

Chapter 08: Recurent Neural Network

Recurrent Neural Networks

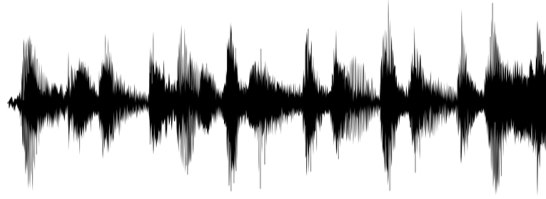


deeplearning.ai

- Coursera course: *Sequence models*
- **Week 1: Recurrent Neural Networks**
- Outline:
 - Why sequence models?
 - Notation
 - Recurrent Neural Network Model
 - Backpropagation through time
 - Different types of RNNs
 - Language model and sequence generation
 - Sampling novel sequences
 - Vanishing gradients with RNNs
 - GRU, LSTM, BiRNN, Deep RNN

Why sequence models?

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

AGCCCCTGTGAGGAACTAG



AG**CCCCTGTGAGGAACT**AG

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**.

Notation:

Motivating example

NLP

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

$\rightarrow \underline{x^{(1)}} \quad x^{(2)} \quad x^{(3)} \quad \dots \quad x^{(t)} \quad \dots \quad x^{(9)}$
 $T_x = 9$

$\rightarrow y:$ $\quad 1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0$
 $\quad y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(9)}$
 $T_y = 9$

$x^{(i)(t)}$

$T_x^{(i)} = 9$

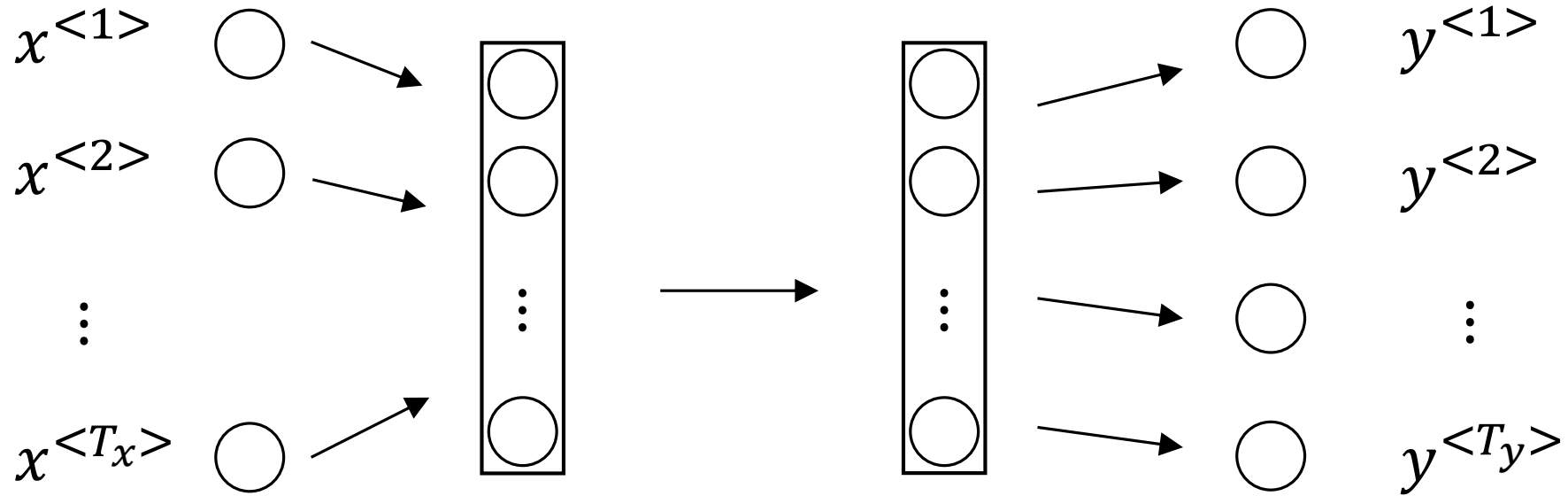
15

$y^{(i)(t)}$
 \uparrow

$T_y^{(i)}$

Recurrent Neural Network Model

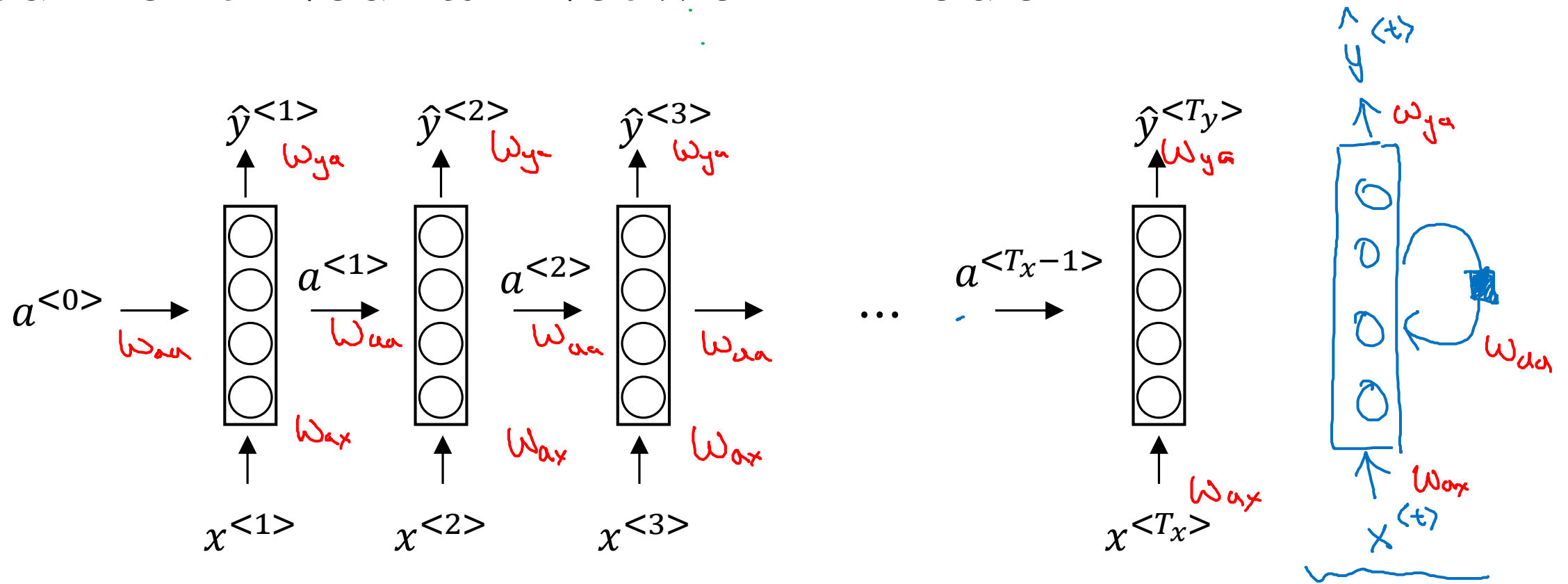
Why not a standard network?



Problems:

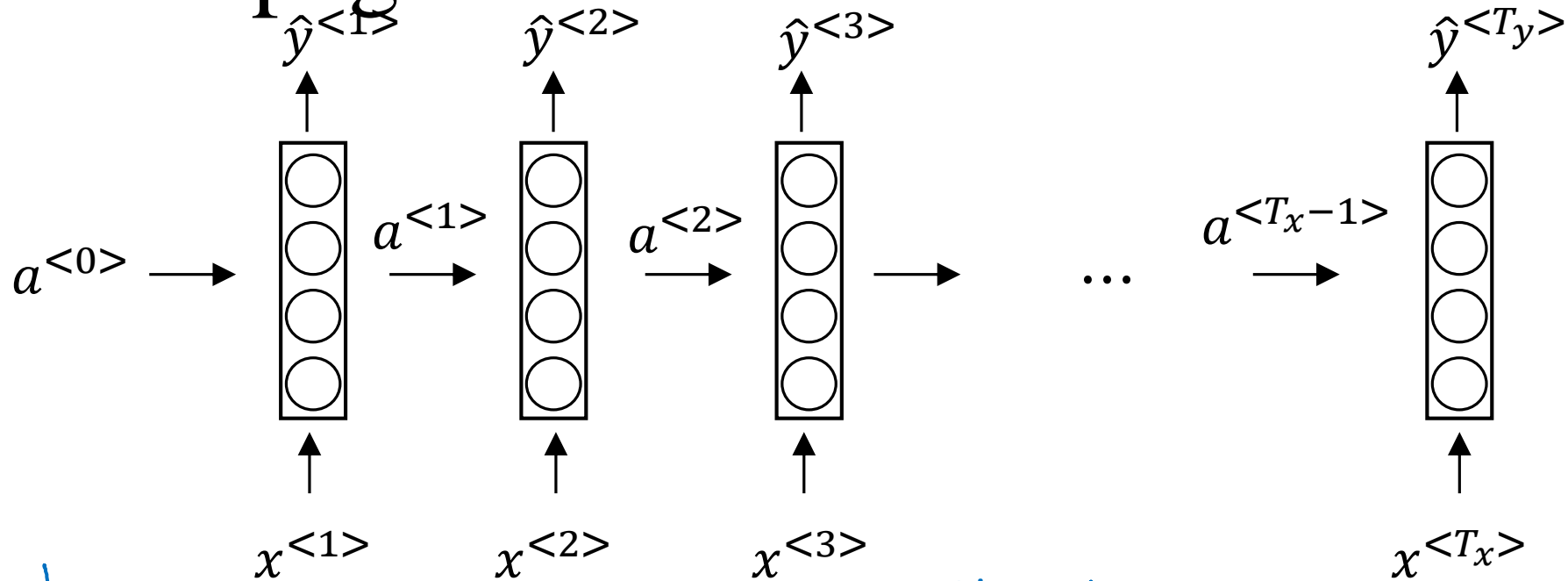
- Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.

Recurrent Neural Network Model



Recurrent Neural Network Model

Forward Propagation



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

$$a^{<1>} = g_1(W_{aa} a^{<0>} + \underline{W_{ax}} x^{<1>} + b_a) \leftarrow \underline{\tanh} / \text{Relu}$$

$$\hat{y}^{<1>} = g_2(W_{ya} a^{<1>} + b_y) \leftarrow \text{Sigmoid}$$

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

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

Recurrent Neural Network Model

Simplified RNN notation

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

Diagram annotations for the equation above:

- A green bracket above the equation indicates a recurrent connection from $a^{<t-1>}$ to $a^{<t>}$.
- A blue box around W_{aa} has an arrow pointing to it from the dimension $(100, 100)$ below.
- A blue box around W_{ax} has an arrow pointing to it from the dimension $(100, 10,000)$ below.
- The dimension 100 is written below $a^{<t-1>}$.
- The dimension $10,000$ is written below $x^{<t>}$.

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

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

Diagram annotations for the equation above:

- An arrow points from $\hat{y}^{<t>}$ to $y^{<t>}$.
- An arrow points up to $y^{<t>}$.
- An arrow points up to W_y .
- An arrow points up to b_y .

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

Diagram annotations for the equation above:

- A green circle around W_a has an arrow pointing to it from the dimension $(100, 10100)$ below.
- A purple box around $[a^{<t-1>}, x^{<t>}]$ has an arrow pointing to it from the dimension 10100 below.

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

Diagram annotations for the equation above:

- A green box around the matrix $\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix}$.
- A blue double-headed arrow below W_{aa} is labeled 100 .
- A blue double-headed arrow below W_{ax} is labeled 10000 .
- The dimension $(100, 10100)$ is written to the right of the matrix.

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

Diagram annotations for the equation above:

- A green circle around the vector $\begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$.
- A purple double-headed arrow to the right of $a^{<t-1>}$ is labeled 100 .
- A purple double-headed arrow to the right of $x^{<t>}$ is labeled 10000 .
- A purple double-headed arrow to the right of the entire vector is labeled 10100 .

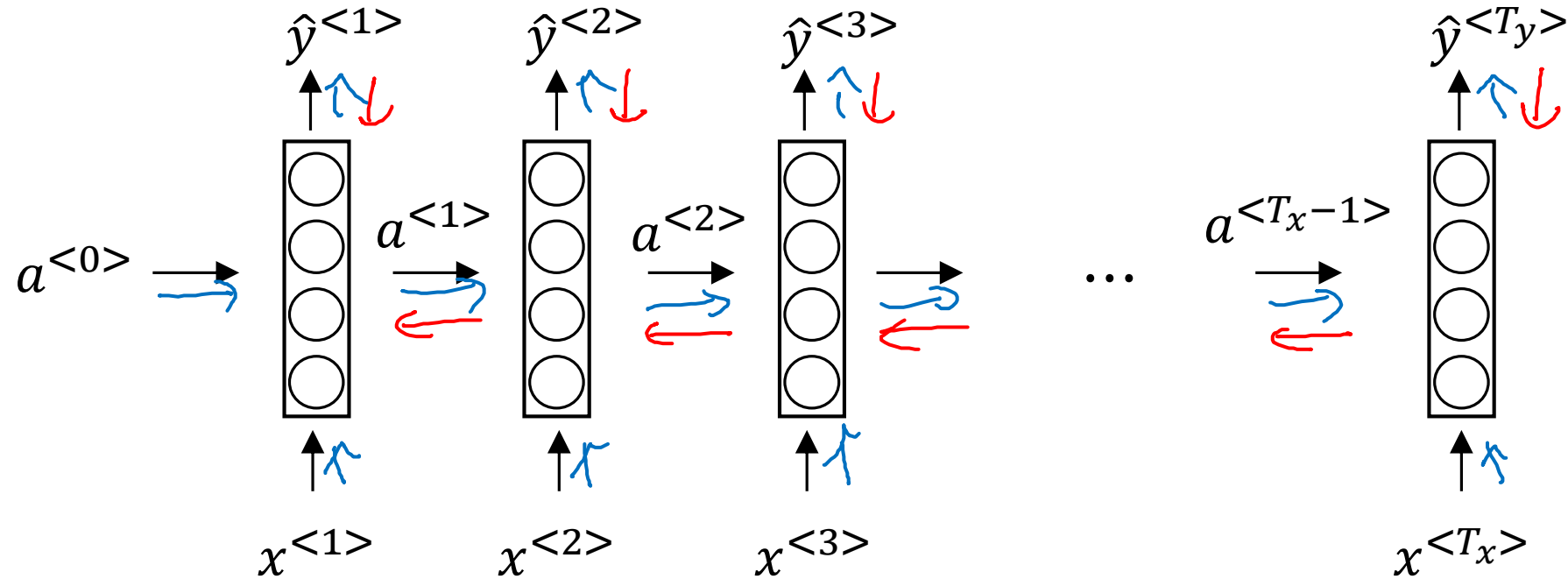
$$\begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$$

Diagram annotations for the equation above:

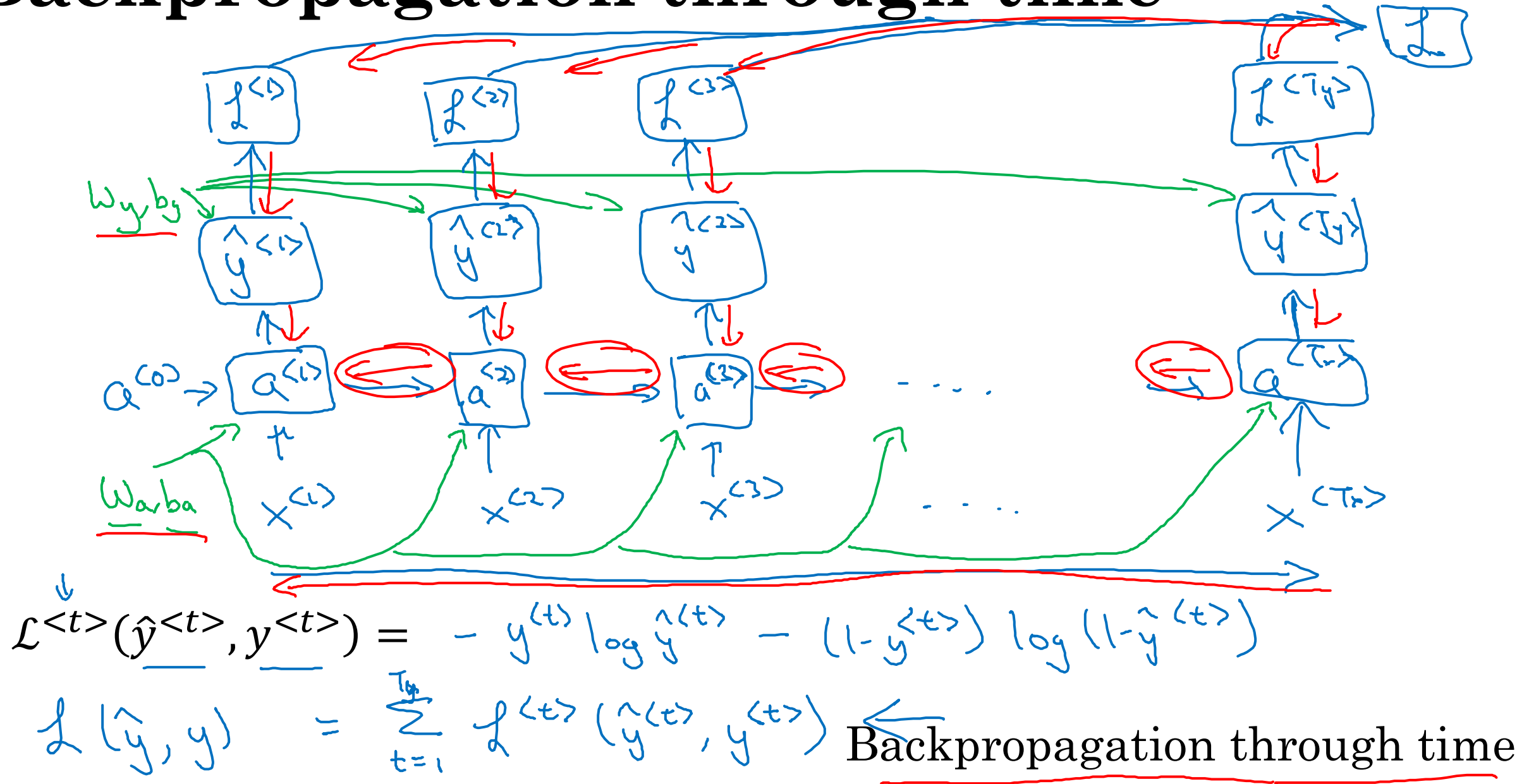
- A green box around the entire equation.

Backpropagation through time

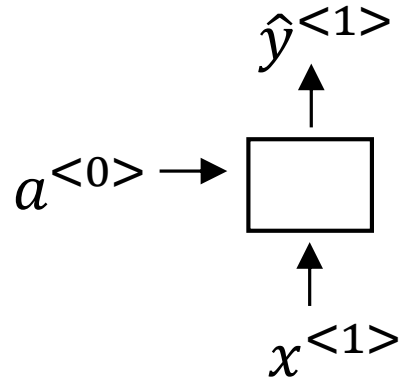
Forward propagation and backpropagation



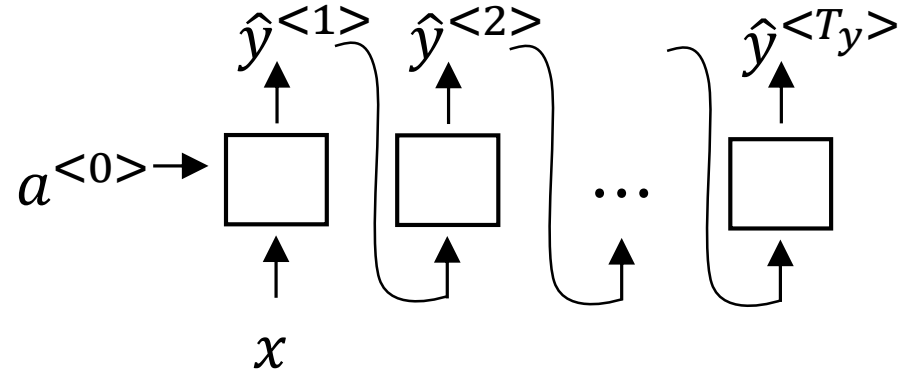
Backpropagation through time



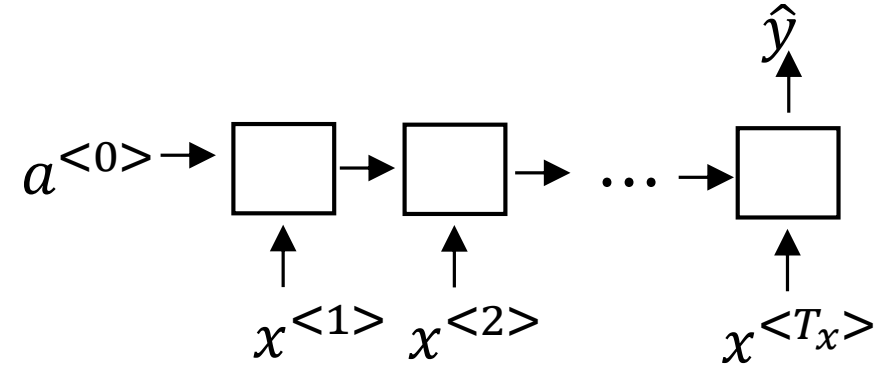
Different types of RNNs



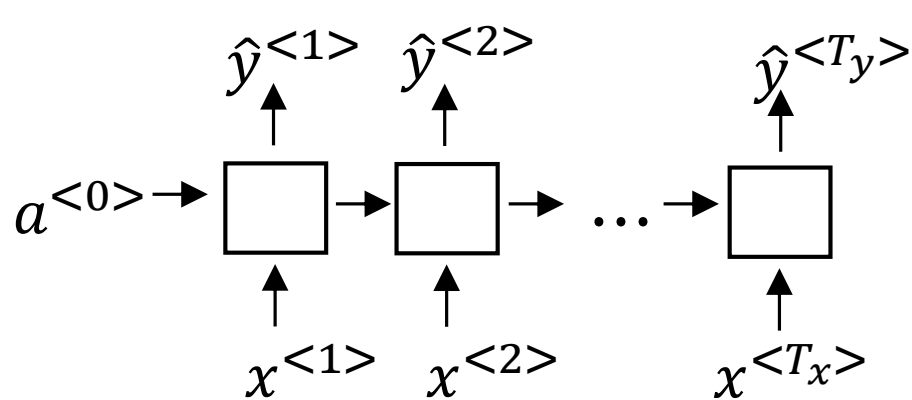
One to one



One to many

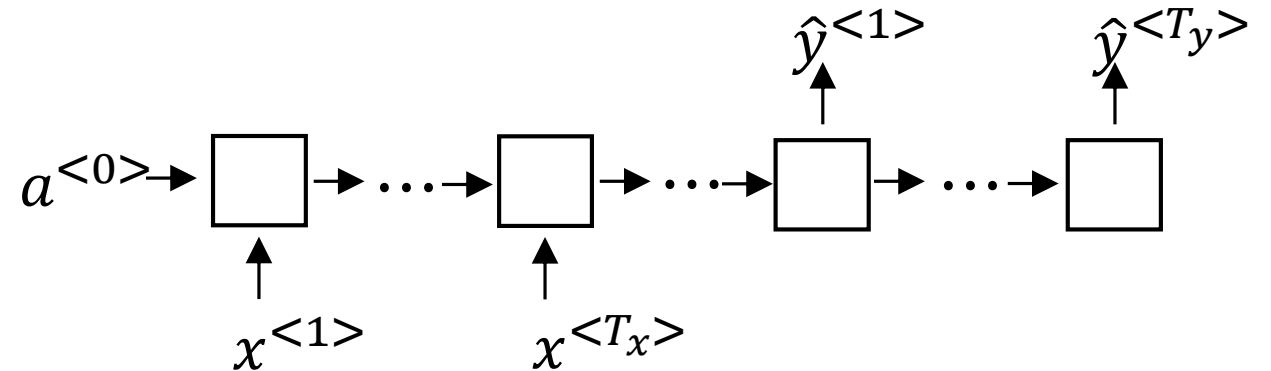


Many to one



Many to many

$T_x = T_y$



Many to many

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 model and sequence generation

Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

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

$y^{(1)}$ $y^{(2)}$ $y^{(3)}$... $y^{(8)}$ $y^{(9)}$
 $x^{(t)} = y^{(t-1)}$

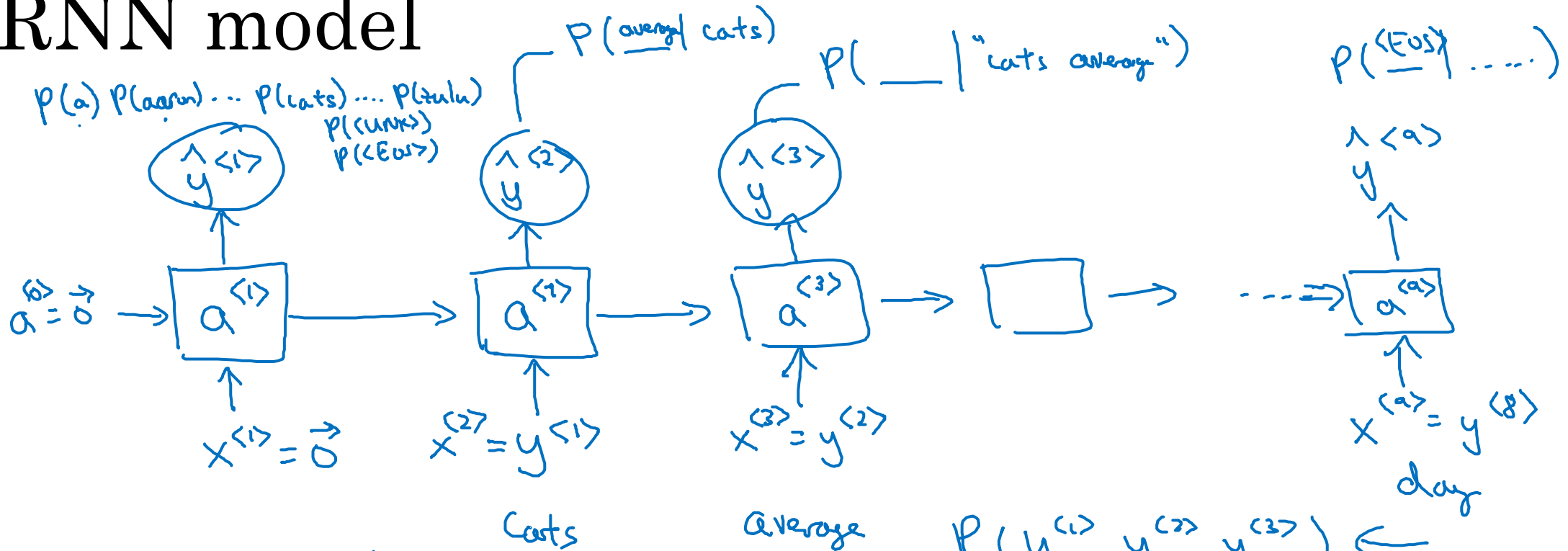
The Egyptian ~~Mau~~ is a breed of cat. $\langle \text{EOS} \rangle$

$\langle \text{UNK} \rangle$

10,000

Language model and sequence generation

RNN model



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

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

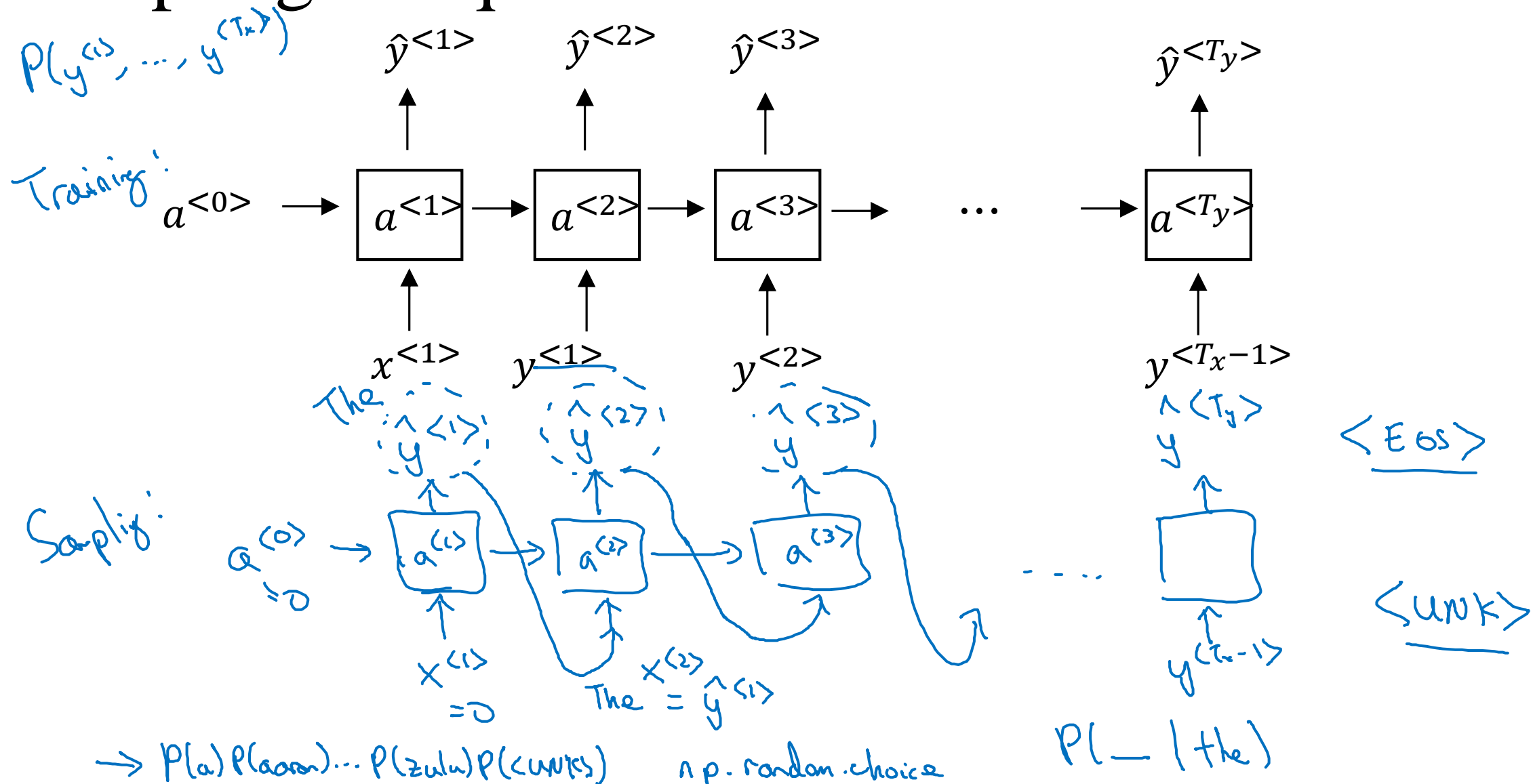
$$\mathcal{L} = \sum_t \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

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

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

Sampling novel sequences

Sampling a sequence from a trained RNN

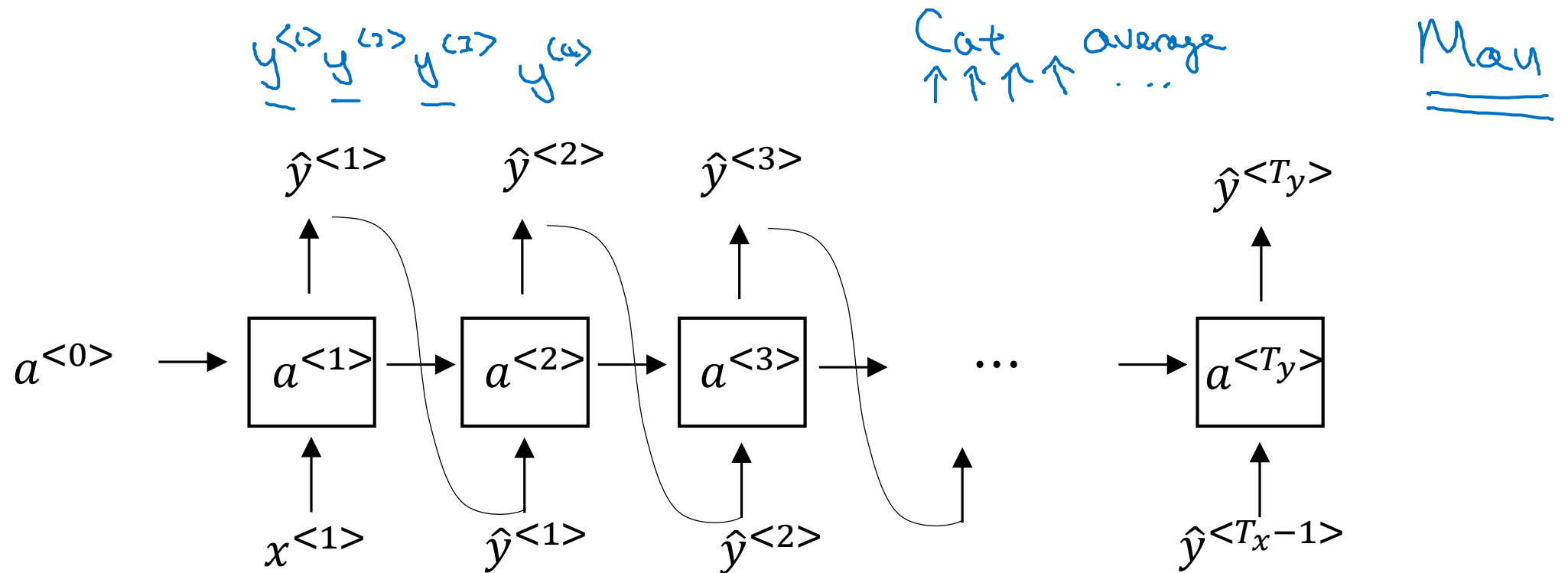


Sampling novel sequences

Character-level language model

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

→ Vocabulary = [a, b, c, ..., z, \backslash , ., , , ;, 0, ..., 9, A, ..., Z]



Sampling novel sequences

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

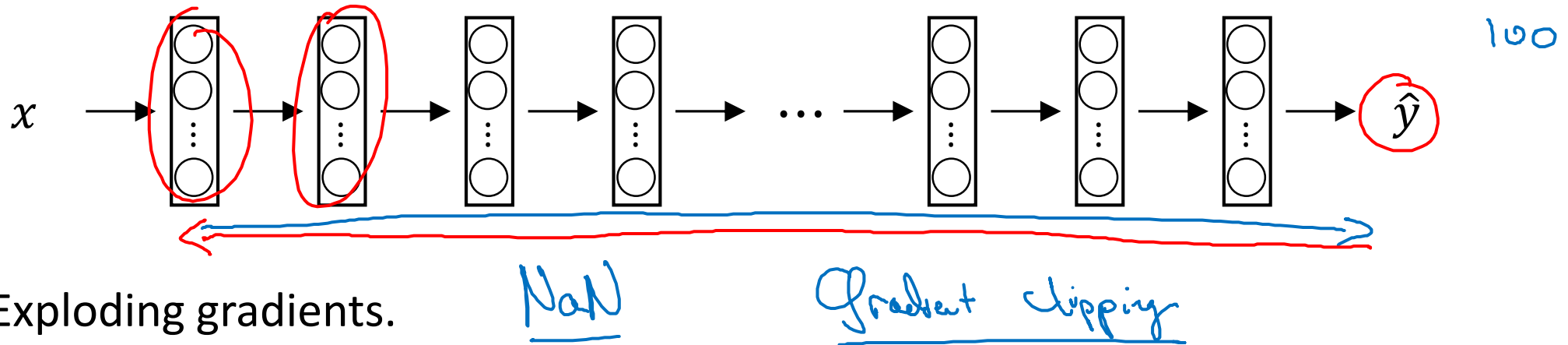
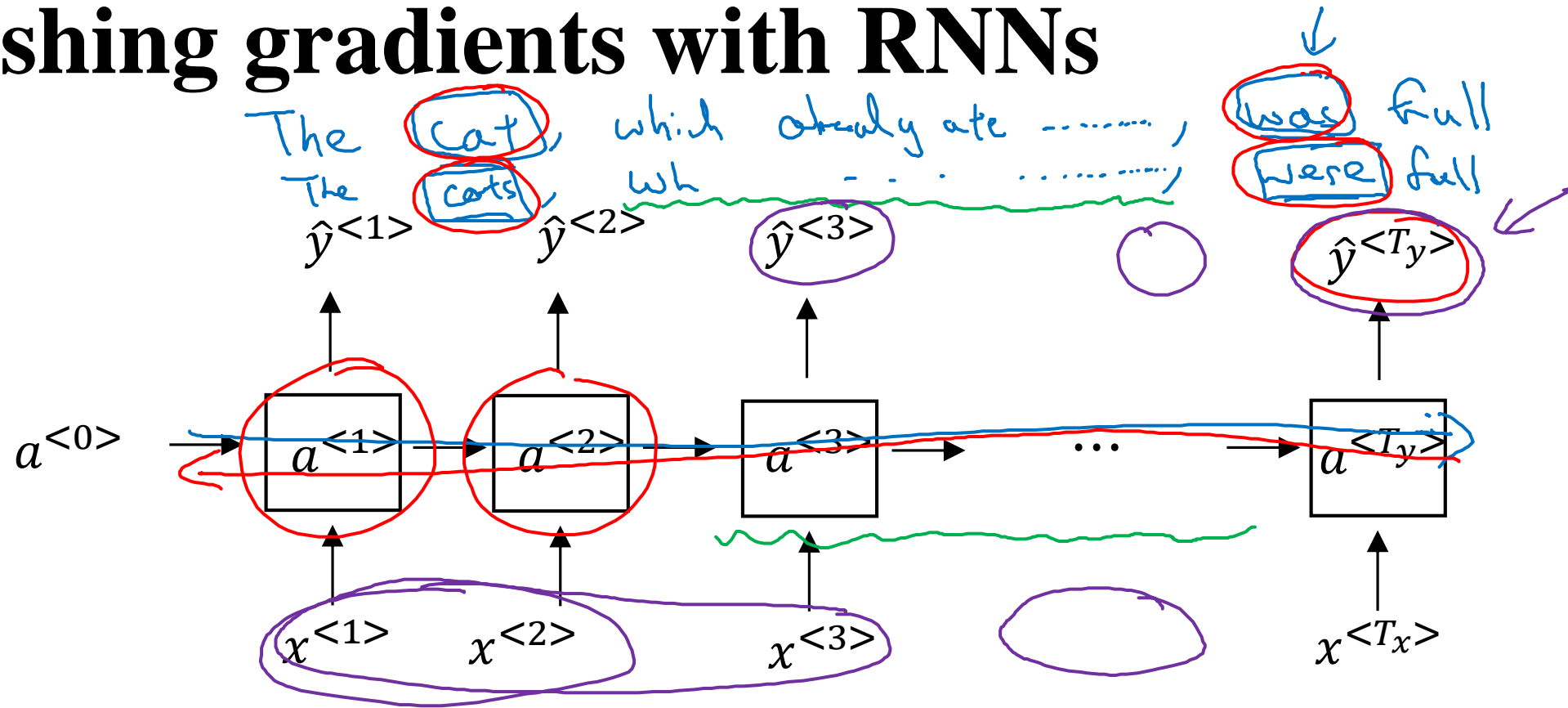
The mortal moon hath her eclipse in love.

And subject of this thou art another this fold.

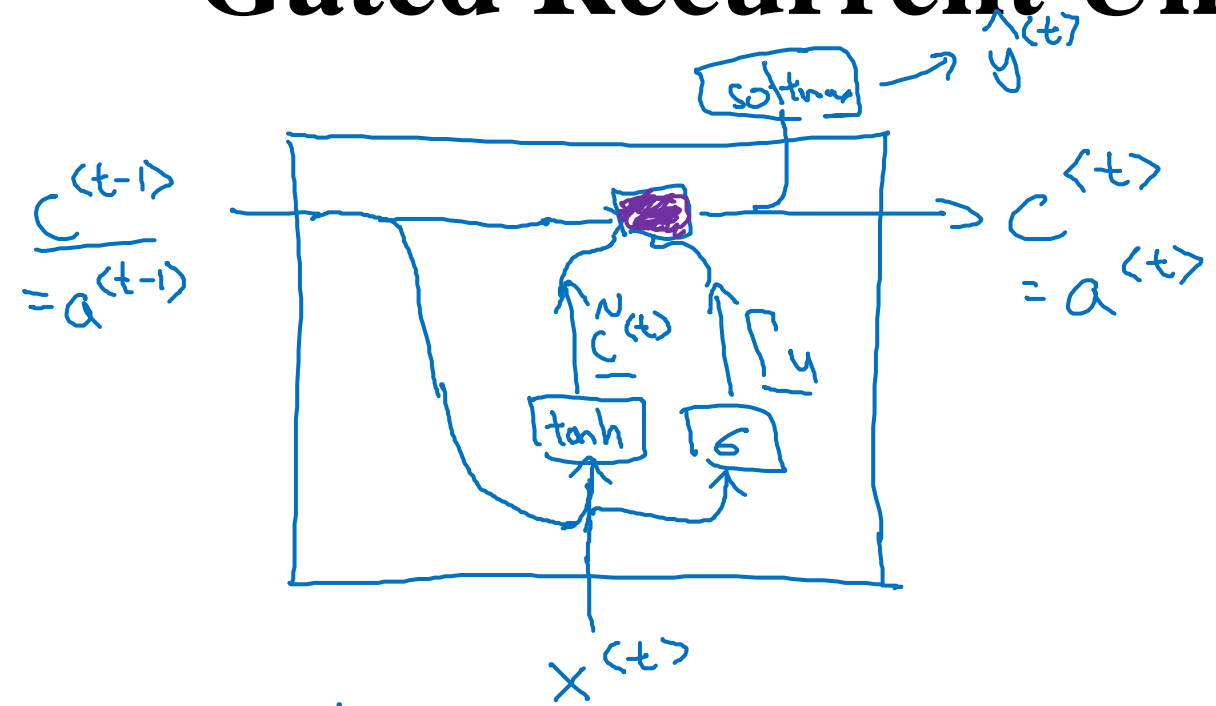
When besser be my love to me see sabl's.

For whose are ruse of mine eyes heaves.

Vanishing gradients with RNNs



Gated Recurrent Unit (GRU) (simplified)



$C = \text{memory cell}$

$$\rightarrow \boxed{C^{(t)}} = \underline{a}^{(t)}$$

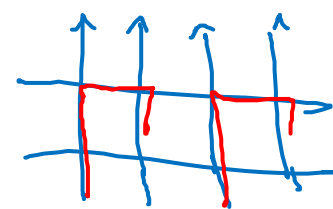
$$\rightarrow \boxed{\tilde{C}^{(t)}} = \tanh(W_c [c^{(t-1)}, x^{(t)}] + b_c)$$

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

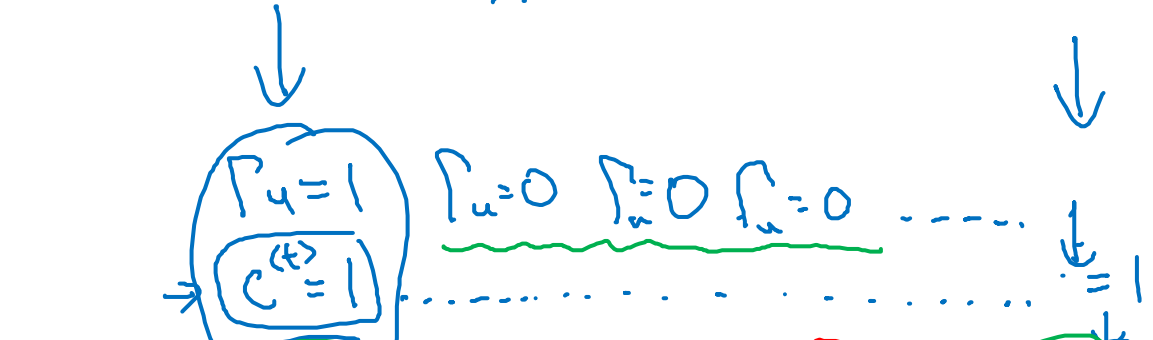
$$\boxed{C^{(t)}} = \underbrace{\Gamma_u}_{\text{update}} * \tilde{C}^{(t)} + (1 - \Gamma_u) * \boxed{C^{(t-1)}}$$

element-wise

$$\Gamma_u = 0.000001$$



Gate



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

Full GRU

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

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

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

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

LSTM

The cat, which ate already, was full.

GRU and LSTM

GRU

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

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

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

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

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

Γ_f

LSTM

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

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

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

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

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

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

LSTM in pictures

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

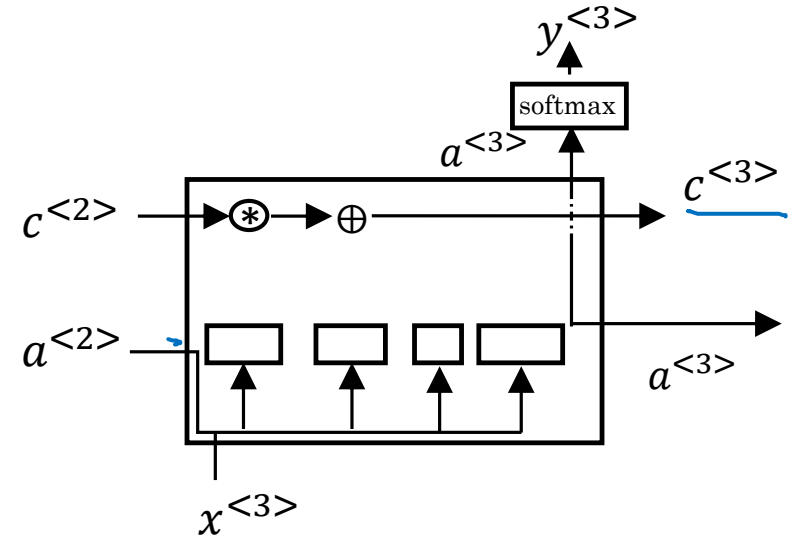
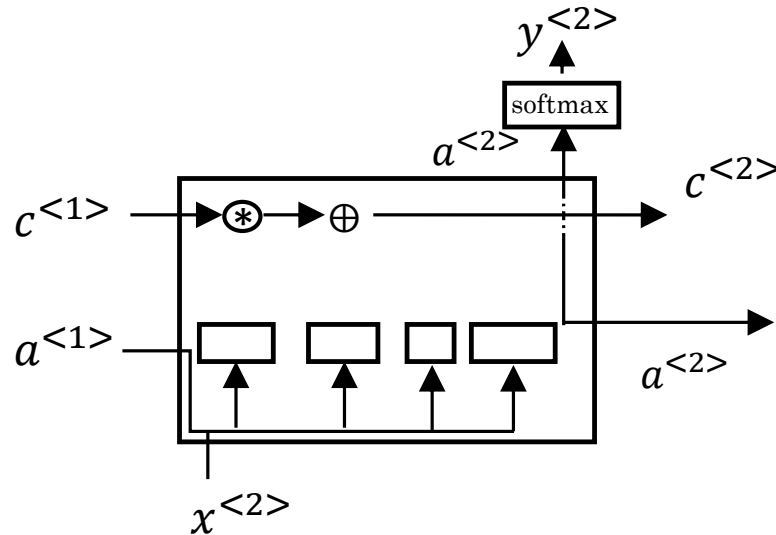
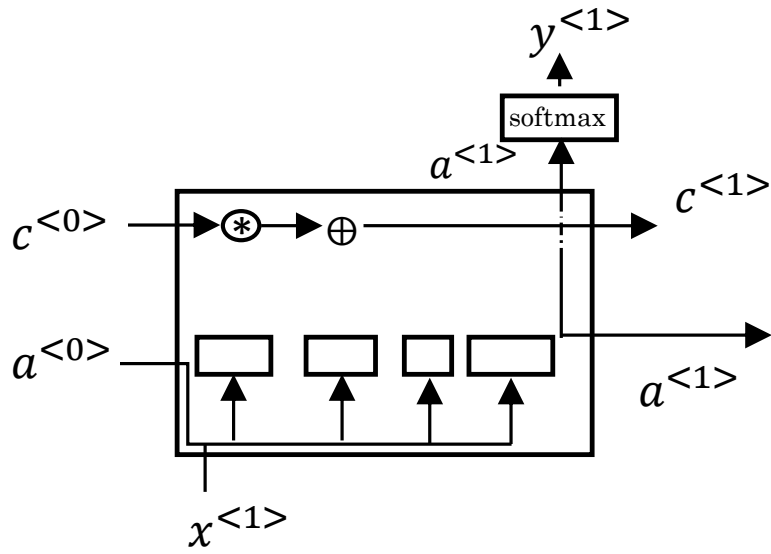
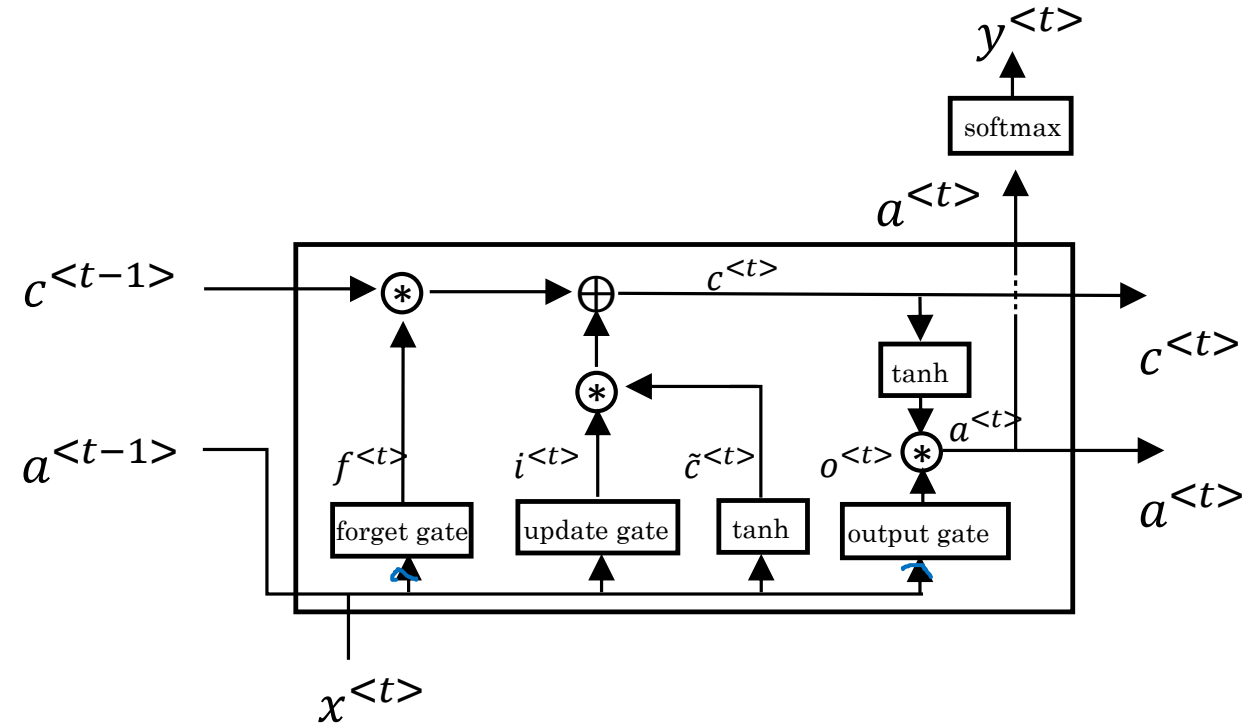
$$\rightarrow \Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

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

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

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

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

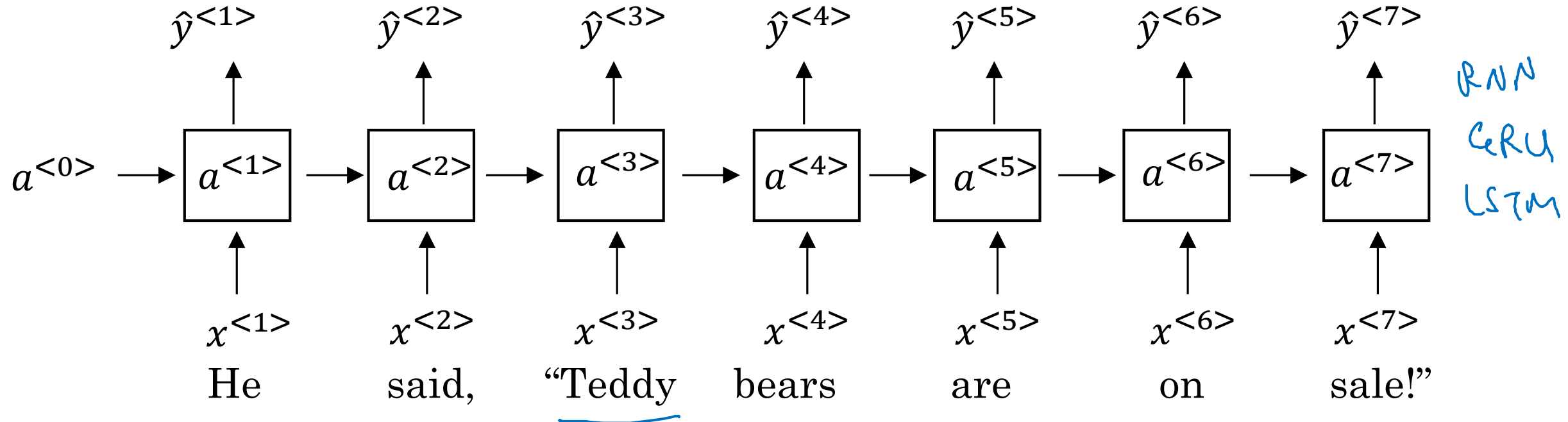


Bidirectional RNN

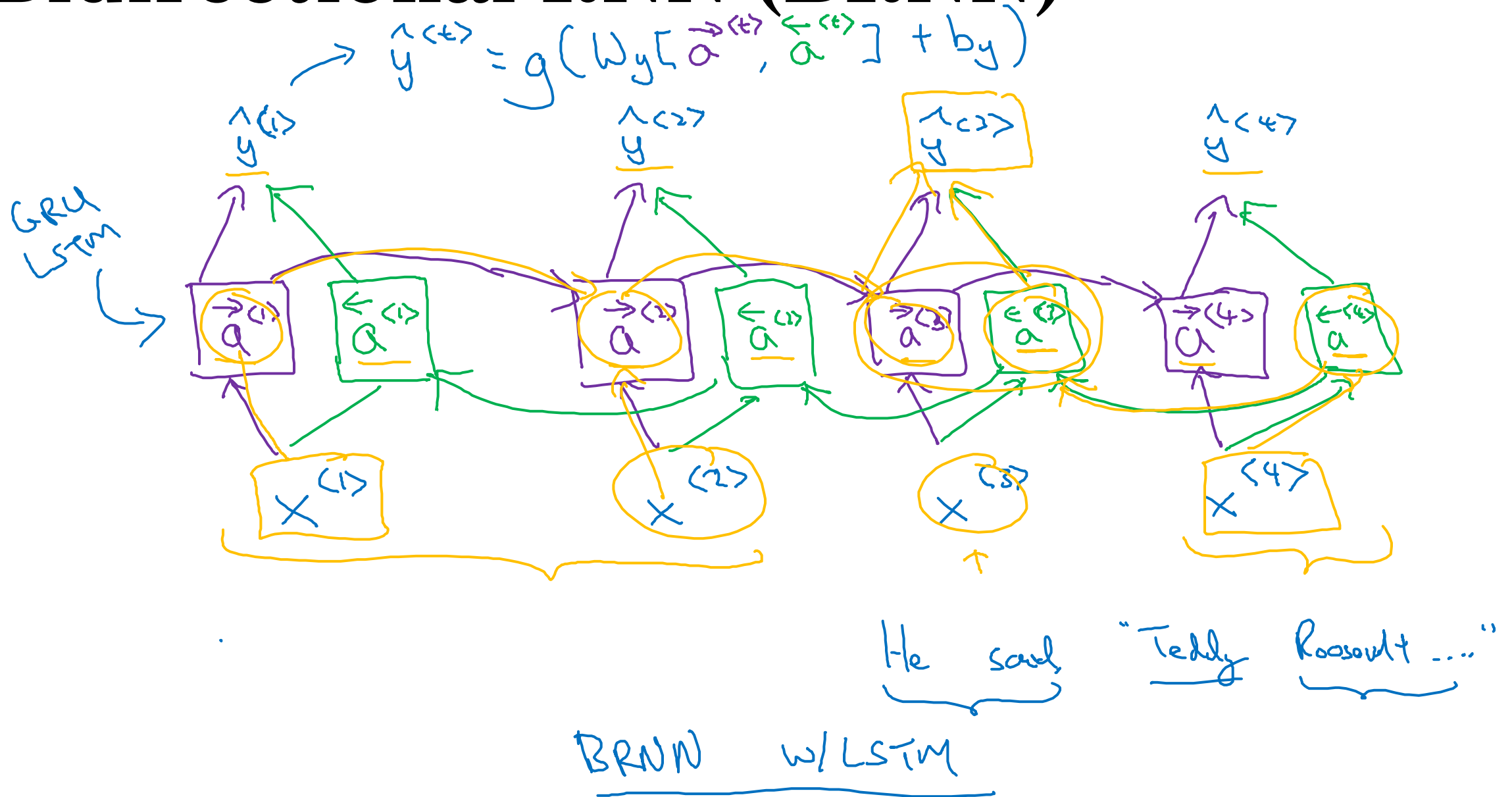
Getting information from the future

He said, “Teddy bears are on sale!”

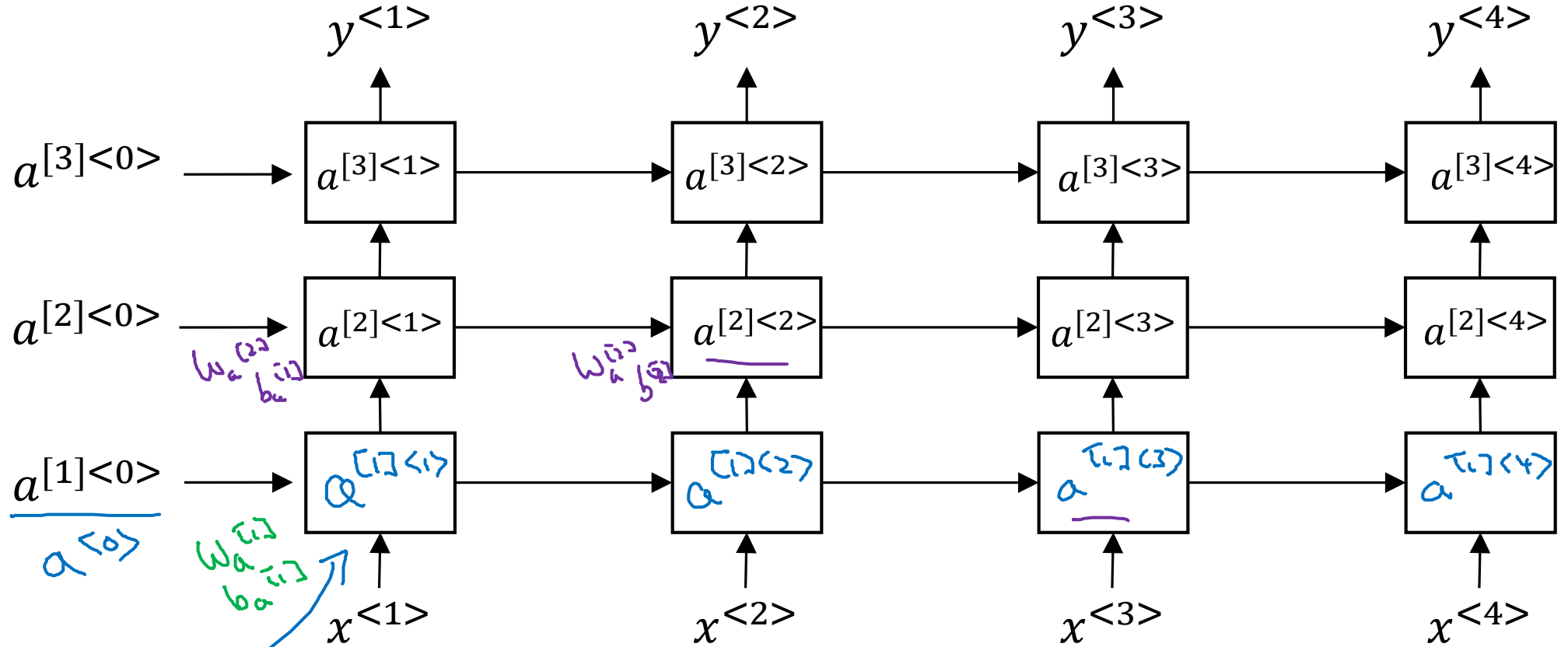
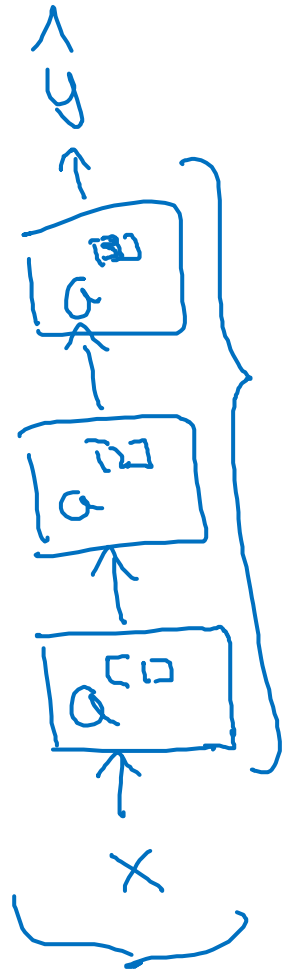
He said, “Teddy Roosevelt was a great President!”



Bidirectional RNN (BRNN)



Deep RNN example



RNN
GRU
LSTM

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