# Recurrent Neural Networks

---
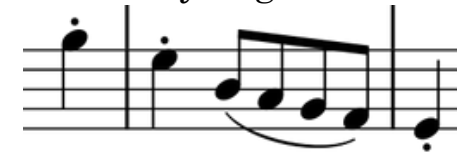
# Why sequence models?

deeplearning.ai

# Examples of sequence data

Speech recognition  → "The quick brown fox jumped over the lazy dog."
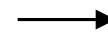
Music generation $\emptyset$ → 

Sentiment classification "There is nothing to like in this movie." → ★☆☆☆☆

DNA sequence analysis → AGCCCCTGTGAGGAACTAG → AGCCCCTGTGAGGAACTAG
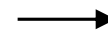
Machine translation Voulez-vous chanter avec moi? → Do you want to sing with me?

Video activity recognition  → Running

Name entity recognition → Yesterday, Harry Potter met Hermione Granger. → Yesterday, Harry Potter met Hermione Granger.

Andrew Ng

Recurrent Neural Networks

---

Notation

deeplearning.ai

# Motivating example

NLP

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

$\rightarrow$ $x^{<1>}$   $x^{<2>}$   $x^{<3>}$   $- - \cdots$   $x^{<t>}$   $- - \cdots$   $x^{<9>}$

$T_x = 9$

$\rightarrow$ y:      1      1      0      1      1      0   0   0   0

$y^{<1>}$   $y^{<2>}$   $y^{<3>}$   $- \cdots$   $y^{<9>}$

$T_y = 9$

$x^{(i)<t>}$        $T_x^{(i)} = 9$      15

$y^{(i)<t>}$  ↑     $T_y^{(i)}$

# Representing words

$x^{<t>}$  $(x, y)$

$x \longrightarrow y$

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

$x^{<1>}$   $x^{<2>}$   $x^{<3>}$   ...   $x^{<t>}$   $x^{<9>}$

Vocabulary

$\begin{bmatrix} a \\ aaron \\ \vdots \\ and \\ \vdots \\ harry \\ potter \\ \vdots \\ zulu \end{bmatrix}$
1 ←
2
367 ←
4075
6830
10,000

<UNK> 10,000

←4075

←6830

←367

10,000

One-hot

# Representing words

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

$x^{<1>}$    $x^{<2>}$    $x^{<3>}$            ...          $x^{<9>}$

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

deeplearning.ai

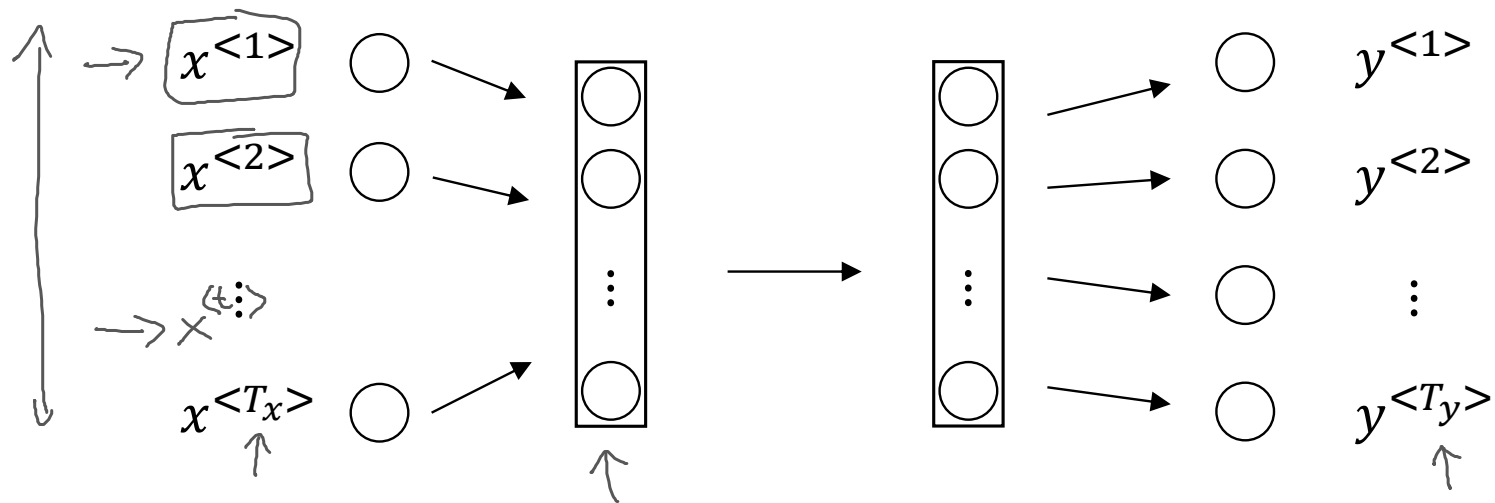Recurrent Neural Networks
_____

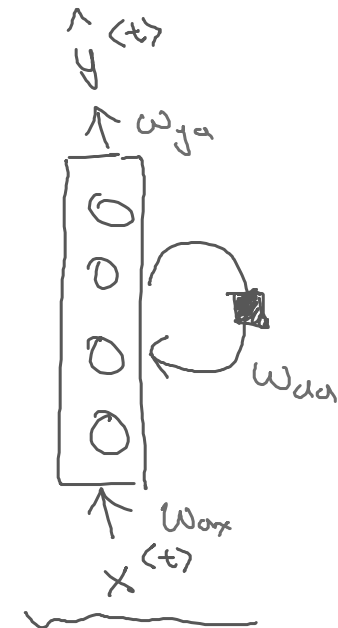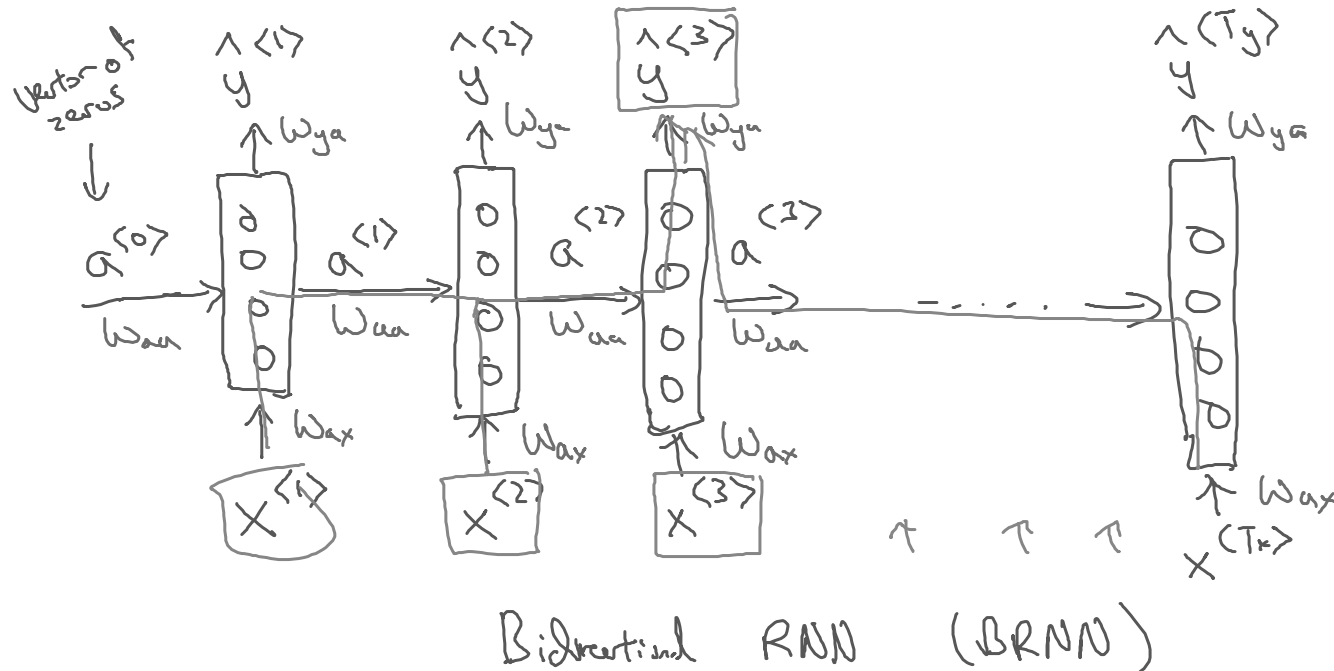Recurrent Neural Network Model

# 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

$T_x = T_y$



Bidirectional RNN (BRNN)

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

He said, "Teddy bears are on sale!"

# Forward Propagation

$a \leftarrow W_{ax} x^{(1)}$

$\hat{y}^{<1>}$     $\hat{y}^{<2>}$     $\hat{y}^{<3>}$     $\hat{y}^{<T_y>}$

$a^{<0>} \rightarrow$   $a^{<1>}$   $a^{<2>}$   ...   $a^{<T_x-1>} \rightarrow$

$x^{<1>}$     $x^{<2>}$     $x^{<3>}$     $x^{<T_x>}$

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

$a^{(1)} = g_1\left(W_{aa} a^{<0>} + W_{ax} x^{(1)} + b_a\right) \leftarrow$ tanh / ReLU

$\hat{y}^{(1)} = g_2\left(W_{ya} a^{(1)} + b_y\right) \leftarrow$ sigmoid

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

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

Andrew Ng

# Simplified RNN notation

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

100      10,000

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

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

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

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

$$100 \left[ W_{aa} \mid W_{ax} \right] = W_a$$

100    10 000

(100, 10100)

$$\left[a^{<t-1>}, x^{<t>}\right] = \left[\begin{array}{c} a^{<t-1>} \\ \hline x^{<t>} \end{array}\right] \begin{array}{c} 100 \\ 10000 \end{array} \right\} 10100$$

$$\left[W_{aa} \mid W_{ax}\right]\left[\begin{array}{c} a^{<t-1>} \\ x^{<t>} \end{array}\right] = W_{aa} a^{<t-1>} + W_{ax} x^{<t>}$$

Andrew Ng

deeplearning.ai

# Recurrent Neural Networks

---

# Backpropagation through time

# Forward propagation and backpropagation



Andrew Ng

# Forward propagation and backpropagation

$W_y, b_y$

$W_a, b_a$

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

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

Backpropagation through time

deeplearning.ai

# Recurrent Neural Networks

---

# Different types of RNNs

# Examples of sequence data

$T_x$     $T_y$

$y$

| | | |
|---|---|---|
| Speech recognition | $x$ [waveform] $\longrightarrow$ | "The quick brown fox jumped over the lazy dog." |
| Music generation | $\emptyset$ $\longrightarrow$ | [musical notes] |
| Sentiment classification | "There is nothing to like in this movie." $\longrightarrow$ | ★☆☆☆☆ |
| DNA sequence analysis | AGCCCCTGTGAGGAACTAG $\longrightarrow$ | AGCCCCTGTGAGGAACTAG |
| Machine translation | Voulez-vous chanter avec moi? $\longrightarrow$ | Do you want to sing with me? |
| Video activity recognition | [images] $\longrightarrow$ | Running |
| Name entity recognition | Yesterday, Harry Potter met Hermione Granger. $\longrightarrow$ | Yesterday, Harry Potter met Hermione Granger. |

Andrew Ng

# Examples of RNN architectures

$T_x = T_y$

$\hat{y}^{(1)}$  $\hat{y}^{(2)}$  $\hat{y}^{(T_y)}$

$a^{(0)}$

$x^{(1)}$  $x^{(2)}$  $\cdots$  $x^{(T_x)}$

Many-to-many

Sentiment classification
$x = \text{text}$
$y = 0/1$    $1 \cdots 5$

$y$

$x^{(1)}$  $x^{(2)}$  $x^{(T_x)}$
There  is  $\cdots\cdots$  movie

Many-to-one

$y$

$x$

One-to-one

Andrew Ng

# Examples of RNN architectures

Music generation

$x \rightarrow y^{<1>} y^{<2>} \ldots y^{<T_y>}$

$\hat{y}^{<1>}$  $\hat{y}^{<2>}$  $\hat{y}^{<3>}$  $\hat{y}^{<T_y>}$

$a^{<0>} \rightarrow$

$x$

One-to-many

$x = \phi$

Machine translation

encoder

$a^{<0>} \rightarrow$

$x^{<1>}$     $x^{<T_x>}$

decoder

$\hat{y}^{<1>}$     $\hat{y}^{<T_y>}$

Many-to-many

Andrew Ng

# Summary of RNN types



One to one

One to many

Many to one

Many to many    $T_x = T_y$

Many to many

Andrew Ng

# Recurrent Neural Networks

## Language model and sequence generation

# What is language modelling?

Speech recognition

The apple and <u>pair</u> salad.

$\longrightarrow$ The apple and <u>pear</u> salad.

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

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

$P(\text{sentence}) = ?$ $\qquad P\left(y^{<1>}, y^{<2>}, \dots, y^{<T_y>}\right)$

# Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

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

$y^{<1>}$  $y^{<2>}$  $y^{<3>}$  $\cdots$  $y^{<8>}$  $y^{<9>}$

$x^{<t>} = y^{<t-1>}$

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

<UNK>

10,000

# RNN model

$P(\text{average} \mid \text{cats})$

$P(\underline{\quad} \mid \text{"cats average"})$

$P(\underline{\langle EOS \rangle} \mid \dots)$

$P(a) \; P(\text{aaron}) \dots P(\text{cats}) \dots P(\text{zulu})$
$P(\langle UNK \rangle)$
$P(\langle EOS \rangle)$

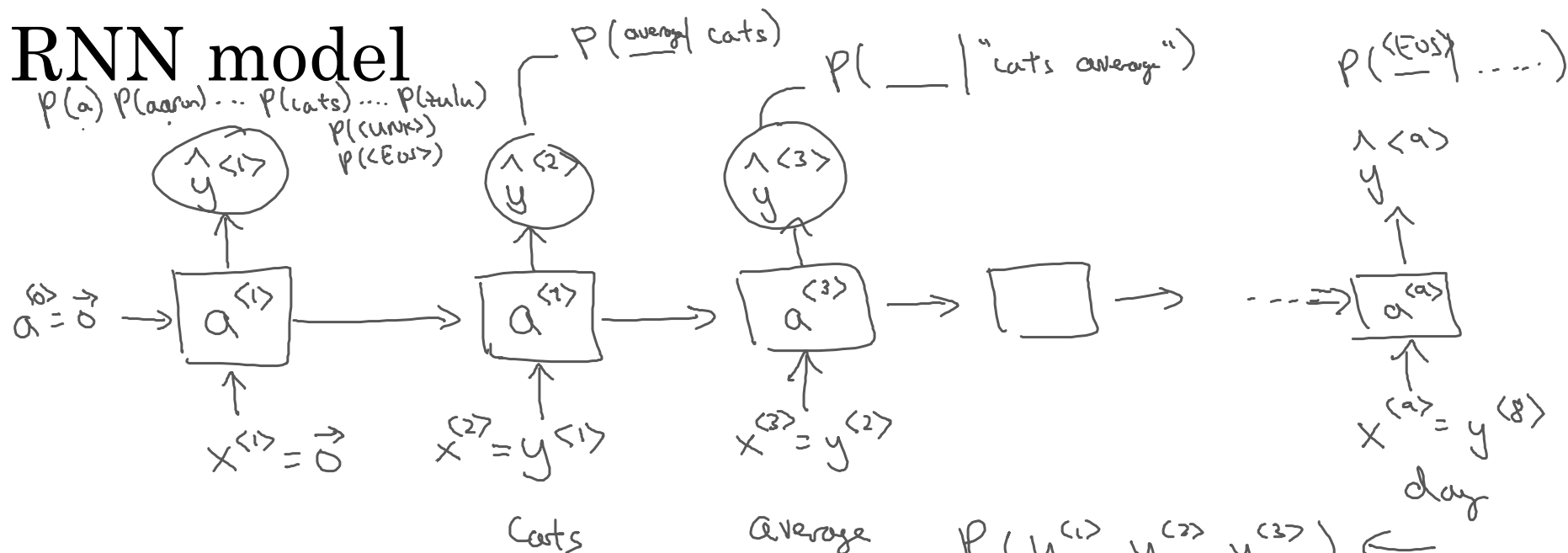$\hat{y}^{\langle 1 \rangle}$

$\hat{y}^{\langle 2 \rangle}$

$\hat{y}^{\langle 3 \rangle}$

$\hat{y}^{\langle 9 \rangle}$

$a^{\langle 0 \rangle} = \vec{0}$ → $a^{\langle 1 \rangle}$ → $a^{\langle 2 \rangle}$ → $a^{\langle 3 \rangle}$ → $\square$ → $\dashrightarrow$ $a^{\langle 9 \rangle}$

$x^{\langle 1 \rangle} = \vec{0}$

$x^{\langle 2 \rangle} = y^{\langle 1 \rangle}$

$x^{\langle 3 \rangle} = y^{\langle 2 \rangle}$

$x^{\langle 9 \rangle} = y^{\langle 8 \rangle}$
day

Cats

average

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

$P(y^{\langle 1 \rangle}, y^{\langle 2 \rangle}, y^{\langle 3 \rangle})$ ←

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

$= P(y^{\langle 1 \rangle}) \; P(y^{\langle 2 \rangle} \mid y^{\langle 1 \rangle})$

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

$P(y^{\langle 3 \rangle} \mid y^{\langle 1 \rangle}, y^{\langle 2 \rangle})$

Recurrent Neural Networks

deeplearning.ai

Sampling novel sequences

# Sampling a sequence from a trained RNN



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

Training:

$a^{<0>} \rightarrow$ boxes with $a^{<1>}, a^{<2>}, a^{<3>}, \ldots, a^{<T_y>}$

Outputs: $\hat{y}^{<1>}, \hat{y}^{<2>}, \hat{y}^{<3>}, \hat{y}^{<T_y>}$

Inputs: $x^{<1>}, y^{<1>}, y^{<2>}, \ldots, y^{<T_x-1>}$

Sampling:

The $\hat{y}^{<1>}$, $\hat{y}^{<2>}$, $\hat{y}^{<3>}$, $\hat{y}^{<T_y>}$

$a^{<0>} = 0 \rightarrow$ $a^{<1>}, a^{<2>}, a^{<3>}$

$x^{<1>} = 0$, $x^{<2>} =$ The $= \hat{y}^{<1>}$

$y^{<T_x-1>}$

$\langle EOS \rangle$

$\langle UNK \rangle$

$\rightarrow P(a) P(aaron) \ldots P(zulu) P(\langle UNK \rangle)$

np.random.choice

$P(\_ | the)$

Andrew Ng

# Character-level language model

$\rightarrow$ Vocabulary = [a, aaron, …, zulu, <UNK>]  $\leftarrow$

$\rightarrow$ Vocabulary = [ a, b, c, …, z, ␣, ., , , ;, 0, …, 9, A, …, Z]

$y^{<1>} y^{<2>} y^{<3>} y^{<4>}$

Cat ↑↑↑↑ average …

Mau



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 besser be my love to me see sabl's.
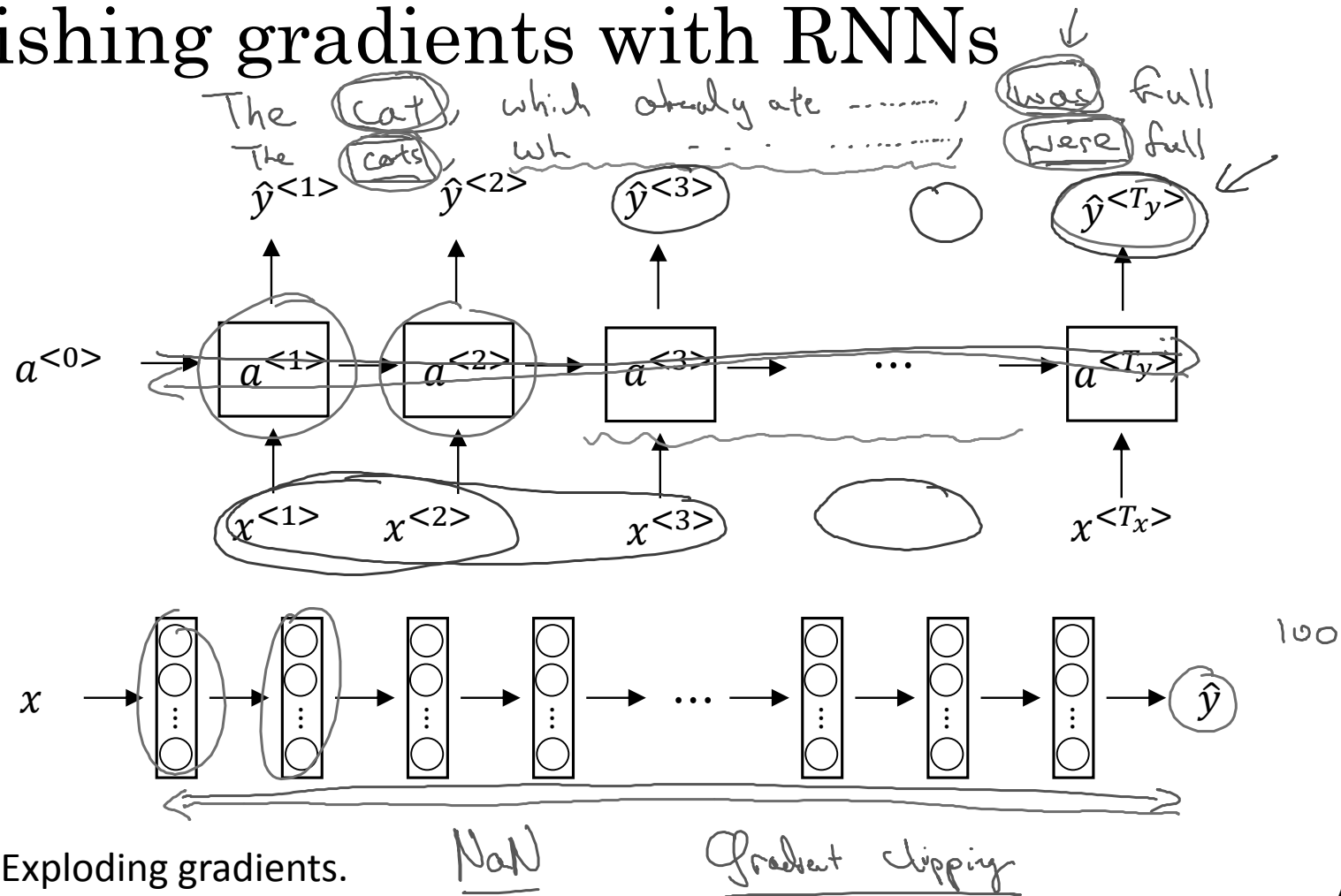
For whose are ruse of mine eyes heaves.

deeplearning.ai

Recurrent Neural Networks

---

Vanishing gradients with RNNs
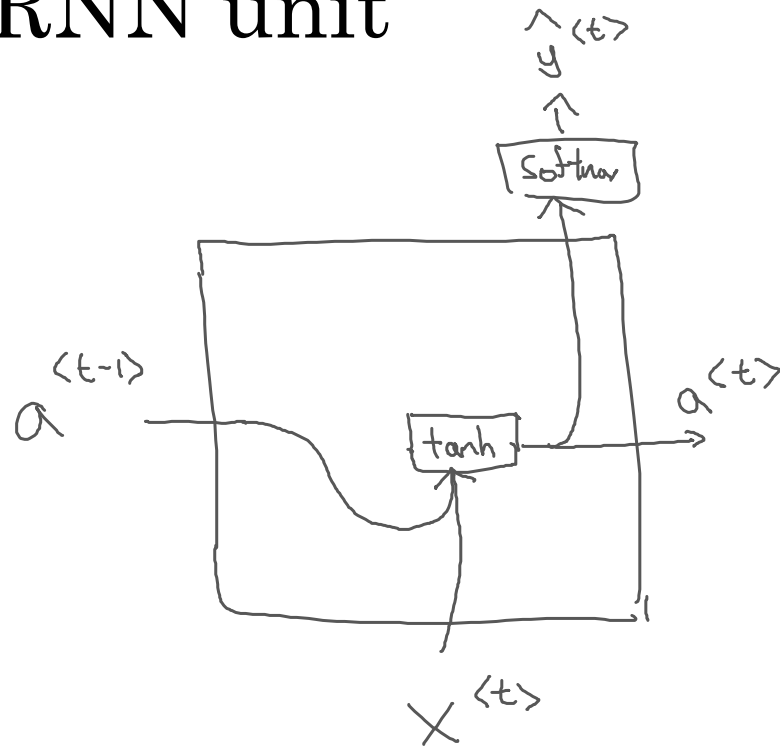
# Vanishing gradients with RNNs

The (Cat), which already ate --------- , (was) full
The (cats) wh - - - -  ......... (were) full



Exploding gradients.   NaN    Gradient clipping

Andrew Ng

deeplearning.ai

# Recurrent Neural Networks

Gated Recurrent Unit (GRU)

# RNN unit

$\hat{y}^{\langle t \rangle}$

Softmax

$a^{\langle t-1 \rangle}$

$a$

tanh

$a^{\langle t \rangle}$

$x^{\langle t \rangle}$

tanh

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

# GRU (simplified)



$c^{\langle t-1 \rangle} = a^{\langle t-1 \rangle}$

softmax $\rightarrow y^{\langle t \rangle}$

$\tilde{c}^{\langle t \rangle}$

tanh    $\sigma$    $\Gamma_u$

$x^{\langle t \rangle}$

$\rightarrow c^{\langle t \rangle} = a^{\langle t \rangle}$

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

$c^{\langle t \rangle} = 1$ ............................

$\rightarrow$ The cat, which already ate ..., was full.

$c$ = memory cell

$\rightarrow \boxed{c^{\langle t \rangle}} = a^{\langle t \rangle}$

$$\tilde{c}^{\langle t \rangle} = \tanh\left(W_c\left[c^{\langle t-1 \rangle}, x^{\langle t \rangle}\right] + b_c\right)$$

$$\Gamma_u = \sigma\left(W_u\left[c^{\langle t-1 \rangle}, x^{\langle t \rangle}\right] + b_u\right)$$

"update"

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

$\Gamma_u = 1$

element-wise

Gate

$\Gamma_u = 0.0000001$

[Cho et al., 2014. On the properties of neural machine translation: Encoder-decoder approaches] $\leftarrow$
[Chung et al., 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling] $\leftarrow$

Andrew Ng

# Full GRU

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

$u$ $\begin{cases} \Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_c) \end{cases}$

$r$

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

LSTM

The cat, which ate already, was full.

Andrew Ng

Recurrent Neural Networks

---

LSTM (long short term memory) unit

deeplearning.ai

# GRU and LSTM

## GRU

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

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

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

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

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

$$\Gamma_f$$

## LSTM

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

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

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

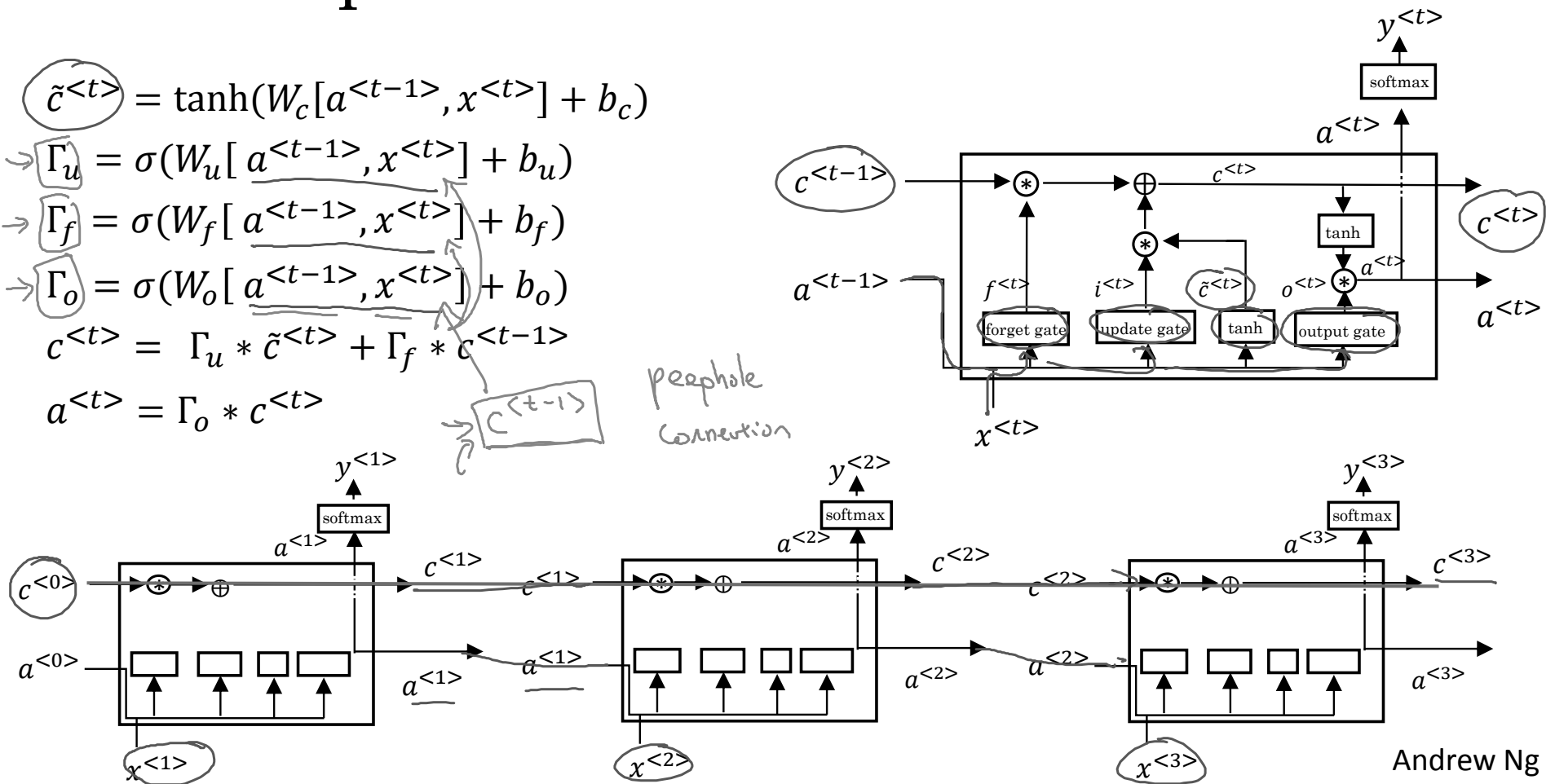$$(output) \quad \Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

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

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

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

Andrew Ng

# LSTM in pictures

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

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

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

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

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

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

peephole connection

$c^{<t-1>}$



Andrew Ng

Recurrent Neural Networks

---

Bidirectional RNN

deeplearning.ai

# Getting information from the future

He said, "Teddy bears are on sale!"

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



RNN
GRU
LSTM

# Bidirectional RNN (BRNN)

$$\hat{y}^{<t>} = g(W_y [\overrightarrow{a}^{<t>}, \overleftarrow{a}^{<t>}] + b_y)$$



GRU
LSTM

Acyclic graph

He said "Teddy Roosevelt ..."

BRNN w/LSTM

Andrew Ng