

# RNN and LSTM

**Zahoor Tanoli (PhD)**  
**CUI Attock**

# What's in it for you?

---

- ▶ What is a Neural Network?
- ▶ Popular Neural Networks
- ▶ Why Recurrent Neural Network?
- ▶ What is a Recurrent Neural Network?
- ▶ How does a RNN work?
- ▶ Vanishing and Exploding Gradient Problem
- ▶ Long Short Term Memory (LSTM)
- ▶ Use case implementation of LSTM



# Introduction to RNN

---

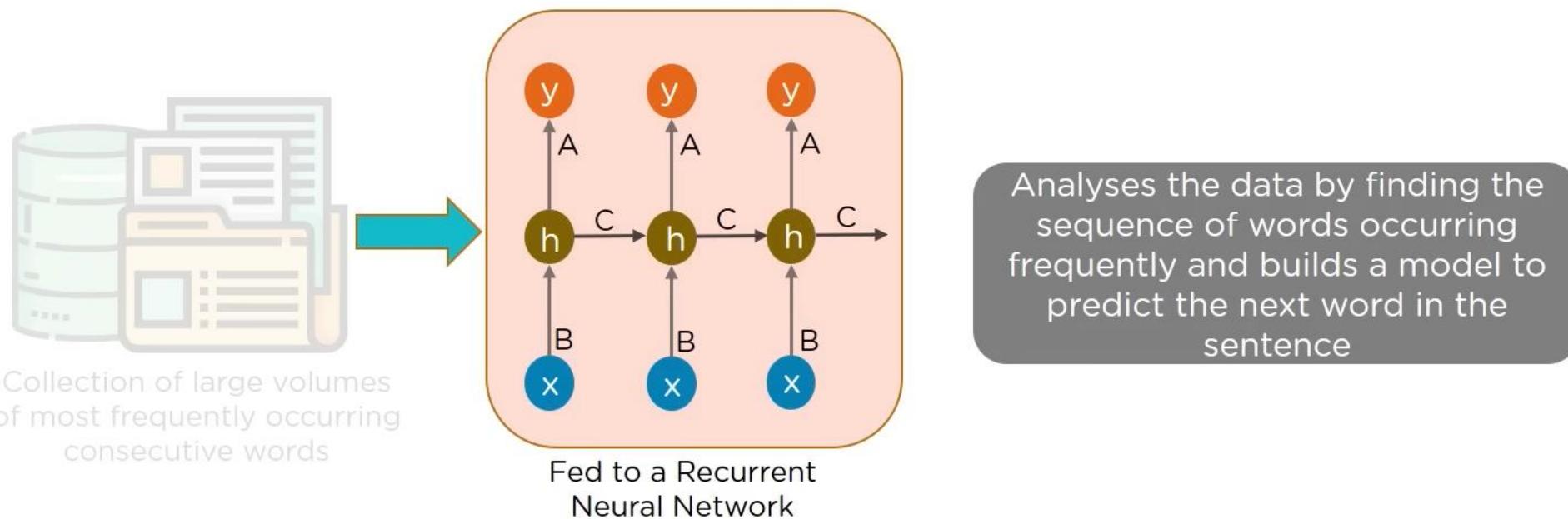
Do you know how Google's autocomplete feature predicts the rest of the words a user is typing?



Collection of large volumes  
of most frequently occurring  
consecutive words

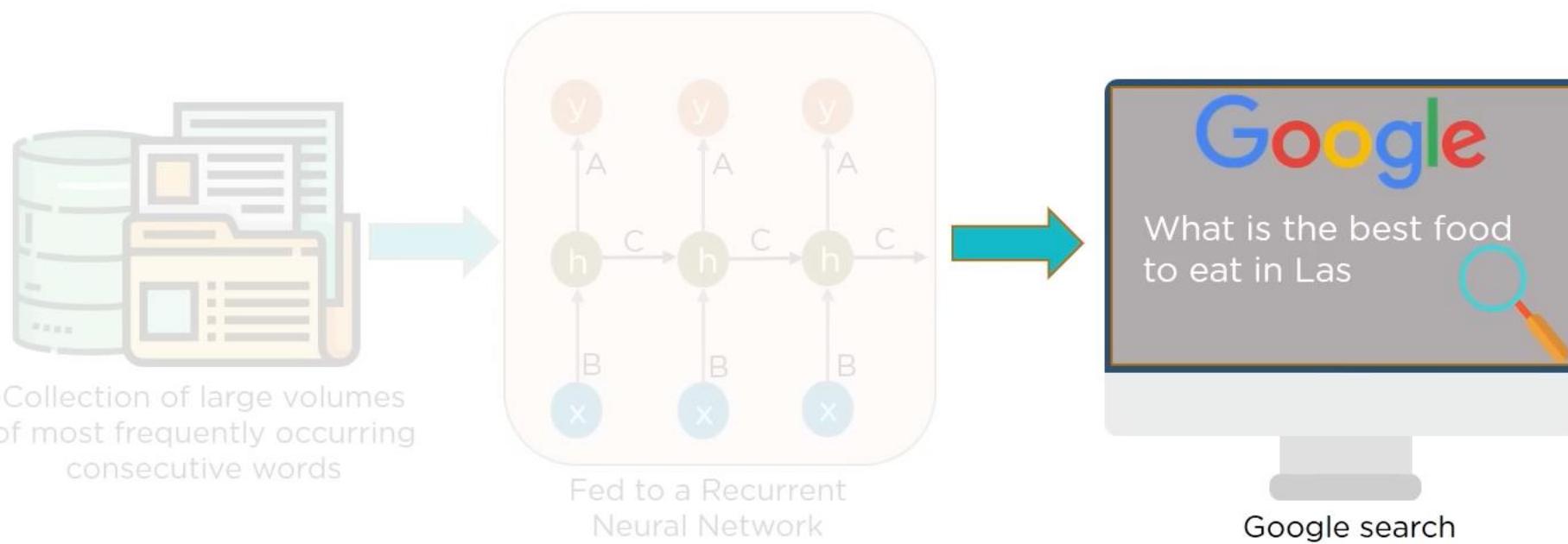
# Introduction to RNN

Do you know how Google's autocomplete feature predicts the rest of the words a user is typing?



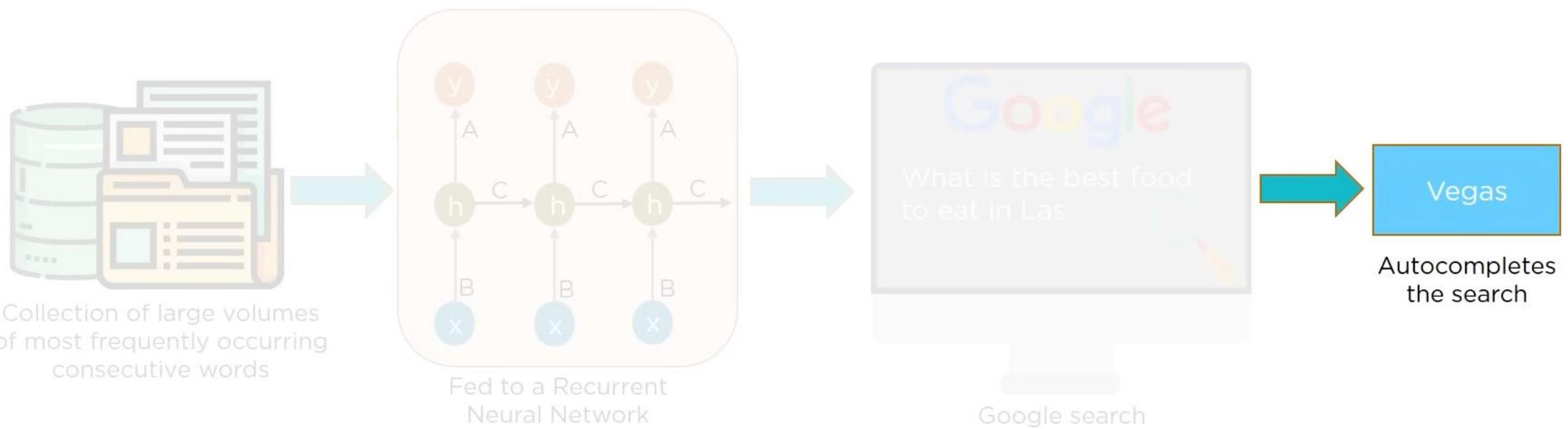
# Introduction to RNN

Do you know how Google's autocomplete feature predicts the rest of the words a user is typing?



# Introduction to RNN

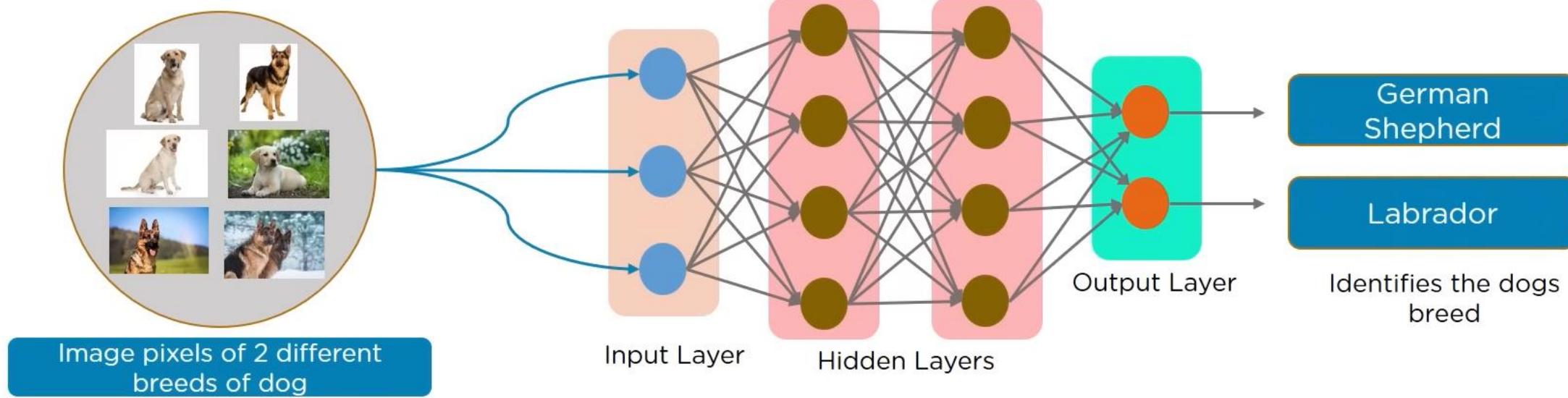
Do you know how Google's autocomplete feature predicts the rest of the words a user is typing?



# What is a Neural Network?



Neural Networks used in Deep Learning, consists of different layers connected to each other and work on the structure and functions of a human brain. It learns from huge volumes of data and uses complex algorithms to train a neural net.

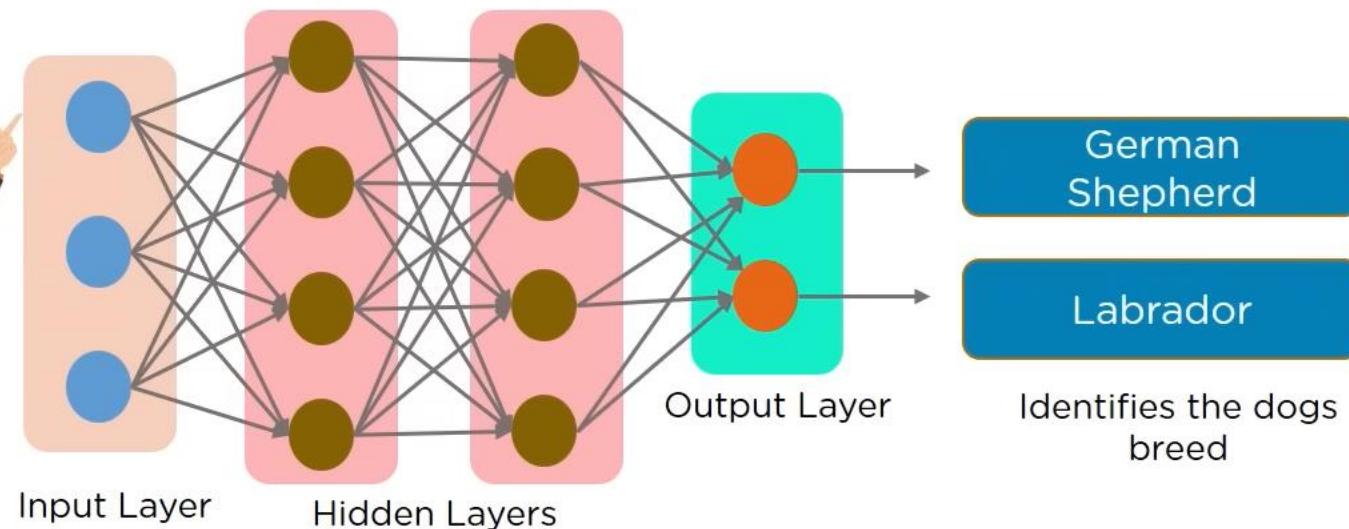


# What is a Neural Network?



Neural Networks used in Deep Learning, consists of different layers connected to each other and work on the structure and functions of a human brain. It learns from huge volumes of data and uses complex algorithms to train a neural net.

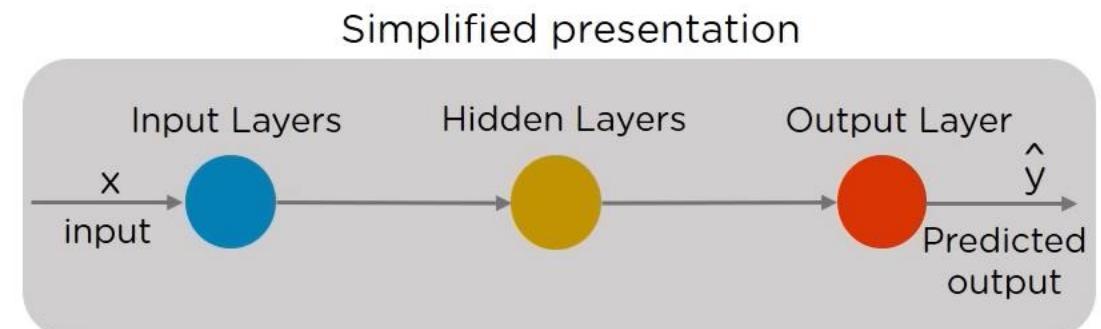
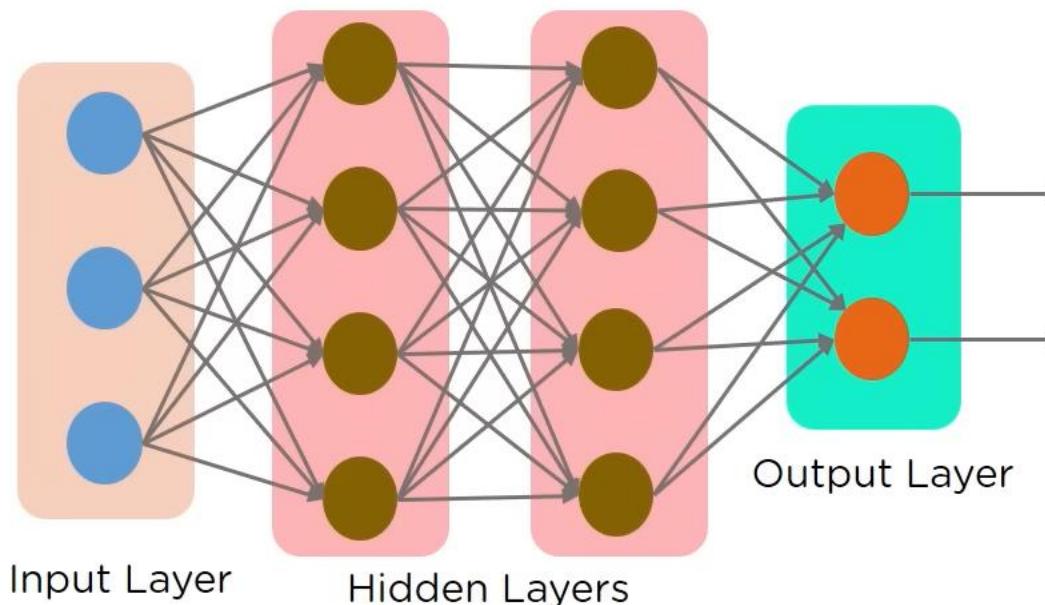
Such networks do not require memorizing the past output



Identifies the dogs breed

# Feed Forward Neural Network

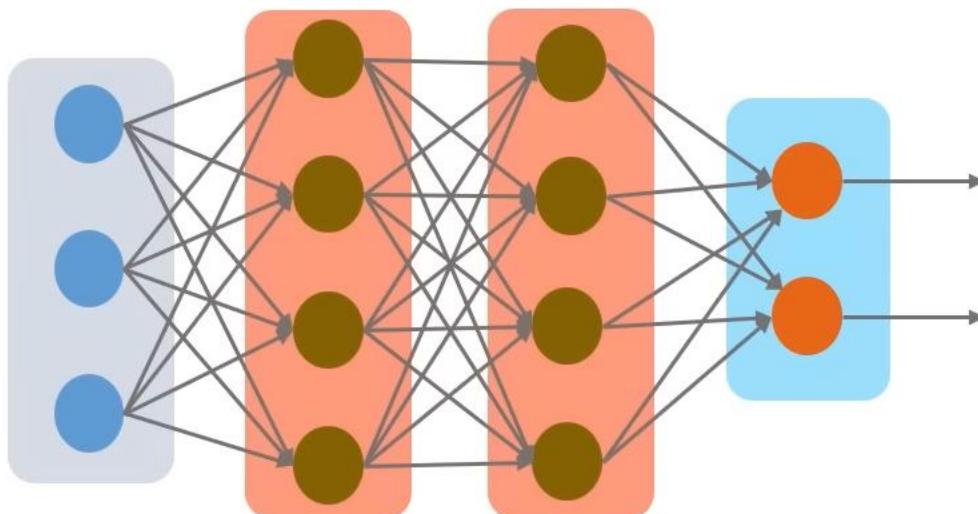
In a Feed-Forward Network , information flows only in forward direction, from the input nodes, through the hidden layers (if any) and to the output nodes. There are no cycles or loops in the network.



- Decisions are based on current input
- No memory about the past
- No future scope

# Why Recurrent Neural Network?

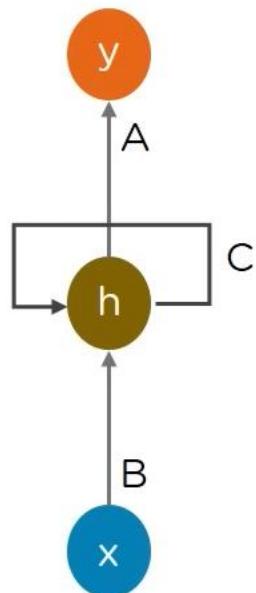
## Issues in Feed Forward Neural Network



- 01 cannot handle sequential data
- 02 considers only the current input
- 03 cannot memorize previous inputs

# Why Recurrent Neural Network?

Solution to Feed Forward Neural Network

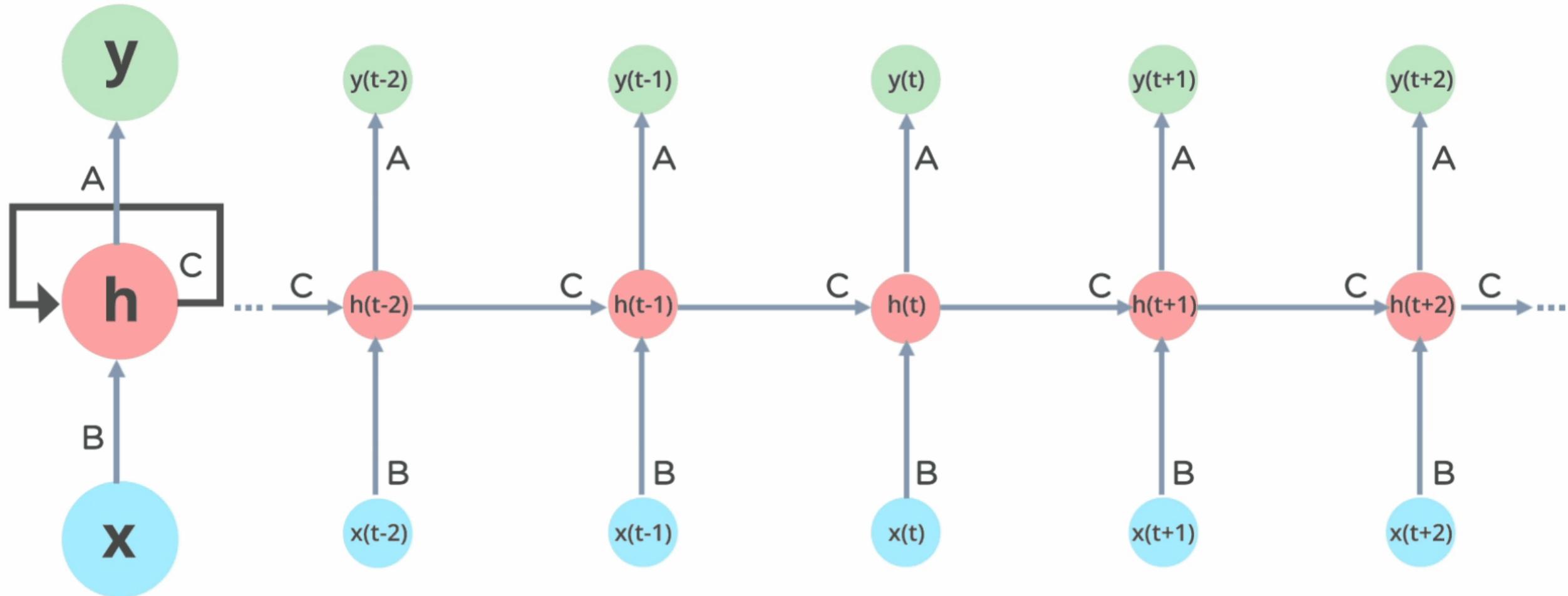


Recurrent Neural  
Network



- 01 can handle sequential data
- 02 considers the current input and also the previously received inputs
- 03 can memorize previous inputs due to its internal memory

# How RNN Looks Like



# Applications of RNN

---



"A Dog catching a ball in mid air"

Image captioning

RNN is used to caption an image by analyzing the activities present in it

# Applications of RNN

---



Time series prediction

Any time series problem like predicting the prices of stocks in a particular month can be solved using RNN

# Applications of RNN

---



Natural Language Processing

Text mining and Sentiment analysis can be carried out using RNN for  
Natural Language Processing

When it rains, look for rainbows.  
When it's dark, look for stars.

Positive Sentiment

# Applications of RNN



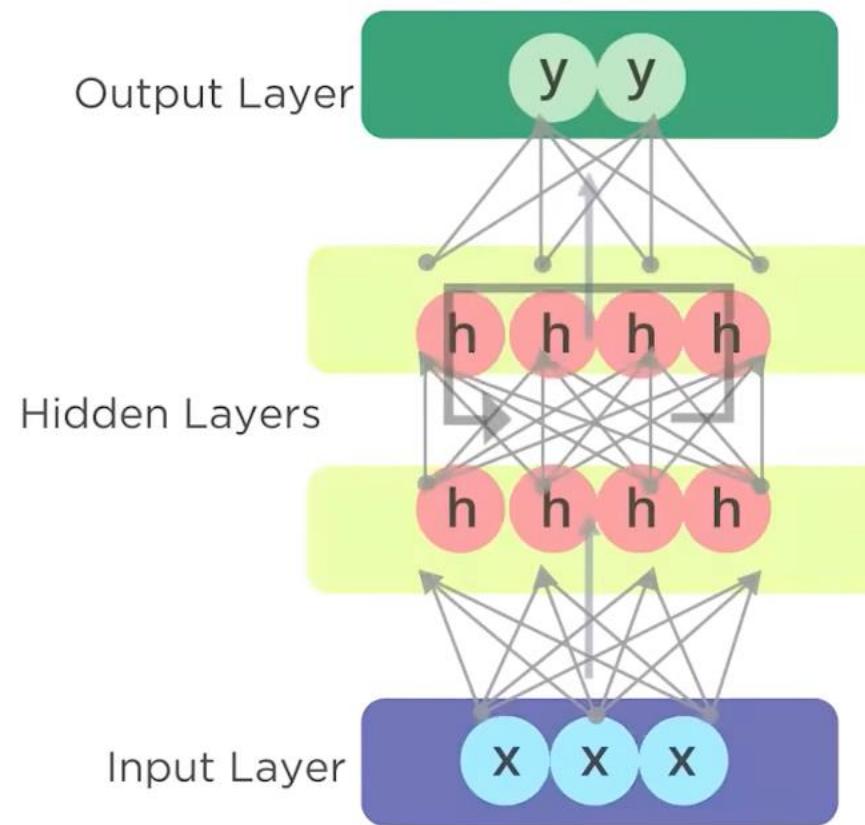
Here the person is speaking in English and it is getting translated into Chinese, Italian, French, German and Spanish languages

Machine Translation

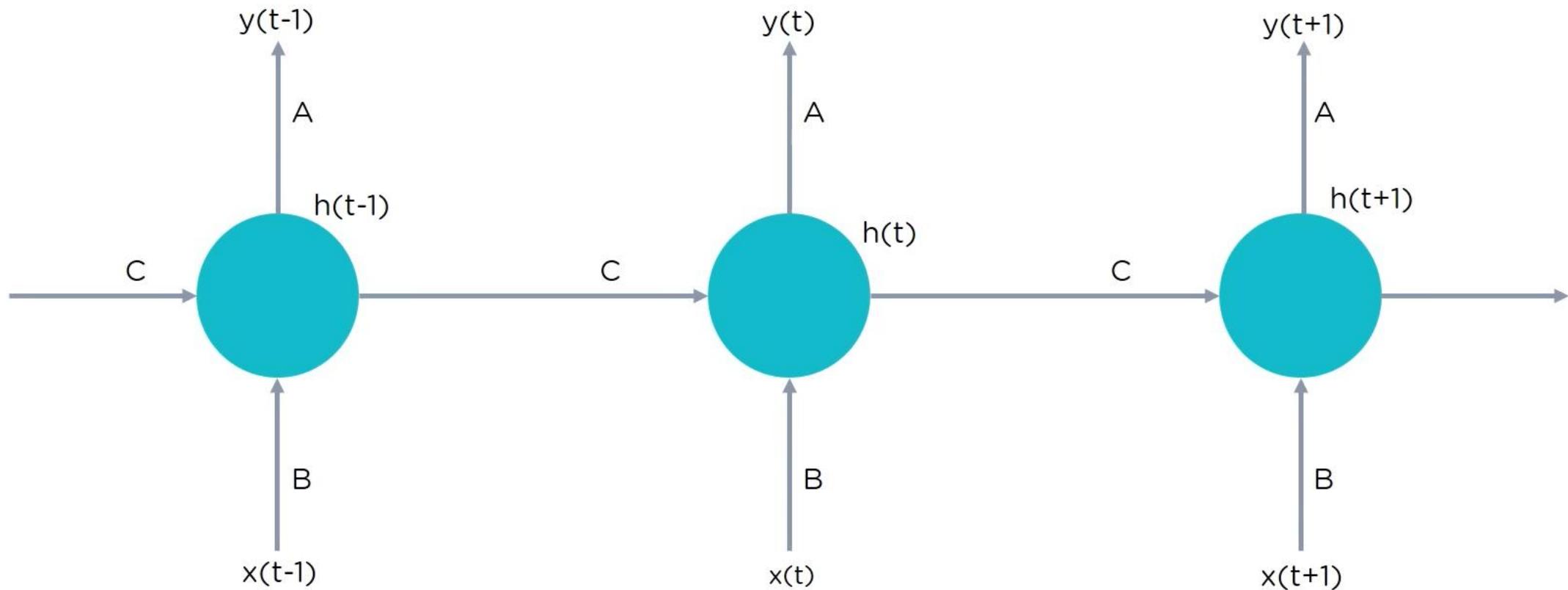
Given an input in one language, RNN can be used to translate the input into different languages as output

# How does a RNN look like?

---



# How does a RNN work?

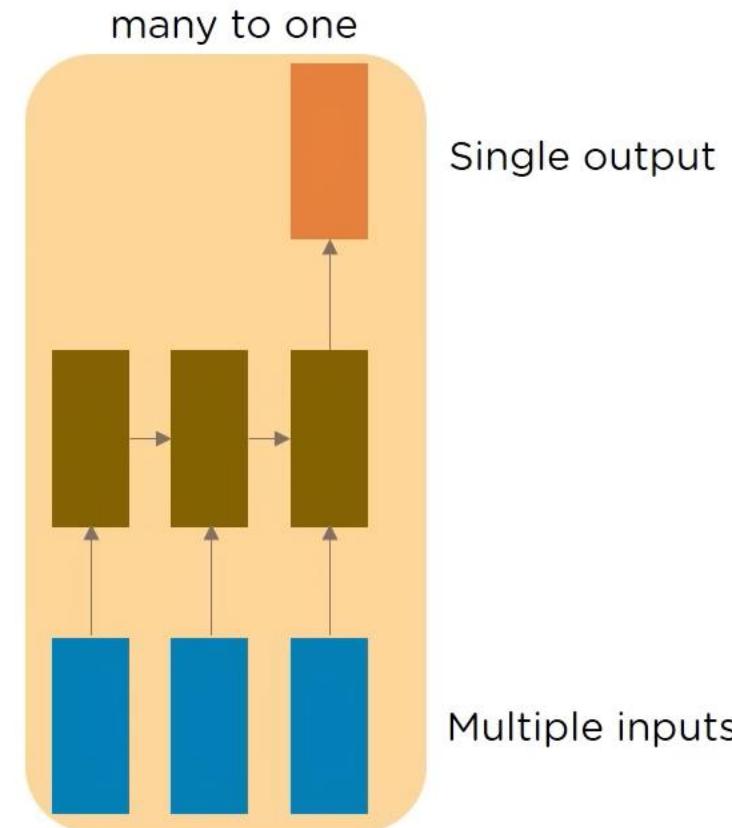


$$h(t) = f_c(h(t-1), x(t))$$

$h(t)$  = new state  
 $f_c$  = function with parameter  $c$   
 $h(t-1)$  = old state  
 $x(t)$  = input vector at time step  $t$

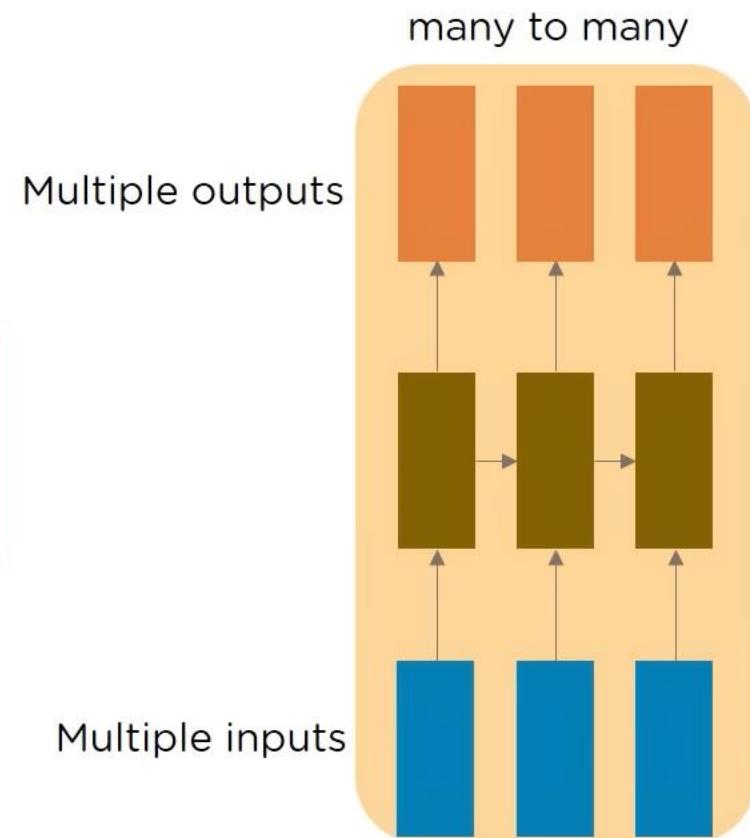
# Types of Recurrent Neural Network

many to one network takes in a sequence of inputs. Example: Sentiment analysis where a given sentence can be classified as expressing positive or negative sentiments

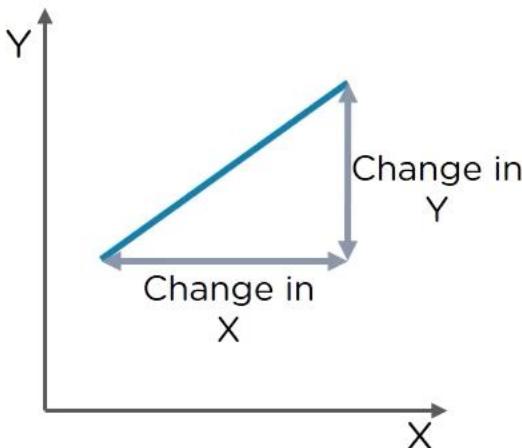


# Types of Recurrent Neural Network

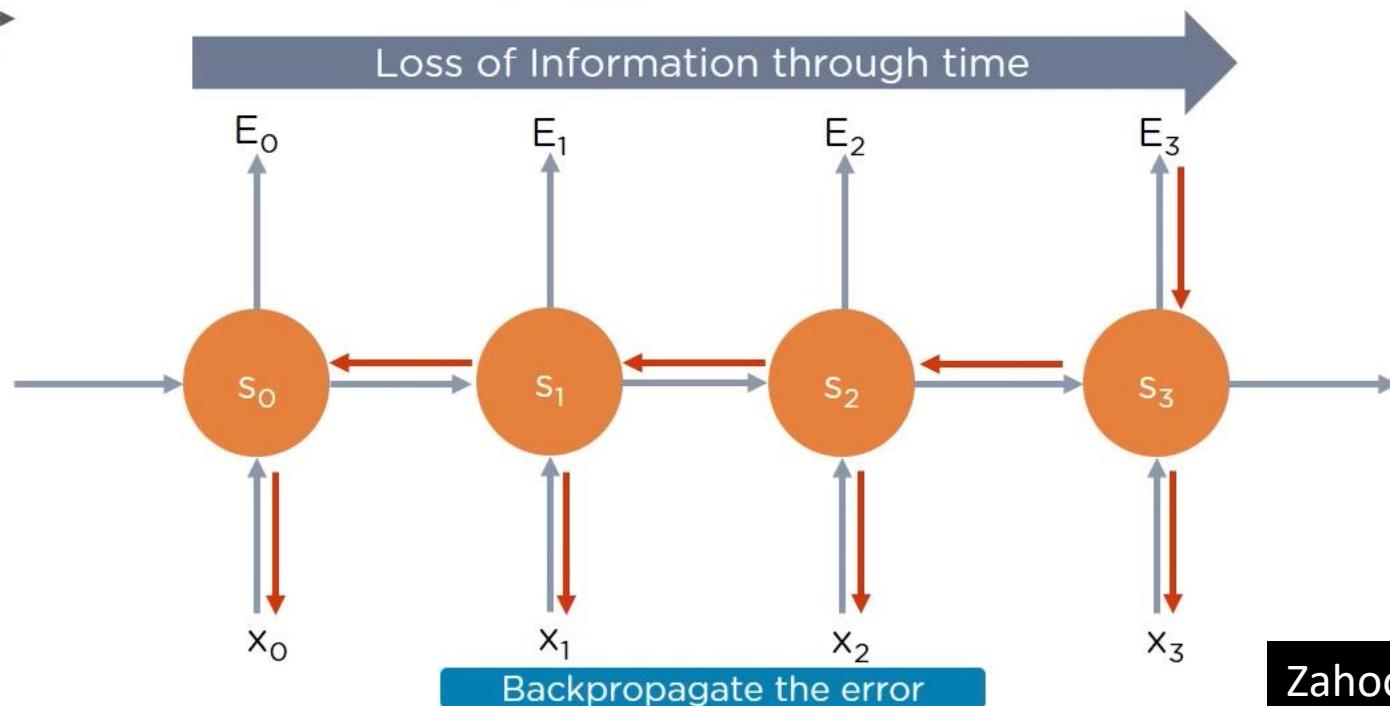
many to many network takes in a sequence of inputs and generates a sequence of outputs. Example: Machine Translation



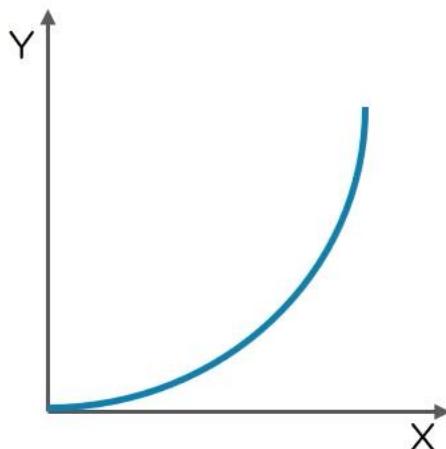
# Vanishing Gradient Problem



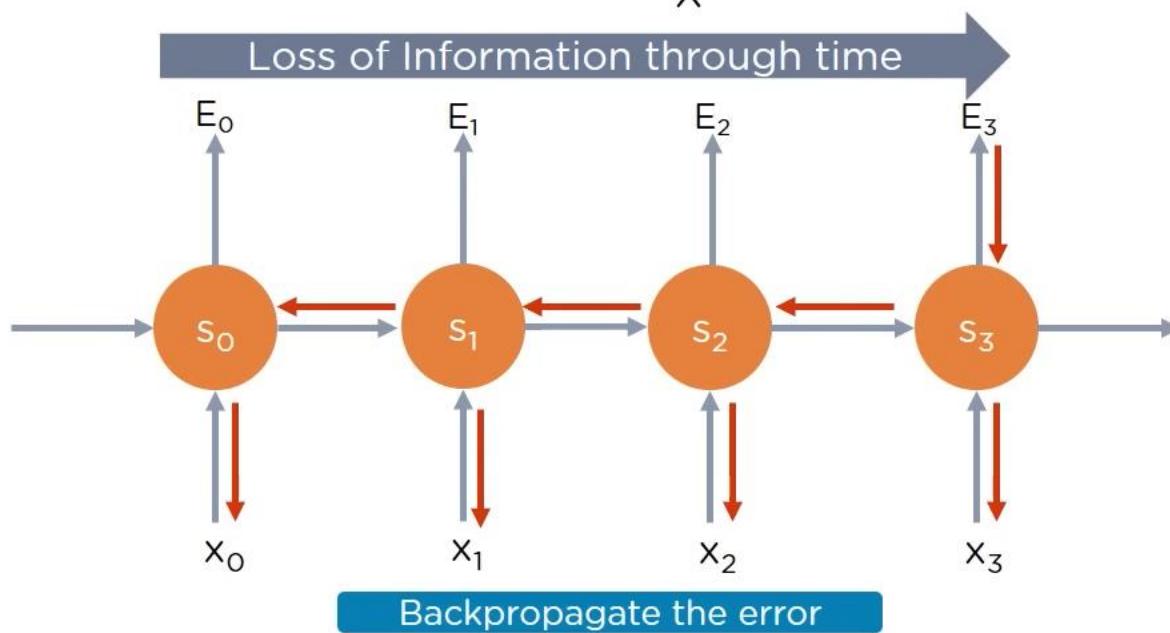
While training a RNN, your slope can be either too small or very large and this makes training difficult. When the slope is too small, the problem is known as *Vanishing gradient*.



# Exploding Gradient Problem



When the slope tends to grow exponentially instead of decaying, this problem is called *Exploding gradient*.



Issues in Gradient Problem

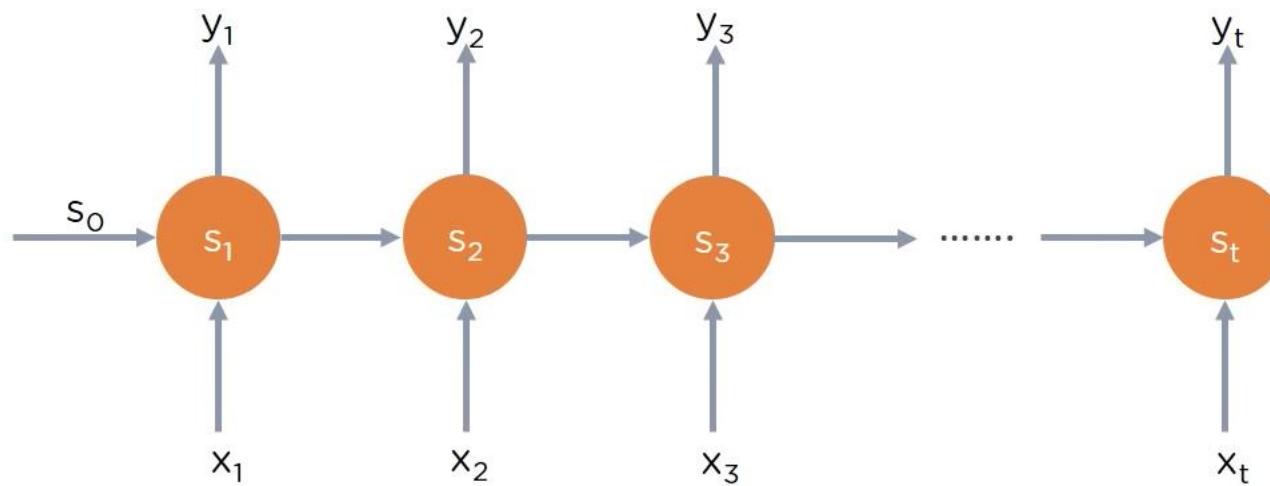
- Long training time
- Poor performance
- Bad accuracy

# Explaining Gradient Problem

Consider the following 2 examples to understand what should be the next word in the sequence:

The person who took my bike and....., was a thief.

The students who got into Engineering with ....., \_\_\_\_\_ from Asia.



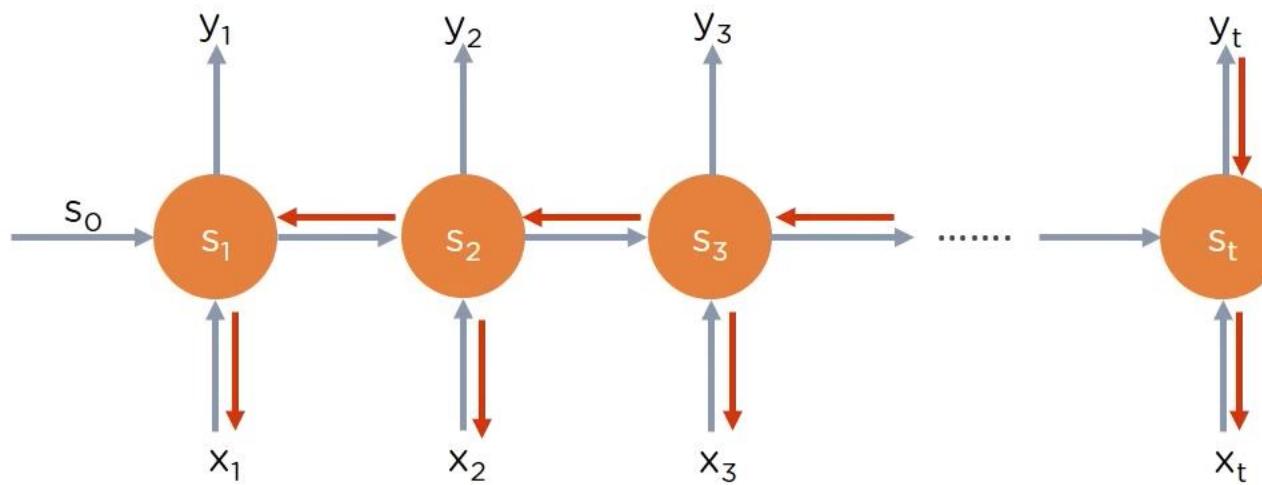
# Explaining Gradient Problem

Consider the following 2 examples:

The person who took my bike and....., a thief.

The students who got into Engineering with ....., were from Asia .

In order to understand what would be the next word in the sequence, the RNN must memorize the previous context whether the subject was singular noun or plural noun



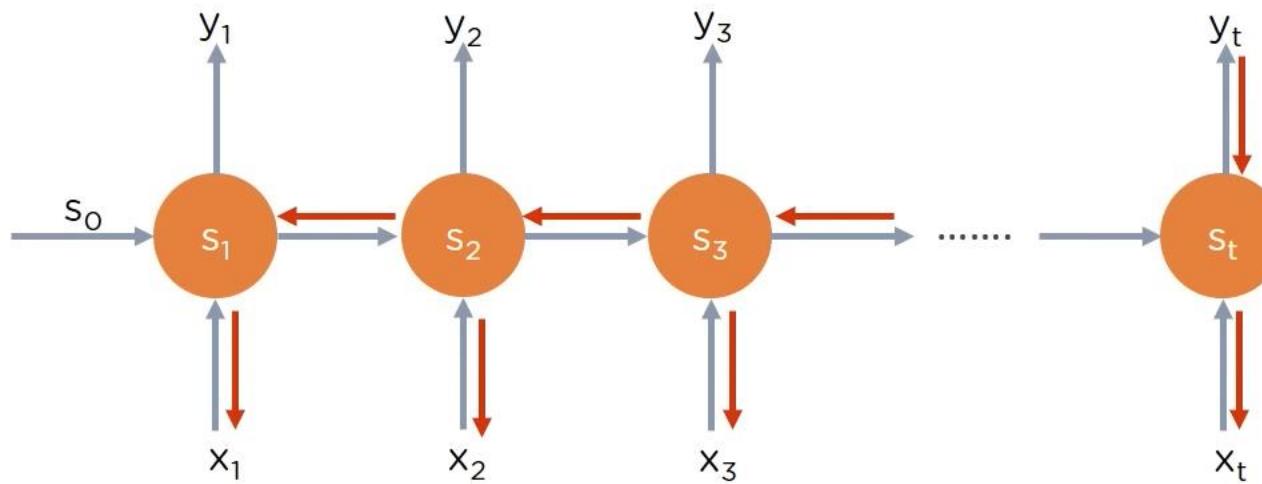
# Explaining Gradient Problem

Consider the following 2 examples:

The person who took my bike and....., was a thief.

The students who got into Engineering with ....., were from Asia .

In order to understand what would be the next word in the sequence, the RNN must memorize the previous context whether the subject was singular noun or plural noun

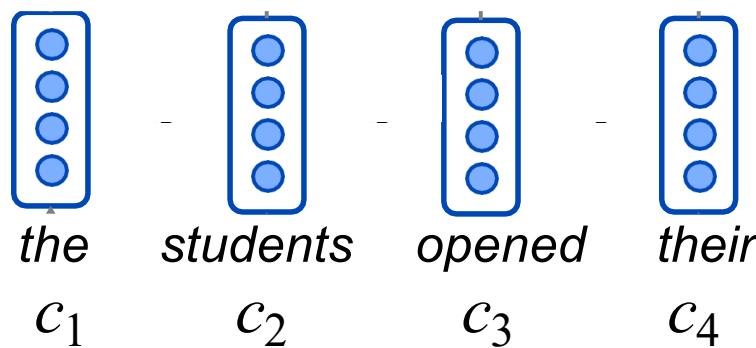


It might be sometimes difficult for the error to backpropagate to the beginning of the sequence to predict what should be the output

# A RNN Language Model

word embeddings  $c_1, c_2, c_3, c_4$

$$e^{(t)} = Ex^{(t)}$$

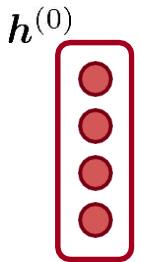


# A RNN Language Model

hidden states

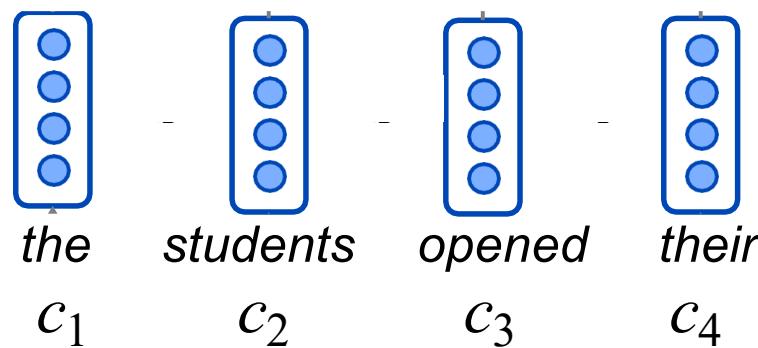
$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$  is initial hidden state!



word embeddings  $c_1, c_2, c_3, c_4$

$$e^{(t)} = E x^{(t)}$$



# A RNN Language Model

hidden states

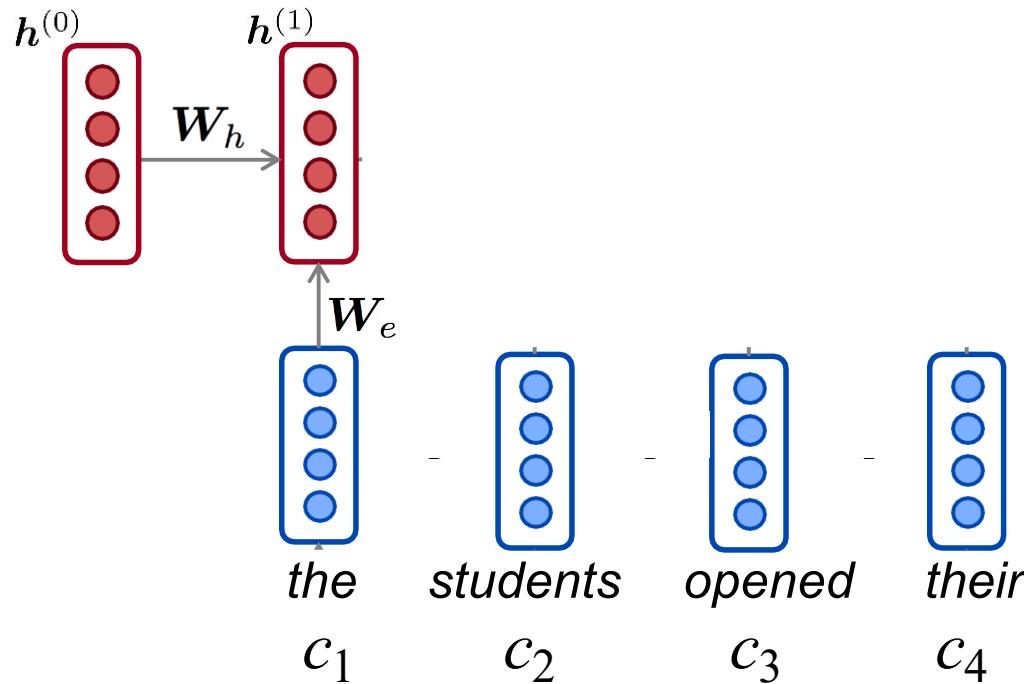
$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$  is initial hidden state!

word embeddings

$c_1, c_2, c_3, c_4$

$$e^{(t)} = Ex^{(t)}$$



# A RNN Language Model

hidden states

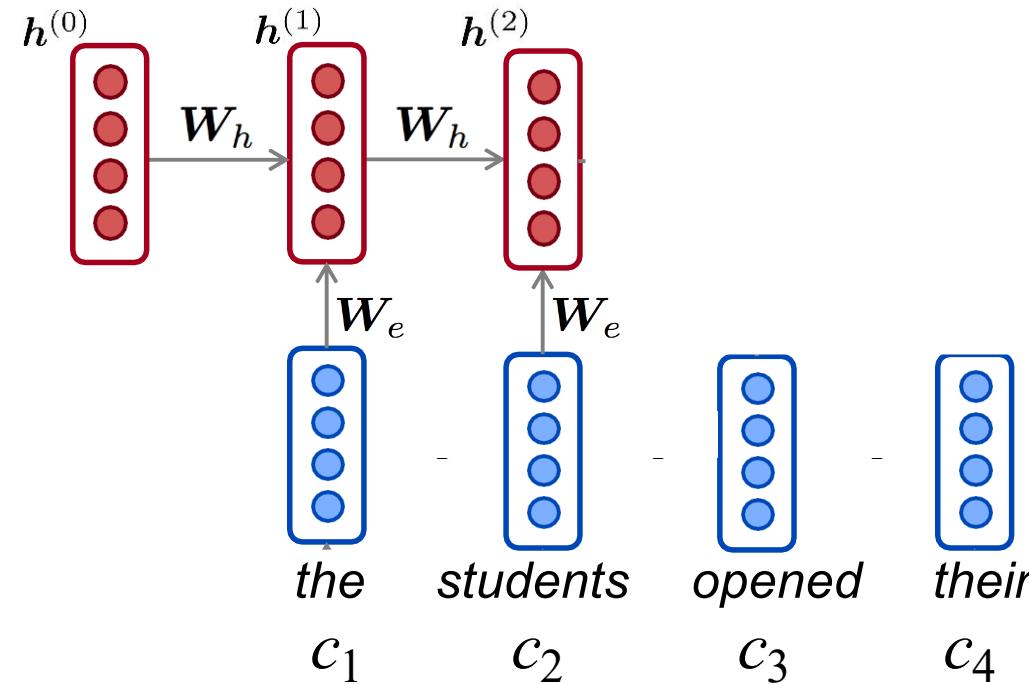
$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$  is initial hidden state!

word embeddings

$c_1, c_2, c_3, c_4$

$$e^{(t)} = E x^{(t)}$$



# A RNN Language Model

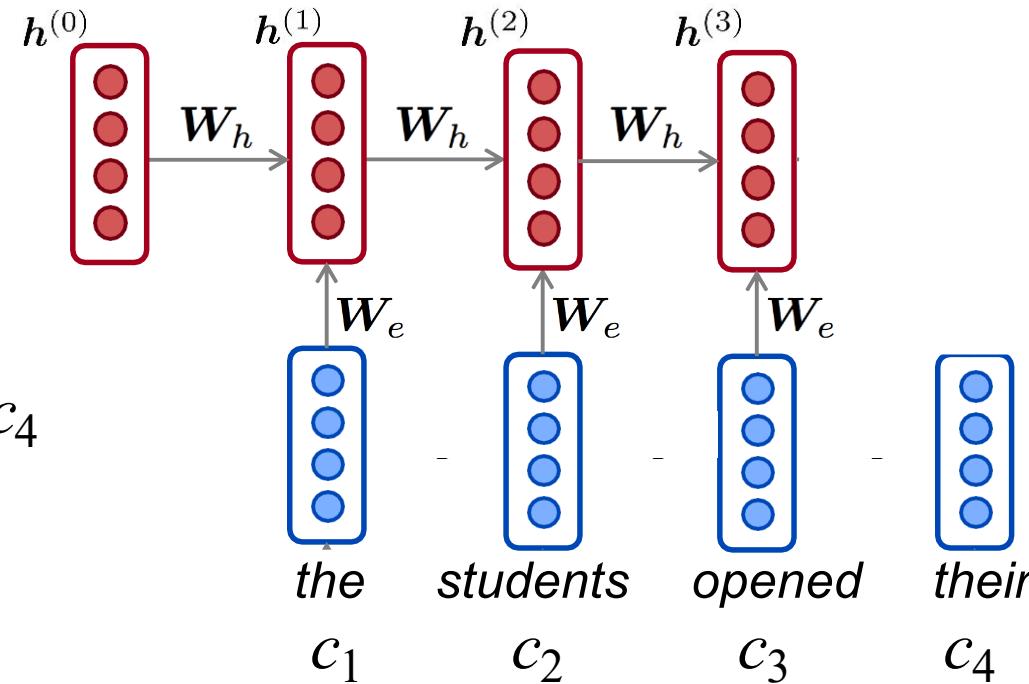
hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$  is initial hidden state!

word embeddings  $c_1, c_2, c_3, c_4$

$$e^{(t)} = Ex^{(t)}$$



# A RNN Language Model

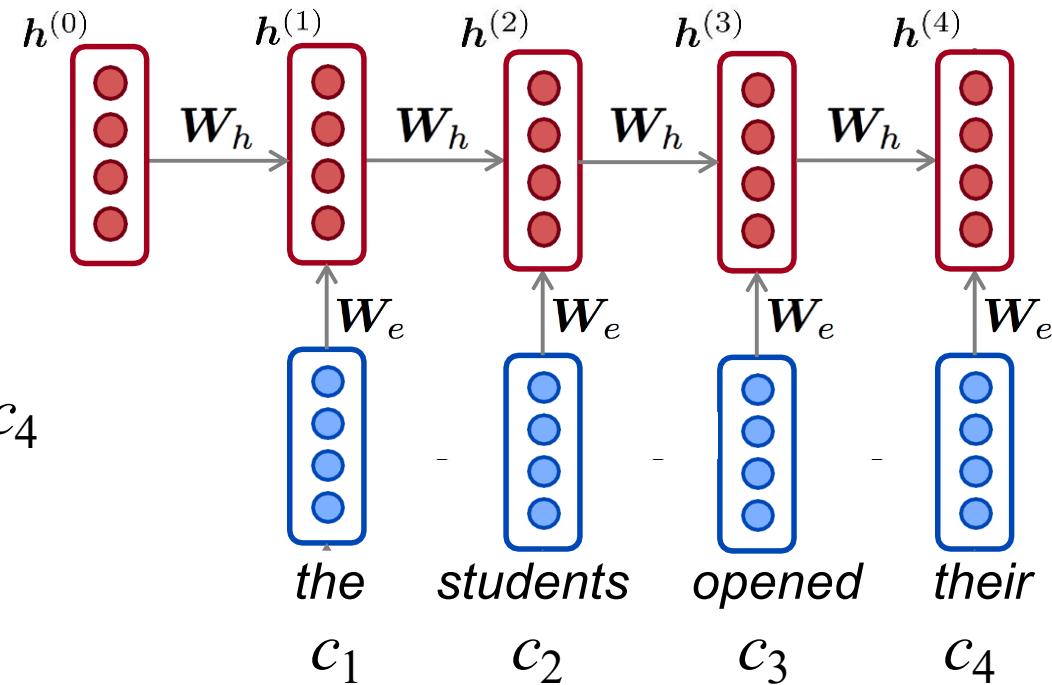
hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$  is initial hidden state!

word embeddings  $c_1, c_2, c_3, c_4$

$$e^{(t)} = Ex^{(t)}$$



# A RNN Language Model

$$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$

output distribution

$$\hat{y} = \text{softmax}(W_2 h^{(t)})$$

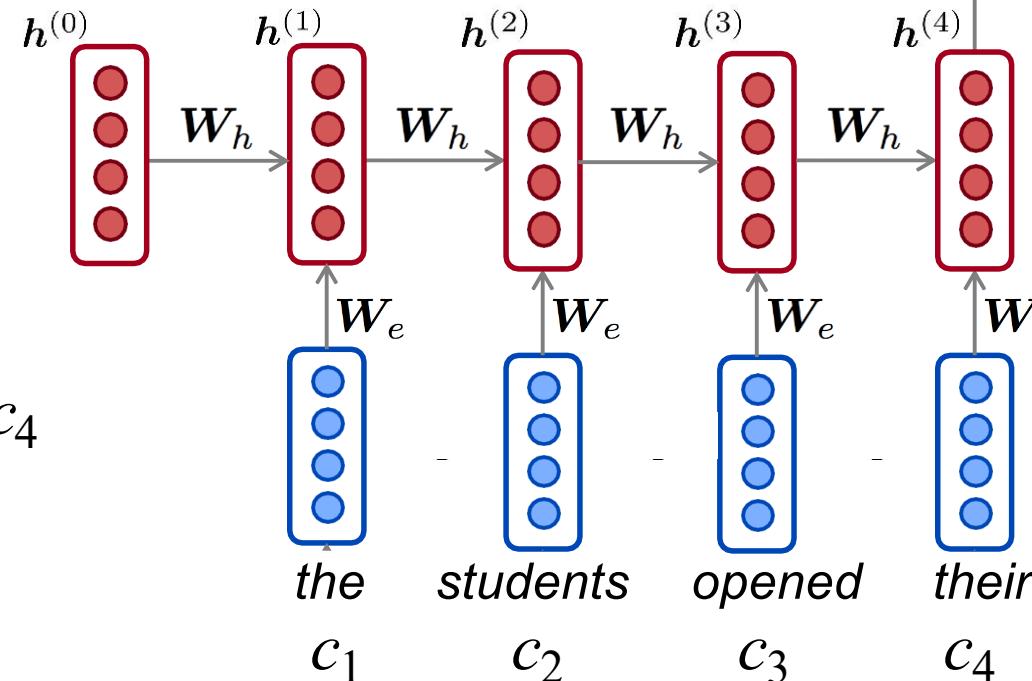
hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

$h^{(0)}$  is initial hidden state!

word embeddings  $c_1, c_2, c_3, c_4$

$$e^{(t)} = E x^{(t)}$$



$$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$

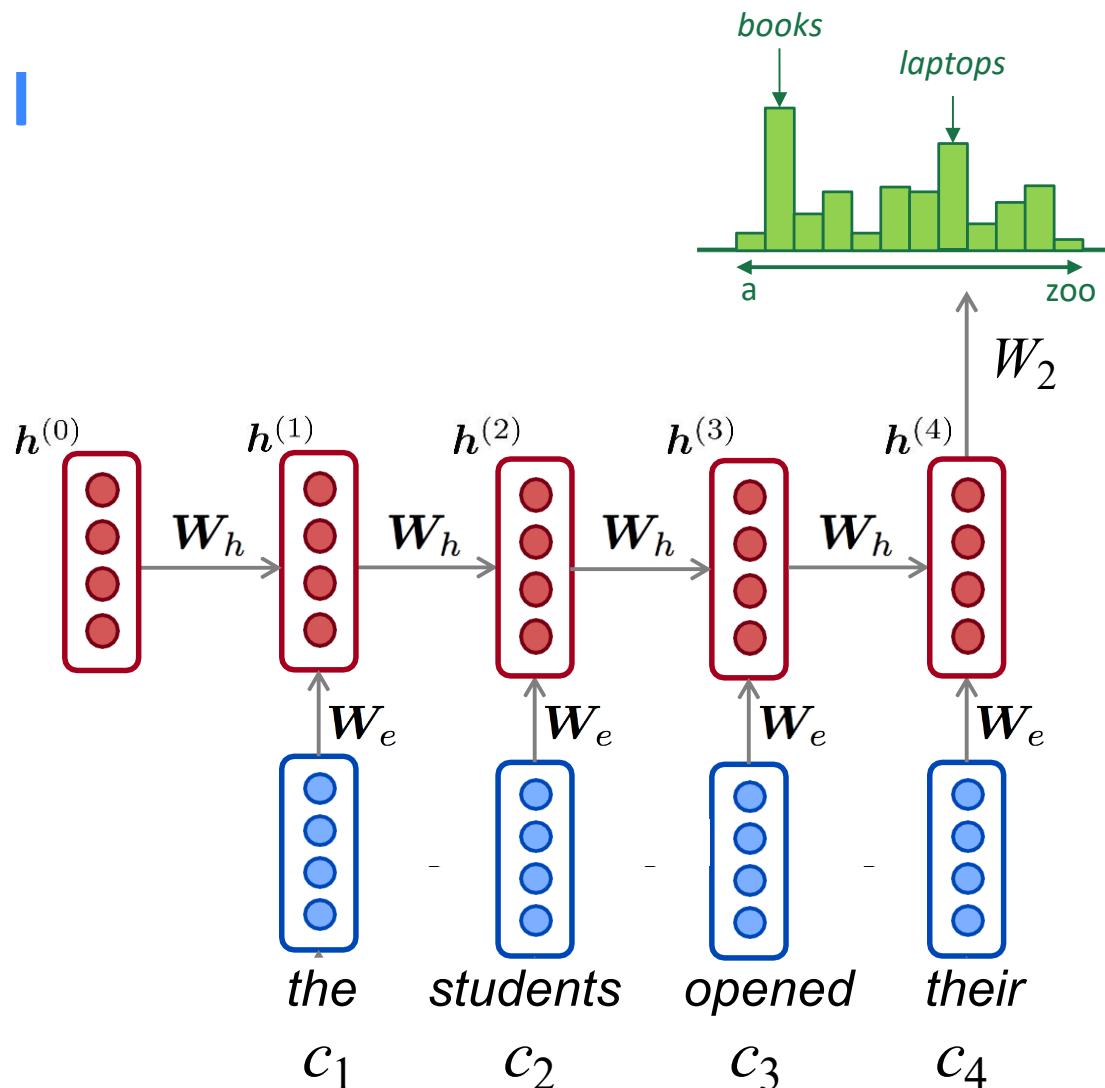
why is this good?

RNN **Advantages:**

- Can process **any length** input
- Model size **doesn't increase** for longer input
- Computation for step  $t$  can (in theory) use information from **many steps back**
- Weights are **shared** across timesteps  $\rightarrow$  representations are shared

RNN **Disadvantages:**

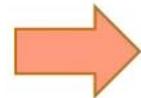
- Recurrent computation is **slow**
- In practice, difficult to access information from **many steps back**



# Long-Term Dependencies

---

Suppose we try to predict  
the last word in the text



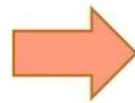
“The clouds are in the sky”



Here we do not need any further context. It's  
pretty clear the last word is going to be “sky”.

# Long-Term Dependencies

Suppose we try to predict  
the last word in the text



"I have been staying in Spain for the last 10  
years... I can speak fluent \_\_\_\_."



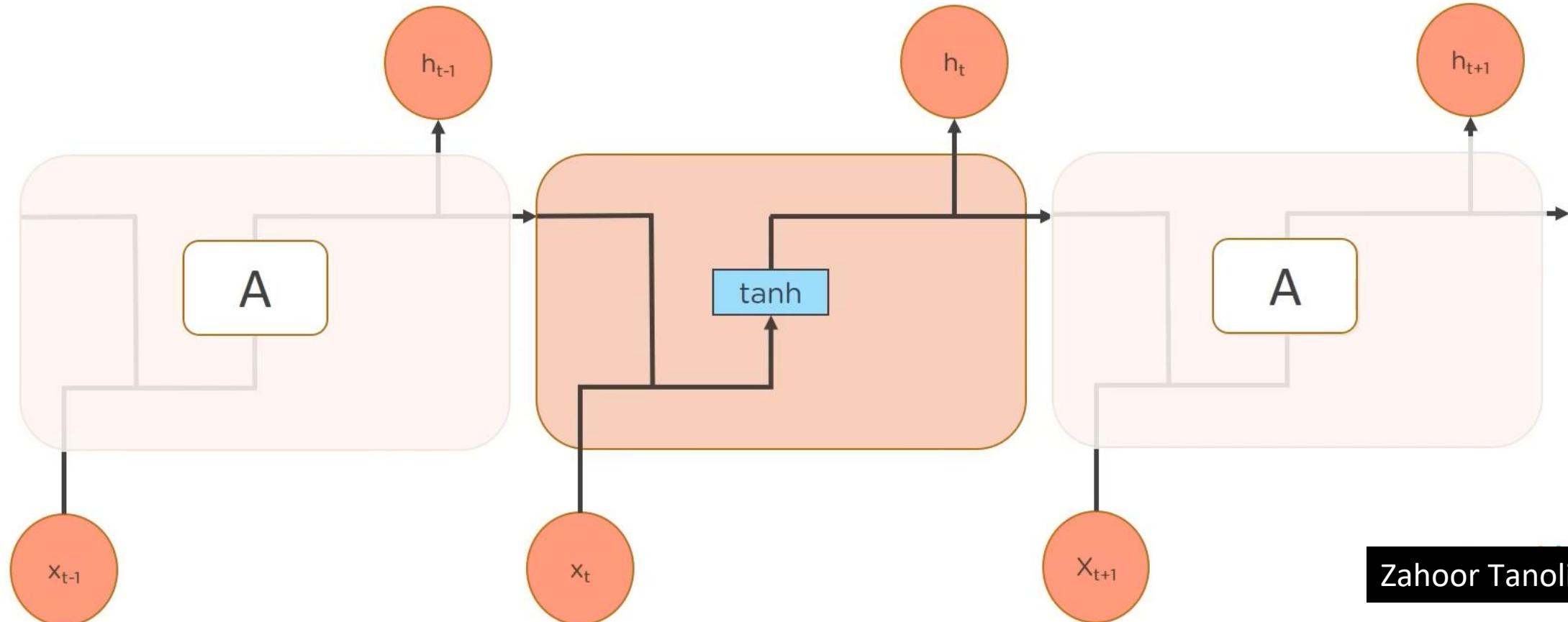
The word we predict will depend on  
the previous few words in context

- Here we need the context of Spain to predict the last word in the text.
- It's possible that the gap between the relevant information and the point where it is needed to become very large.
- LSTMs help us solve this problem.

# Long Short-Term Memory Networks

LSTMs are special kind of Recurrent Neural Networks, capable of learning long-term dependencies. Remembering information for long periods of time is their default behavior.

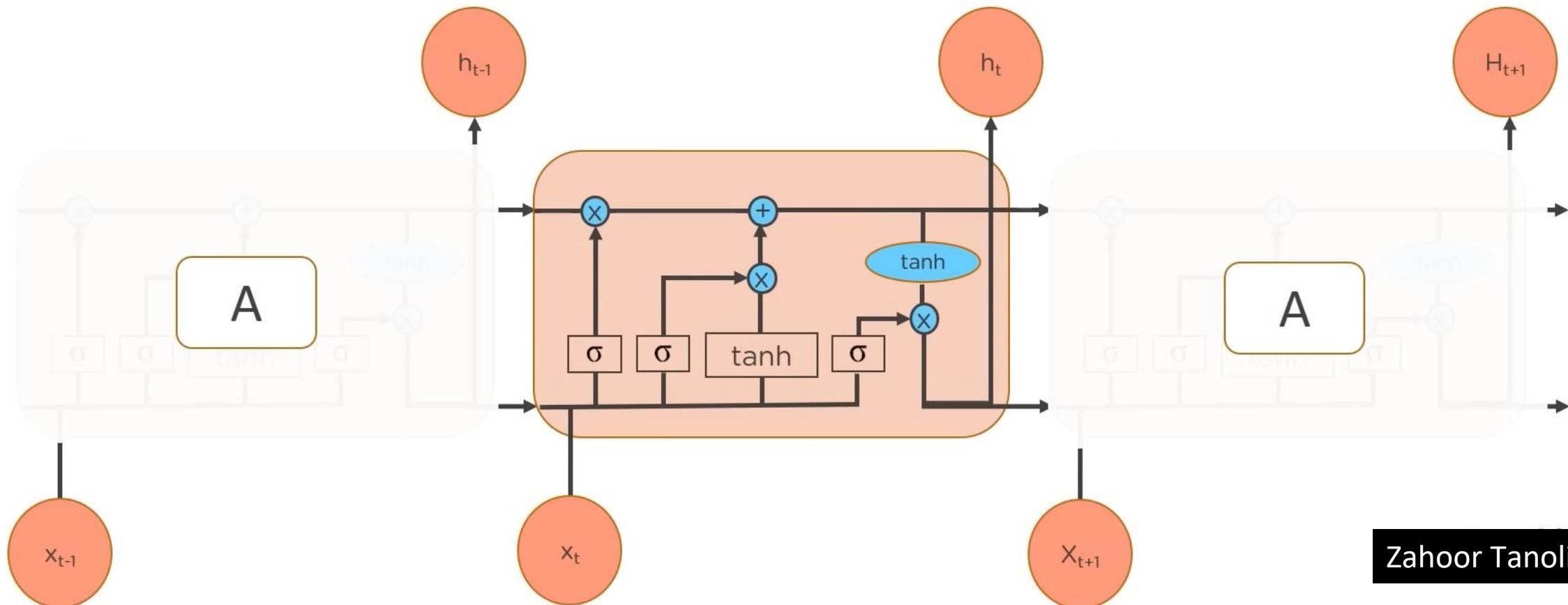
All recurrent neural networks have the form of a chain of repeating modules of neural network. In standard RNNs, this repeating module will have a very simple structure, such as a single tanh layer.



# Long Short-Term Memory Networks

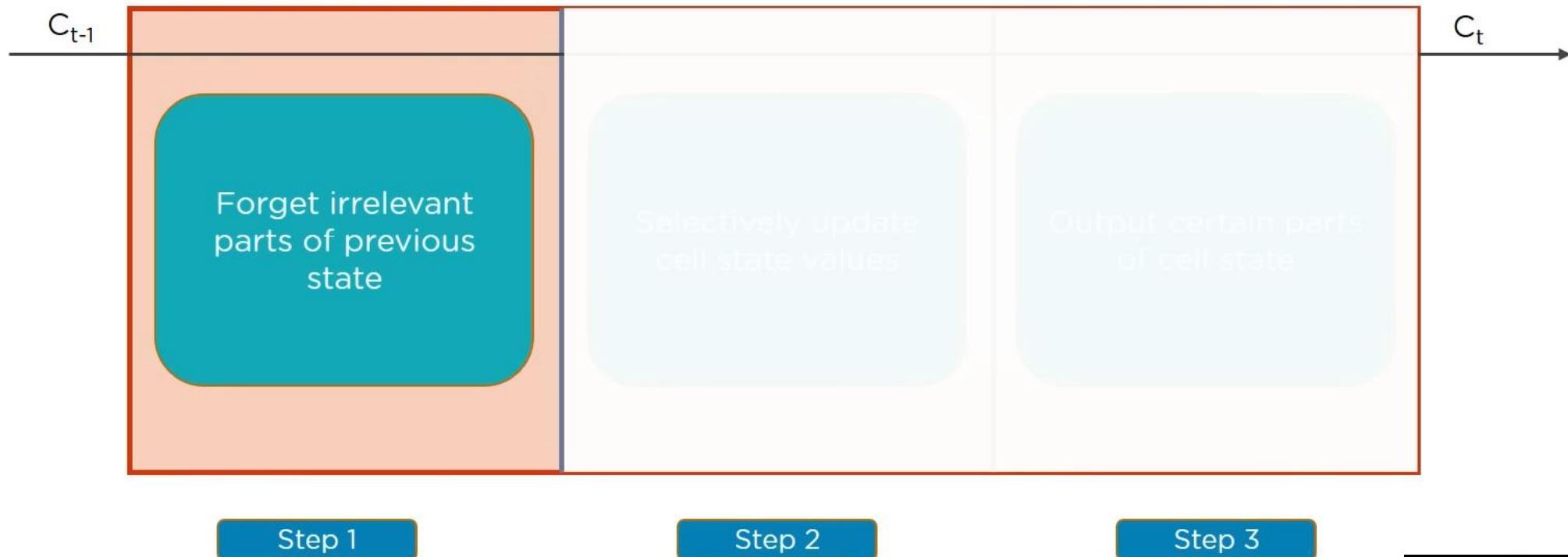
LSTMs are special kind of Recurrent Neural Networks, capable of learning long-term dependencies. Remembering information for long periods of time is their default behavior.

LSTMs also have a chain like structure, but the repeating module has a different structure. Instead of having a single neural network layer, there are four interacting layers communicating in a very special way.



# Long Short-Term Memory Networks

3 step process of LSTMs



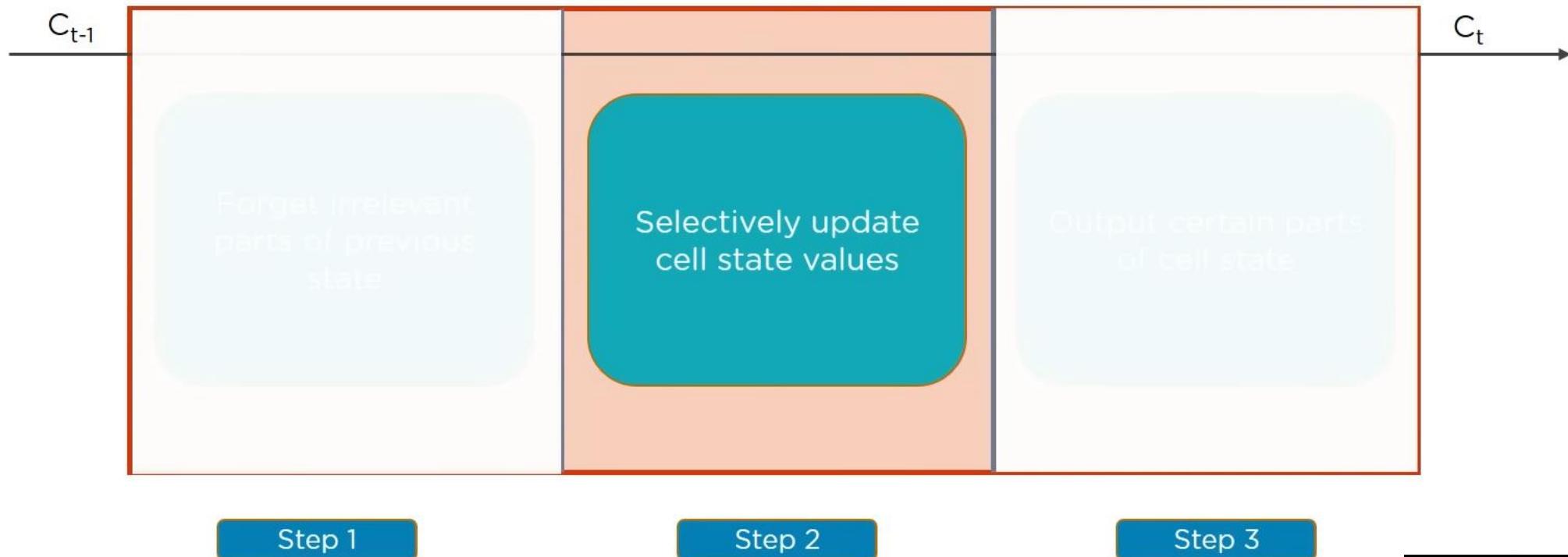
Step 1

Step 2

Step 3

# Long Short-Term Memory Networks

3 step process of LSTMs



Step 1

Step 2

Step 3

# Working of LSTMs

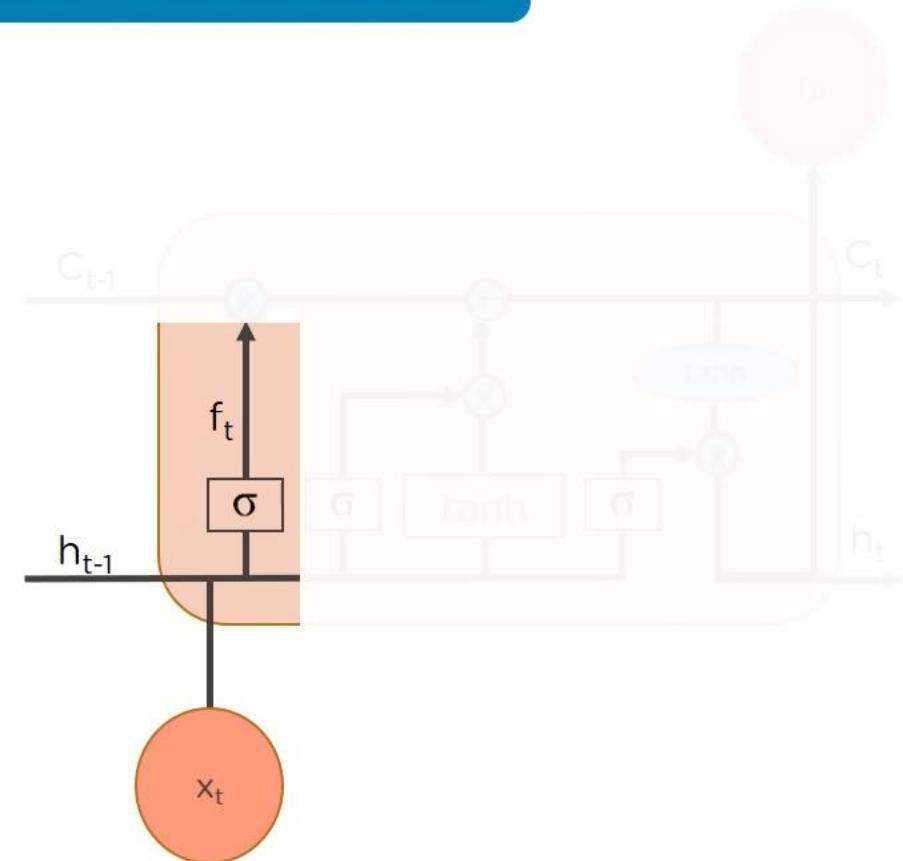
Step-1

Decides how much of the past it should remember

First step in the LSTM is to decide which information to be omitted in from the cell in that particular time step. It is decided by the sigmoid function. It looks at the previous state ( $h_{t-1}$ ) and the current input  $x_t$  and computes the function.

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

$f_t$  = forget gate  
Decides which information to delete that is not important from previous time step



# Working of LSTMs

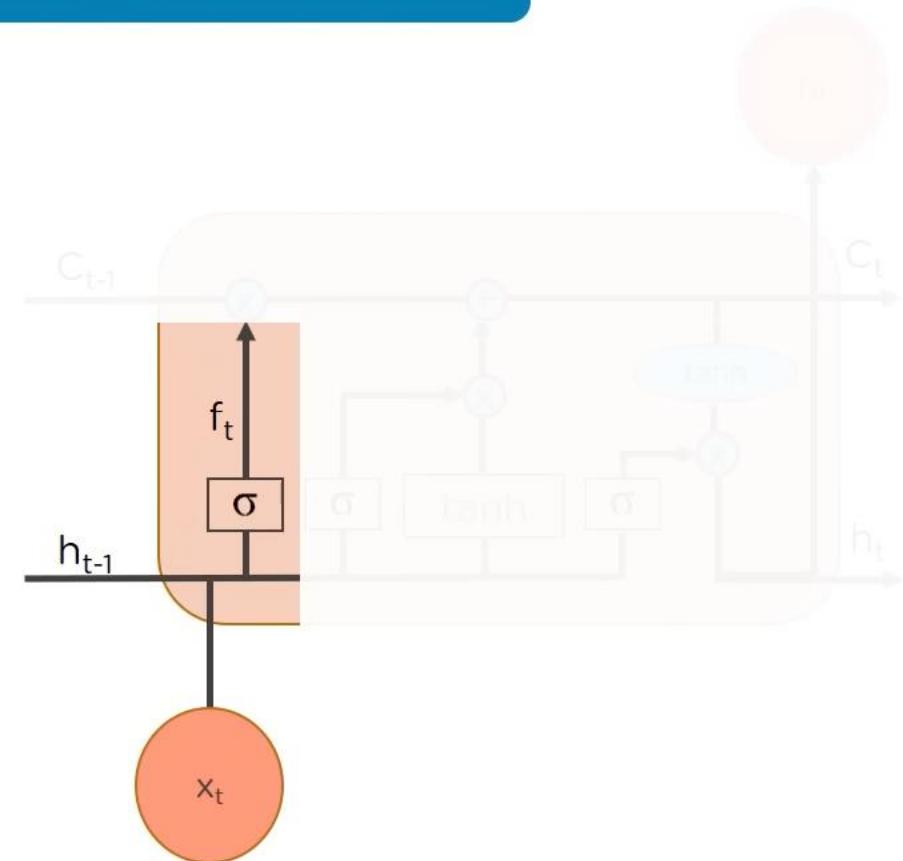
Step-1

Decides how much of the past it should remember

First step in the LSTM is to decide which information to be omitted in from the cell in that particular time step. It is decided by the sigmoid function. It looks at the previous state ( $h_{t-1}$ ) and the current input  $x_t$  and computes the function.

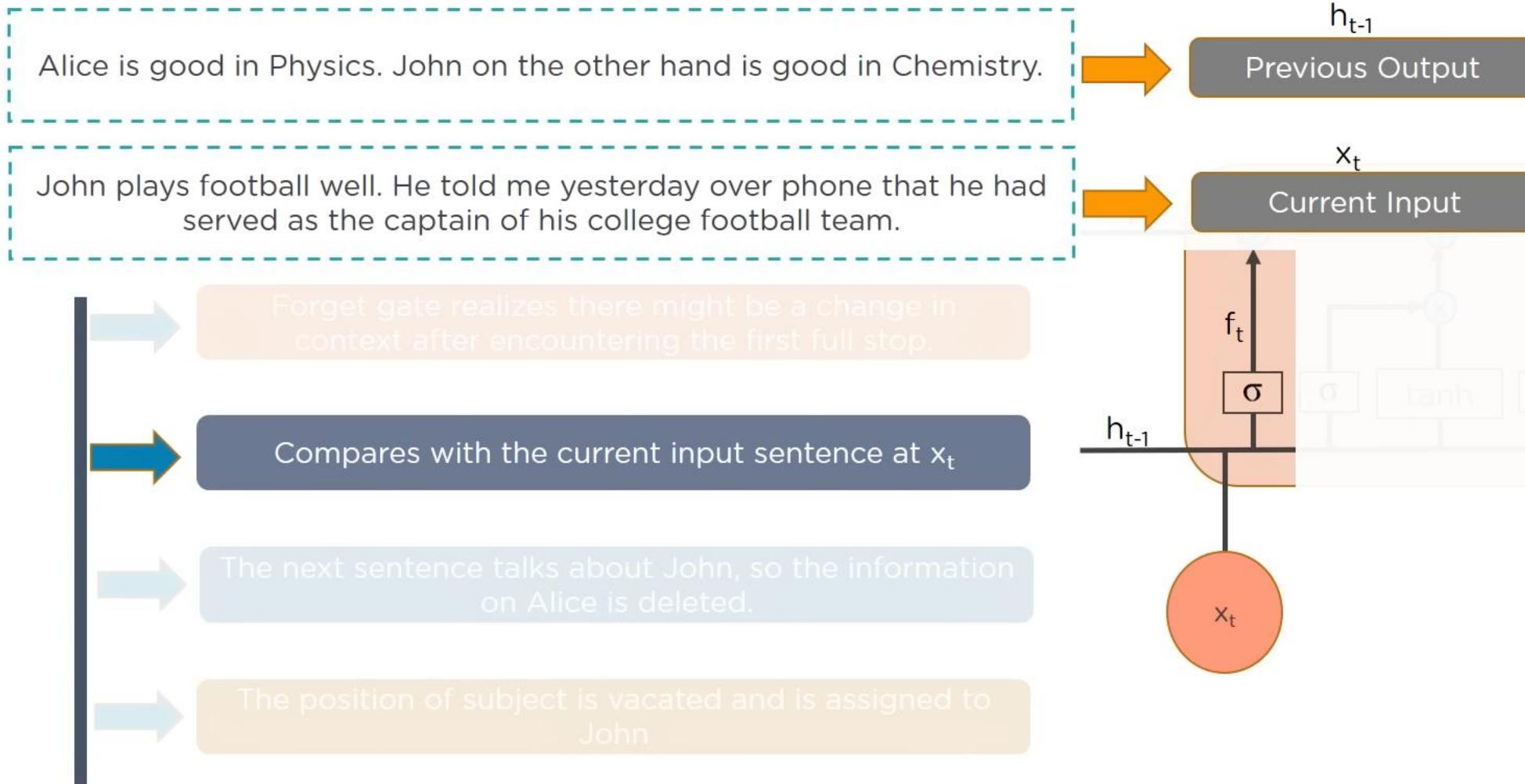
$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

$f_t$  = forget gate  
Decides which information to delete that is not important from previous time step



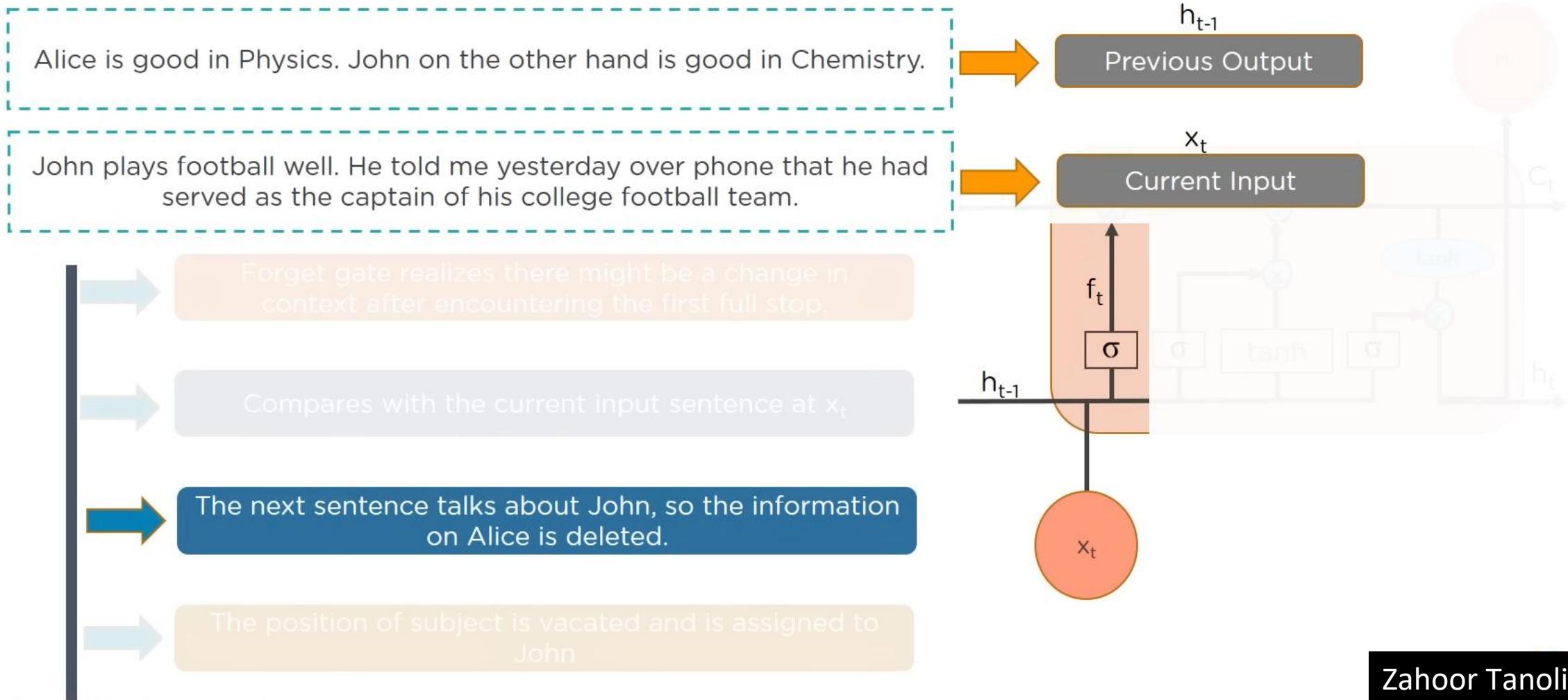
# Working of LSTMs

Consider an LSTM is fed with the following inputs from previous and present time step :



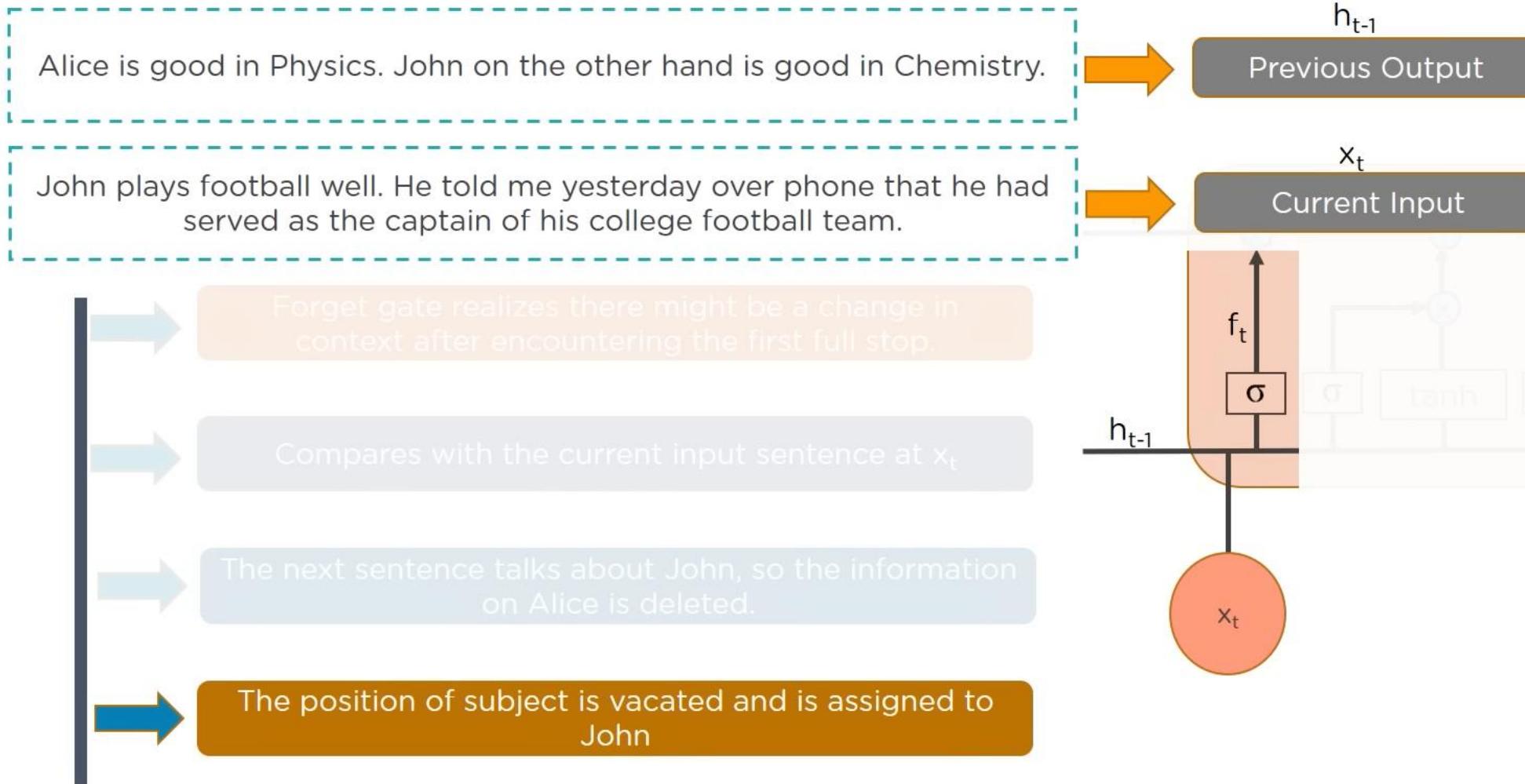
# Working of LSTMs

Consider an LSTM is fed with the following inputs from previous and present time step :



# Working of LSTMs

Consider an LSTM is fed with the following inputs from previous and present time step :



# Working of LSTMs

Step-2

Decides how much should this unit add to the current state

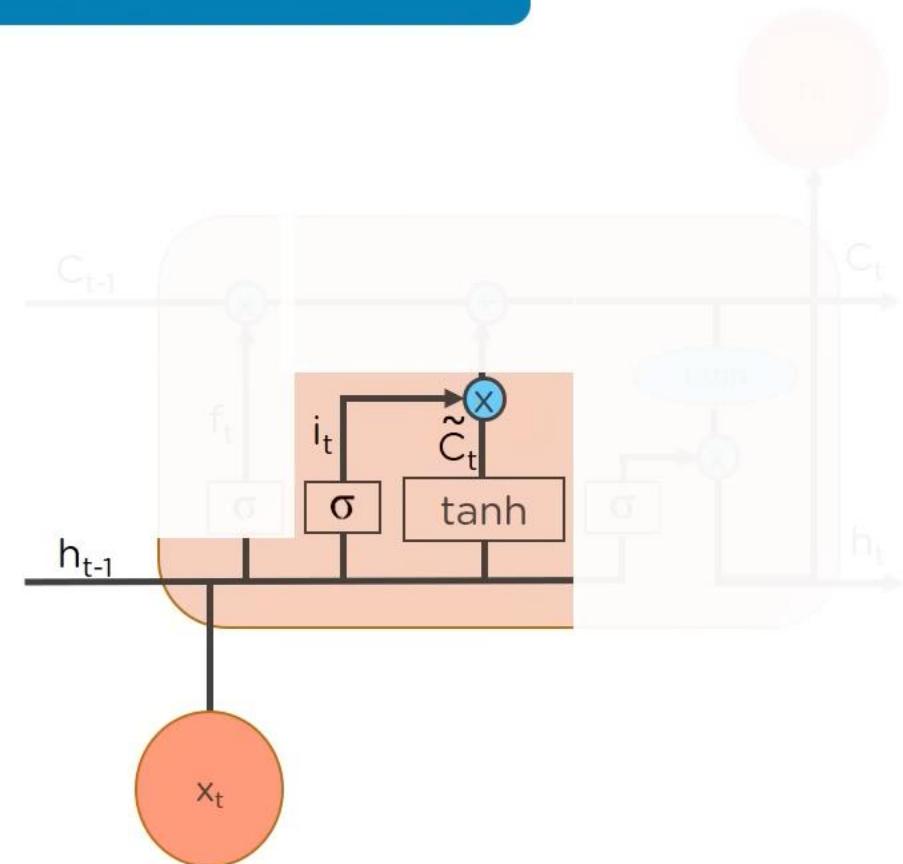
In the second layer, there are 2 parts. One is the sigmoid function and the other is the tanh. In the **sigmoid** function, it decides which values to let through(0 or 1). **tanh** function gives the weightage to the values which are passed deciding their level of importance(-1 to 1).

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

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

$i_t$  = input gate

Determines which information to let through based on its significance in the current time step



# Working of LSTMs

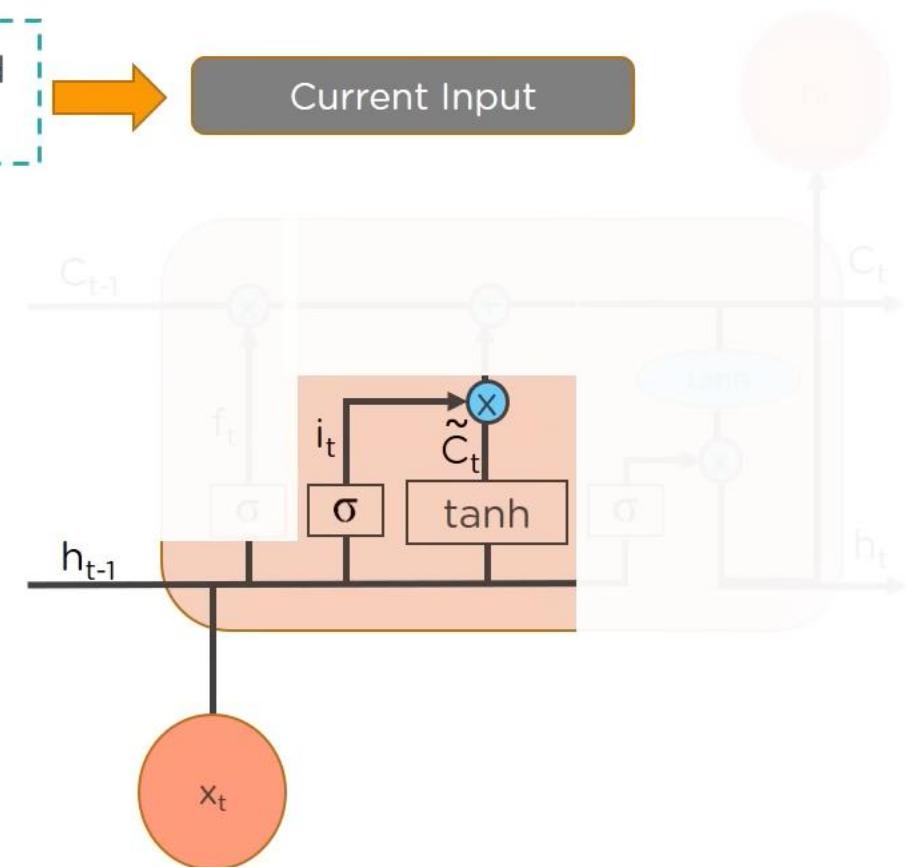
Consider the current input at  $x_t$

John plays football well. He told me yesterday over phone that he had served as the captain of his college football team.

Current Input



Input gate analyses  
the important  
information



# Working of LSTMs

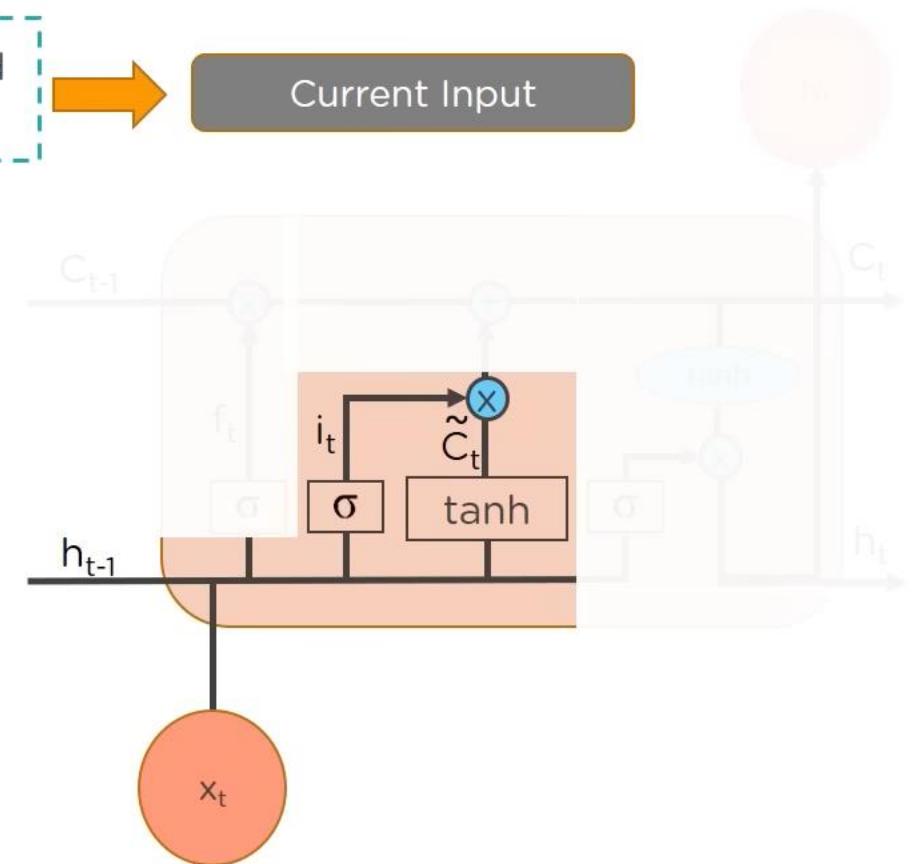
Consider the current input at  $x_t$

John plays football well. He told me yesterday over phone that he had served as the captain of his college football team.

Current Input



John plays football  
and he was the  
captain of his  
college team is  
important



# Working of LSTMs

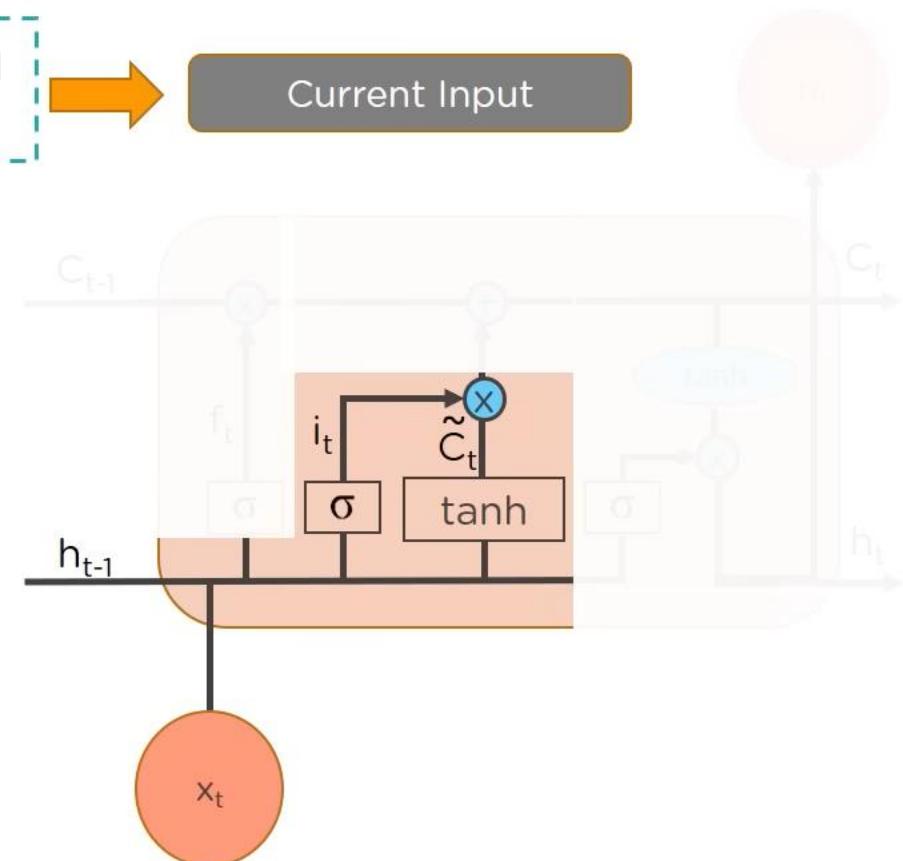
Consider the current input at  $x_t$

John plays football well. He told me yesterday over phone that he had served as the captain of his college football team.

Current Input



He told me over phone yesterday is less important, hence it is forgotten



# Working of LSTMs

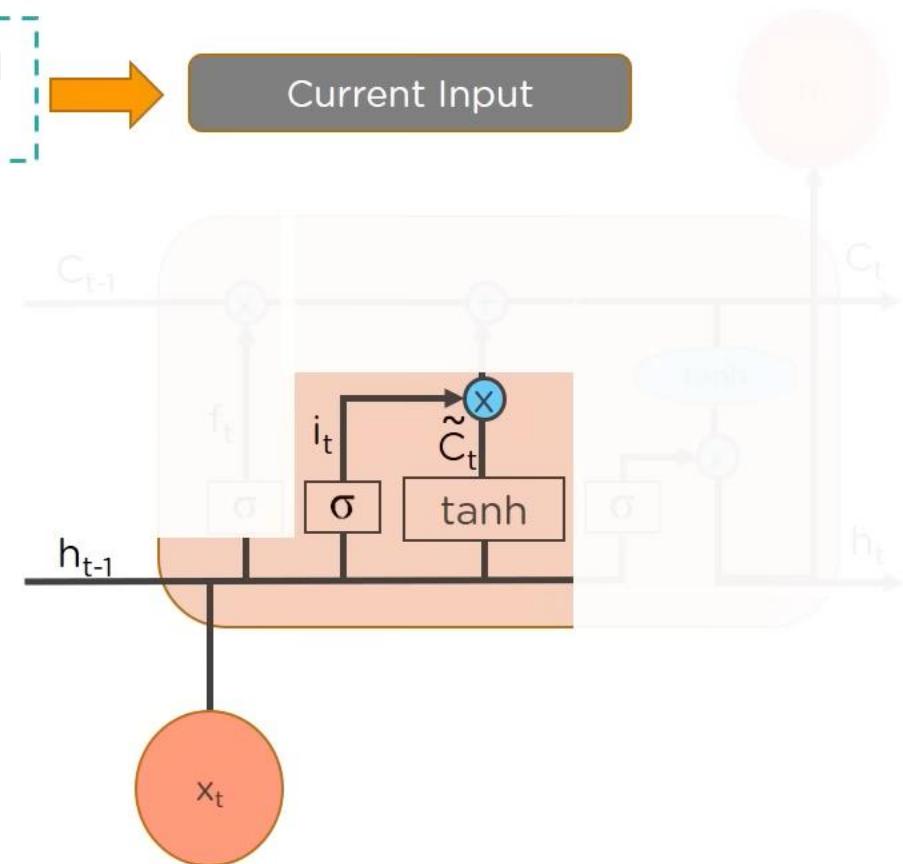
Consider the current input at  $x_t$

John plays football well. He told me yesterday over phone that he had served as the captain of his college football team.

Current Input



This process of adding some new information can be done via the input gate



# Working of LSTMs

Step-3

Decides what part of the current cell state makes it to the output

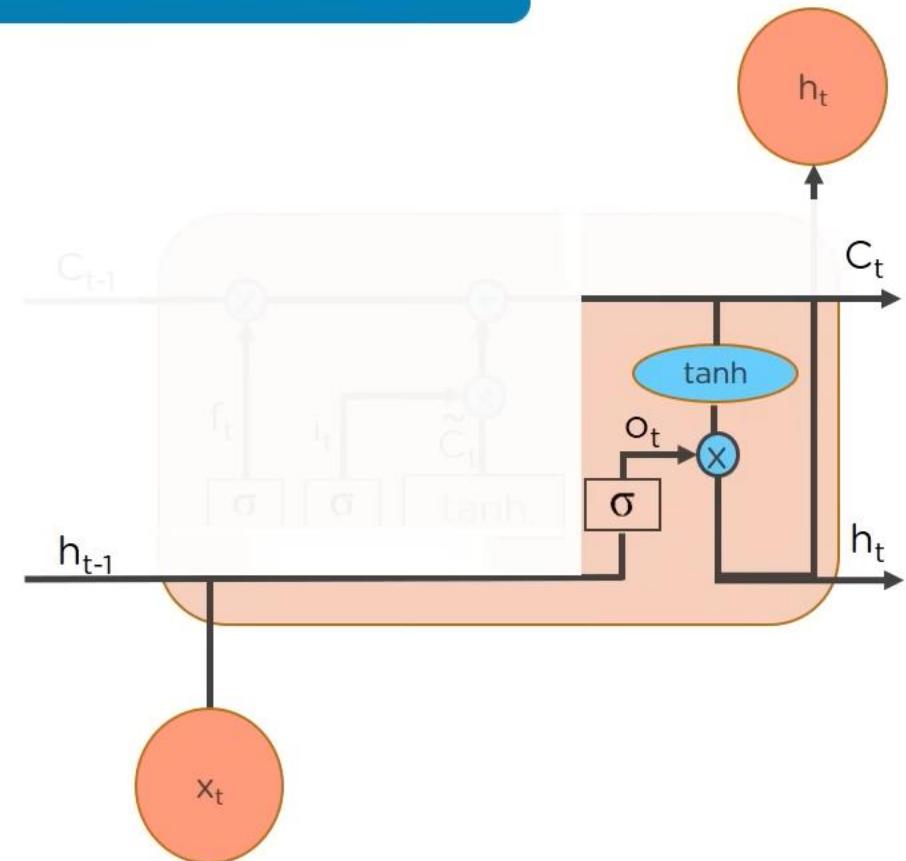
The third step is to decide what will be our output. First, we run a sigmoid layer which decides what parts of the cell state make it to the output. Then, we put the cell state through tanh to push the values to be between -1 and 1 and multiply it by the output of the sigmoid gate.

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

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

$o_t$  = output gate

Allows the passed in information to impact the output in the current time step



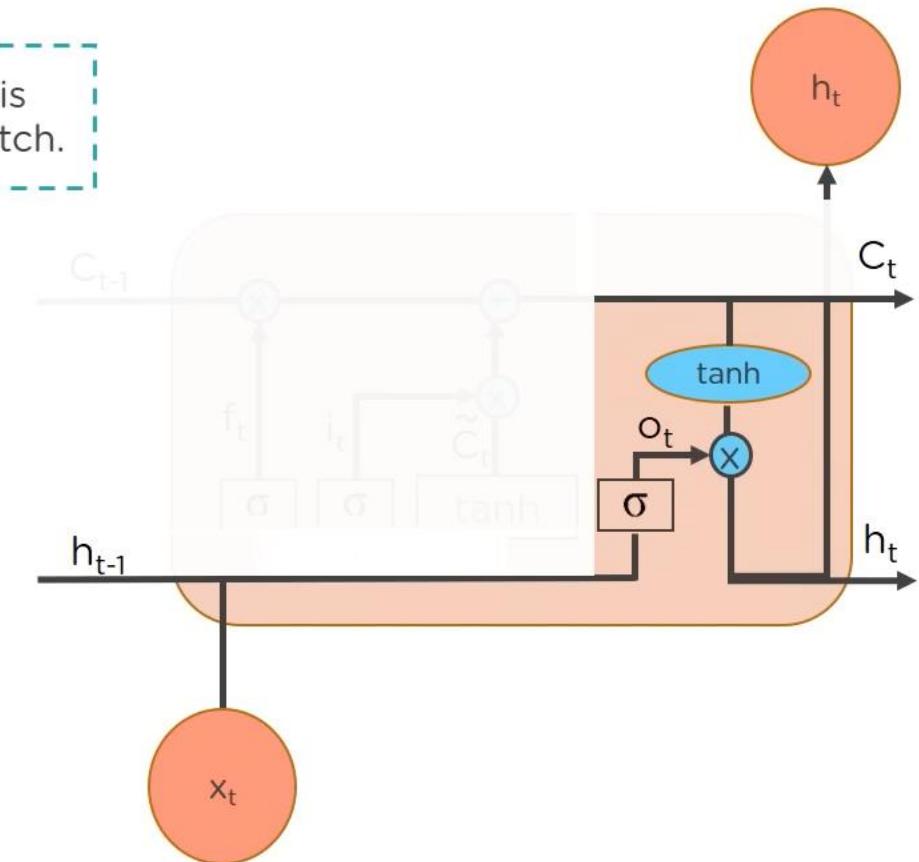
# Working of LSTMs

Let's consider this example to predicting the next word in the sentence:

John played tremendously well against the opponent and won for his team. For his contributions, brave \_\_\_\_\_ was awarded player of the match.



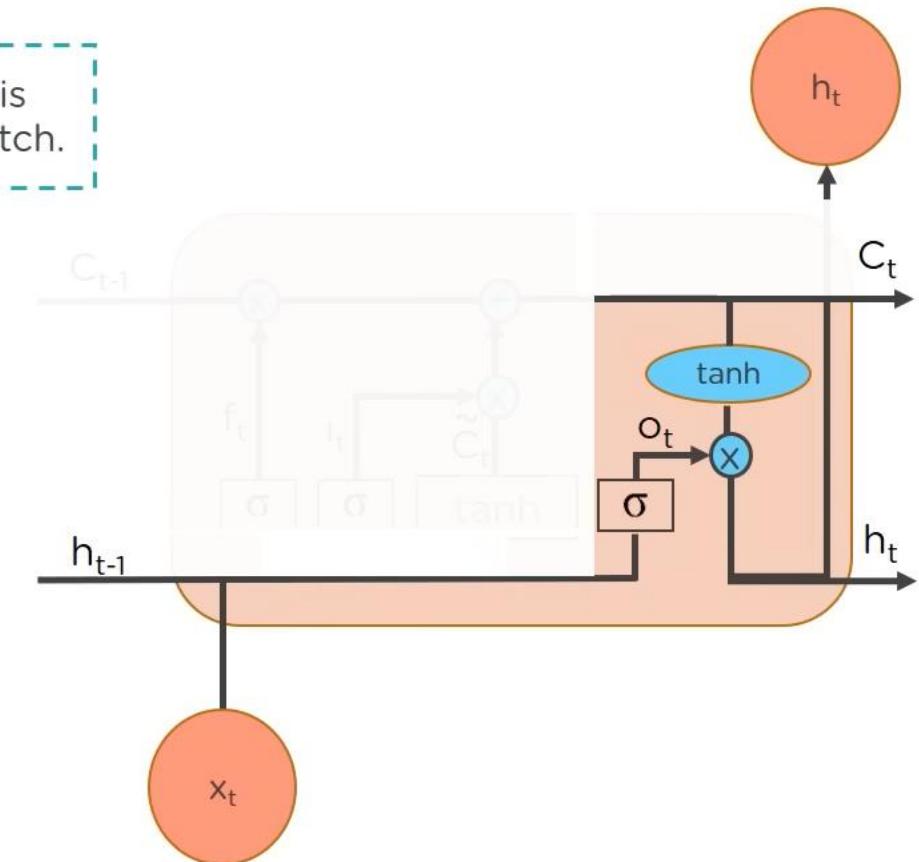
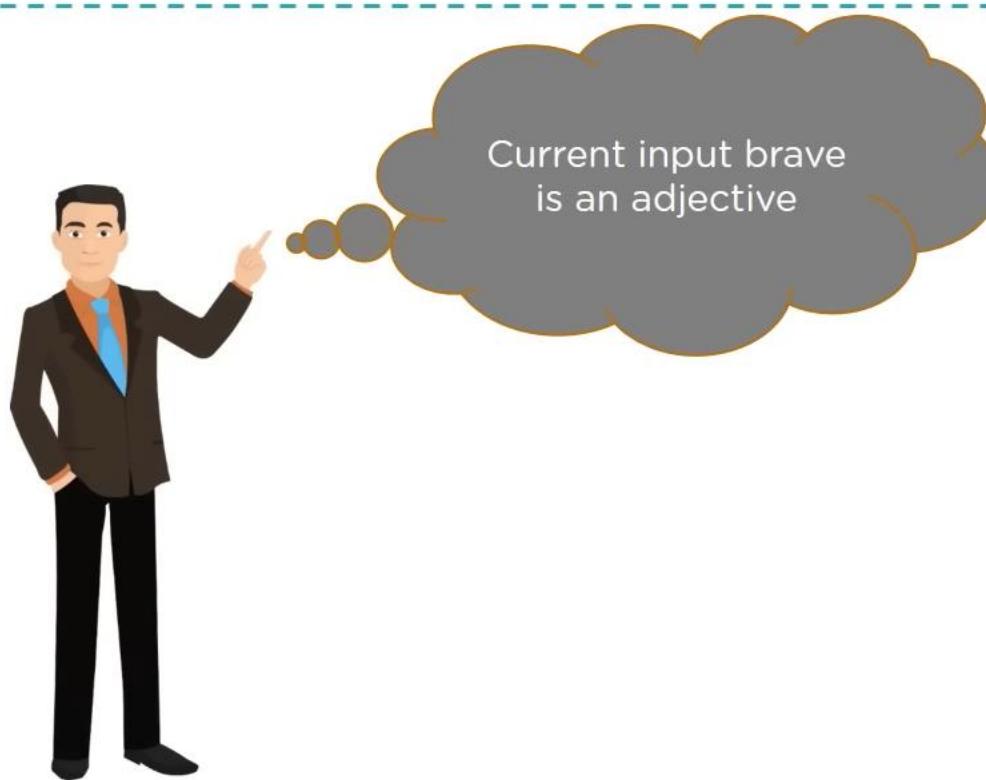
There could be a lot of choices for the empty space



# Working of LSTMs

Let's consider this example to predicting the next word in the sentence:

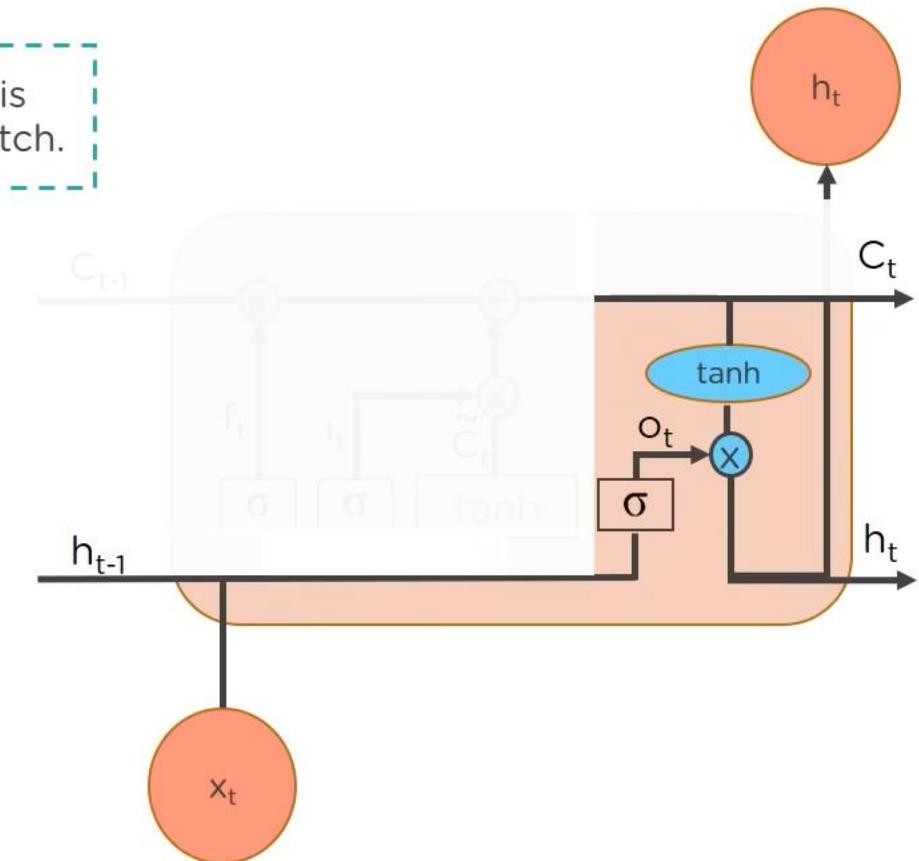
John played tremendously well against the opponent and won for his team. For his contributions, brave \_\_\_\_\_ was awarded player of the match.



# Working of LSTMs

Let's consider this example to predicting the next word in the sentence:

John played tremendously well against the opponent and won for his team. For his contributions, brave \_\_\_\_\_ was awarded player of the match.



# Working of LSTMs

Let's consider this example to predicting the next word in the sentence:

John played tremendously well against the opponent and won for his team. For his contributions, brave \_\_\_\_\_ was awarded player of the match.

