

There are a lot of different sequence models

# Recurrent Neural Networks

---

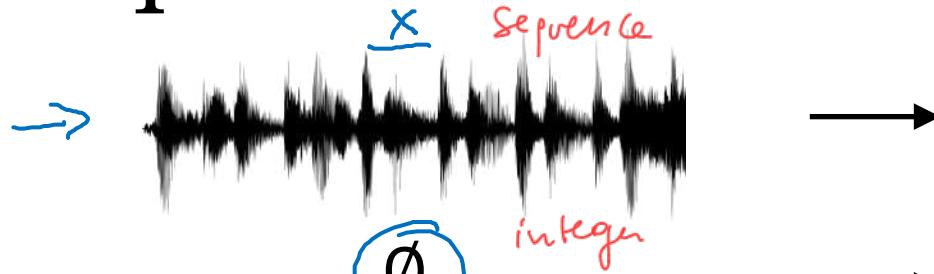


deeplearning.ai

Why sequence  
models?

# Examples of sequence data

Speech recognition

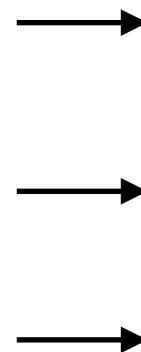


Music generation



Sentiment classification

"There is nothing to like  
in this movie."

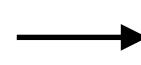


DNA sequence analysis → AGCCCCTGTGAGGAAC TAG

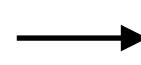


Machine translation

Voulez-vous chanter avec  
moi?



Video activity recognition



Name entity recognition

Yesterday, Harry Potter  
met Hermione Granger.



Yesterday, Harry Potter  
met Hermione Granger.  
Andrew Ng



deeplearning.ai

# Recurrent Neural Networks

---

## Notation

# Motivating example

NLP

$(x, y)$  = example and label

$T_x^{(i)}$  and  $T_y^{(i)}$  = no. of tokens in i-th sample

$x^{(i), t}$  = t-th token in i-th input sample

x:

Harry Potter) and (Hermione Granger) invented a new spell.

$\rightarrow \underline{x^{(1)}}$   $x^{(2)}$   $x^{(3)}$  - - -  $x^{(t)}$  - - -  $x^{(q)}$

$$T_x = q$$

$\rightarrow y:$

| | 0 | | 0 0 0 0  
 $y^{(1)}$   $y^{(2)}$   $y^{(3)}$  - - -  $y^{(q)}$

$$T_y = q$$

$x^{(i), t}$

$y^{(i), t}$

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

$$15$$

$$T_y^{(i)}$$

# Representing words

$x^{<t>}$

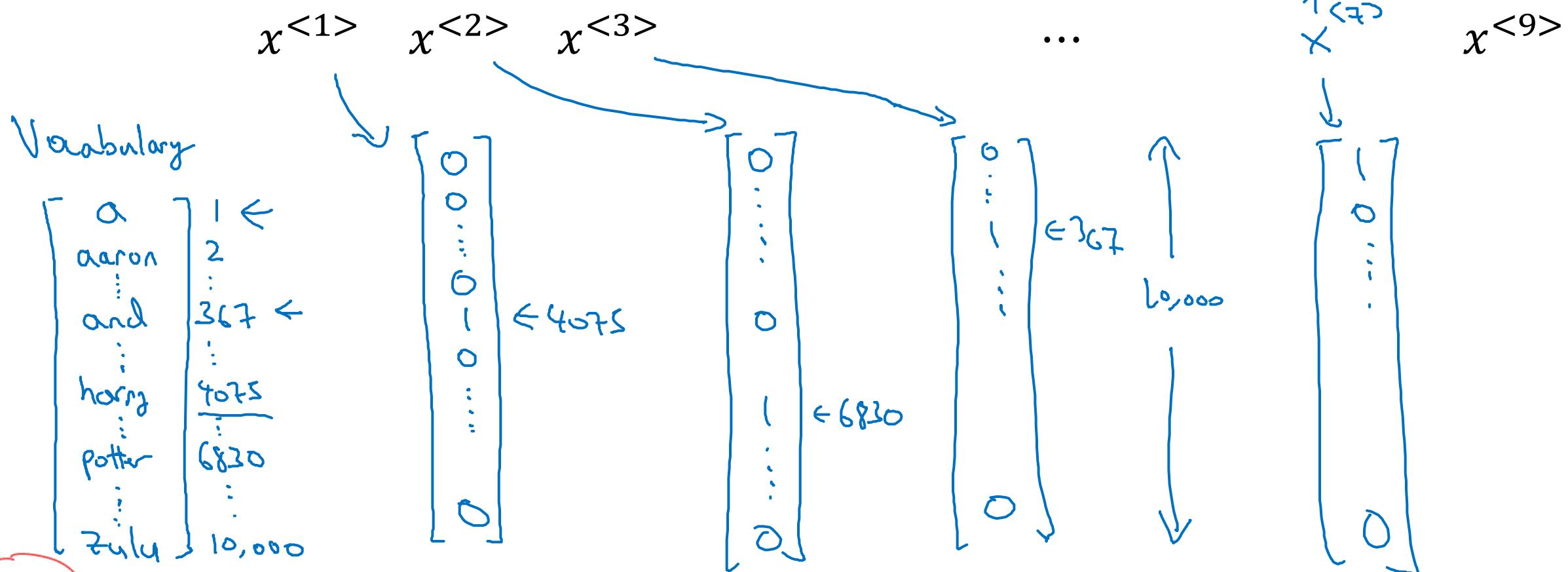
$x \rightarrow y$

$(x, y)$

In NLP, one way to encode a sample (sentence) is by one-hot encoding its tokens based on a vocabulary.

$x:$

Harry Potter and Hermione Granger invented a new spell.



10,000 → Vocabulary size in commercial applications is ~50k. Exceptional models move to 1M words.

<UNK>

→ Vocabulary entry to describe unknown words

One-hot

# 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



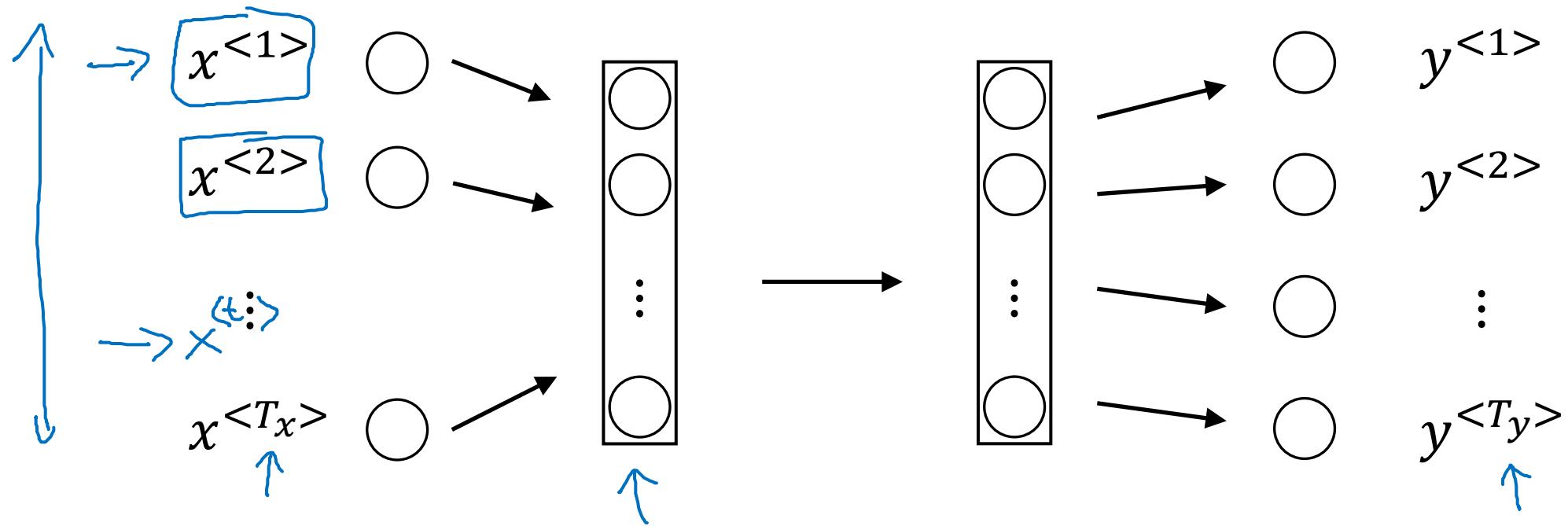
deeplearning.ai

# Recurrent Neural Networks

---

## Recurrent Neural Network Model

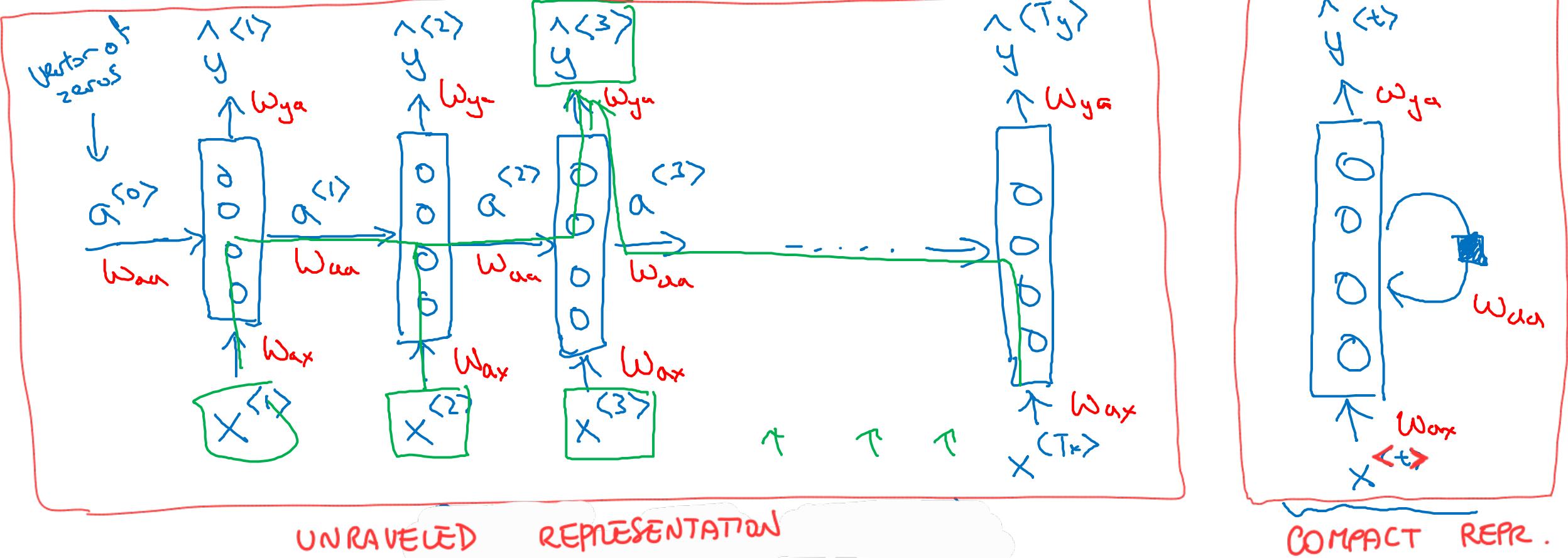
# Why not a standard network?



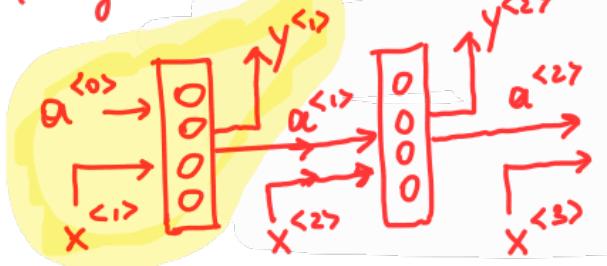
## Problems:

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

# Recurrent Neural Networks



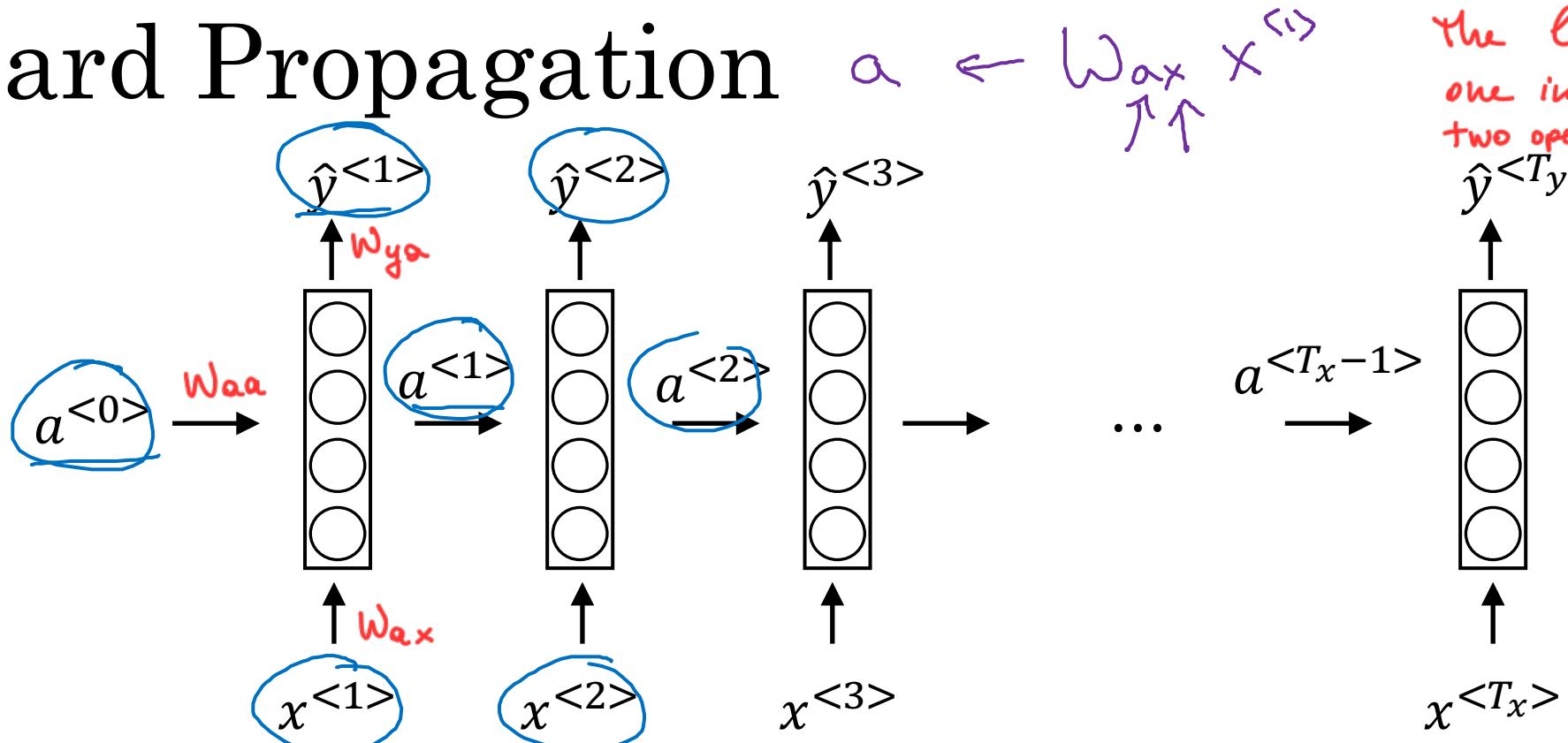
I imagine it like this (unraveled):



in yellow is the network

- the basic structure can be represented in unrolled (left) and compact (right) way.
- look at the compact representation. This is a RNN with 1 layer and 4 hidden units.  
 (The non-RNN equivalent would be  $x^{<i>} \rightarrow \boxed{0} \rightarrow y^{<i>} = g(W_x x^{<i>} + b_x)$ ). The RNN peculiarity is that each hidden unit feeds back to itself: the activation of unit  $j$  at time  $k$  at time  $i$ ,  $a^{[k]}_{\{j\}<i>}$ , becomes one input of the same unit at time  $i+1$ .  
depends on the number of neurons  
 NOTE THAT the unit has an "hidden state"  $\tilde{z}^{<i>} = W_x x^{<i>} + b_x$  and its dimension  $\tilde{z}$  differs in general from the dimension of  $x^{<i>}$ . With one-hot encoding  $x$  can be in  $\{0,1\}^{10000}$  but if the unit has only  $h$  neurons then the hidden state is also  $h$ -dimensional ( $W$  is  $4 \times 10000$ ).
- In the unrolled representation the same network is replicated horizontally at different points in time. The tail of the network is at the bottom, the head at the top.

# Forward Propagation



The layer (it's only one in this case) computes two operations ( $a^{<i>}$  and  $y^{<i>}$ ). Note that in non-RNNs the layer performed only one equation  $y^{<i>} = g(a^{<i>})$ . Here two ops are needed because the activation is more complicated (not just  $Wx$  but  $g(W_x x + W_a a + b_a)$ )

$$a^{<0>} = \vec{0}.$$

$$\underline{a^{<i>} = g_i(W_{aa} a^{<0>} + \underline{W_{ax}} x^{<i>} + b_a)} \quad \leftarrow \text{tanh / ReLU}$$

$$\underline{\hat{y}^{<i>} = g_2(W_{ya} a^{<i>} + b_y)} \quad \leftarrow \text{Sigmoid}$$

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

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

# Simplified RNN notation

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

Dimensions:  
 $W_{aa}$ :  $(100, 100)$   
 $x^{(t)}$ :  $(100, 10,000)$

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

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

$W_y$ :  $(100, 100)$   
 $a^{(t)}$ :  $(100)$   
 $b_y$ :  $(100)$

by stacking the layer's inputs at time  $t$ , one gets the simplified notation:

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

$W_a$ :  $(100, 100)$   
 $[a^{(t-1)}, x^{(t)}]$ :  $(100, 100)$   
 $b_a$ :  $(100)$

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



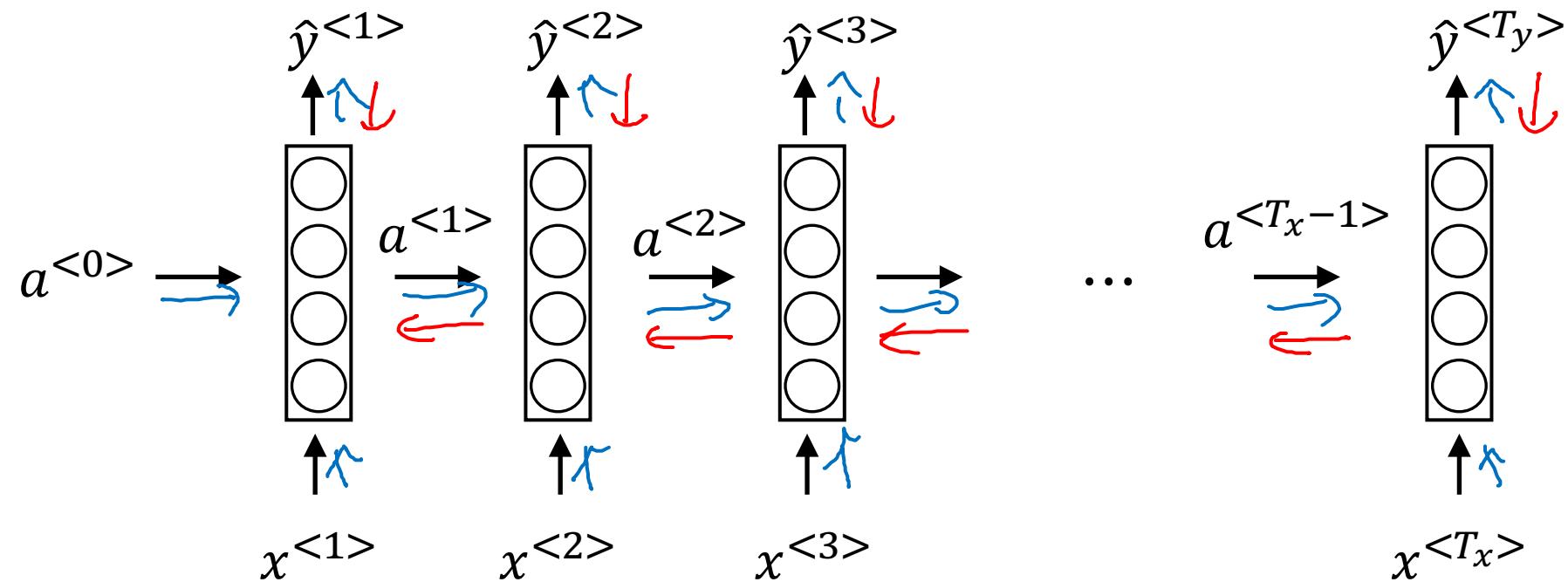
deeplearning.ai

# Recurrent Neural Networks

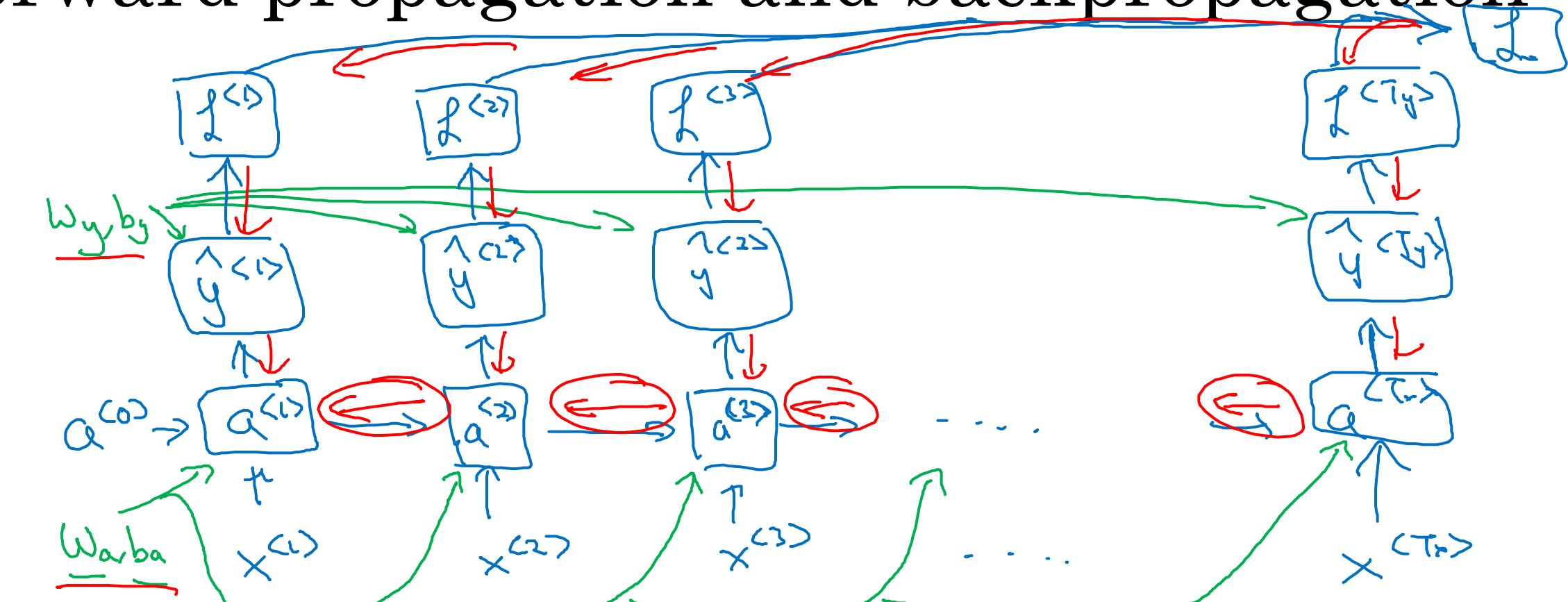
---

Backpropagation  
through time

# Forward propagation and backpropagation



# Forward propagation and backpropagation



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

$$\text{OVERALL LOSS} = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

LOSS FOR ELEMENT T OF THE INPUT SEQUENCE

The loss for a certain input sequence is the summed loss of the single inputs of the sequence.

**Backpropagation through time**



deeplearning.ai

# Recurrent Neural Networks

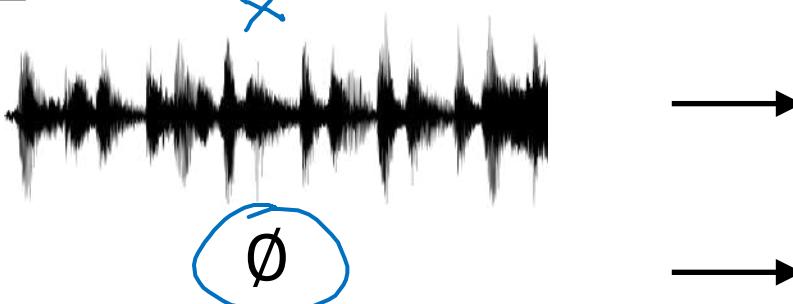
---

## Different types of RNNs

Different architecture to tackle different types of problems. The main difference between problems is the number of inputs and outputs.

# Examples of sequence data

Speech recognition



$$T_x \quad T_y$$

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



AGCCCCTGTGAGGAAC 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.

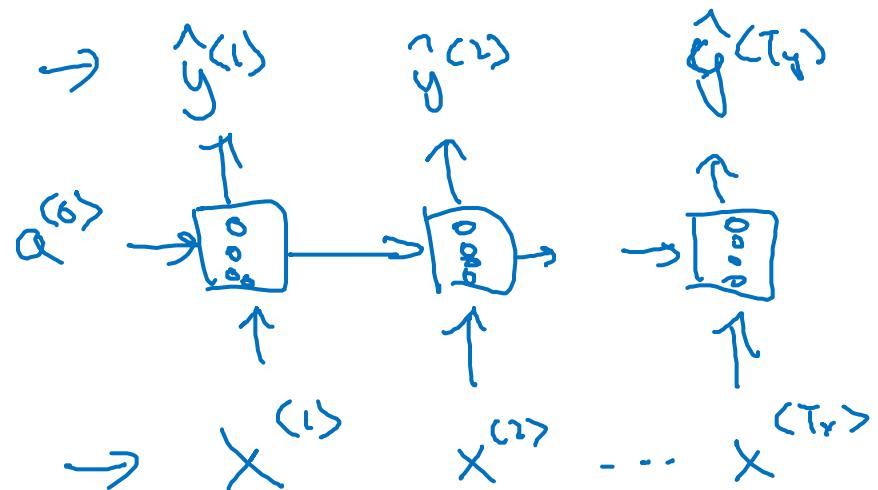


Yesterday, Harry Potter  
met Hermione Granger.

Andrew Ng

# Examples of RNN architectures

$$T_x = T_y$$



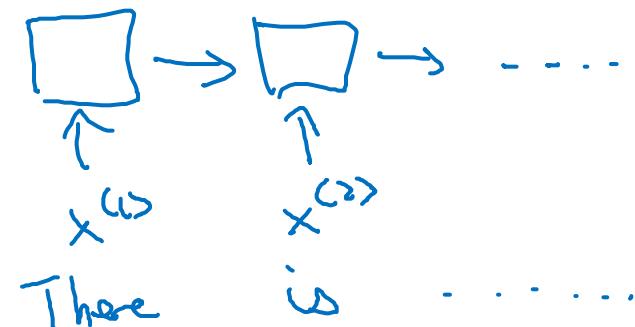
Many-to-many

{ many inputs and many outputs  
no inputs = no outputs

Sentiment classification

$$x = \text{text}$$

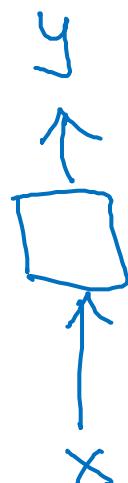
$$y = 0/1 \quad 1 \dots 5$$



Many-to-one

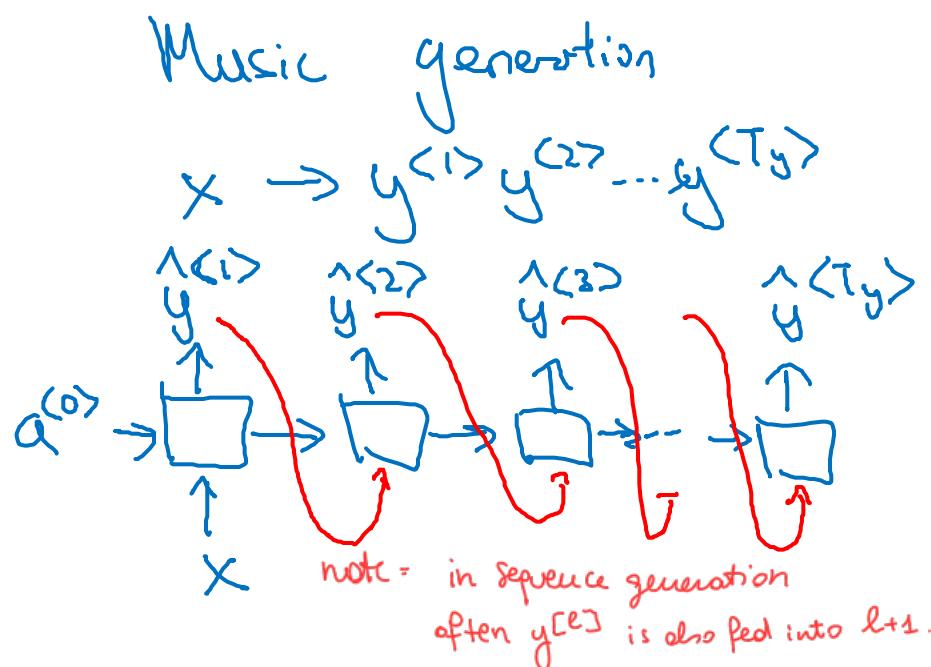
{ many inputs  
only one output

In this frame,  
a standard NN  
can be defined  
as one-to-one



One-to-one

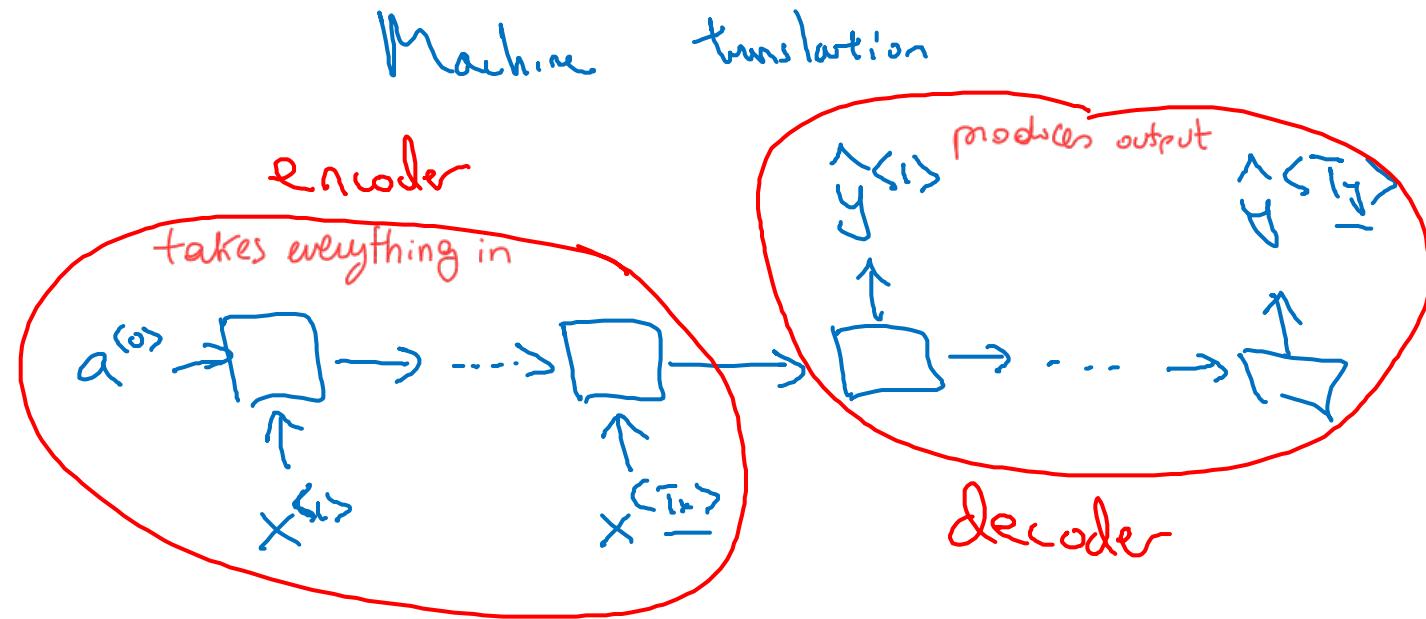
# Examples of RNN architectures



One-to-many

$$x = \phi$$

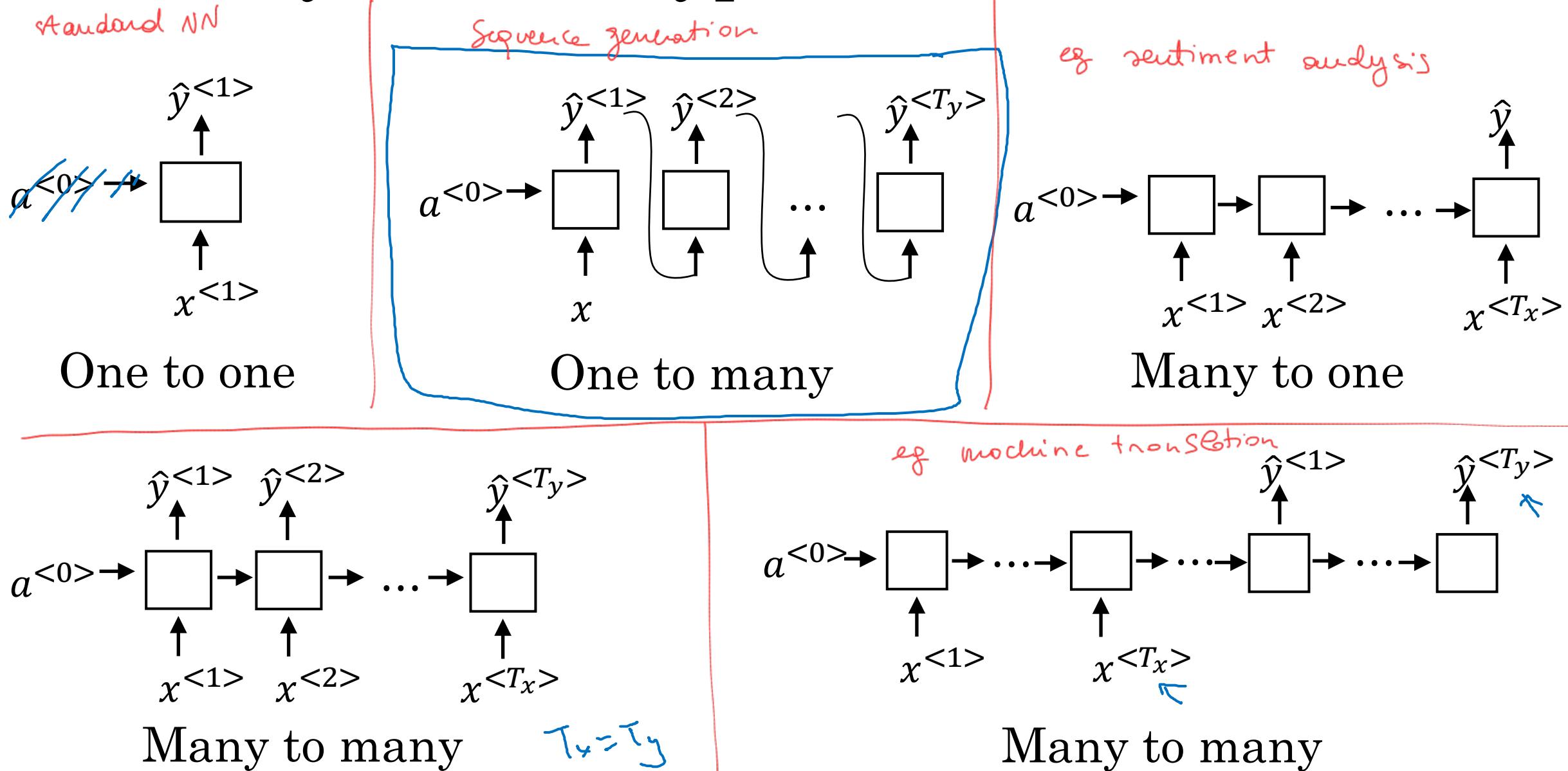
A FINAL TYPE OF ARCHITECTURE IS THE TRANSFORMER



Many - to - many

different number of inputs and outputs

# Summary of RNN types





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.

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

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

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

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

language modelling = establishing the probability of a sequence of words for a certain language.

# Language modelling with an RNN

The first step is to tokenize the sentence = mapping each element of the sentence with an individual token (element of the vocabulary) of the language.

Training set: large corpus of english text.

Tokenize

The vocabulary can include special tokens like <UNK>, <EOL>, <·>, ...

Cats average 15 hours of sleep a day.  $\downarrow$  <EOS>

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

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

...

$y^{(8)}$        $y^{(9)}$

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

<UNK>

10,000

## language modelling task

- DEFINITION: determining the probability of a series of words in a certain language.
- STEP 1: TOKENIZATION AND MAPPING: Isolating the relevant elements in the input sentence (words, relevant punctuation EOL if desired, etc), that is the "tokens", and mapping each of these with an entry of the vocabulary of the given language.
- STEP 2: PREDICTING PROBABILITY of the sequence of words.  
(Andrew Ng provided a somewhat confusing explanation).  
For each token in the sequence, the network outputs in order

$$P(1^{\text{st}} \text{tok}), P(2^{\text{nd}} \text{tok} | 1^{\text{st}} \text{tok}), \dots, P(t^{\text{th}} \text{tok} | (t-1)^{\text{th}} \text{tok})$$

that is

$$\hat{y}^{<1>} = P(y^{<1>}), \hat{y}^{<2>} = P(y^{<2>} | y^{<1>}), \dots, \hat{y}^{<t>} = P(y^{<t>} | y^{<1>} , y^{<2>} , \dots , y^{<t-1>}).$$

The probability of the sentence is obtained as the product of those probabilities (... although this assumes independence of the elements in the sequence...).

$$\begin{aligned} P(y^{<1>} , \hat{y}^{<2>} , \dots , \hat{y}^{<t>}) &= P(y^{<1>}) \cdot P(\hat{y}^{<2>} | y^{<1>}) \cdot \dots \cdot P(\hat{y}^{<t>} | y^{<1>} , y^{<2>} , \dots , y^{<t-1>}) \\ &= \hat{y}^{<1>} \cdot \hat{y}^{<2>} \cdot \dots \cdot \hat{y}^{<t>} \end{aligned}$$

- ACTUALLY: the NN does not return only one probability but rather a softmax over all the elements of the vocabulary

$$\hat{y}^{<1>} = [P(v_1), P(v_2), \dots, P(v_{10000})], \hat{y}^{<2>} = [P(v_1 | y^{<1>}), \dots, P(v_{10000} | y^{<1>})], \hat{y}^{<3>} = \dots$$

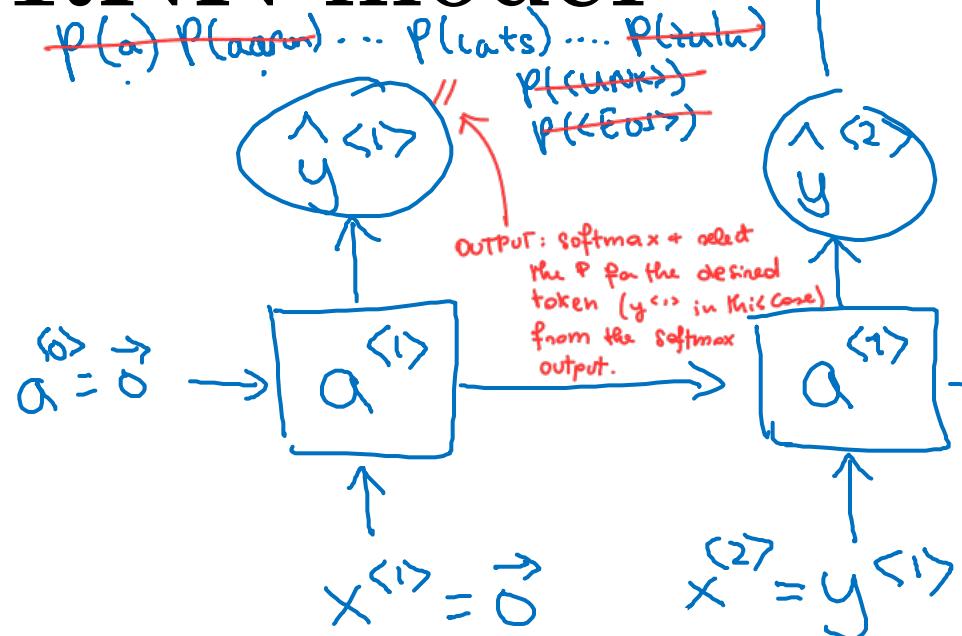
(although Andrew Ng writes  $\hat{y}^{<t>} = \sqrt[10000]{\max([P(v_i | y^{<1>} , \dots , y^{<t-1>}])_{i=1}^{10000})}$ .)

- the NN's cost function uses a cross entropy loss that compares the predicted softmax to the ground truth probability distib of the desired word wnt vocabulary.

example The cat is on the table

$$y^{<2>} = \text{cat} \dots \text{ectrely} \dots = \begin{bmatrix} \frac{1}{10000} & \dots & \frac{1}{10000} \end{bmatrix}_{\text{parts of Aaron}}^{\text{prob of cat}}, \hat{y}^{<2>} = [10^{-13} \dots 10^{-4} 0.6 10^{-2} \dots 10^{-6}], L^{<2>} = \sum_{i=1}^{10000} y_i^{<2>} \log \hat{y}_i^{<2>}$$

# RNN model

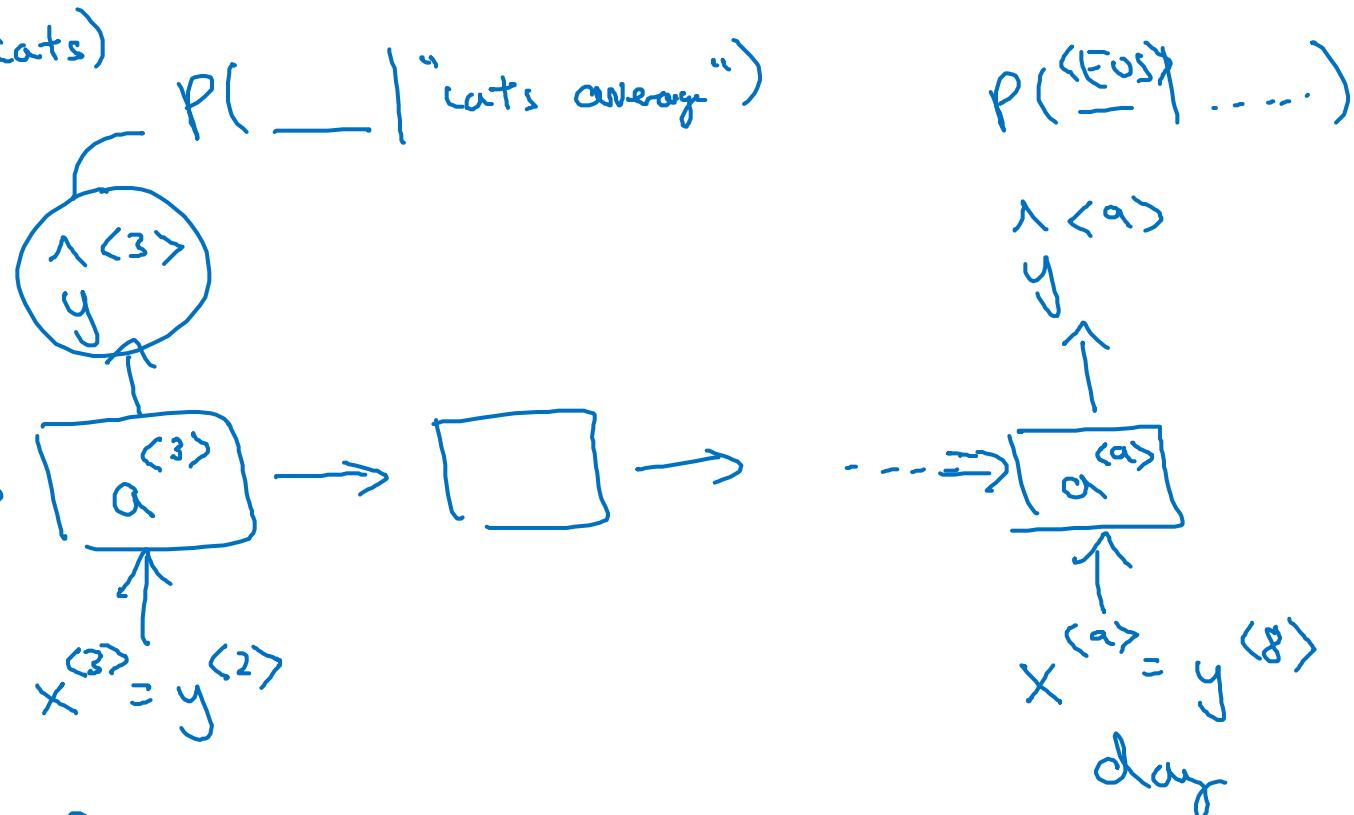


Cats average 15 hours of sleep a day. <EOS>

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

$\mathcal{L} = \sum_t \mathcal{L}(\hat{y}^{(t)}, y^{(t)})$

*LOSS FUNCTION*



$$\begin{aligned}
 P(y^{(1)}, y^{(2)}, y^{(3)}) &\leftarrow \\
 &= \frac{P(y^{(1)}) P(y^{(2)} | y^{(1)})}{P(y^{(3)} | y^{(1)}, y^{(2)})}
 \end{aligned}$$

Once the NN has been fitted to a language model, it can be used for sampling a sequence:

- interpret the output:  $\text{output}^{<t>} = \text{vocab. entry corresponding to } \max(\hat{y}^{<t>})$ . e.g. "cat"
- use the interpreted input as new entry in the NN:  $x^{<t+1>} = \text{output}^{<t>}$
- terminate the generation after a fixed number of tokens  
or when <EOL> is generated.

- implement fall-back strategies for when <UNK> is generated.



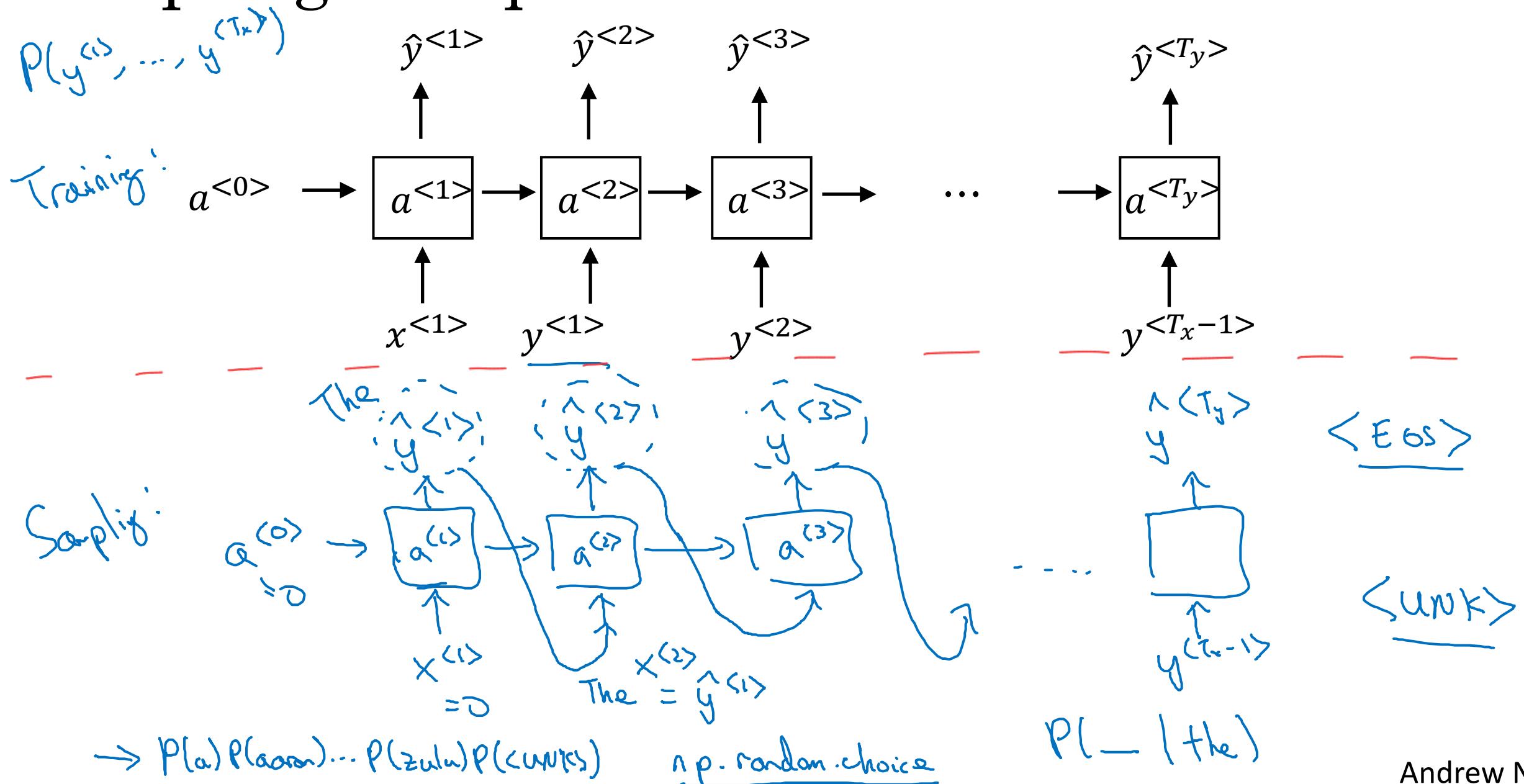
deeplearning.ai

# Recurrent Neural Networks

---

## Sampling novel sequences

# Sampling a sequence from a trained RNN



# Character-level language model

- (as opposed to language-level language model). One advantage is not requiring an `<UNK>` character, but the big disadvantage is having to be trained on much longer sequences.

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

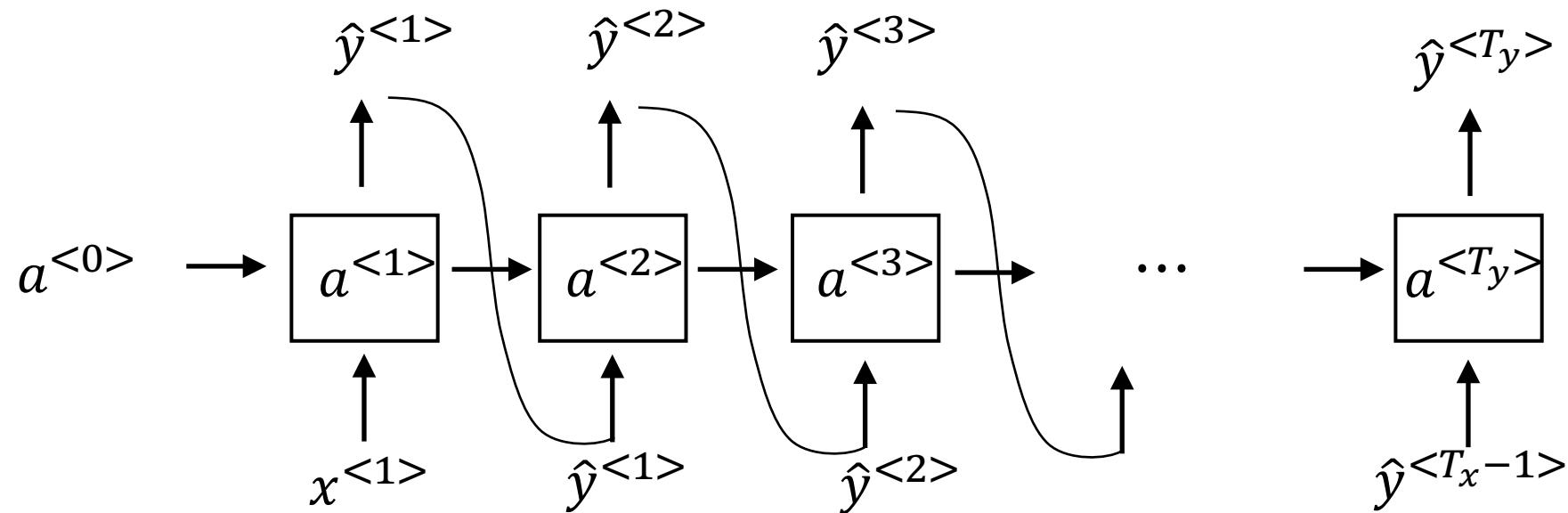
One advantage is not repiving an  character, but the big  disadvantage is having to be trained on much longer sequences.

$\rightarrow \text{Vocabulary} = [a, b, c, \dots, z, \cup, \circ, \rightarrow, ;, \circlearrowleft, \circlearrowright, \circlearrowuparrow, \circlearrowdownarrow, A, \dots, Z]$

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

Cat average  
↑ ↑ ↑ ↑ ...

May



# 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.



deeplearning.ai

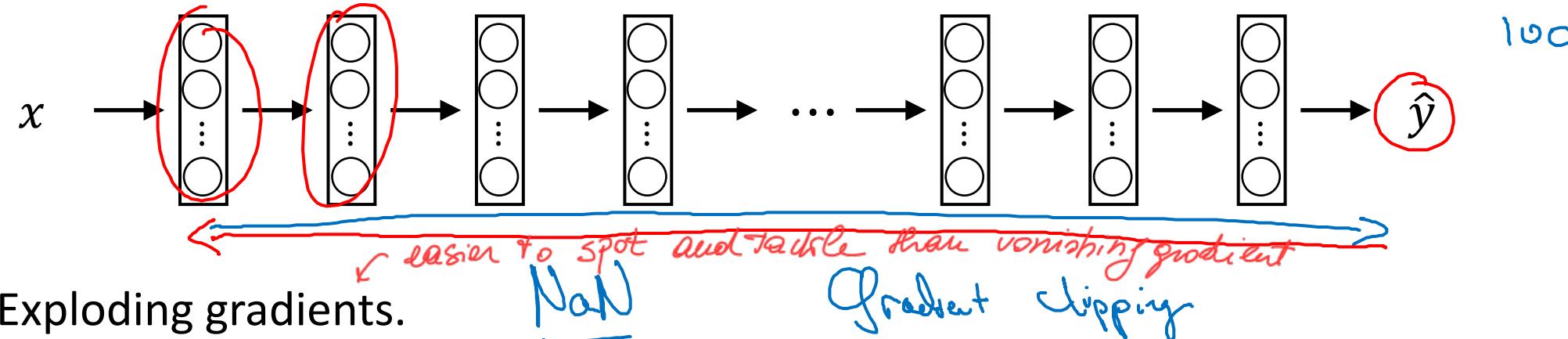
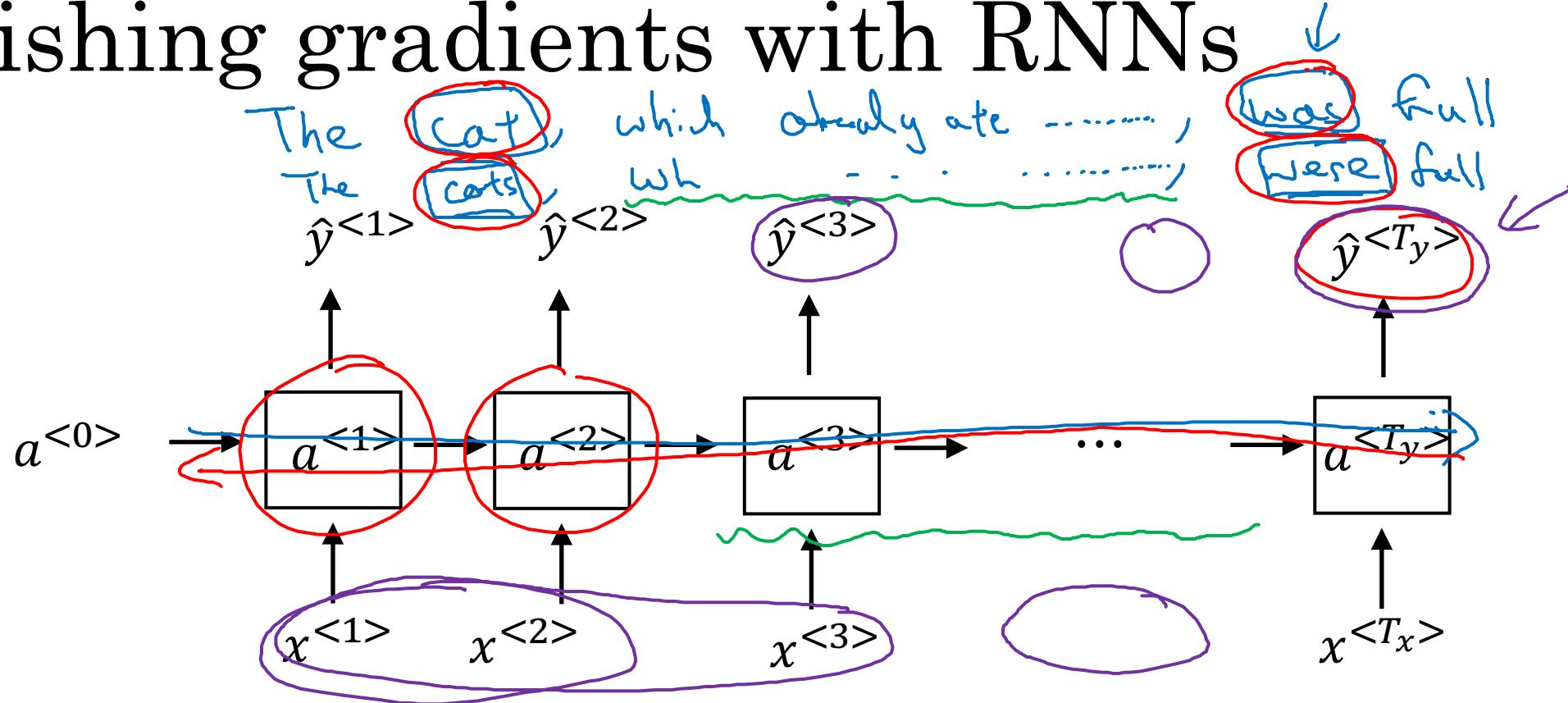
# Recurrent Neural Networks

---

## Vanishing gradients with RNNs

RNNs are prone to vanishing gradients in case they are used to handle long sequences.  
=> RNNs are not at capturing long term dependencies

# Vanishing gradients with RNNs



A GRU is a modification of the standard Recurrent Unit that produces outputs (activations and updates) with some "adaptive momentum". The activation at time  $t$  is the weighted sum of the hidden state  $c^{<t>}$  of the unit and the old activation  $c^{<t-1>}$  using a "gate function" as weight. The gate function  $\Gamma$  depends from

$c^{<t-1>}$  and  $x^{<t>}$ .

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

The idea is that  $\Gamma$  allows to keep track of relevant tokens across very long sentences.  $\Gamma$  is set (automatically) close to 1 when the input is important for the sentence and then is near zero for the remaining terms until another important word appears. When  $\Gamma$  is zero, it forces the output to strongly relate to the word for which it was last set to 1.

"The cat, which ate the soup and the apple, was fat and the child ..."

$\Gamma = 0$	1	0	0	0	0	0	0	0	0	0	1	1
--------------	---	---	---	---	---	---	---	---	---	---	---	---

# Gated Recurrent

deeplearning.ai

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

$$n^{<t>} = \tanh(W_a * [c^{<t-1>}, x^{<t>}] + b_a)$$

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

NOTE:

$c^{<t>} = n^{<t>}$  for GRU but  
this is not the case for LSTM units.

NOTE:

$\tilde{c}^{<t>}$

update gate: in more advanced GRUs there are more types of gates.

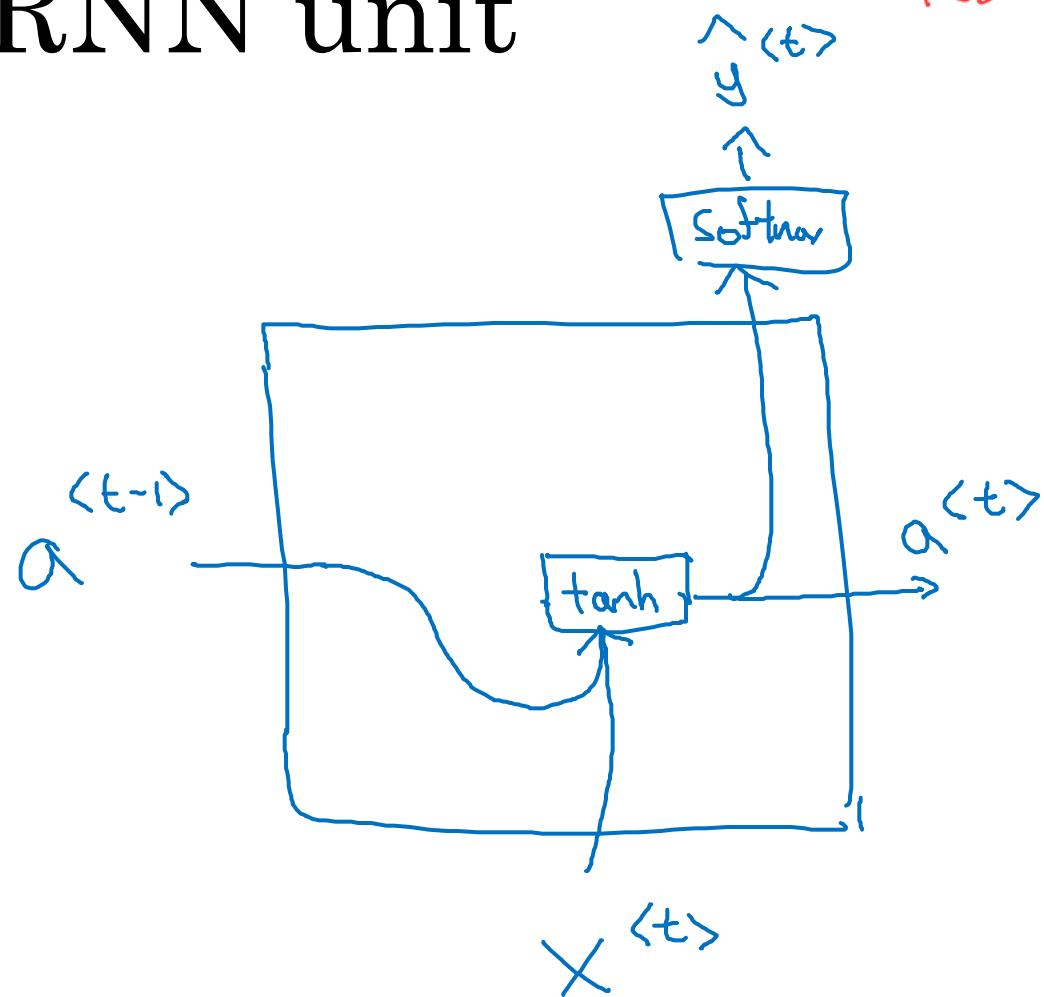
An example is the relevance gate  $\Gamma_r$ , which helps forgetting the previous hidden state (unseen).

# Recurrent Neural Networks

NOTE: GRU came often LSTM, as a simplification of LSTM and often work as well.

# Unit (GRU)

# RNN unit



TLDR: differently from RV, the activation  $c^{(t)}$  is mixed with the old activations using a gate function as weight. The gate function allows to maintain the influence of old activations high for a long time. In practice it separates the hidden state ( $\tilde{c}$ ) from the actual activation  $c$ .

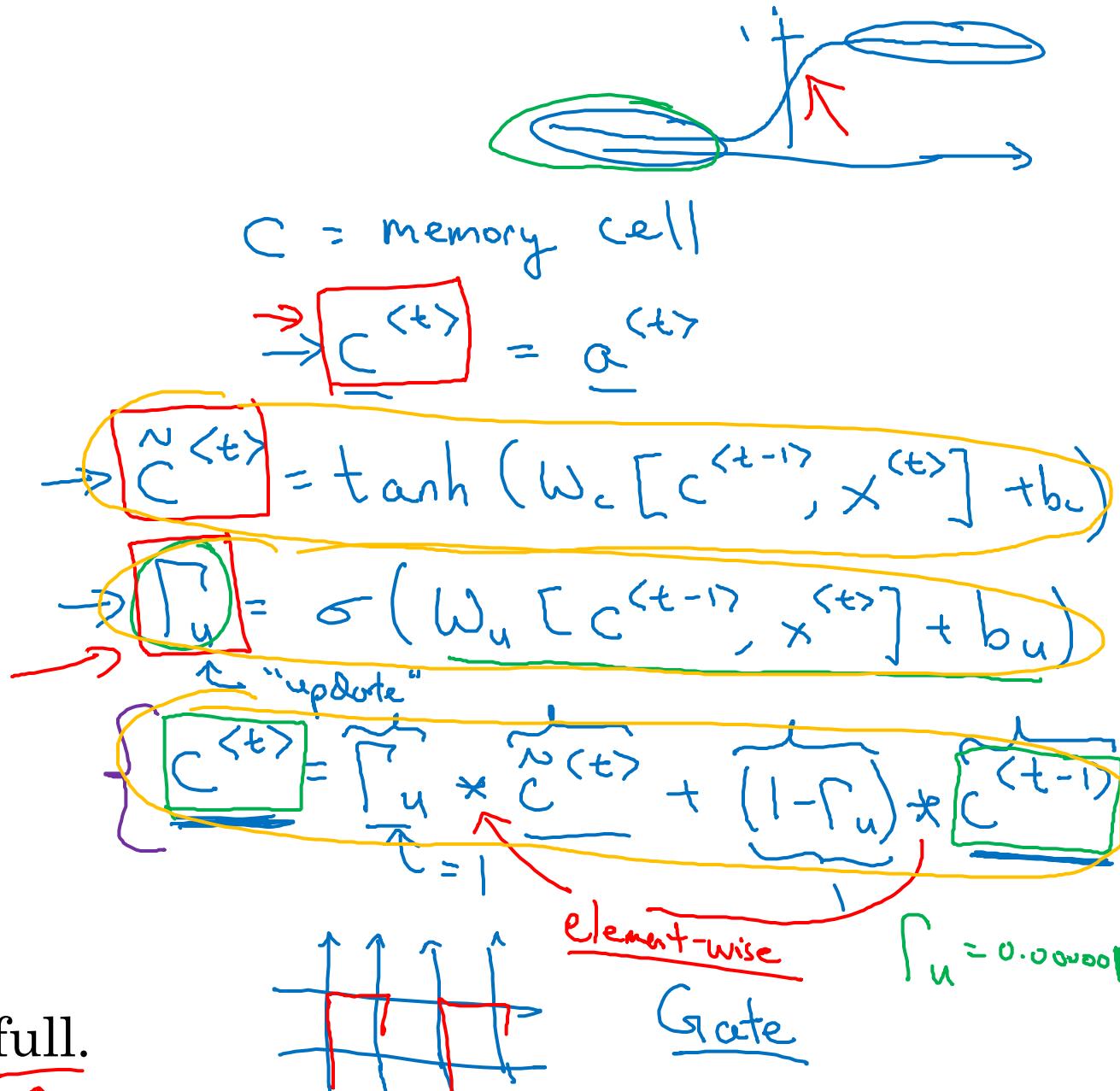
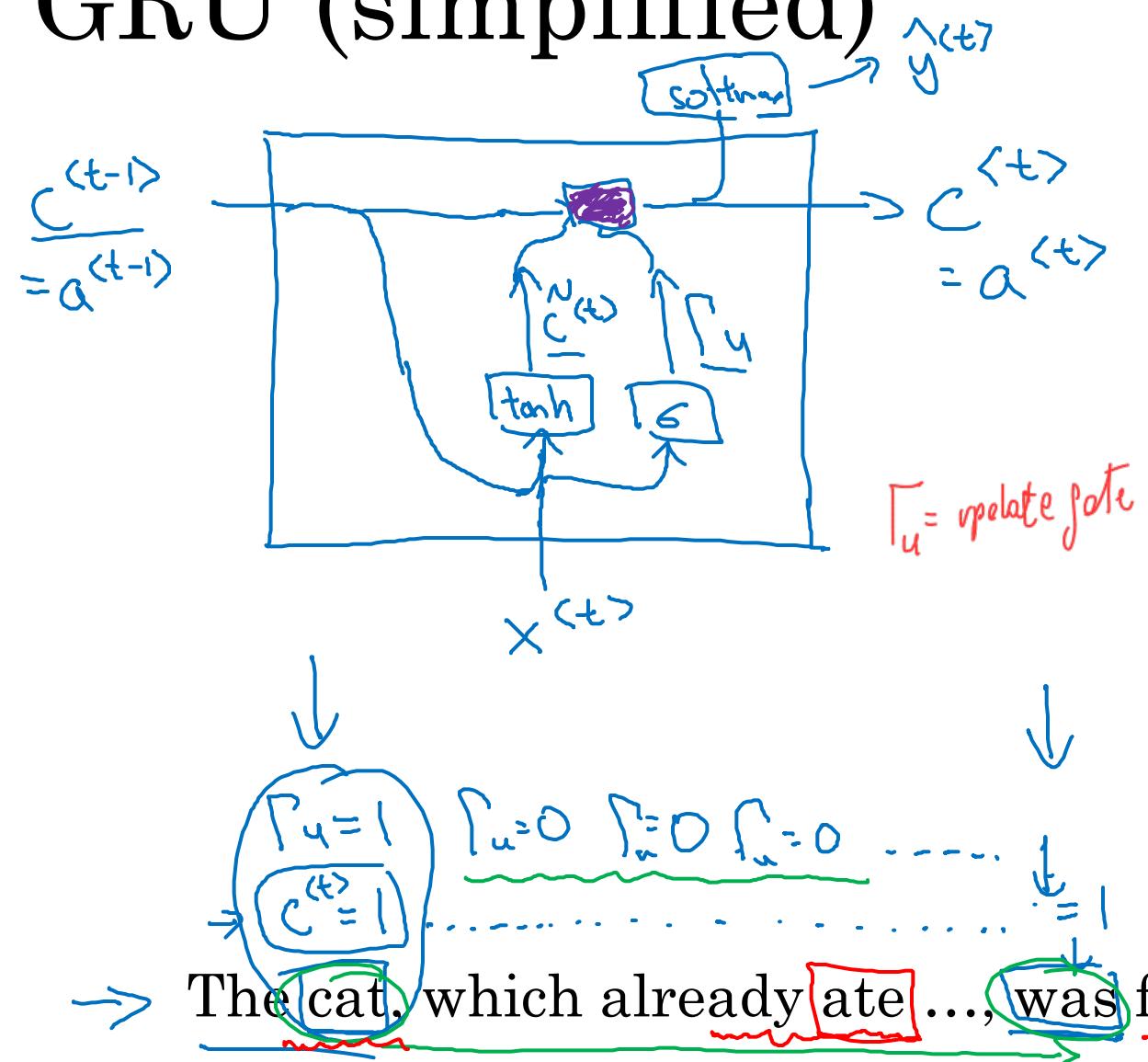
$$\underline{a}^{(t)} = g(W_a[\underline{a}^{(t-1)}, \underline{x}^{(t)}] + b_a)$$

$\text{tanh}$

↓

— ↑ —

# GRU (simplified)



[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

$\Gamma_r$  = relevance gate

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r^{<t-1>} \cdot \underline{\underline{x}}^{<t>}] + b_c) \leftarrow \tanh([\Gamma_r \cdot c^{<t-1>} \cdot \underline{\underline{x}}^{<t>}] + b_c)$$
$$\begin{cases} \Gamma_u = \sigma(W_u[c^{<t-1>} \cdot \underline{\underline{x}}^{<t>}] + b_u) \\ \Gamma_r = \sigma(W_r[c^{<t-1>} \cdot \underline{\underline{x}}^{<t>}] + b_r) \end{cases}$$

$\Gamma_r$  ← IGNORED IN LSTMS !

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

↑  
alternative notation in literature

The cat, which ate already, was full.

LSTM units are a more general version of GRU that have 3 gates (input, forget, output) instead of just one (update). The idea is to create a separation between the input of the cell and the output in the form of an extra memory cell  $c^{<t>}$  called "cell state", which is distinct from the output  $a^{<t>} = \tilde{c}^{<t>}$  (in GRU the output =  $a^{<t>} = c^{<t>} = \Gamma \tilde{c}^{<t>} + (1-\Gamma) c^{<t-1>}$ ). So in GRU we have hidden state  $\tilde{c}$  and activation  $c=a$ , while in LSTM we have hidden state, cell state and output activation.

- LSTM is more general because has 2 gates  $i_f$  that regulate the

information flow in and out of the cell state  $c$  (while GRU only has one  $\Gamma_u$  to regulate the input of new info  $\tilde{c}$  into the activation  $c=a$ ).

- Moreover LSTM has a dedicated Forgetting gate  $f_f$  to regulate the contribution of memory (old cell state  $c^{<t-1>}$ ) to the current cell state. (Actually they call it forget gate, but "memory gate" would be more precise).



deeplearning.ai

# Recurrent Neural Networks

TLDR : • more general/flexible version of GRU (also older). SEE FORMULAS IN 2 SLIDES

- it has 3 gates (input = update, output, forget) instead of only one (update).
- separate "cell state" memory cell from hidden state and output.
- 2 gates (in, out) regulate information flow through cell state.
- 1 gate (update) regulate influence of old cell state on new one.

## LSTM (long short term memory) unit

- in general one should try both GRU (lighter) and LSTM. If no time to check, pick LSTM.

# GRU and LSTM

GRU vs LSTM

## GRU

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

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

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

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

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

$\Gamma_f$

## LSTM

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

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

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

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

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

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

# 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

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

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

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

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

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

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

# LSTM in pictures

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

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

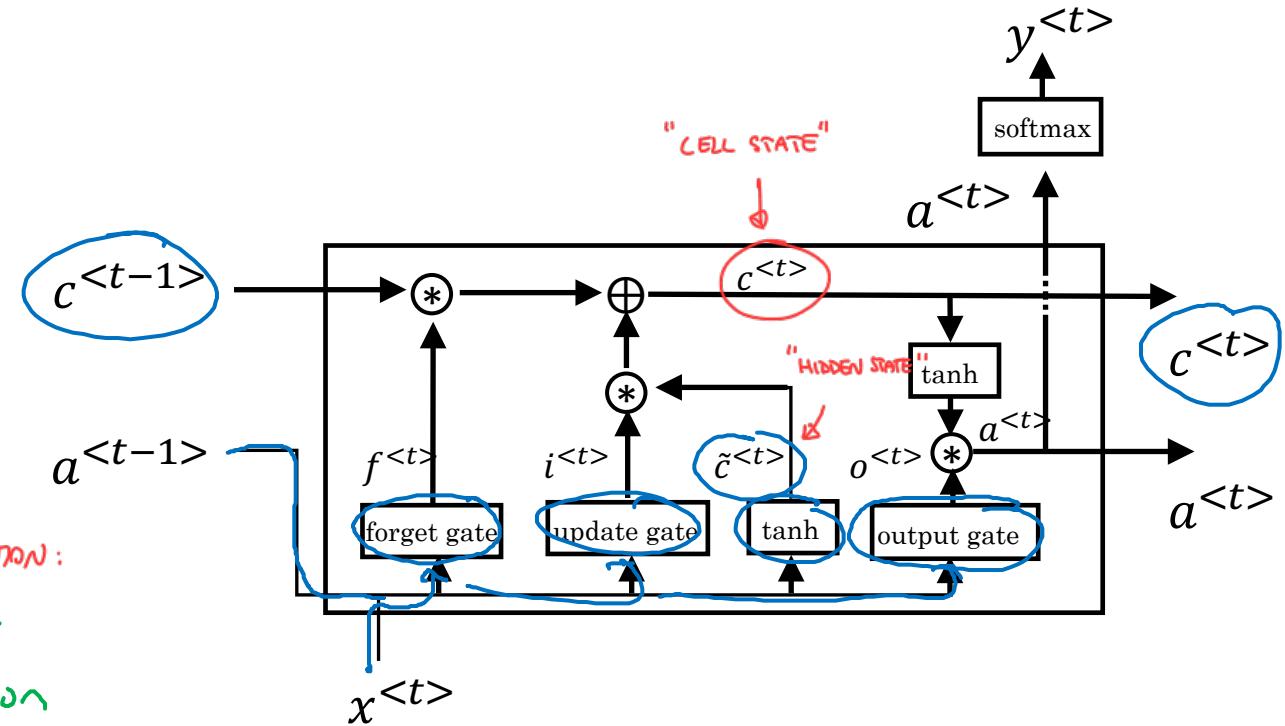
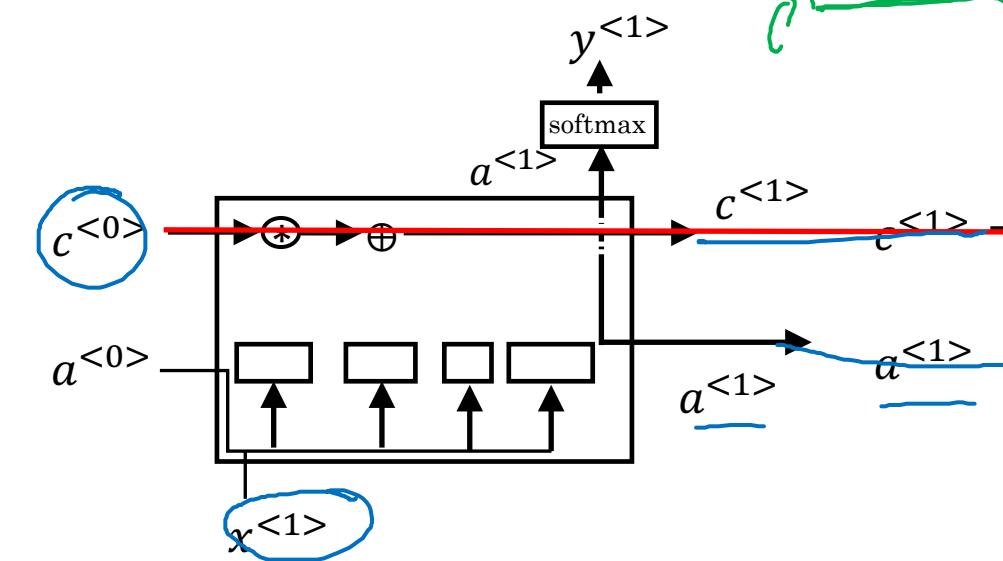
$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

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

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

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

*COMMON VARIATION: peephole connection*



Andrew Ng

BRNN are used to make predictions also based to information in the future (for example useful for named entity recognition). They do so:

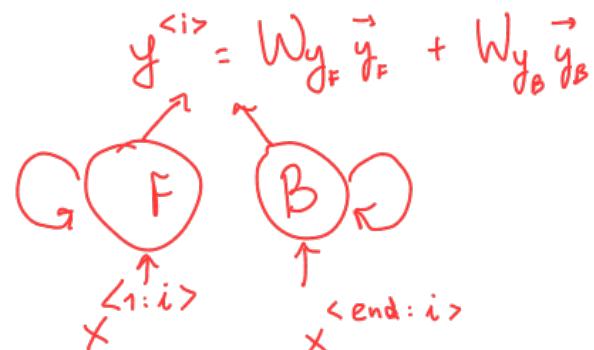
- by having access to the entire input sequence (no online)



- by having one "forward" recurrent unit that processes the input from the beginning and one "backward" recurrent unit that processes the input from the end token. These two units produce a unique output by weighting forward and backward prediction (with learnt weights).

# Recurrent Neural Networks

## Bidirectional RNN

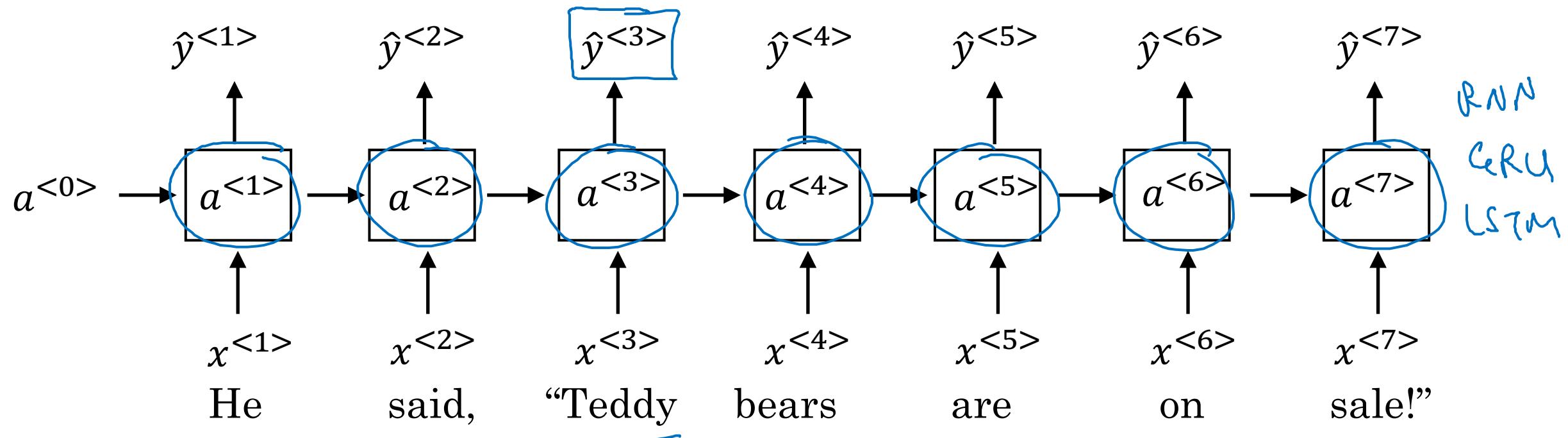


# Getting information from the future

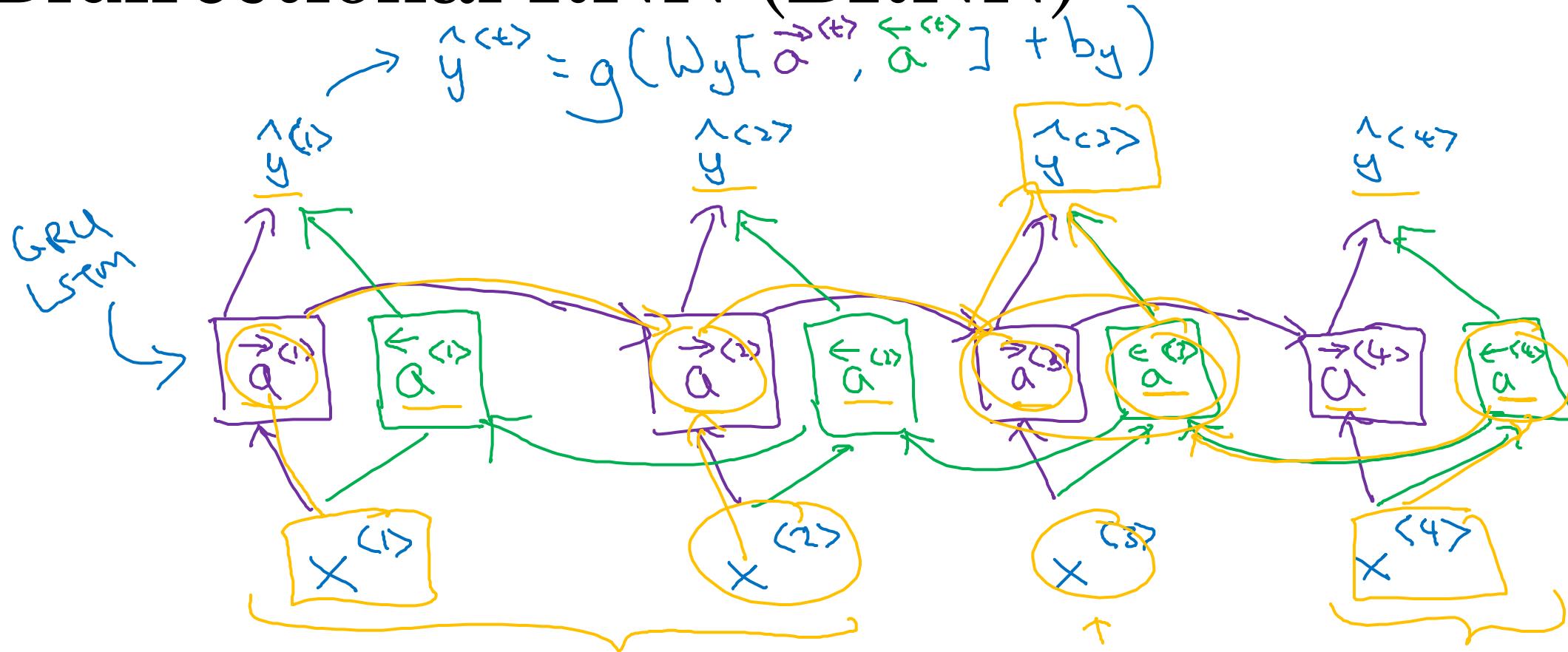
He said, “Teddy bears are on sale!”

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

*Application = named entity recognition*



# Bidirectional RNN (BRNN)



Acyclic graph

BRNN w/ LSTM

He said,

"Teddy Roosevelt ..."

Obtained stacking recurrent units (RU, GRU, LSTM units).

- "Vertical" stacking from bottom to top ; usual "horizontal" time connection.
- Similarly to before, each



unit has 2 inputs. Input units like  $a^{[1] < i>}$  have inputs  $x^{< i>}$  and  $a^{[1] < i-1>}$

# Recurrent Neural Networks

Hidden units like  $a^{[2] < i>}$  have inputs  $a^{[1] < i>}$  and  $a^{[2] < i-1>}$ .

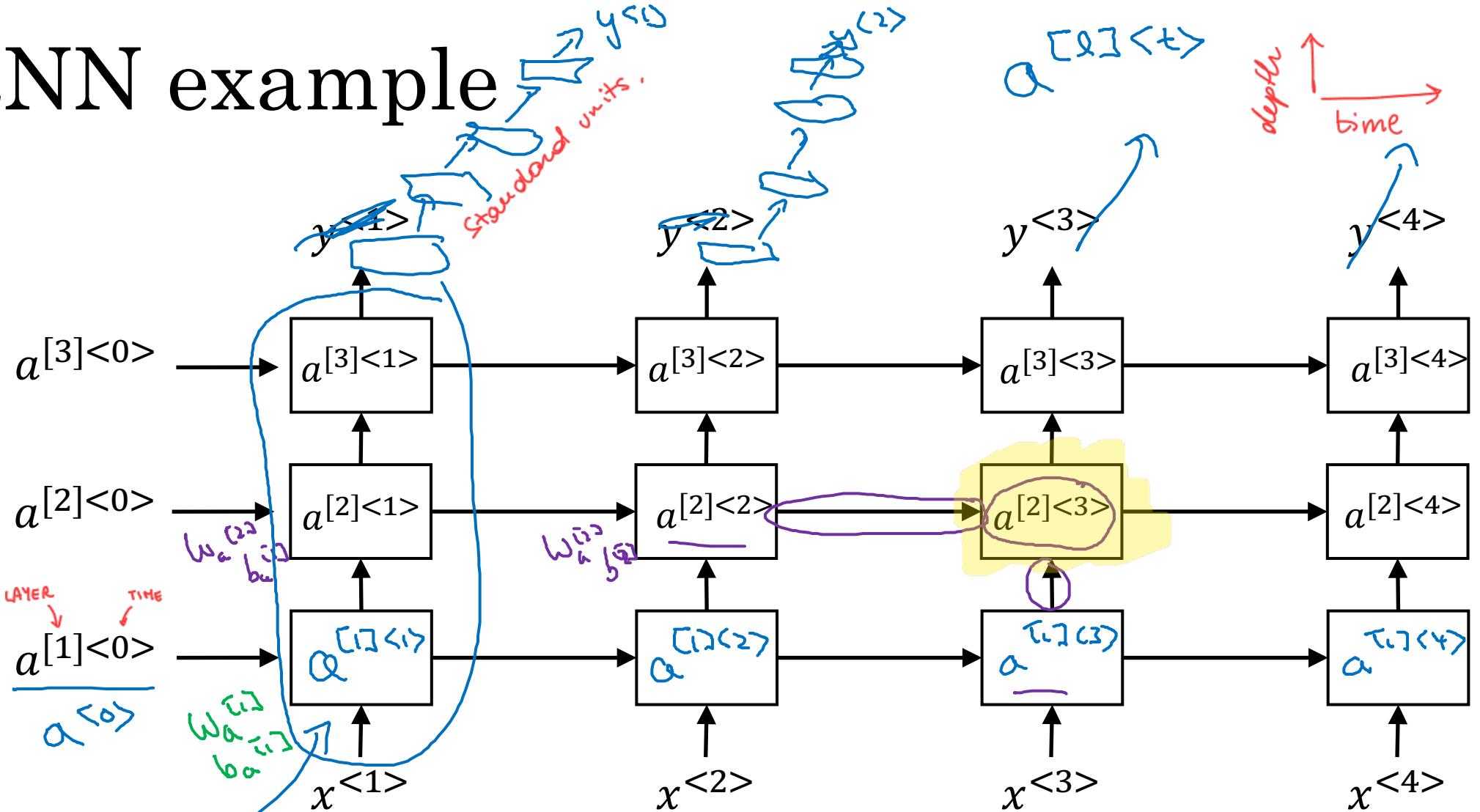
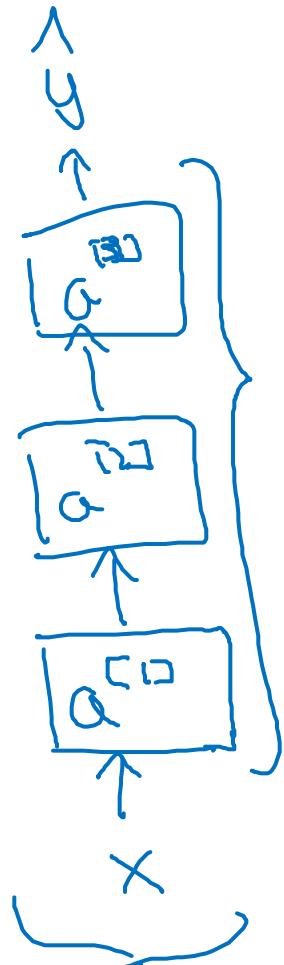
---

## Deep RNNs

deeplearning.ai

It's unusual to find very deep RNNs, usually 3 layers at most. One possibility is to stack few recurrent layers and several standard units.

# Deep RNN example



$RNN$        $a^{[2]<3>} = g(W_a^{[2]} [a^{[1]<2>}, a^{[1]<3>}] + b_a^{[2]})$

$GRU$

$LSTM$

$BRNN$