# Recurrent Neural Networks (RNN)

# Model Selection and Regularization

Speech recognition

 → "The quick brown fox jum ped over the lazy dog."

Music generation

∅ → 

Sentiment classification

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

DNA sequence analysis

AGCCCCTGTGAGGAACTAG → AG<span style="color:red">CCCCTGTGAGGAACT</span>AG

Machine translation

Voulez-vous chanter avec moi? → Do you want to sing with me?
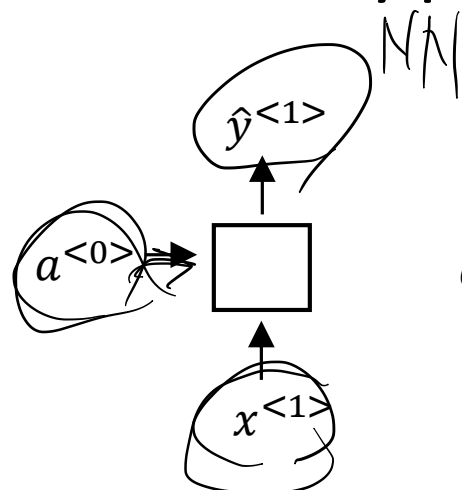
Video activity recognition

 → Running

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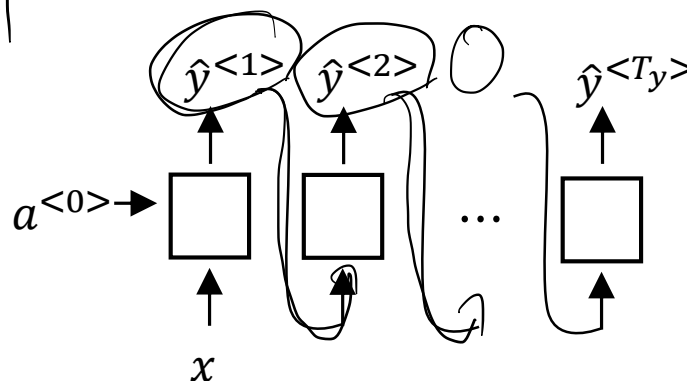
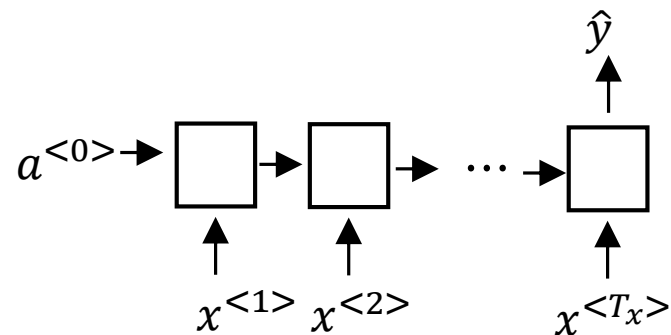Yesterday, Harry Potter met Hermione Granger. → Yesterday, Harry Potter met Hermione Granger.

# RNN types

NN

$\hat{y}^{<1>}$
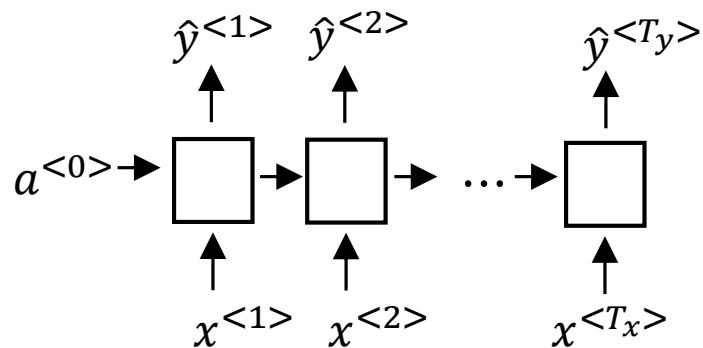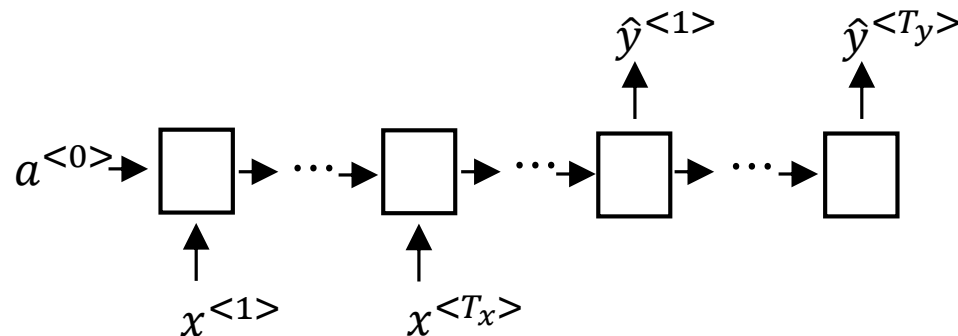
$a^{<0>}$

$x^{<1>}$

One to one

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

$a^{<0>}$

$\cdots$

$x$

One to many

$\hat{y}$

$a^{<0>}$   $\cdots$

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

Many to one

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

$a^{<0>}$   $\cdots$

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

Many to many

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

$a^{<0>}$   $\cdots$   $\cdots$   $\cdots$

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

Many to many

# Motivation

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

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

y: 1 1 0 1 1 0 0 0 0

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

Word 2 Vec

# Why not a standard neural net?

$x^{<1>}$ ◯ →

$x^{<2>}$ ◯ →

⋮

$x^{<T_x>}$ ◯ →

→

→ ◯ $y^{<1>}$

→ ◯ $y^{<2>}$

→ ◯ ⋮

→ ◯ $y^{<T_y>}$
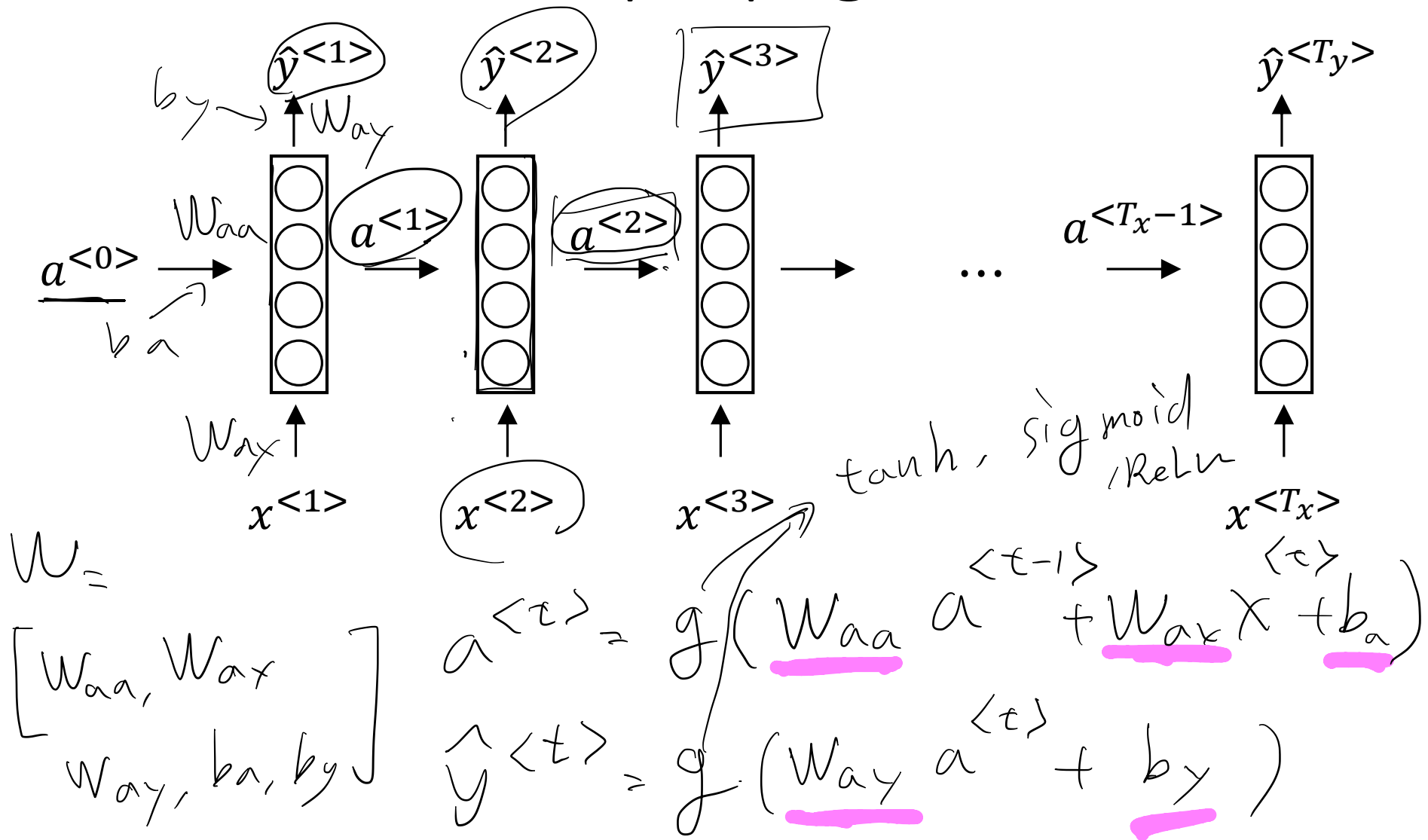
- Problems
  - Inputs, outputs can be different lengths in different examples
  - Inputs, outputs can be different lengths in different examples

# RNN- Forward propagation



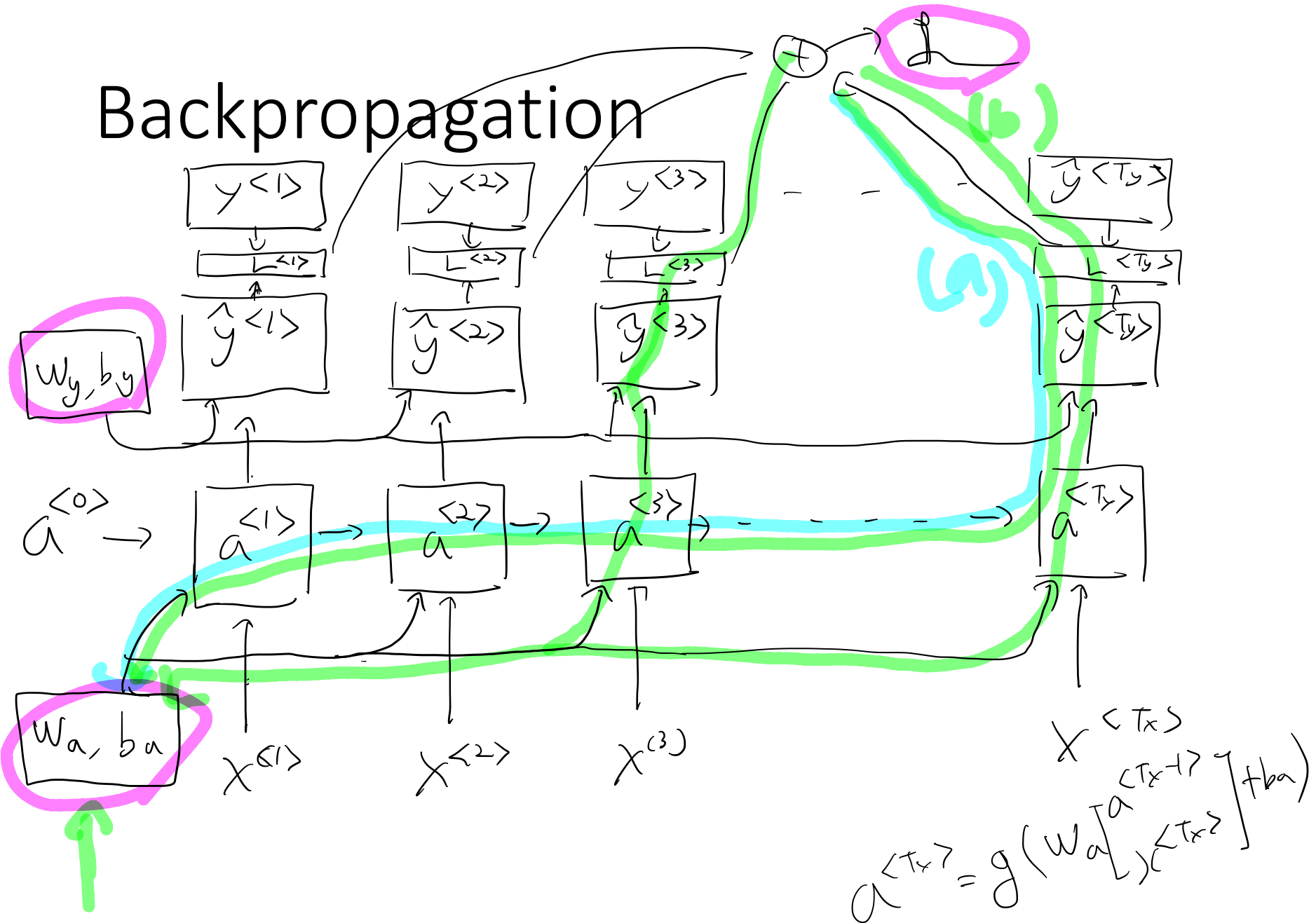$\hat{y}^{<1>}$  $\hat{y}^{<2>}$  $\hat{y}^{<3>}$  $\hat{y}^{<T_y>}$

$b_y$  $W_{ay}$

$W_{aa}$  $a^{<1>}$  $a^{<2>}$  $a^{<T_x-1>}$

$a^{<0>}$  $b_a$

$W_{ax}$

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

tanh, sigmoid, ReLu

$W =$

$\begin{bmatrix} W_{aa}, W_{ax} \\ W_{ay}, b_a, b_y \end{bmatrix}$

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

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

# RNN- Forward propagation

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

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

$$\frac{a^{<t>}}{h^{<t>}} = g\left(W_a \cdot \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} + b_a\right)$$

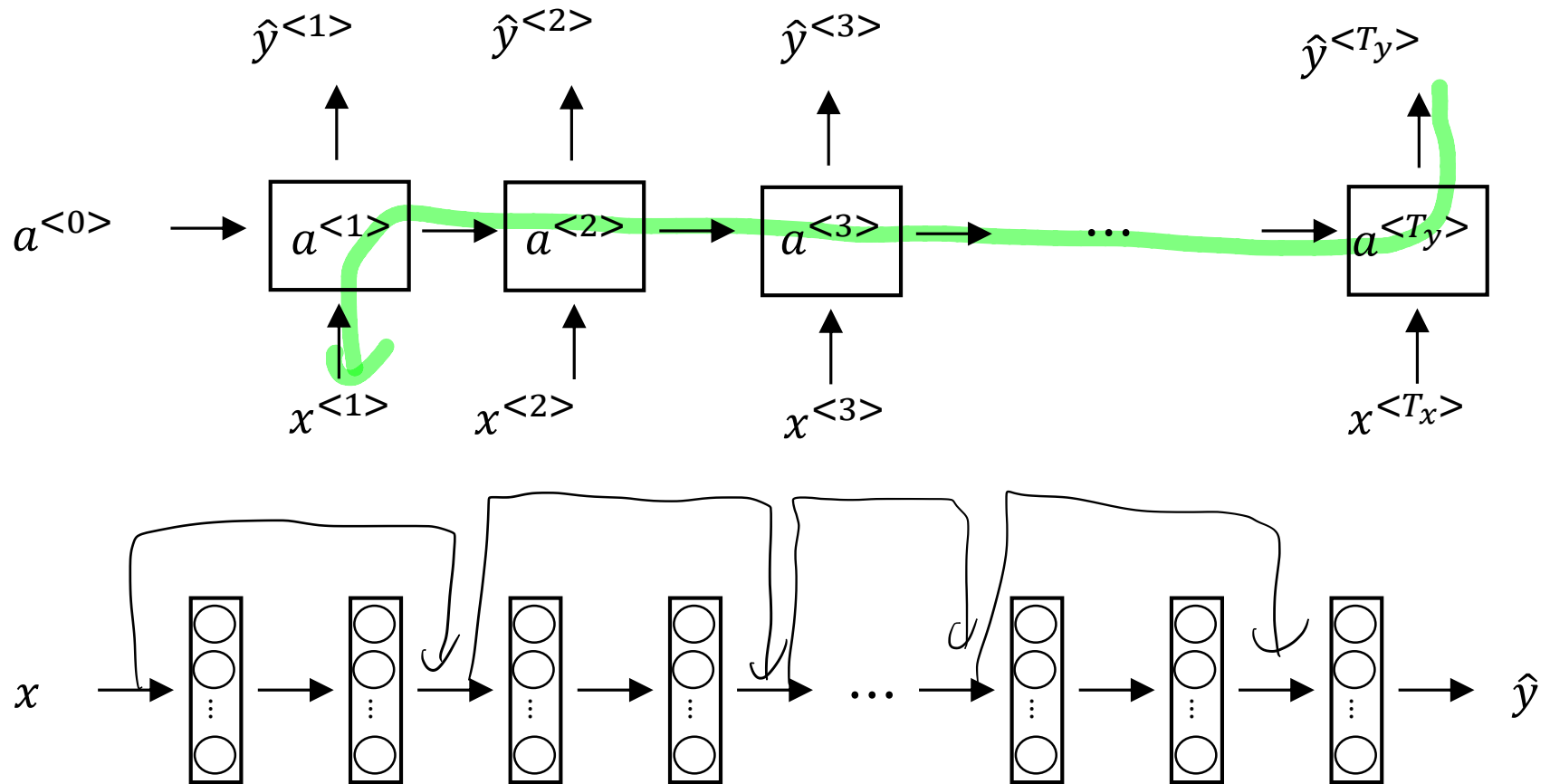$$\hat{y}^{<t>} = g\left(W_y \, a^{<t>} + b_y\right)$$

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

# Backpropagation

$y^{<1>}$

$y^{<2>}$

$y^{<3>}$

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

$L^{<1>}$

$L^{<2>}$

$L^{<3>}$

$L^{<T_y>}$

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

$\hat{y}^{<2>}$

$\hat{y}^{<3>}$

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

$W_{y}, b_y$

$a^{<0>} \longrightarrow$

$a^{<1>}$

$a^{<2>}$

$a^{<3>}$

$a^{<T_y>}$

$W_a, b_a$

$x^{<1>}$

$x^{<2>}$

$x^{(3)}$

$x^{<T_x>}$

$$a^{<T_x>} = g\left(W_a\begin{bmatrix} a^{<T_x-1>} \\ x^{<T_x>} \end{bmatrix} + b_a\right)$$

# Backpropagation

$$dW_a =$$

$$\nabla_{W_a} \mathcal{L} =$$

$$= \sum_{i=1}^{T_y} \nabla_{W_a} L^{\langle e \rangle} = \sum_{i=1}^{T_y} \frac{\partial L^{\langle i \rangle}}{\partial \hat{y}^{\langle i \rangle}} \cdot \boxed{\frac{\partial \hat{y}^{\langle i \rangle}}{\partial W_a}}$$

$$\frac{\partial L^{\langle i \rangle}}{\partial W_a}$$

$$= \sum_{i=1}^{T_y} \frac{\partial L^{\langle i \rangle}}{\partial \hat{y}^{\langle i \rangle}} \boxed{\sum_{j=1}^{i} \frac{\partial \hat{y}^{\langle i \rangle}}{\partial a^{\langle j \rangle}}} \cdot \frac{\partial a^{\langle j \rangle}}{\partial W_a}$$

$$= \sum_{i=1}^{T_y} \frac{\partial L^{\langle i \rangle}}{\partial \hat{y}^{\langle i \rangle}} \sum_{j=1}^{i} \boxed{\frac{\partial \hat{y}^{\langle i \rangle}}{\partial a^{\langle i \rangle}} \cdot \frac{\partial a^{\langle i \rangle}}{\partial a^{\langle i-1 \rangle}} \cdots \frac{\partial a^{\langle j+1 \rangle}}{\partial a^{\langle j \rangle}}} \frac{\partial a^{\langle j \rangle}}{\partial W_a}$$
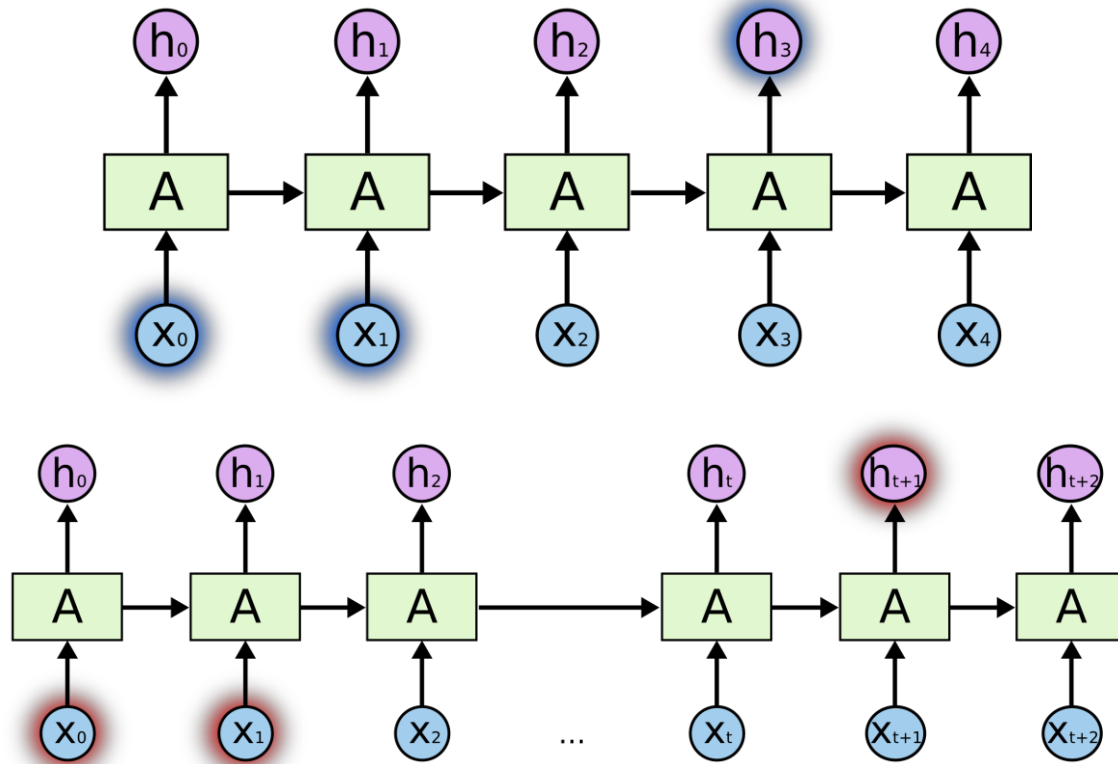
# Vanishing Gradient

# Long-term dependency

LSTM!

GRU!



- RNN cannot learn the long-term dependency in the bellow

# LSTM



$g^{(t-1)}$

$g(W_a \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} + b_a)$

$y$

tanh

$W_y, b_y$

$h_{t-1}$   $h_t$   $h_{t+1}$

RNN

tanh

$W_a, b_a$

$x_{t-1}$   $x_t$   $x_{t+1}$

sigmoid.

$f(x) = \dfrac{e^x}{1+e^x} \in [0,1]$

$h_{t-1}$   $h_t$   $h_{t+1}$

LSTM

$\sigma$   $\sigma$   tanh   $\sigma$

switch.

$x_{t-1}$   $x_t$   $x_{t+1}$

# LSTM

Cell state : long term memory



$C_{t-1}$
$C_t$
$f_t$  $i_t$  $\tilde{C}_t$  $o_t$
tanh
$h_t$
$\sigma$  $\sigma$  tanh  $\sigma$
$h_{t-1}$
$h_t$
$x_t$

a conveyor belt of information

$$\sigma(x) = \frac{e^x}{1+e^x}$$

# LSTM



$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] \ + \ b_f \right)$$

$$i_t = \sigma \left( W_i \cdot [h_{t-1}, x_t] \ + \ b_i \right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] \ + \ b_C)$$

# LSTM



$$(1 - f_t)$$

$$\downarrow$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

# Variants on LSTM



$$f_t = \sigma\left(W_f \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_t] \;+\; b_f\right)$$
$$i_t = \sigma\left(W_i \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_t] \;+\; b_i\right)$$
$$o_t = \sigma\left(W_o \cdot [\boldsymbol{C_t}, h_{t-1}, x_t] \;+\; b_o\right)$$

$$C_t = f_t * C_{t-1} + (\boldsymbol{1 - f_t}) * \tilde{C}_t$$

# GRU



$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$