

الشبكات العصبية

محاضرة 3

# Recurrent Neural Networks

د. فادي متوج

# Language model and sequence generation

# What is language modelling?



## Speech recognition

The apple and pair salad. The

apple and pear salad.

$$P(\text{The apple and pair salad}) =$$

$$3.2 \times 10^{-13}$$

$$P(\text{The apple and pear salad}) =$$

$$5.7 \times 10^{-10}$$

$$P(\text{Sentence}) = ?$$

$$P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$$

# Language modelling with an RNN



Training set: large corpus of english text.

Tokenize

Cats average 15 hours of sleep a day.

↓ <EOS>

$y^{(1)}$   $y^{(2)}$   $y^{(3)}$  ...  $y^{(8)}$   $y^{(9)}$

$$x^{(t)} = y^{(t-1)}$$

The Egyptian Mau is a bread of cat. <EOS>

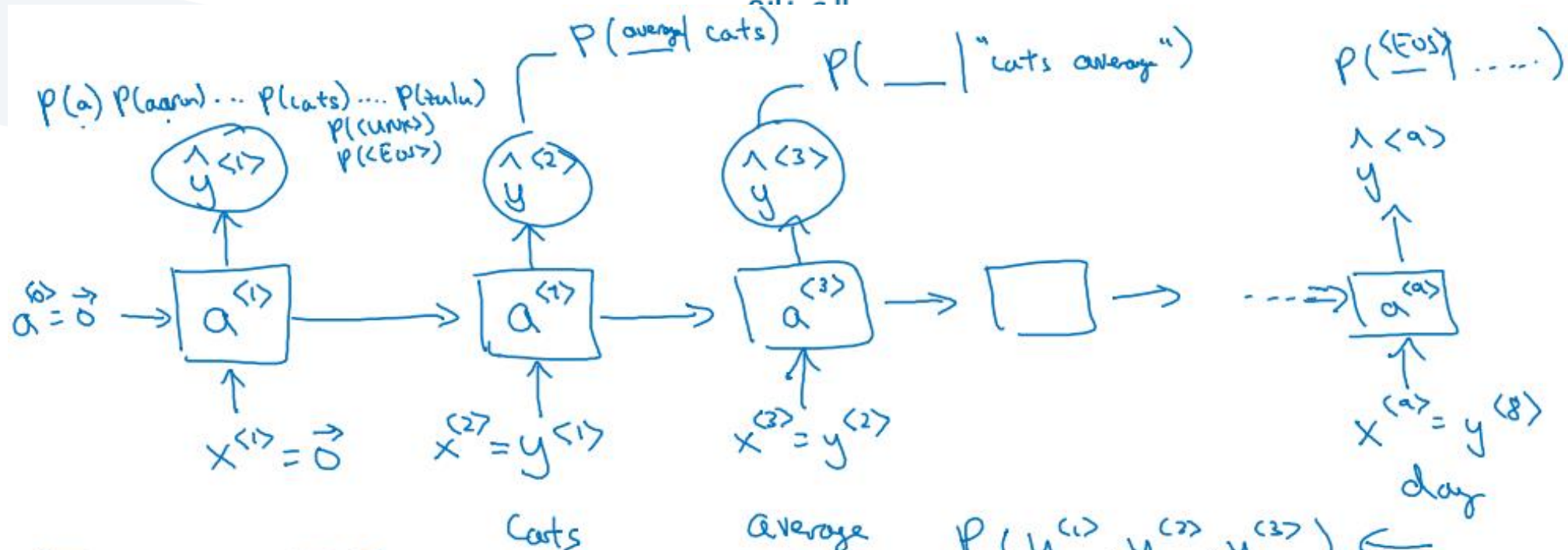
10,000

<UNK>

# RNN model



جامعة  
منارة



→ Cats average 15 hours of sleep a day. <EOS>

$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$\mathcal{L} = \sum \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

$$P(y^{<1>}, y^{<2>}, y^{<3>}) \leftarrow$$

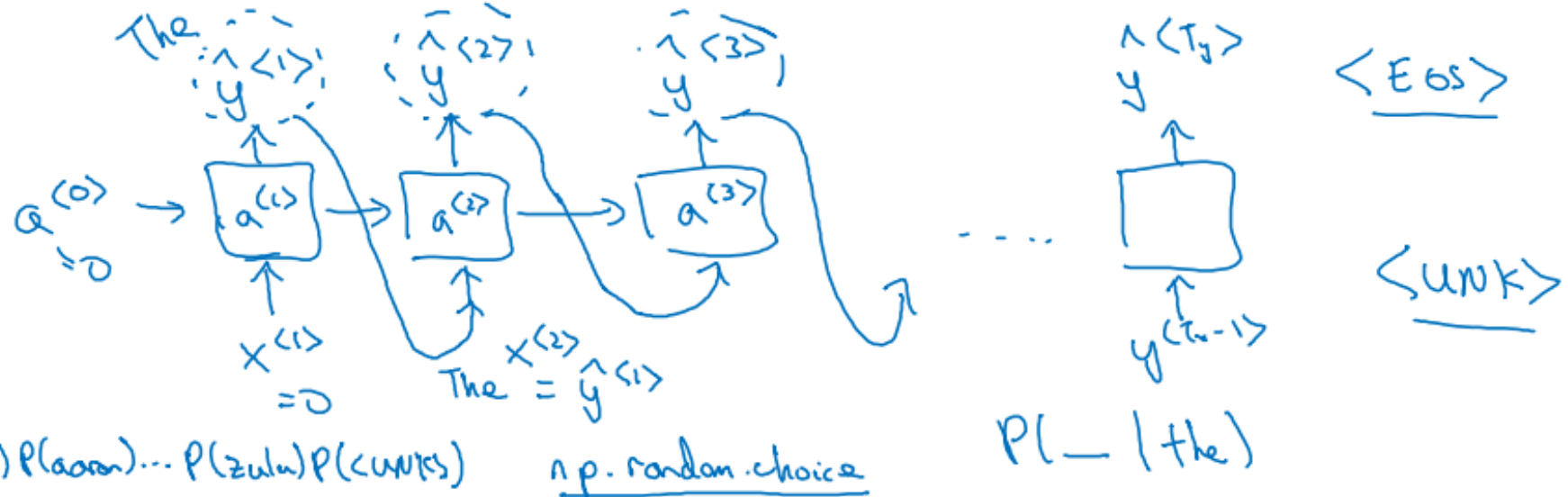
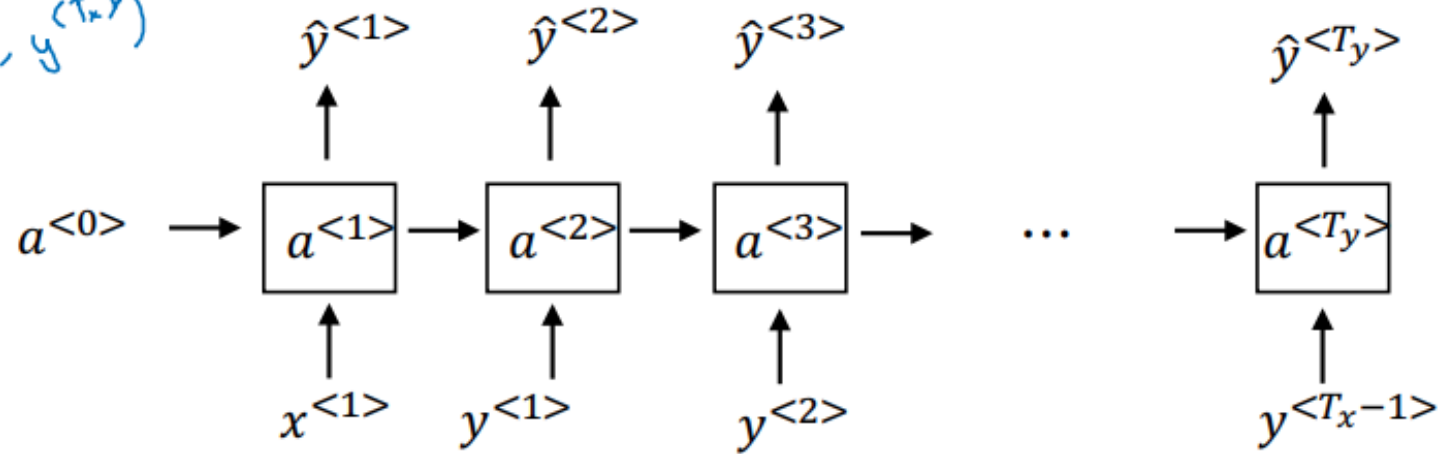
$$= \frac{P(y^{<1>}) P(y^{<2>} | y^{<1>})}{P(y^{<3>} | y^{<1>}, y^{<2>})}$$

# Sampling novel sequences

# Sampling a sequence from a trained RNN

$P(y^{(1)}, \dots, y^{(T_x)})$

Training:



# Character-level language model

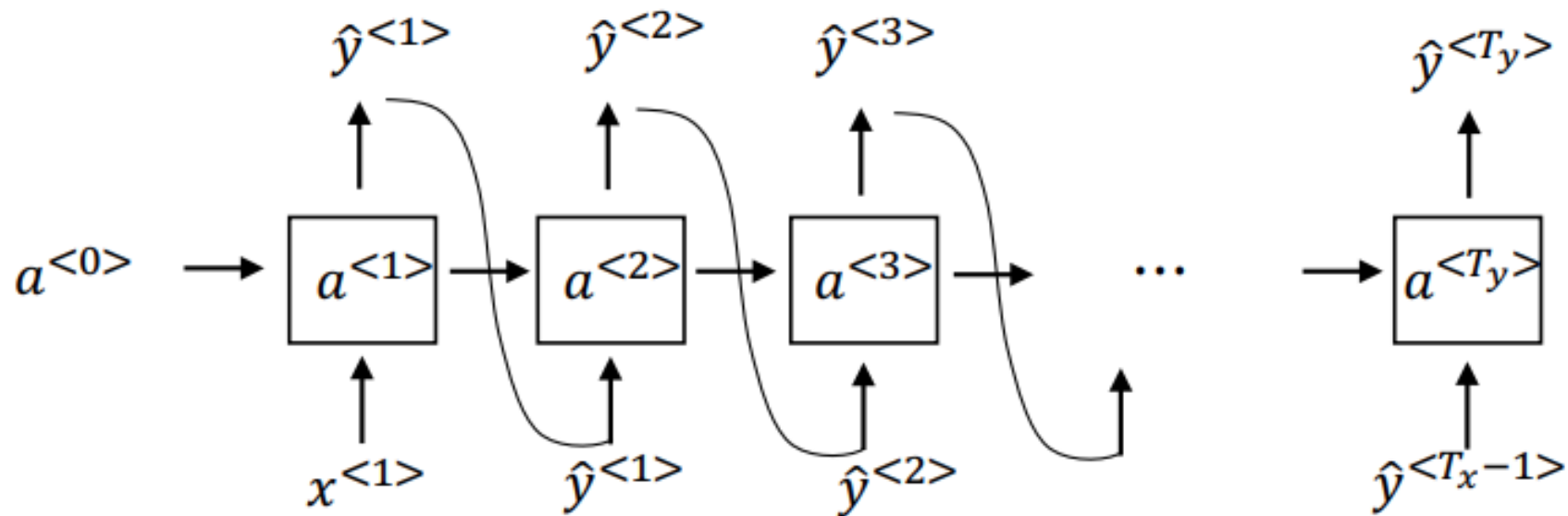
Vocabulary = [a, aaron, ..., zulu, <UNK>].

→ Vocabulary = [a, b, c, ..., z, ٠, ١, ٢, ٣, ٤, ٥, ٦, ٧, ٨, ٩, A, ..., Z]

$y^{(1)}$   $y^{(2)}$   $y^{(3)}$   $y^{(4)}$

Cat  
↑↑↑↑  
Average  
...

Man





# Recurrent Neural Networks



## News

President Enrique Peña Nieto, announced  
sench's sulk former coming football langston  
paring.

"I was not at all surprised," said high langston.

"Concussion epidemic", to be examined.

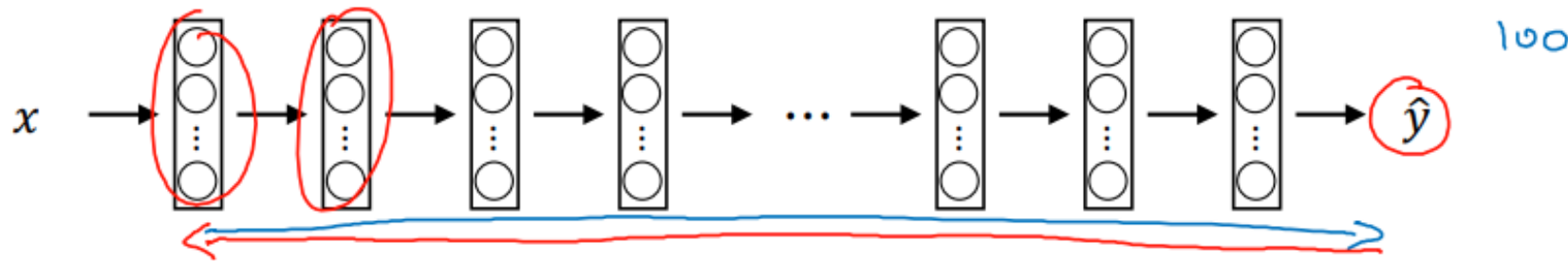
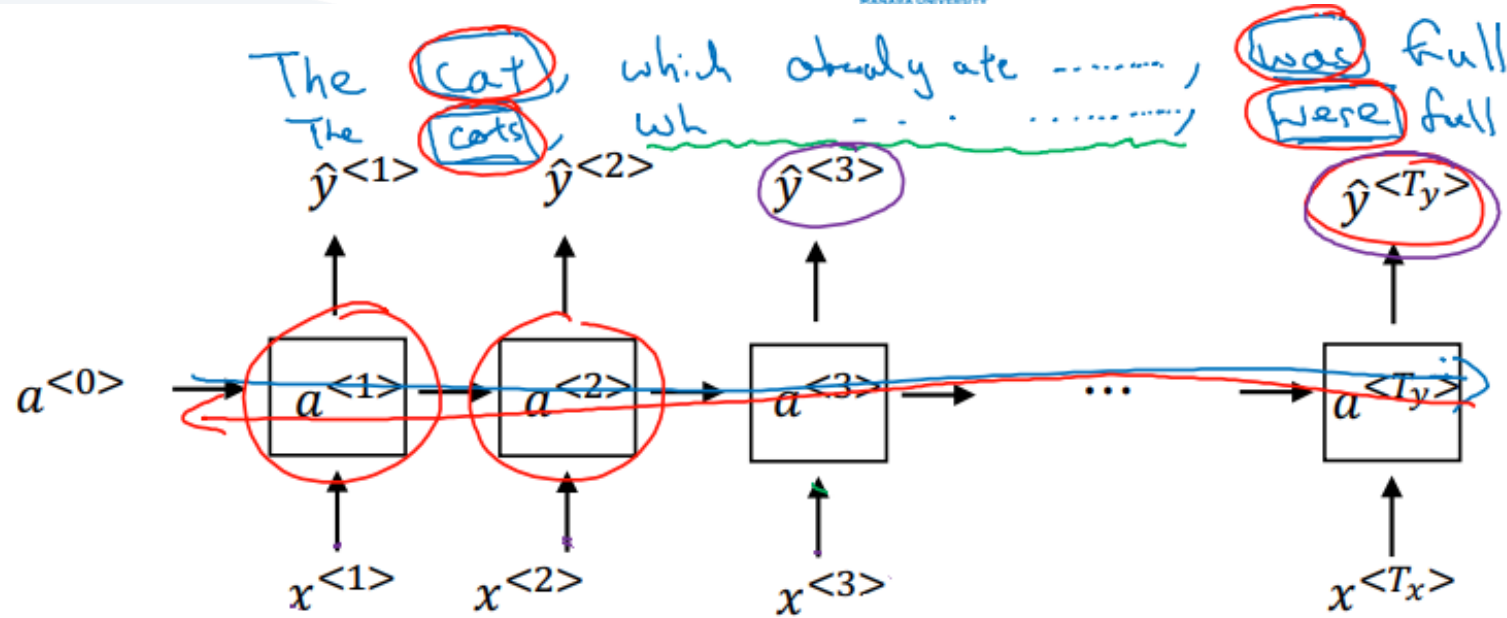
The gray football the told some and this has on  
the uefa icon, should money as.

## Shakespeare

The mortal moon hath her eclipse in love.  
And subject of this thou art another this fold.  
When better be my love to me see sabl's.  
For whose are ruse of mine eyes heaves.

# Vanishing gradients with RNNs

# Vanishing gradients with RNNs



Exploding gradients.

Gradient clipping

# LSTM (long short term memory) unit

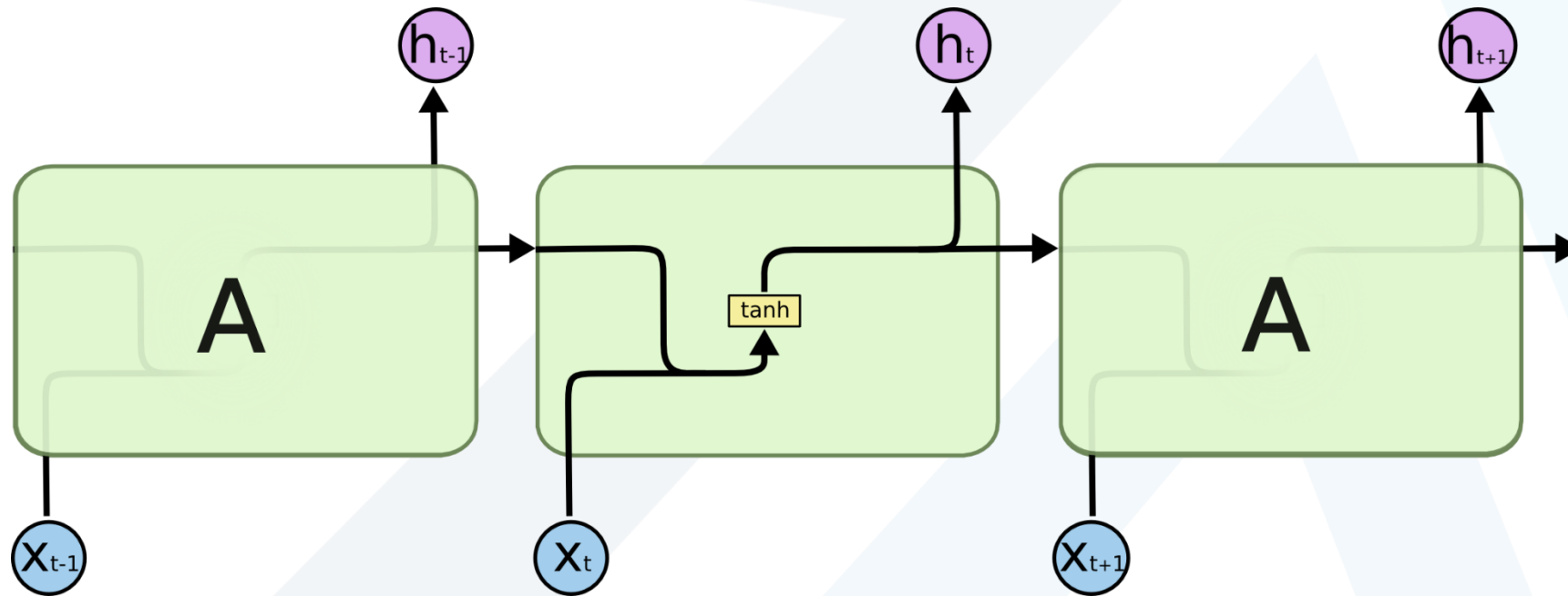
## Simplified RNN notation



- Consider trying to predict the last word in the text “I grew up in France... I speak fluent French.” Recent information suggests that the next word is probably the name of a language, but if we want to narrow down which language, we need the context of France, from further back. It’s entirely possible for the gap between the relevant information and the point where it is needed to become very large.
- Unfortunately, as that gap grows, RNNs become unable to learn to connect the information.
- Long Short Term Memory networks – usually just called “LSTMs” – are a special kind of RNN, capable of learning long-term dependencies.

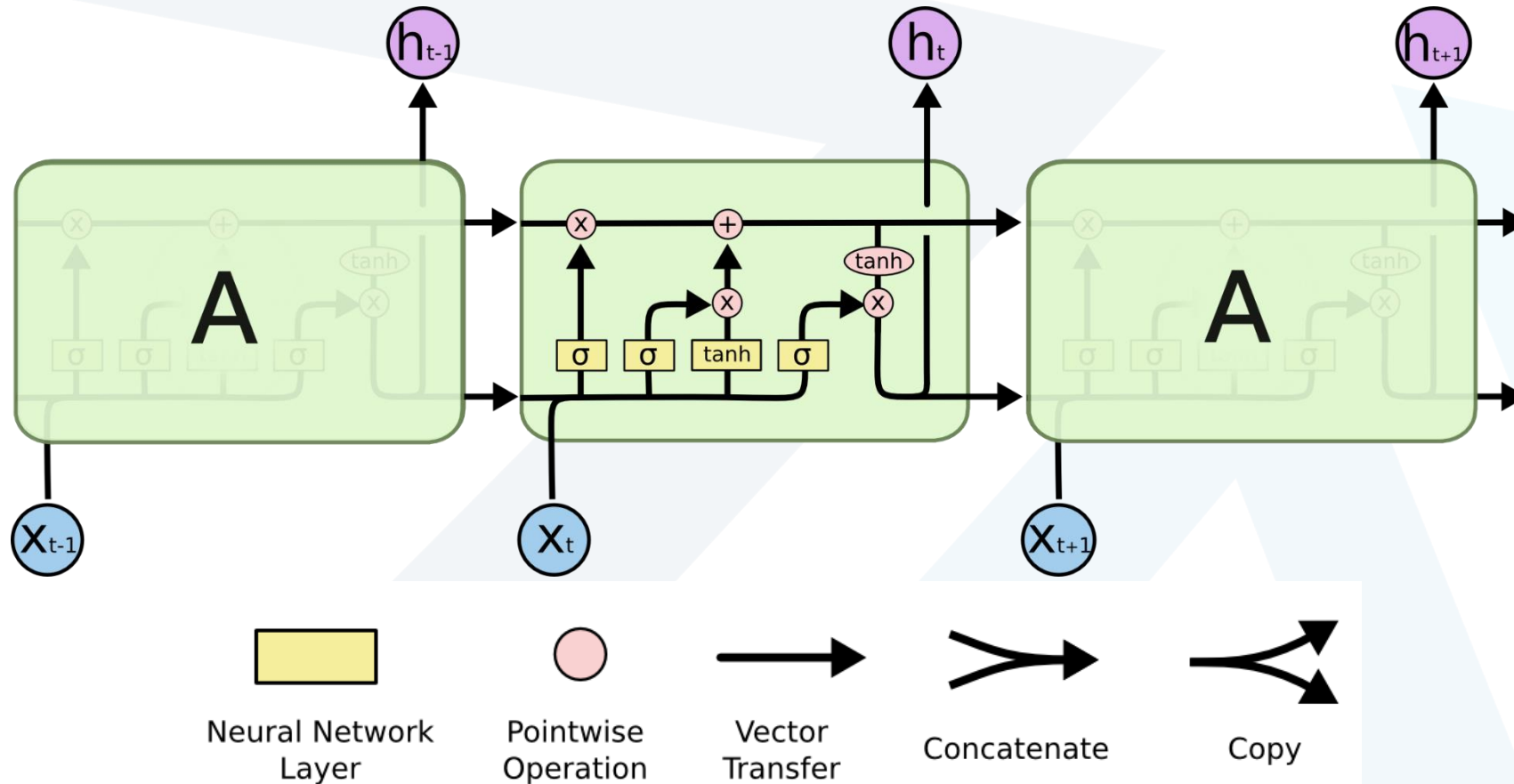
## LSTMs structure

- 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.



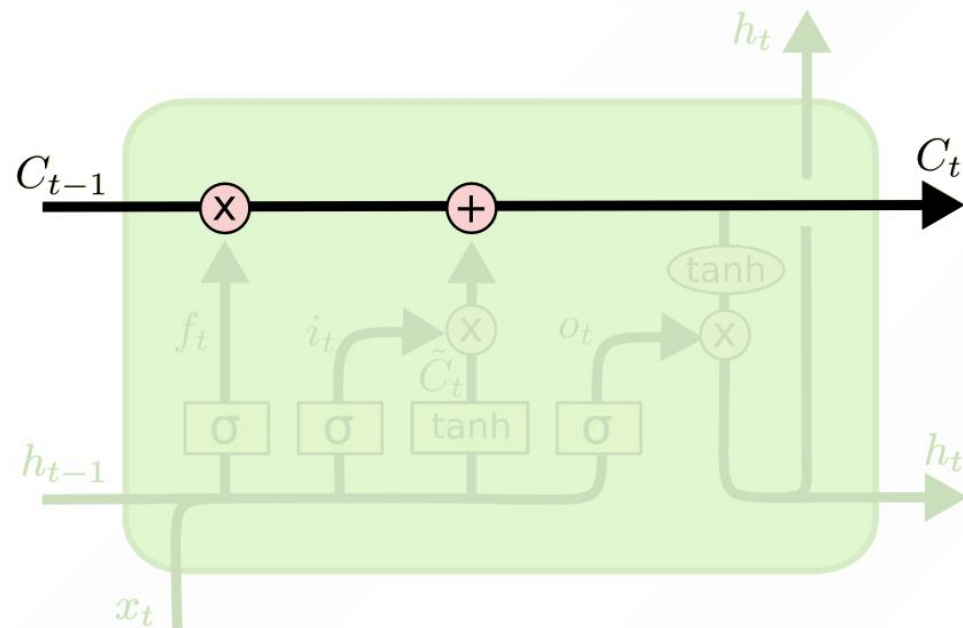
# LSTMs structure

- 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.



# The Core Idea Behind LSTMs

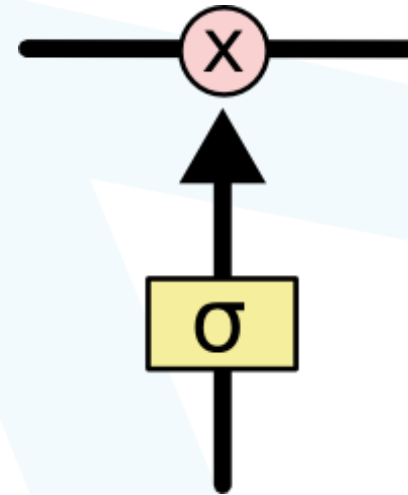
- The key to LSTMs is the cell state, the horizontal line running through the top of the diagram.
- The cell state is kind of like a conveyor belt. It runs straight down the entire chain. It's very easy for information to just flow along it unchanged.





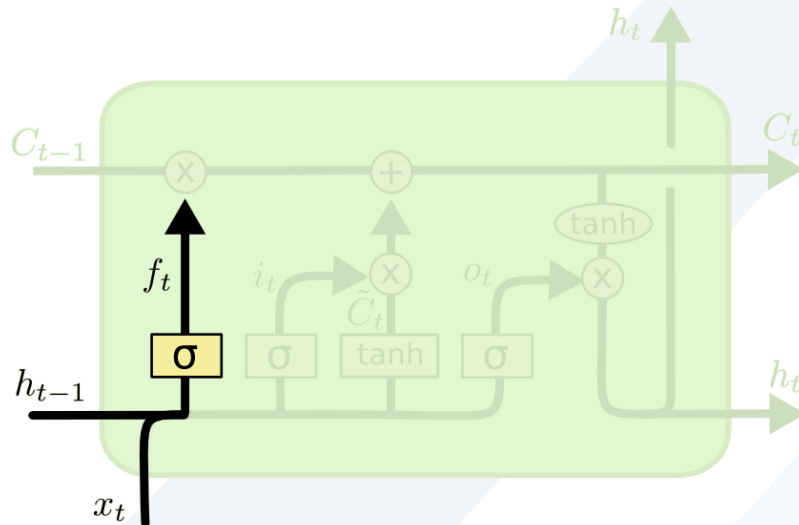
# The Core Idea Behind LSTMs

- The LSTM does have the ability to remove or add information to the cell state, carefully regulated by structures called gates.
- Gates are a way to optionally let information through. They are composed out of a sigmoid neural net layer and a pointwise multiplication operation.
- The sigmoid layer outputs numbers between zero and one, describing how much of each component should be let through. A value of zero means “let nothing through,” while a value of one means “let everything through!”
- An LSTM has three of these gates, to protect and control the cell state.



# Step-by-Step LSTM Walk Through

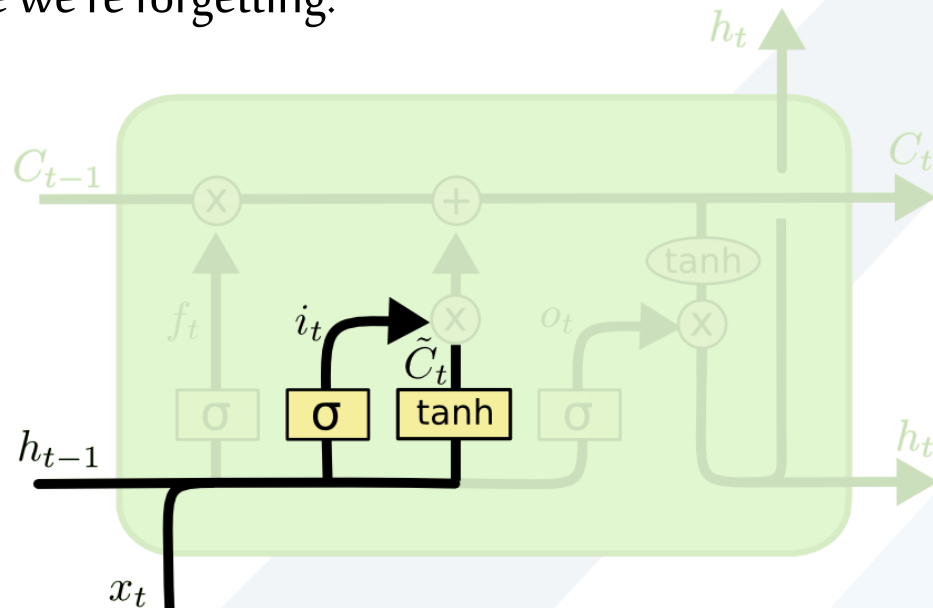
- The first step in our LSTM is to decide what information we're going to throw away from the cell state. This decision is made by a sigmoid layer called the "forget gate layer." It looks at  $h_{t-1}$  and  $x_t$ , and outputs a number between 0 and 1 for each number in the cell state  $C_{t-1}$ . A 1 represents "completely keep this" while a 0 represents "completely get rid of this."
- Let's take an example of a language model trying to predict the next word based on all the previous ones. In such a problem, the cell state might include the gender of the present subject, so that the correct pronouns can be used. When we see a new subject, we want to forget the gender of the old subject.



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

# Step-by-Step LSTM Walk Through

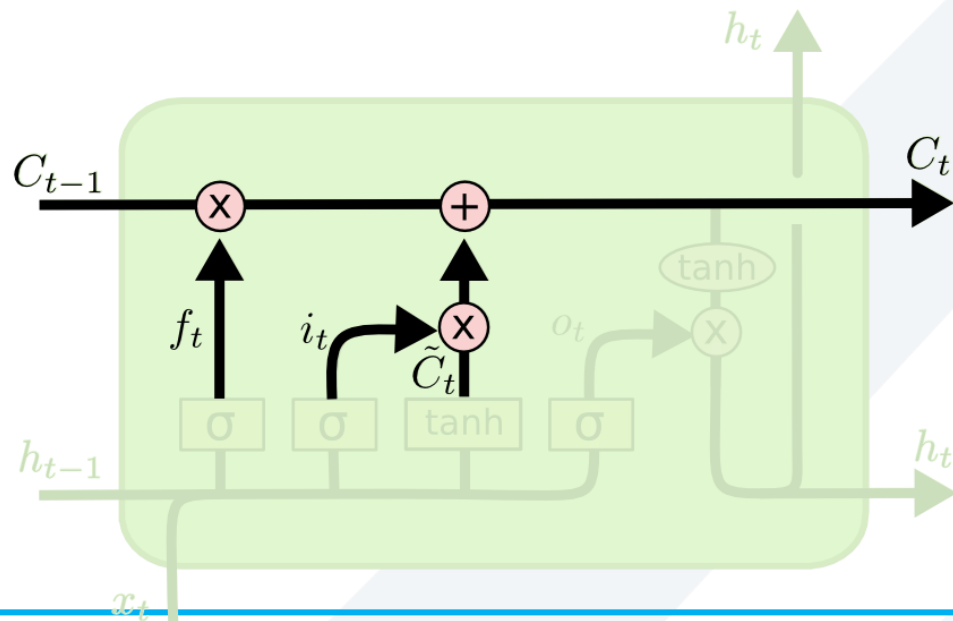
- The next step is to decide what new information we're going to store in the cell state. This has two parts. First, a sigmoid layer called the "input gate layer" decides which values we'll update. Next, a tanh layer creates a vector of new candidate values,  $\tilde{C}_t$ , that could be added to the state. In the next step, we'll combine these two to create an update to the state.
- In the example of our language model, we'd want to add the gender of the new subject to the cell state, to replace the old one we're forgetting.



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

# Step-by-Step LSTM Walk Through

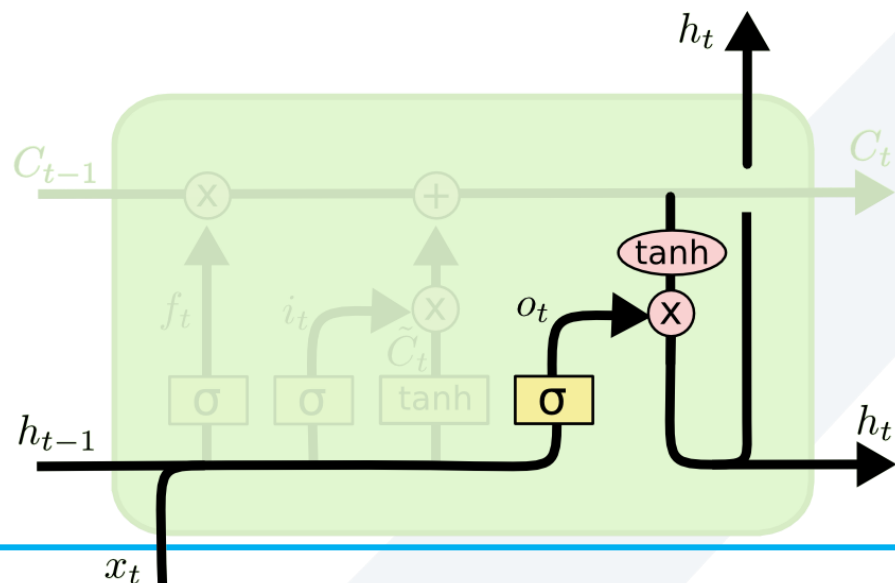
- It's now time to update the old cell state,  $C_{t-1}$ , into the new cell state  $C_t$ . The previous steps already decided what to do, we just need to actually do it.
- We multiply the old state by  $f_t$ , forgetting the things we decided to forget earlier. Then we add  $i_t * \tilde{C}_t$ . This is the new candidate values, scaled by how much we decided to update each state value.
- In the case of the language model, this is where we'd actually drop the information about the old subject's gender and add the new information, as we decided in the previous steps.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# Step-by-Step LSTM Walk Through

- Finally, we need to decide what we're going to output. This output will be based on our cell state, but will be a filtered version. First, we run a sigmoid layer which decides what parts of the cell state we're going to 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, so that we only output the parts we decided to.
- For the language model example, since it just saw a subject, it might want to output information relevant to a verb, in case that's what is coming next. For example, it might output whether the subject is singular or plural, so that we know what form a verb should be conjugated into if that's what follows next.

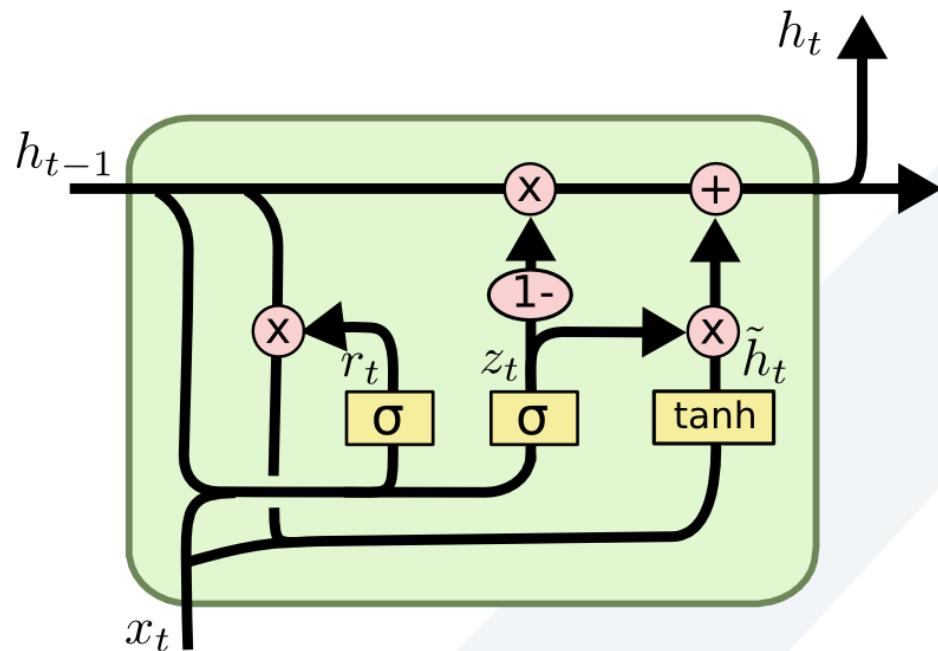


$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

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

# Gated Recurrent Unit(GRU)

- A variation on the LSTM is the **Gated Recurrent Unit**, or **GRU**. It combines the forget and input gates into a single “**update gate**.” It also merges the cell state and hidden state, and makes some other changes. The resulting model is simpler than standard LSTM models, and has been growing increasingly popular.



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

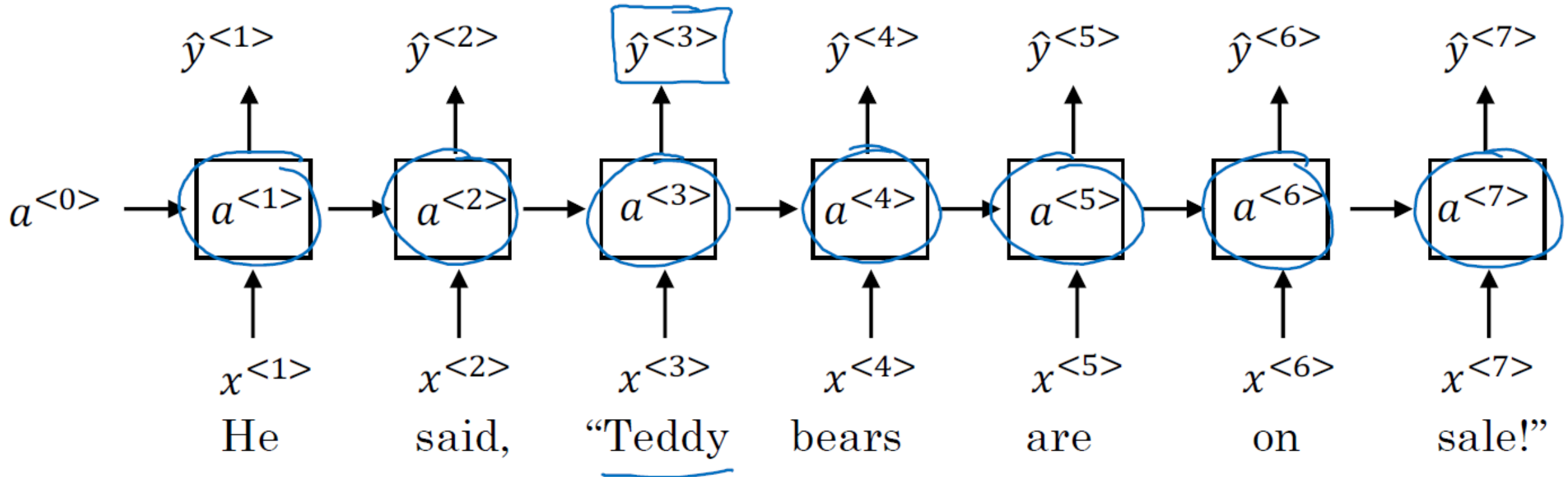
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Bidirectional RNN

# Getting information from the future

He said, "Teddy bears are on sale!"

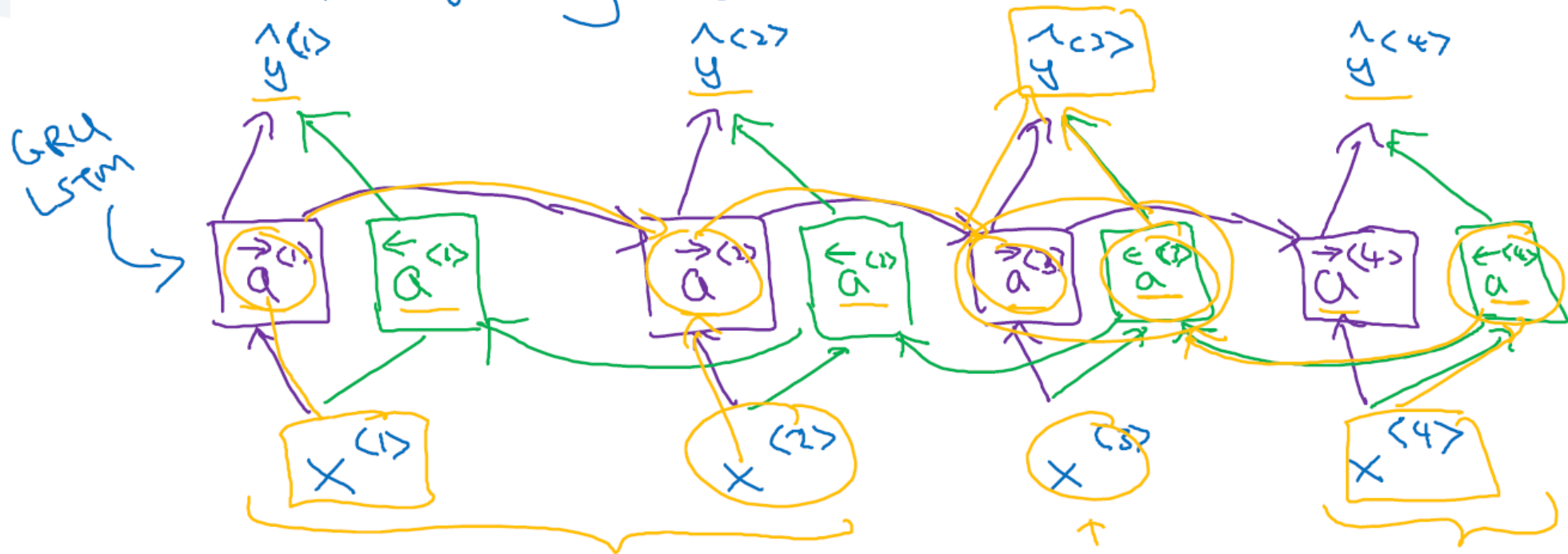
He said, "Teddy Roosevelt was a great President!"





# Bidirectional RNN (BRNN)

$$\hat{y}^{(t)} = g(W_y [a^{\rightarrow(t)}, a^{\leftarrow(t)}] + b_y)$$



Acyclic graph

He said "Teddy Roosevelt ..."

BRNN w/ LSTM

# Deep RNNs

# Deep RNN example

