



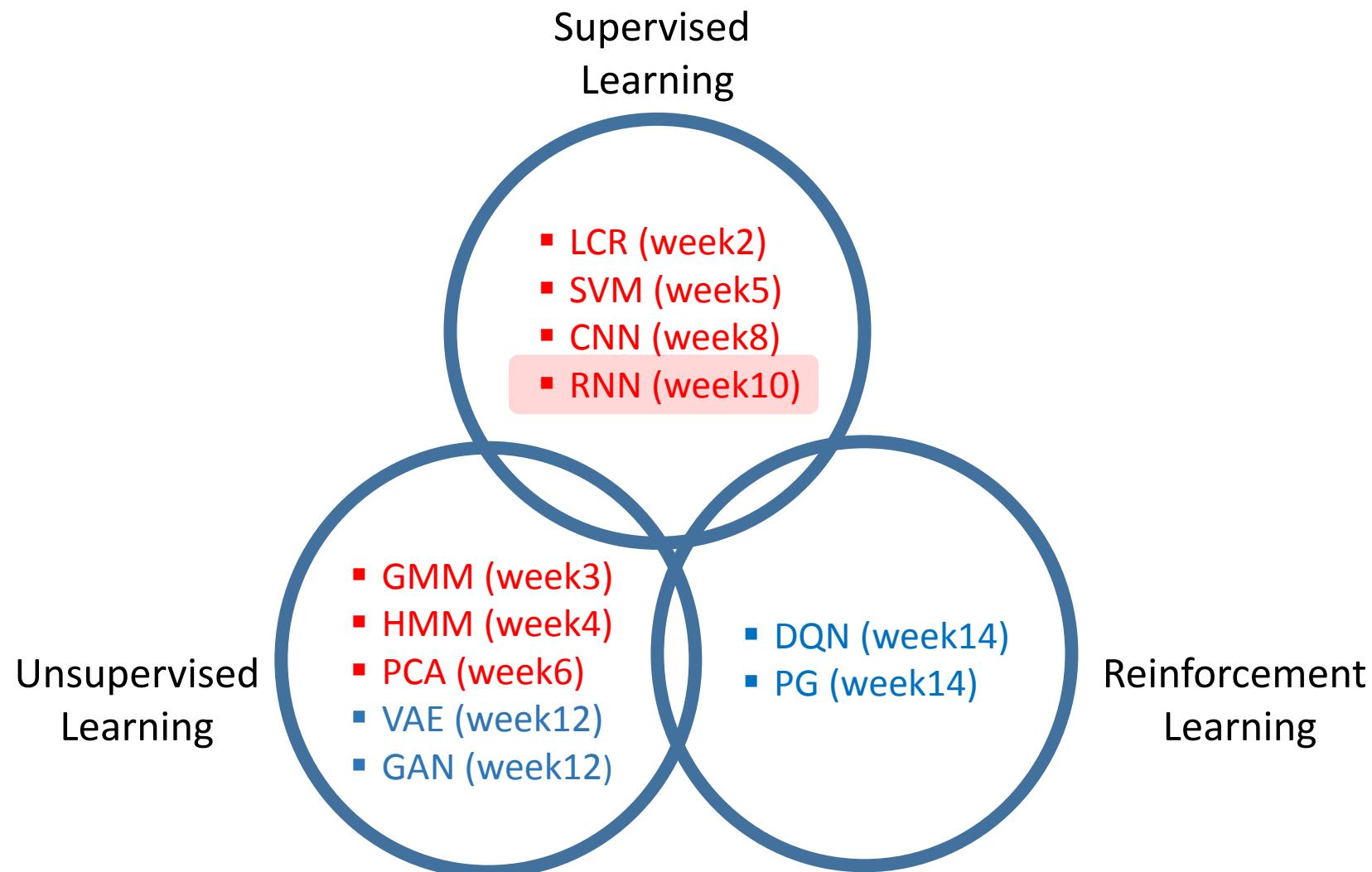
Practical Machine Learning

Lecture 10 Recurrent Neural Networks (RNN) & Long Short Term Memory (LSTM)

Dr. Suyong Eum



Where we are



You are going to learn

- ❑ Recurrent Neural Networks (RNN)
- ❑ Long Short Term Memory (LSTM)
- ❑ Backpropagation in RNN and LSTM

Recurrent Neural Networks (RNN)

Recurrent Neural Network (RNN)

- ❑ Feed forward neural networks

- Independent
- Fixed length

- ❑ Recurrent Neural Networks

- Temporal dependencies
- Variable sequence length

Some interesting applications

1. Music composition

- <http://people.idsia.ch/~juergen/blues/>



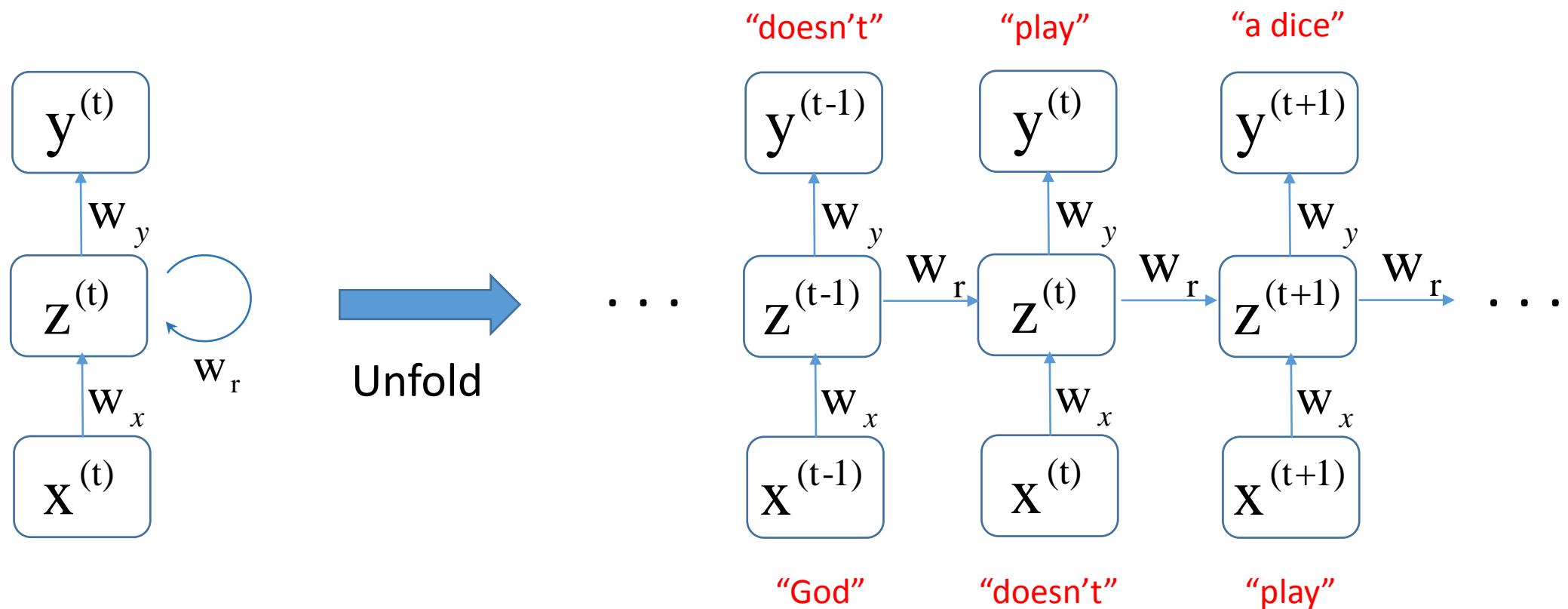
2. Writing a poem

- <https://github.com/dvictor/lstm-poetry>

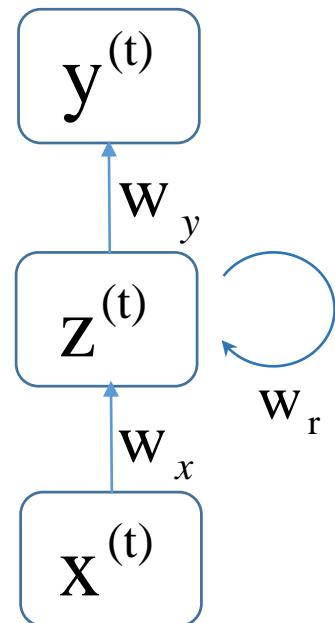
A butterfly in the sun

Just because I know that I should leave this heart for you
You said I was falling apart
I wish I were you
I wanted you to know how I feel
I could have settled it all
It's time to go and do it big and you can be my side
I can't believe it when I see you
I'm lost in the world and I can't see you cry
I'm asking you to love me then let me go
I can't stop this way

Recurrent Neural Network (RNN)

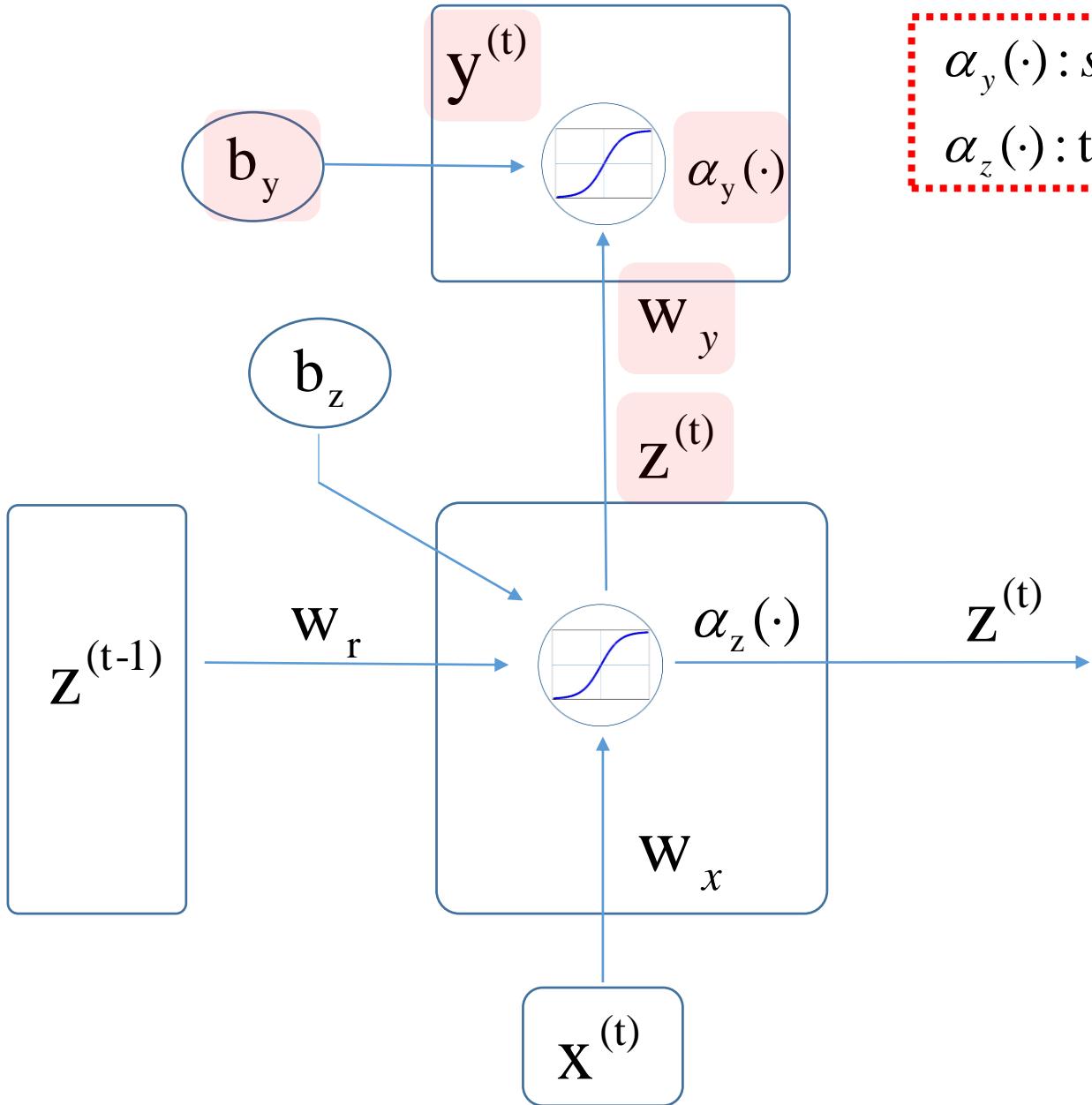


Recurrent Neural Network (RNN): inner structure



$$z^{(t)} = \alpha_z(w_x x^{(t)} + w_r z^{(t-1)} + b_z)$$

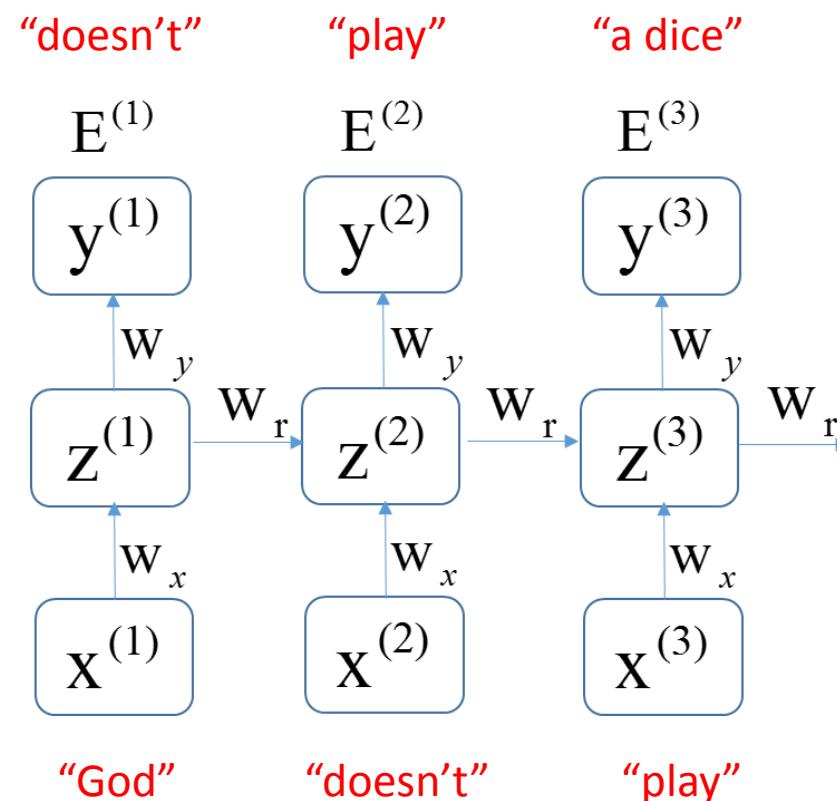
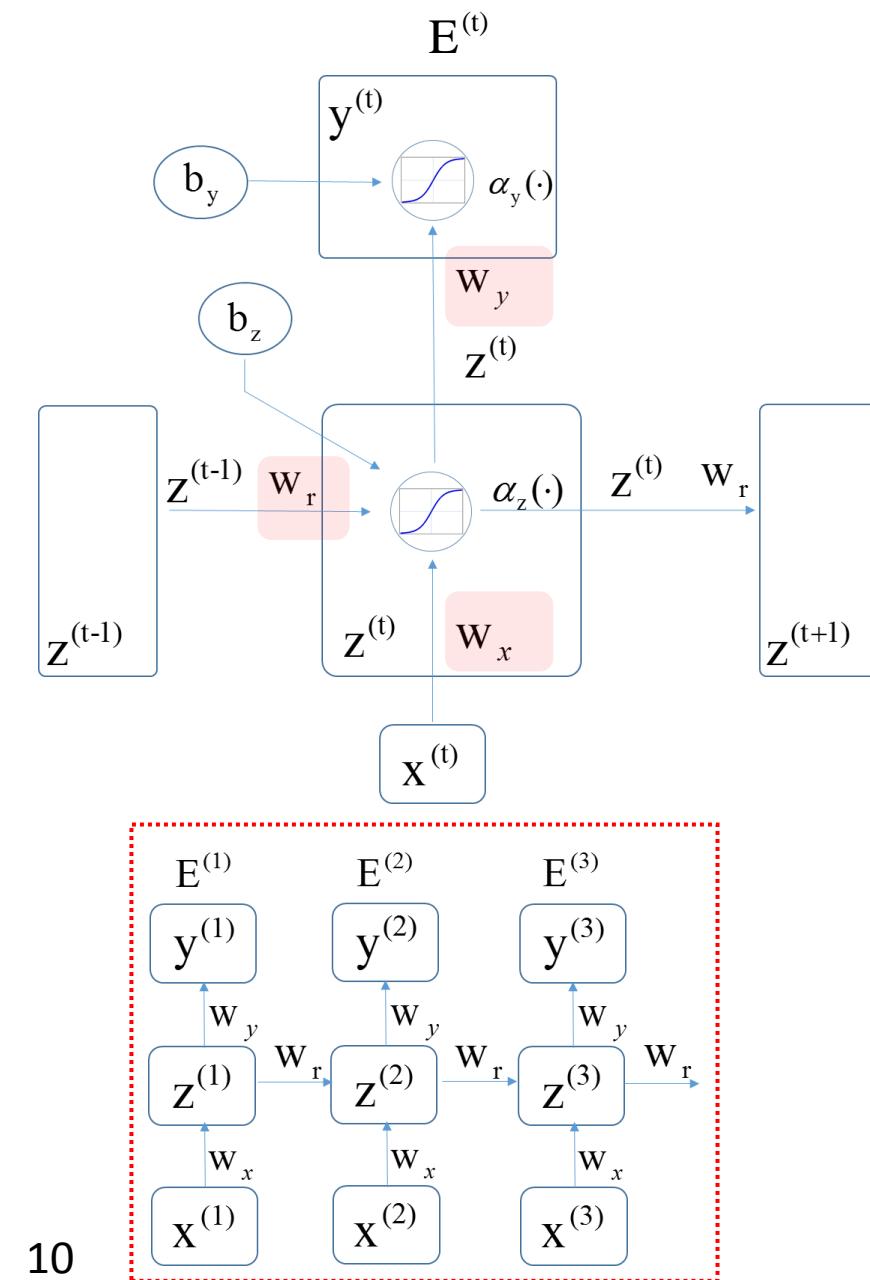
$$y^{(t)} = \alpha_y(w_y z^{(t)} + b_y)$$



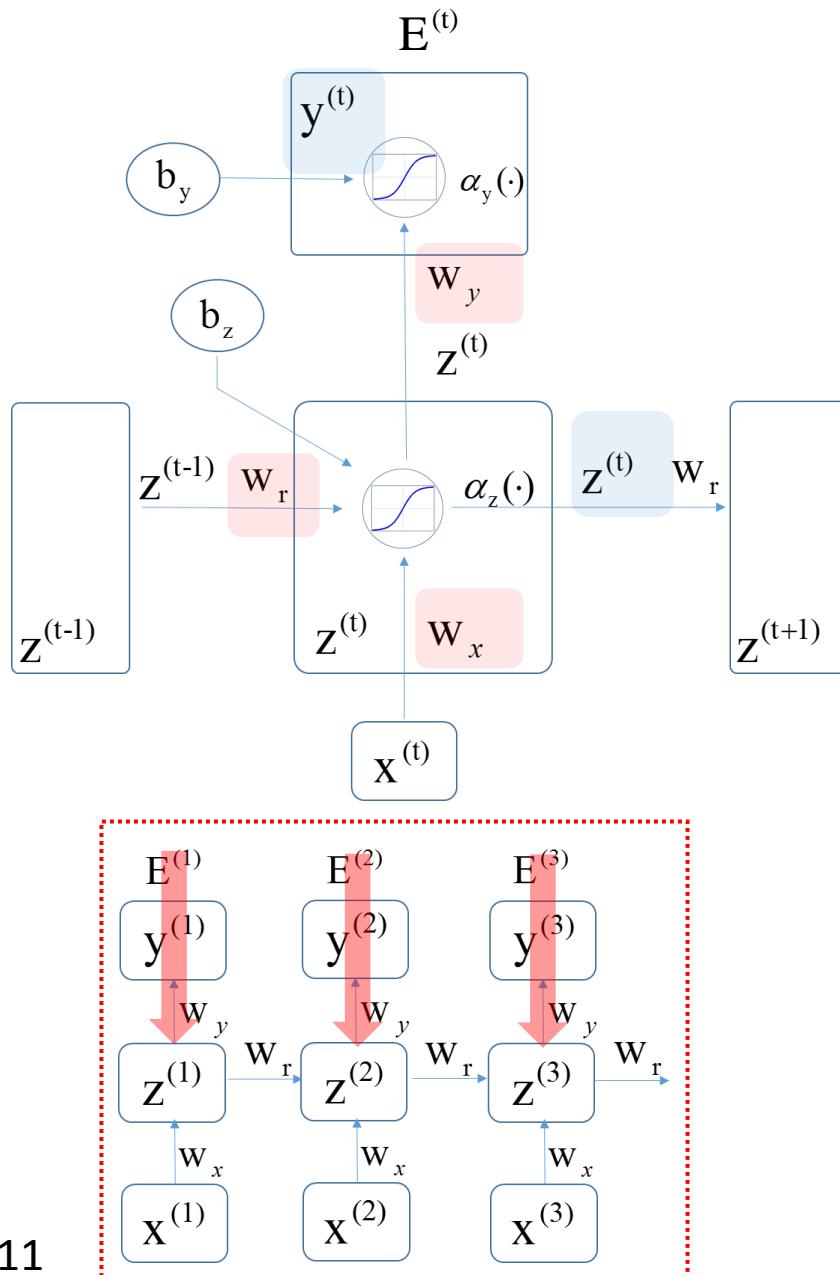
$\alpha_y(\cdot)$: soft max
 $\alpha_z(\cdot)$: tanh

Backpropagation through time (BPTT) in RNN

Recurrent Neural Network (RNN): Backpropagation through time (BPTT)



Recurrent Neural Network (RNN): Backpropagation through time (BPTT)



$$\begin{aligned} z^{(t)} &= \alpha_z(w_x x^{(t)} + w_r z^{(t-1)} + b_z) \\ y^{(t)} &= \alpha_y(w_y z^{(t)} + b_y) \end{aligned}$$

$$\begin{aligned} s_z^{(t)} &= w_x x^{(t)} + w_r z^{(t-1)} + b_z \\ s_y^{(t)} &= w_y z^{(t)} + b_y \end{aligned}$$

$$\frac{\partial E}{\partial w_y} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_y} = \frac{\partial E^{(1)}}{\partial w_y} + \frac{\partial E^{(2)}}{\partial w_y} + \frac{\partial E^{(3)}}{\partial w_y}$$

$$\frac{\partial E}{\partial w_r} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_r} = \frac{\partial E^{(1)}}{\partial w_r} + \frac{\partial E^{(2)}}{\partial w_r} + \frac{\partial E^{(3)}}{\partial w_r}$$

$$\frac{\partial E}{\partial w_x} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_x} = \frac{\partial E^{(1)}}{\partial w_x} + \frac{\partial E^{(2)}}{\partial w_x} + \frac{\partial E^{(3)}}{\partial w_x}$$

$$\frac{\partial E}{\partial w_y} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_y} = \frac{\partial E^{(1)}}{\partial w_y} + \frac{\partial E^{(2)}}{\partial w_y} + \frac{\partial E^{(3)}}{\partial w_y}$$

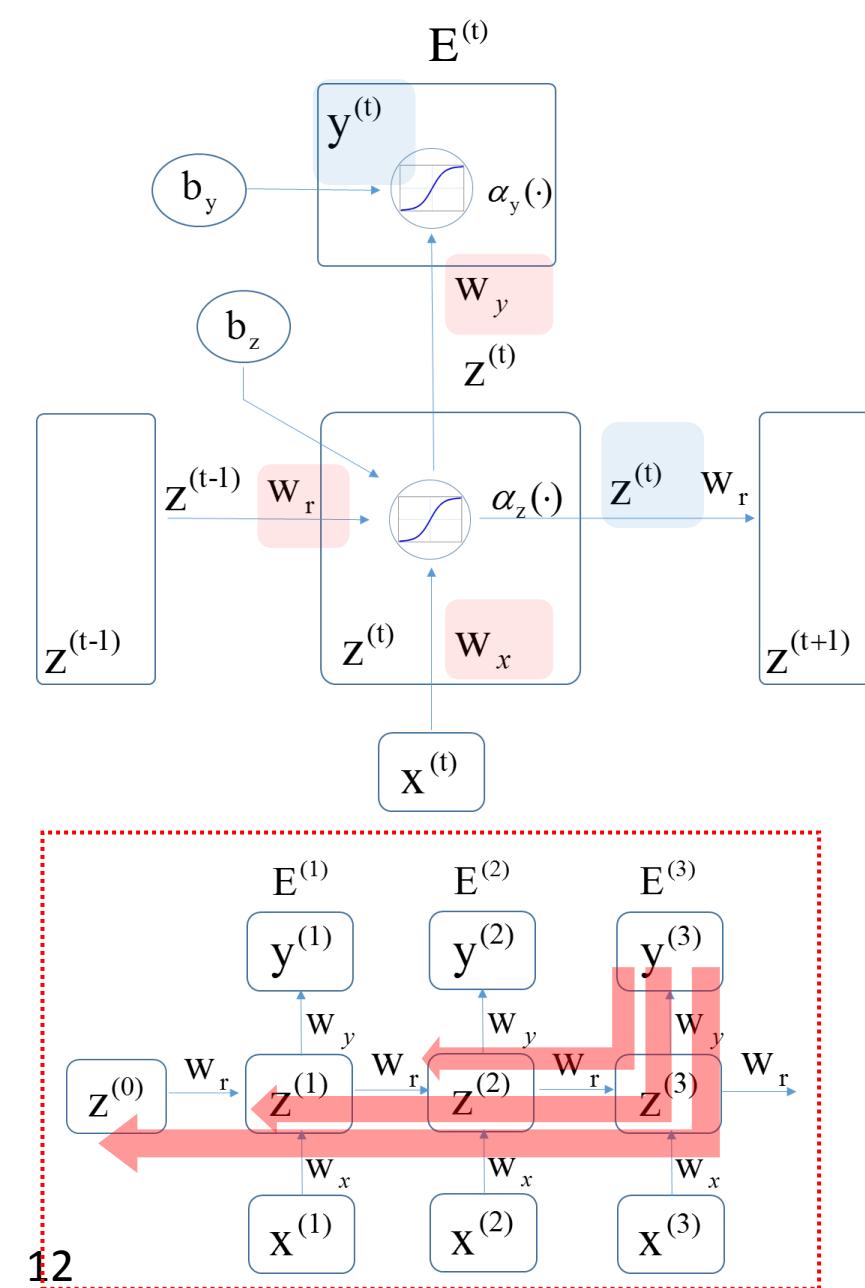
$$\frac{\partial E^{(1)}}{\partial w_y} = \frac{\partial E^{(1)}}{\partial y^{(1)}} \frac{\partial y^{(1)}}{\partial s_y^{(1)}} \frac{\partial s_y^{(1)}}{\partial w_y}$$

$$\frac{\partial E^{(2)}}{\partial w_y} = \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial w_y}$$

$$\frac{\partial E^{(3)}}{\partial w_y} = \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial w_y}$$

→ Linear input function
 → Activation function
 → Error function

Recurrent Neural Network (RNN): Backpropagation through time (BPTT)



$$z^{(t)} = \alpha_z(w_x x^{(t)} + w_r z^{(t-1)} + b_z)$$

$$y^{(t)} = \alpha_y(w_y z^{(t)} + b_y)$$

$$s_z^{(t)} = w_x x^{(t)} + w_r z^{(t-1)} + b_z$$

$$s_y^{(t)} = w_y z^{(t)} + b_y$$

$$\frac{\partial E}{\partial w_y} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_y} = \frac{\partial E^{(1)}}{\partial w_y} + \frac{\partial E^{(2)}}{\partial w_y} + \frac{\partial E^{(3)}}{\partial w_y}$$

$$\frac{\partial E}{\partial w_r} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_r} = \frac{\partial E^{(1)}}{\partial w_r} + \frac{\partial E^{(2)}}{\partial w_r} + \frac{\partial E^{(3)}}{\partial w_r}$$

$$\frac{\partial E}{\partial w_x} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_x} = \frac{\partial E^{(1)}}{\partial w_x} + \frac{\partial E^{(2)}}{\partial w_x} + \frac{\partial E^{(3)}}{\partial w_x}$$

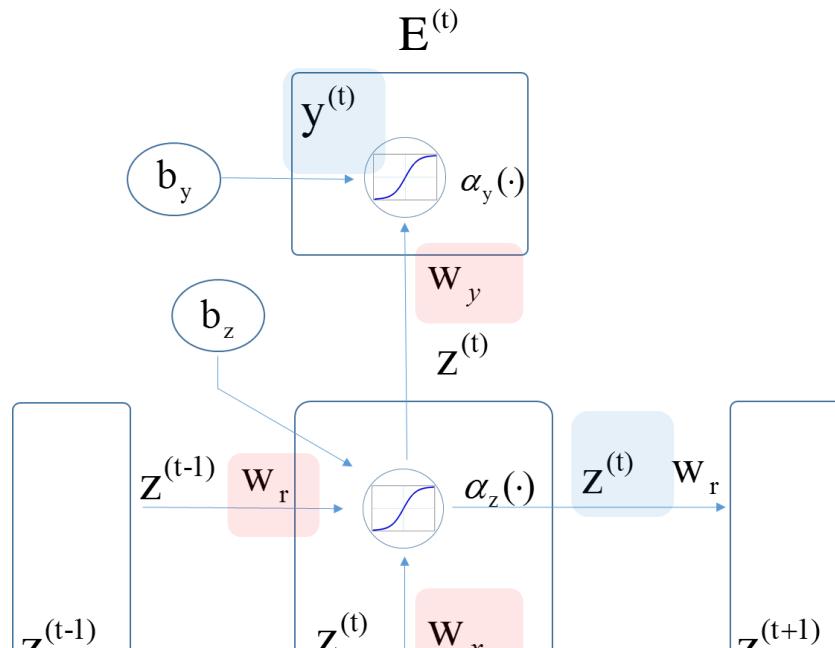
$$\frac{\partial E}{\partial w_r} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_r} = \frac{\partial E^{(1)}}{\partial w_r} + \frac{\partial E^{(2)}}{\partial w_r} + \frac{\partial E^{(3)}}{\partial w_r}$$

$$\sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_r} = \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial w_r} \quad t = 3$$

$$+ \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial s_z^{(2)}} \frac{\partial s_z^{(2)}}{\partial w_r} \quad t = 2$$

$$+ \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial s_z^{(2)}} \frac{\partial s_z^{(2)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial s_z^{(1)}} \frac{\partial s_z^{(1)}}{\partial w_r} \quad t = 1$$

Recurrent Neural Network (RNN): Backpropagation through time (BPTT)



$$\begin{aligned} z^{(t)} &= \alpha_z(w_x x^{(t)} + w_r z^{(t-1)} + b_z) \\ y^{(t)} &= \alpha_y(w_y z^{(t)} + b_y) \\ s_z^{(t)} &= w_x x^{(t)} + w_r z^{(t-1)} + b_z \\ s_y^{(t)} &= w_y z^{(t)} + b_y \end{aligned}$$

$$\frac{\partial E}{\partial w_y} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_y} = \frac{\partial E^{(1)}}{\partial w_y} + \frac{\partial E^{(2)}}{\partial w_y} + \frac{\partial E^{(3)}}{\partial w_y}$$

$$\frac{\partial E}{\partial w_r} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_r} = \frac{\partial E^{(1)}}{\partial w_r} + \frac{\partial E^{(2)}}{\partial w_r} + \frac{\partial E^{(3)}}{\partial w_r}$$

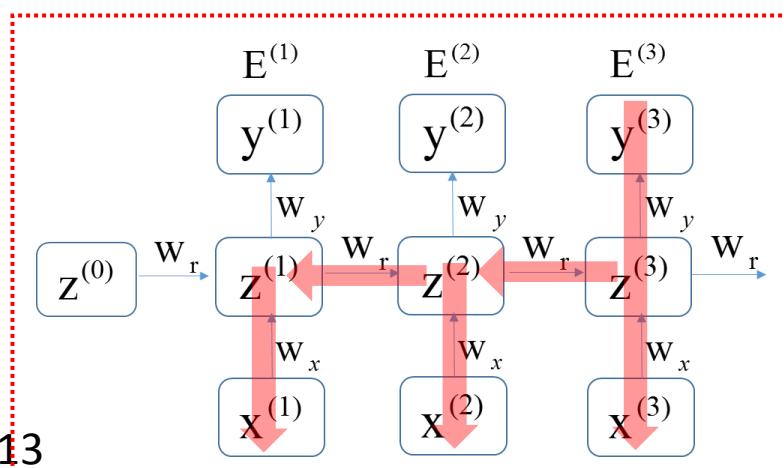
$$\frac{\partial E}{\partial w_x} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_x} = \frac{\partial E^{(1)}}{\partial w_x} + \frac{\partial E^{(2)}}{\partial w_x} + \frac{\partial E^{(3)}}{\partial w_x}$$

$$\frac{\partial E}{\partial w_x} = \sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_x} = \frac{\partial E^{(1)}}{\partial w_x} + \frac{\partial E^{(2)}}{\partial w_x} + \frac{\partial E^{(3)}}{\partial w_x}$$

$$\sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_r} = \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial w_x} \quad t = 3$$

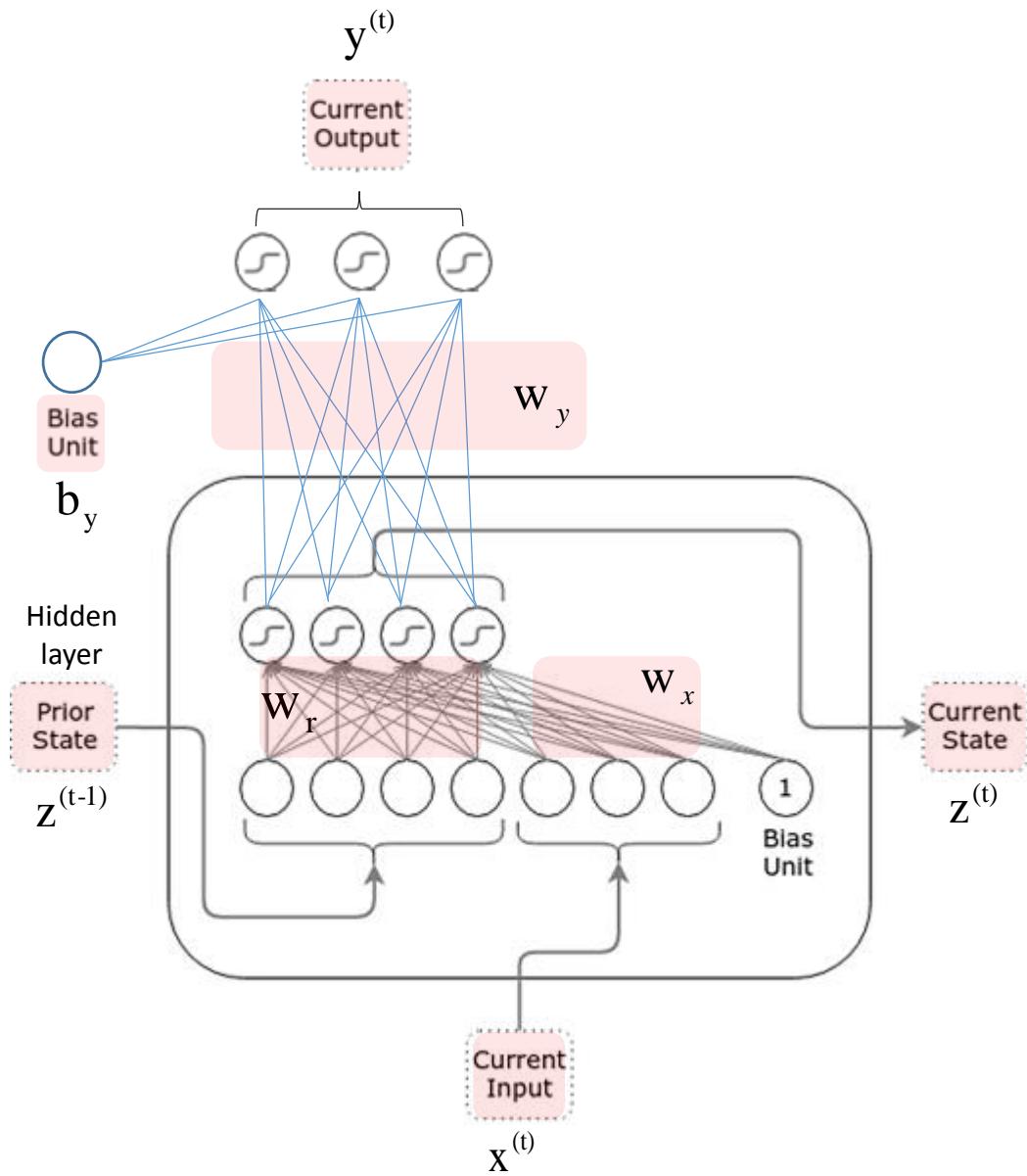
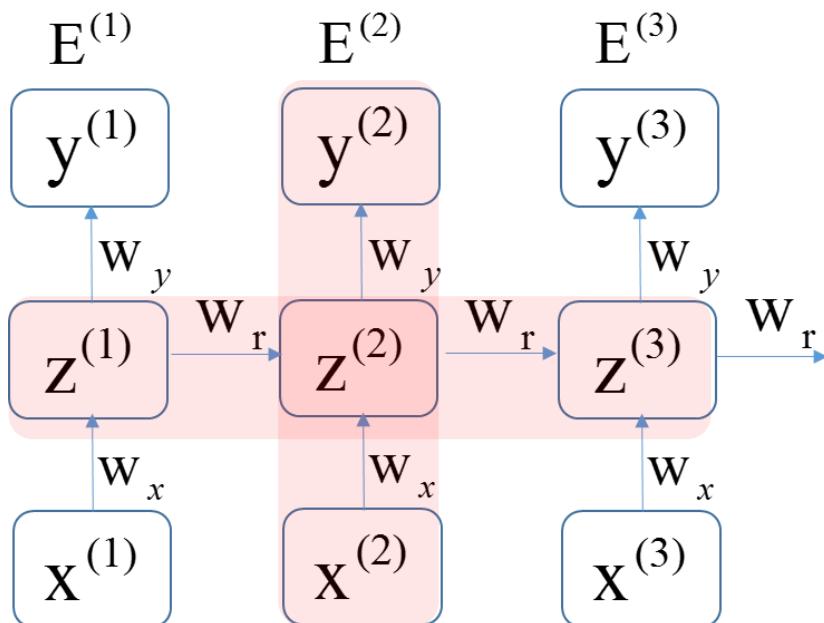
$$+ \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial s_z^{(2)}} \frac{\partial s_z^{(2)}}{\partial w_x} \quad t = 2$$

$$+ \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial s_z^{(2)}} \frac{\partial s_z^{(2)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial s_z^{(1)}} \frac{\partial s_z^{(1)}}{\partial w_x} \quad t = 1$$

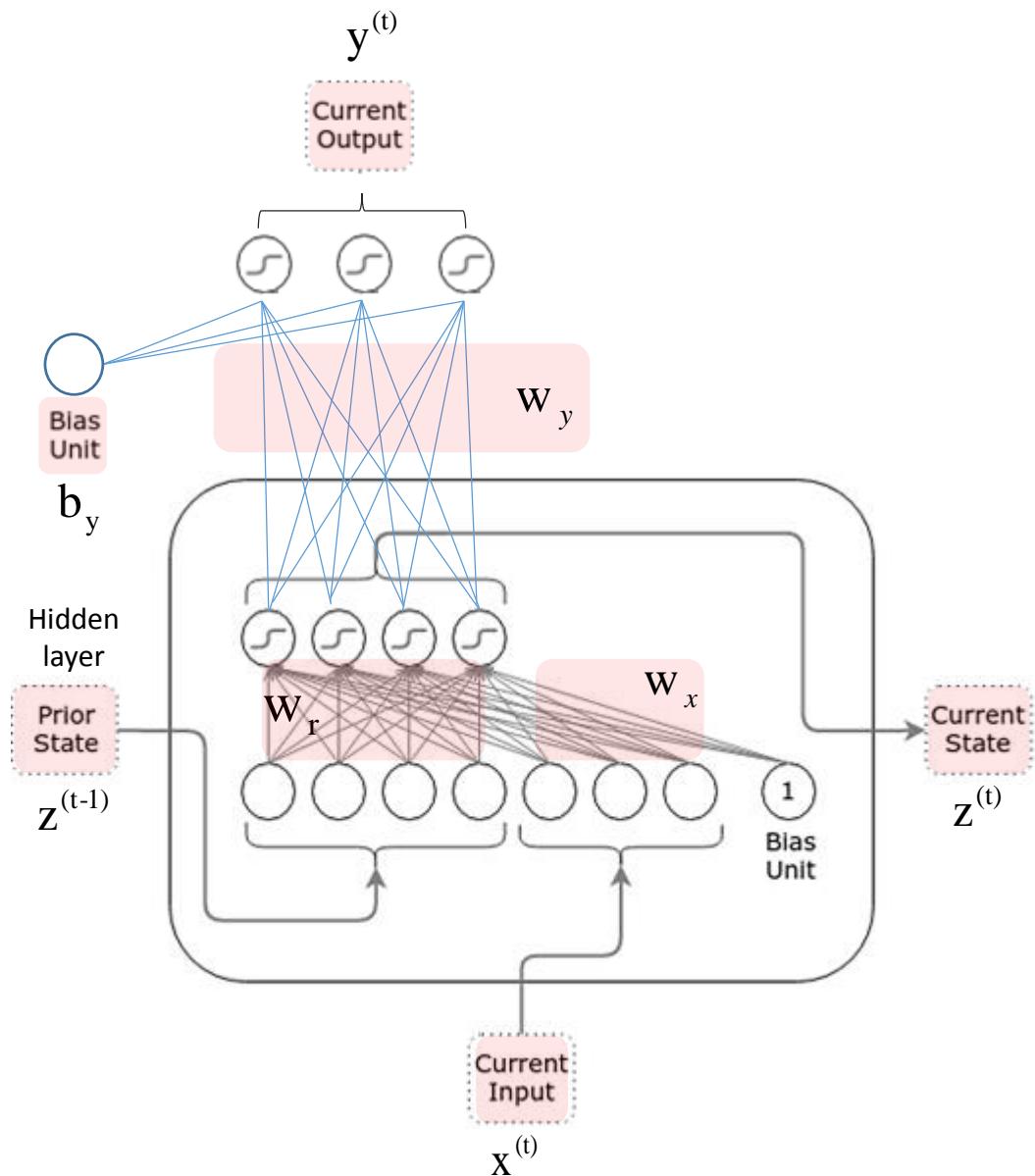


Connectivity of neurons in vanilla RNN component

Operation of RNN: Connectivity of neurons in vanilla RNN component



Operation of RNN: Connectivity of neurons in vanilla RNN component



where:

$$\begin{array}{lll} z^{(t-1)}, z^{(t)} \in \mathbb{R}^n & W_r \in \mathbb{R}^{n \times n} & b_z \in \mathbb{R}^n \\ x^{(t)} \in \mathbb{R}^m & W_x \in \mathbb{R}^{n \times m} & b_y \in \mathbb{R}^k \\ y^{(t)} \in \mathbb{R}^k & W_y \in \mathbb{R}^{k \times n} & \end{array}$$

$$z^{(t)} = \alpha_z (W_x x^{(t)} + W_r z^{(t-1)} + b_z)$$

$$(n \times 1) = (n \times m) (m \times 1) + (n \times n) (n \times 1) + (n \times 1)$$

$$y^{(t)} = \alpha_y (W_y z^{(t)} + b_y)$$

$$(k \times 1) = (k \times n) (n \times 1) + (k \times 1)$$

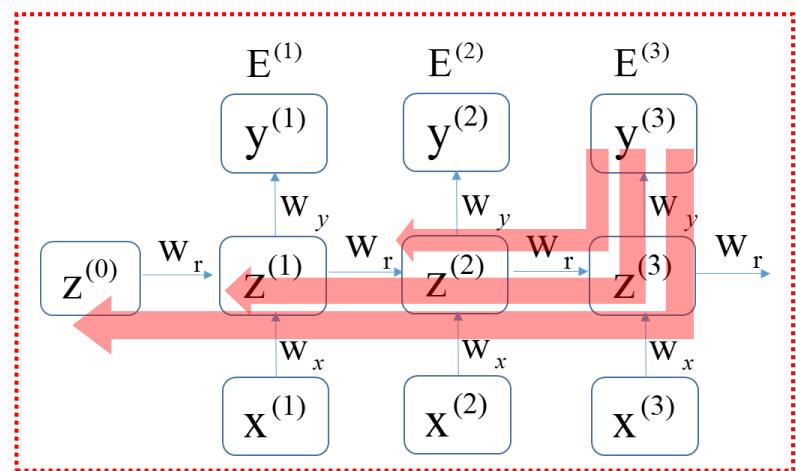
Long Short Term Memory (LSTM)

Long Short Term Memory (LSTM) network

- Long Short Term Memory (LSTM) architecture was motivated to overcome the problem: error is not back-propagated properly to the end of RNN architecture.

- Hochreiter and Schmidhuber (1997) proposed the Long Short-Term Memory (LSTM) cell which includes “*a memory unit*”:
 - 1) A cell with a number of components that together act similar to a memory cell.
 - 2) Inside one cell, multiple layers called “*gates*” are used.
 - 1) *Forget gate*
 - 2) *Input gate*
 - 3) *Output gate*

RNN using BPTT: Vanishing and exploding problems



$$\sum_{t=1}^3 \frac{\partial E^{(t)}}{\partial w_r} = \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial w_r}$$

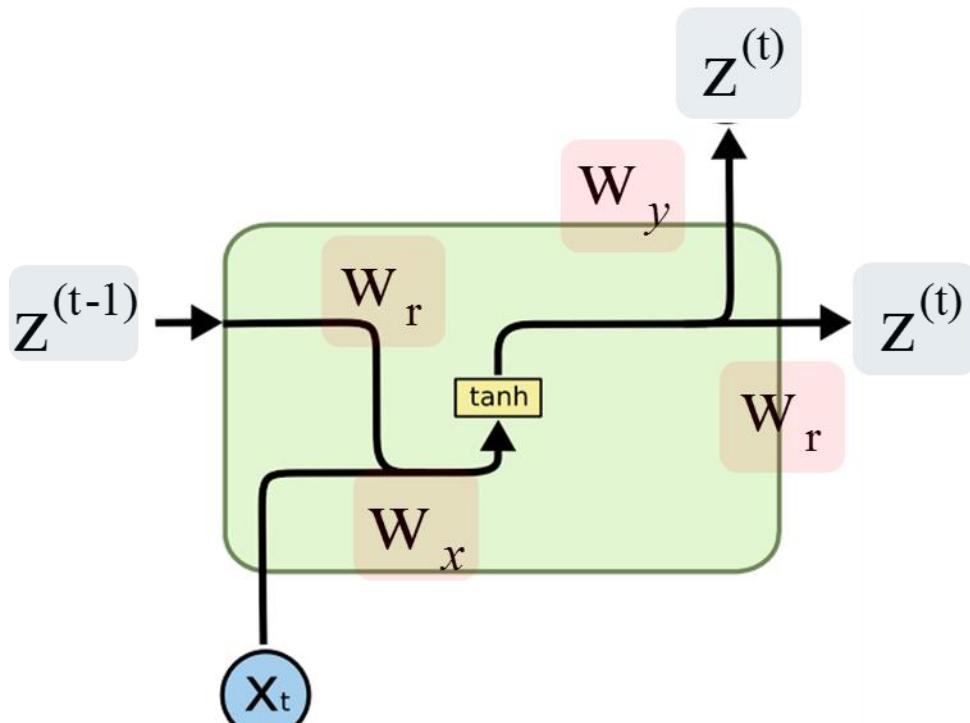
$$+ \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial s_z^{(2)}} \frac{\partial s_z^{(2)}}{\partial w_r}$$

$$+ \frac{\partial E^{(3)}}{\partial y^{(3)}} \frac{\partial y^{(3)}}{\partial s_y^{(3)}} \frac{\partial s_y^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial s_z^{(3)}} \frac{\partial s_z^{(3)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial s_z^{(2)}} \frac{\partial s_z^{(2)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial s_z^{(1)}} \frac{\partial s_z^{(1)}}{\partial w_r}$$

Activation function (tanh)

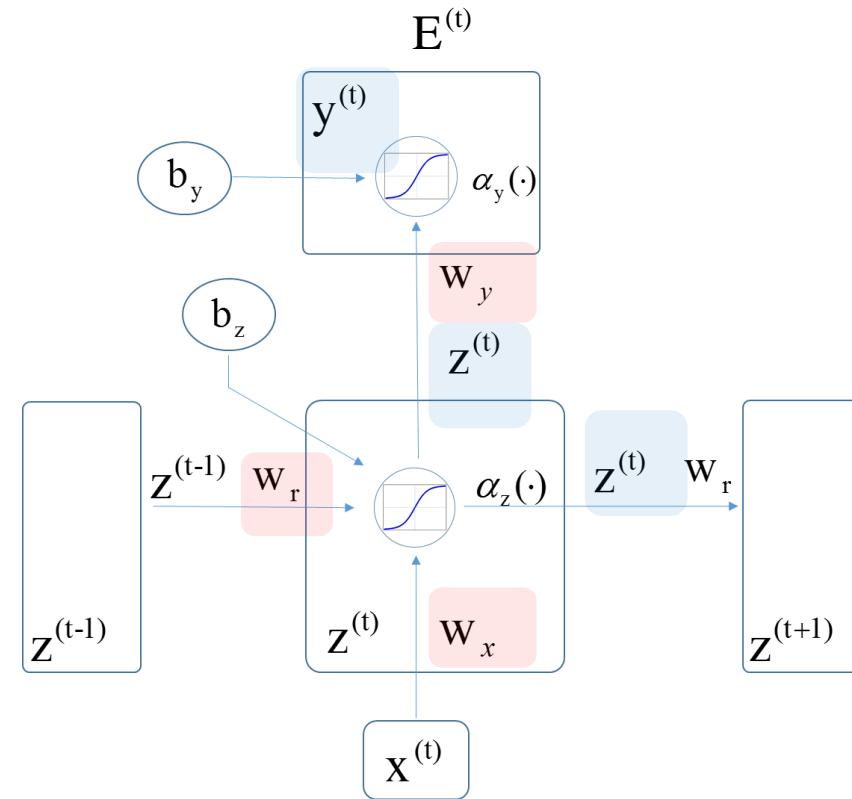
$$\frac{d\sigma(a)}{da} = 1 - \sigma(a)^2$$

Long Short Term Memory (LSTM) network



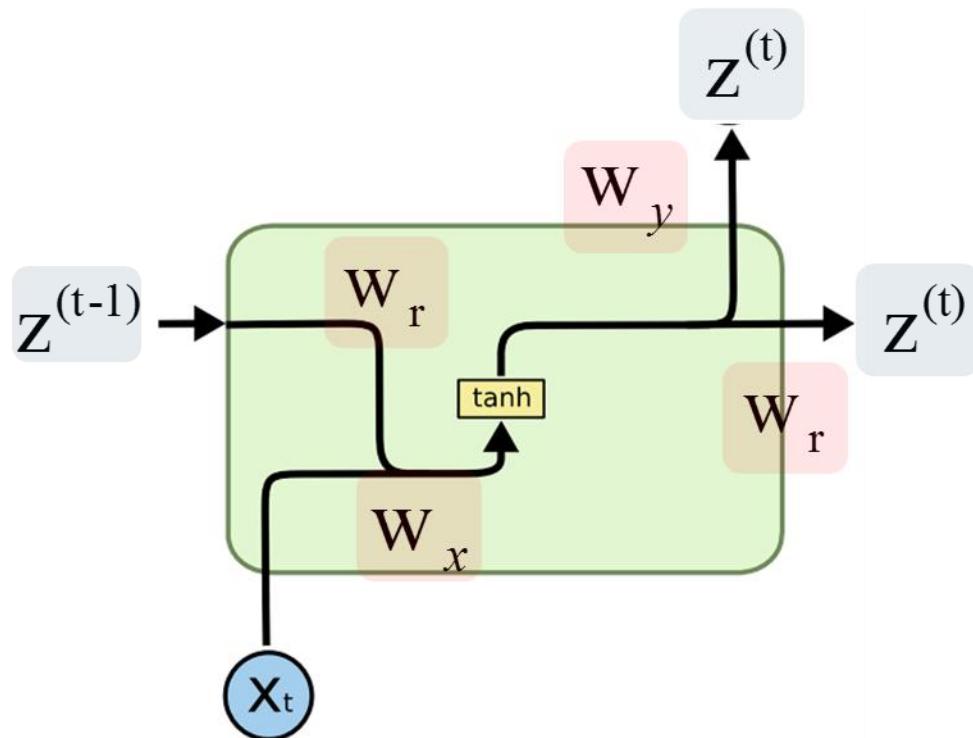
$$Z^{(t)} = \alpha_z (W_x X^{(t)} + W_r Z^{(t-1)} + b_z)$$

$$y^{(t)} = \alpha_y (W_y Z^{(t)} + b_y)$$



RNN

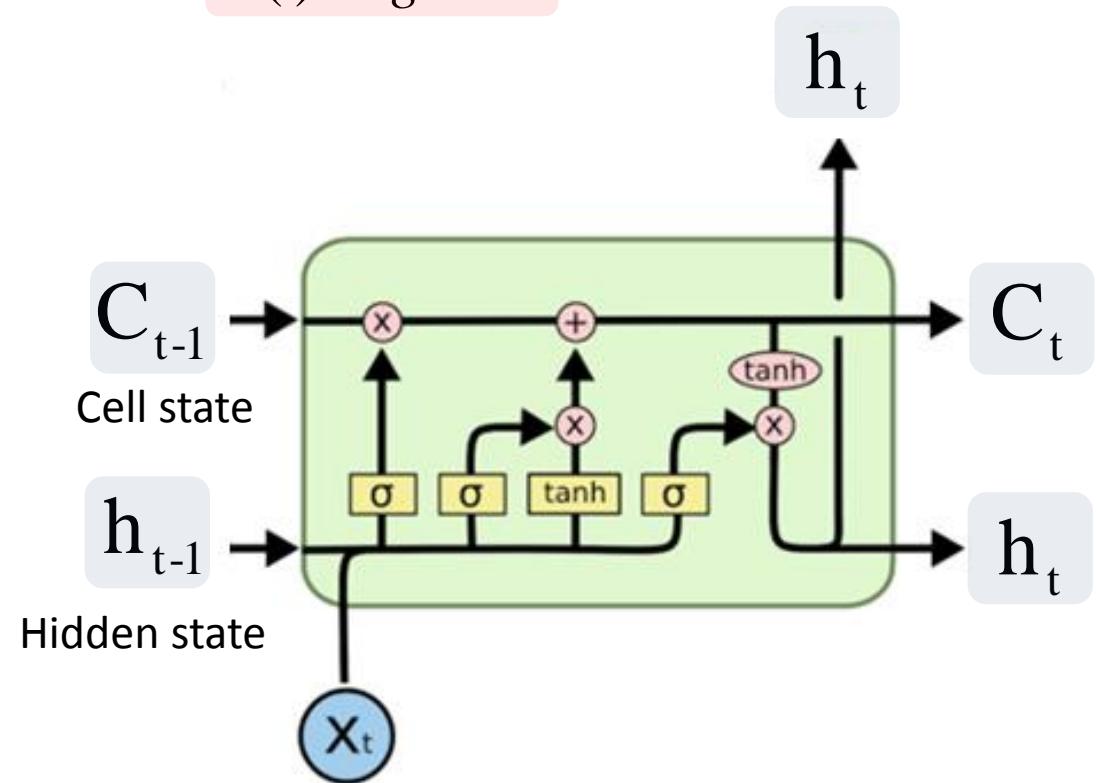
Long Short Term Memory (LSTM) network



$$z^{(t)} = \alpha_z (w_x x^{(t)} + w_r z^{(t-1)} + b_z)$$

$$y^{(t)} = \alpha_y (w_y z^{(t)} + b_y)$$

$\sigma(\cdot)$: sigmoid

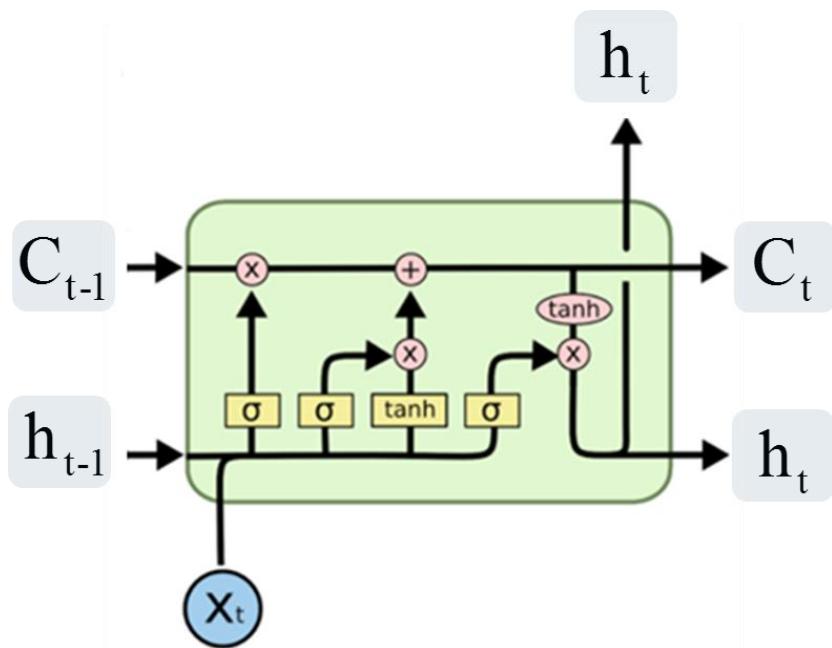


- Cell state: a memory of LSTM cell
- Hidden state: an output of this cell

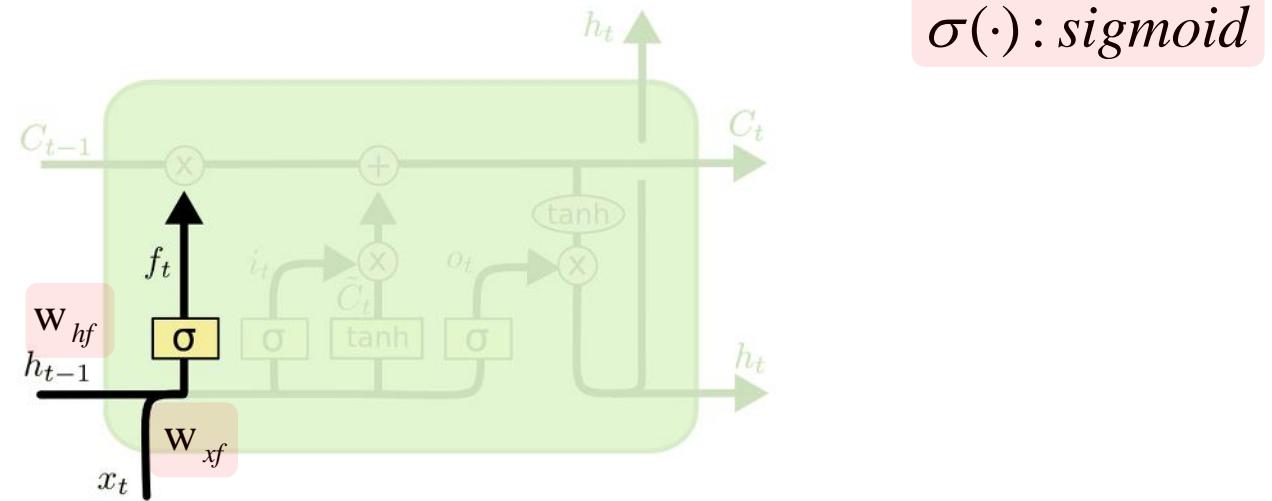
RNN

LSTM

Long Short Term Memory (LSTM) network: forget gate



LSTM

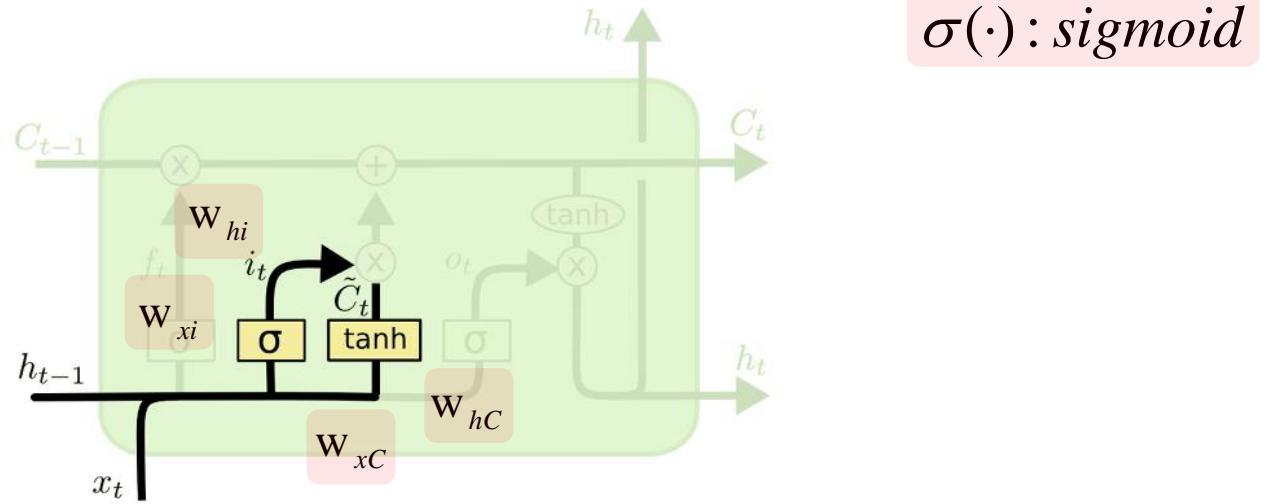
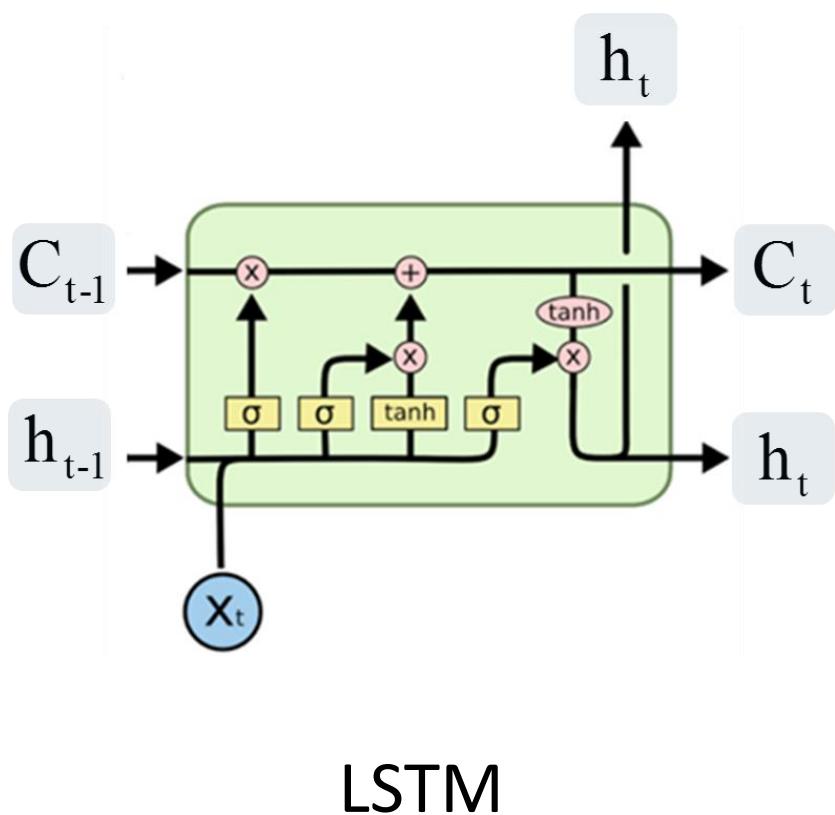


$$f_t = \text{sigm}(w_{xf} x_t + w_{hf} h_{t-1} + b_f)$$

❑ Forget gate layer

- Decide how much " C_{t-1} " is forgotten?
- If " f_t " is zero, forget " C_{t-1} " completely.
- If " f_t " is one, do not forget " C_{t-1} " at all.

Long Short Term Memory (LSTM) network: input gate

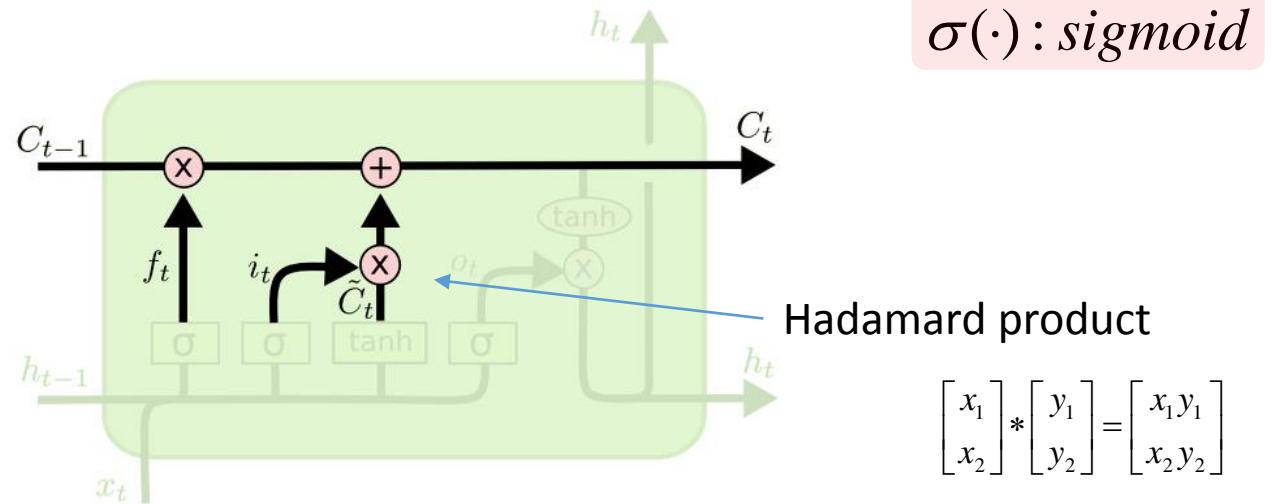
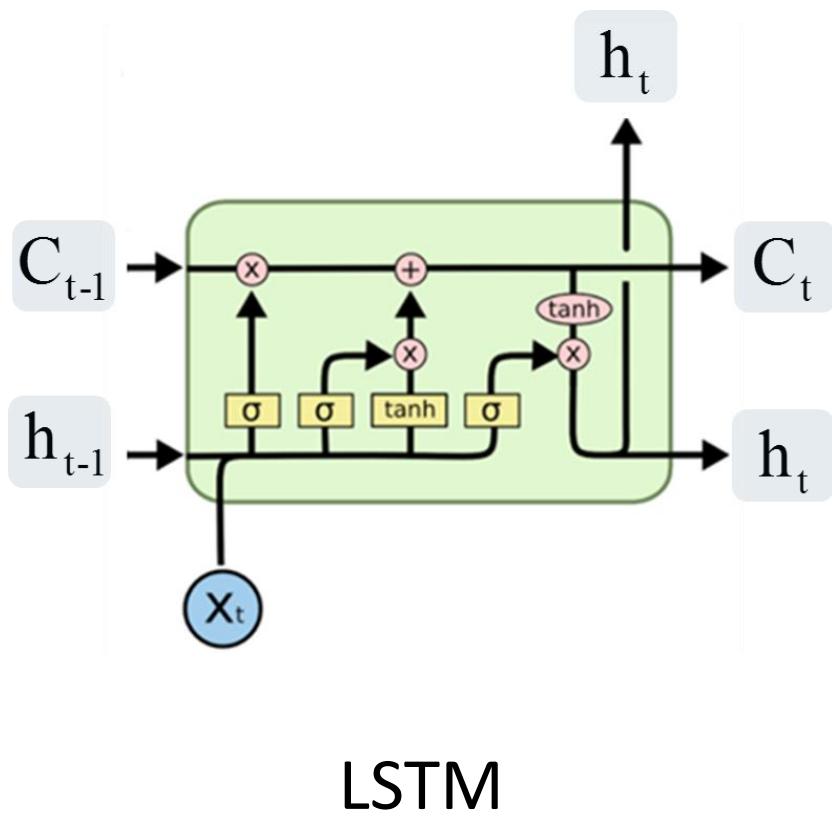


$$i_t = \text{sigm}(w_{xi}x_t + w_{hi}h_{t-1} + b_i)$$

$$\tilde{C}_t = \tanh(w_{xC}x_t + w_{hC}h_{t-1} + b_C)$$

- ❑ *Input gate layer*
 - decide how much “ \tilde{C}_t ” is forgotten?
- ❑ *tanh layer:*
 - *decide which value is updated.*

Long Short Term Memory (LSTM) network: update cell



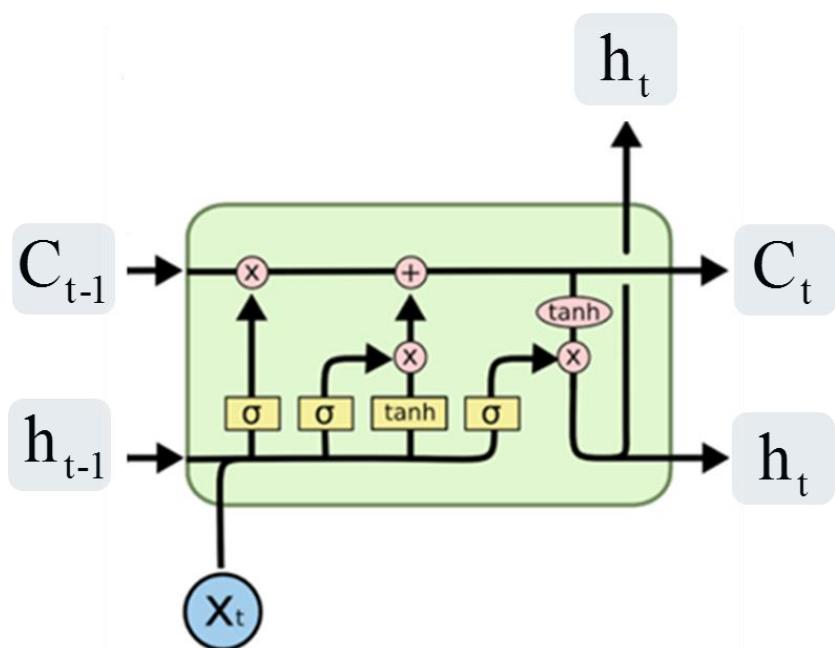
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

past
memory

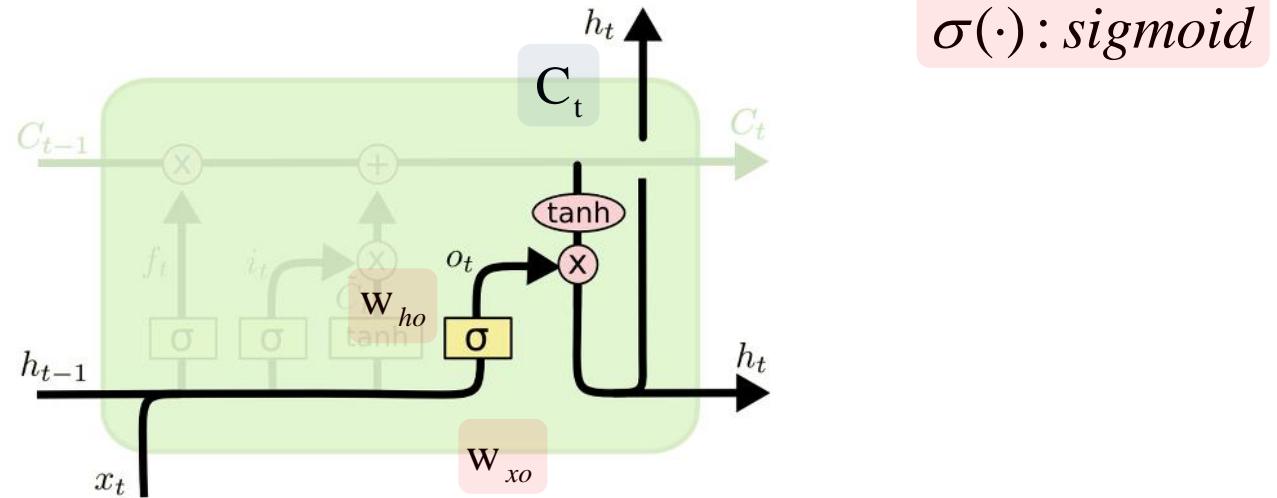
present
memory

- Update the output cell state of “ C_t ” by adding the past cell state of “ C_{t-1} ” to the present cell state of “ \tilde{C}_t ”

Long Short Term Memory (LSTM) network: output gate



LSTM



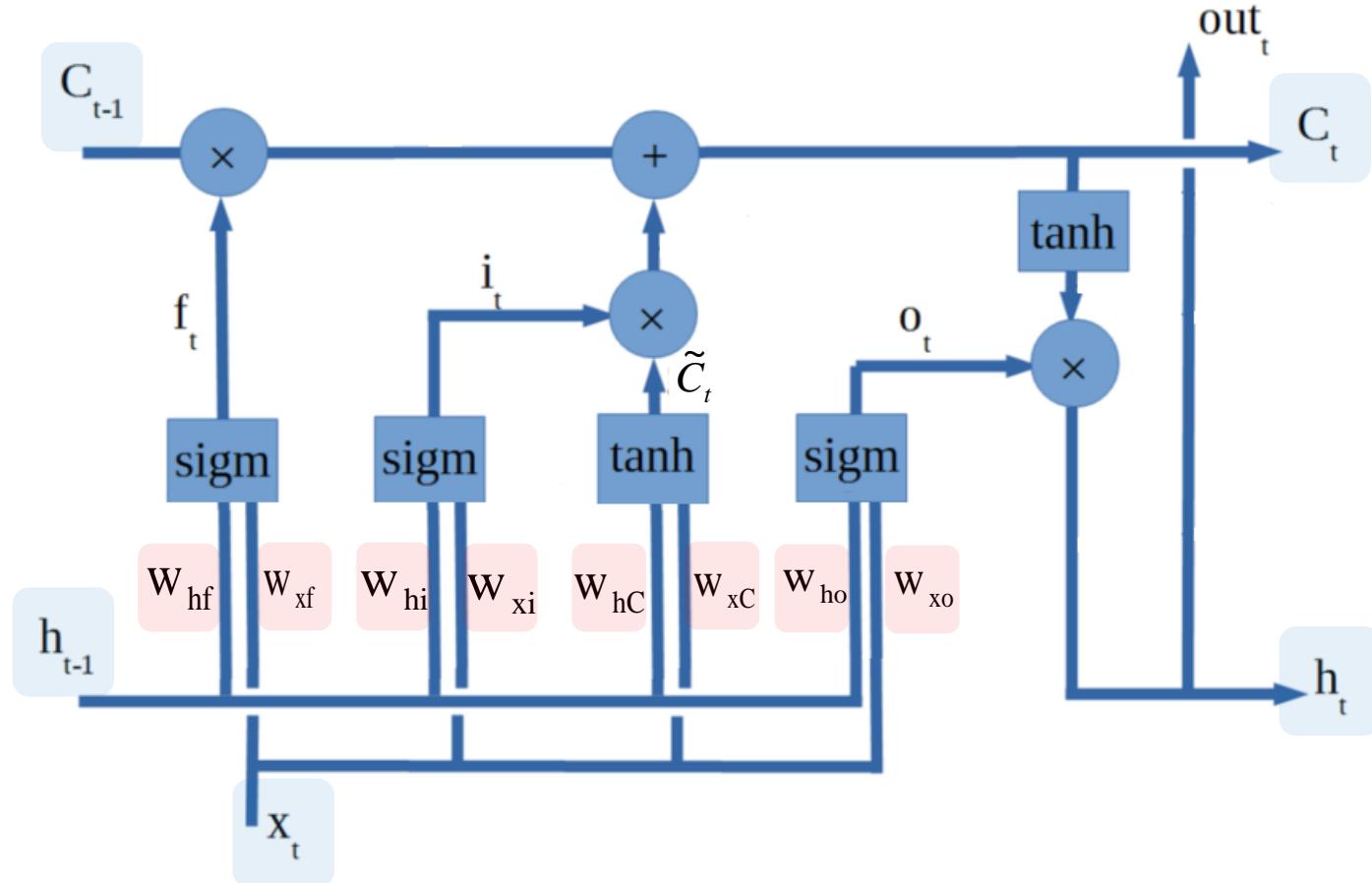
$$o_t = \text{sigm}(\mathbf{w}_{xo} \mathbf{x}_t + \mathbf{w}_{ho} \mathbf{h}_{t-1} + b_o)$$

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

□ Output gate layer

- The output cell state is put through $\tanh()$ and rescaled by the output of the sigmoid function.

Long Short Term Memory (LSTM) network: summary



$$f_t = \text{sigm}(W_{xf}x_t + W_{hf}h_{t-1} + b_f)$$

$$i_t = \text{sigm}(W_{xi}x_t + W_{hi}h_{t-1} + b_i)$$

$$o_t = \text{sigm}(W_{xo}x_t + W_{ho}h_{t-1} + b_o)$$

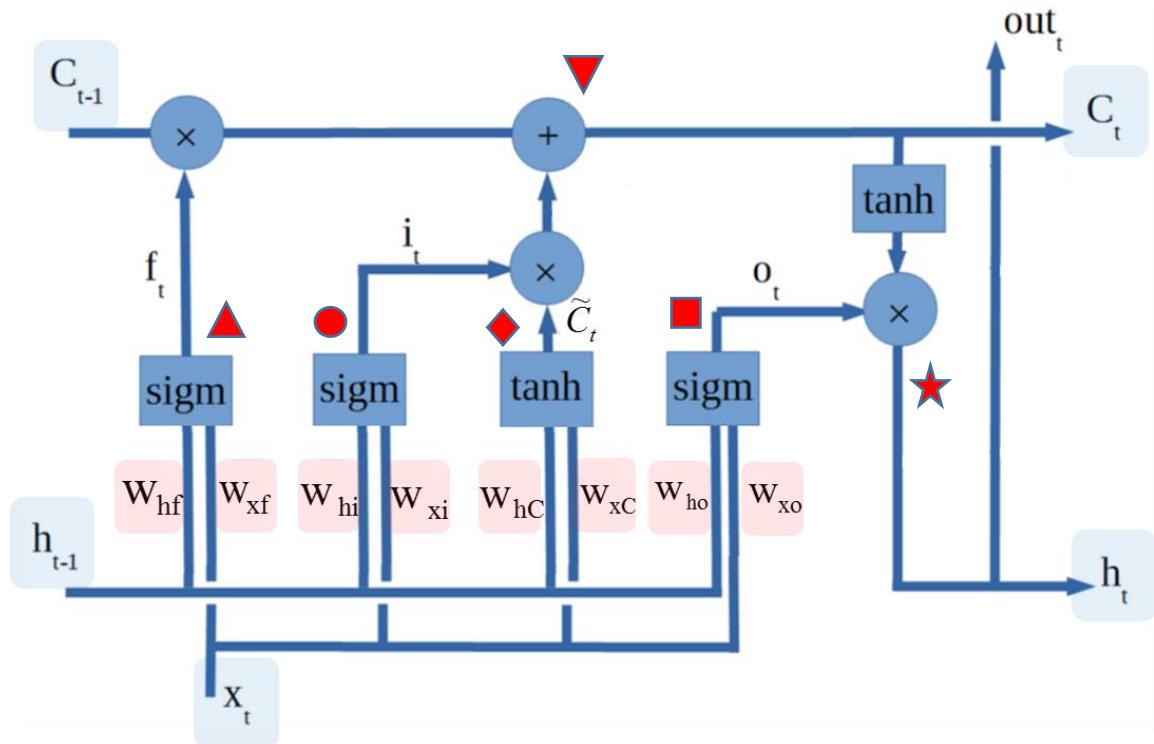
$$\tilde{C}_t = \tanh(W_{xC}x_t + W_{hC}h_{t-1} + b_C)$$

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

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

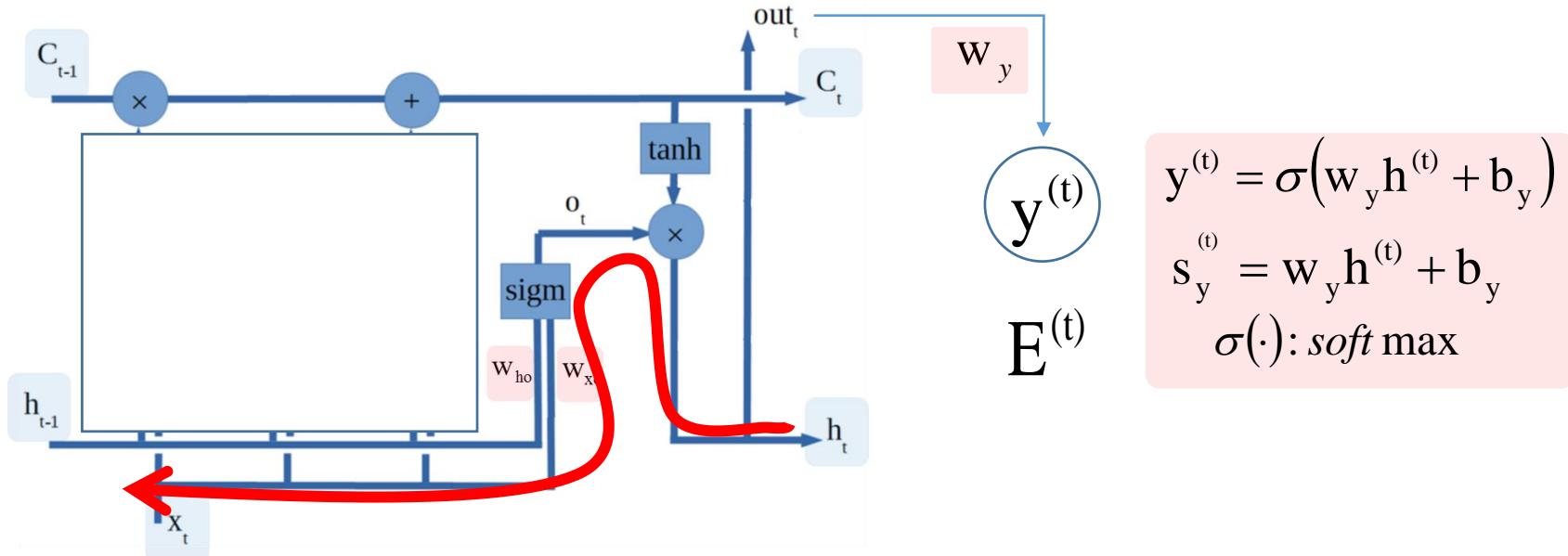
LSTM Backpropagation

Long Short Term Memory (LSTM) network: Backpropagation

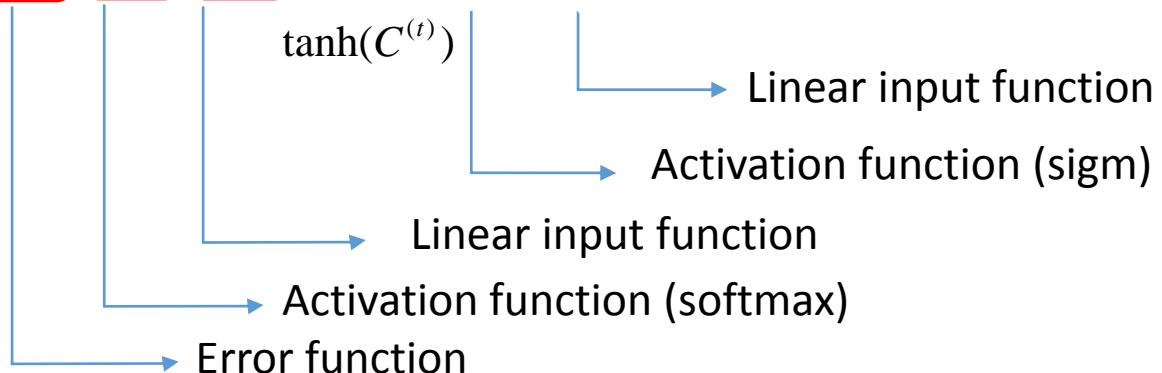


- ▲ $a_f^{(t)} = W_{xf} X^{(t)} + W_{hf} h^{(t-1)}$
 $b_f^{(t)} = \text{sigm}(a_f^{(t)})$
- $a_i^{(t)} = W_{xi} X^{(t)} + W_{hi} h^{(t-1)}$
 $b_i^{(t)} = \text{sigm}(a_i^{(t)})$
- $a_o^{(t)} = W_{xo} X^{(t)} + W_{ho} h^{(t-1)}$
 $b_o^{(t)} = \text{sigm}(a_o^{(t)})$
- ◆ $a_c^{(t)} = W_{xc} X^{(t)} + W_{hc} h^{(t-1)}$
 $b_c^{(t)} = \tanh(a_c^{(t)})$
- ▼ $C^{(t)} = b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)}$
- ★ $h^{(t)} = b_o^{(t)} * \tanh(C^{(t)})$

LSTM network: Backpropagation – output gate (t=1)



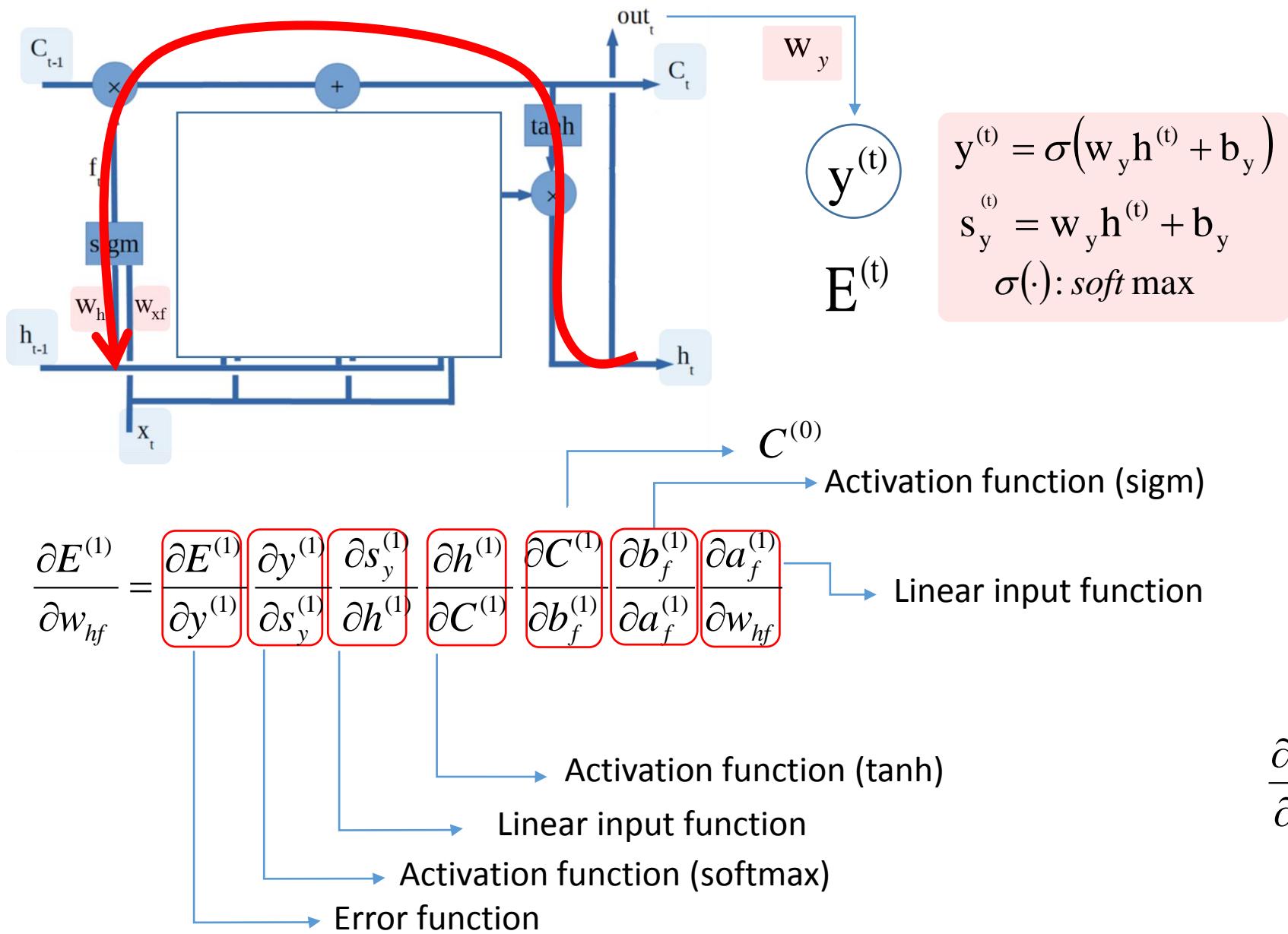
$$\frac{\partial E^{(1)}}{\partial w_{ho}} = \frac{\partial E^{(1)}}{\partial y^{(1)}} \frac{\partial y^{(1)}}{\partial s_y^{(1)}} \frac{\partial s_y^{(1)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial b_o^{(1)}} \frac{\partial b_o^{(1)}}{\partial a_o^{(1)}} \frac{\partial a_o^{(1)}}{\partial w_{ho}}$$



$$\frac{\partial E^{(1)}}{\partial w_{xo}} = ?$$

$$\begin{aligned}
 a_f^{(t)} &= w_{xf} x^{(t)} + w_{hf} h^{(t-1)} \\
 b_f^{(t)} &= \text{sigm}(a_f^{(t)}) \\
 a_i^{(t)} &= w_{xi} x^{(t)} + w_{hi} h^{(t-1)} \\
 b_i^{(t)} &= \text{sigm}(a_i^{(t)}) \\
 a_o^{(t)} &= w_{xo} x^{(t)} + w_{ho} h^{(t-1)} \\
 b_o^{(t)} &= \text{sigm}(a_o^{(t)}) \\
 a_c^{(t)} &= w_{xc} x^{(t)} + w_{hc} h^{(t-1)} \\
 b_c^{(t)} &= \tanh(a_c^{(t)}) \\
 C^{(t)} &= b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)} \\
 h^{(t)} &= b_o^{(t)} * \tanh(C^{(t)})
 \end{aligned}$$

LSTM network: Backpropagation – forget gate (t=1)



$$a_f^{(t)} = w_{xf}x^{(t)} + w_{hf}h^{(t-1)}$$

$$b_f^{(t)} = \text{sigm}(a_f^{(t)})$$

$$a_i^{(t)} = w_{xi}x^{(t)} + w_{hi}h^{(t-1)}$$

$$b_i^{(t)} = \text{sigm}(a_i^{(t)})$$

$$a_o^{(t)} = w_{xo}x^{(t)} + w_{ho}h^{(t-1)}$$

$$b_o^{(t)} = \text{sigm}(a_o^{(t)})$$

$$a_c^{(t)} = w_{xc}x^{(t)} + w_{hc}h^{(t-1)}$$

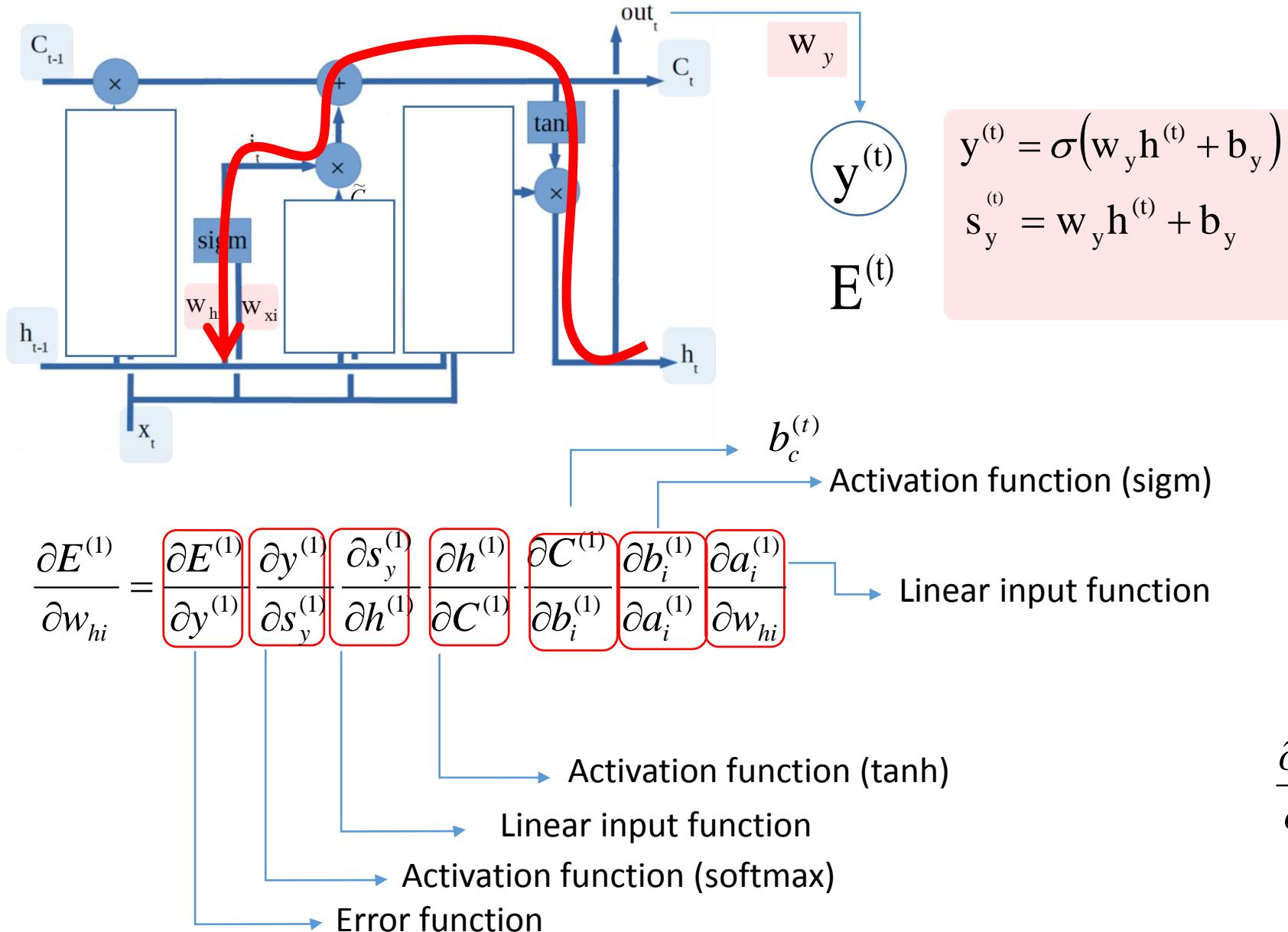
$$b_c^{(t)} = \tanh(a_c^{(t)})$$

$$C^{(t)} = b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)}$$

$$h^{(t)} = b_o^{(t)} * \tanh(C^{(t)})$$

$$\frac{\partial E^{(1)}}{\partial w_{xf}} = ?$$

LSTM network: Backpropagation – input gate (t=1)



$$a_f^{(t)} = w_{xf}x^{(t)} + w_{hf}h^{(t-1)}$$

$$b_f^{(t)} = \text{sigm}(a_f^{(t)})$$

$$a_i^{(t)} = w_{xi}x^{(t)} + w_{hi}h^{(t-1)}$$

$$b_i^{(t)} = \text{sigm}(a_i^{(t)})$$

$$a_o^{(t)} = w_{xo}x^{(t)} + w_{ho}h^{(t-1)}$$

$$b_o^{(t)} = \text{sigm}(a_o^{(t)})$$

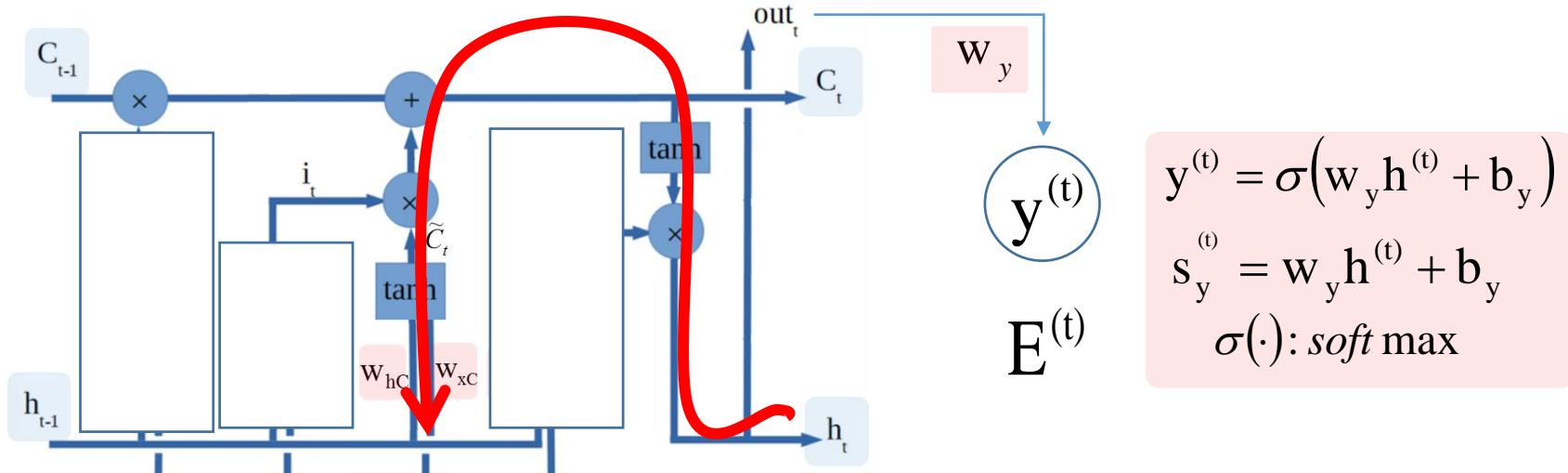
$$a_c^{(t)} = w_{xc}x^{(t)} + w_{hc}h^{(t-1)}$$

$$b_c^{(t)} = \tanh(a_c^{(t)})$$

$$C^{(t)} = b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)}$$

$$h^{(t)} = b_o^{(t)} * \tanh(C^{(t)})$$

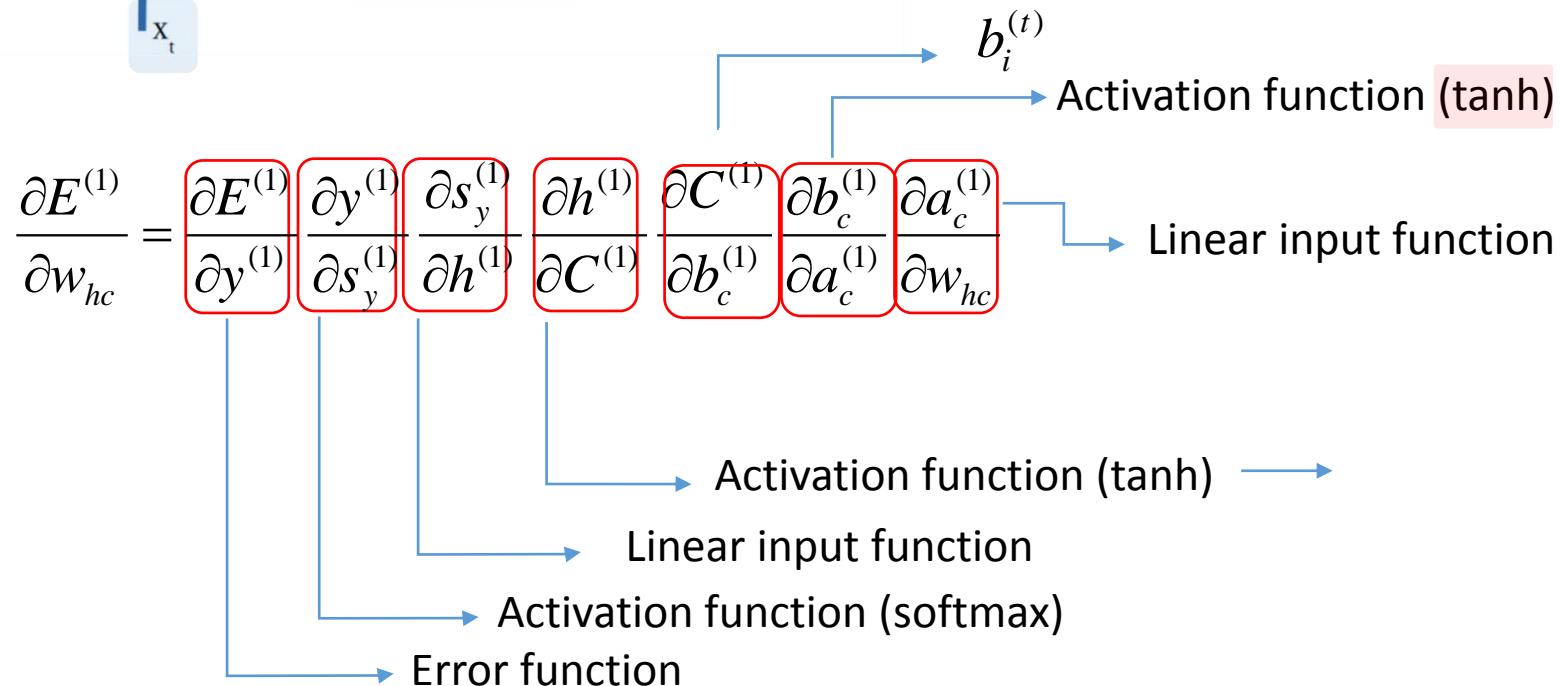
LSTM network: Backpropagation – block input (t=1)



$$y^{(t)} = \sigma(w_y h^{(t)} + b_y)$$

$$s_y^{(t)} = w_y h^{(t)} + b_y$$

$$\sigma(\cdot): \text{soft max}$$



$$\frac{\partial E^{(1)}}{\partial w_{xc}} = ?$$

$$a_f^{(t)} = w_{xf}x^{(t)} + w_{hf}h^{(t-1)}$$

$$b_f^{(t)} = \text{sigm}(a_f^{(t)})$$

$$a_i^{(t)} = w_{xi}x^{(t)} + w_{hi}h^{(t-1)}$$

$$b_i^{(t)} = \text{sigm}(a_i^{(t)})$$

$$a_o^{(t)} = w_{xo}x^{(t)} + w_{ho}h^{(t-1)}$$

$$b_o^{(t)} = \text{sigm}(a_o^{(t)})$$

$$a_c^{(t)} = w_{xc}x^{(t)} + w_{hc}h^{(t-1)}$$

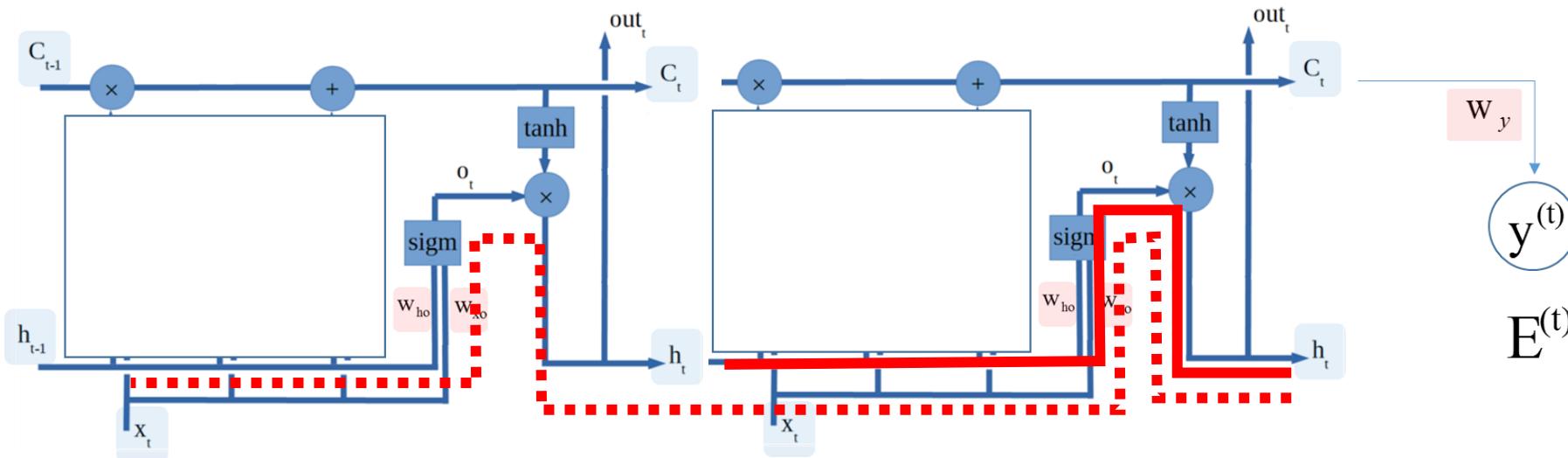
$$b_c^{(t)} = \tanh(a_c^{(t)})$$

$$C^{(t)} = b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)}$$

$$h^{(t)} = b_o^{(t)} * \tanh(C^{(t)})$$

LSTM Backpropagation (t=2)

LSTM network: Backpropagation – output gate (t=2)



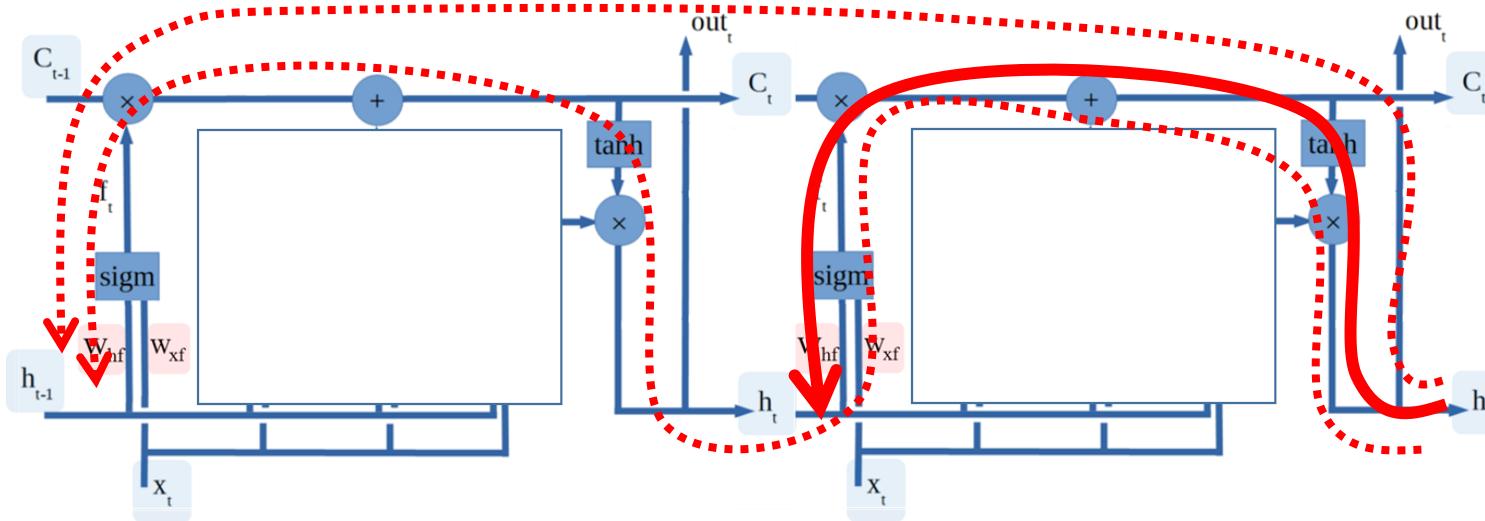
$$\begin{aligned}
 y^{(t)} &= \sigma(w_y h^{(t)} + b_y) \\
 s_y^{(t)} &= w_y h^{(t)} + b_y \\
 \sigma(\cdot) &: \text{soft max}
 \end{aligned}$$

$$\begin{aligned}
 \sum_{t=1}^2 \frac{\partial E^{(t)}}{\partial w_{ho}} &= \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial b_o^{(2)}} \frac{\partial b_o^{(2)}}{\partial a_o^{(2)}} \frac{\partial a_o^{(2)}}{\partial w_{ho}} \\
 &+ \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial b_o^{(2)}} \frac{\partial b_o^{(2)}}{\partial a_o^{(2)}} \frac{\partial a_o^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial b_o^{(1)}} \frac{\partial b_o^{(1)}}{\partial a_o^{(1)}} \frac{\partial a_o^{(1)}}{\partial w_{ho}}
 \end{aligned}$$

Activation function (sigm) Activation function (sigm)

$a_f^{(t)} = w_{xf}x^{(t)} + w_{hf}h^{(t-1)}$
$b_f^{(t)} = \text{sigm}(a_f^{(t)})$
$a_i^{(t)} = w_{xi}x^{(t)} + w_{hi}h^{(t-1)}$
$b_i^{(t)} = \text{sigm}(a_i^{(t)})$
$a_o^{(t)} = w_{xo}x^{(t)} + w_{ho}h^{(t-1)}$
$b_o^{(t)} = \text{sigm}(a_o^{(t)})$
$a_c^{(t)} = w_{xc}x^{(t)} + w_{hc}h^{(t-1)}$
$b_c^{(t)} = \tanh(a_c^{(t)})$
$C^{(t)} = b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)}$
$h^{(t)} = b_o^{(t)} * \tanh(C^{(t)})$

LSTM network: Backpropagation – forget gate (t=2)



$$\begin{aligned} \sum_{t=1}^2 \frac{\partial E^{(t)}}{\partial w_{hf}} &= \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \frac{\partial C^{(2)}}{\partial b_f^{(2)}} \frac{\partial b_f^{(2)}}{\partial a_f^{(2)}} \frac{\partial a_f^{(2)}}{\partial w_{hf}} \\ &\quad + \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \frac{\partial C^{(2)}}{\partial b_f^{(2)}} \frac{\partial b_f^{(2)}}{\partial a_f^{(2)}} \frac{\partial a_f^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial C^{(1)}} \frac{\partial C^{(1)}}{\partial b_f^{(1)}} \frac{\partial b_f^{(1)}}{\partial a_f^{(1)}} \frac{\partial a_f^{(1)}}{\partial w_{hf}} \\ &\quad + \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \frac{\partial C^{(2)}}{\partial b_f^{(1)}} \frac{\partial b_f^{(1)}}{\partial a_f^{(1)}} \frac{\partial a_f^{(1)}}{\partial w_{hf}} \end{aligned}$$

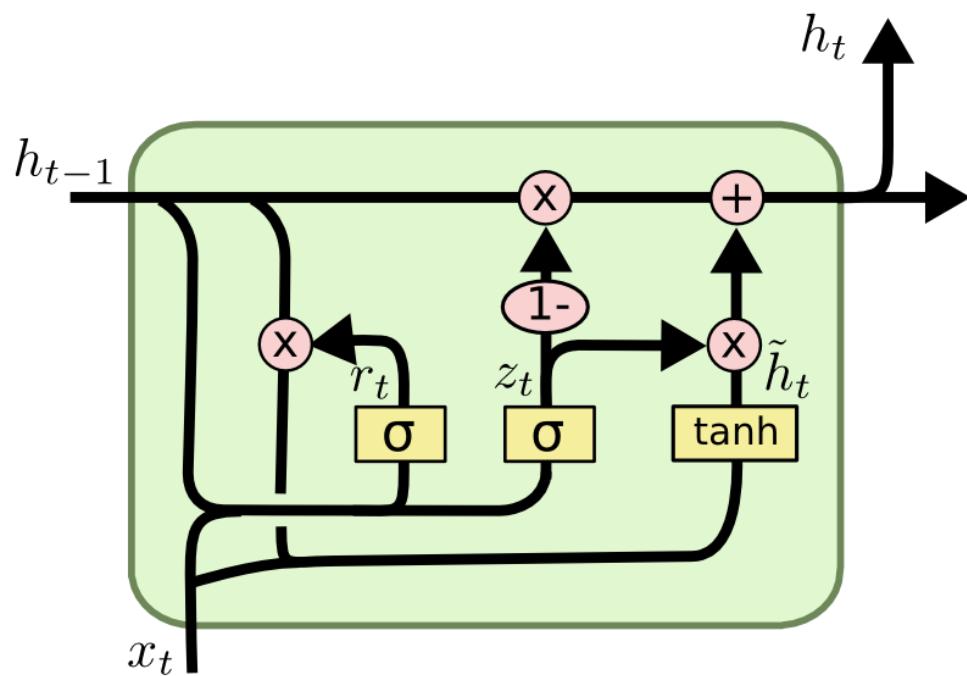
$a_f^{(t)} = w_{xf}x^{(t)} + w_{hf}h^{(t-1)}$
$b_f^{(t)} = \text{sigm}(a_f^{(t)})$
$a_i^{(t)} = w_{xi}x^{(t)} + w_{hi}h^{(t-1)}$
$b_i^{(t)} = \text{sigm}(a_i^{(t)})$
$a_o^{(t)} = w_{xo}x^{(t)} + w_{ho}h^{(t-1)}$
$b_o^{(t)} = \text{sigm}(a_o^{(t)})$
$a_c^{(t)} = w_{xc}x^{(t)} + w_{hc}h^{(t-1)}$
$b_c^{(t)} = \tanh(a_c^{(t)})$
$C^{(t)} = b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)}$
$h^{(t)} = b_o^{(t)} * \tanh(C^{(t)})$

LSTM network: Backpropagation – input gate and block gate (t=2)

- ❑ Backpropagations of input gate and block gate are similar to that of forget gate in the previous slide.
- ❑ Please, refer to the backup slides.

Gated Recurrent Unit (GRU)

- ❑ Simpler than LSTM and so training is faster,
- ❑ Cell state (C) is replaced by hidden state (h),
- ❑ GRU has two gates: update gate (z), reset gate (r).



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

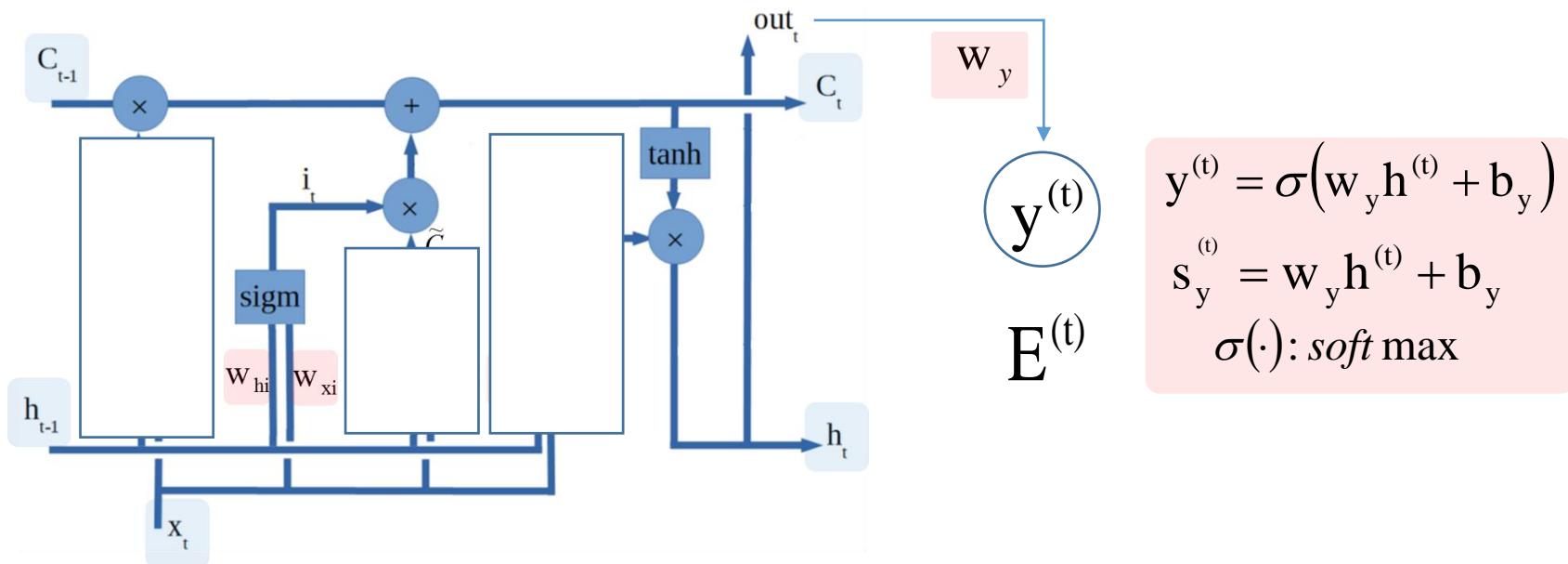
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

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

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

Backup Slides

LSTM network: Backpropagation – input gate (t=2)



$$y^{(t)} = \sigma(w_y h^{(t)} + b_y)$$

$$s_y^{(t)} = w_y h^{(t)} + b_y$$

$$\sigma(\cdot): \text{soft max}$$

$$\frac{\partial E^{(2)}}{\partial w_{hi}} = \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \frac{\partial C^{(2)}}{\partial b_i^{(2)}} \frac{\partial b_i^{(2)}}{\partial a_i^{(2)}} \frac{\partial a_i^{(2)}}{\partial w_{hi}}$$

$$+ \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \frac{\partial C^{(2)}}{\partial b_i^{(2)}} \frac{\partial b_i^{(2)}}{\partial a_i^{(2)}} \boxed{\frac{\partial h^{(1)}}{\partial C^{(1)}}} \boxed{\frac{\partial C^{(1)}}{\partial b_i^{(1)}}} \boxed{\frac{\partial b_i^{(1)}}{\partial a_i^{(1)}}} \boxed{\frac{\partial a_i^{(1)}}{\partial w_{hi}}}$$

$$+ \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \boxed{\frac{\partial C^{(2)}}{\partial C^{(1)}}} \boxed{\frac{\partial C^{(1)}}{\partial b_i^{(1)}}} \boxed{\frac{\partial b_i^{(1)}}{\partial a_i^{(1)}}} \boxed{\frac{\partial a_i^{(1)}}{\partial w_{hi}}}$$

$$a_f^{(t)} = w_{xf} x^{(t)} + w_{hf} h^{(t-1)}$$

$$b_f^{(t)} = \text{sigm}(a_f^{(t)})$$

$$a_i^{(t)} = w_{xi} x^{(t)} + w_{hi} h^{(t-1)}$$

$$b_i^{(t)} = \text{sigm}(a_i^{(t)})$$

$$a_o^{(t)} = w_{xo} x^{(t)} + w_{ho} h^{(t-1)}$$

$$b_o^{(t)} = \text{sigm}(a_o^{(t)})$$

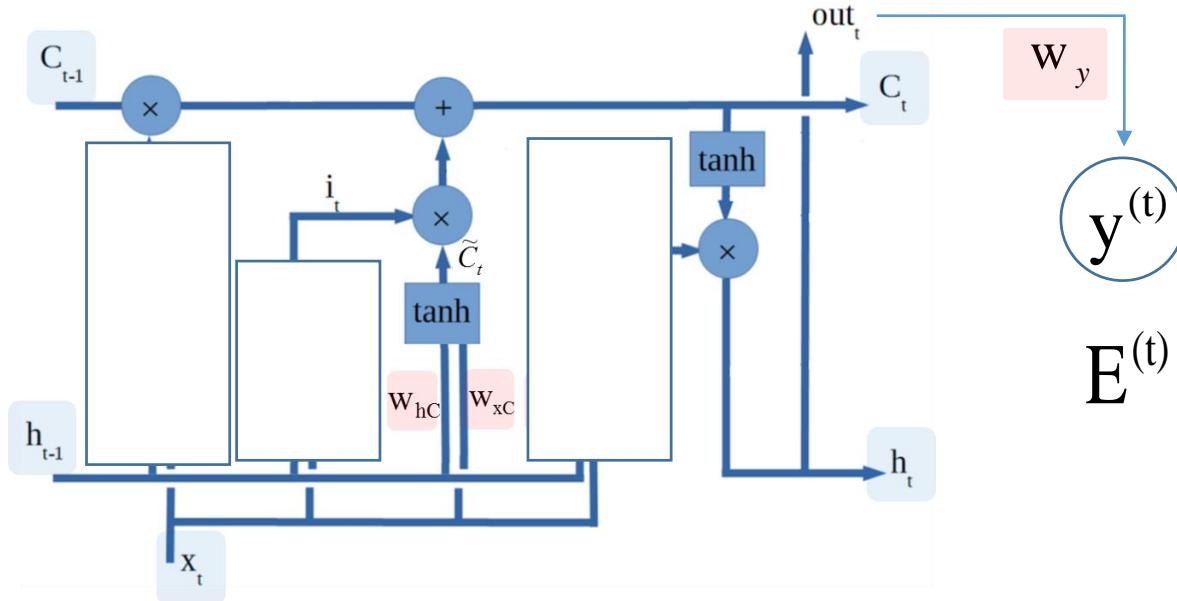
$$a_c^{(t)} = w_{xc} x^{(t)} + w_{hc} h^{(t-1)}$$

$$b_c^{(t)} = \tanh(a_c^{(t)})$$

$$C^{(t)} = b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)}$$

$$h^{(t)} = b_o^{(t)} * \tanh(C^{(t)})$$

LSTM network: Backpropagation – block input (t=2)



$$y^{(t)} = \sigma(w_y h^{(t)} + b_y)$$

$$s_y^{(t)} = w_y h^{(t)} + b_y$$

$$\sigma(\cdot): \text{soft max}$$

$$\frac{\partial E^{(2)}}{\partial w_{hc}} = \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \frac{\partial C^{(2)}}{\partial b_c^{(2)}} \frac{\partial b_c^{(2)}}{\partial a_c^{(2)}} \frac{\partial a_c^{(2)}}{\partial w_{hc}}$$

$$+ \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \frac{\partial C^{(2)}}{\partial b_c^{(2)}} \frac{\partial b_c^{(2)}}{\partial a_c^{(2)}} \frac{\partial a_c^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial C^{(1)}} \frac{\partial C^{(1)}}{\partial b_c^{(1)}} \frac{\partial b_c^{(1)}}{\partial a_c^{(1)}} \frac{\partial a_c^{(1)}}{\partial w_{hc}}$$

$$+ \frac{\partial E^{(2)}}{\partial y^{(2)}} \frac{\partial y^{(2)}}{\partial s_y^{(2)}} \frac{\partial s_y^{(2)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial C^{(2)}} \frac{\partial C^{(2)}}{\partial C^{(1)}} \frac{\partial C^{(1)}}{\partial b_c^{(1)}} \frac{\partial b_c^{(1)}}{\partial a_c^{(1)}} \frac{\partial a_c^{(1)}}{\partial w_{hc}}$$

$$a_f^{(t)} = w_{xf} x^{(t)} + w_{hf} h^{(t-1)}$$

$$b_f^{(t)} = \text{sigm}(a_f^{(t)})$$

$$a_i^{(t)} = w_{xi} x^{(t)} + w_{hi} h^{(t-1)}$$

$$b_i^{(t)} = \text{sigm}(a_i^{(t)})$$

$$a_o^{(t)} = w_{xo} x^{(t)} + w_{ho} h^{(t-1)}$$

$$b_o^{(t)} = \text{sigm}(a_o^{(t)})$$

$$a_c^{(t)} = w_{xc} x^{(t)} + w_{hc} h^{(t-1)}$$

$$b_c^{(t)} = \tanh(a_c^{(t)})$$

$$C^{(t)} = b_f^{(t)} * C^{(t-1)} + b_i^{(t)} * b_c^{(t)}$$

$$h^{(t)} = b_o^{(t)} * \tanh(C^{(t)})$$

Operation of LSTM: weight update

$$\begin{bmatrix} \frac{\partial E}{\partial w_{xo}} \\ \frac{\partial E}{\partial w_{xf}} \\ \frac{\partial E}{\partial w_{xi}} \\ \frac{\partial E}{\partial w_{xc}} \end{bmatrix} = \begin{bmatrix} \frac{\partial E^{(1)}}{\partial w_{xo}} + \frac{\partial E^{(2)}}{\partial w_{xo}} + \cdots + \frac{\partial E^{(t)}}{\partial w_{xo}} \\ \frac{\partial E^{(1)}}{\partial w_{xf}} + \frac{\partial E^{(2)}}{\partial w_{xf}} + \cdots + \frac{\partial E^{(t)}}{\partial w_{xf}} \\ \frac{\partial E^{(1)}}{\partial w_{xi}} + \frac{\partial E^{(2)}}{\partial w_{xi}} + \cdots + \frac{\partial E^{(t)}}{\partial w_{xi}} \\ \frac{\partial E^{(1)}}{\partial w_{xc}} + \frac{\partial E^{(2)}}{\partial w_{xc}} + \cdots + \frac{\partial E^{(t)}}{\partial w_{xc}} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial E}{\partial w_{ho}} \\ \frac{\partial E}{\partial w_{hf}} \\ \frac{\partial E}{\partial w_{hi}} \\ \frac{\partial E}{\partial w_{hc}} \end{bmatrix} = \begin{bmatrix} \frac{\partial E^{(1)}}{\partial w_{ho}} + \frac{\partial E^{(2)}}{\partial w_{ho}} + \cdots + \frac{\partial E^{(t)}}{\partial w_{ho}} \\ \frac{\partial E^{(1)}}{\partial w_{hf}} + \frac{\partial E^{(2)}}{\partial w_{hf}} + \cdots + \frac{\partial E^{(t)}}{\partial w_{hf}} \\ \frac{\partial E^{(1)}}{\partial w_{hi}} + \frac{\partial E^{(2)}}{\partial w_{hi}} + \cdots + \frac{\partial E^{(t)}}{\partial w_{hi}} \\ \frac{\partial E^{(1)}}{\partial w_{hc}} + \frac{\partial E^{(2)}}{\partial w_{hc}} + \cdots + \frac{\partial E^{(t)}}{\partial w_{hc}} \end{bmatrix}$$

where:

$$\begin{array}{lll} h^{(t-1)}, h^{(t)} \in R^n & w_{xo} \in R^{n \times m} & w_{ho} \in R^{n \times n} \\ C^{(t)} \in R^n & w_{xf} \in R^{n \times m} & w_{hf} \in R^{n \times n} \\ x^{(t)} \in R^m & w_{xi} \in R^{n \times m} & w_{hi} \in R^{n \times n} \\ w_{xc} \in R^{n \times m} & & w_{hc} \in R^{n \times n} \end{array}$$

$$\begin{bmatrix} w_{xo}^{new} \\ w_{xf}^{new} \\ w_{xi}^{new} \\ w_{xc}^{new} \end{bmatrix} = \begin{bmatrix} w_{xo}^{old} - \lambda \frac{\partial E}{\partial w_{xo}} \\ w_{xf}^{old} - \lambda \frac{\partial E}{\partial w_{xf}} \\ w_{xi}^{old} - \lambda \frac{\partial E}{\partial w_{xi}} \\ w_{xc}^{old} - \lambda \frac{\partial E}{\partial w_{xc}} \end{bmatrix}$$

$$\begin{bmatrix} w_{ho}^{new} \\ w_{hf}^{new} \\ w_{hi}^{new} \\ w_{hc}^{new} \end{bmatrix} = \begin{bmatrix} w_{ho}^{old} - \lambda \frac{\partial E}{\partial w_{ho}} \\ w_{hf}^{old} - \lambda \frac{\partial E}{\partial w_{hf}} \\ w_{hi}^{old} - \lambda \frac{\partial E}{\partial w_{hi}} \\ w_{hc}^{old} - \lambda \frac{\partial E}{\partial w_{hc}} \end{bmatrix}$$

λ : Learning rate