# Recurrent Neural Networks

Directed Reading Course (Part Two)

Arman Afrasiyabi
Electrical Engineering and Computer Engineering Department
Université Laval

UNIVERSITÉ LAVAL

# Content

- **Motivation**

- **Vanilla Recurrent Neural Network (RNNs)**

- **Bidirectional RNN (Bi-RNN)**

- **Problems of Vanilla RNNs**

- **Gated RNNs**

    - Gated Recurrent Units (GRU)

    - Long Short-Term Memory(LSTM)

- **Examples:**

    - Encoder-Decoder RNNs

    - Attention based RNNs

# Motivation

The problem of language model

$$P(\mathrm{w}_1, \ldots, \mathrm{w}_m)$$

Traditional methods

$$p(\mathrm{w}_2 \mid \mathrm{w}_1) = \frac{\mathrm{count}(\mathrm{w}_1, \mathrm{w}_2)}{\mathrm{count}(\mathrm{w}_1)}$$

$$p(\mathrm{w}_3 \mid \mathrm{w}_1, \mathrm{w}_2) = \frac{\mathrm{count}(\mathrm{w}_1, \mathrm{w}_2, \mathrm{w}_3)}{\mathrm{count}(\mathrm{w}_1, \mathrm{w}_2)}$$

$\vdots$

# Motivation

Markov assumption

$$P(\mathrm{w}_1, ..., \mathrm{w}_m) = \prod_{i=1}^{m} P(\mathrm{w}_i \mid \mathrm{w}_1, ..., \mathrm{w}_{i-1})$$

$$\approx \prod_{i=1}^{m} P(\mathrm{w}_i \mid \mathrm{w}_{i-(n-1)}, ..., \mathrm{w}_{i-1})$$

# Motivation

Markov assumption

$$P(\mathrm{w}_1,...,\mathrm{w}_m) = \prod_{i=1}^{m} P(\mathrm{w}_i \mid \mathrm{w}_1,...,\mathrm{w}_{i-1})$$

$$\approx \prod_{i=1}^{m} P(\mathrm{w}_i \mid \mathrm{w}_{i-(n-1)},...,\mathrm{w}_{i-1})$$

Necessary...

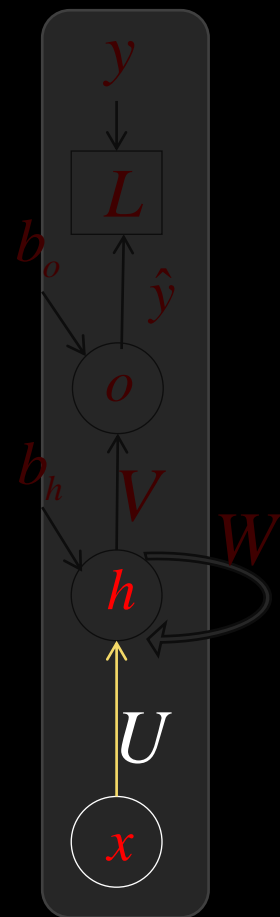Incorrect in many cases!

# Recurrent Neural Networks (RNNs)

RNNs are a type of neural networks with the following goal:

$$\hat{P}\left( x_{t+1} = \hat{y} \mid x_t, ..., x_1 \right)$$
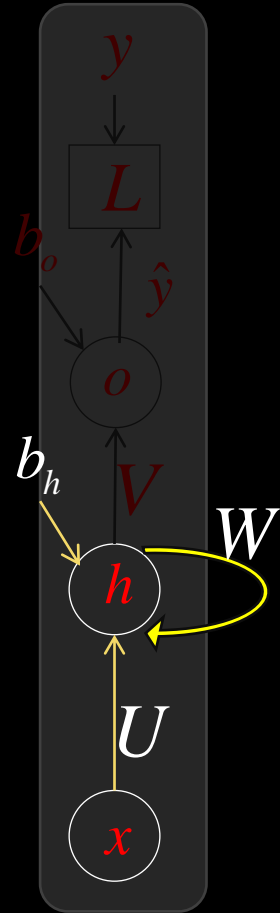
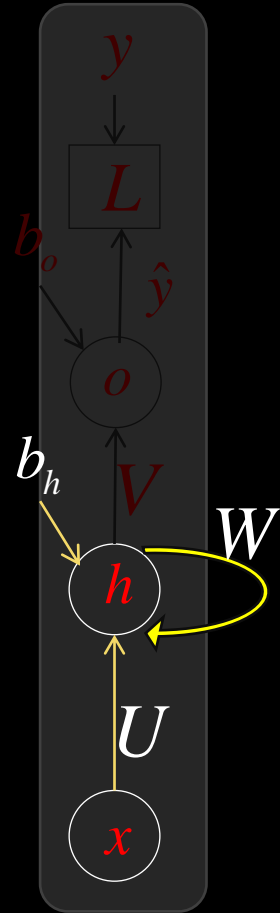# Vanilla Recurrent Neural Networks

# Forward Pass



$$Ux^{(t)}$$

# Forward Pass



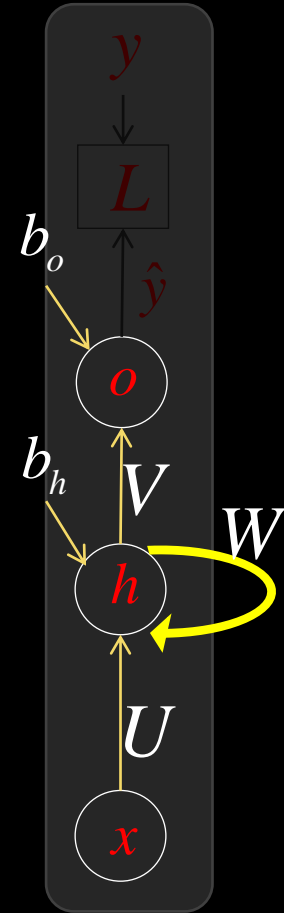$$h^{(t)} = \varphi\left(Ux^{(t)} + Wh^{(t-1)} + b_h\right)$$

# Forward Pass



$$h^{(t)} = \varphi\left(Ux^{(t)} + Wh^{(t-1)} + b_h\right)$$

$$(n \times 1) = \varphi\left((n \times m)(m \times 1) + (n \times n)(n \times 1) + (n \times 1)\right)$$
$$= \varphi\left((n \times 1) + (n \times 1) + (n \times 1)\right)$$
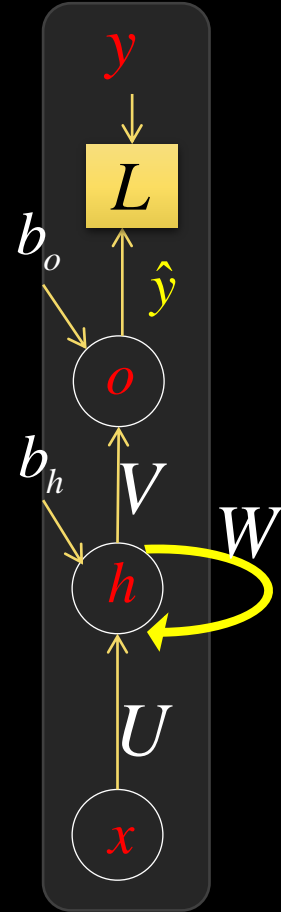$$= \varphi\left(n \times 1\right)$$

# Forward Pass



$$h^{(t)} = \varphi\left(Ux^{(t)} + Wh^{(t-1)} + b_h\right)$$

$$o^{(t)} = Vh^{(t)} + b_o$$

# Forward Pass

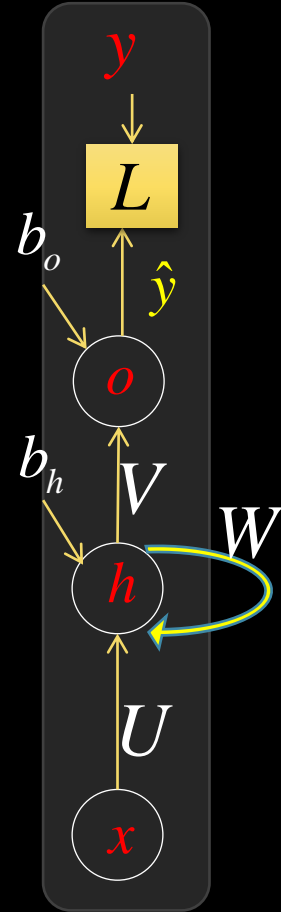$$h^{(t)} = \varphi\left(Ux^{(t)} + Wh^{(t-1)} + b_h\right)$$

$$o^{(t)} = Vh^{(t)} + b_o$$

Example:

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$
$$L = \log-\text{likelihood}$$

# Forward Pass



$$h^{(t)} = \varphi\left(Ux^{(t)} + Wh^{(t-1)} + b_h\right)$$

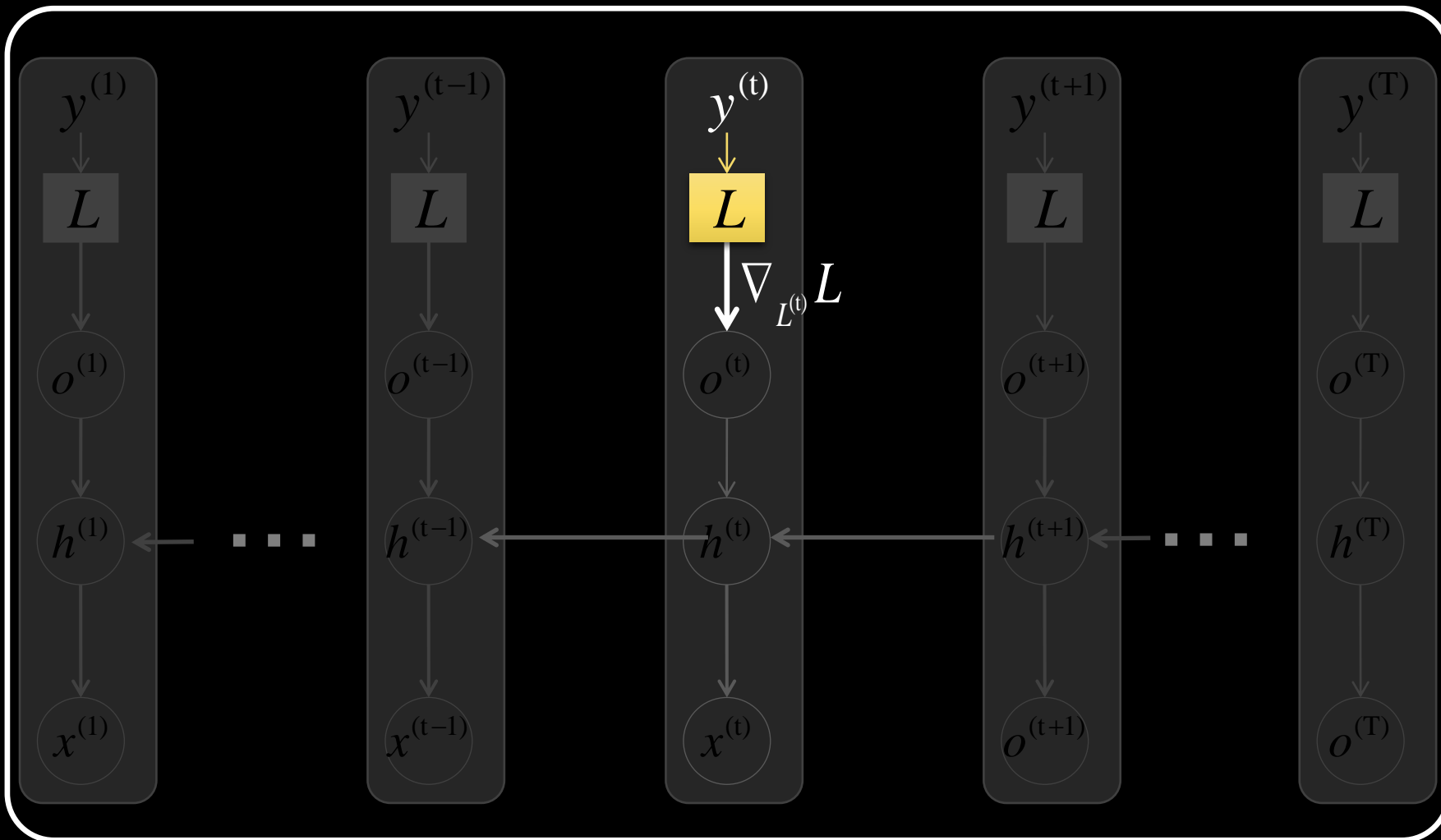$$o^{(t)} = Vh^{(t)} + b_o$$

Example:

$$\hat{y}^{(t)} = \text{softmax}(o_t)$$
$$L = \log-\text{likelihood}$$

In some cases $\hat{y}^{(t)}$ is optional ☺
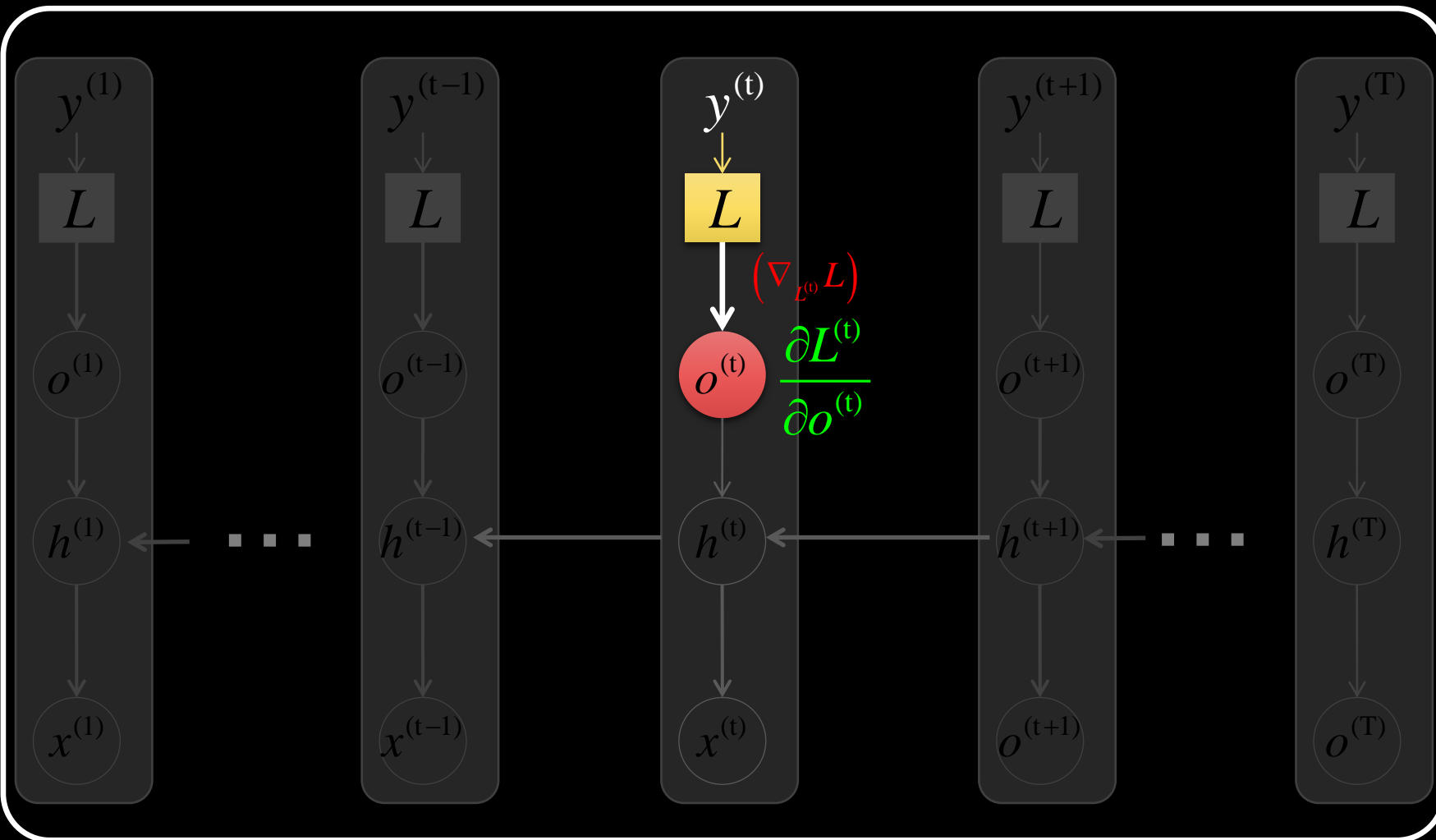
# Back-Propagation Through Time (BPTT)

$$\nabla_{L^{(t)}} L = \frac{\partial L}{\partial L^{(t)}}$$

# Back-Propagation Through Time (BPTT)

$$\nabla_{L^{(t)}} L = \frac{\partial L}{\partial L^{(t)}}$$

$$\nabla_{o^{(t)}} L = \left( \nabla_{L^{(t)}} L \right) \frac{\partial L^{(t)}}{\partial o^{(t)}}$$
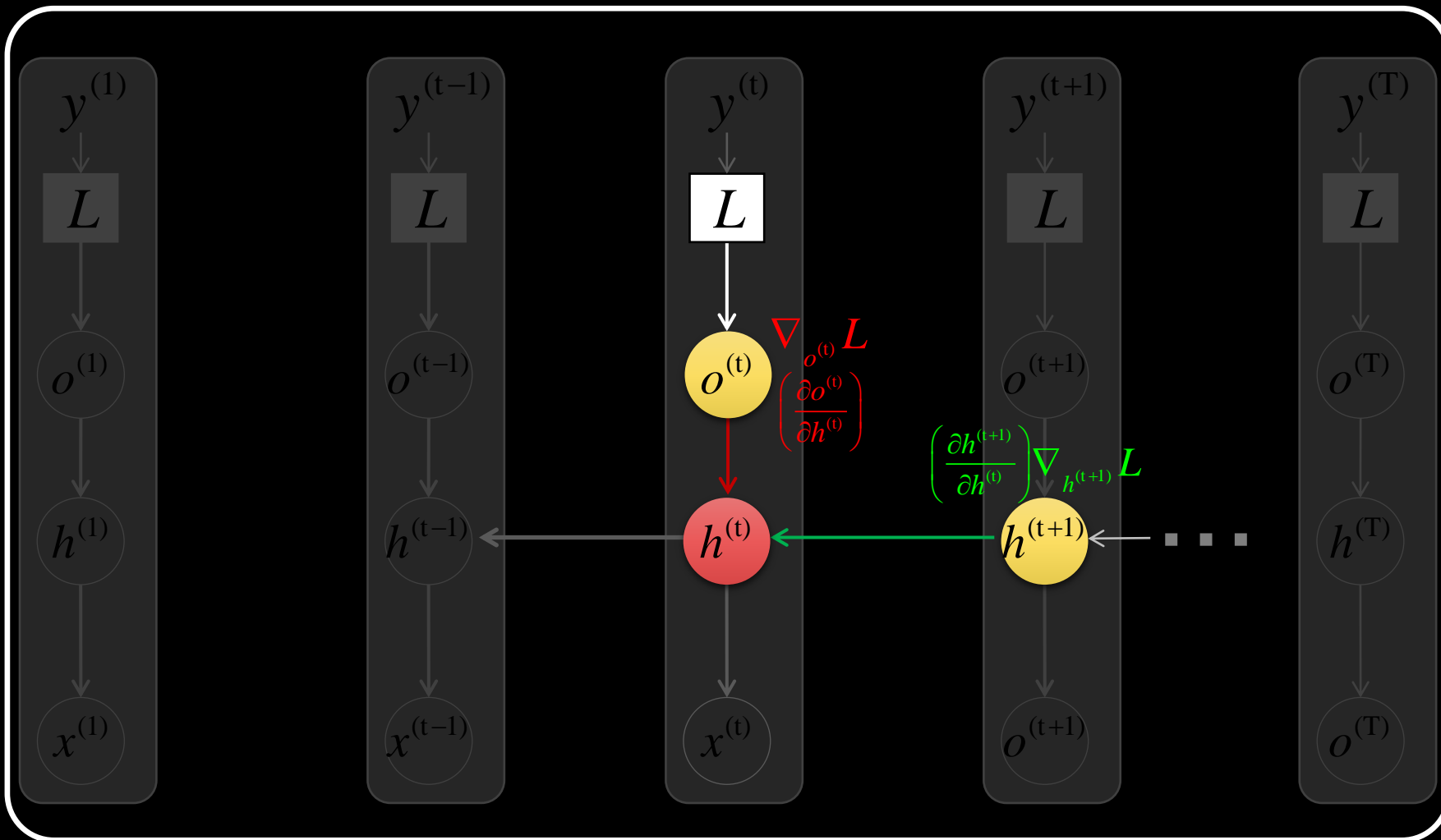
# Back-Propagation Through Time (BPTT)

$$\nabla_{L^{(t)}} L = \frac{\partial L}{\partial L^{(t)}}$$

$$\nabla_{o^{(t)}} L = \left( \nabla_{L^{(t)}} L \right) \frac{\partial L^{(t)}}{\partial o_i^{(t)}}$$

$$\nabla_{h^{(t)}} L = \left( \frac{\partial o^{(t)}}{\partial h^{(t)}} \right) \nabla_{o^{(t)}} L + \left( \frac{\partial h^{(t+1)}}{\partial h^{(t)}} \right) \nabla_{h^{(t+1)}} L$$

# Back-Propagation Through Time (BPTT)

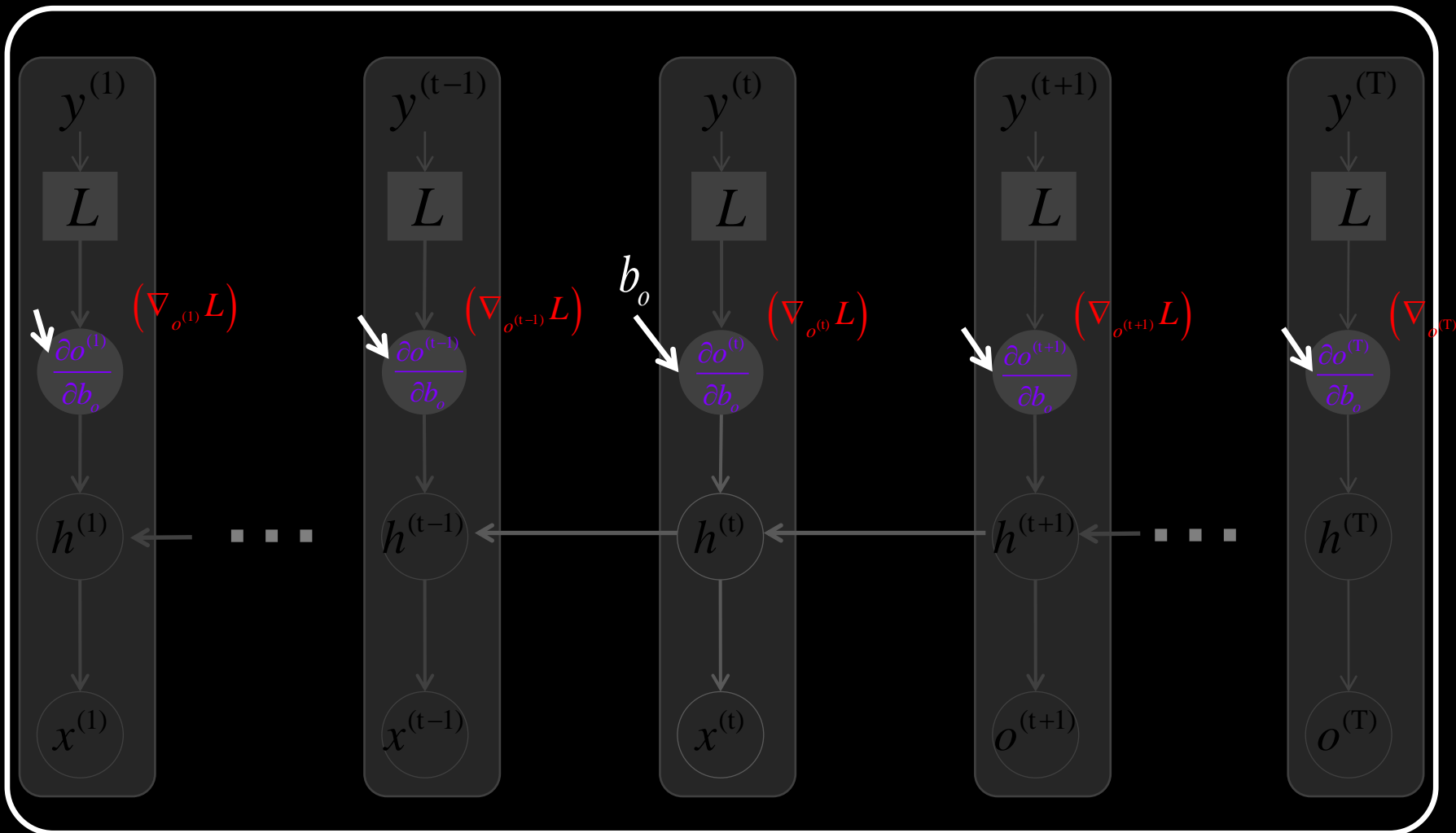$$\nabla_{L^{(t)}} L = \frac{\partial L}{\partial L^{(t)}}$$

$$\nabla_{o^{(t)}} L = \left(\nabla_{L^{(t)}} L\right) \frac{\partial L^{(t)}}{\partial o_i^{(t)}}$$

$$\nabla_{h^{(t)}} L = \left(\frac{\partial h^{(t+1)}}{\partial h^{(t)}}\right) \nabla_{h^{(t+1)}} L + \left(\frac{\partial o^{(t)}}{\partial h^{(t)}}\right) \nabla_{o^{(t)}} L$$

$$\nabla_{b_o} L = \sum_{t=1}^{T} \left(\frac{\partial o^{(t)}}{\partial b_o}\right) \left(\nabla_{o^{(t)}} L\right)$$

biases

Weights

# Back-Propagation Through Time (BPTT)



$$\nabla_{L^{(t)}} L = \frac{\partial L}{\partial L^{(t)}}$$

$$\nabla_{o^{(t)}} L = \left( \nabla_{L^{(t)}} L \right) \frac{\partial L^{(t)}}{\partial o_i^{(t)}}$$
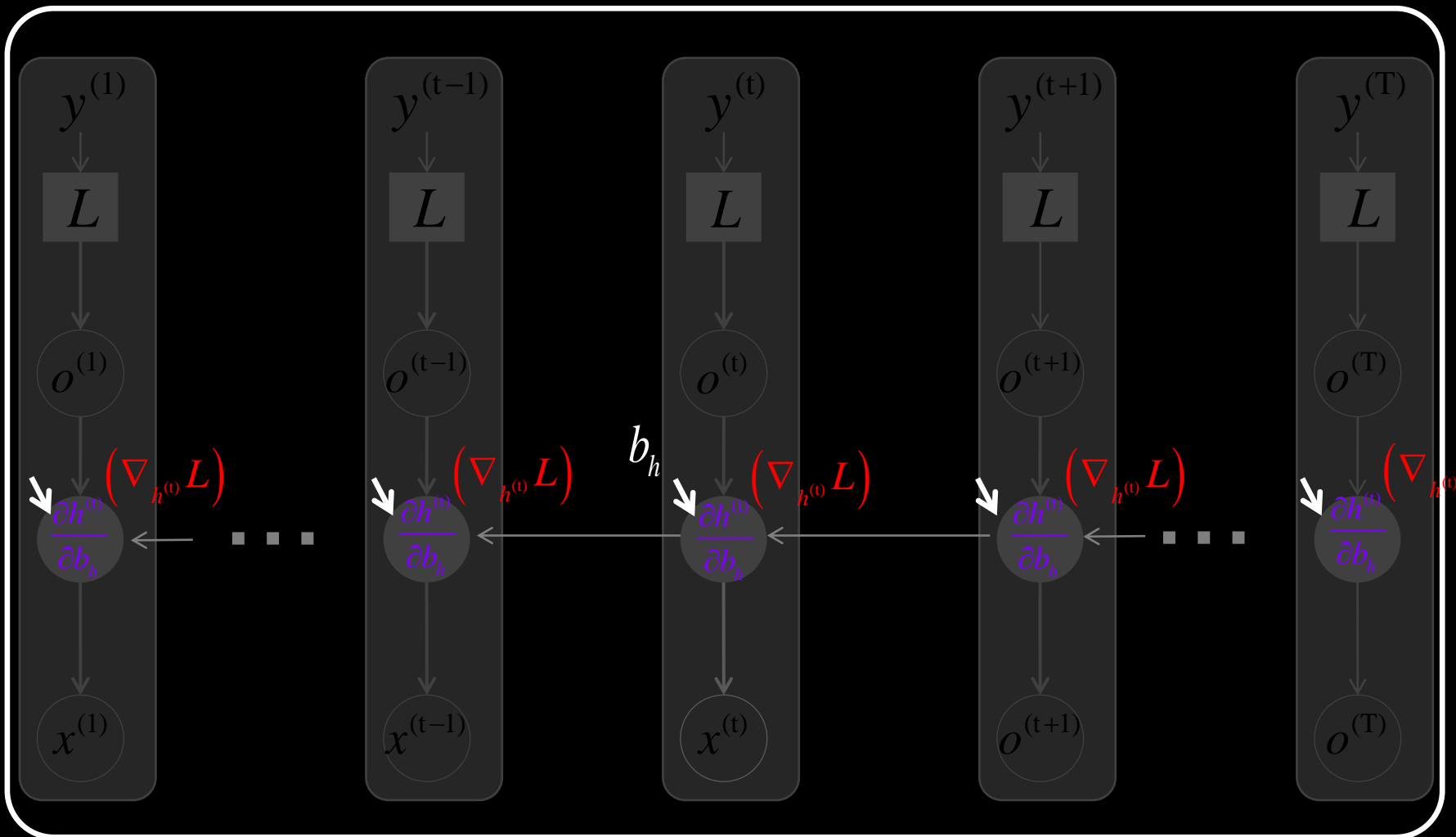
$$\nabla_{h^{(t)}} L = \left( \frac{\partial h^{(t+1)}}{\partial h^{(t)}} \right) \nabla_{h^{(t+1)}} L + \left( \frac{\partial o^{(t)}}{\partial h^{(t)}} \right) \nabla_{o^{(t)}} L$$

$$\nabla_{b_o} L = \sum_{t=1}^{T} \left( \frac{\partial o^{(t)}}{\partial b_o} \right) \left( \nabla_{o^{(t)}} L \right)$$

$$\nabla_{b_h} L = \sum_{t=1}^{T} \left( \frac{\partial h^{(t)}}{\partial b_h} \right) \left( \nabla_{h^{(t)}} L \right)$$

# Back-Propagation Through Time (BPTT)



$$\nabla_{L^{(t)}} L = \frac{\partial L}{\partial L^{(t)}}$$

$$\nabla_{o^{(t)}} L = \left( \nabla_{L^{(t)}} L \right) \frac{\partial L^{(t)}}{\partial o_i^{(t)}}$$

$$\nabla_{h^{(t)}} L = \left( \frac{\partial h^{(t+1)}}{\partial h^{(t)}} \right) \nabla_{h^{(t+1)}} L + \left( \frac{\partial o^{(t)}}{\partial h^{(t)}} \right) \nabla_{o^{(t)}} L$$
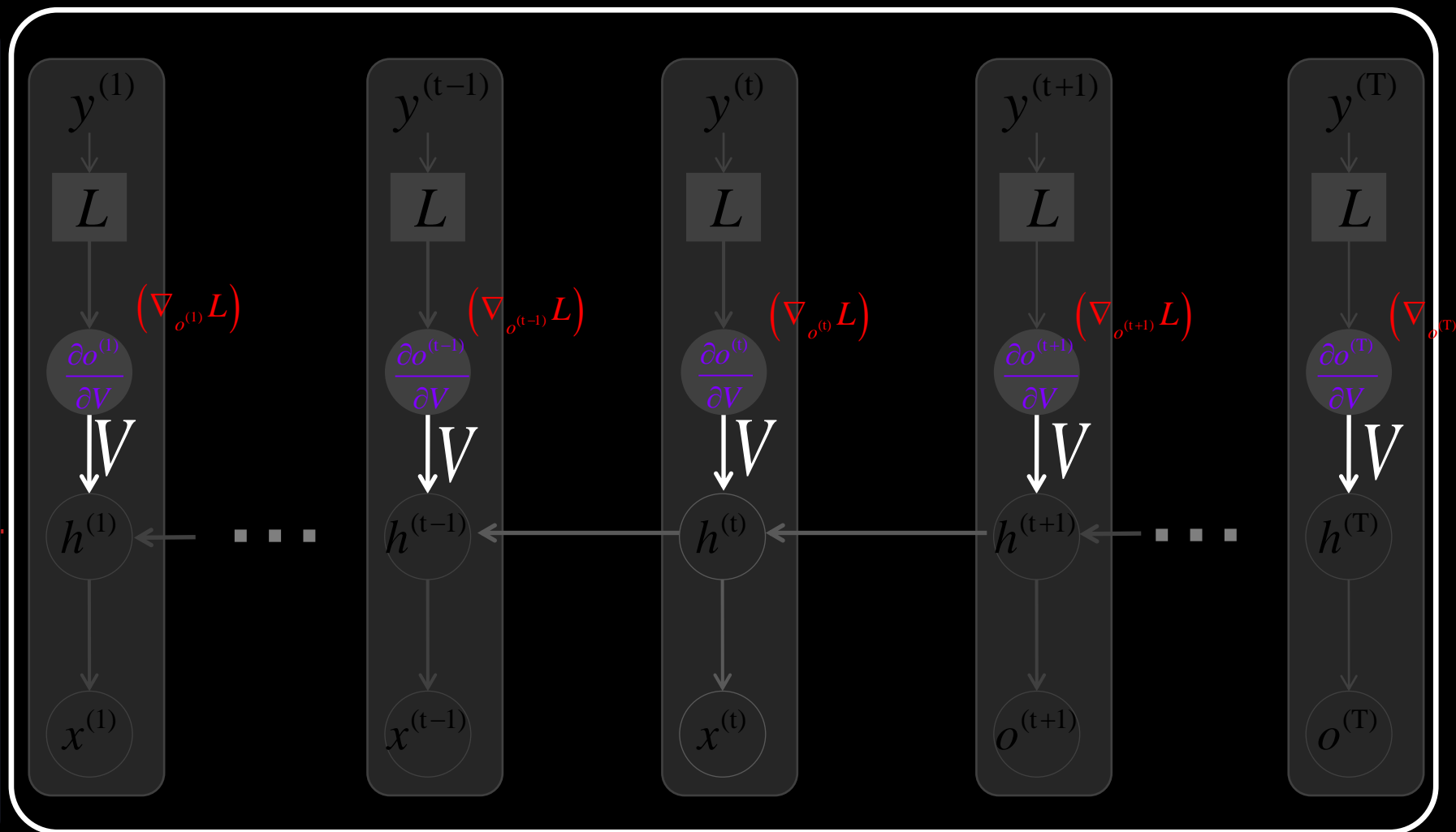
$$\nabla_{b_o} L = \sum_{t=1}^{T} \frac{\partial o^{(t)}}{\partial b_o} \left( \nabla_{o^{(t)}} L \right)$$

$$\nabla_{b_h} L = \sum_{t=1}^{T} \left( \frac{\partial h^{(t)}}{\partial b_h} \right) \left( \nabla_{h^{(t)}} L \right)$$

$$\nabla_V L = \sum_{t=1}^{T} \left( \frac{\partial o^{(t)}}{\partial V} \right) \left( \nabla_{o^{(t)}} L \right)$$

biases

Weights

# Back-Propagation Through Time (BPTT)

$$\nabla_{L^{(t)}} L = \frac{\partial L}{\partial L^{(t)}}$$

$$\nabla_{o^{(t)}} L = \left(\nabla_{L^{(t)}} L\right) \frac{\partial L^{(t)}}{\partial o_i^{(t)}}$$

$$\nabla_{h^{(t)}} L = \left(\frac{\partial h^{(t+1)}}{\partial h^{(t)}}\right) \nabla_{h^{(t+1)}} L + \left(\frac{\partial o^{(t)}}{\partial h^{(t)}}\right) \nabla_{o^{(t)}} L$$

$$\nabla_{b_o} L = \sum_{t=1}^{T} \frac{\partial o^{(t)}}{\partial b_o} \left(\nabla_{o^{(t)}} L\right)$$
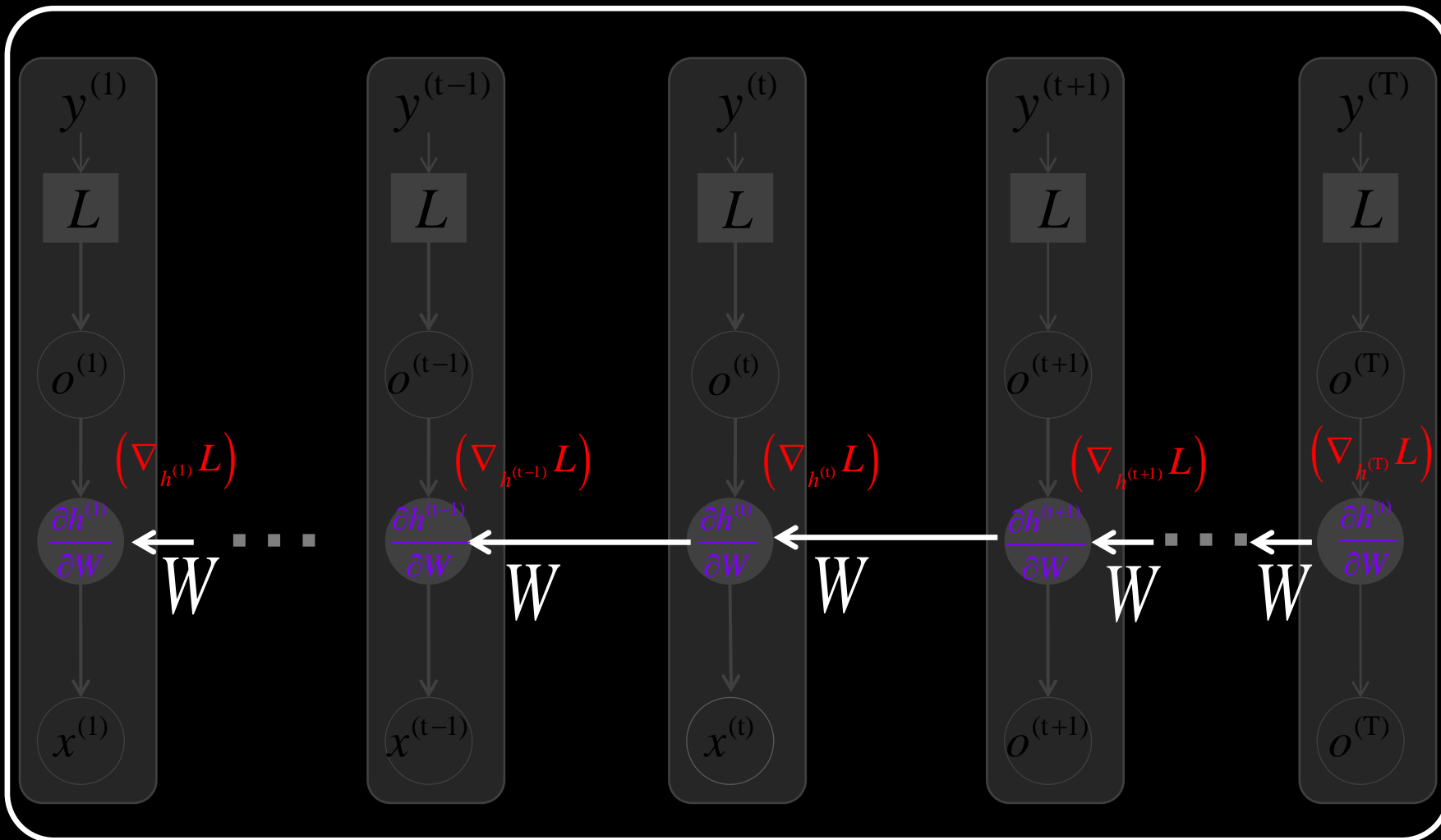
$$\nabla_{b_h} L = \sum_{t=1}^{T} \left(\frac{\partial h^{(t)}}{\partial b_h}\right) \left(\nabla_{h^{(t)}} L\right)$$
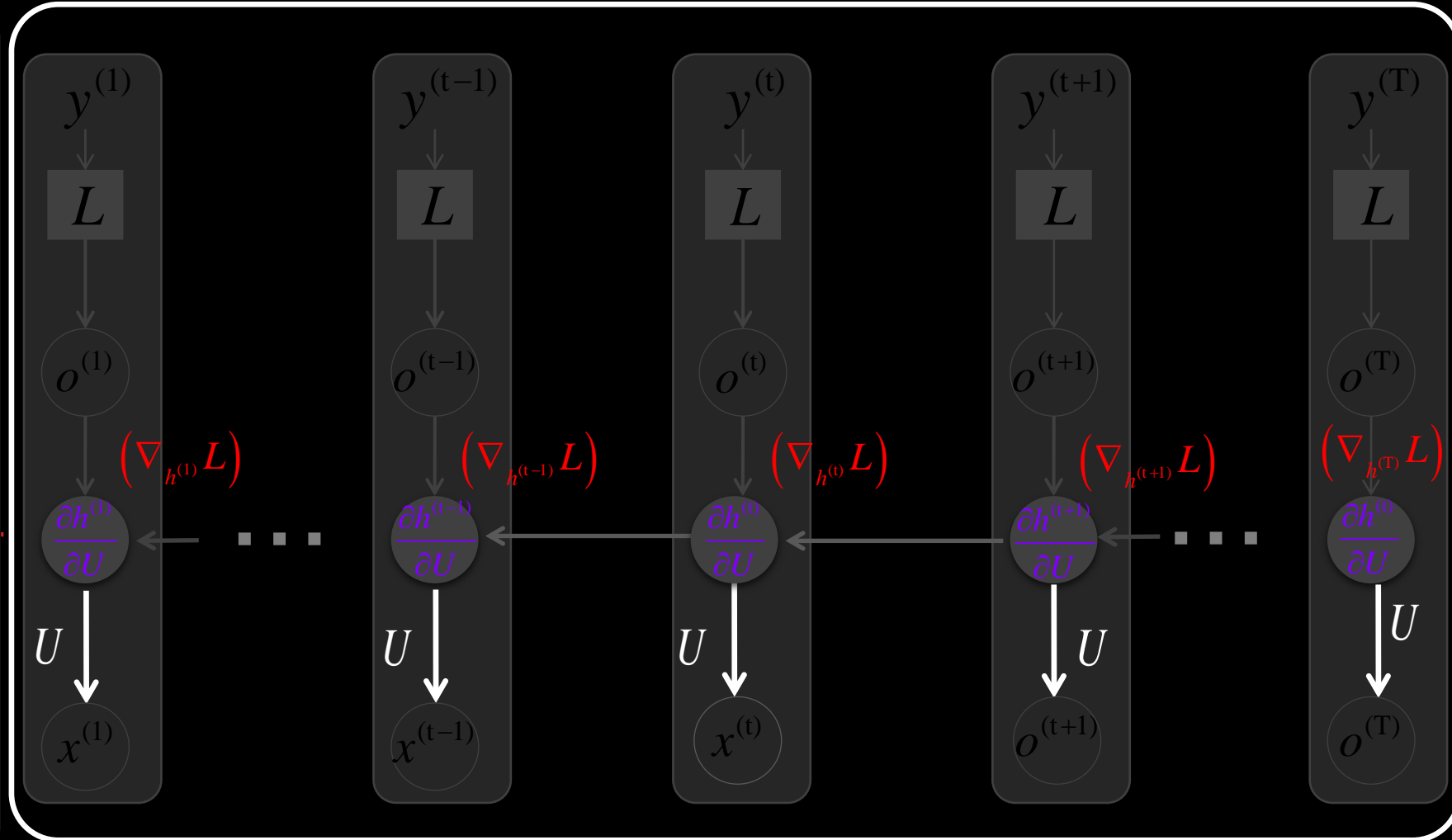
$$\nabla_V L = \sum_{t=1}^{T} \left(\frac{\partial o^{(t)}}{\partial V}\right) \left(\nabla_{o^{(t)}} L\right)$$

$$\nabla_W L = \sum_{t=1}^{T} \left(\frac{\partial h^{(t)}}{\partial W}\right) \left(\nabla_{h^{(t)}} L\right)$$

biases

Weights

# Back-Propagation Through Time (BPTT)

# RNN: one to one



Vanilla Neural Networks

# RNN: one to many



Image Captioning
image -> sequence of words

# RNN: many to one



Sentiment Classification
sequence of words -> sentiment

# RNN: many to many



Machine Translation
seq of words -> seq of words

# RNN: many to many



Video classification on frame level

# Bidirectional RNNs

Vanilla RNNs
        looks into the one side of sequence (left or past) to predict the next output.

Bidirectional RNNs (BiRNNs)
        can focus on both past and future (right side of the sequence).

# Bidirectional RNNs



$$\vec{h}_t = f\left(\vec{U}x_t + \vec{W}h_{t-1} + \vec{b}_h\right)$$

$$\overleftarrow{h}_t = f\left(\overleftarrow{U}x_t + \overleftarrow{W}h_{t+1} + \overleftarrow{b}_h\right)$$

$$\hat{y} = g(V h_t + b_o) = g(V[\vec{h}_t ; \overleftarrow{h}_t] + b_o)$$

# Deep Bidirectional RNNs



$$\vec{h}_t^{(i)} = f\left(\vec{U}^{(i)} h_t^{(i-1)} + \vec{W} h_{t-1}^{(i)} + \vec{b}_h^{(i)}\right)$$

$$\overleftarrow{h}_t^{(i)} = f\left(\overleftarrow{U}^{(i)} h_t^{(i-1)} + \overleftarrow{W} h_{t-1}^{(i)} + \overleftarrow{b}_h^{(i)}\right)$$

$$\hat{y}_{(t)} = g(V\, h_t + b_o) = g(U[\vec{h}_t^{(L)} ; \overleftarrow{h}_t^{(L)}] + b_o)$$

# Two Problems of Vanilla RNNS

A) Information Morphing

B) Vanishing/Exploding of the Gradient

# Problem 1: Information morphing

Suppose $x^{(3)}$ is important and we must remember it.



Goal:
$$h_3 = h_4 = \ldots = h_{t+1}$$

# Problem 1: Information morphing

$h_{(4)}$ contains this information.



$$h^{(4)} = \varphi\left(W\, h_3 + b\right)$$

# Problem 1: Information morphing

But, this information is washing away as time precedes...



$$h^{(4)} = \varphi\left(W\, h_3 + b\right) \qquad h^{(t)} = \varphi\left(W\, h_{t-1} + b\right)$$

# Problem 2: Exploding Gradient

If the weights are big, the gradients grow exponentially.

# Problem 2: Vanishing Gradient

If the weights are small, the gradients shrink exponentially.

# Two obvious solutions

- Better weight initialization (instead of random).

- Rectified Linear Unit (RLU) as activation function.

# Gradient clipping

Simple heuristic solution that clips gradients to a small number whenever they explode.



[Thomas Mikolov et al. , Pascanu, Razvan, Tomas Mikolov, and Yoshua Bengio, 2013]

# Questions?

# Can We Change the Structure?

Let's reformulated our vanilla RNN in the following form:

$$
\begin{pmatrix} h^{(t)} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix} = function \left[ \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}, \begin{pmatrix} x^{(t)} \\ \text{input} \\ \text{at time } t \end{pmatrix} \right]
$$

# Can We Change the Structure?

Let's reformulated our vanilla RNN in the following form:

$$\begin{pmatrix} h^{(t)} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix} = function \left[ \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}, \begin{pmatrix} x^{(t)} \\ \text{input} \\ \text{at time } t \end{pmatrix} \right]$$

Memory

The Source of problem

# Memory's Properties

Memory should has three properties:

I) Forgettablility

II) Writability

III) Readability

These ideas of memory is used to develop robust RNNS such as GRU and LSTM.

These ideas of memory is used to develop robust gated RNN cells such as GRU and LSTM.

Before, explaining GRU
Let's talk about its intuition.

# Linear operation on Memory

Memory should be a linear combination of Forgettablility and Writability.

$$
\begin{pmatrix} h^{(t)} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix} = \begin{pmatrix} \text{Clearning} \\ \text{Memory} \end{pmatrix} + \begin{pmatrix} \text{Writing} \\ \text{Memory} \end{pmatrix}
$$

# Clearing and Writing on Memory

We need two separate parts to clear the unnecessary information from the memory and writing the new information.

$$\begin{pmatrix} h^{(t)} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix} = \begin{pmatrix} \text{Clearning} \\ \text{Memory} \end{pmatrix} + \begin{pmatrix} \text{Writing} \\ \text{Memory} \end{pmatrix}$$

$$\begin{pmatrix} \text{Clearning} \\ \text{Memory} \end{pmatrix} = \begin{pmatrix} f^{(t)} \\ \text{forget gate} \\ \text{binary vector} \end{pmatrix} \otimes \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}$$

$$\begin{pmatrix} \text{Writing} \\ \text{Memory} \end{pmatrix} = \begin{pmatrix} s^{(t)} \\ \text{store gate} \\ \text{binary vector} \\ \text{at time } t \end{pmatrix} \otimes \begin{pmatrix} \tilde{h}^{(t)} \\ \text{candidate} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix}$$

$$f^{(t)}, s^{(t)}, \tilde{h}^{(t)}$$

# Clearing Memory

$$
\begin{pmatrix} \text{Clearning} \\ \text{Memory} \end{pmatrix} = \begin{pmatrix} \color{red}{f^{(t)}} \\ \text{forget gate} \\ \text{binary vector} \end{pmatrix} \otimes \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}
$$

$$
\color{red}{f^{(t)}} = \sigma_f \left[ \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}, \begin{pmatrix} x^{(t)} \\ \text{input} \\ \text{at time } t \end{pmatrix} \right]
$$

# Writing on Memory

$$\begin{pmatrix} \text{Writing} \\ \text{Memory} \end{pmatrix} = \begin{pmatrix} \textcolor{red}{s^{(t)}} \\ \text{store gate} \\ \text{binary vector} \\ \text{at time } t \end{pmatrix} \otimes \begin{pmatrix} \textcolor{red}{\tilde{h}^{(t)}} \\ \text{candidate} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix}$$

$$\textcolor{red}{s^{(t)}} = \sigma_s \left[ \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}, \begin{pmatrix} x^{(t)} \\ \text{input} \\ \text{at time } t \end{pmatrix} \right]$$

$$\begin{pmatrix} \textcolor{red}{\tilde{h}^{(t)}} \\ \text{candidate} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix} = \varphi \left[ \begin{pmatrix} \textcolor{green}{r^{(t)}} \\ \text{read gate} \\ \text{at time } t \end{pmatrix} \otimes \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}, \begin{pmatrix} x^{(t)} \\ \text{input} \\ \text{at time } t \end{pmatrix} \right]$$

# Writing on Memory

$$
\begin{pmatrix} \text{Writing} \\ \text{Memory} \end{pmatrix} = \overset{\textcolor{red}{s^{(t)}}}{\begin{pmatrix} \text{store gate} \\ \text{binary vector} \\ \text{at time } t \end{pmatrix}} \times \overset{\textcolor{red}{\tilde{h}^{(t)}}}{\begin{pmatrix} \text{candidate} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix}}
$$

$$
\textcolor{red}{s^{(t)}} = \sigma_f \left[ \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}, \begin{pmatrix} x^{(t)} \\ \text{input} \\ \text{at time } t \end{pmatrix} \right]
$$

$$
\overset{\textcolor{red}{\tilde{h}^{(t)}}}{\begin{pmatrix} \text{candidate} \\ \text{hid. state} \\ \text{at time } t \end{pmatrix}} = \tanh \left[ \overset{\textcolor{green}{r^{(t)}}}{\begin{pmatrix} \text{read gate} \\ \text{at time } t \end{pmatrix}} \otimes \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}, \begin{pmatrix} x^{(t)} \\ \text{input} \\ \text{at time } t \end{pmatrix} \right]
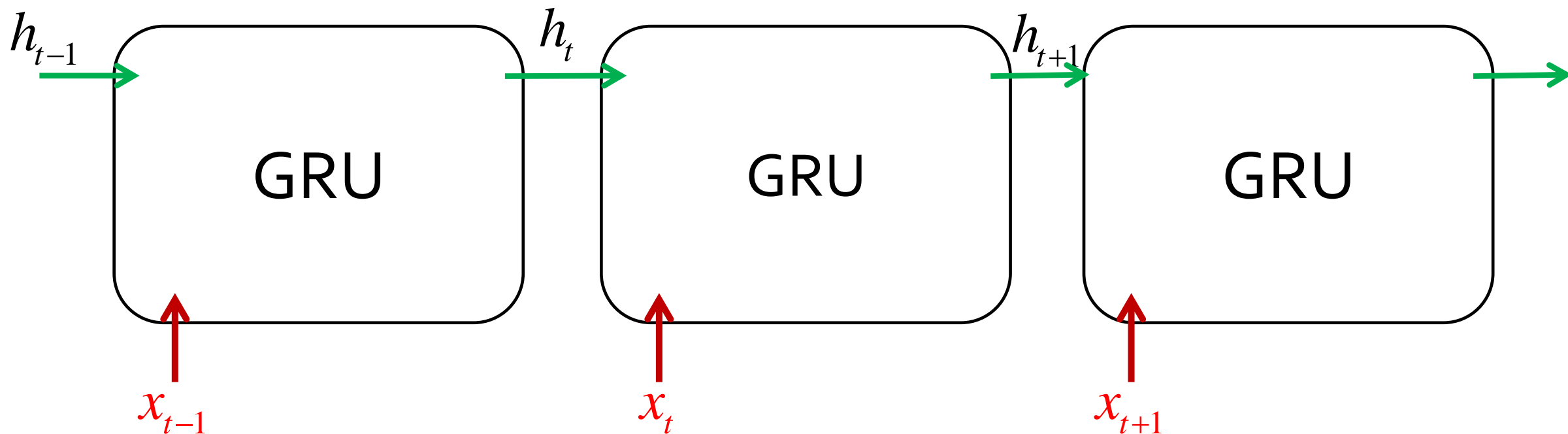$$

$$
\textcolor{red}{r^{(t)}} = \sigma_f \left[ \begin{pmatrix} h^{(t-1)} \\ \text{hid. state} \\ \text{at time } t-1 \end{pmatrix}, \begin{pmatrix} x^{(t)} \\ \text{input} \\ \text{at time } t \end{pmatrix} \right]
$$

# Gated Recurrent Unit-(GRU)

# GRU: Forget gate



$$f_t = \sigma(\mathbf{U}_f \mathbf{x}_t + \mathbf{W}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

# GRU: Forget gate and Store gate



$$f_t = \sigma(\mathrm{U}_f \mathrm{x}_t + \mathrm{W}_f \mathrm{h}_{t-1} + \mathrm{b}_f)$$

$$\mathrm{s}_t = 1 - f_t$$

# GRU: Reading gate



$$f_t = \sigma(\mathbf{U}_f \mathbf{x}_t + \mathbf{W}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$s_t = 1 - f_t$$

$$r_t = \sigma(\mathbf{U}_r \mathbf{x}_t + \mathbf{W}_r \mathbf{h}_{t-1} + \mathbf{b}_r)$$

# GRU: Candidate Hidden State



$$f_t = \sigma(\mathrm{U}_f \mathrm{x}_t + \mathrm{W}_f \mathrm{h}_{t-1} + \mathrm{b}_f)$$

$$\mathrm{s}_t = 1 - f_t$$

$$r_t = \sigma(\mathrm{U}_r \mathrm{x}_t + \mathrm{W}_r \mathrm{h}_{t-1} + \mathrm{b}_r)$$

$$\tilde{h}_t = \tanh\big(\mathrm{W}(\mathrm{r}_t \otimes \mathrm{h}_{t-1}) + Ux_t + b\big)$$

# GRU: Hidden State



$$f_t = \sigma(\mathrm{U}_f \mathrm{x}_t + \mathrm{W}_f \mathrm{h}_{t-1} + \mathrm{b}_f)$$

$$\mathrm{s}_t = 1 - f_t$$

$$r_t = \sigma(\mathrm{U}_r \mathrm{x}_t + \mathrm{W}_r \mathrm{h}_{t-1} + \mathrm{b}_r)$$

$$\tilde{h}_t = \tanh\big(\mathrm{W}(\mathrm{r}_t \otimes \mathrm{h}_{t-1}) + Ux_t + b\big)$$

# GRU: Memory Update



$$f_t = \sigma(\mathrm{U}_f \mathrm{x}_t + \mathrm{W}_f \mathrm{h}_{t-1} + \mathrm{b}_f)$$

$$\mathrm{s}_t = 1 - f_t$$

$$r_t = \sigma(\mathrm{U}_r \mathrm{x}_t + \mathrm{W}_r \mathrm{h}_{t-1} + \mathrm{b}_r)$$

$$\tilde{h}_t = \tanh\big(\mathrm{W}(\mathrm{r}_t \otimes \mathrm{h}_{t-1}) + Ux_t + b\big)$$

$$h_t = \big(f_t \otimes \mathrm{h}_{t-1}\big) + \big(\mathrm{s}_t \otimes \tilde{h}_t\big)$$

# Long Short-Term Memory (LSTM)

# LSTM- Forget Gate



$$f^{(t)} = \sigma(U_f x^{(t)} + W_f h^{(t-1)} + b_f)$$

# LSTM- Store Gate



$$f^{(t)} = \sigma(\mathrm{U}_f \mathrm{x}^{(t)} + \mathrm{W}_f \mathrm{h}^{(t-1)} + \mathrm{b}_f)$$

$$i^{(t)} = \sigma(\mathrm{U}_i \mathrm{x}^{(t)} + \mathrm{W}_i \mathrm{h}^{(t-1)} + \mathrm{b}_i)$$

$$\tilde{c}^{(t)} = \tanh\left(\mathrm{W} \mathrm{h}^{(t-1)} + U \mathrm{x}^{(t)} + b\right)$$

# LSTM- Memory Update



$$f^{(t)} = \sigma(U_f x^{(t)} + W_f h^{(t-1)} + b_f)$$

$$i^{(t)} = \sigma(U_i x^{(t)} + W_i h^{(t-1)} + b_i)$$

$$\tilde{c}^{(t)} = \tanh\left(W h^{(t-1)} + U x^{(t)} + b\right)$$

$$c^{(t)} = \left(f^{(t)} \circ \tilde{c}^{(t-1)}\right) + \left(i \circ \tilde{c}^{(t)}\right)$$

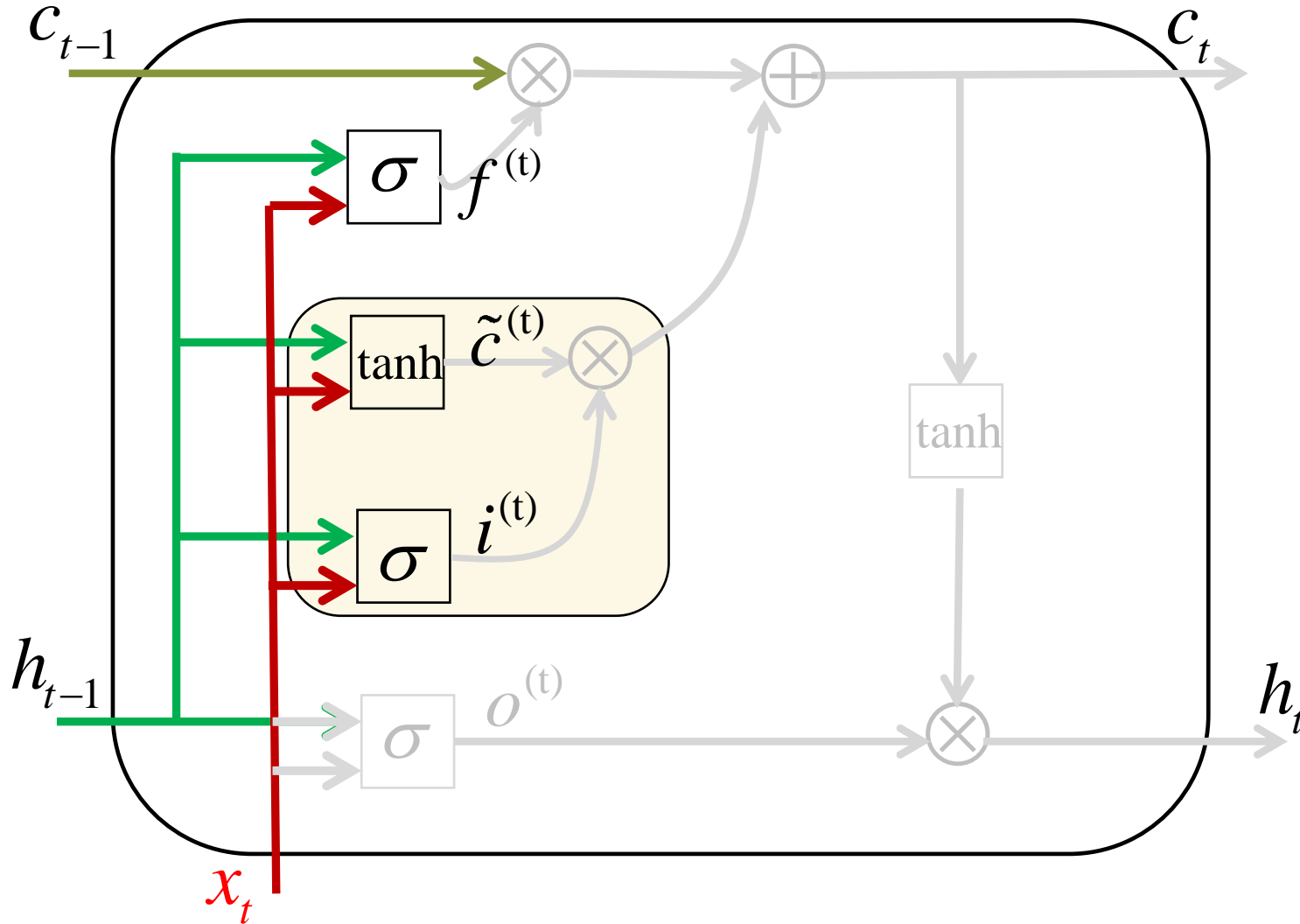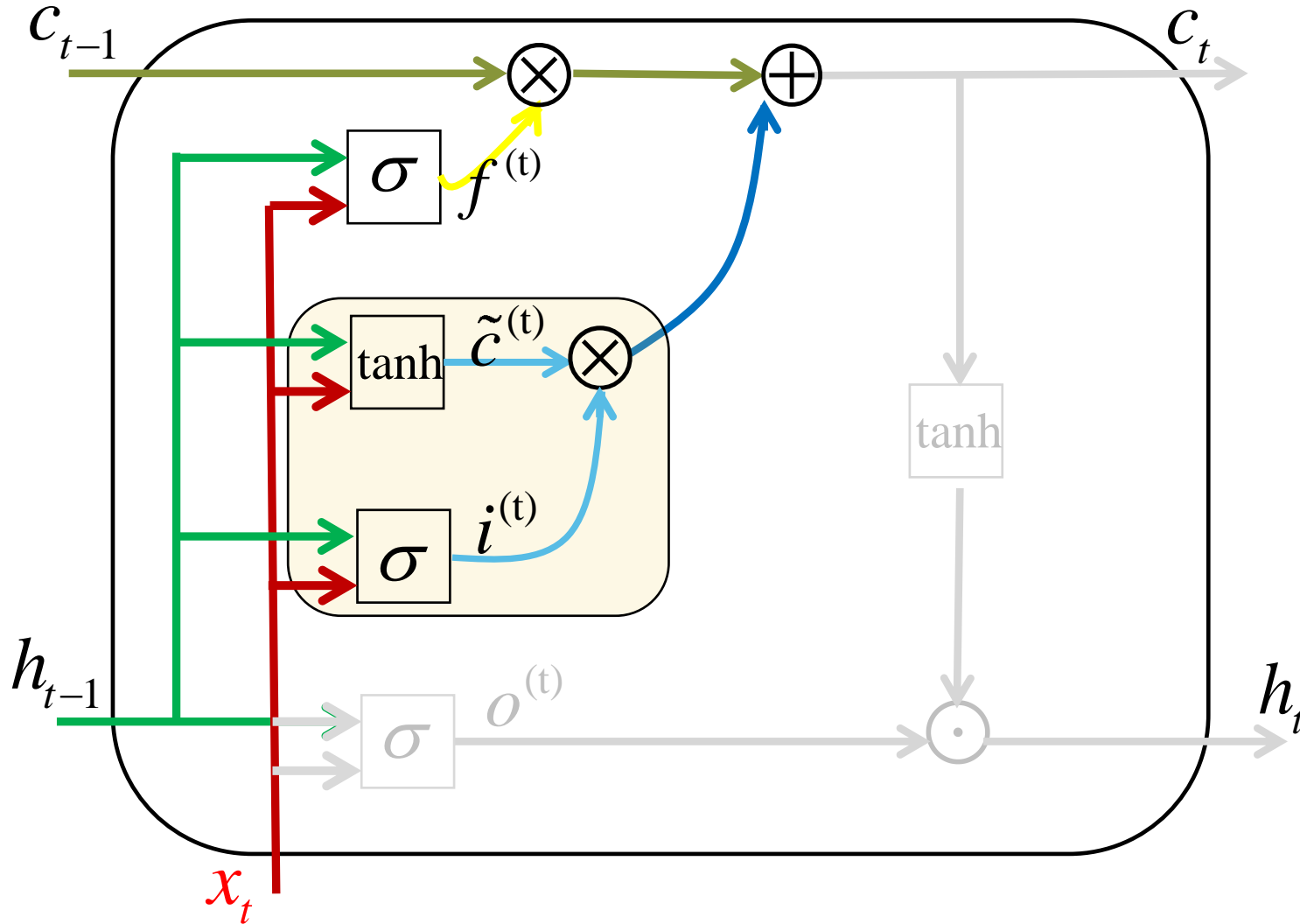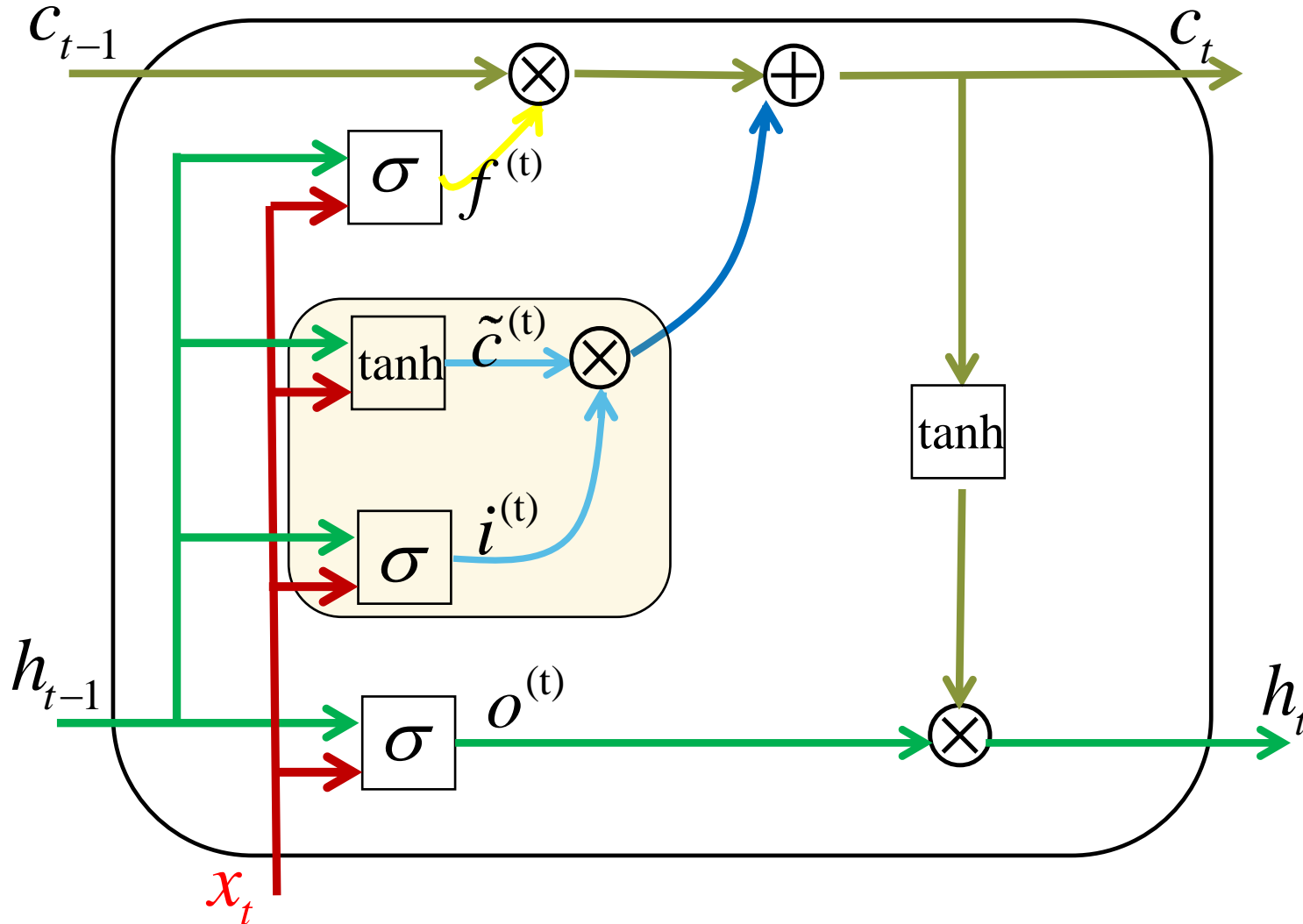# LSTM- Output Gate



$$f^{(t)} = \sigma(U_f x^{(t)} + W_f h^{(t-1)} + b_f)$$

$$i^{(t)} = \sigma(U_i x^{(t)} + W_i h^{(t-1)} + b_i)$$

$$\tilde{c}^{(t)} = \tanh\left(W h^{(t-1)} + U x^{(t)} + b\right)$$

$$c^{(t)} = \left(f^{(t)} \circ \tilde{c}^{(t-1)}\right) + \left(i \circ \tilde{c}^{(t)}\right)$$

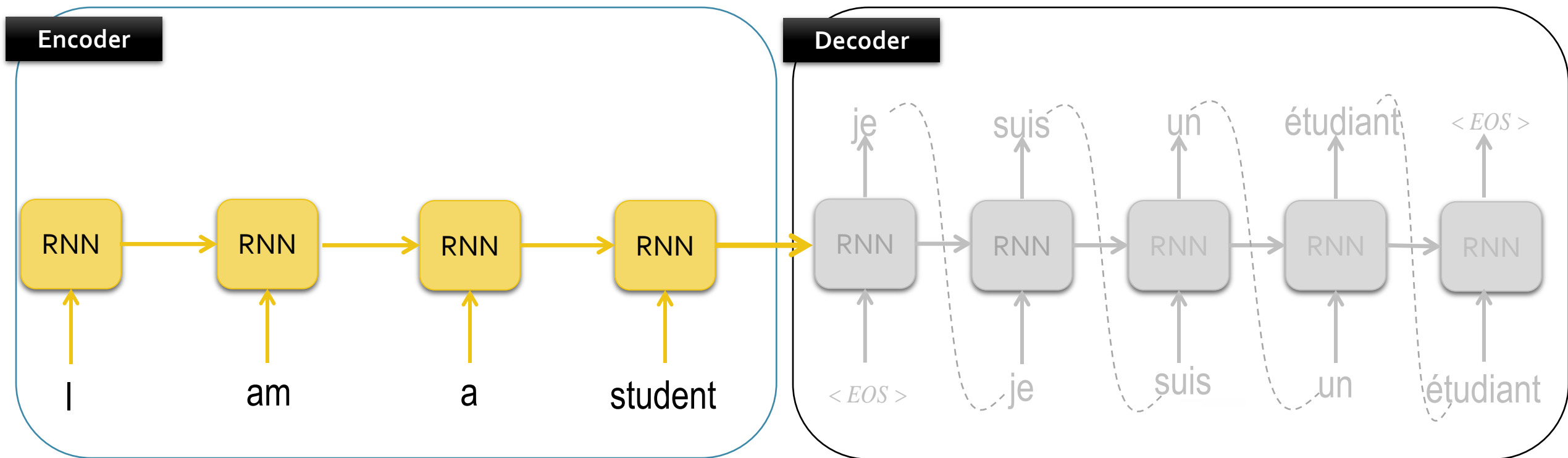$$o^{(t)} = \sigma(U_o x^{(t)} + W_o h^{(t-1)} + b_o)$$

$$h^{(t)} = o^{(t)} \circ \tanh(c^{(t)})$$

$$y^{(t)} = h^{(t)}$$

# Questions?

# Some examples of RNN

Sequence to Sequence

Encoder

RNN RNN RNN RNN

I am a student

Decoder

je suis un étudiant < EOS >

RNN RNN RNN RNN RNN

< EOS > je suis un étudiant

Sequence to Sequence

Sequence to Sequence

Deep Sequence to Sequence

Encoder

Decoder

$y_1 \rightarrow y_2 ... \rightarrow < EOS >$

SoftMax

$h_{L,1}$ $h_{L,2}$ $h_{L,N}$ $s_{L,1}$ $s_{L,2}$ $s_{L,N}$

$h_{2,1}$ $h_{2,2}$ $h_{2,N}$ $s_{2,1}$ $s_{2,2}$ $s_{2,N}$

$h_{1,1}$ $h_{1,2}$ $h_{1,N}$ $s_{1,1}$ $s_{1,2}$ $s_{1,N}$

$x_0$ $x_1$ $x_N$ $< EOS >$ $y_1$ $y_M$

- The last hidden state summarizes the entire input sentence into C.

[D Bahdanau, K Cho, Y Bengio, 2015]

[D Bahdanau, K Cho, Y Bengio, 2015]

**Decoder**

**Encoder**

$\hat{y}_{t-1}$

$\hat{y}_t$

$< EOS >$

g

**g**

g

$S_{t-1}$

RNN

$S_t$

RNN

$S_{T_y}$

RNN

$c_t$

$\hat{y}_{t-1}$

$\alpha_{t,1}$

$\alpha_{t,2}$

$\alpha_{t,T}$

Bi-RNN

$\leftarrow h_1$

$\rightarrow h_1$

Bi-RNN

$\leftarrow h_2$

$\leftarrow h_2$

$\leftarrow h_{T_x}$

$\leftarrow h_{T_x}$

Bi-RNN

$x_0$

$x_1$

$x_{T_x}$

Attention based Encoder-Decoder

$$h_j = \left[ \rightarrow h_j ; \leftarrow h_j \right]_{concat}$$

Decoder

Encoder

$\hat{y}_{t-1}$   $\hat{y}_t$   $< EOS >$

g   g   g

$s_{t-1}$   $s_t$   $s_T$

$\hat{y}_{t-1}$

$c_t$

$\alpha_{t,1}$   $\alpha_{t,2}$   $\alpha_{t,T}$

Bi-RNN   $\leftarrow h_1$   Bi-RNN   $\leftarrow h_2$   $\leftarrow h_{T_x}$   Bi-RNN

$\rightarrow h_1$   $\rightarrow h_2$   $\rightarrow h_{T_x}$

$x_0$   $x_1$   $x_{T_x}$
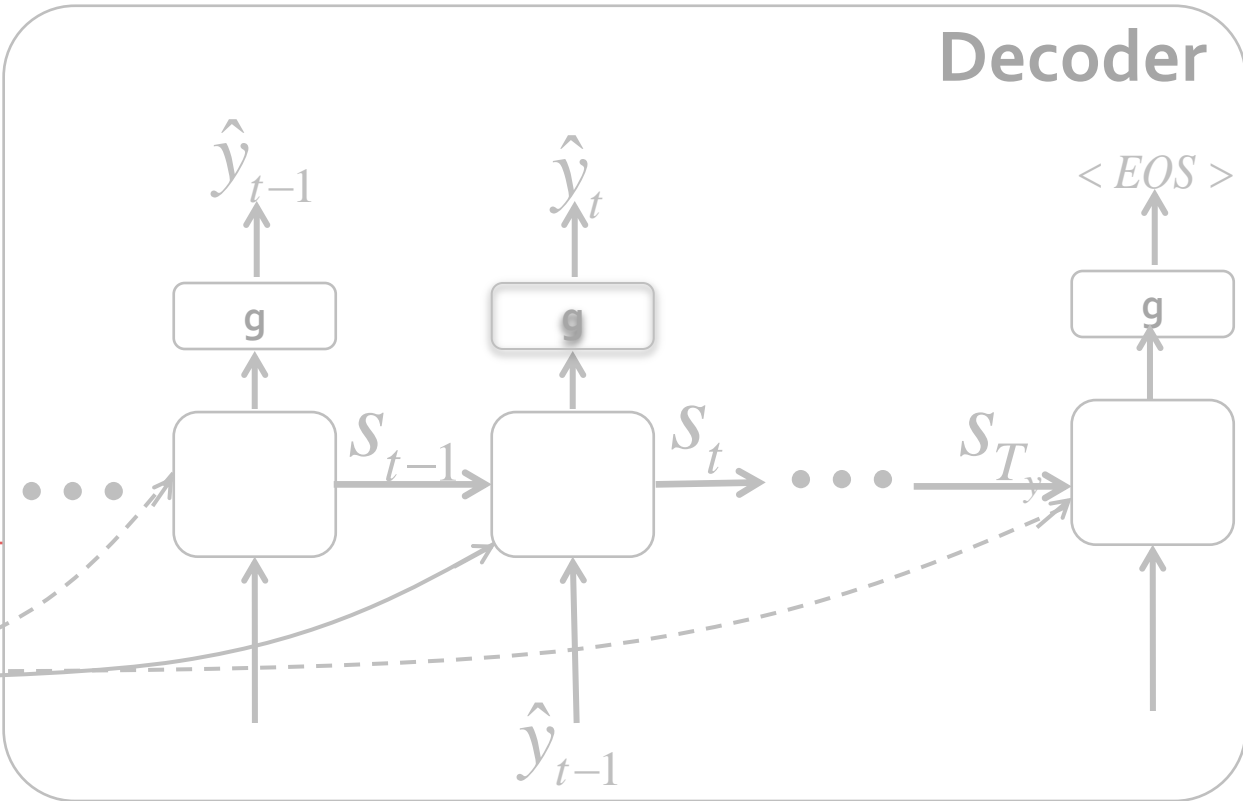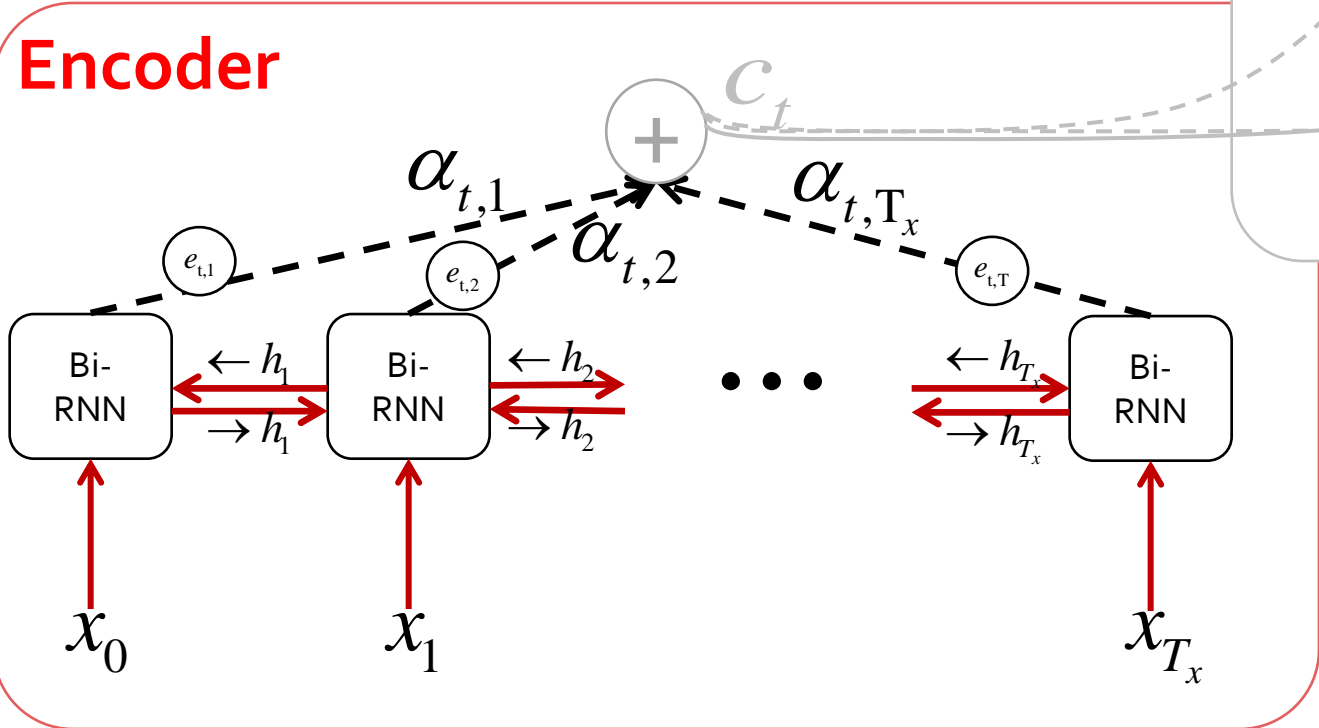
Attention based Encoder-Decoder

$$h_j = \left[ \rightarrow h_j ; \leftarrow h_j \right]_{concat}$$

$$\alpha_{i,j} = \frac{exp(e_{i,j})}{\sum_{k=1}^{T_x} exp(e_{i,k})}$$

$$e_{i,j} = a(s_{i-1}, h_j)$$
$$= v_a^{\mathrm{T}} \tanh(W_a s_{i-1} + U_a h_j)$$

Decoder

$\hat{y}_{t-1}$    $\hat{y}_t$    $<EOS>$

g    g    g

$S_{t-1}$    $S_t$    $S_T$

$c_t$

$\hat{y}_{t-1}$

Encoder

$\alpha_{t,1}$    $\alpha_{t,T_x}$
$\alpha_{t,2}$

$e_{t,1}$    $e_{t,2}$    $e_{t,T}$

Bi-RNN    $\leftarrow h_1$    Bi-RNN    $\leftarrow h_2$    $\cdots$    $\leftarrow h_{T_x}$    Bi-RNN
$\rightarrow h_1$    $\rightarrow h_2$    $\rightarrow h_{T_x}$

$x_0$    $x_1$    $x_{T_x}$

$$a_t(s) = \frac{e^{\text{score}(s)}}{\sum_{s'} e^{\text{score}(s')}}$$



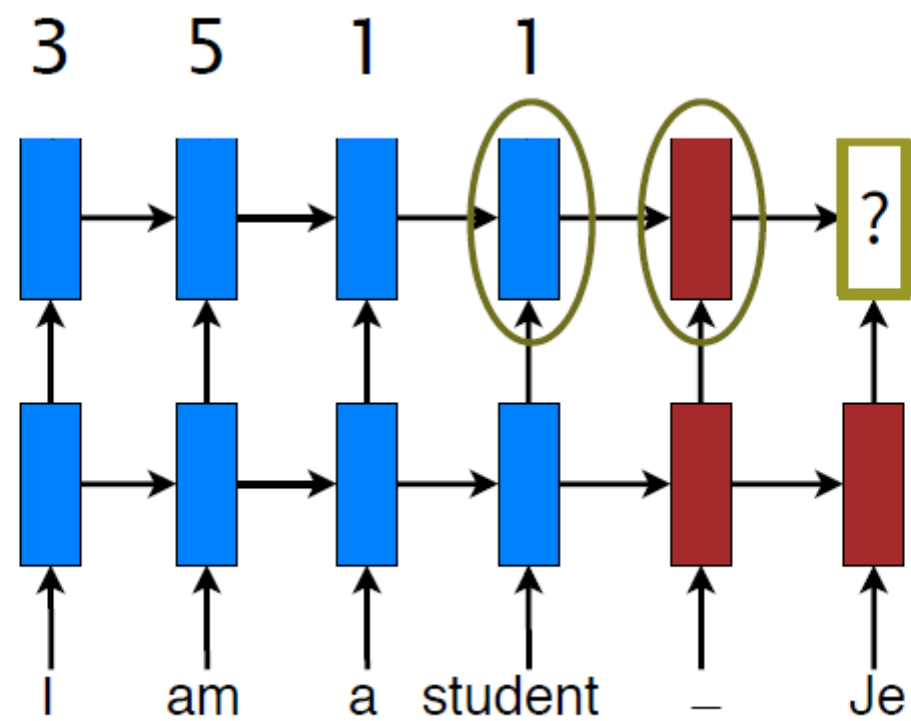$a_t$  0.3 0.5 0.1 0.1

I   am   a   student   _   Je

Attention based Encoder-Decoder

**Decoder**

$h_j = \left[ \rightarrow h_j ; \leftarrow h_j \right]_{concat}$

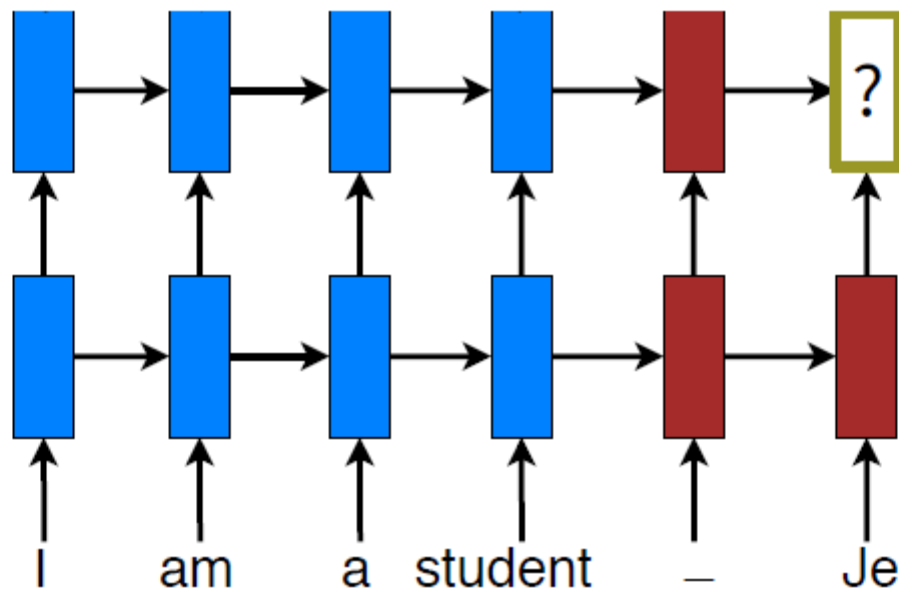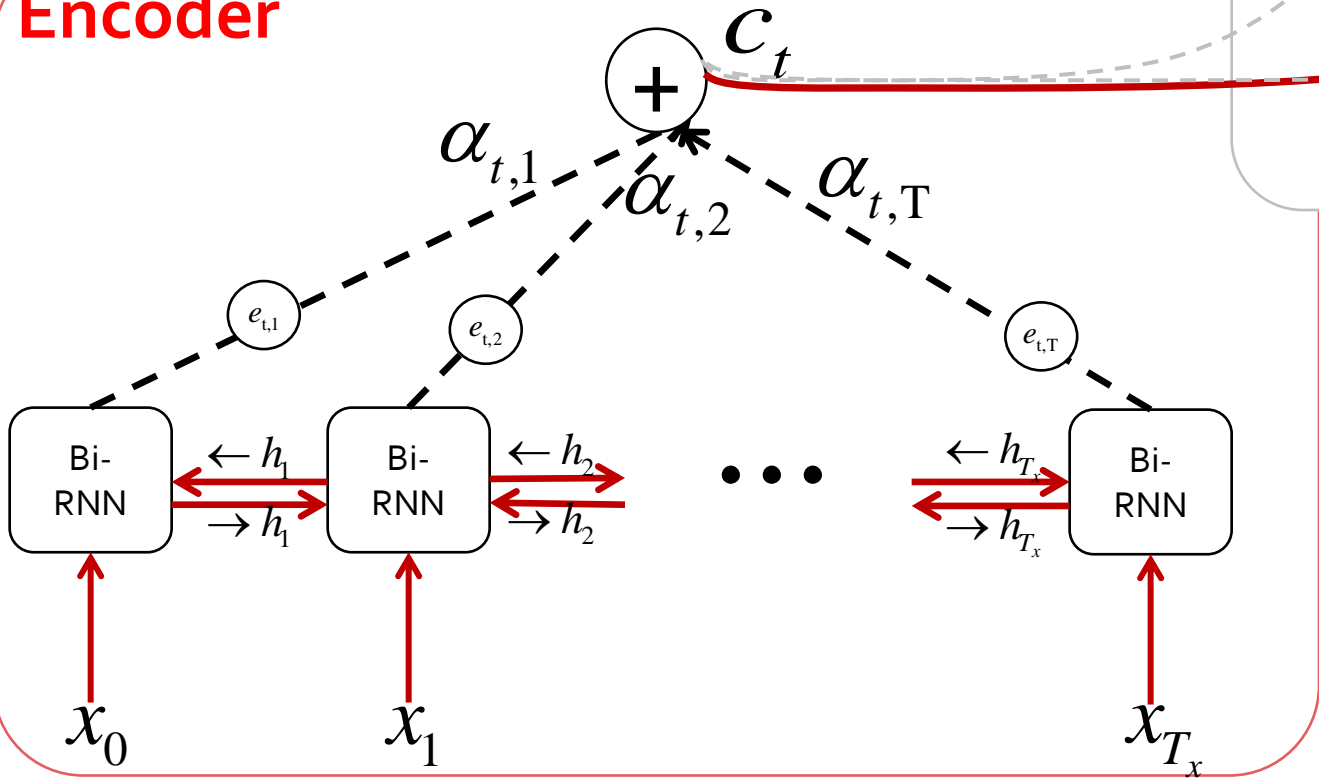$\alpha_{i,j} = \dfrac{exp(e_{i,j})}{\sum_{k=1}^{T_x} exp(e_{i,k})}$

$e_{i,j} = a(s_{i-1}, h_j)$

$\quad = v_a^T \tanh(W_a s_{i-1} + U_a h_j)$

$c_i = \sum_{k=1}^{T_x} \alpha_{i,j} h_j$

$\hat{y}_{t-1}$  $\hat{y}_t$  $<EOS>$

g  g  g

$S_{t-1}$  $S_t$  $S_{T_x}$

$\hat{y}_{t-1}$

**Encoder**

$c_t$

$+$

$\alpha_{t,1}$  $\alpha_{t,2}$  $\alpha_{t,T}$

$e_{t,1}$  $e_{t,2}$  $e_{t,T}$

Bi-RNN  $\leftarrow h_1$  Bi-RNN  $\leftarrow h_2$  $\cdots$  $\leftarrow h_{T_x}$  Bi-RNN

$\rightarrow h_1$  $\rightarrow h_2$  $\rightarrow h_{T_x}$

$x_0$  $x_1$  $x_{T_x}$

Attention based Encoder-Decoder

$$p(\mathrm{y}_i \mid \mathrm{y}_1, \mathrm{y}_2, ..., \mathrm{y}_{i-1}, \mathrm{X}) = g(\mathrm{y}_{i-1}, \mathrm{s}_i, \mathrm{c}_i)$$

Decoder

Encoder

Attention based Encoder-Decoder

$$p(\mathrm{y}_i \mid \mathrm{y}_1, \mathrm{y}_2, ..., \mathrm{y}_{i-1}, \mathrm{X}) = \mathrm{g}(\mathrm{y}_{i-1}, \mathrm{s}_i, \mathrm{c}_i)$$

$$s_i = f(s_{i-1}, \hat{\mathrm{y}}_{i-1}, \mathrm{c}_i)$$

**Decoder**

$\hat{y}_{t-1}$

$\hat{y}_t$

$< EOS >$

g

g

g

RNN $\quad S_{t-1}$ RNN $\quad S_t$ $\quad S_{T_y}$ RNN

$\hat{y}_{t-1}$

**Encoder**

$c_t$

$\alpha_{t,1}$ $\quad +$ $\quad \alpha_{t,T}$

$\alpha_{t,2}$

$e_{t,1}$ $\quad e_{t,2}$ $\quad e_{t,T}$

Bi-RNN $\quad \leftarrow h_1$ $\quad$ Bi-RNN $\quad \leftarrow h_2$ $\quad \cdots \quad \leftarrow h_{T_x}$ $\quad$ Bi-RNN

$\rightarrow h_1$ $\qquad \rightarrow h_2$ $\qquad \rightarrow h_{T_x}$

$x_0$ $\qquad x_1$ $\qquad x_{T_x}$

# Thank **You!**

Questions?