



deeplearning.ai

Recurrent Neural Networks

RNN have revolutionized speech recognition, natural language processing and other areas, here we learn how to build these models

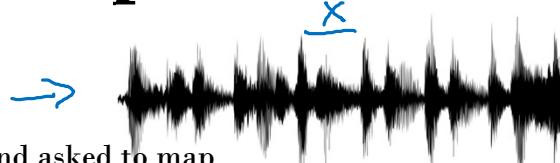
Why sequence models?

Examples of sequence data

Speech recognition

you are given an audio clip input x and asked to map

it to a transcript y , both input and output are sequence data



y
“The quick brown fox jumped
over the lazy dog.”

Music generation

this is another example with sequence data, only output is a sequence, input can be empty set or single integer



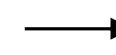
Sentiment classification

“There is nothing to like
in this movie.”



DNA sequence analysis

\rightarrow AGCCCCCTGTGAGGAAC TAG



AGCCCCTGTGAGGAAC TAG

this part of the dna sequence correspond
to a protein

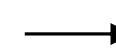
Do you want to sing with
me?

Plasmid sequence AP4 Oligo

Running

Machine translation

Voulez-vous chanter avec
moi?

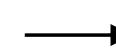


Do you want to sing with
me?

Plasmid sequence AP4 Oligo

Running

Video activity recognition



Plasmid sequence AP4 Oligo

Running

you are given a sequence of frames and you are
asked to recognize the activity

Name entity recognition

\rightarrow Yesterday, Harry Potter
met Hermione Granger.



Yesterday, Harry Potter
met Hermione Granger.

Andrew Ng

you are given a sentence and you are asked
to identify the people

so its supervised learning and there are different type of sequences
in some examples x and y have same lenght, sometime only one is a sequence.



deeplearning.ai

Recurrent Neural Networks

that's start by defining a notation to define these sequence models

Notation

Motivating example

say you want to input a sequence like x and we want y to tell where are peoples names. this problem is called Named Entity Recognition, and its used by search engines. Some also tell where name starts and where it ends.

$x:$ Harry Potter) and Hermione Granger invented a new spell.

$\rightarrow x^{<1>} \quad x^{<2>} \quad x^{<3>}$

these index the words

$\dots \quad x^{<t>} \quad \dots \quad x^{<n>}$

$$T_x = 9$$

$\rightarrow y:$

| | 0 | |
 $y^{<1>} \quad y^{<2>} \quad y^{<3>}$

same in indexing output

$\dots \quad 0 \quad 0 \quad 0 \quad 0 \quad y^{<n>}$

$$T_y = 9$$

that i there indicates
the training example

$x^{(i)<t>} \quad y^{(i)<t>}$

$$T_x^{(i)} = 9$$

15

$$T_y^{(i)}$$

these indicate the lenght of the input and
output of the x and y of the example i

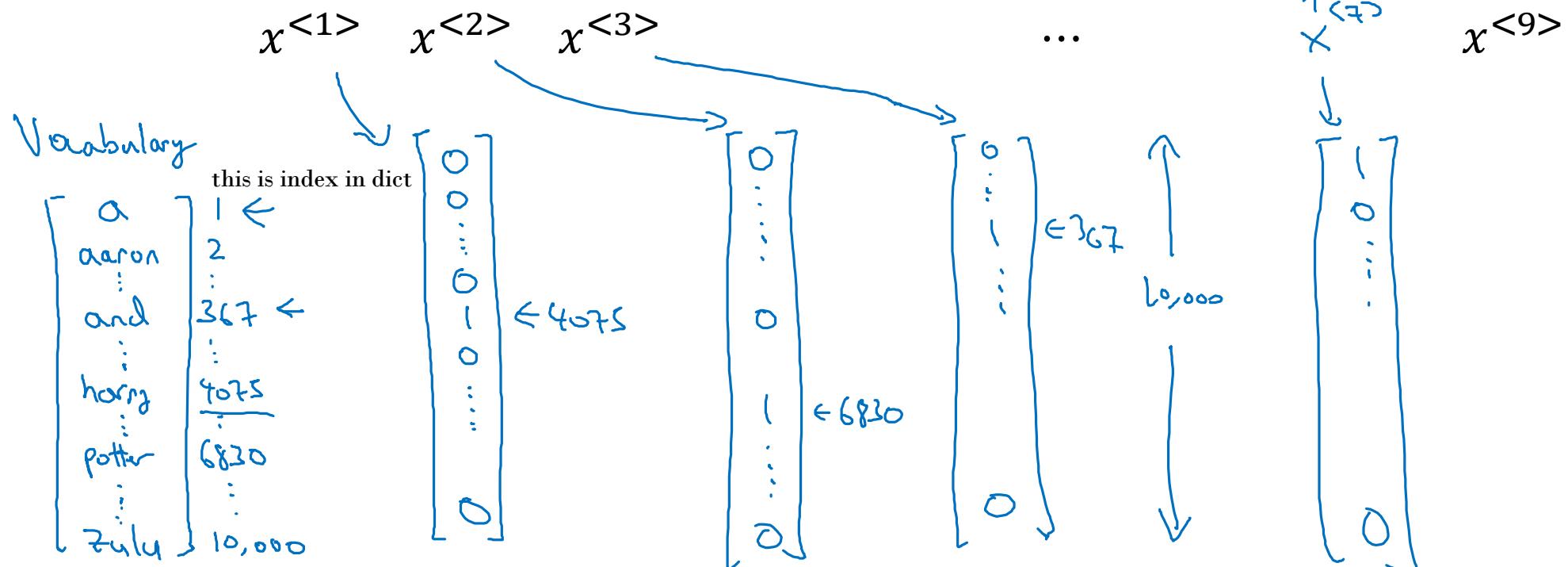
Andrew Ng

how to represent the words in a sentence, you take a vocabulary

Representing words

$$x^{(t)} \rightarrow y^{(t)}$$

$x:$ Harry Potter and Hermione Granger invented a new spell.



One-hot

then we use one hot encoding to represent Harry which will be all 0 vector except a 1 in index where we find harry in dict each of this vectors will be of size 10000 just like dict

Andrew Ng

this for words not in dictionary

Representing words

x: Harry Potter and Hermione Granger invented a new spell.
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<9>}$

And = 367
Invented = 4700
A = 1
New = 5976
Spell = 8376
Harry = 4075
Potter = 6830
Hermione = 4200
Gran... = 4000

Andrew Ng



deeplearning.ai

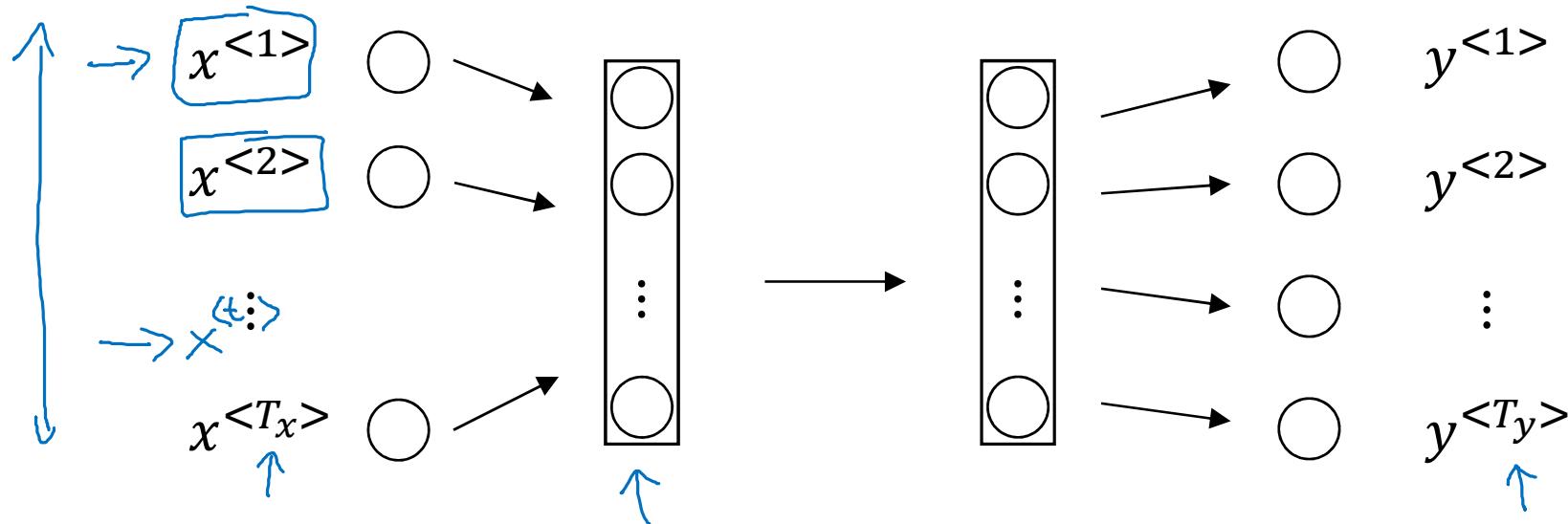
Now we see how to learn the mapping from x to y with the neural network

Recurrent Neural Networks

Recurrent Neural Network Model

Why not a standard network?

feeding the 9 inputs to a NN and then get those 9 y its not a good idea for the problems below



Problems:

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

because u have examples with different sentences that contain different amount of words

Andrew Ng

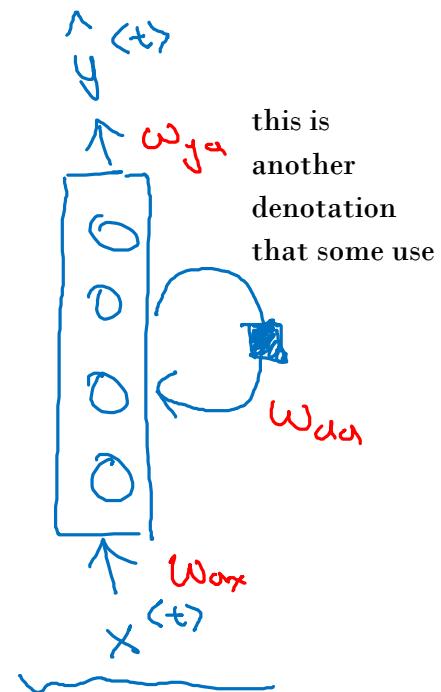
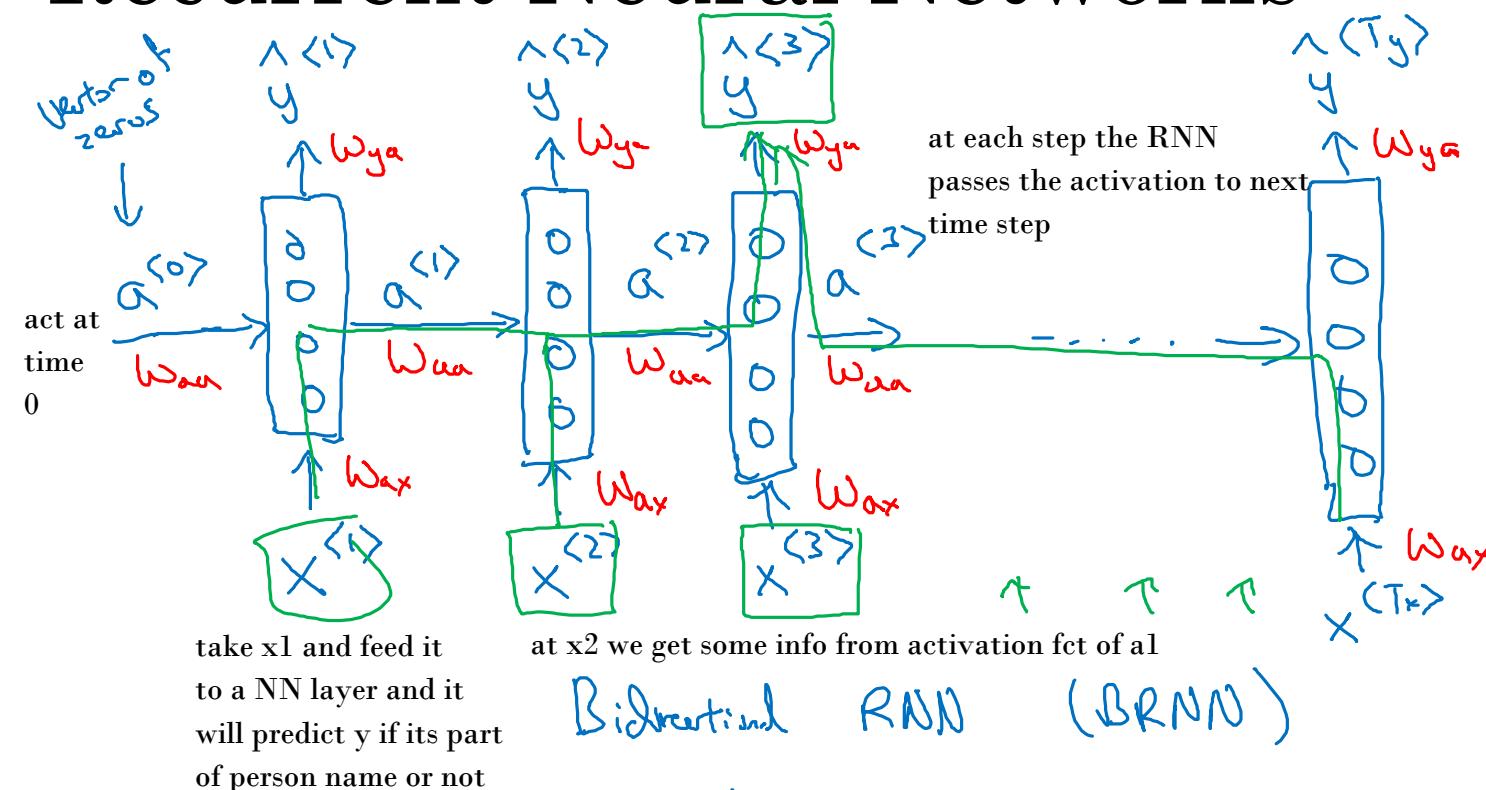
A RNN does not have these two disadvantages

What is a recurrent NN? that's build one up

Recurrent Neural Networks

$$T_x = T_y$$

in this case they are equal,
architecture change a bit if not equal



He said, “Teddy Roosevelt was a great President.”

He said, “Teddy bears are on sale!”

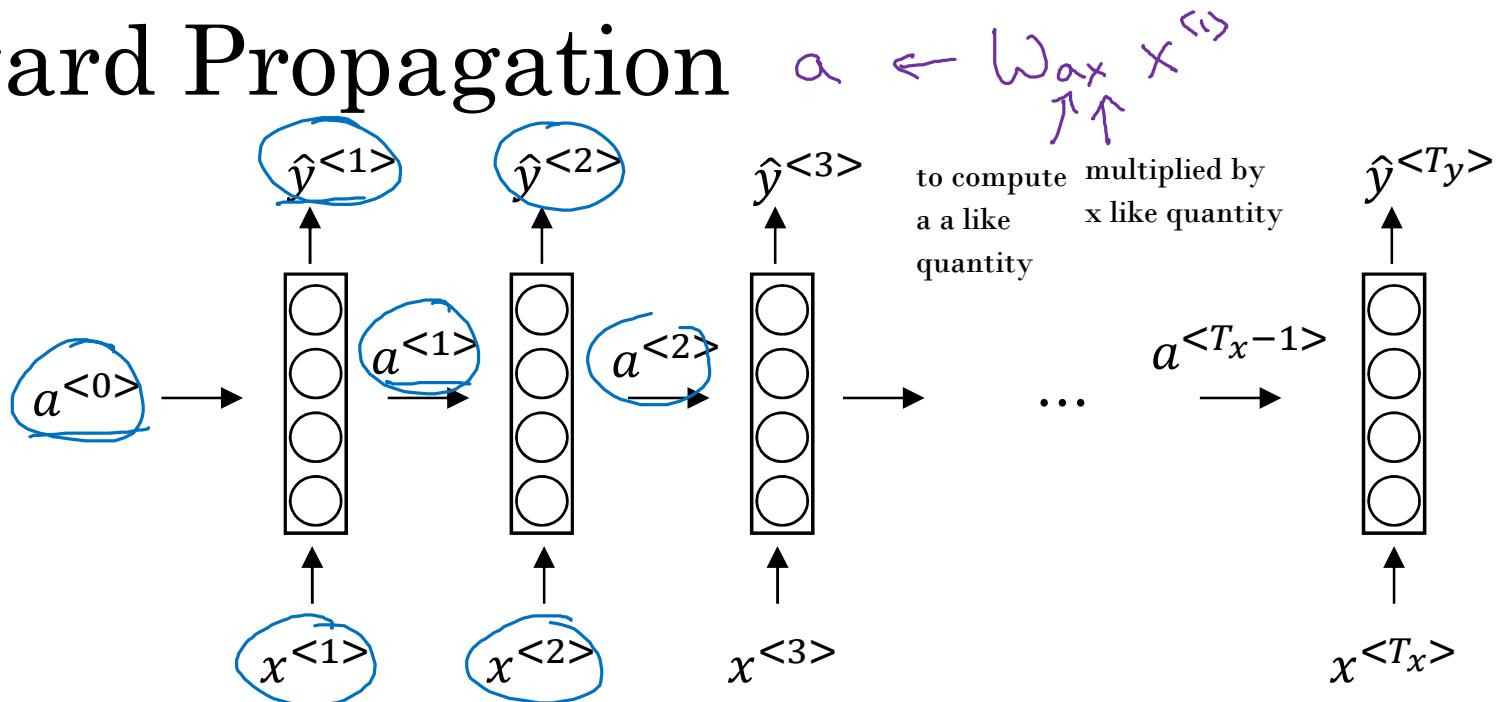
you can't tell whether Teddy is a name by just looking at first 3 words. We will address this when we talk about bidirectional recurrent NN
the parameters W_{ax} are shared. we describe in detail next slide

one weakness is that it uses info only about previous words not next ones

Andrew Ng

that see what are the calculations that the RNN does

Forward Propagation



$$a^{(t)} = g_1(W_{aa} a^{(t-1)} + W_{ax} x^{(t)} + b_a) \leftarrow \text{tanh / ReLU}$$

$$\hat{y}^{(t)} = g_2(W_{ya} a^{(t)} + b_y) \leftarrow \text{Sigmoid}$$

depends on type of output u have

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

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

this is more generally

Andrew Ng

Simplified RNN notation

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

we write this simpler way

$(100, 100)$ 100
 $(100, 10, 000)$ $(100, 10, 000)$

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

$$y^{(t)} = g(W_y a^{(t)} + b_y)$$

↑ ↑ ↑

Here you have the simplified notation

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

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

$(100, 10, 000)$

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

100 $10,000$ $10,000$

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

Andrew Ng

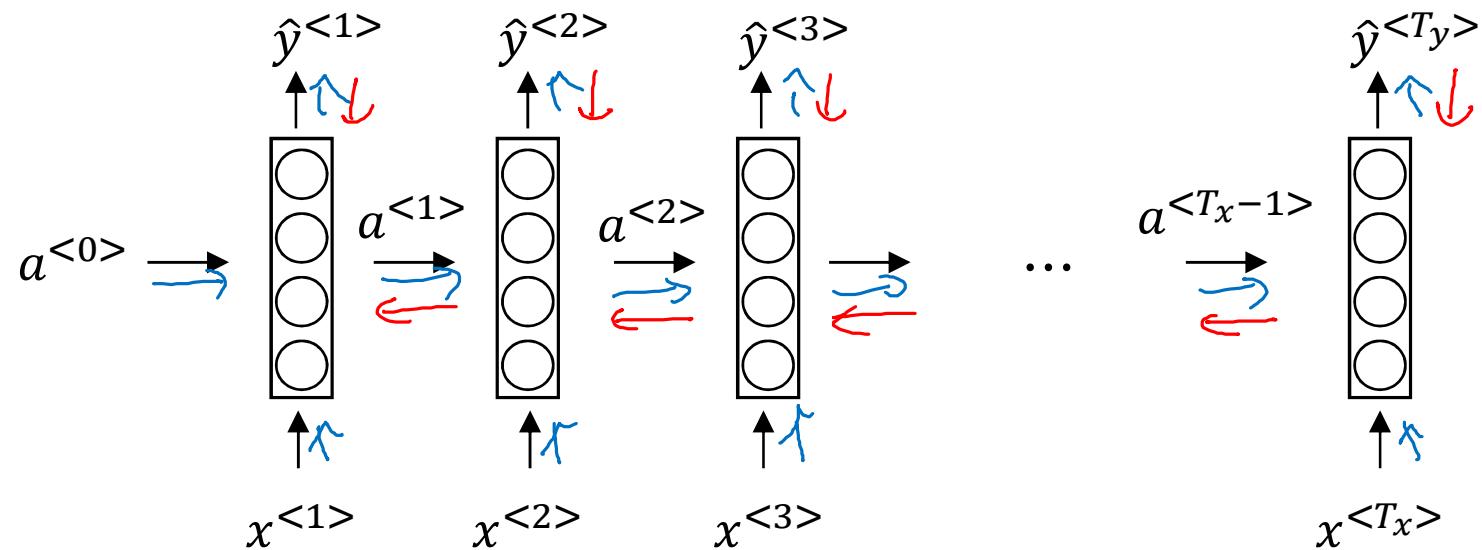


deeplearning.ai

Recurrent Neural Networks

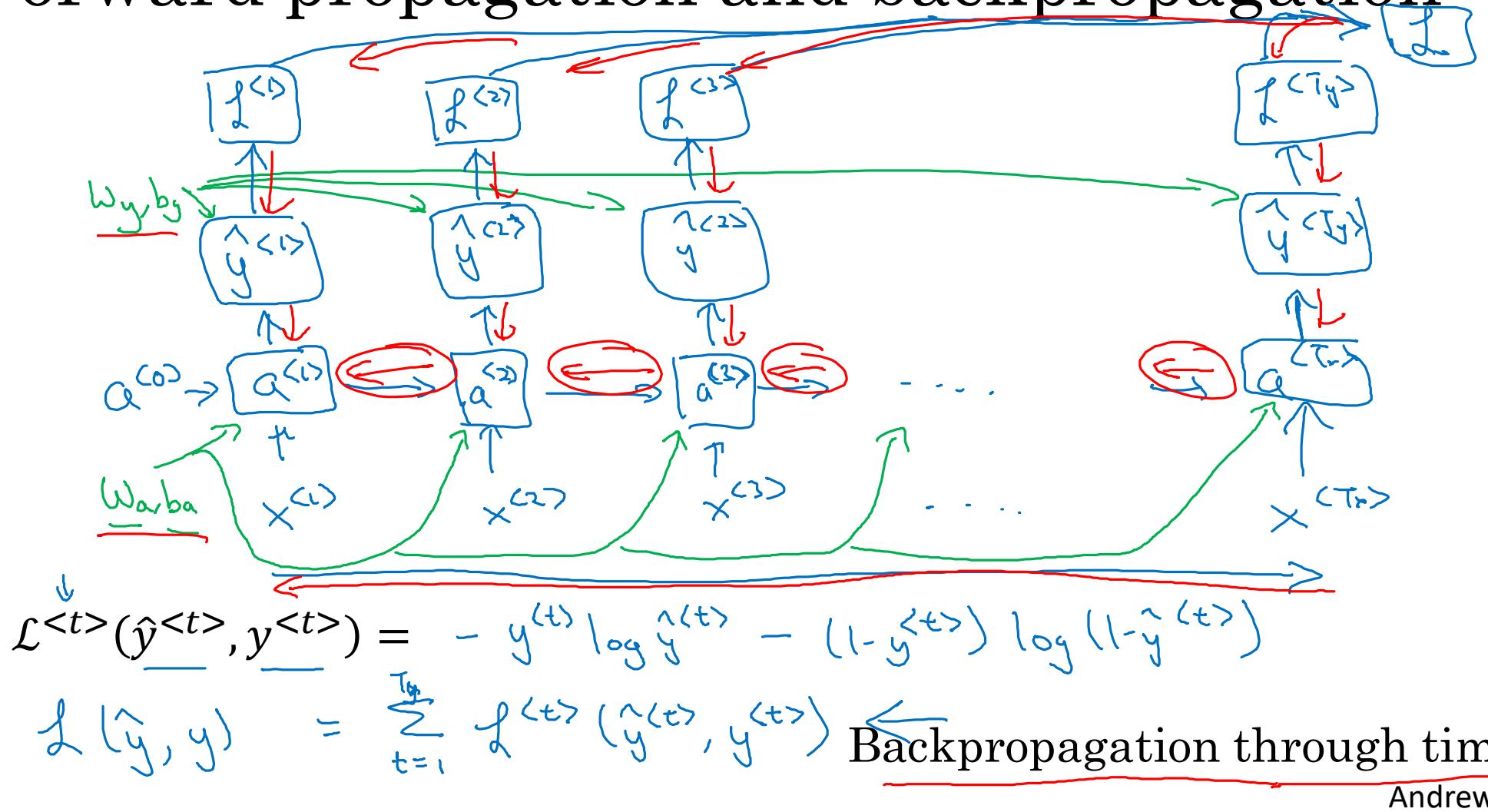
Backpropagation through time

Forward propagation and backpropagation



Andrew Ng

Forward propagation and backpropagation



Andrew Ng



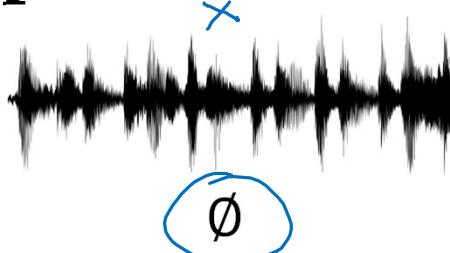
deeplearning.ai

Recurrent Neural Networks

Different types of RNNs

Examples of sequence data

Speech recognition



not allways these two are equal

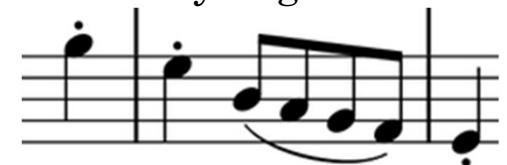
T_x

T_y

y

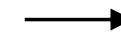
"The quick brown fox jumped over the lazy dog."

Music generation



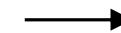
Sentiment classification

"There is nothing to like
in this movie."



DNA sequence analysis

AGCCCCTGTGAGGAAC TAG



AG~~CCCCTGTGAGGAAC~~ TAG

Machine translation

Voulez-vous chanter avec
moi?



Do you want to sing with
me?

Video activity recognition



Running

Name entity recognition

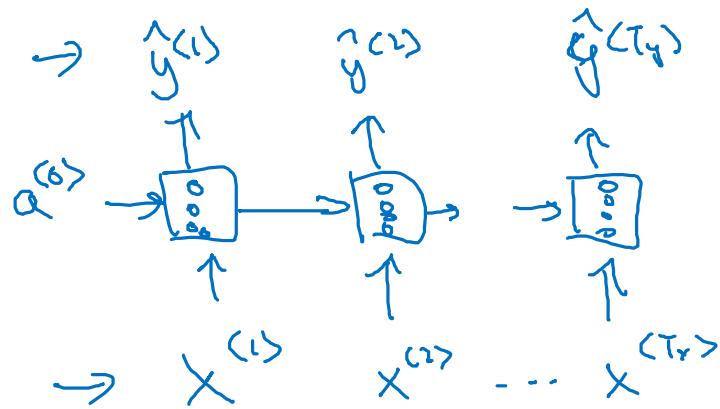
Yesterday, Harry Potter
met Hermione Granger.



Yesterday, ~~Harry Potter~~
met ~~Hermione Granger~~.
Andrew Ng

Examples of RNN architectures

$$T_x = T_y$$



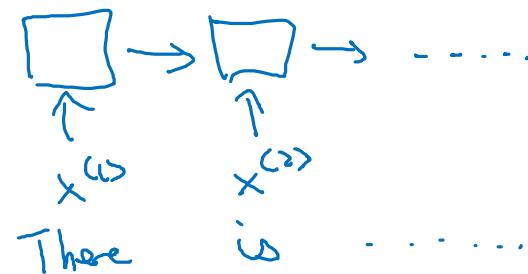
Many-to-many

many inputs many outputs

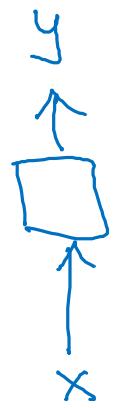
Sentiment classification

$x = \text{text}$ this is the text input, the review

$y = 0/1 \quad 1 \dots 5$ review of a movie positive/negative or starts from 1 to 5



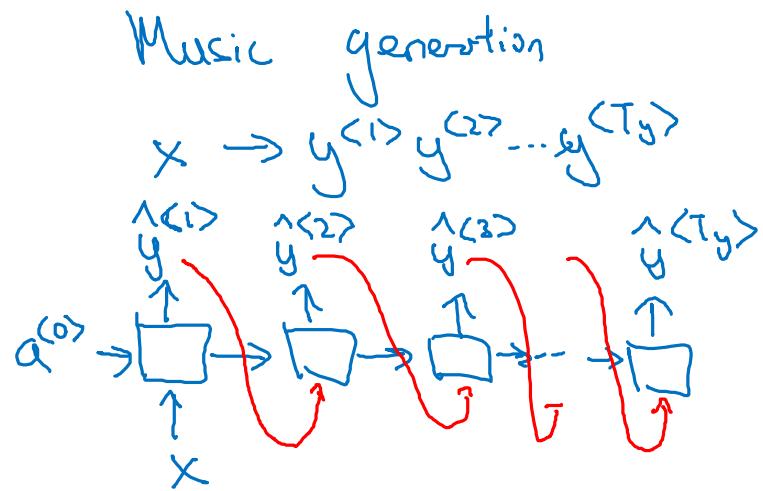
Many-to-one



One-to-one

this is NN that
Andrew Ng
we did in first
two courses

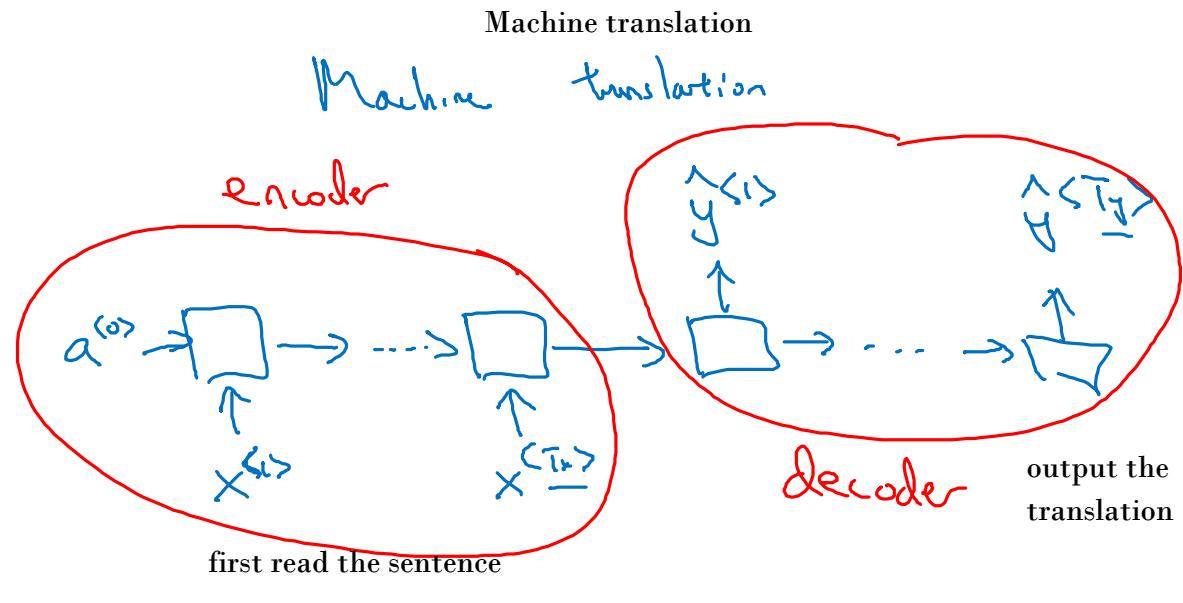
Examples of RNN architectures



One-to-many

$$x = \phi$$

an example of one to many is music generation,
and we can implement thi in the exercises. Input can
be an integer identifying the gener of music u want, or u can have
no input just input of 0

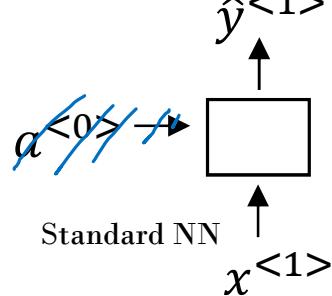


Many - to - many

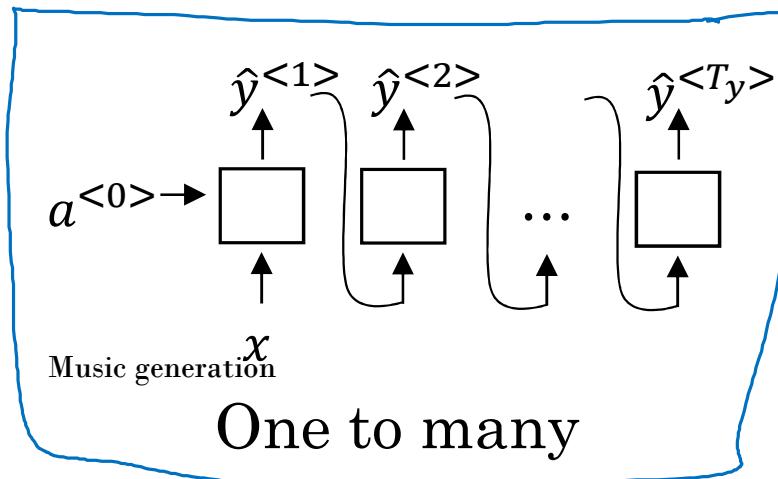
when input length and output length are different.

Andrew Ng

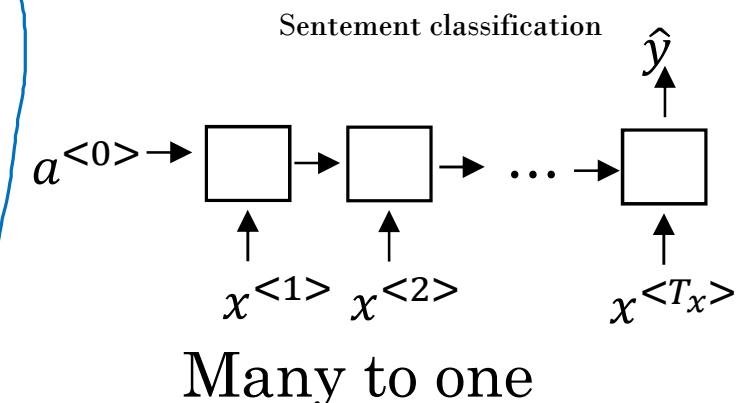
Summary of RNN types



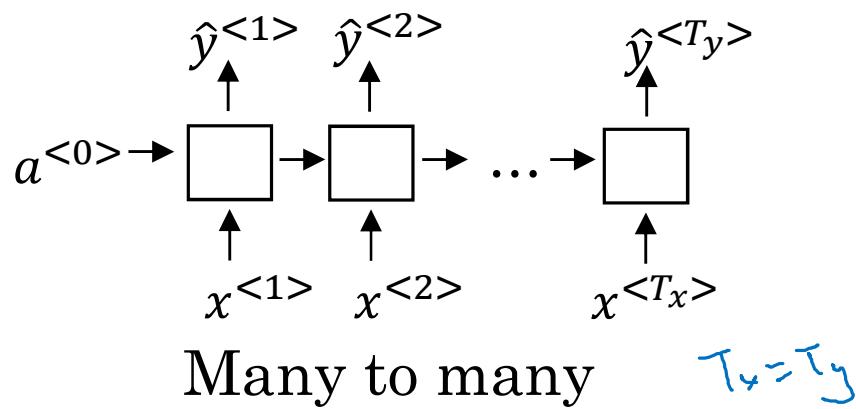
One to one



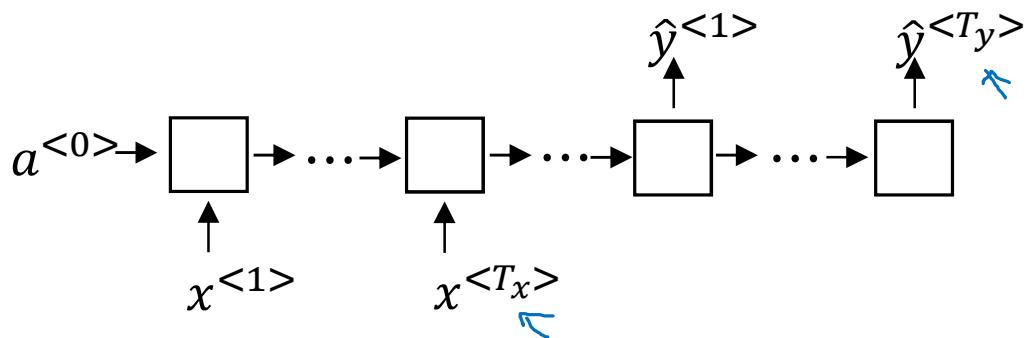
One to many



Many to one



Many to many



Many to many

Machine translation

Andrew Ng



deeplearning.ai

Recurrent Neural Networks

Language model and sequence generation

What is language modelling?

Speech recognition

The apple and pair salad.
→ The apple and pear salad.

These two sound very similar but a speech recognit would hear the second one.

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

A probability model tells the chances of these two models.

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

this is much more likely

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

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

What a language model does is tells the probability of that sentence.

Andrew Ng

Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

Cats average 15 hours of sleep a day. \downarrow $\langle \text{EOS} \rangle$

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

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

...

$y^{(8)}$

$y^{(9)}$

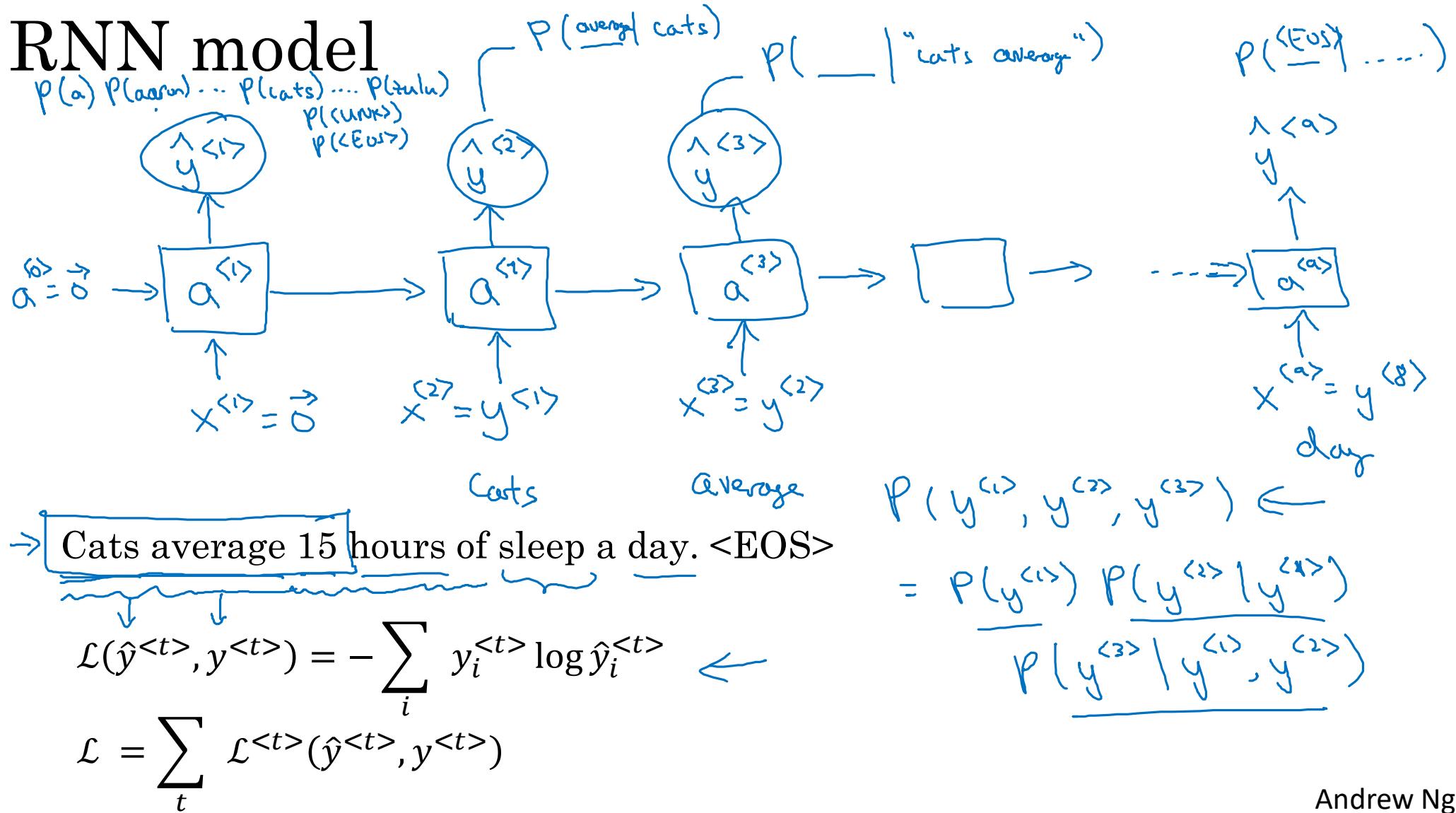
The Egyptian ~~Mau~~ is a bread of cat. $\langle \text{EOS} \rangle$

10,000

$\langle \text{UNK} \rangle$

Andrew Ng

RNN model



Andrew Ng



deeplearning.ai

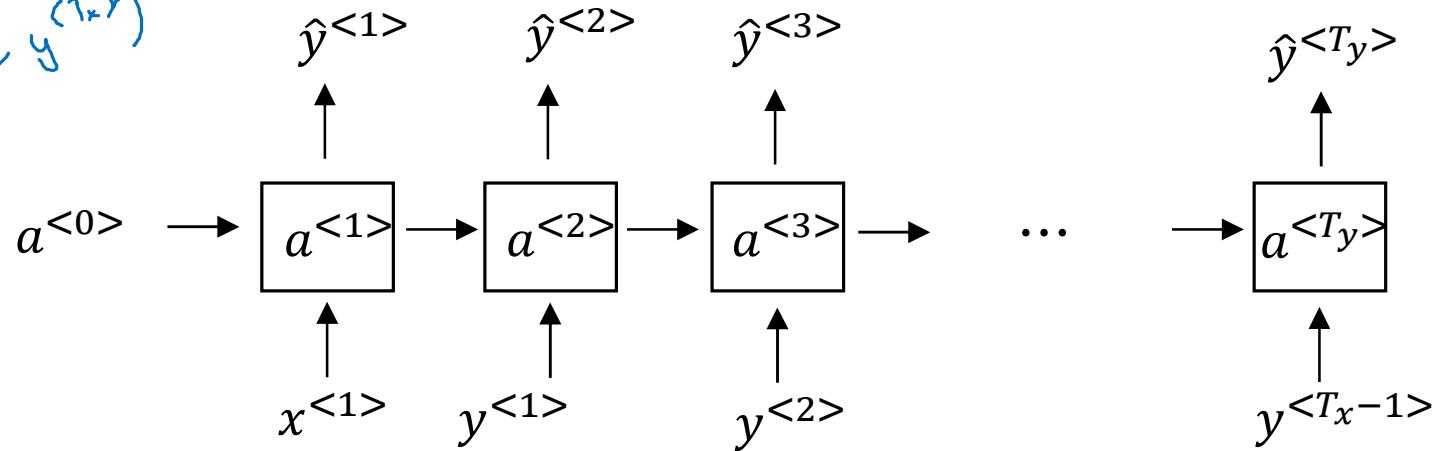
Recurrent Neural Networks

Sampling novel sequences

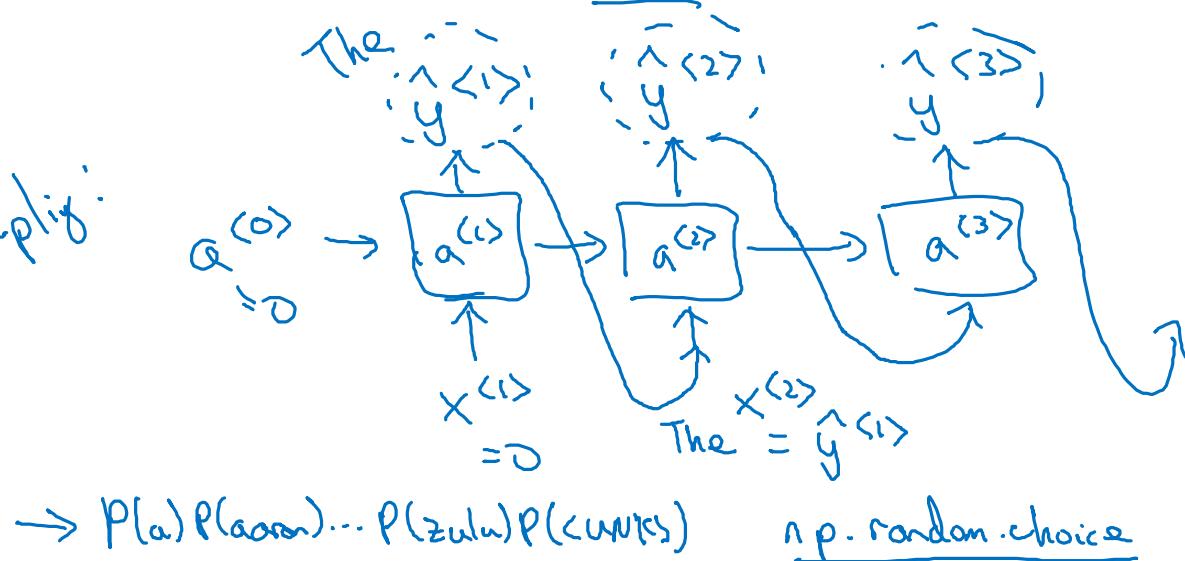
Sampling a sequence from a trained RNN

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

Training:



Sampling:



$$\rightarrow p(a) p(a|a_0) \dots p(z|a_{T_y}) p(c|z)$$

n.p. random.choice

$P(- | \text{the})$

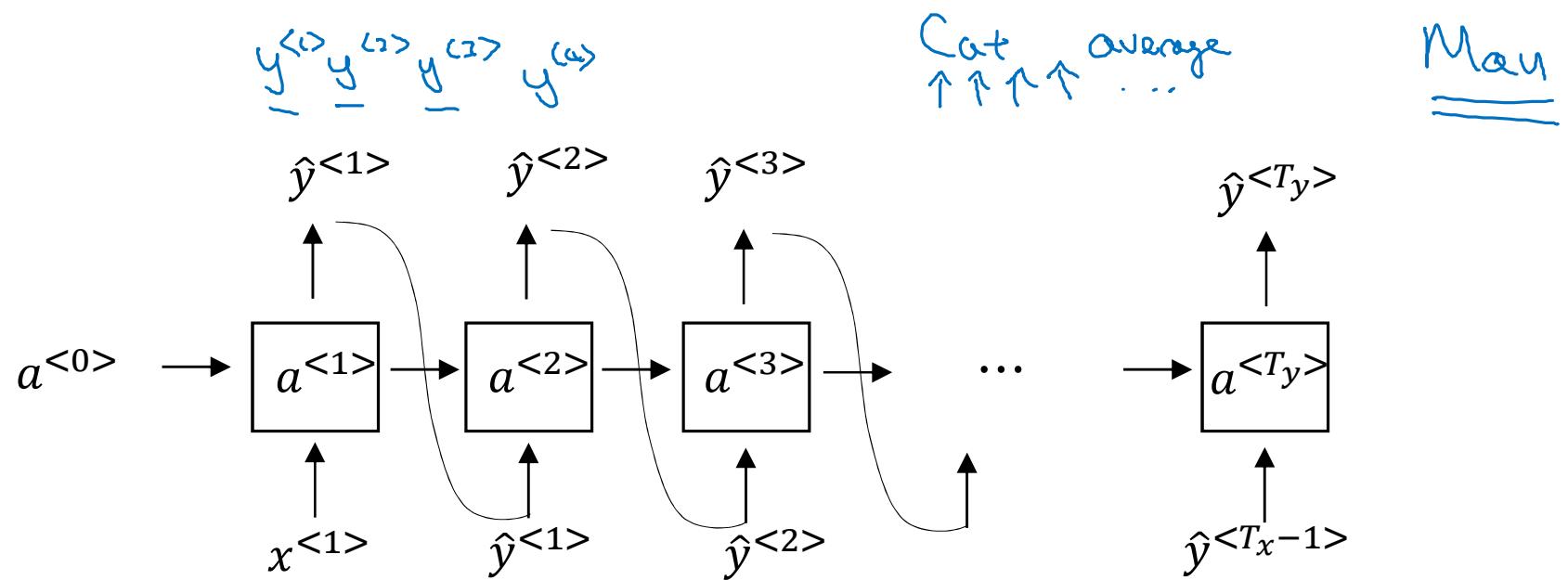
<Eos>
<unk>

Andrew Ng

Character-level language model

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

$\rightarrow \text{Vocabulary} = [a, b, c, \dots, z, \cup, \circ, \rightarrow, ;, i, o, \dots, q, A, \dots, \tilde{z}]$



Andrew Ng

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.

Andrew Ng

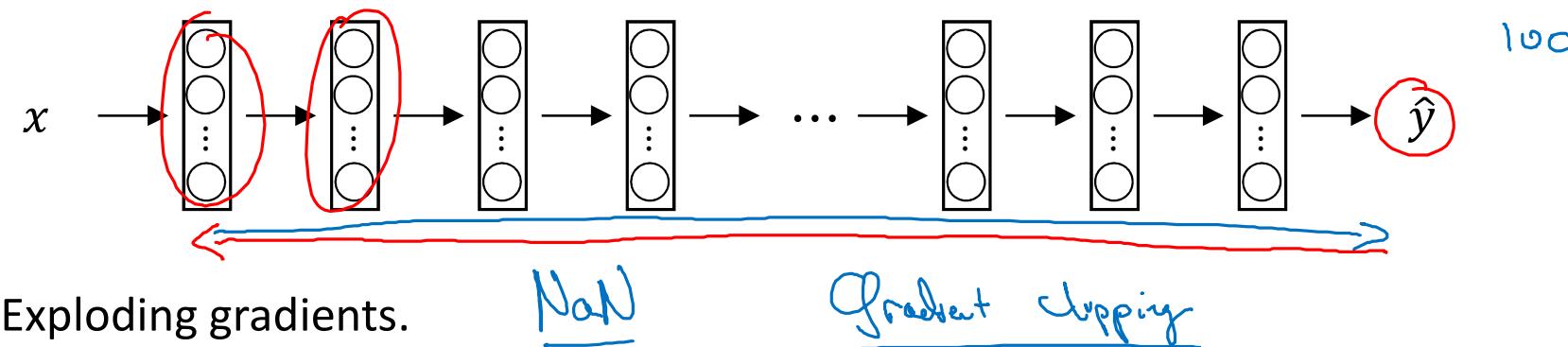
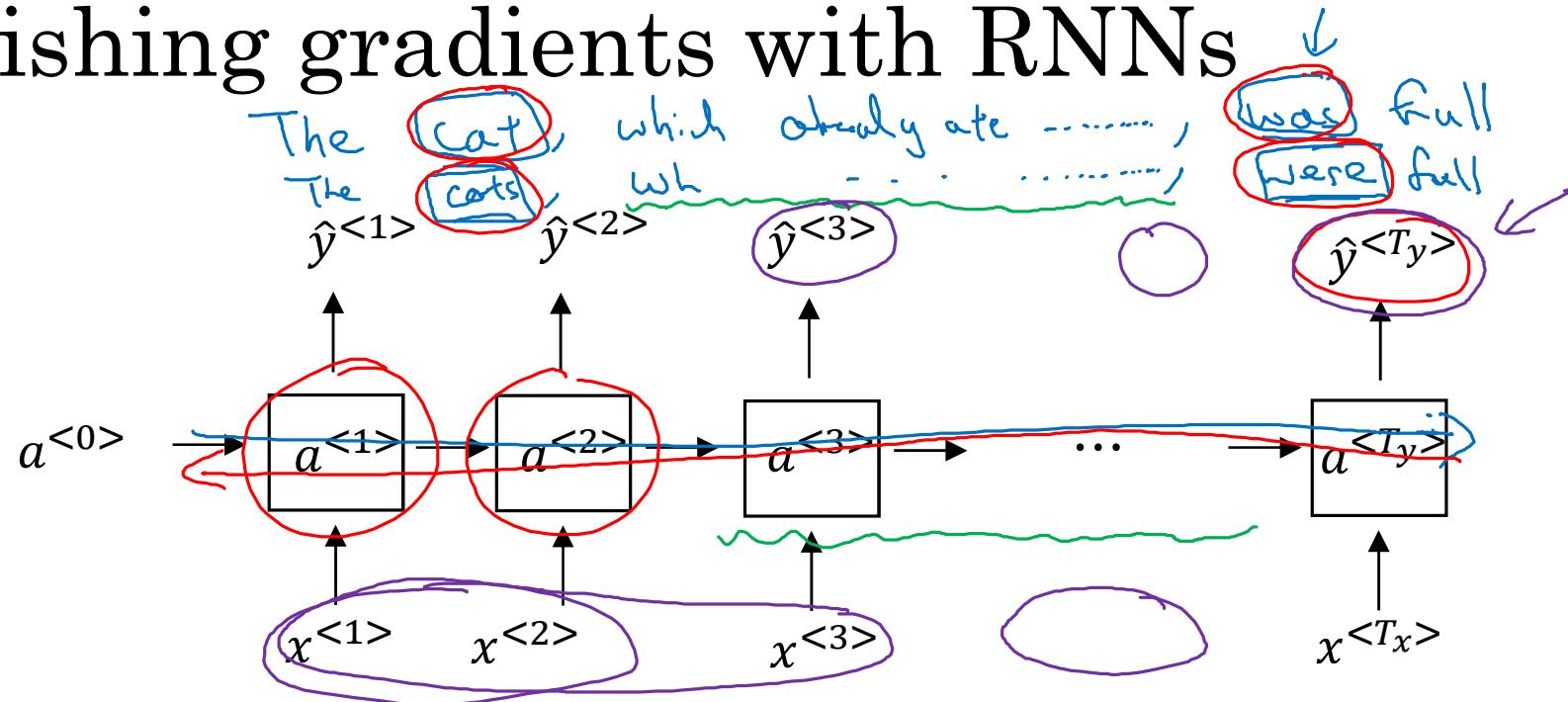


deeplearning.ai

Recurrent Neural Networks

Vanishing gradients with RNNs

Vanishing gradients with RNNs



Andrew Ng

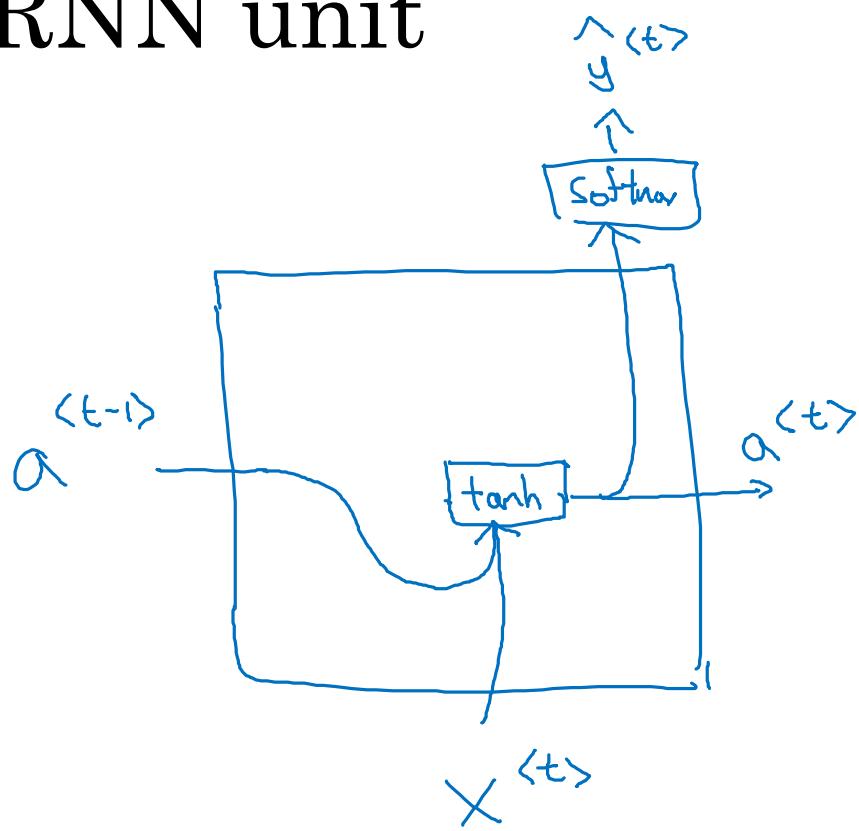


deeplearning.ai

Recurrent Neural Networks

Gated Recurrent Unit (GRU)

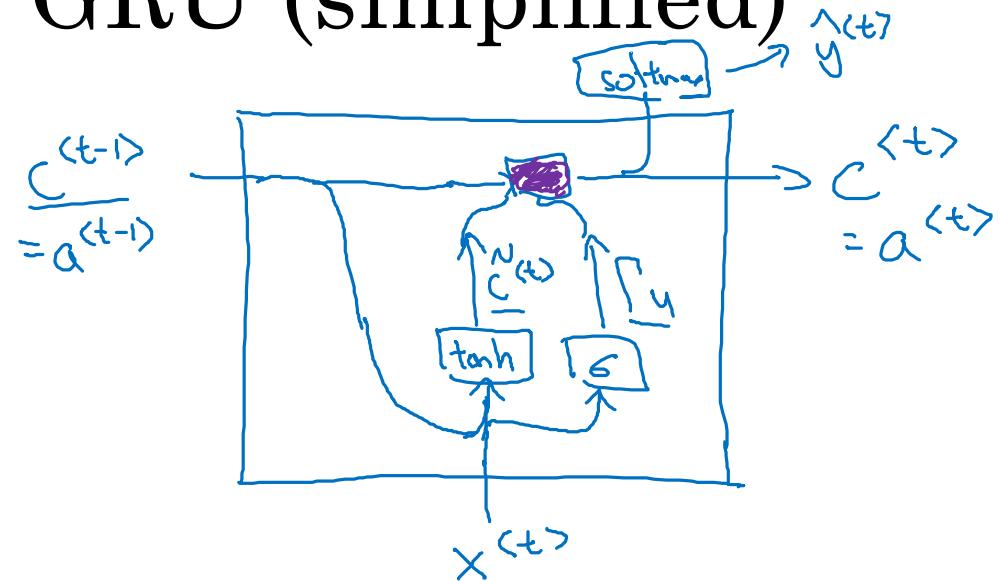
RNN unit



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

Andrew Ng

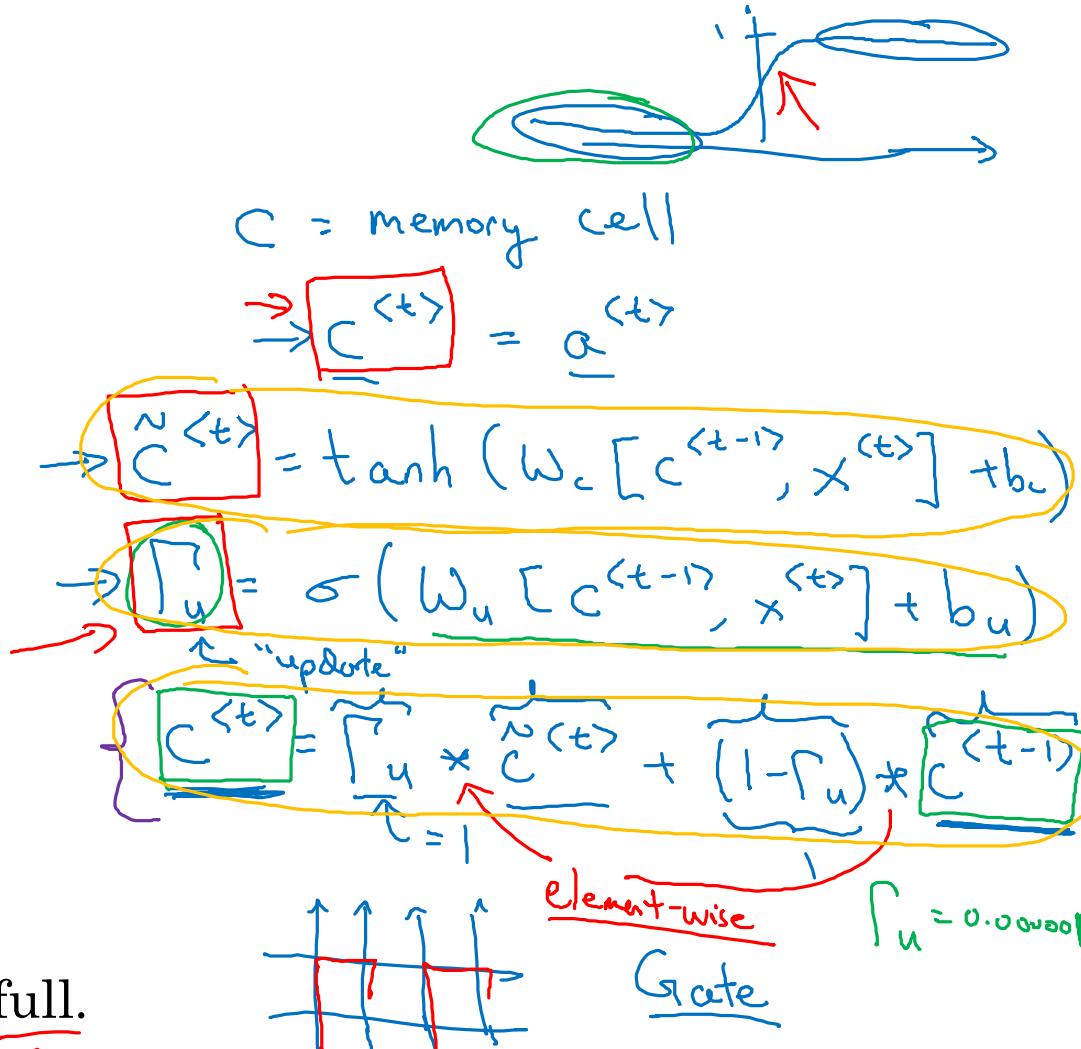
GRU (simplified)



$\Gamma_u = 1$ $\Gamma_u = 0$ $\Gamma_u = 0$ $\Gamma_u = 0$... $\Gamma_u = 1$

$c^{(t)} = 1$

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



[Cho et al., 2014. On the properties of neural machine translation: Encoder-decoder approaches]

[Chung et al., 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling]

Andrew Ng

Full GRU

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

$$\begin{matrix} u \\ r \end{matrix} \quad \left\{ \begin{array}{l} \Gamma_u = \sigma(W_u [c^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_r = \sigma(W_r [c^{<t-1>}, x^{<t>}] + b_r) \end{array} \right.$$

LSTM

$$h \quad c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

The cat, which ate already, was full.

Andrew Ng



deeplearning.ai

Recurrent Neural Networks

LSTM (long short term memory) unit

GRU and LSTM

GRU

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

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

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

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

$$\underline{a}^{<t>} = \underline{c}^{<t>}$$

LSTM

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

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

(update)

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

(forget)

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

(output)

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

$$\underline{a}^{<t>} = \underline{\Gamma_o} * \underline{c}^{<t>}$$

[Hochreiter & Schmidhuber 1997. Long short-term memory]

Andrew Ng

LSTM units

GRU

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

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

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

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

$$a^{<t>} = c^{<t>}$$

LSTM

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

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

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

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

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

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

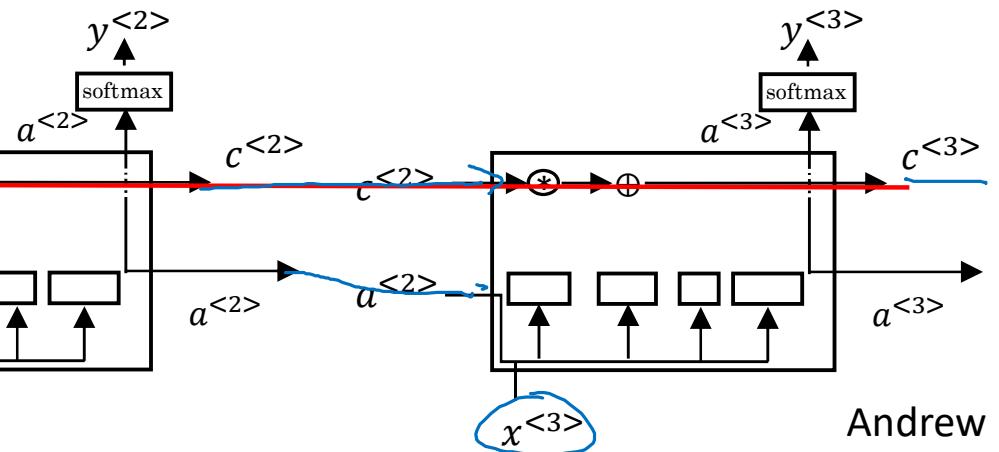
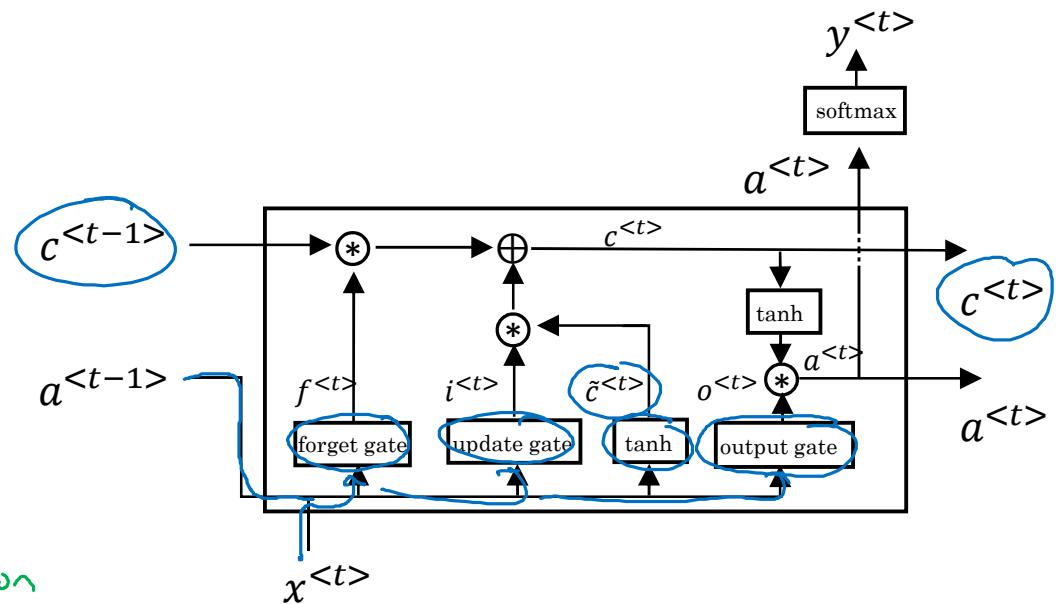
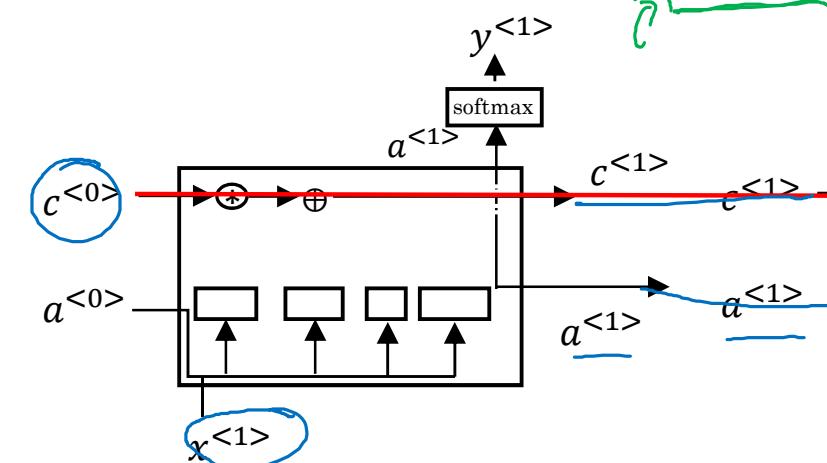
[Hochreiter & Schmidhuber 1997. Long short-term memory]

Andrew Ng

LSTM in pictures

$$\begin{aligned}\tilde{c}^{<t>} &= \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_f &= \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f) \\ \Gamma_o &= \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o) \\ c^{<t>} &= \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>} \\ a^{<t>} &= \Gamma_o * c^{<t>}\end{aligned}$$

peephole connection



Andrew Ng



deeplearning.ai

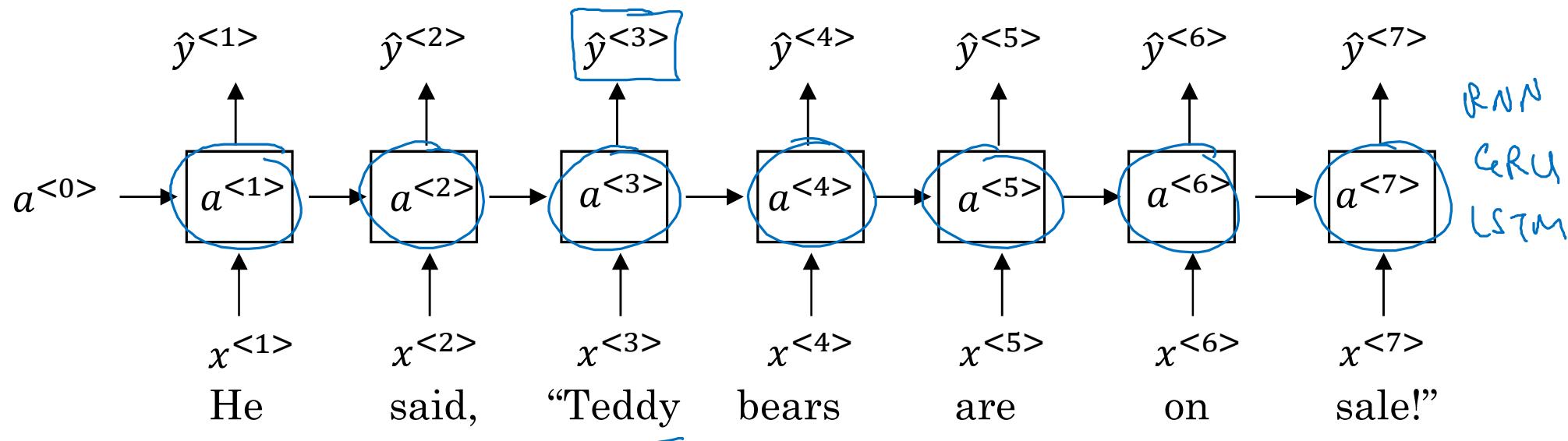
Recurrent Neural Networks

Bidirectional RNN

Getting information from the future

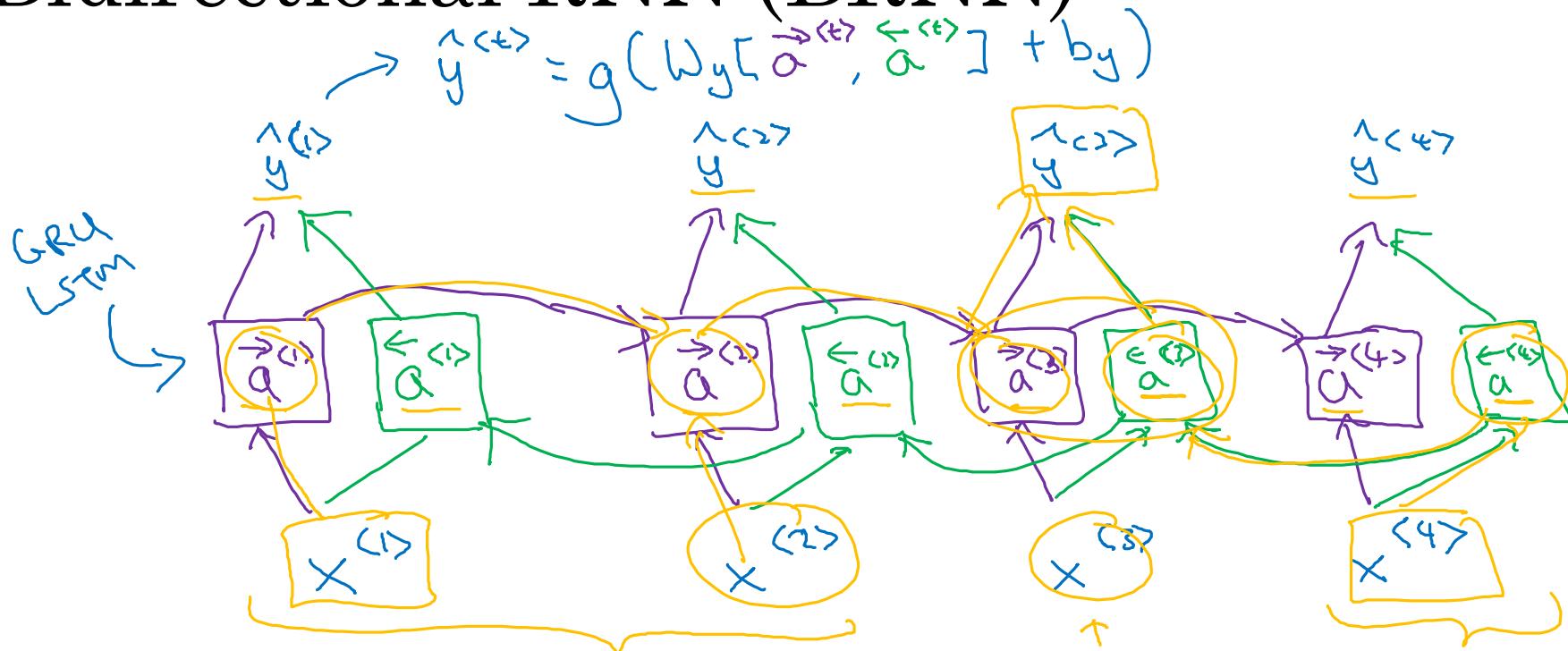
He said, “Teddy bears are on sale!”

He said, “Teddy Roosevelt was a great President!”



Andrew Ng

Bidirectional RNN (BRNN)



Acyclic graph

BRNN w/ LSTM

He said

"Teddy Roosevelt ..."

Andrew Ng

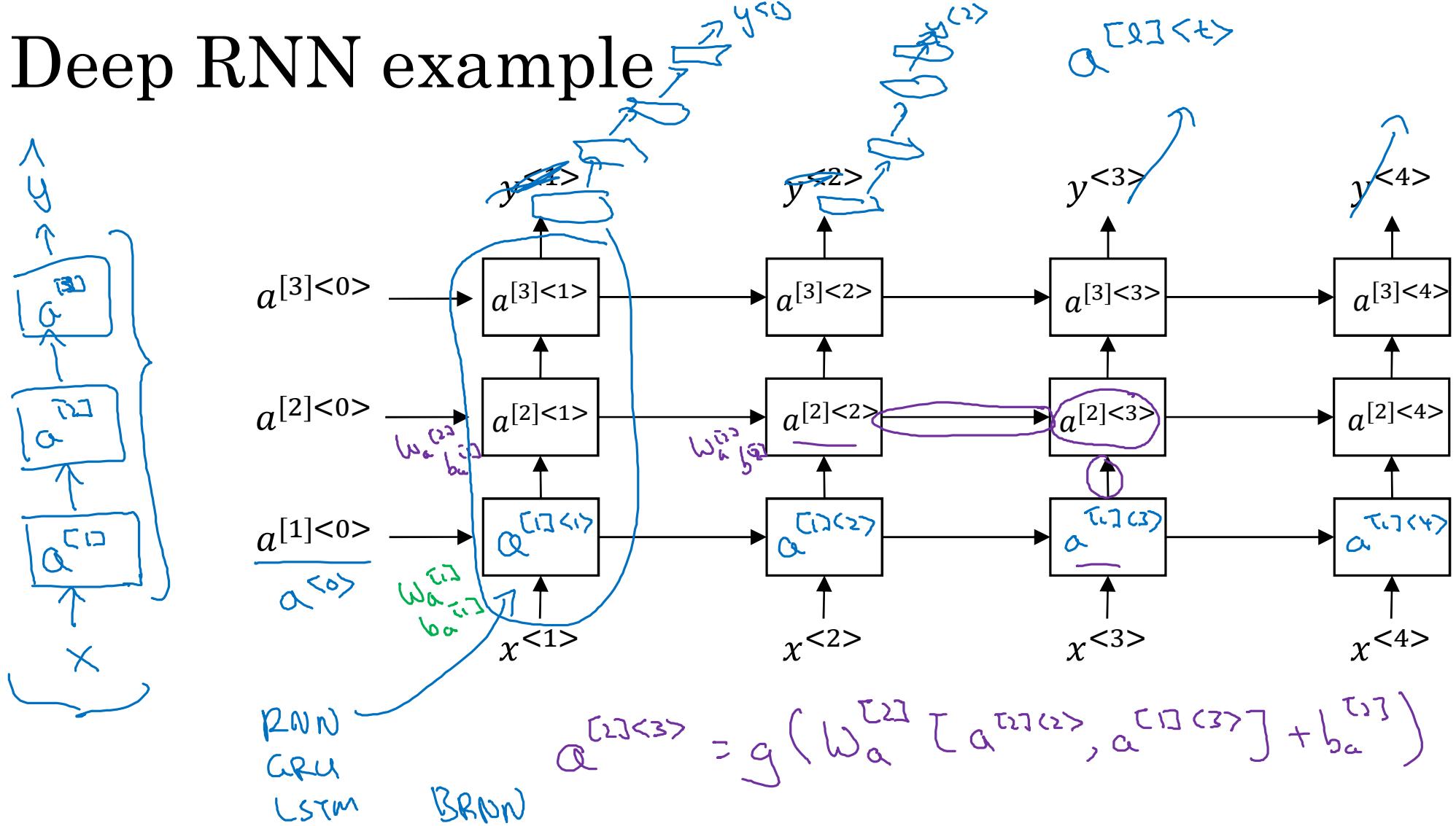


deeplearning.ai

Recurrent Neural Networks

Deep RNNs

Deep RNN example



Andrew Ng