



deeplearning.ai

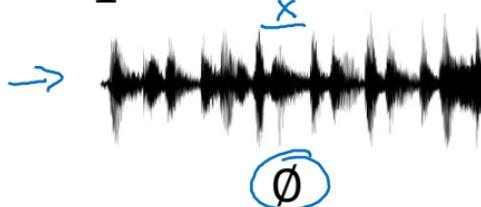
# Recurrent Neural Networks

---

## Why sequence models?

# Examples of sequence data

Speech recognition



"The quick brown fox jumped  
over the lazy dog."

Music generation



Sentiment classification

"There is nothing to like  
in this movie."



DNA sequence analysis → AGCCCCCTGTGAGGAACTAG

AGCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter avec  
moi?



Do you want to sing with  
me?

Video activity recognition

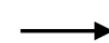


Running

Name entity recognition



Yesterday, Harry Potter  
met Hermione Granger.



Yesterday, **Harry Potter**  
met **Hermione Granger**.

Andrew Ng

## Notation used in Sequence Models

An motivating example

Goal: Find persons' name in the sentence below. (Name Entity Recognition)

Input  $X$ : Harry Potter and Hermione Granger invented a new spell.  
 $T_x = 9$        $X^{(1)}$      $X^{(2)}$      $X^{(3)}$     ...       $X^{(t)}$     ...       $X^{(9)}$

Desired output  $y$ : 1 1 0 1 1 0 0 0 0  
 $T_y = 9$        $y^{(1)}$      $y^{(2)}$      $y^{(3)}$     ...       $y^{(t)}$     ...       $y^{(9)}$

The  $\langle t \rangle$  is used to index the position of a word in a sequence.

$T_x$ : the length of the input sequence.  $T_y$ : the length of the output sequence.

$T_x$  and  $T_y$  are not necessarily the same value.

$X^{(i)\langle t \rangle}$ : the  $t$ -th element in the input sequence of the training sample  $i$ .

$y^{(i)\langle t \rangle}$ : the  $t$ -th element in the output sequence of the training sample  $i$ .

$T_x^{(i)}$ : the input sequence length for the training sample  $i$ .

$T_y^{(i)}$ : the output sequence length for the training sample  $i$ .

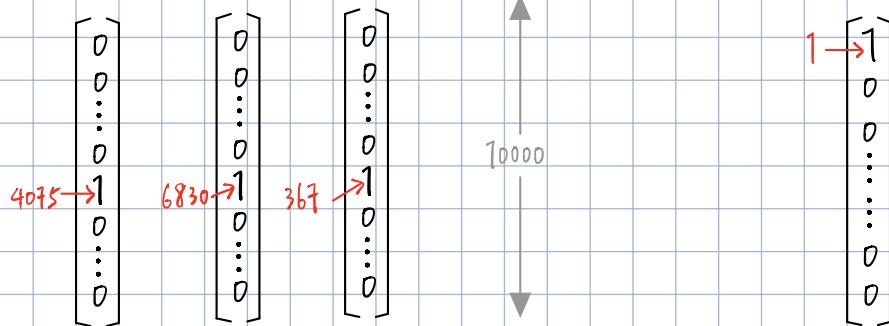
Vocabulary: Build a list with 10000 most used English words. What if encountering a word

a	1	that is not in the vocabulary. I can create a new token <UNK>.
aaron	2	
:		
and	367	
:		
harry	4075	
:		
potter	6830	
:		
zulu	10,000	

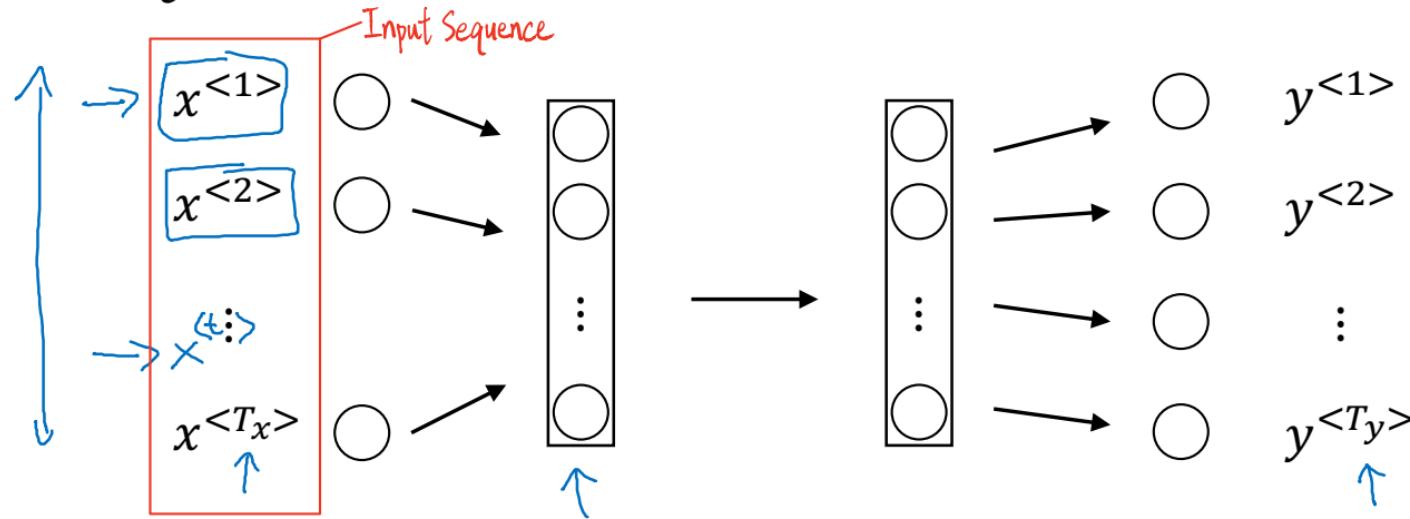
## Representing Words

Input  $X$ : Harry Potter and Hermione Granger invented a new spell.  
 $X^{(1)}$      $X^{(2)}$      $X^{(3)}$     ...     $X^{(t)}$     ...     $X^{(n)}$

One-Hot representation:



# Why not a standard network?



Problems: Every sentence has different lengths.

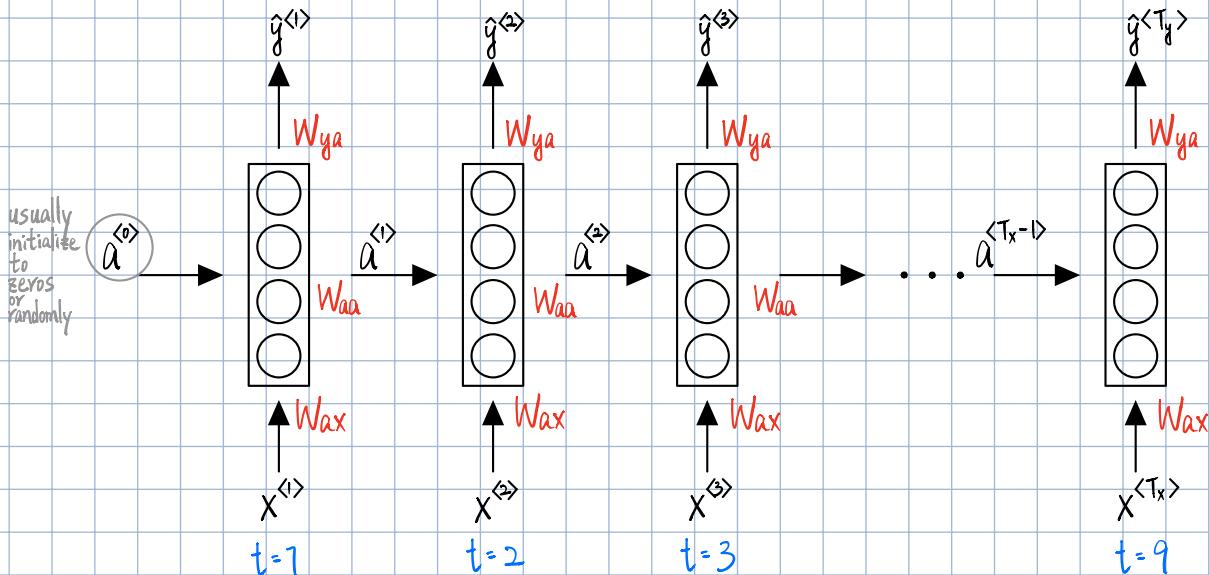
- - Inputs, outputs can be different lengths in different examples.
- - Doesn't share features learned across different positions of text.

# Recurrent Neural Network Model

Input  $X$ : Harry Potter and Hermione Granger invented a new spell.  $T_x = T_y$

$$X^{(1)} \quad X^{(2)} \quad X^{(3)} \quad \dots \quad X^{(t)} \quad \dots \quad X^{(8)} \quad X^{(9)}$$

Previously,  $X^{(i)}$  can be seen as the  $t$ -th element in the input sequence of the training sample  $i$ . Here, it can also be seen as the time step  $t$ .



$a^{(1)}$ : the activation value of the time step 1.  $a^{(2)}$ : the activation value of the time step 2.

At each time step  $t$ , the RNN passes on activation to the next time step for it to use.

The RNN scans through the data from left to right. The parameters,  $W_{ax}$ ,  $W_{aa}$ ,  $W_{ya}$ , it uses for each time step are shared.

In this RNN, when making the prediction for  $\hat{y}^{(3)}$ , it gets the information not only from  $X^{(3)}$  but also from  $X^{(1)}$  and  $X^{(2)}$ .

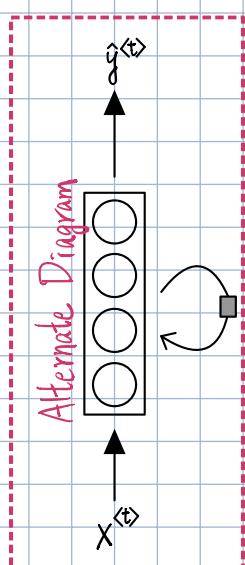
**Weakness:** This RNN only uses the information that is earlier in the sequence to make a prediction. In particular, when predicting  $\hat{y}^{(3)}$ , it doesn't use information about words  $X^{(4)}, X^{(5)}$  and so on.

E.g.:

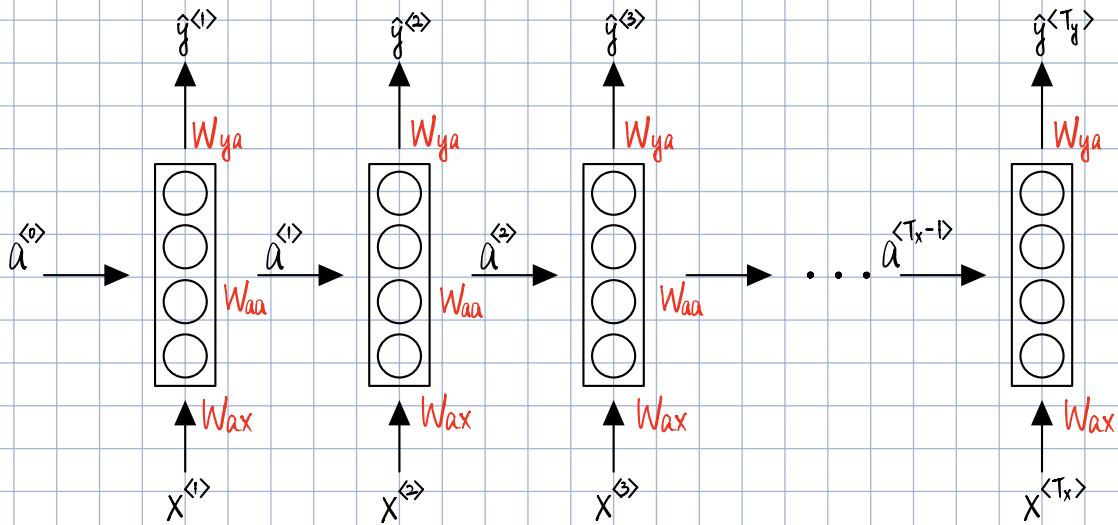
He said, "Teddy Roosevelt was a great president."

He said, "Teddy bears are on sale."

Using unidirectional RNN is just impossible to say if Teddy is a name or not.



## Forward Propagation



$$a^{(0)} = \vec{0} \quad a^{(t)} = f_1(W_{aa}a^{(t-1)} + W_{ax}x^{(t)} + b_a) \quad \text{The activation function } f_1 \text{ is often tanh or ReLU}$$

$$\hat{y}^{(t)} = f_2(W_{ya}a^{(t)} + b_y)$$

$f_2$ : sigmoid for binary class. (Name Entity recognition)  
softmax for k class.

$$a^{(t)} = f_1(W_{aa}a^{(t-1)} + W_{ax}x^{(t)} + b_a)$$

$$\hat{y}^{(t)} = f_2(W_{ya}a^{(t)} + b_y)$$

$$\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = W_{aa}a^{(t-1)} + W_{ax}x^{(t)}$$

Simplifying RNN notation

$$a^{(t)} = f_1(W_a[a^{(t-1)}, x^{(t)}] + b_a)$$

Stack  $W_{aa}$  and  $W_{ax}$  horizontally  $[W_{aa} \mid W_{ax}] = W_a$

$$\hat{y}^{(t)} = f_2(W_ya^{(t)} + b_y)$$

$[a^{(t-1)}, x^{(t)}]$  means stack vertically

E.g.:  $a^{(t-1)}$  is  $(100, .)$ .  $x^{(t)}$  is  $(10000, .)$

$$\begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}$$

Then  $W_a$  is  $(100, 100)$  and  $W_x$  is  $(100, 10000)$

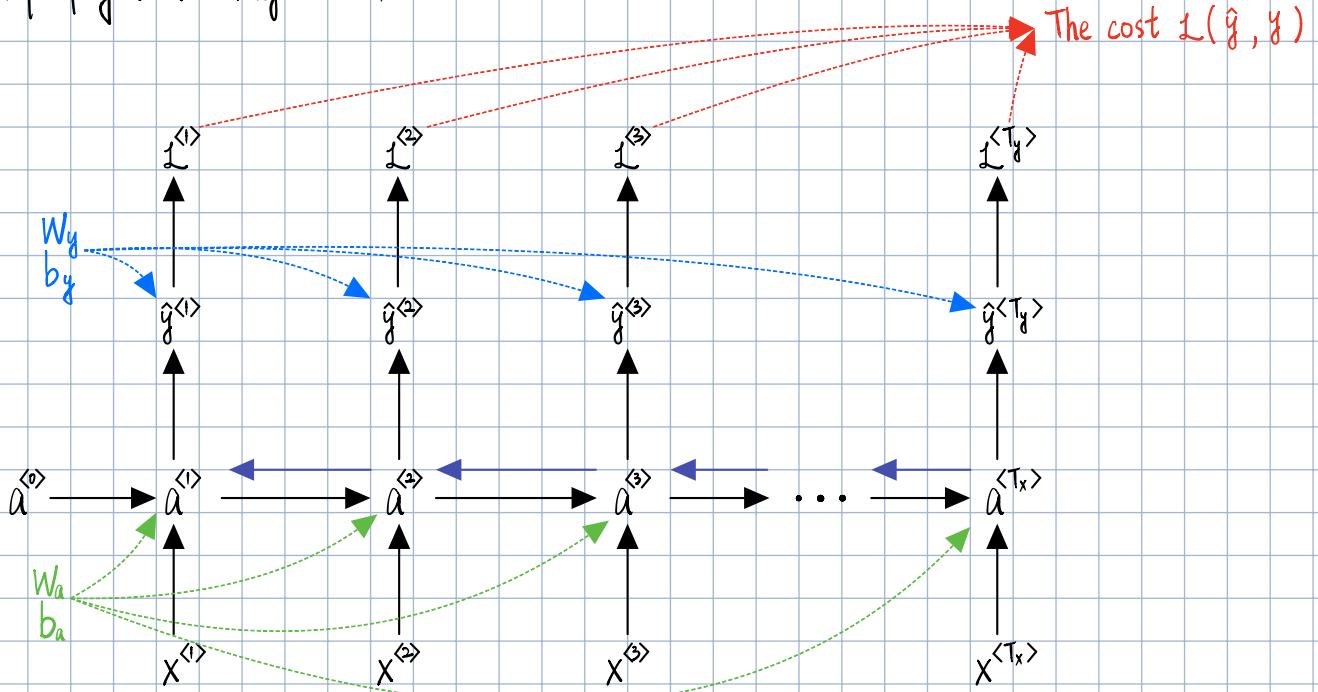
$$100 \begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} = W_a$$

100      10000      (100, 100)

$$[a^{(t-1)}, x^{(t)}] = \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix}$$

$\uparrow$   
100  
 $\downarrow$   
10000

## Backpropagation through Time



$$L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

$$\text{The cost } L(\hat{y}, y) = \sum_{t=1}^{T_x} L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

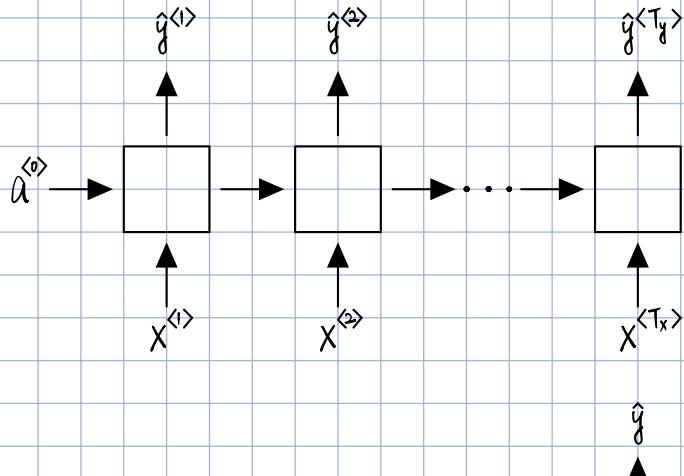
Backpropagation through Time

# Different Types of RNNs

Inspired by Andrej Karpathy "The Unreasonable Effectiveness of Recurrent Neural Networks"

Many-to-many.  $T_x = T_y$

Eg: Name Entity Recognition

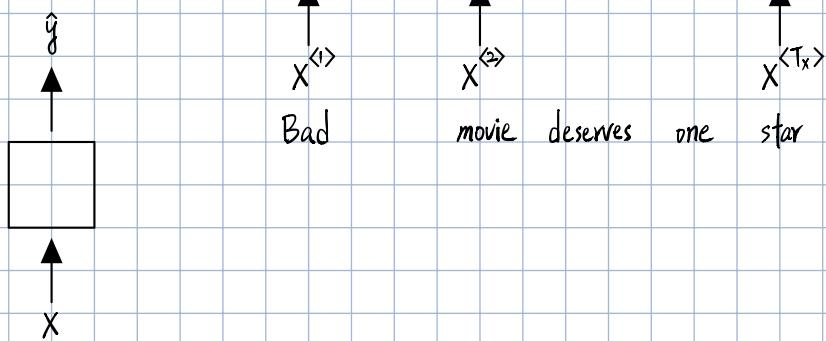


Many-to-one.  $T_x = T_y$

Eg: Sentiment Classification

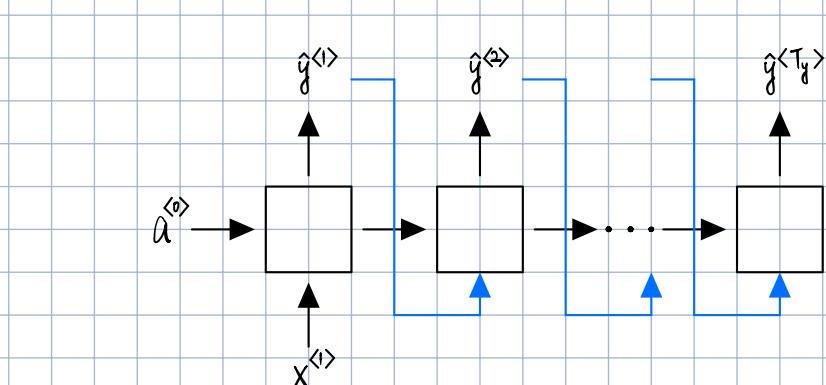
X: text

y: 1...5



One-to-one.

The standard neural network.



One-to-many .

Eg: Music generation

X: genre, p

y: music

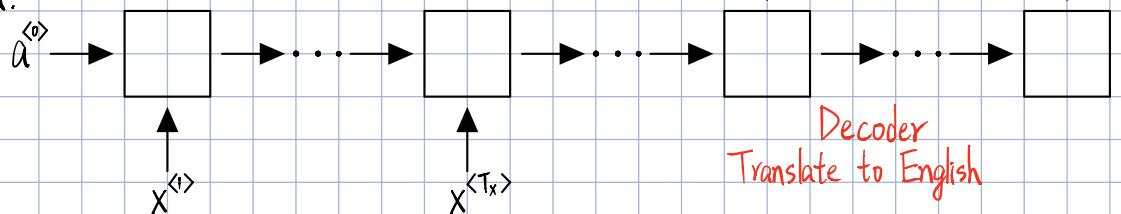
Many-to-many.  $T_x \neq T_y$

Encoder  
Read German as input

Eg: Machine translation.

X: German

y: English



# Language Model and Sequence Generation

What is language model?

The output of speech recognition:

The apple and pair salad.

Which sentence has higher probability?  $P(\text{sentence})$

The apple and pear salad.

$$P(\text{The apple and pair salad.}) = 3.2 \times 10^{-13}$$

$$P(\text{The apple and pear salad.}) = 5.7 \times 10^{-10}$$

Given a sentence, a sequence of words  $y^{(1)}, y^{(2)}, \dots, y^{(T_k)}$ , a language model will estimate the probability of the particular sequence of words.  $P(y^{(1)}, y^{(2)}, \dots, y^{(T_k)})$

Now, the question is how to build a language model? Let's build a language model with RNN.

Training set: large corpus of English text

Give a sentence: Cats average 15 hours of sleep a day. <EOS>

Tokenize :  $y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(8)} \quad y^{(9)}$

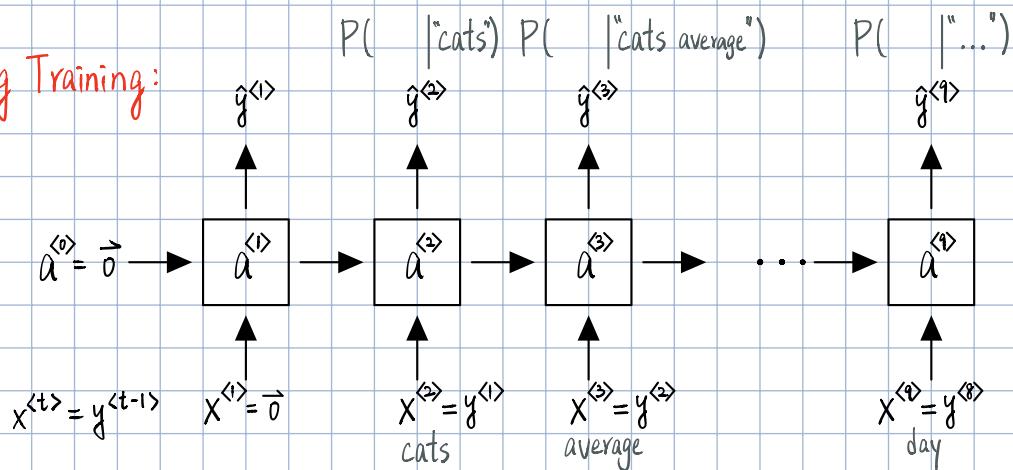
End of Sentence

<UNK> unknown word. Not in vocabulary.  
The Egyptian Mau is a breed of cat. <EOS>

RNN model:

Using "Cats average 15 hours of sleep a day. <EOS>" as training sample.

During Training:



$\hat{y}^{(1)}$ : the prob. of the first prediction being any one of the word in the vocabulary.

$\hat{y}^{(2)}$ : the prob. of the 2nd prediction being any one of the word in the vocabulary given cats.

$$\mathcal{L}(\hat{y}^{(t)}, y^{(t)}) = -\sum_i y_i^{(t)} \log \hat{y}_i^{(t)} \quad \text{softmax loss function}$$

$$\text{The cost } \mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}(\hat{y}^{(t)}, y^{(t)})$$

During Inference:

$$y^{(1)} \quad y^{(2)} \quad y^{(3)}$$

Given first three words, "cats", "average", "15", it can predict the next word.

$$P(y^{(1)}, y^{(2)}, y^{(3)}) = P(y^{(1)}) P(y^{(2)} | y^{(1)}) P(y^{(3)} | y^{(1)}, y^{(2)})$$

It calculates the probability of "cats average 15" being a sentence.



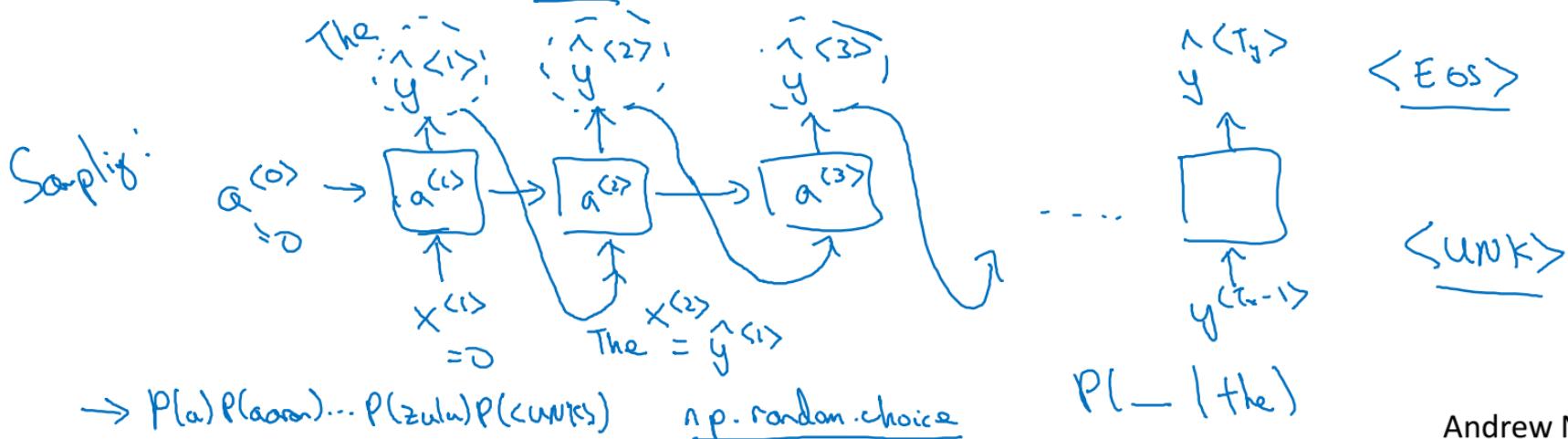
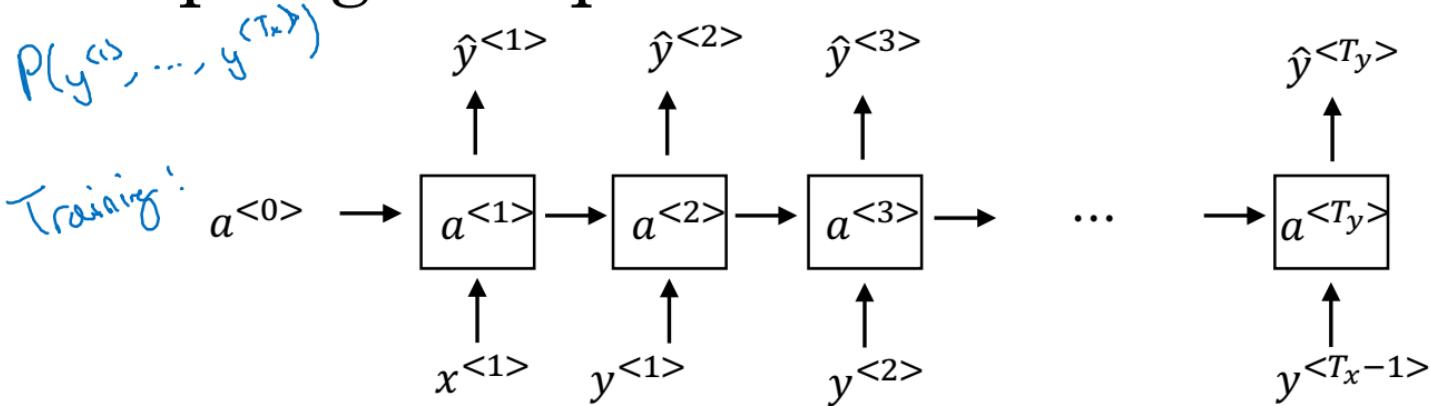
deeplearning.ai

# Recurrent Neural Networks

---

## Sampling novel sequences

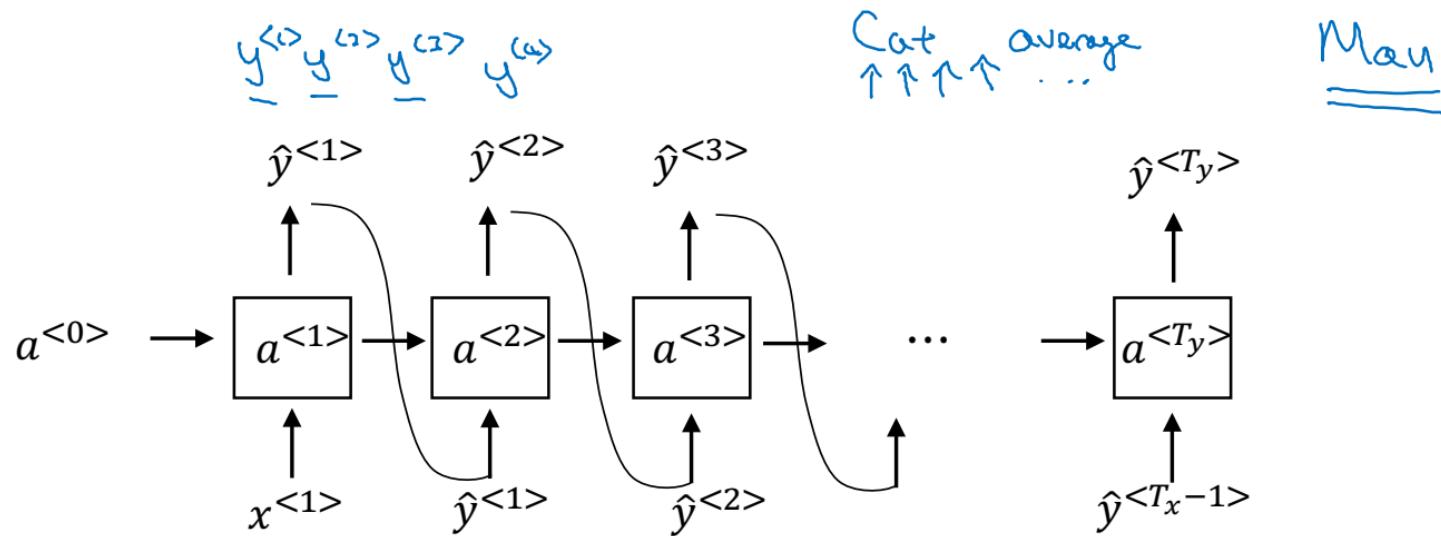
# Sampling a sequence from a trained RNN



# Character-level language model

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

$\rightarrow \text{Vocabulary} = [a, b, c, \dots, z, \cup, \circ, \rightarrow, ;, 0, \dots, 9, A, \dots, Z]$



# Sequence generation

## News

President enrique peña nieto, announced  
sench's sulk former coming football langston  
paring.

“I was not at all surprised,” said hich 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 lesser be my love to me see sabl’s.

For whose are ruse of mine eyes heaves.

## Vanishing Gradients with RNNs

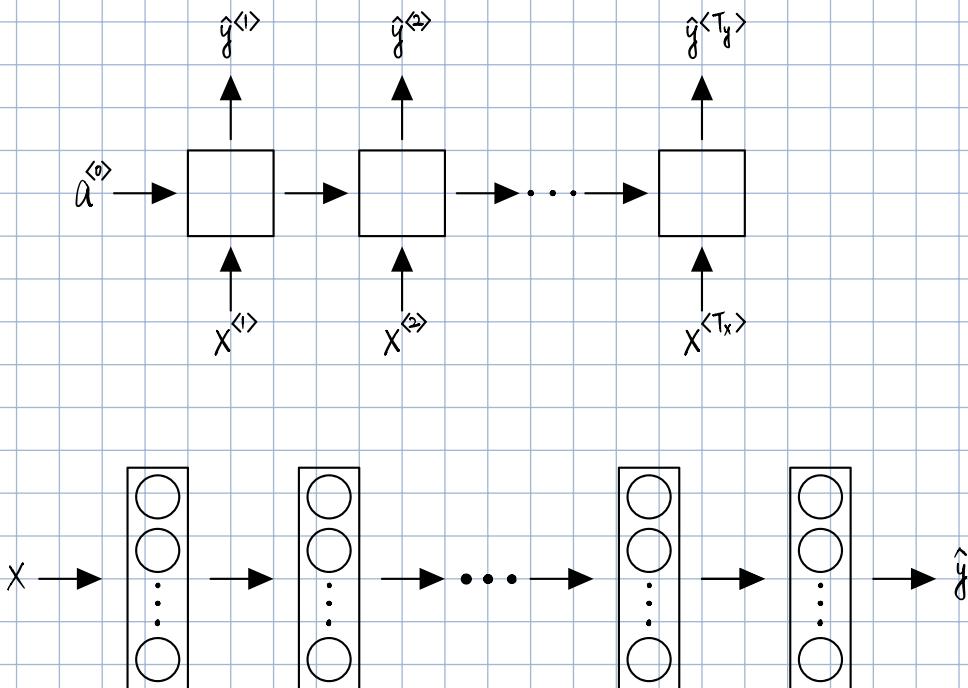
Example sentences:

This can be arbitrary long-

The cat, which already ate apple, pear, banana, ..., was full.

The cats, which already ate apple, pear, banana, ..., were full.

The basic RNN is not good at capturing very long term dependency.



Vanishing gradients is hard to solve.

In contrast, exploding gradients is easier to solve.

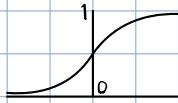
Apply gradient clipping to clip very large gradients > threshold.

# Gated Recurrent Unit (GRU)

RNN unit:

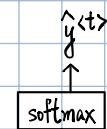
$$\hat{a}^{(t)} = g(W_a[a^{(t-1)}, x^{(t)}] + b_a)$$

$\sigma$ : sigmoid function



$a^{(t-1)}$

$x^{(t)}$



GRU (simplified) C: memory cell  $C^{(t)} = a^{(t)}$

$$\tilde{C}^{(t)} = \tanh(W_c[C^{(t-1)}, x^{(t)}] + b_c)$$

$$\text{Gate } \Gamma_u = \sigma(W_u[C^{(t-1)}, x^{(t)}] + b_u)$$

Imagine  $\Gamma_u$  is either 0 or 1 most of the time.

$$C^{(t)} = \Gamma_u * \tilde{C}^{(t)} + (1 - \Gamma_u) * C^{(t-1)}$$

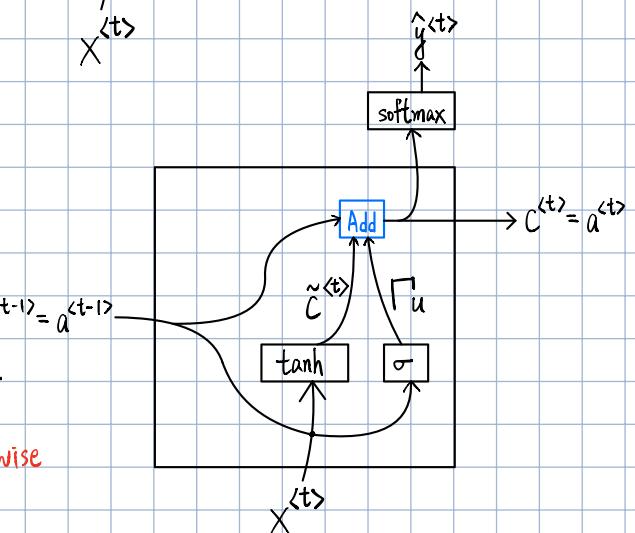
\*: element-wise

If  $\Gamma_u$  is very close to 0, eg,  $\Gamma_u = 10^{-6}$ ,  $C^{(t)}$  will be assigned to  $C^{(t-1)}$ , which can mitigate the vanishing gradient problem. The value of  $C^{(t)}$  is maintained.

Intuition:  $\Gamma_u$  is like a binary gate bit 

0
1

, which can control if the information is updated. The element-wise multiplication can keep some bits as constant while updating other bits.



The cat, which already ate ..., was full.

Full GRU

$$\tilde{C}^{(t)} = \tanh(W_c[\Gamma_r * C^{(t-1)}, x^{(t)}] + b_c)$$

RELEVANCE

$$\Gamma_u = \sigma(W_u[C^{(t-1)}, x^{(t)}] + b_u)$$

update gate

$$\Gamma_r = \sigma(W_r[C^{(t-1)}, x^{(t)}] + b_r)$$

relevance gate

$$C^{(t)} = \Gamma_u * \tilde{C}^{(t)} + (1 - \Gamma_u) * C^{(t-1)}$$

Read "On the properties of neural machine translation: Encoder-decoder approaches"

"Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling"

# Long Short Term Memory (LSTM)

Read 1997 "Long Short Term Memory"

Read Christopher Olah's "Understanding LSTM Networks"

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

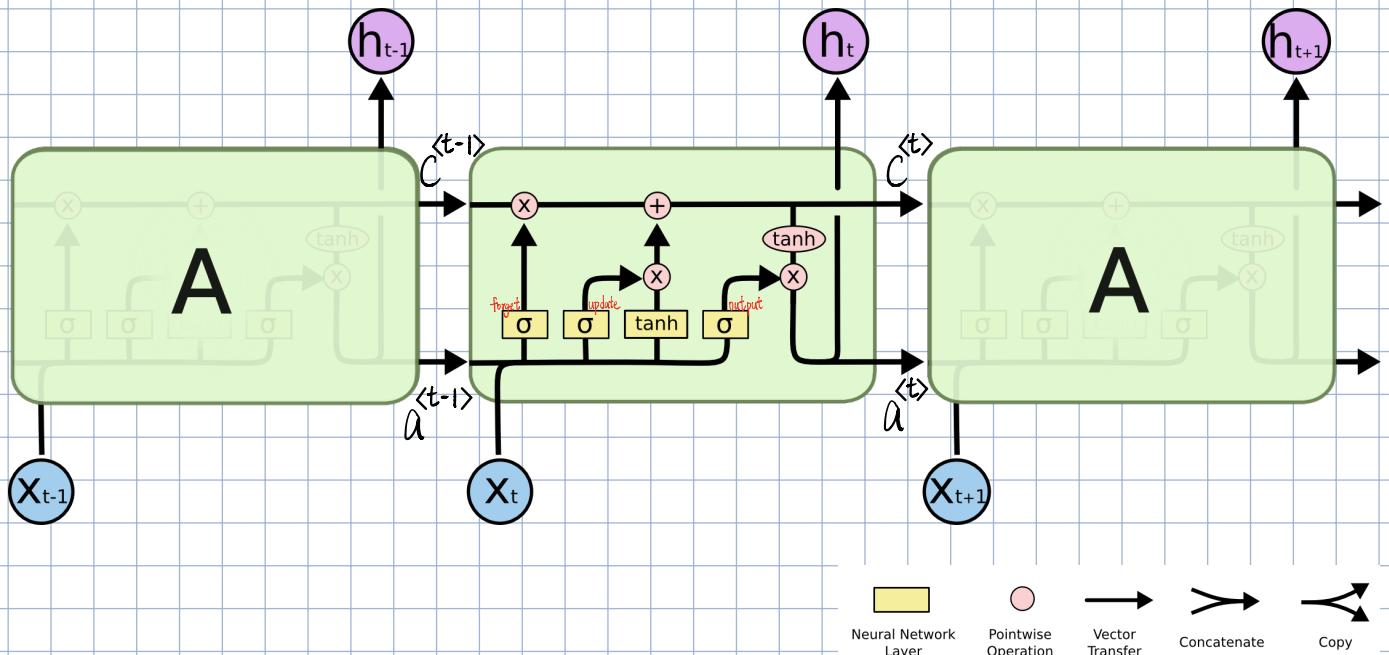
*update gate*  $\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$

*forget gate*  $\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$

*output gate*  $\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh(c^{<t>})$$



# Bidirectional RNN (BRNN)

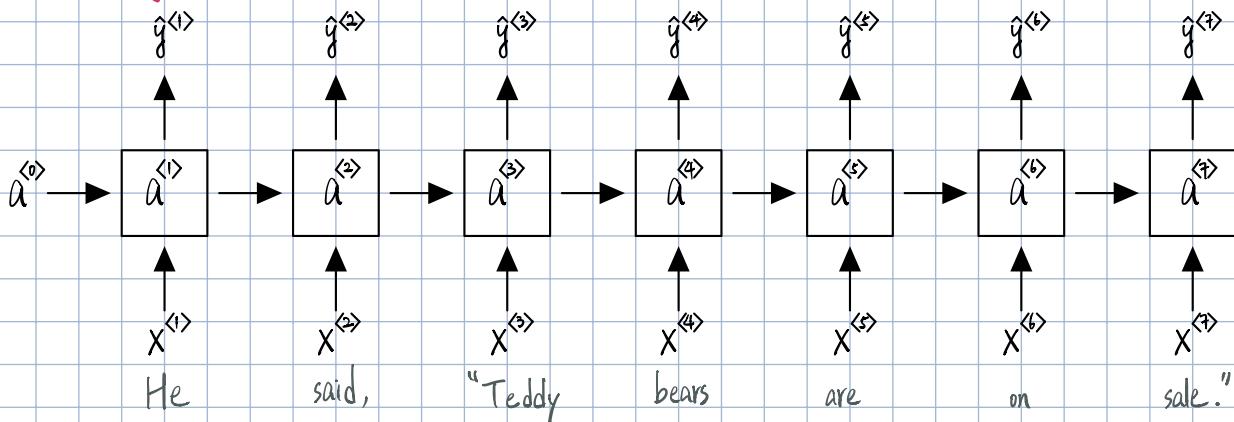
We need to get information from the future to predict the current output.

E.g.:

He said, "Teddy Roosevelt was a great president."

He said, "Teddy bears are on sale."

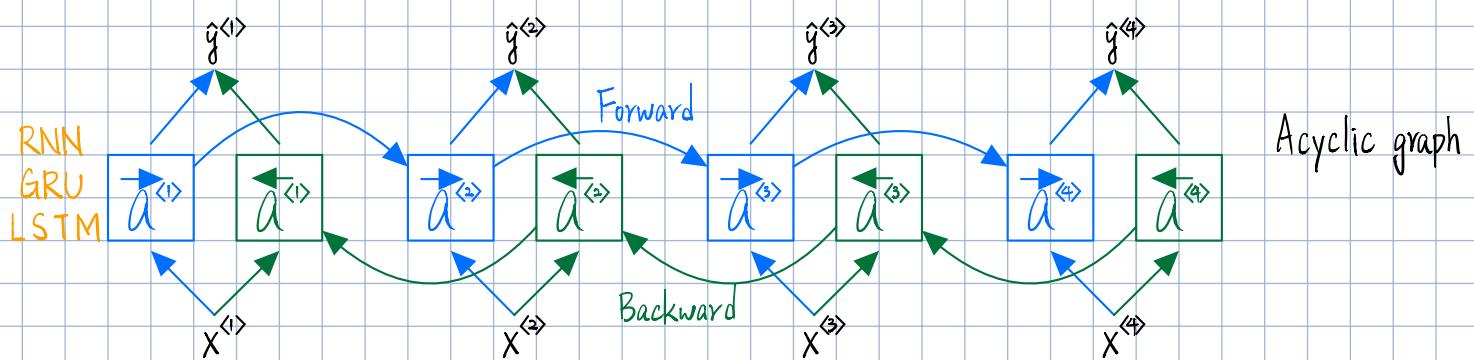
Using unidirectional RNN is just impossible to say if Teddy is a name or not.



Unidirectional structure doesn't work well regardless of using standard RNN, GRU or LSTM or not.

## Bidirectional RNN

$$\text{Prediction: } \hat{y}^{(t)} = g(W_y [a^{(t)}, \bar{a}^{(t)}] + b_y)$$

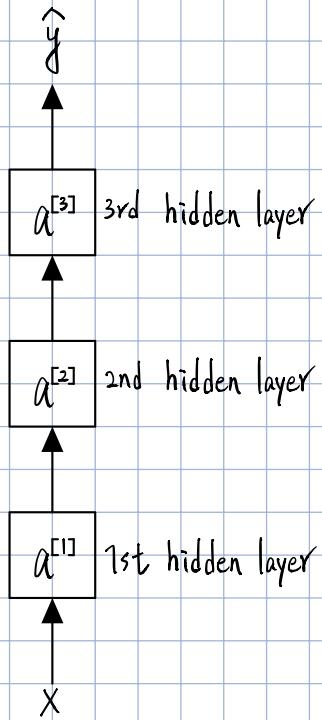


Notice that Forward and Backward connection are both forward propagation. Backward connection is not backpropagation. In other words, forward propagation composes of both Forward and Backward connection.

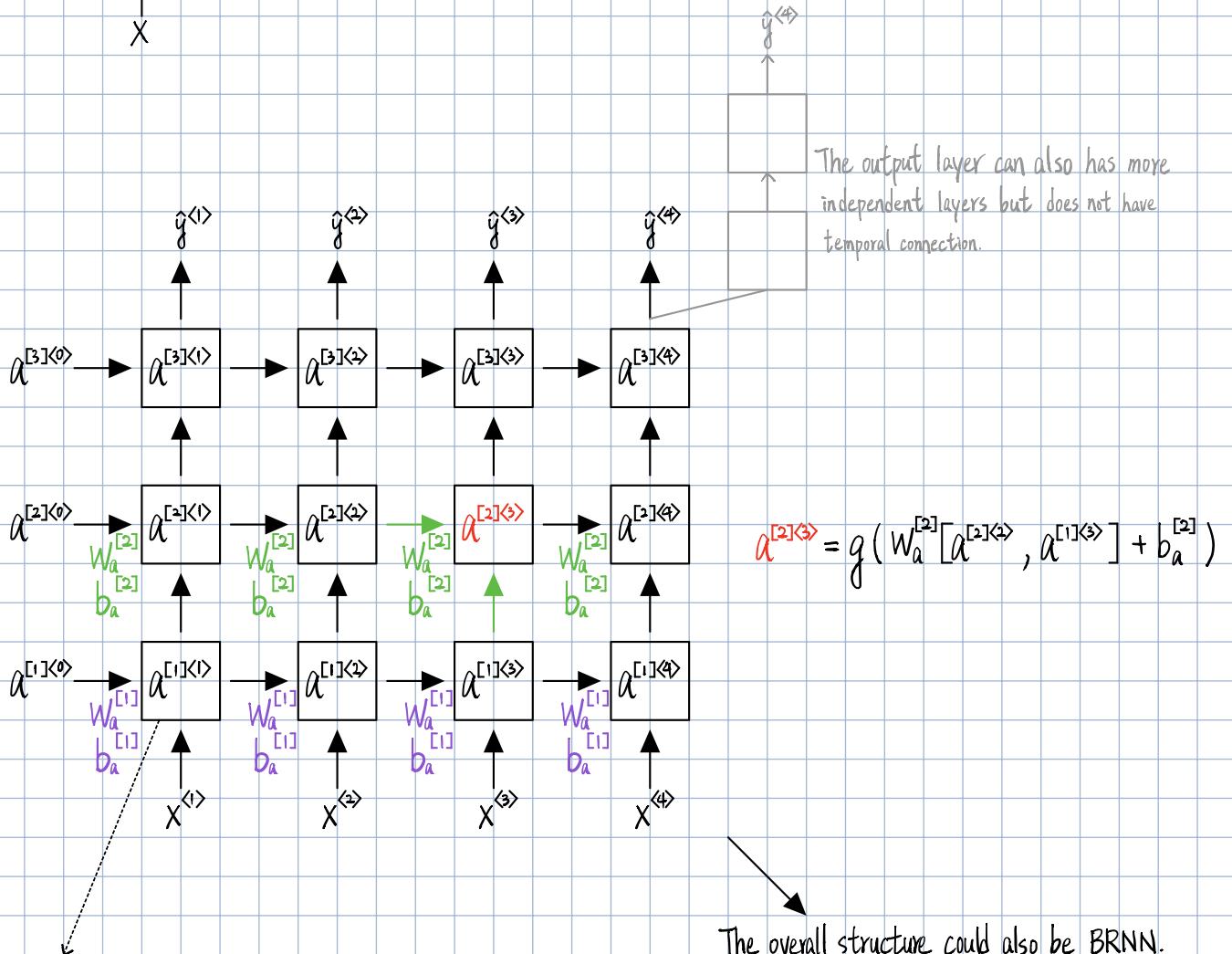
In prediction, to compute  $\hat{y}^{(t)}$ , information from  $x^{(1)}$  to  $x^{(t)}$  are used.

# Deep RNNs

Recall



$a^{[l](t)}$ :  $l$  is associated with the layer  $l$ .  
 $t$  is associated with the time  $t$ .



Blocks could be standard RNNs, GRU, LSTM

The overall structure could also be BRNN.