

NLP in Deep learning

Text data \rightarrow Vector \rightarrow Numerical Representations

① BPE

② BOW

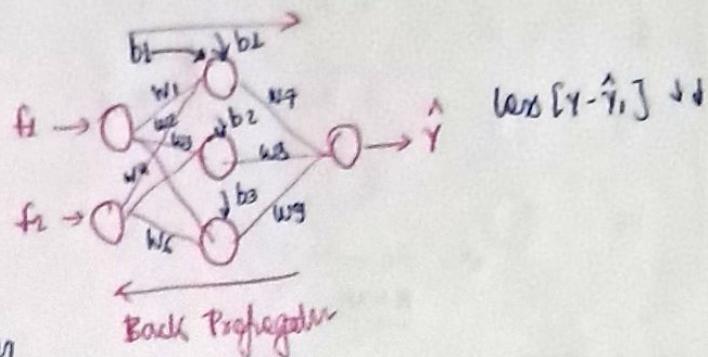
③ Tf-IDF

④ Word2Vec, AvgWord2Vec & Sentence Analysis, Text Classifier

⑤ ANN \rightarrow Artificial Neural Network \rightarrow Classifier / Regressor - Tabular Data

Eg:- House price Predictor

f_1	f_2	y
House Size	No. of Room	Price
-	-	-



⑥ Sequence of Data Doesn't Matter

⑦ CNN (Convolutional Neural Network) \rightarrow Images, Video Frames

Eg \rightarrow Image Classification, Object Detection

⑧ Data \rightarrow Sequential Data. I/p \rightarrow Juice

① Text Generation \rightarrow This is a Apple \rightarrow juice

② Chatbot Conversation \rightarrow Q & A \rightarrow Ans

(Sequential Data \rightarrow The food is good \rightarrow good the food is \rightarrow Meaning may change)

③ Language Translation \rightarrow English \rightarrow French

④ Auto Suggestion \rightarrow LinkedIn, Gmail \rightarrow Auto suggestion

⑤ Sales Data \rightarrow Date Time \rightarrow Sales \rightarrow Sale forecast

Can we solve this using ANN?

NLP in Deep L [Given A \rightarrow LLM, multiModel]

① Simple RNN \rightarrow LSTM / GRU RNN \rightarrow Bidirectional RNN \rightarrow Encoder

Dataset of sentiment analysis

Text O/p

The food is good 1
The food is bad 0
The food is not good 0

Transform \leftarrow Self Attention \leftarrow Decoder

ANN

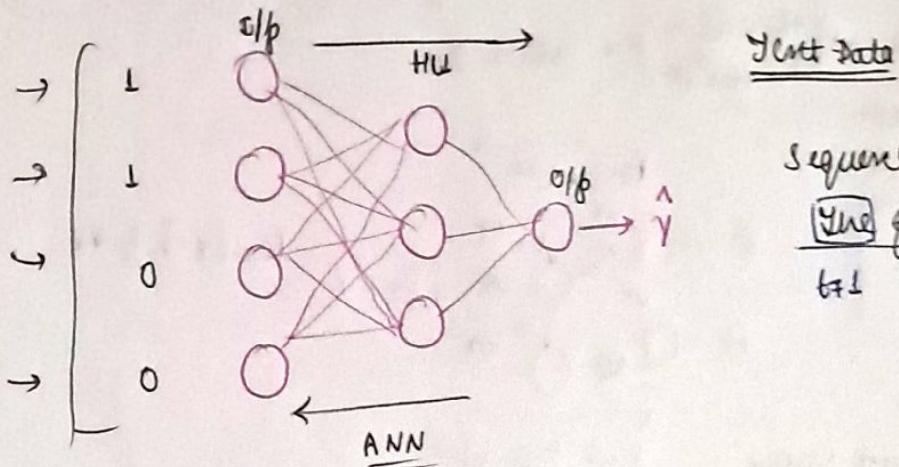
Text \rightarrow Text \rightarrow Vector

Vocabulary Size \rightarrow 4

BOW, TF-IDF

Text Data \rightarrow Sequence of Information is Imp
Meaning of the Sentence is less.

feed good bad not
 S1 [1 1 0 0]
 S2 [1 0 1 0]
 S3 [1 0 0 1]



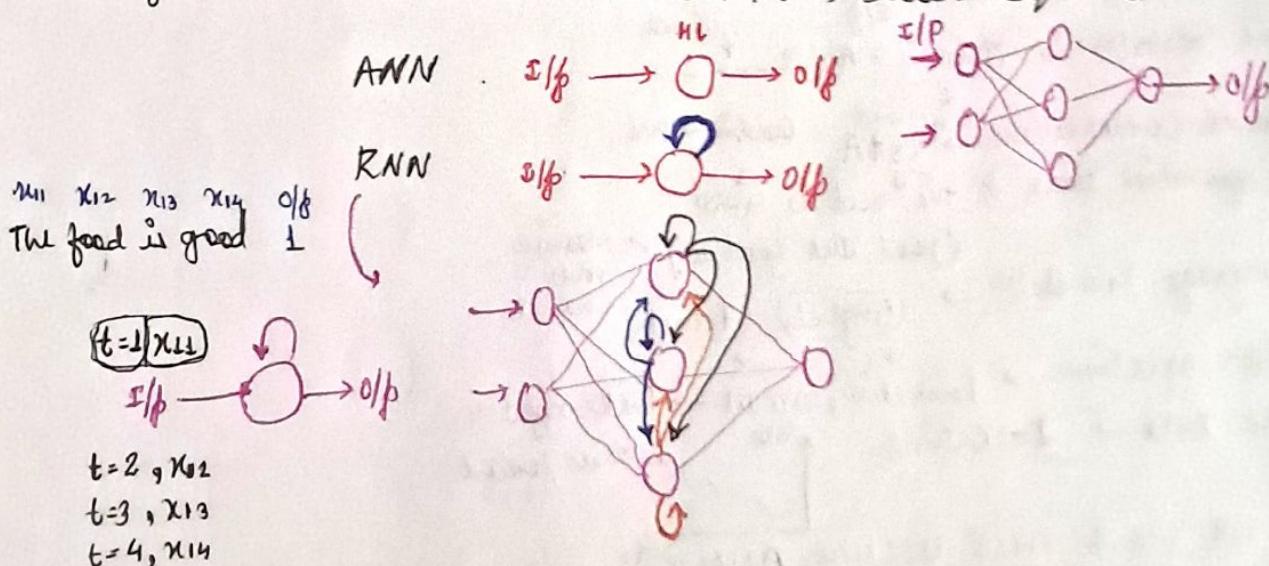
Sequence of Information

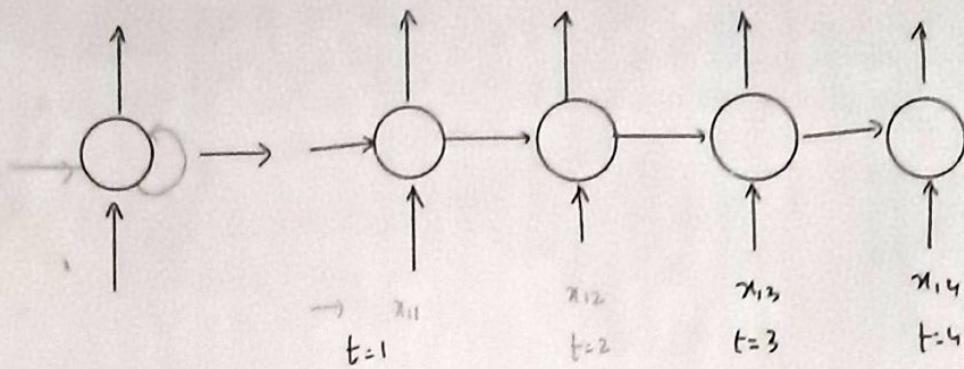
The food is good
t=1 t=2 t=3 t=4

Google Translate \rightarrow English \rightarrow French

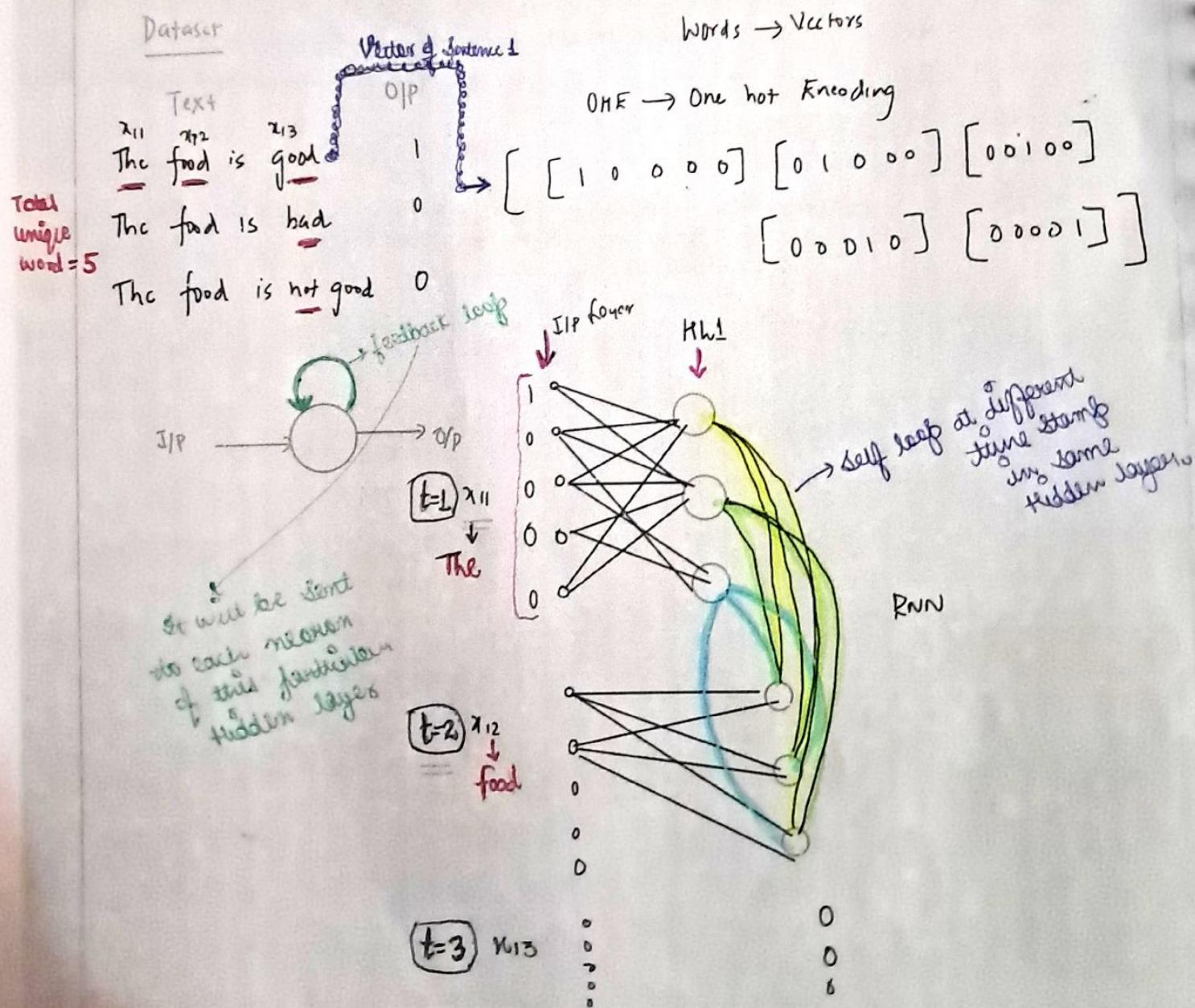
The --- \downarrow ---

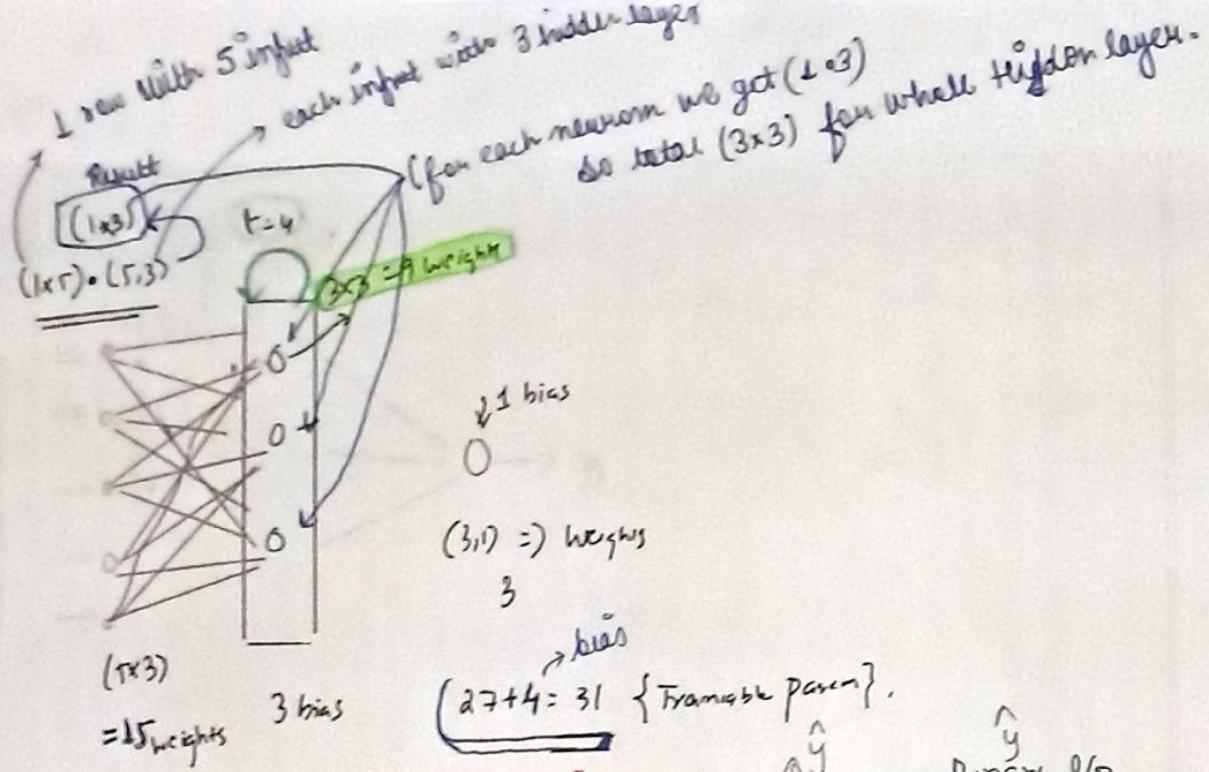
④ Simple RNN \rightarrow Recurrent Neural N/w \rightarrow Data \rightarrow Sequence



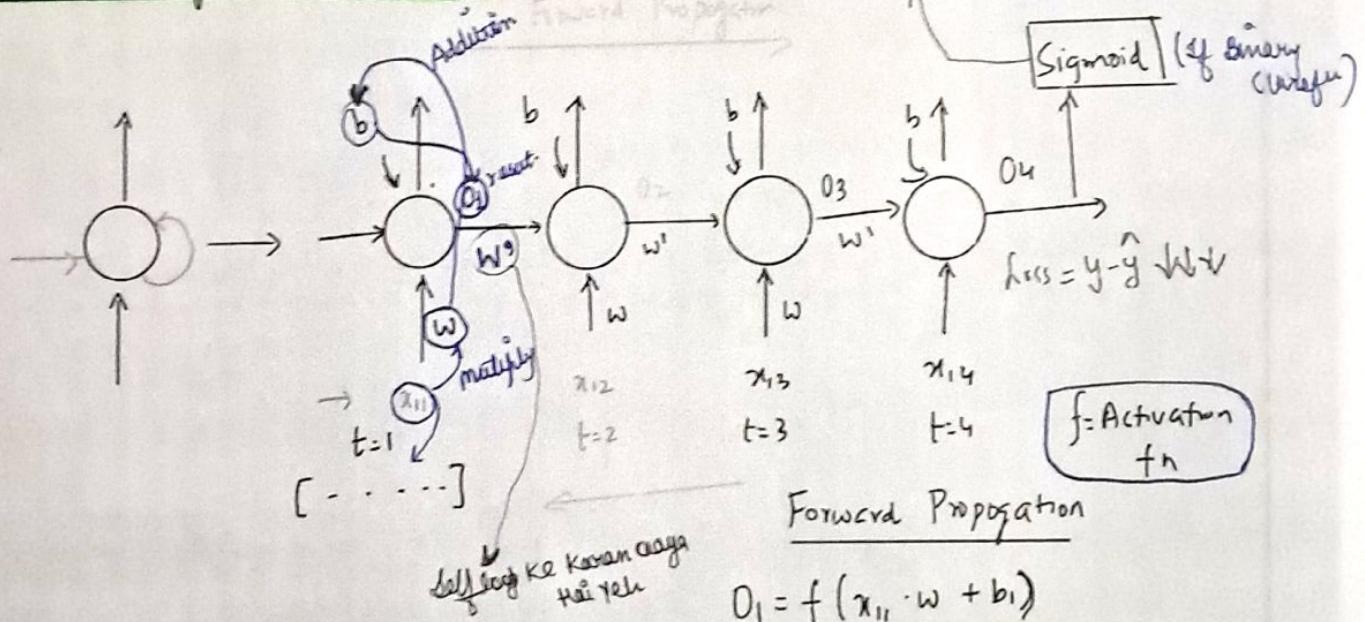


Working Of Simple RNN With Forward Propagation





Forward Propagation With Time (This whole diagram is same hidden layer)



Dataxt

The food is good
 $x_{11} \quad x_{12} \quad x_{13} \quad x_{14}$

O/P

1

$$O_1 = f(x_{11} \cdot w + b_1)$$

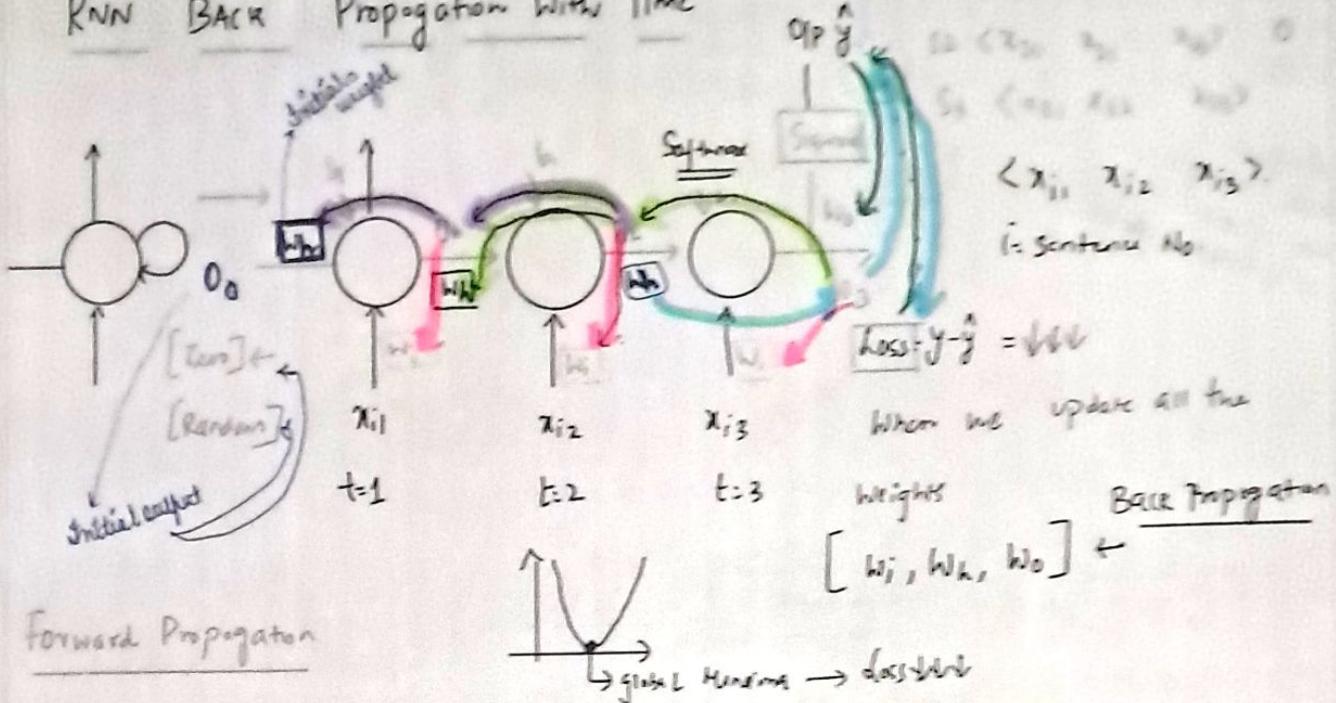
$$O_2 = f(x_{12} \cdot w + O_1 \cdot w + b_1)$$

$$O_3 = f(x_{13} \cdot w + O_2 \cdot w + b_1)$$

Any activation func

This step happens
 Bcoz of w_0 we
 get from
 feedback loop

RNN BACK Propagation with Time



Forward Propagation

$$O_1 = f(x_{i1} * w_i + h_0 * w_h + b_1)$$

$$O_2 = f(x_{i2} * w_i + O_1 * w_h + b)$$

$$O_3 = f(x_{i3} * w_i + O_2 * w_h + b)$$

$$\hat{y} = \sigma(O_3 + w_o)$$

① Update w_o

Chain Rule of Derivative

$$\frac{\partial L}{\partial w_{oold}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{oold}}$$

Softmax
 w_o update
topo.

$$w_{oold} = w_{oold} - \eta \frac{\partial L}{\partial w_{oold}}$$

Back Propagation

② Update w_h [Hidden layer weights] \rightarrow Time Stamps

$$w_{hnew} = w_{hold} - \eta \frac{\partial L}{\partial w_{hhold}}$$

$t=1, 2, 3$

Backward Propagation with Time

\rightarrow Update $[w_i, w_h, w_o]$

Weight update formula

$$w_{hnew} = w_{hold} - \eta \frac{\partial L}{\partial w_{hold}}$$

Derivative, slope of
GRADIENT DESCENT
to calculate

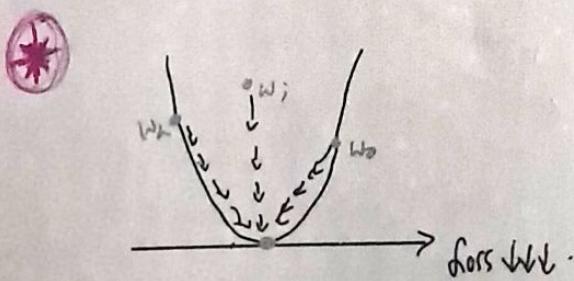
$$\Rightarrow \frac{\partial L}{\partial w_{h,old}} = \left[\frac{\partial L}{\partial y} * \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial w_h} \right] + \left[\frac{\partial L}{\partial y} * \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial w_h} \right] + \left[\frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_h} \right]$$

In One Back Propagation All Hidden weight (w_h) Updated

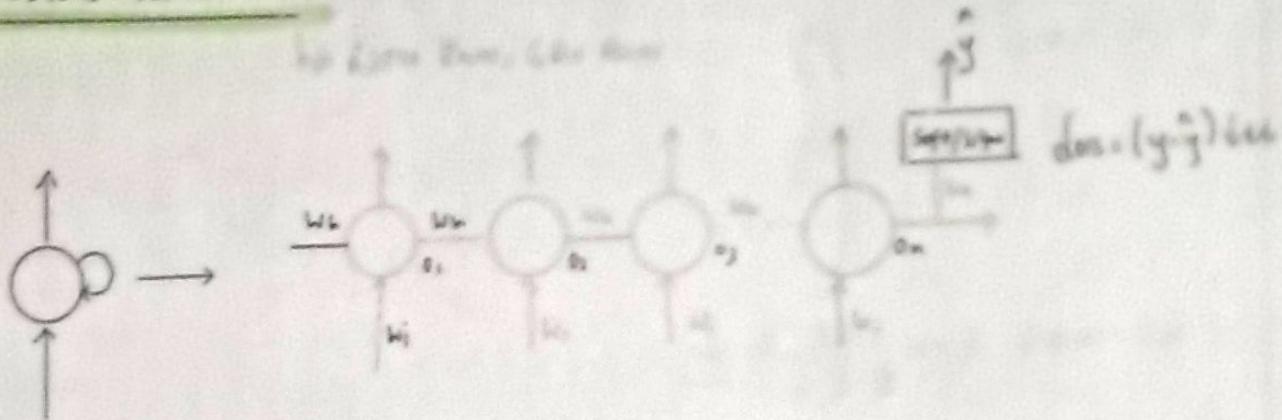
③ Updating weights w_i → Timestamp

$$w_{i,new} = w_{i,old} - \eta \left[\frac{\partial L}{\partial w_{i,old}} \right] \quad \{ \text{update } w_i \text{ weights} \}$$

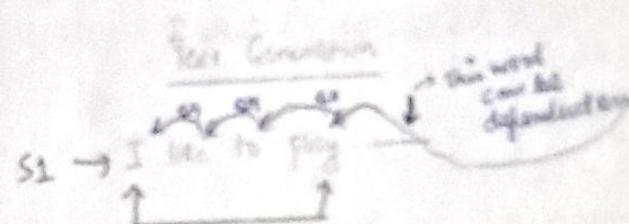
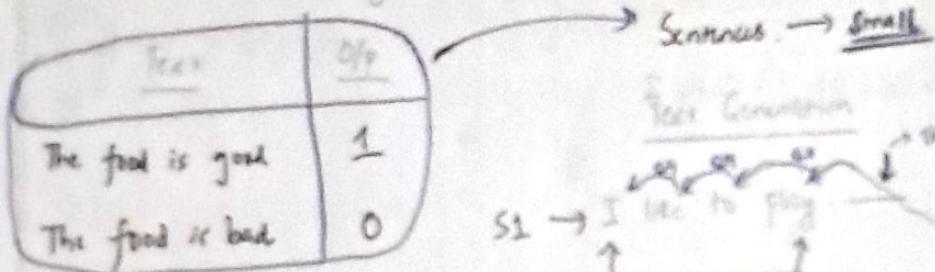
$$\frac{\partial L}{\partial w_{i,old}} = \left[\frac{\partial L}{\partial y} * \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial w_{i,old}} \right] + \left[\frac{\partial L}{\partial y} * \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial w_{i,old}} \right] + \left[\frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{i,old}} \right]$$



Problems With RNN



ANN \rightarrow Vanishing Gradient Problem



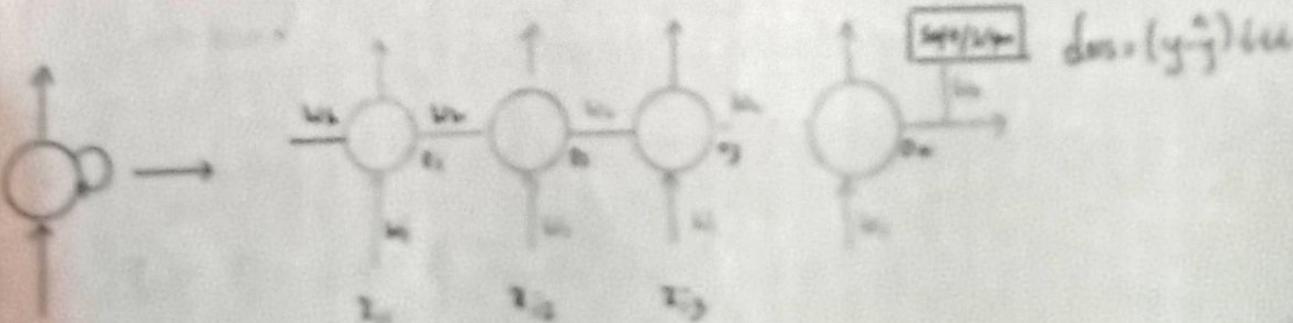
Sequence is
long

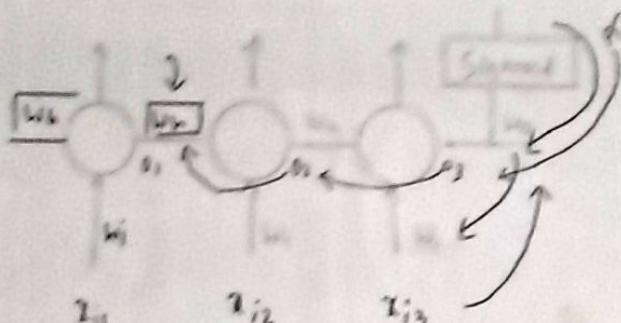
{ My Name is KRISH and I like sports like
CRICKET, VOLLEYBALL AND ALSO LIKE TO MAKE

if u have sentence
with ≈ 100 words. So at $t=100$ we will find last word, for \rightarrow ~~long~~ happy

① Long term dependency cannot be captured by RNN

provide Accuracy
↓
Dependency





outputs w_h, w_i, w_o

$$w_{new} = w_{old} - \eta \frac{\partial L}{\partial w_{old}}$$

$x_{i1} \quad x_{i2} \quad x_{i3}$
 $t=1 \quad t=2 \quad t=3$

for example if we want it for
 $t=1$ ← it for

$$\frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial w_o} \leftarrow$$

$t=2$

$$\frac{\partial L}{\partial w_{h, old}} = \left[\frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_{h, old}} \right] +$$

for word = Only 3

length of sentence \Rightarrow 50 words

$t=1$

$t=2$

$$\frac{\partial L}{\partial w_{h, old}} = \left[\frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial o_{50}} \cdot \frac{\partial o_{50}}{\partial o_{49}} \left[\frac{\partial o_{49}}{\partial o_{48}} \frac{\partial o_{48}}{\partial w_{h, old}} \right] + \frac{0.25 \times 0.02 \times 0.01 \times 0.015}{\frac{\partial o_{49}}{\partial o_{48}}} \right]$$

↓

Small
Value

≈ 0

$$\frac{\partial o_3}{\partial o_2} = \frac{\partial \sigma}{\partial o_2} = \frac{\partial \sigma}{\partial (x_{i3} * w_i + o_2 * w_h + b)} \leftarrow o_3 = \sigma(x_{i3} * w_i + o_2 * w_h + b)$$

\downarrow

$\frac{\partial \sigma}{\partial o_2}$ (Because differentiation is wrt o_2)

$\frac{\partial \sigma}{\partial (x_{i3} * w_i + o_2 * w_h + b)}$ $\leftarrow o_3 = \sigma(x_{i3} * w_i + o_2 * w_h + b)$

$\frac{\partial \sigma}{\partial o_2}$ constant wrt o_2
 ≈ 0.25

$\frac{\partial \sigma}{\partial (x_{i3} * w_i + o_2 * w_h + b)}$ Derivative of sigmoid
 $0 - 0.25$

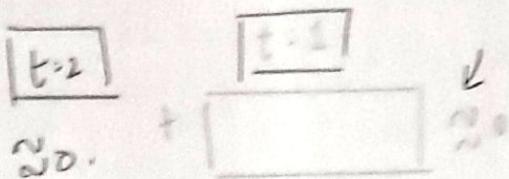
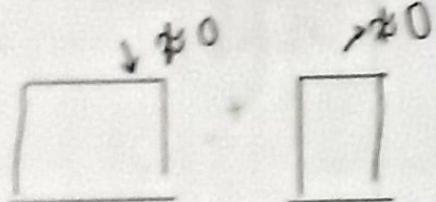
$$= \sigma'(1 + w_h) \Rightarrow [0 - 0.25] \sigma'(w_h)$$

$\downarrow t=1 \Rightarrow$ The word is not participating to update the weights
 Value. (Because weight to change w_h is 0)

Chain Rule is Big

Conclusion
 At $t=0$ whole \hat{y} from $t=0$ update h_0
 But $t=1$ whole \hat{y} update h_1 h_0
 So, output will depend on
 $t=0$ word to predict
 depending on $t=0$ word

$$\frac{\partial h}{\partial w_{word}} = \frac{\partial h}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \hat{o}_{t=0}} + \frac{\partial \hat{o}_{t=0}}{\partial w_{word}}$$



To solve this

① Recu, Leaky Recu \rightarrow

② LSTM RNN \rightarrow long Short term Memory RNN } \Rightarrow Simple RNN

③ GRU RNN \rightarrow

LSTM RNN

→ Long Short Term Memory [LSTM]

RNN → Long Term Dependence → Vanishing Gradient Problem

① RNN → Problem? ✓

② Why LSTM RNN? ✓ Basic Representation.

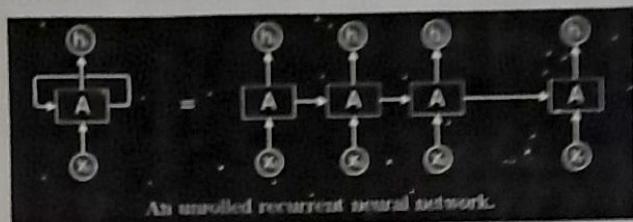
③ How LSTM RNN Works

→ Long Term Memory ✓ ↵
→ Short Term Memory ✓ ↵

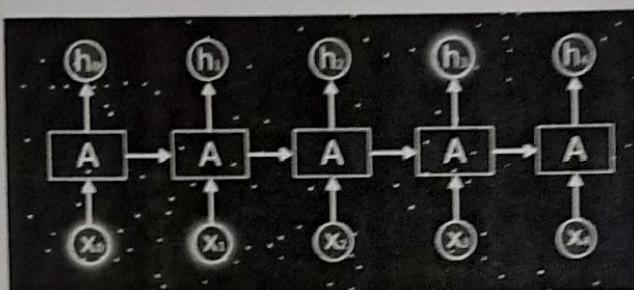
④ LSTM Architecture

⑤ Working of LSTM RNN

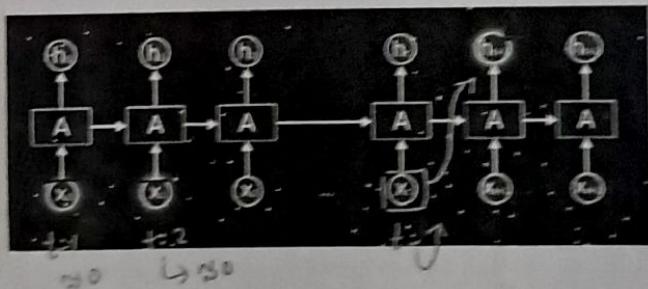
Problems with RNN \rightarrow Long Term Dependency



Vanishing Gradient Problem



Gap is 1cm



Huge gap \rightarrow Long Term Dependency

0-0.25

0-1

RNN \rightarrow Long Term Dependency \rightarrow Vanishing Gradient Problem

Chain Rule \rightarrow ≈ 0

Short term dependency

Task

Next word in a sentence

The color of the sky is blue
= further context

Huge gap O/p \leftarrow (over \rightarrow some other sentence)

I grew up in India \cdots I speak

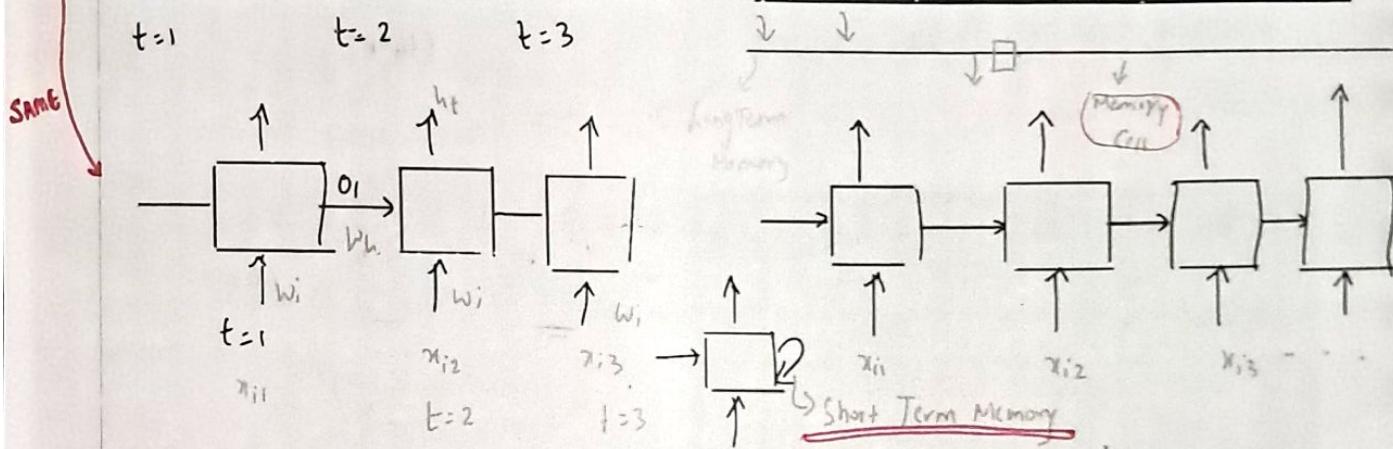
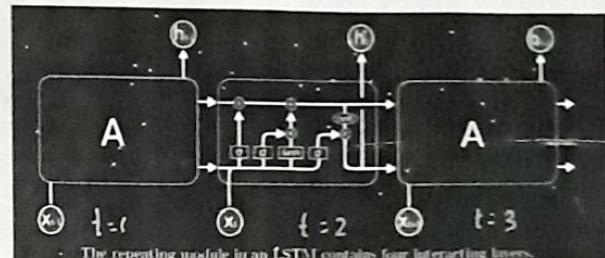
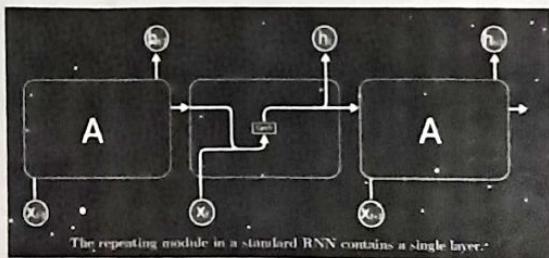
fluent \leftarrow \rightarrow language

Name of language

further context \rightarrow dependent on Sentence #2 \downarrow
Long term dependency

Basic Representation of RNN And LSTM RNN

LSTM RNN

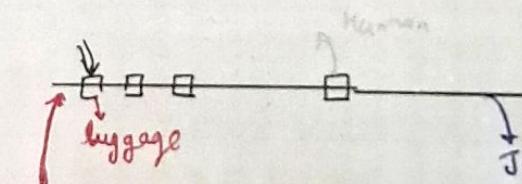


LSTM RNN

- Long Term Memory
- Short Term Memory

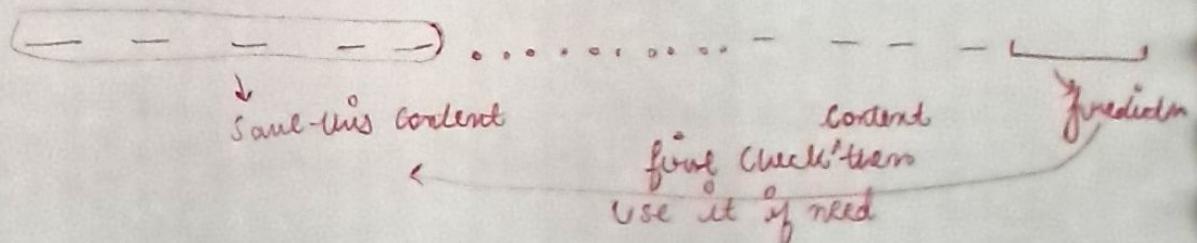
Airport

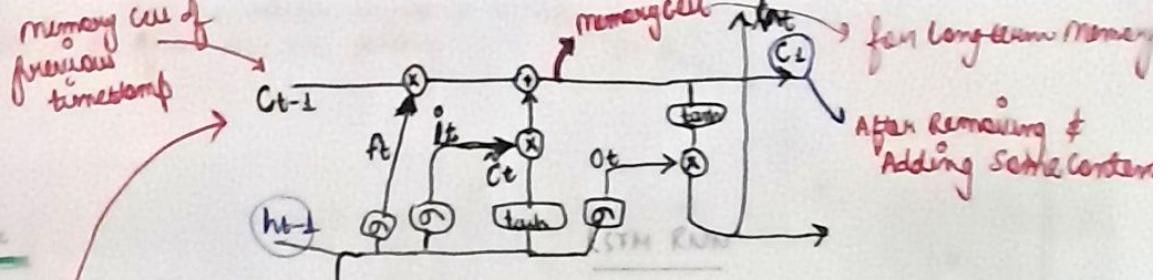
Convoyana But : luggage



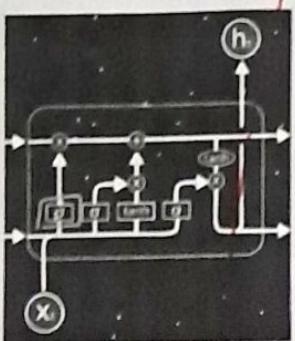
Iske upar luggage r khe
hai cheating ke liye

Jo cheez kaam ki
hai usse RNN
wara tha do.



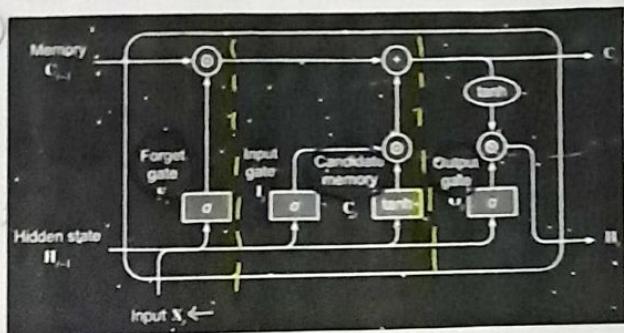


LSTM Architecture



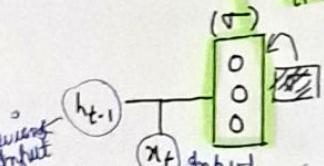
Basic Architecture

LSTM Have Two Loop

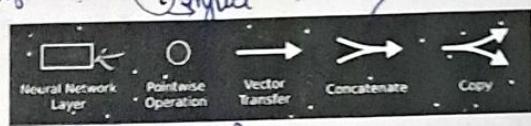


3 Gate

{Neural New Layer}



combine two vector



Notation

$$v_1 = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \xrightarrow{\text{tanh}}$$

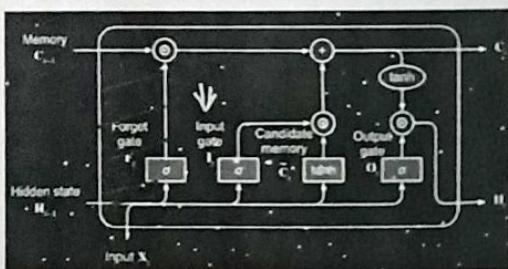
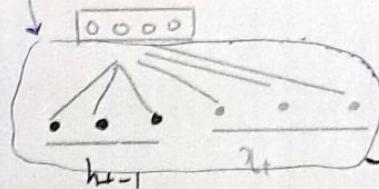
$$\times = [4 \ 10 \ 18]$$

$$+ = [5 \ 7 \ 9]$$

$$\text{tanh} = \begin{bmatrix} \text{tanh}(1) & \text{tanh}(2) \\ \text{tanh}(3) \end{bmatrix}$$

Combining 2 vectors

$$h_{t-1} = [1 \ 2 \ 3]$$



Forget Gate

Text

Next Word

$x_1 \ x_2 \ x_3 \ x_4 \ y_5$

$$\sigma [w_f \cdot [h_{t-1}, x_t] + b_f]$$

$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

h_{t-1} = Hidden state of previous time stamp

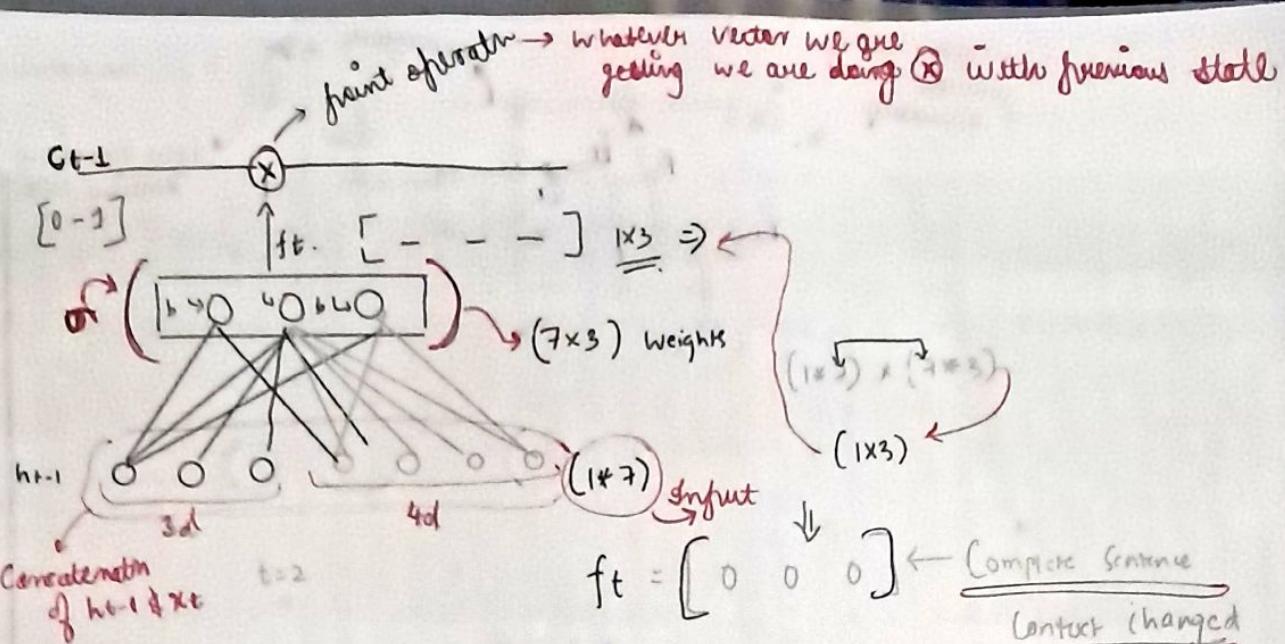
x_t = Word passed as i/p in the current time stamp

$$[0 \ 2 \ 4 \ 1] \quad [4 \ 5 \ 1 \ 2] \quad \dots$$

$$x_t \quad x_{t+1}$$

$$h_{t-1} \quad [1 \ 2 \ 4] \quad [3d]$$

$$c_{t-1} \quad [4 \ 2 \ 1] \quad [3d]$$



$$\textcircled{1} \quad c_{t-1} = [6 \ 8 \ 9] \otimes [0 \ 0 \ 0]$$

previous state

$$= [0 \ 0 \ 0] \leftarrow \text{Removing all the previous context}$$

$$\textcircled{2} \quad c_{t-1} = [6 \ 8 \ 9] \otimes [1 \ 1 \ 1] = [6 \ 8 \ 9]$$

$$\textcircled{3} \quad c_{t-1} = \begin{bmatrix} 6 \\ 8 \\ 9 \end{bmatrix} \otimes \begin{bmatrix} 0.5 & 1 & 0.5 \end{bmatrix} = \begin{bmatrix} 3 \\ 8 \\ 4.5 \end{bmatrix}$$

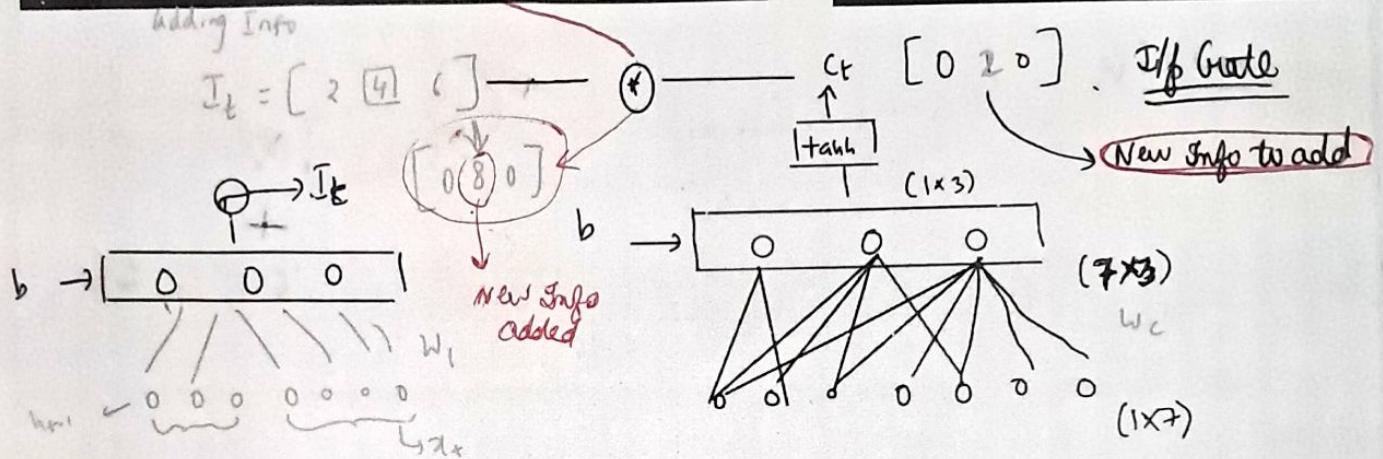
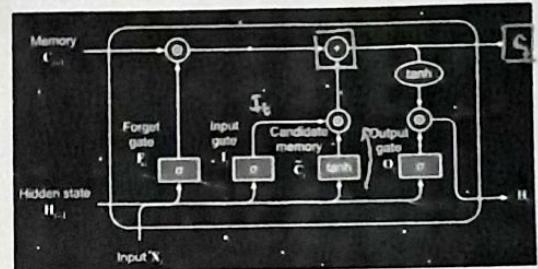
Conclusion \therefore Based on the context \rightarrow Forget gate will let go some information or will not let go some info { Forgetting }.

② Input Gate And Candidate Memory

64

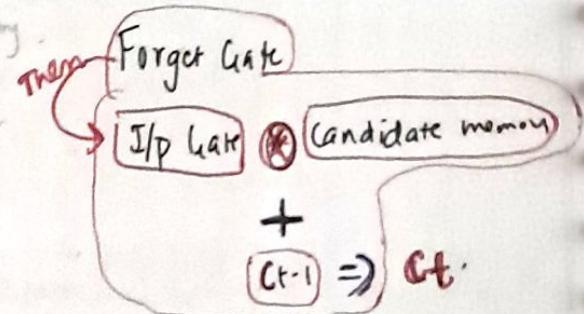
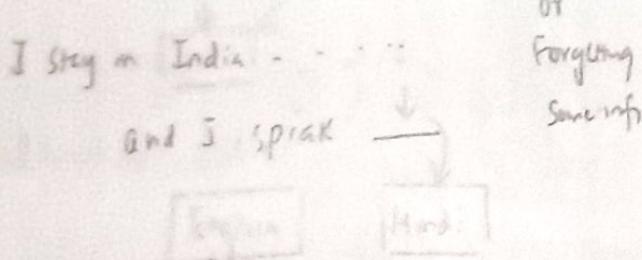
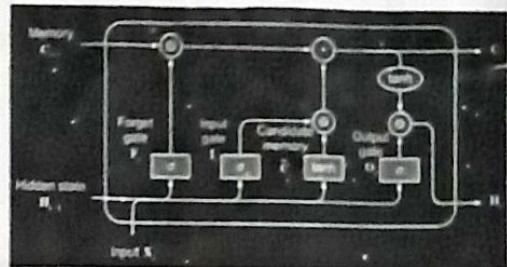
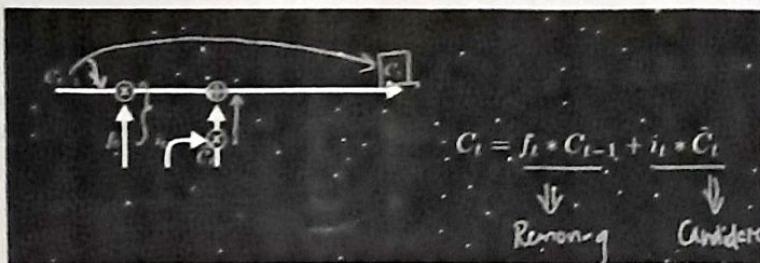
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

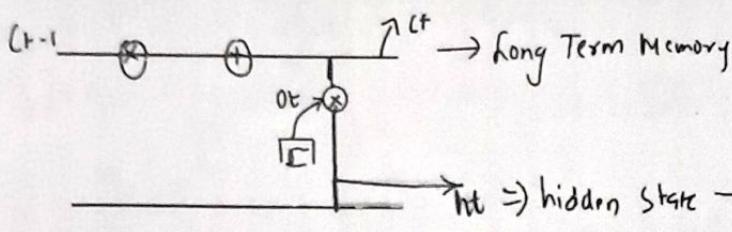
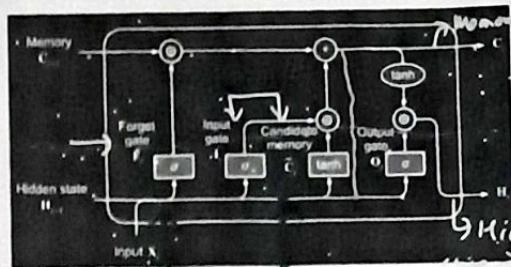
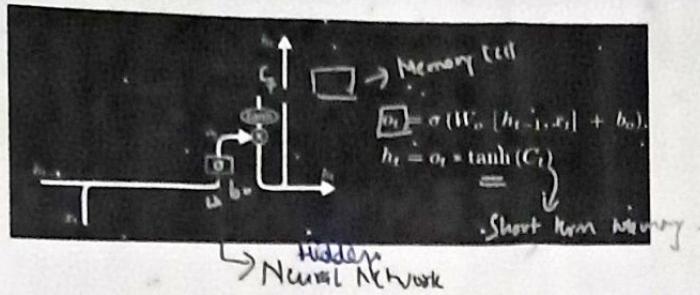


Context = If any information needed to be added in the memory $c_{t-1} \rightarrow$ The information will be added

④ Combine forget & Input Gate

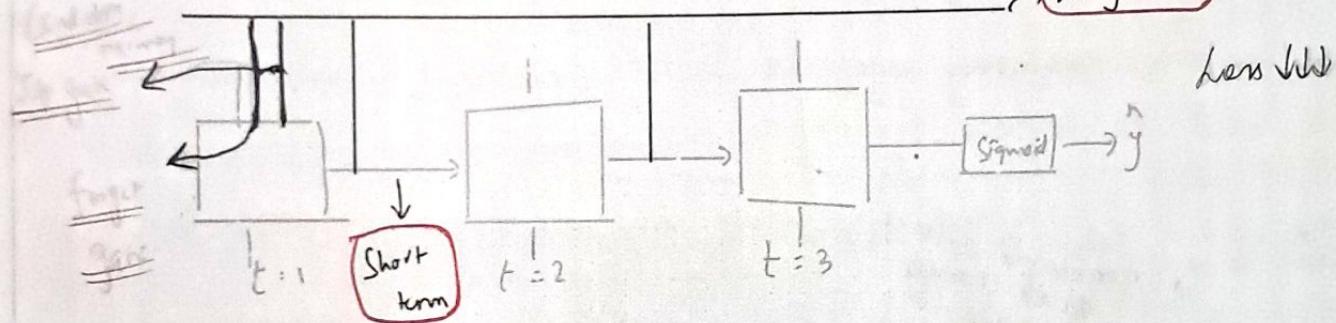


Output gate LSTM RNN



Forget Add Info \rightarrow Context
Some Info \rightarrow Context \rightarrow Memory Cell

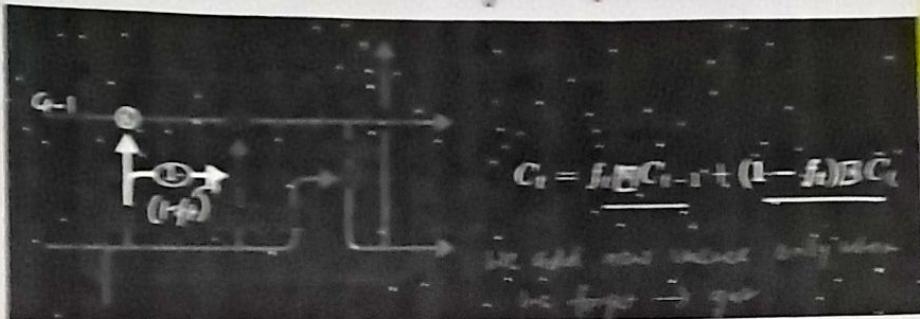
Long Term



$[W_i, W_c, W_o] \rightarrow$ Updating \leftarrow Back Propagation

GRU RNN \Rightarrow LSTM Variant

④ Another variation \rightarrow Coupling forget and zip gate



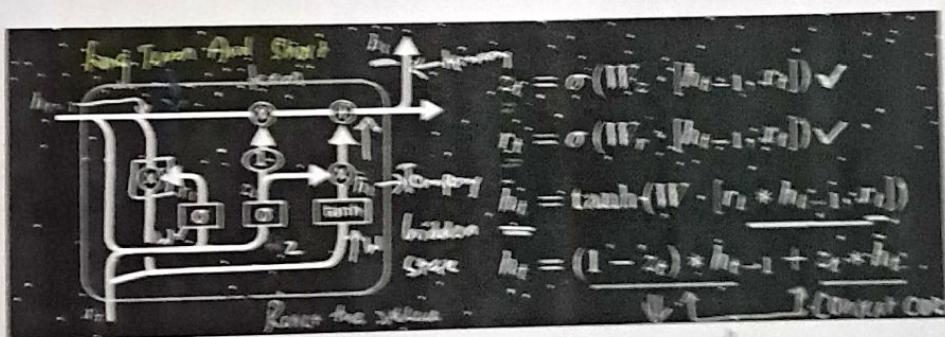
Goal: we only forget when we are going to add something in its place we only add new values to the state when we forget something else.

[Instead of separately deciding what to forget and what we should add now Info, we make the decision together.]

Another variant

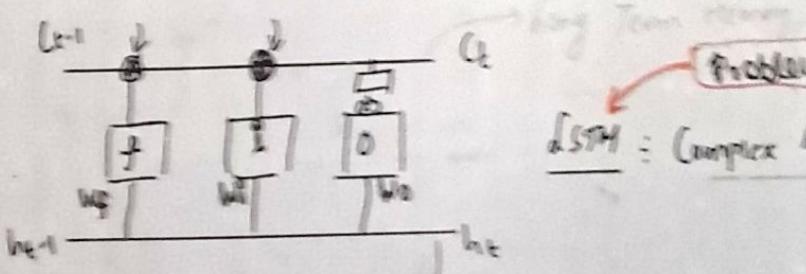
④ GRU \rightarrow Gated Recurrent Unit [Cho, et al [2014]]

1990 \rightarrow LSTM
2000 - Variants
2014 \rightarrow GRU



$z_t \Rightarrow$ Update Gate
 $r_t \Rightarrow$ Reset Gate
 $\tilde{h}_t \Rightarrow$ Temporary hidden state

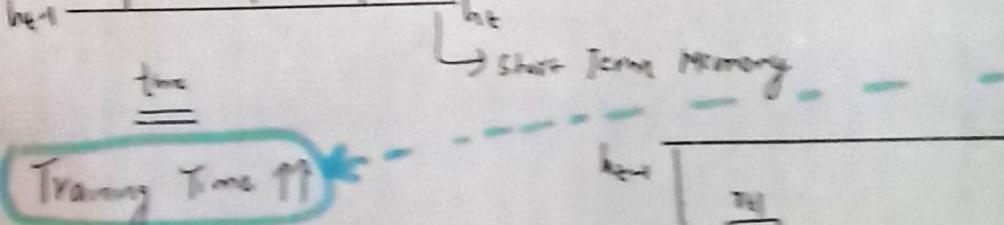
bias
[w_f, w_r, w_o]



LSTM: Complex Architecture

problem + candidate memory

O/P



Reset gate = $r_t =$

\rightarrow Resetting some info from $h_{t-1} \Rightarrow$ memory \rightarrow LSTM + STM

$$h_{t-1} = [0.6 \quad 0.5 \quad 0.3 \quad 0.9]$$

\otimes

$$x_t = []$$

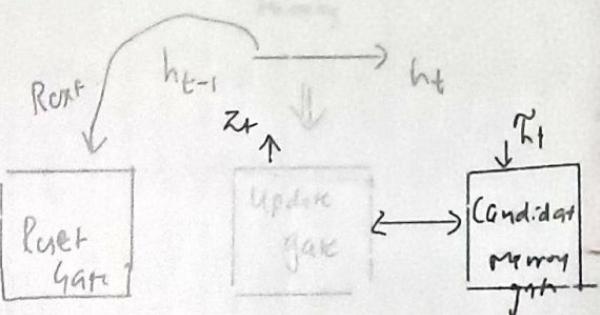
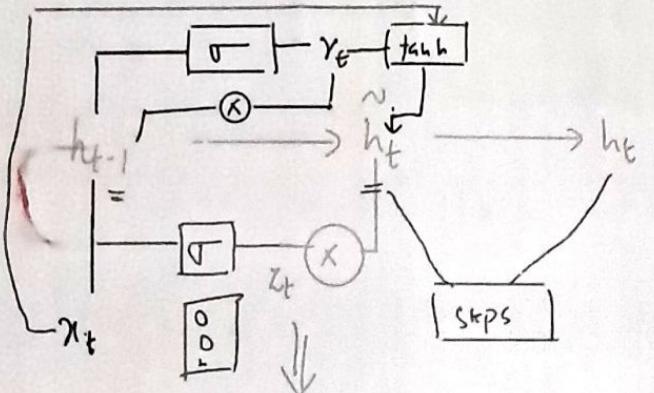
\downarrow

\downarrow

\downarrow

\downarrow

$$x_t \rightarrow [0.12 \quad 0.20 \quad 0.24 \quad 0.18] \leftarrow \text{Rectifying} \rightarrow \text{Context}$$



What Context Info needs to be Added

\Downarrow

Candidate hidden state. [current context]

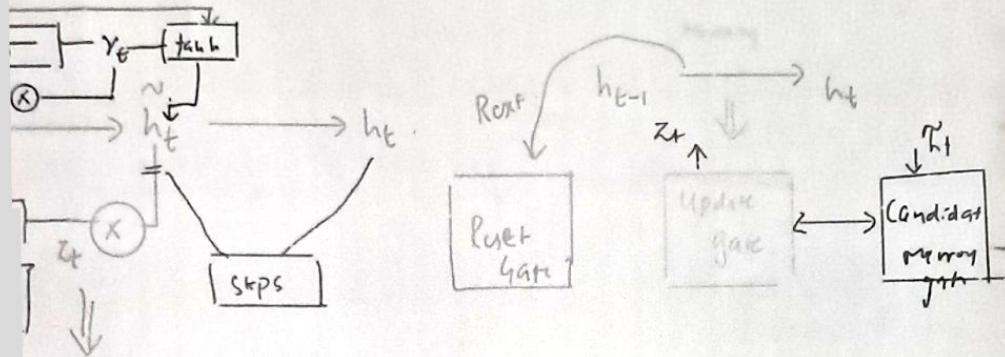
\Downarrow
Imp \rightarrow Add Info

$h_t \Rightarrow$ New Info

$$0.5 \quad 0.3 \quad 0.9]$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$0.20 \quad 0.24 \quad 0.18] \leftarrow \text{Reaching} \rightarrow \text{Context}$$



What Context Info needs to be Added



Candidate hidden state [Current Context]

\downarrow
Imp \rightarrow Add Info

$h_t \Rightarrow$ New Info

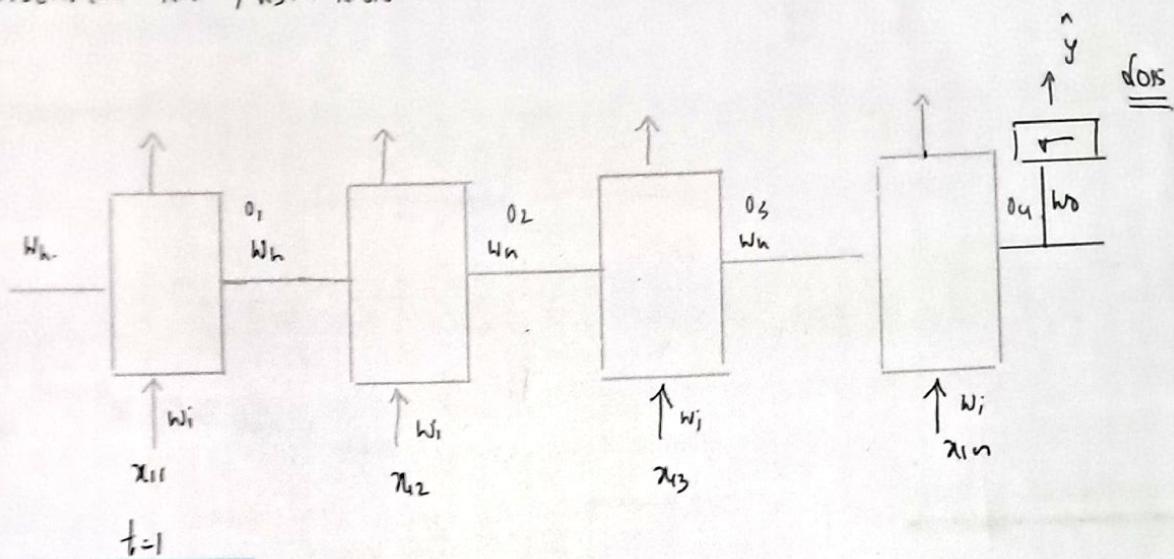
\equiv

Bidirectional RNN

① Simple RNN \rightarrow practical Implement \rightarrow Embedding layers

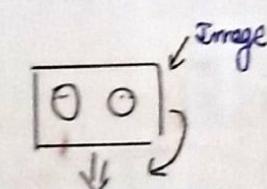
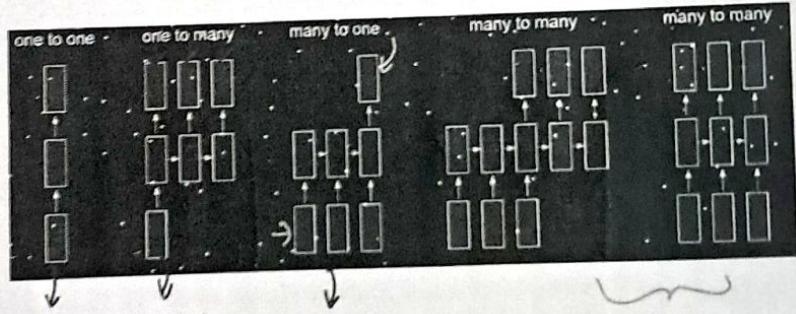
② LSTM, GRU variants RNN \rightarrow Practical Implement

③ Bidirectional RNN / LSTM RNN



Types Of RNN

- ① One to Many RNN
- ② Many to One RNN
- ③ Many to Many RNN
- ④ One to One RNN



Eg:

Eg: Image \rightarrow Captioning

There is a dog and cat

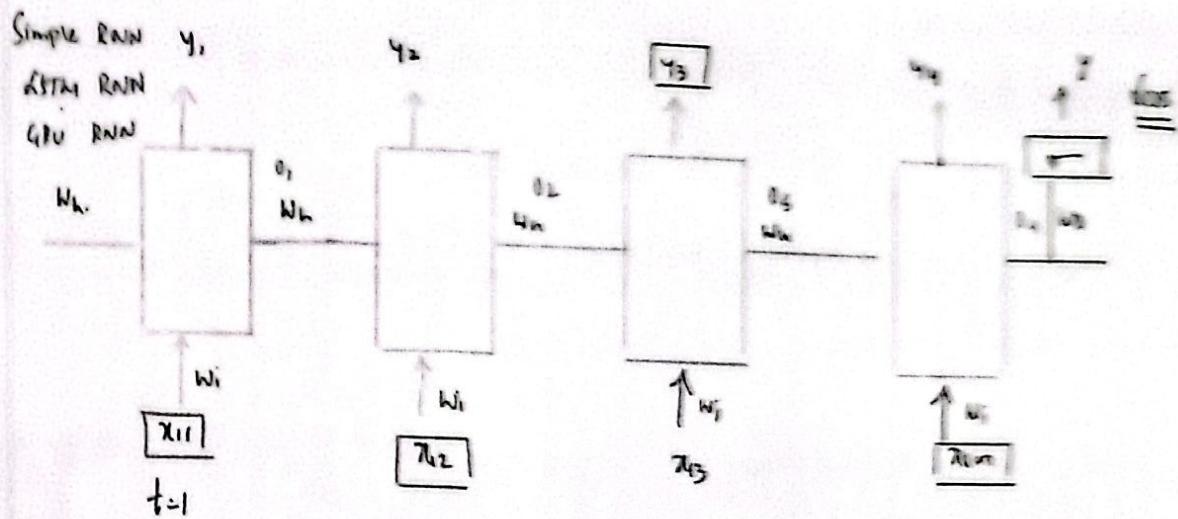
A cat eating food \rightarrow o eating

O/p \rightarrow picture

Many I/P
Many O/P

Text input on google

Language Translation



Text Example = KRISH

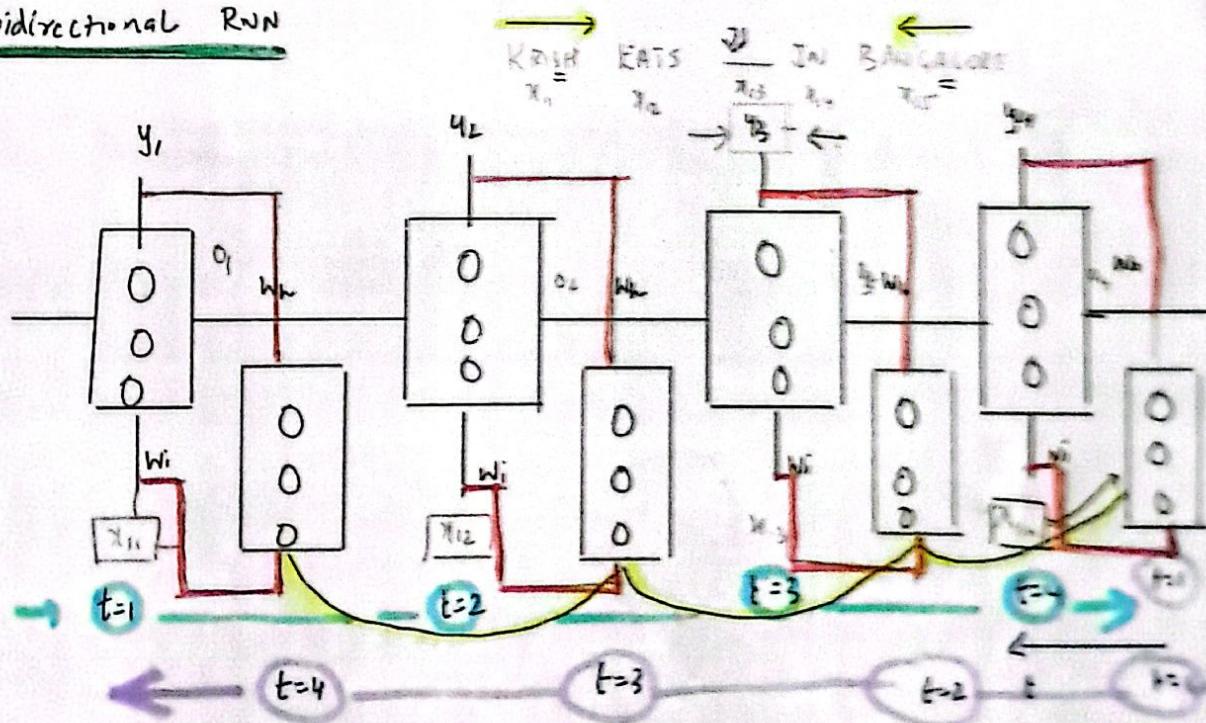
DOG IN BENGAL=CAT

KRISH KATS

PIZZA PAZIS

Defend and Scream

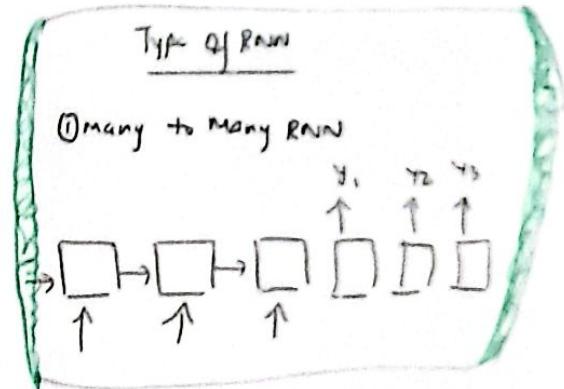
Bidirectional RNN



① Forward propagation \rightarrow Equation \leftarrow

Encoder And Decoder

- ① Simple RNN \rightarrow Vanishing Gradient Problem
- ② LSTM RNN \rightarrow Long Short Term Memory.
- ③ GRU RNN
- ④ Bidirectional RNN \leftarrow If we have context with further words



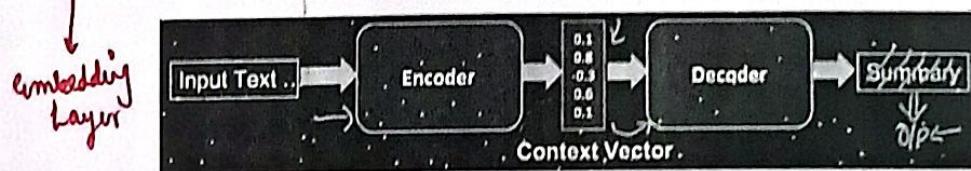
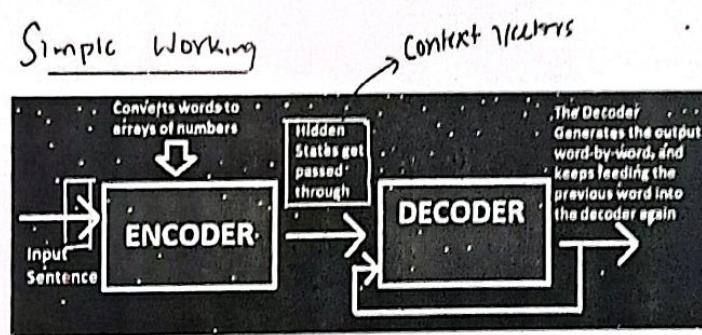
Encoder And Decoder?

Eg: One language To other

English \rightarrow French

Eg: Linked chat \rightarrow Hi, how are you?

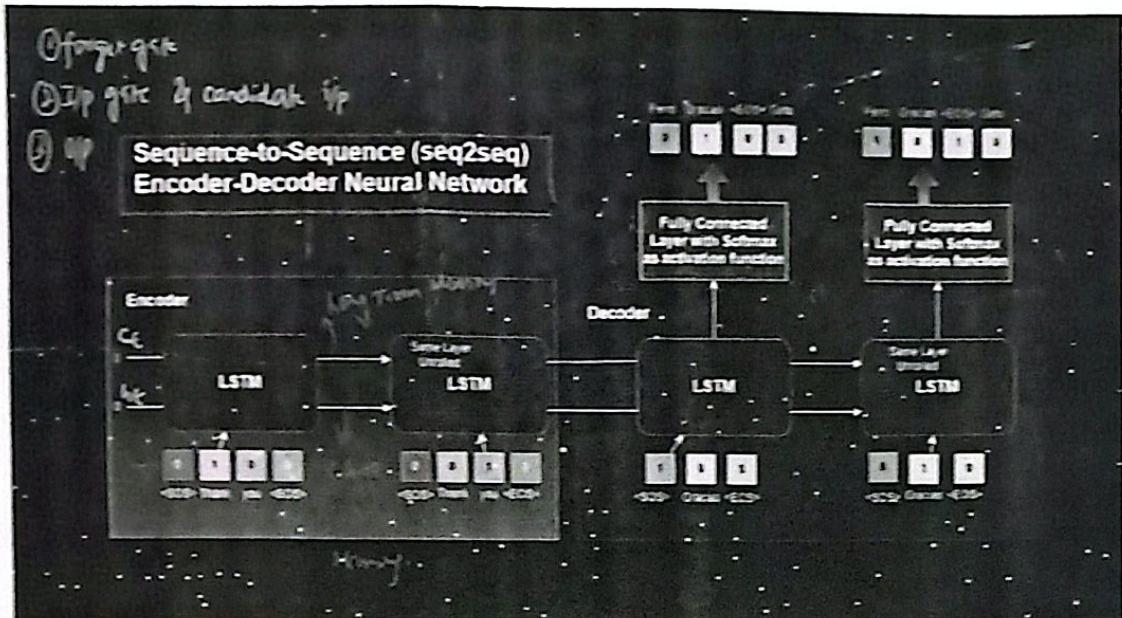
Simple Working



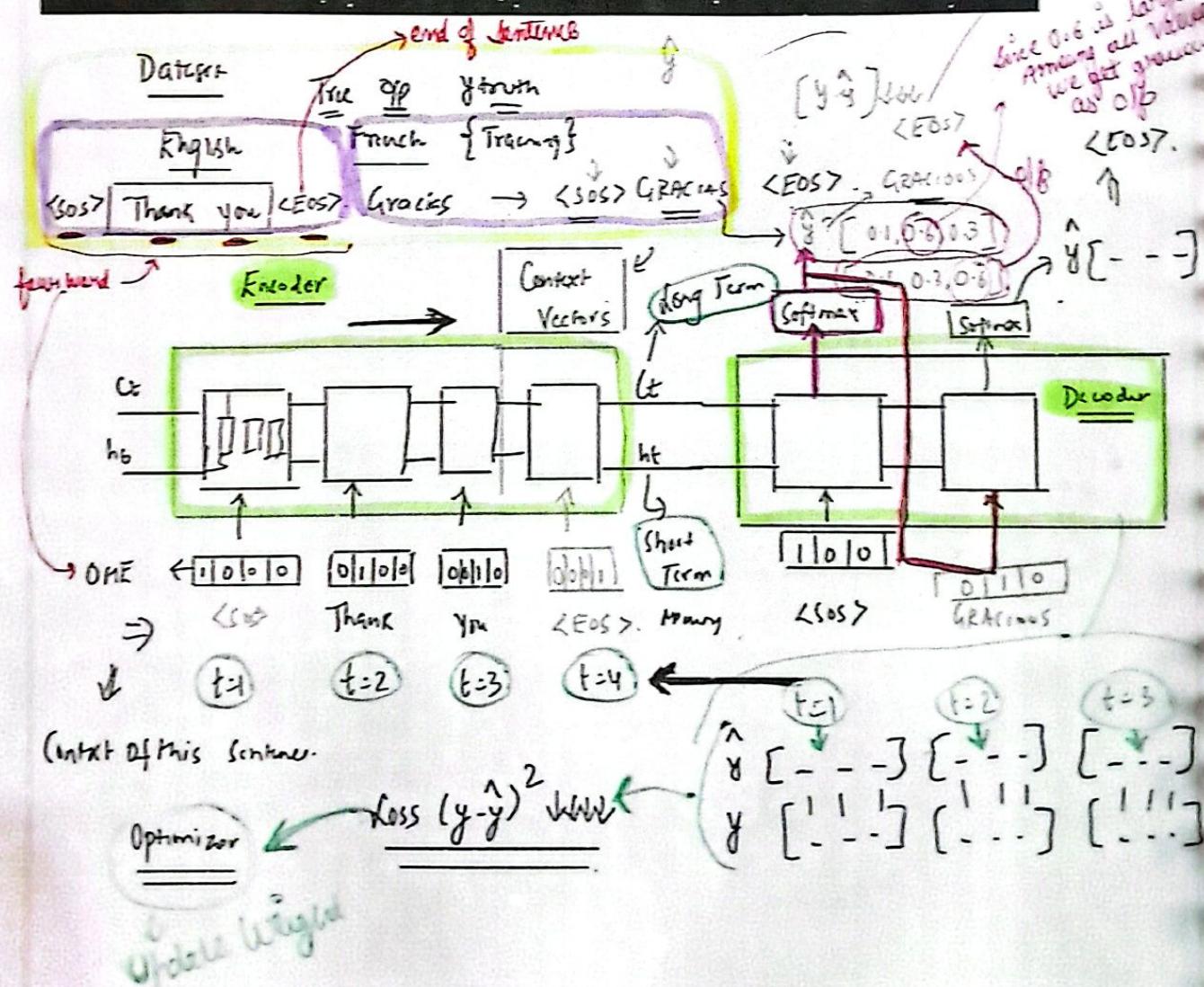
- Convert
- ① Encoder \Rightarrow I/p \Rightarrow Context Vector
 - ② Decoder \Rightarrow O/P

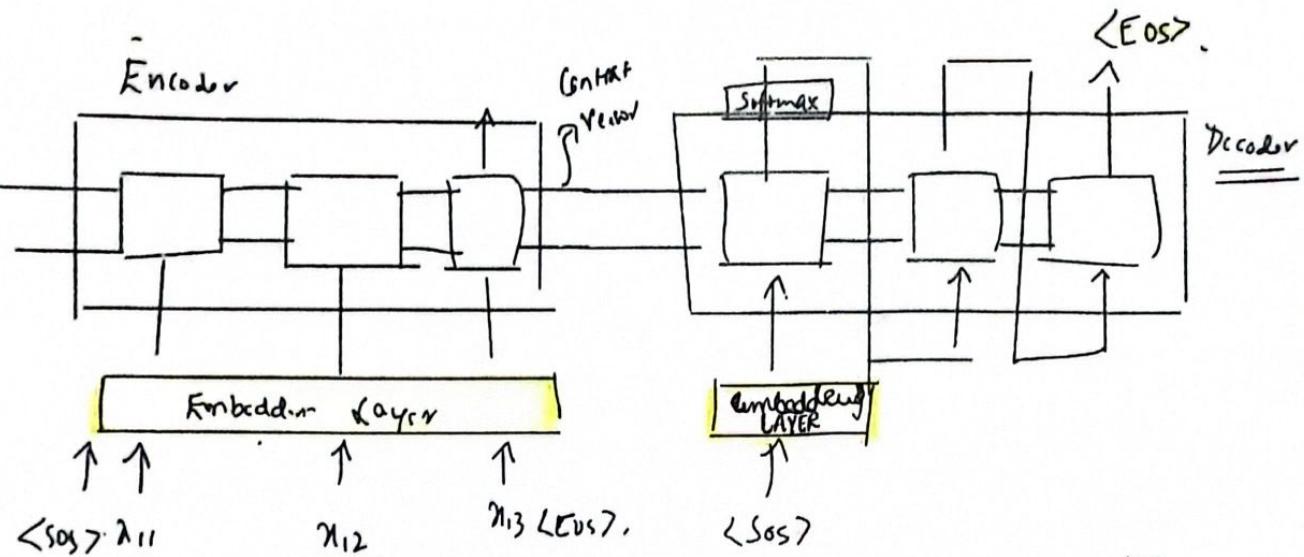
- Usecase
- ① Language Translation
 - ② Text Generation
 - ③ Text Suggestion

Row \rightarrow converting constraint problem



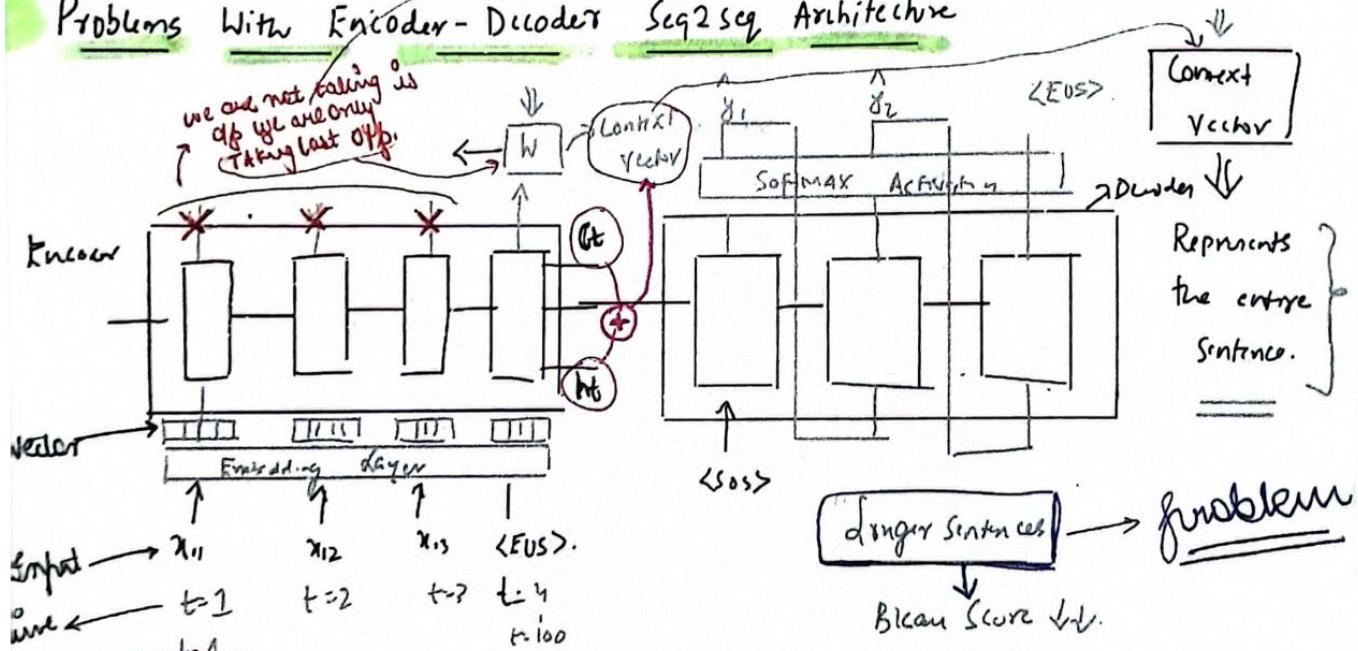
we will not go
further with



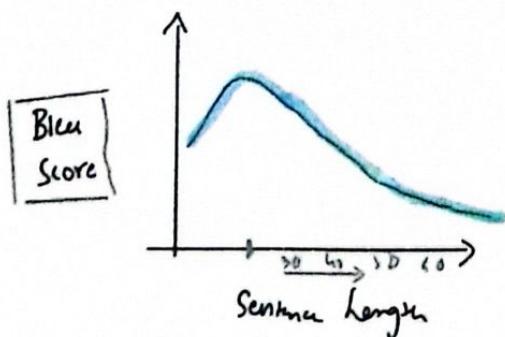


Sentence will have more info about last word's

Problems With Encoder-Decoder Seq2Seq Architecture



Researchers: Sentences of varying length

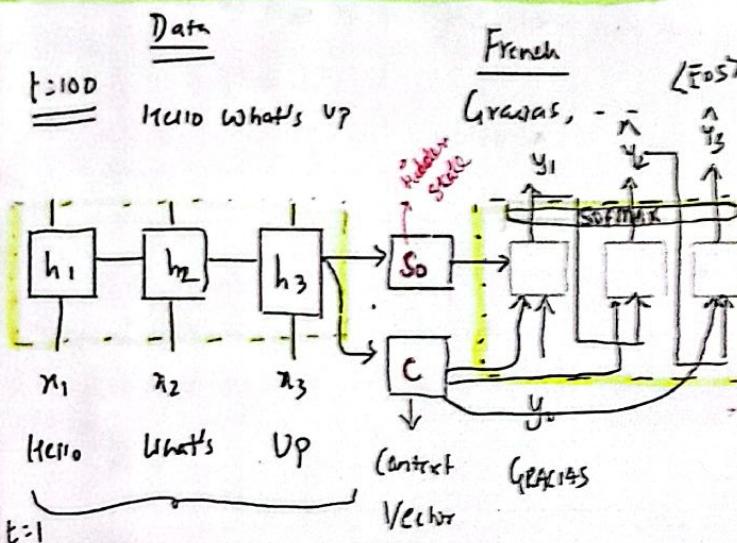


⇒ Seq to Seq Dots

⑧ Attention Mechanism \rightarrow Seq2Seq Network

longer paragraph \rightarrow {Context Vector} + {Context}

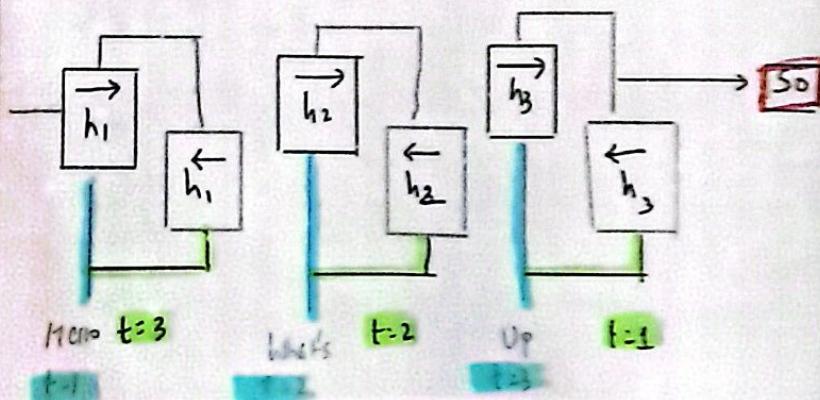
Attention Mechanism | Seq2Seq Networks



Encoder Decoder Architecture

Attention Mechanism

<https://erdem.pl/2021/05/introduction-to-attention-mechanism>



3. LEARNING TO PREDICT AND TRANSLATE

In this section we propose a novel architecture for neural machine translation that consists of a bidirectional RNN as an encoder (Sec. 3.2) and a decoder (Sec. 3.3) that translates through a sequence-to-sequence using decoding autoencoders (Sec. 3.4).

3.1. Decoding Given a Document

This new model architecture will learn to predict words in a sequence given a document. The input sequence is $x = (x_1, \dots, x_n)$ and the output sequence is $y = (y_1, \dots, y_n)$.

where x is an RNN hidden state sequence computed by

$$x_t = f(x_{t-1}, h_t, a_t)$$

It should be noted that unlike the ordinary encoder-decoder approach (Eq. (1)), here the probability is conditioned on the context vector c for each target word y_t .

The context vector c depends on a sequence of annotations h_1, \dots, h_n to which an encoder maps the input sentence. The annotation h_t contains information about the words in the input sentence x with a strong focus on the part preceding the t -th word of the input sequence. We will go in detail how the attention mechanism is computed in the next section.

The context vector c_t is often computed as a weighted sum of the annotations h_j :

$$c_t = \sum_{j=1}^n \alpha_{tj} h_j$$

The weight α_{tj} of each annotation h_j is computed by

$$\alpha_{tj} = \frac{\exp(z_{tj})}{\sum_{i=1}^n \exp(z_{ti})}$$

where

$$z_{ti} = \phi(h_i, a_t, h_t)$$

