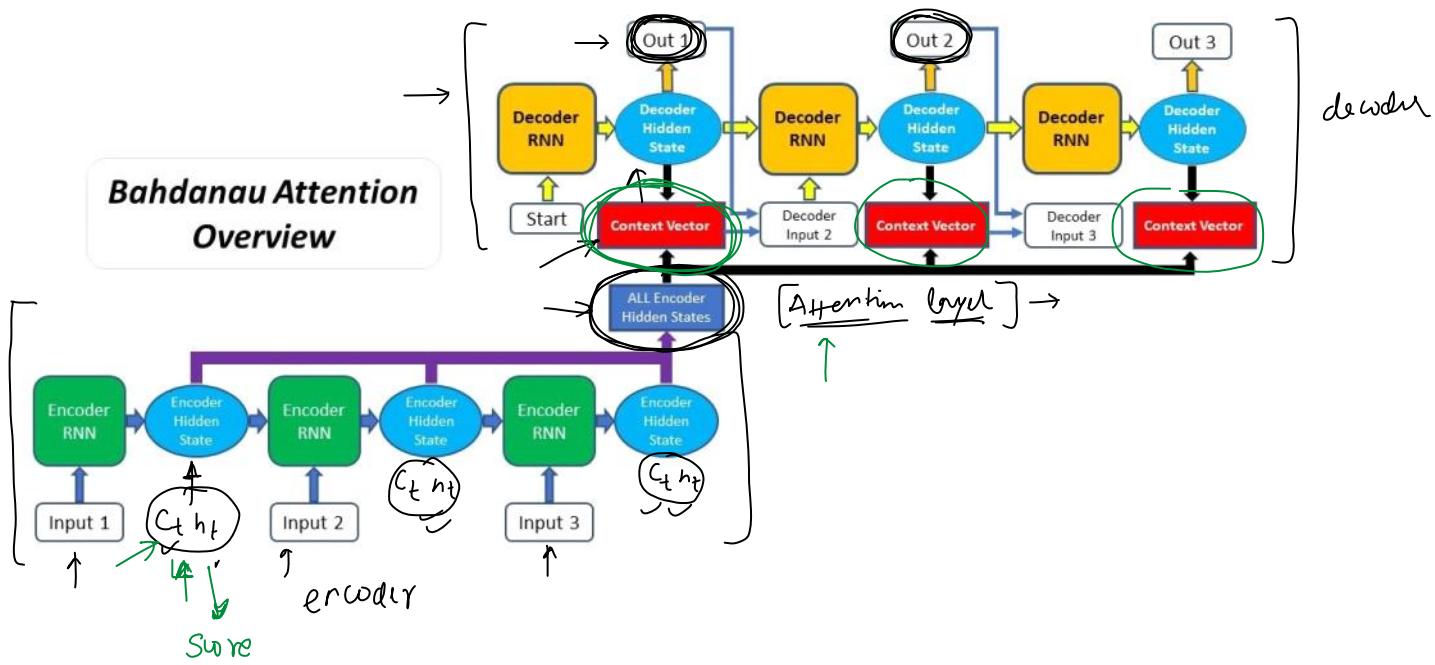
1 INTRODUCTION

Neural machine translation is a newly emerging approach to machine translation, recently proposed by Kalchbrenner and Blunsom (2013), Sutskever et al. (2014) and Cho et al. (2014b). Unlike the traditional phrase-based translation system (see, e.g., Koehn et al., 2003) which consists of many small sub-components that are tuned separately, neural machine translation attempts to build and train a single, large neural network that reads a sentence and outputs a correct translation.

Most of the proposed neural machine translation models belong to a family of *encoder-decoders* (Sutskever et al., 2014; Cho et al., 2014a), with an encoder and a decoder for each language, or involve a language-specific encoder applied to each sentence whose outputs are then compared (Hermann and Blunsom, 2014). An encoder neural network reads and encodes a source sentence into a fixed-length vector. A decoder then outputs a translation from the encoded vector. The whole encoder-decoder system, which consists of the encoder and the decoder for a language pair, is jointly trained to maximize the probability of a correct translation given a source sentence.

A potential issue with this encoder-decoder approach is that a neural network needs to be able to compress all the necessary information of a source sentence into a fixed-length vector. This may make it difficult for the neural network to cope with long sentences, especially those that are longer than the sentences in the training corpus. Cho et al. (2014b) showed that indeed the performance of a basic encoder-decoder deteriorates rapidly as the length of an input sentence increases.





Stage 3 - Transformers
20 November 2023 12:18

{ computational complexity }

m words

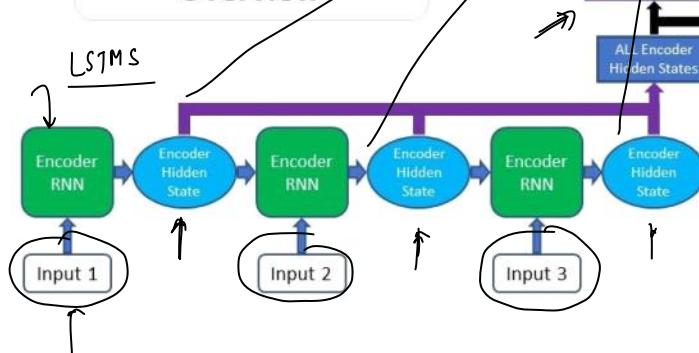
2015-2017

np

m words

after

Bahdanau Attention Overview



n words

sequential order

en wieder diewel

parallel processing

2017

Attention Is All You Need

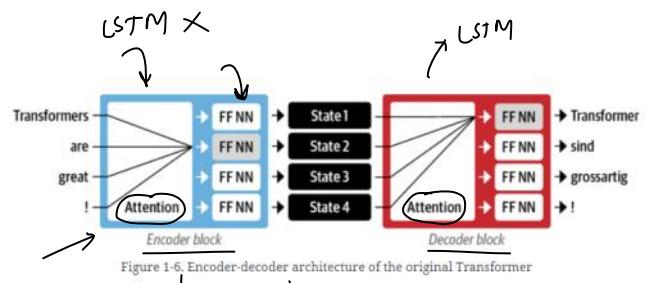
Ashish Vaswani*	Noam Shazeer*	Niki Parmar*	Jakob Uszkoreit*
Google Brain	Google Brain	Google Research	Google Research
avaswani@google.com	noam@google.com	nikip@google.com	usz@google.com

Llion Jones*	Aidan N. Gomez* [†]	Lukasz Kaiser*
Google Research	University of Toronto	Google Brain
llion@google.com	aidan@cs.toronto.edu	lukaszkaiser@google.com

Ilia Polosukhin* [‡]
ilia.ilia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.



LSTM / RNN cell

Attention

Self-attention

stage

arch

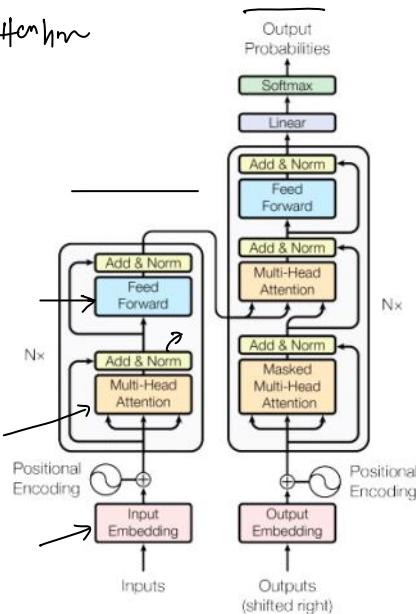
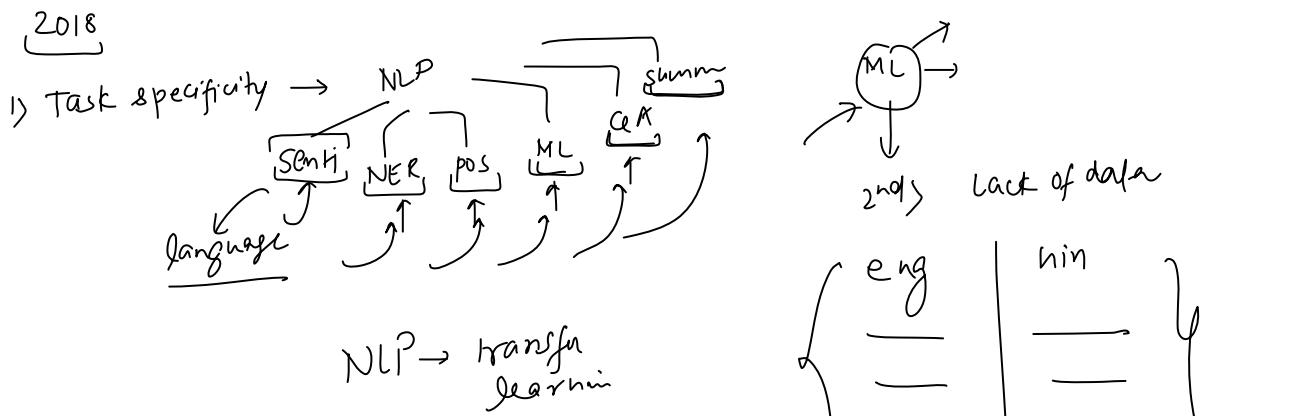
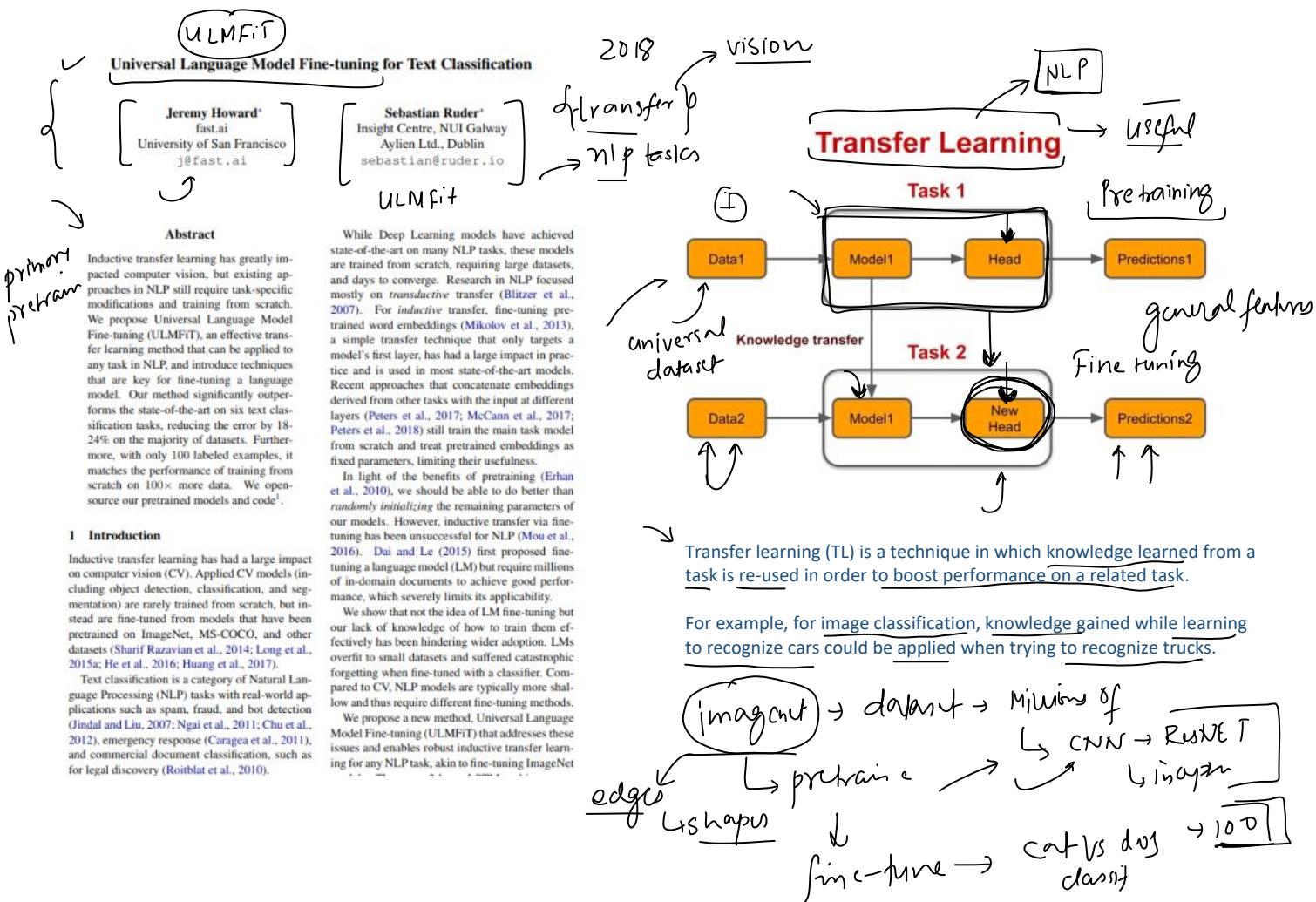
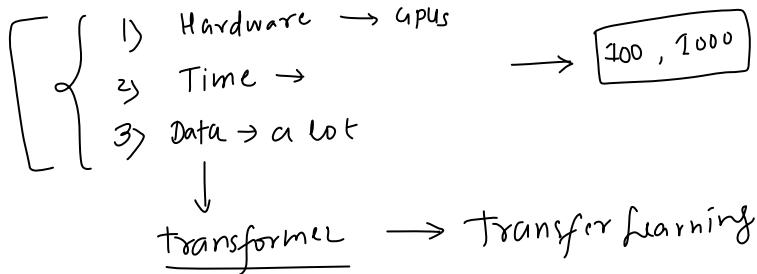
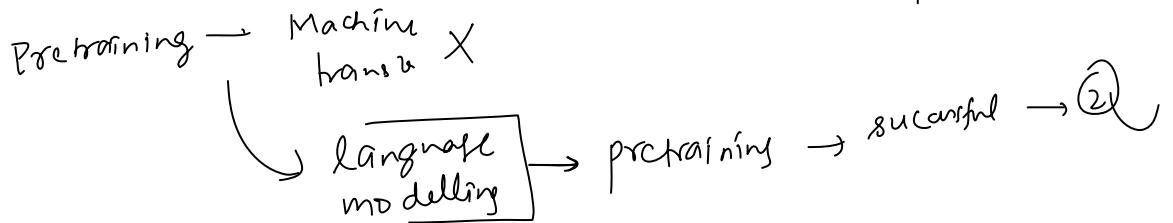
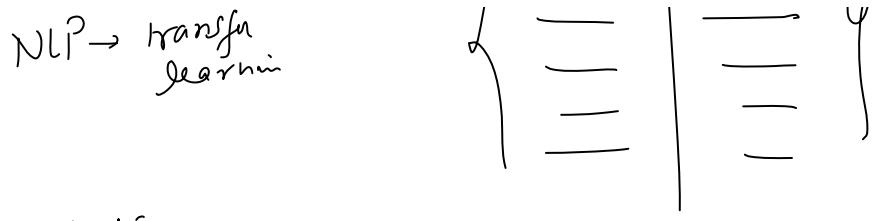


Figure 1: The Transformer - model architecture.

Stage 4 - Transfer Learning

20 November 2023 15:39





NLP task → NLP/PL model next word pred
 I live in India. and the capital is New Delhi

1) Rich feature learning
Language modeling as a Pretraining task
 The hotel was exceptionally clean, yet the service was bad ↓
pathetic

→ know trans
 ↓
 text classif / ques. | textsum) NLP / PLM

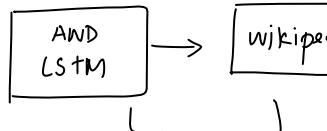
mt (kew → supervised
 eng | hin labeled
 → unsupervised task

2) Huge avail of data
 pdf → dataset
 labelling

fine tuning

[ULMFIE]

X transformer



Unsupervised
 pretrain
 Language
 modeling

classifier

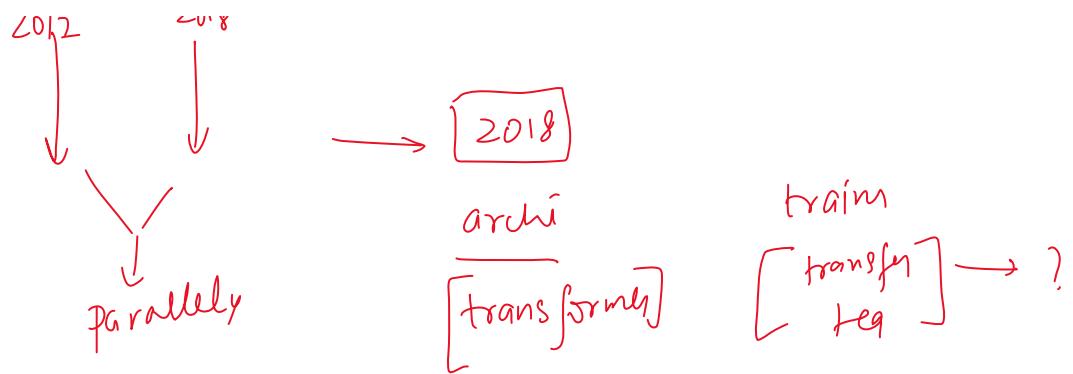
imdb
 yelp
 new dataset

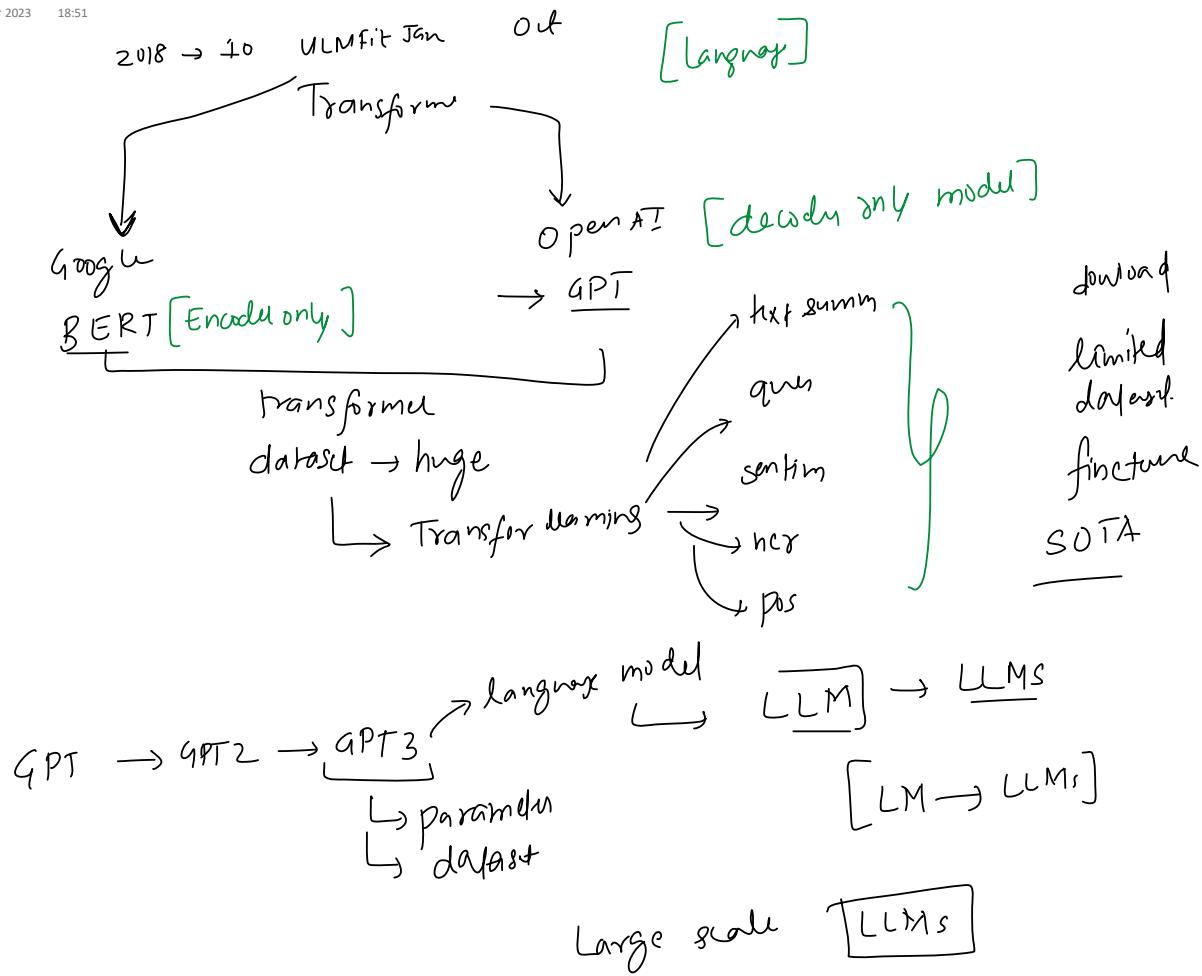
Scratch → 10000 rows
 100 row → better →

model
 ↓
 test

State of the art

2012 2018



Qualities of LLMs

1) Data → billions → GPT3 → 45 TBs
 ↗ book, websites, internet
 ↗ diversity → bias

2) Hardware → Cluster of GPU → GPT3 → Supercomputer → 100s NVIDIA GPU

3) Training → days to wccs

4) Cost → hardware + elec + infra + experiments → individual
 ↗ millions → companies, govt, institutes

4) energy consumption
 ↗ GPT3 → ...