# RNN Types

① One-to-One / Vanilla Mode: Vanilla mode of processing w/o RNN

fixed size i/p to fixed sized o/p.

eg: Image Classification



5

MNIST
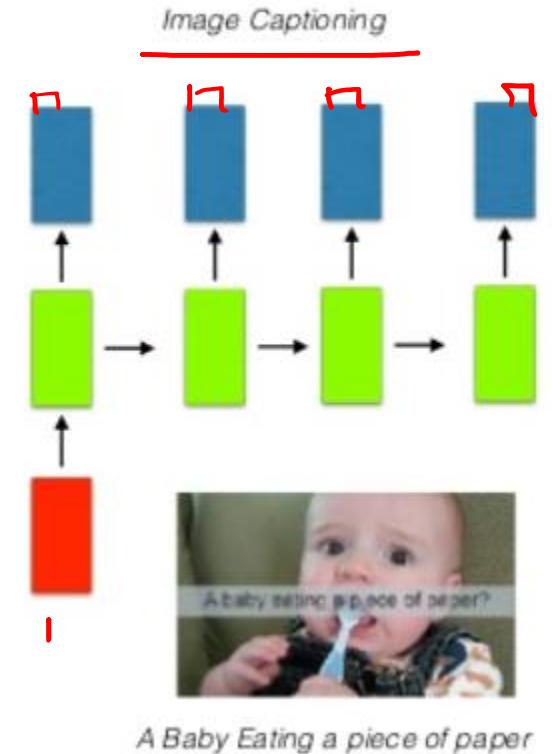
# RNN Types

② One — to — Many!

Sequence O/p

Image Captioning!

Input: An image

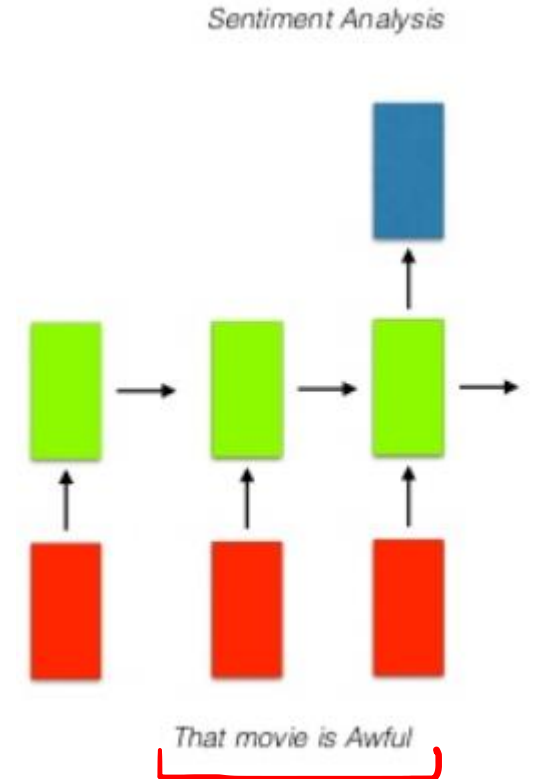O/p: Sentence of words



Image Captioning

A Baby Eating a piece of paper

# RNN Types

③ Many -to-one:

Sequence input

eg: Sentiment analysis where a given sentence is
classified as expressing +ve or -ve sentiment.

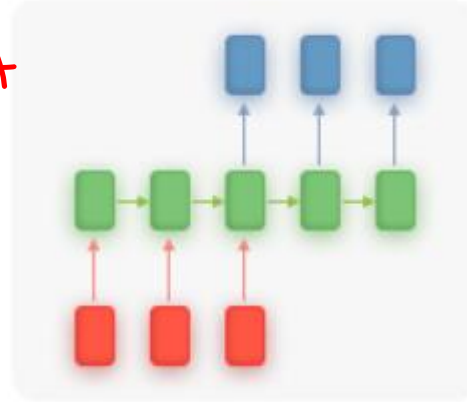Sentiment Analysis



That movie is Awful

# RNN Types

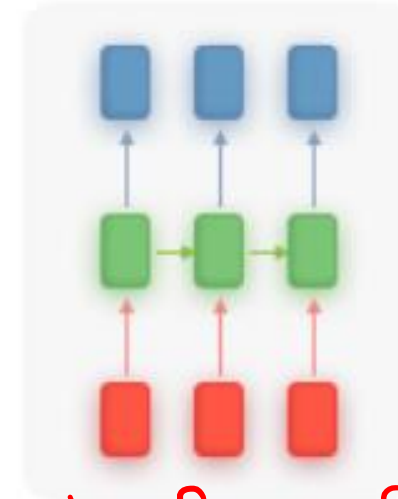**① Many to Many:-**

Sequence input & Sequence Input

eg: Google Translate

Machine Translation

Video Classification where we wish to label each frame of the video.
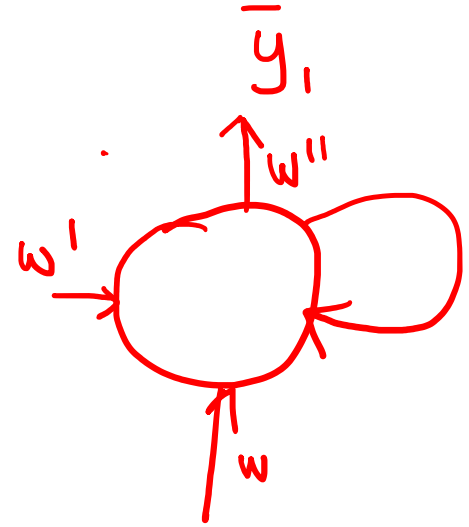
Many to many



Many to many

# Why RNN? Used popularly in NLP problems.

**NLP:** Input Data

X = text format ⟩ Preprocessing
Processed X = Number format ⟨ + Bag of Words (CountVectorizer)

Naive Bayes Algo.



① Bag of Words ⟩ Sequence info is lost.
② TF-IDF ⟩
③ WORD2VEC ⟩

|   | love | playing | cricket | fun |
|---|------|---------|---------|-----|
| 0 | 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 | 0 |
| 2 | 0 | 1 | 0 | 1 |
| Count: | 2 | 3 | 1 | 1 |

0   love play
1   love play cricket
2   play fun

Descending order of Count

$0 = [\underset{x_0}{1}, \underset{x_1}{1}, \underset{x_2}{0}, \underset{x_3}{0}]$

$1 = [1, 1, 1, 0]$

$2 = [\underset{x_0}{1}, \underset{x_1}{0}, \underset{x_2}{0}, \underset{x_3}{1}]$

→ Sequence info is LOST.

|   | play | love | cricket | fun |
|---|------|------|---------|-----|
| 0 | 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 | 0 |
| 2 | 1 | 0 | 0 | 1 |

# Why RNN?

Time Series Prob

up/down/same

exponential smoothing

eg: Image Captioning,
Google Translate
Sentiment analysis for restaurant reviews
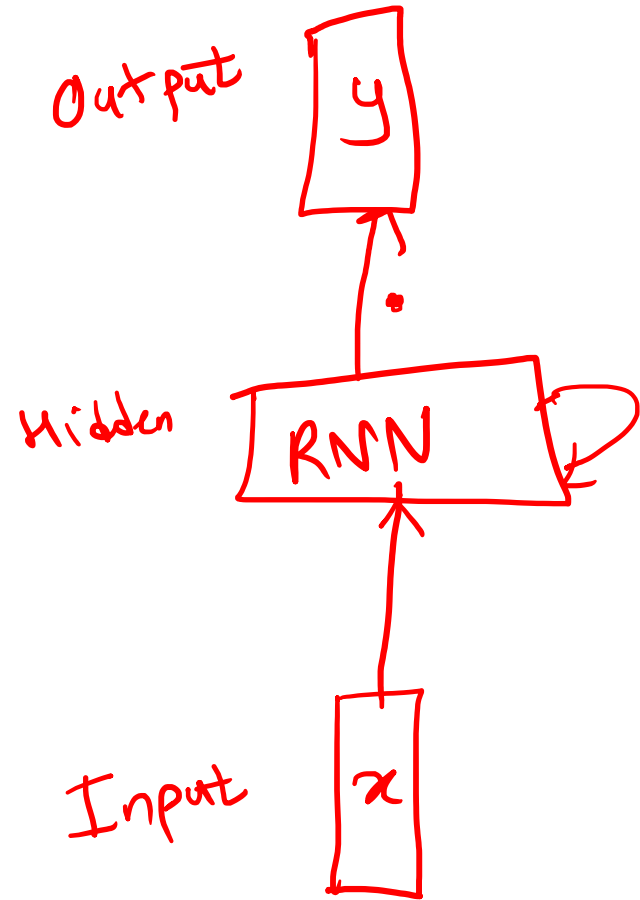Amazon alexa
Google Assistant, Siri/ Microsoft Cortana
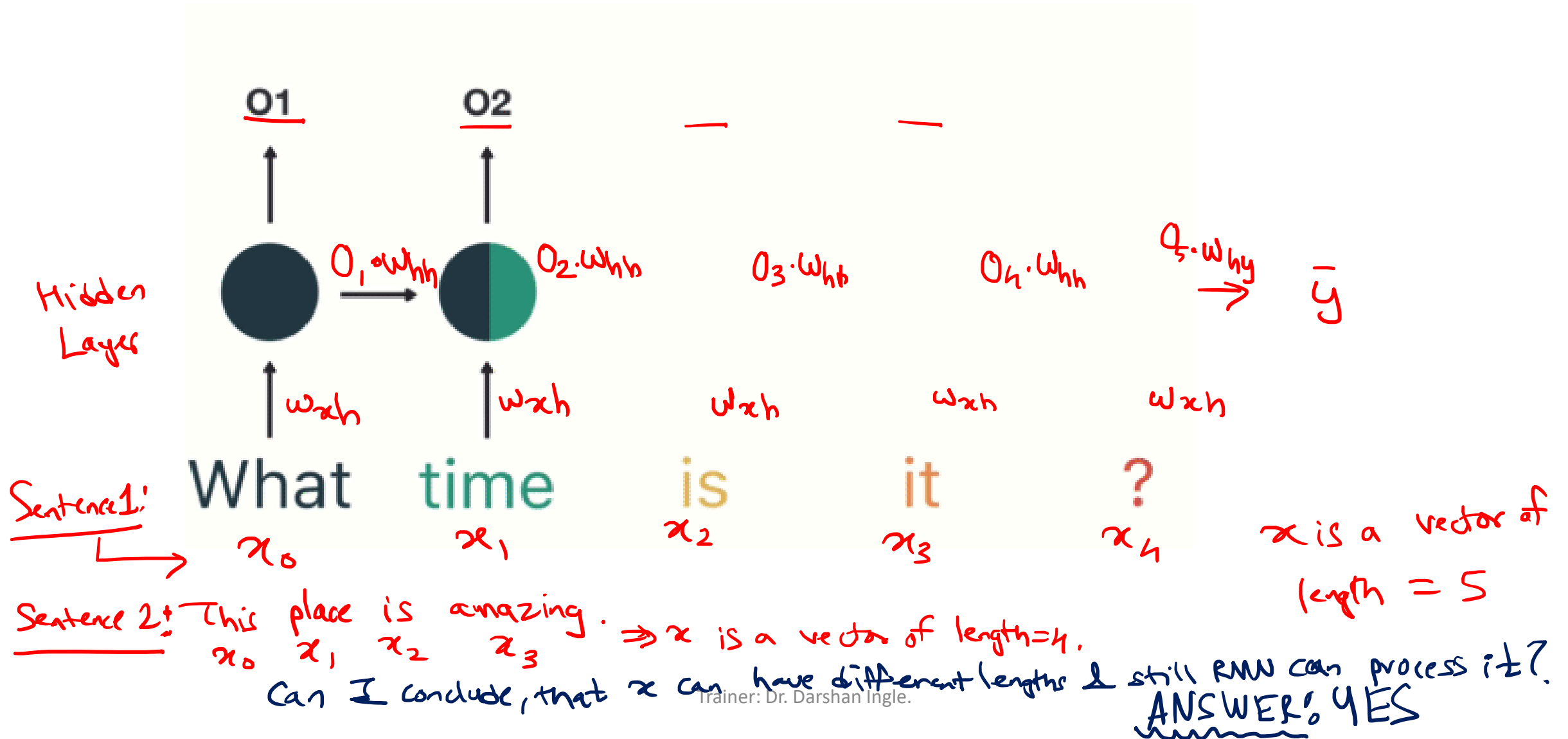
Imp: Sequence info.

# Summary

Simple form of Vanilla RNN

Output

$y$

Hidden

RNN

Input

$x$

$$h_t = f_w\left(h_{t-1}, x_t\right)$$

$$\therefore h_t = \tanh\left(W_{hh} \cdot h_{t-1} + W_{xh} \cdot x_t\right)$$

$$y_t = W_{hy} \cdot h_t$$

# Summary

$O1$     $O2$    —    —

Hidden Layer

$O_1 \cdot W_{hh}$    $O_2 \cdot W_{hb}$    $O_3 \cdot W_{hb}$    $O_h \cdot W_{hh}$    $O_5 \cdot W_{hy}$ $\Rightarrow \bar{y}$

$W_{xh}$    $W_{xh}$    $W_{xh}$    $W_{xh}$    $W_{xh}$

**What**    **time**    **is**    **it**    **?**

Sentence1:   $x_0$    $x_1$    $x_2$    $x_3$    $x_4$    $x$ is a vector of length = 5

Sentence 2: This place is amazing. $\Rightarrow x$ is a vector of length=4.
$x_0$   $x_1$   $x_2$   $x_3$

Can I conclude, that $x$ can have different lengths & still RNN can process it?
ANSWER: YES

# RNN Architecture (Fwd Prop over Time)

F.P. over time

$\omega \approx \omega_{xh}$
$\omega' \approx \omega_{hh}$
$\omega'' \approx \omega_{hy}$

$t=1$   $t=2$   $t=3$   $t=4$

128 neurons

Unroll

Softmax / Sigmoid

$\bar{y}_1$

$O_1$   $\omega'$   $O_2$   $\omega'$   $O_3$   $\omega'$   $O_4$   $\omega''$   A.F.   $\bar{y}_1$

$\bar{y}_1$ $O_t$ $O_t$ A.F. $x_{it}$

$\omega$ $x_{11}$   $\omega$ $x_{12}$   $\omega$ $x_{13}$   $\omega$ $x_{14}$

Sentence:   The   ambience   is   fantastic.
            $x_{11}$   $x_{12}$   $x_{13}$   $x_{14}$
            $x_1$      $x_2$     $x_{13}$   $x_{14}$

$Loss = (\hat{y} - y)$ ↓ reduce

$$O_1 = f[(x_{11} \cdot \omega + b) + (O_0 \cdot \omega' + b')]$$

$$O_2 = f[(x_{12} \cdot \omega + b) + (O_1 \cdot \omega' + b')]$$

$$O_3 = f[(x_{13} \cdot \omega + b) + (O_2 \cdot \omega' + b')]$$

$$O_4 = f[(x_{14} \cdot \omega + b) + (O_3 \cdot \omega' + b')]$$

Trainer: Dr. Darshan Ingle.

# RNN Architecture (Fwd Prop over Time)

# RNN Architecture (Back Prop over Time)



$$Loss = (\hat{y} - y) \downarrow reduce$$
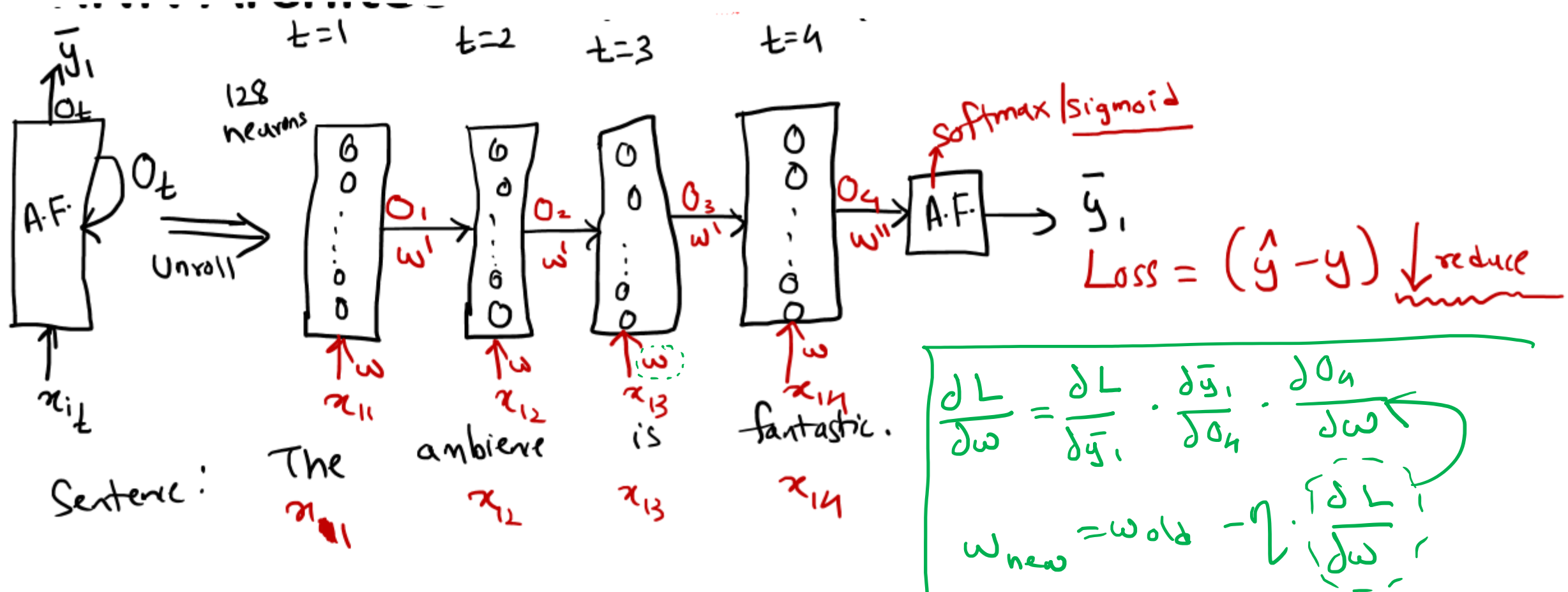
$\frac{\partial L}{\partial \bar{y}_1}$ can be computed directly w/o using chain rule. ∴ SORTED.

$$\frac{\partial L}{\partial w''} = \frac{\partial L}{\partial \bar{y}_1} \cdot \frac{\partial \bar{y}_1}{\partial w''}$$

$$\therefore w''_{new} = w''_{old} - \eta \cdot \frac{\partial L}{\partial w''}$$

# RNN Architecture (Back Prop over Time)

# Problems with RNN

① Vanishing Gradient Problem:

tanh() most common

Der. of tanh → $[0,1]$

A.F. in RNN → $[-1, +1]$

② Exploding Gradient problem: It is due to the higher assigned weights

# Problems with RNN

Due to these problems, we cannot do a lot of B.P.,
∴ we cannot consider older states i.e. we cannot process a
lot of sequence data.

∴ We need a better or Modified RNN Model.

✳ Solution: <u>1997</u> : Long Short Term Memory ( <u>LSTM</u> ). ✳

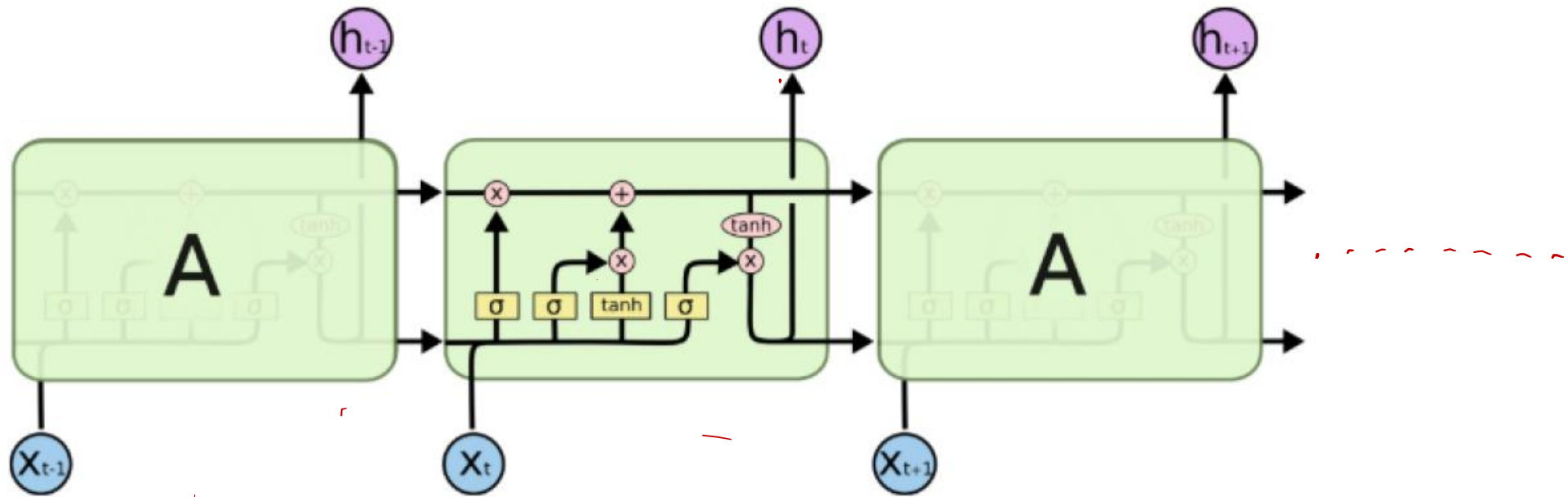eg1: I live in Gujrat & I know Gujrati ⟹ RNN can process it.

eg2: I live in Gujrat & I am a businessman and <u>I know ___?</u>
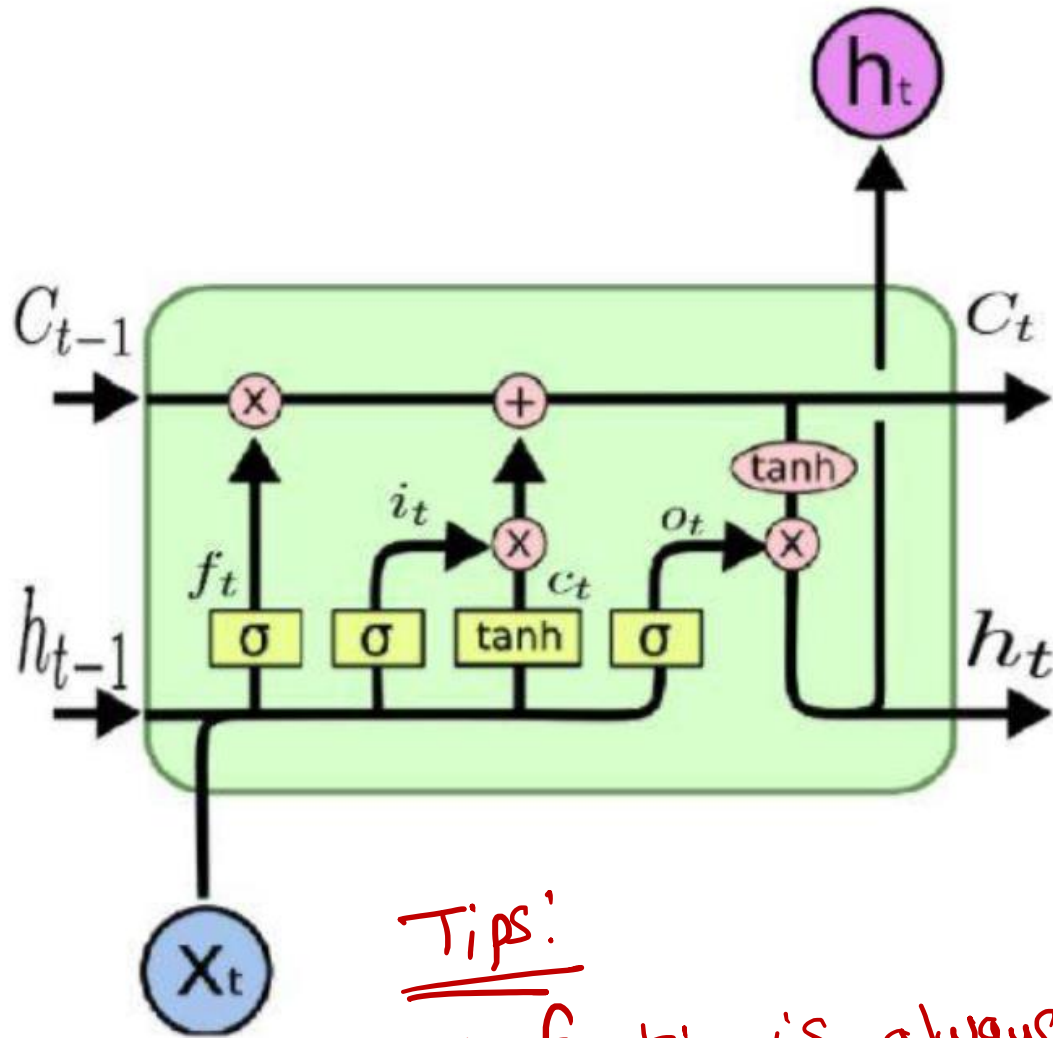⟹ RNN fails
bez the context is far away.

far

# LSTM Architecture

① We have not studied entire NLP course. ② LSTM requires a good background of NLP. ③ LSTM is easy if you compare this to learn car driving.



The repeating module in an LSTM contains four interacting layers.

① f : Forget gate: whether to erase the cell data

② i : Input gate: whether to write to the cell

③ g : Gate gate: How much to write to the cell?

④ o: Output gate: How much to reveal cell.

Trainer: Dr. Darshan Ingle.

C : Cell State
h : hidden State
f : forget gate
i : input gate
o : output gate
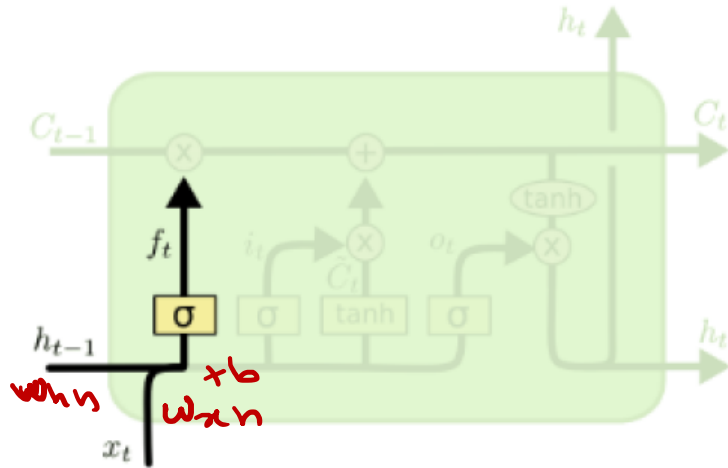
Tips:
① Sigmoid function is always used to filter out information.
② tanh function is used for adding new or modifying existing information.
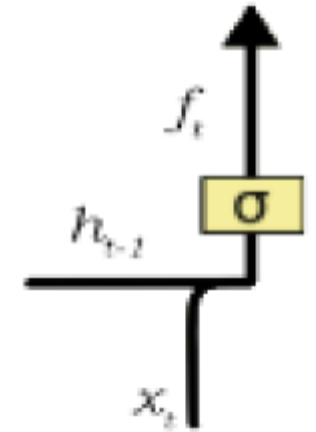
Trainer: Dr. Darshan Ingle.

# Forget Gate

After getting o/p from $h_{t-1}$ i.e. prev. state, Forget gate helps me to take decision about what must be removed from $h_{t-1}$ & thus keep only relevant info.
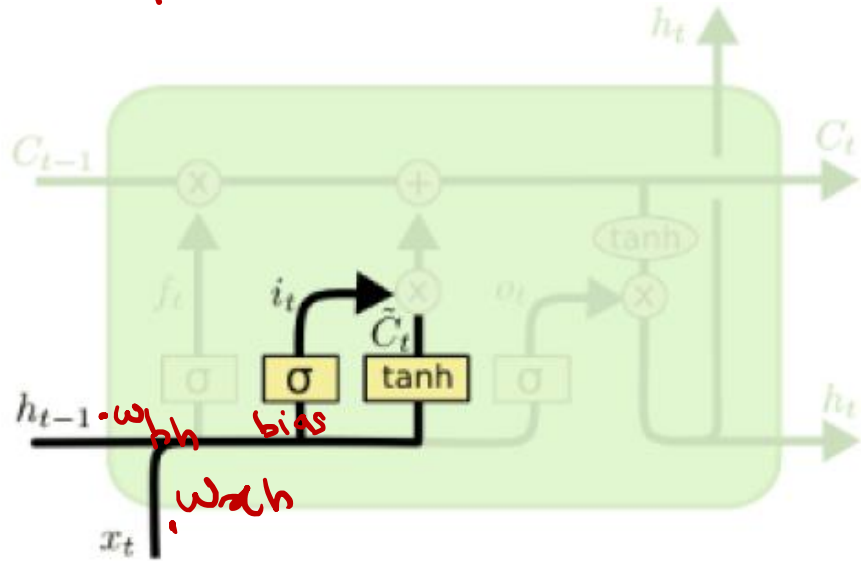


$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] \; + \; b_f \right)$$

# Forget Gate Example

Bob is a nice person. Dan on the other hand is evil.

$h_{t-1}$ $\sigma$ $f_t$

$x_t$

# Input Gate

In input gate, we decide to add new info from present input to our present cell state scaled by how much we wist to add them.



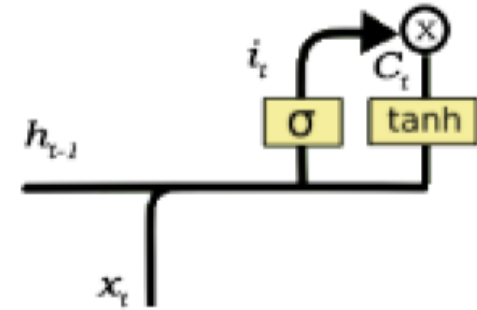$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] + b_i\right)$$

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

Sigmoid layer decides which value shud be updated & tanh create a vector for new info to be added to present cell state.

# Input Gate Example

Bob, knows swimming. He told me over the phone that he
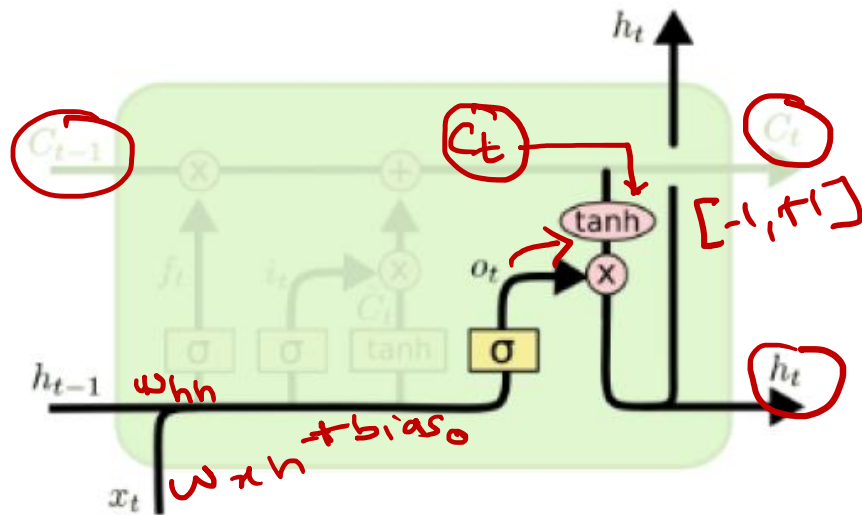had served the navy for 4 long years.

ignored

added to cell state



$$\begin{bmatrix} 0.002 \\ \downarrow \\ 0.1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 5 \\ 6 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.004 \\ 6 \\ 0.3 \end{bmatrix}$$

# Output Gate

we decide finally what to output from our cell state (which will be done by sigmoid).
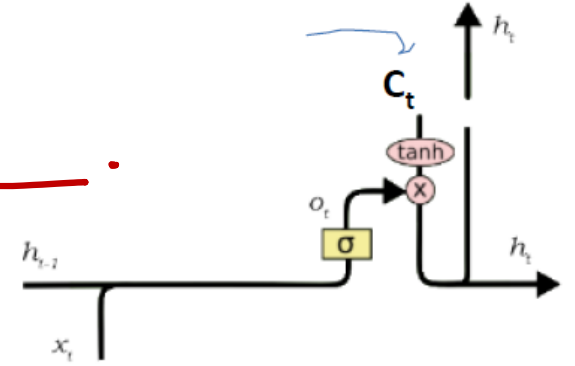


$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

annotations on figure: $C_{t-1}$, $C_t$, $[-1, +1]$, $h_t$, $w_{hh}$, $w \times h + bias_o$

# Output Gate Example

Bob fought single handedly with the enemy and
died for his country. For his contributions brave____?____.
                                                   ~~~~
                                              adjective
                                              i.e. to
                                         describe a noun

# Additional Resource

- https://colah.github.io/posts/2015-08-Understanding-LSTMs/

- https://www.analyticsvidhya.com/blog/2017/12/fundamentals-of-deep-learning-introduction-to-lstm/

- https://towardsdatascience.com/illustrated-guide-to-lstms-and-gru-s-a-step-by-step-explanation-44e9eb85bf21

- https://medium.com/@aidangomez/let-s-do-this-f9b699de31d9