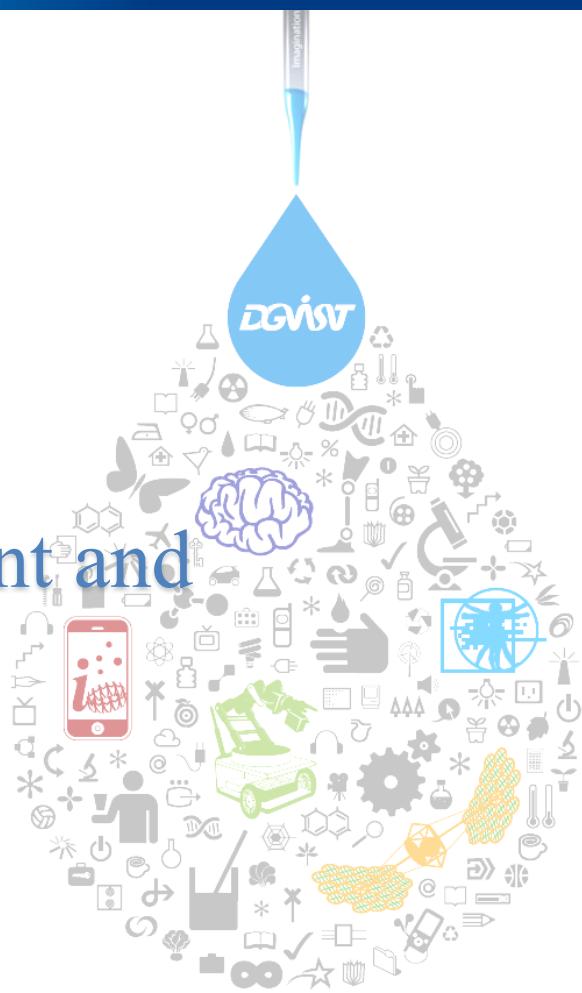




# Deep Learning Seminar

## CH 10. Sequence modeling: Recurrent and Recursive Nets – Part 1

2017-10-25  
Eunjeong Yi



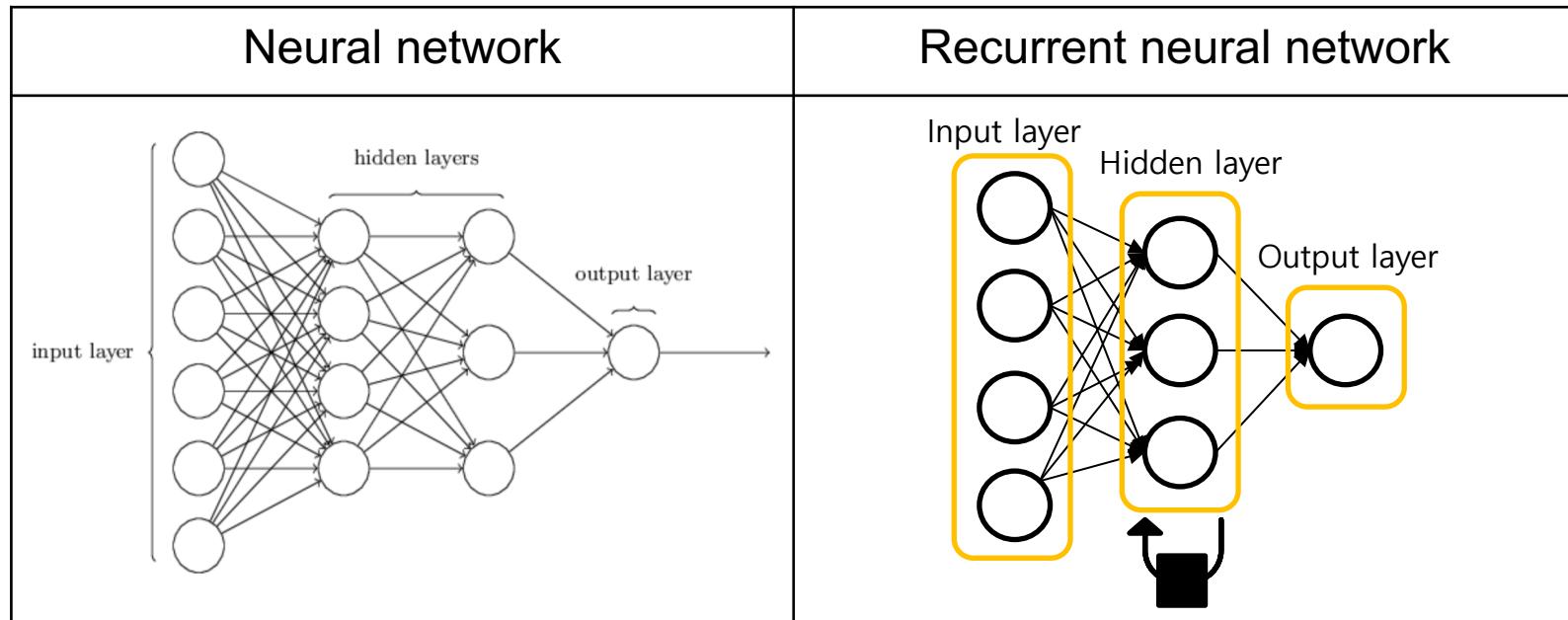
# Chapter 10. Sequence modeling: Recurrent and Recursive Nets

## ■ Part 1

- 10.1 Unfolding Computational Graphs**
- 10.2 Recurrent Neural Networks**
- 10.3 Bidirectional RNNs**
- 10.4 Encoder-Decoder Sequence-to-Sequence Architectures**
- 10.5 Deep Recurrent Networks**
- 10.6 Recursive Neural Networks**

# RNN(Recurrent Neural Network)

- Specialized for processing a sequence data
- Sharing parameters across different parts of a model
- Time step index refer to the position in sequence



Michael Nielsen. "CHAPTER 1 Using neural nets to recognize handwritten digits"  
<http://neuralnetworksanddeeplearning.com/chap1.html>

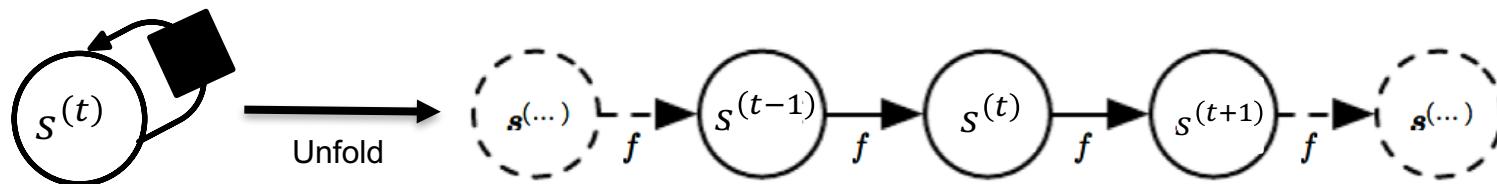
# Unfolding Computational Graphs

## ■ Computational graph

- a way to formalize the structure of a set of computations
- Ex)  $s^{(t)} = f(s^{(t-1)}; \theta)$

## ■ Unfolding equation

$$\begin{aligned}s^{(t)} &= f(s^{(t-1)}; \theta) \\&= f(f(s^{(t-2)}; \theta)) \\&= f(f(f(s^{(t-3)}; \theta))) \\&\vdots \\&= f(f(f(\cdots(f(s^{(1)}; \theta))))\end{aligned}$$

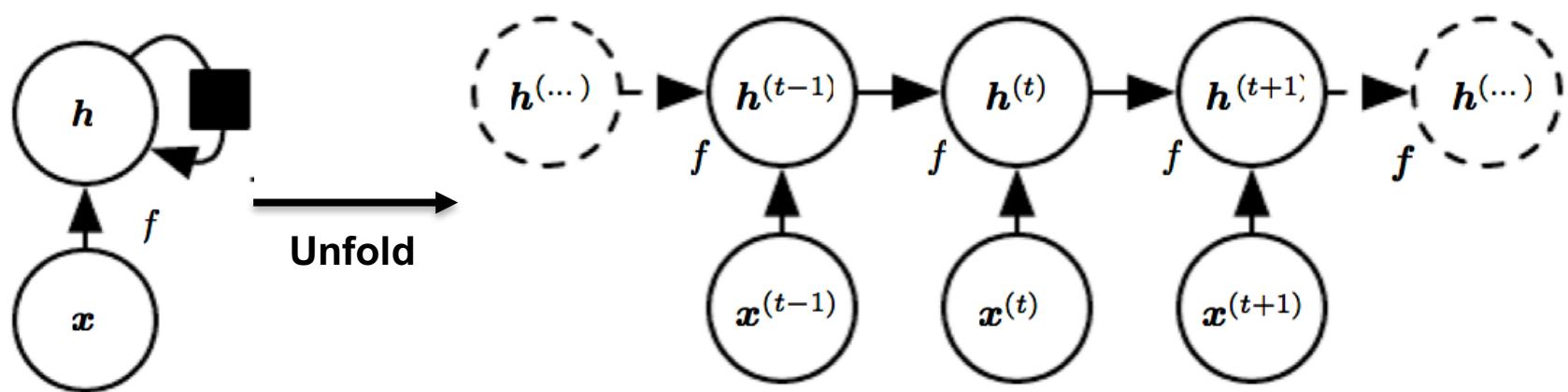


# Unfolding recurrent network

- $\mathbf{h}^{(t)} = \mathbf{g}^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)}) = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \theta)$

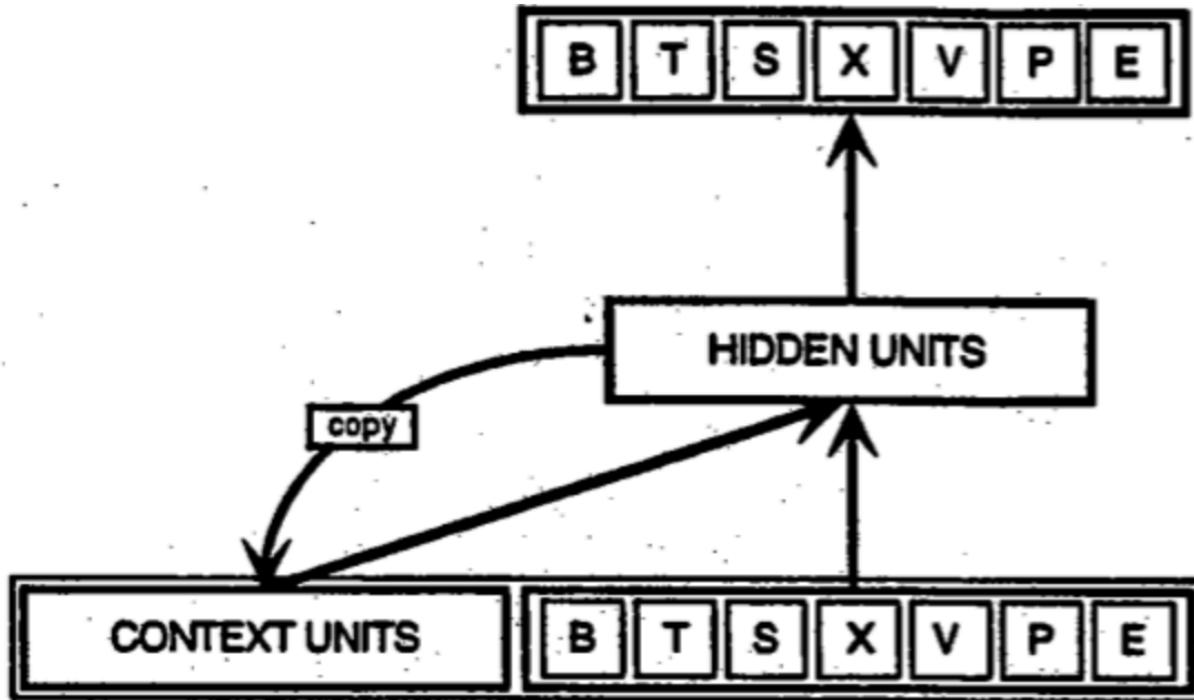
## ■ Advantage

- To use the transition function  $f$  with the same parameters at every time step.



# Recurrent Neural Networks

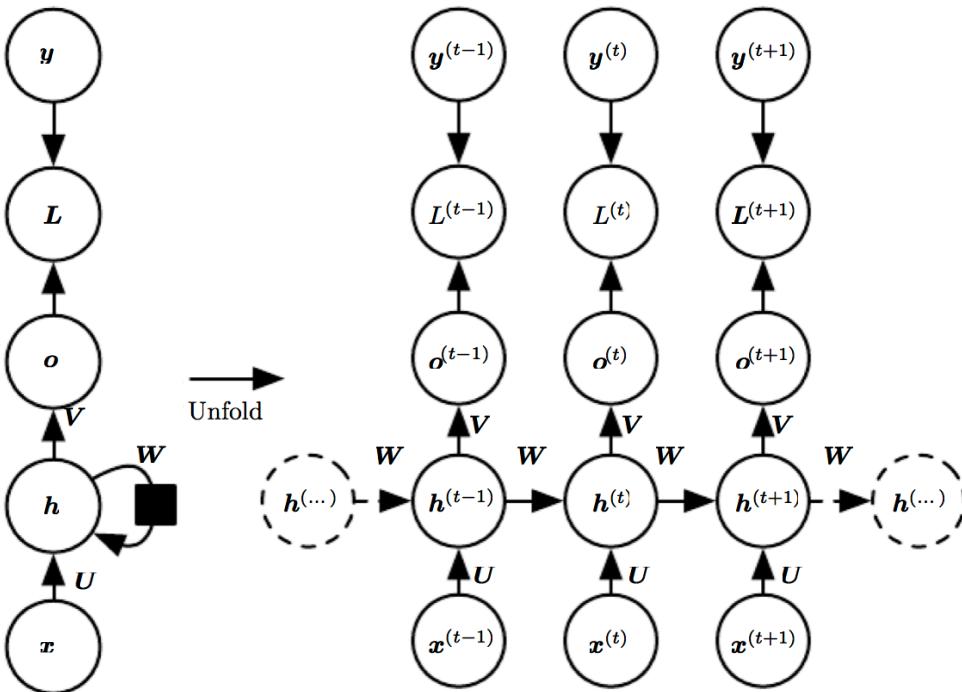
- Receiving input data and context unit in hidden layer
- Feedback structure that seems like memory



Elman, J. L. Finding structure in time. In *Cognitive Science*, 1990

# Recurrent Neural Networks

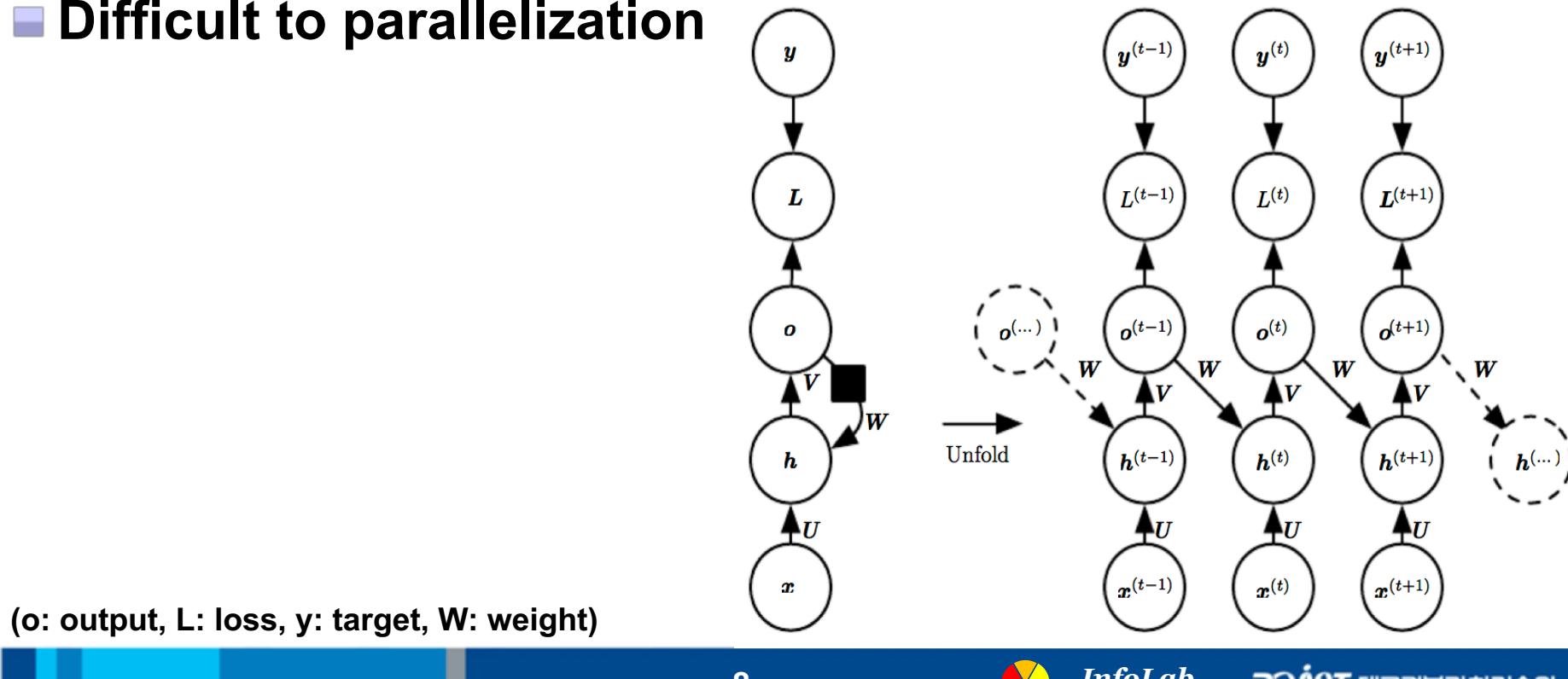
- $\mathbf{h}^{(t)} = f(\mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)})$
  - $\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)}$
  - $\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)})$
- } Forward propagation



( $\mathbf{o}$ : output,  $\mathbf{L}$ : loss,  $\mathbf{y}$ : target,  $\mathbf{W}$ : weight)

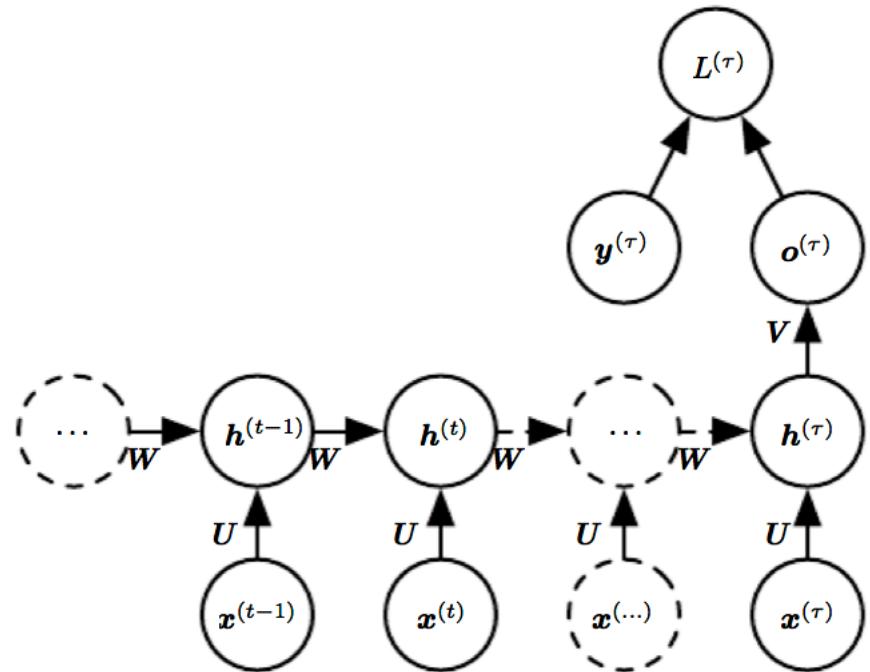
# Recurrent Neural Networks

- Produce an output at each time step and have recurrent connections only from the output at one time step to the hidden units at the next time step
- Lack important information from the past
- Difficult to parallelization



# Recurrent Neural Networks

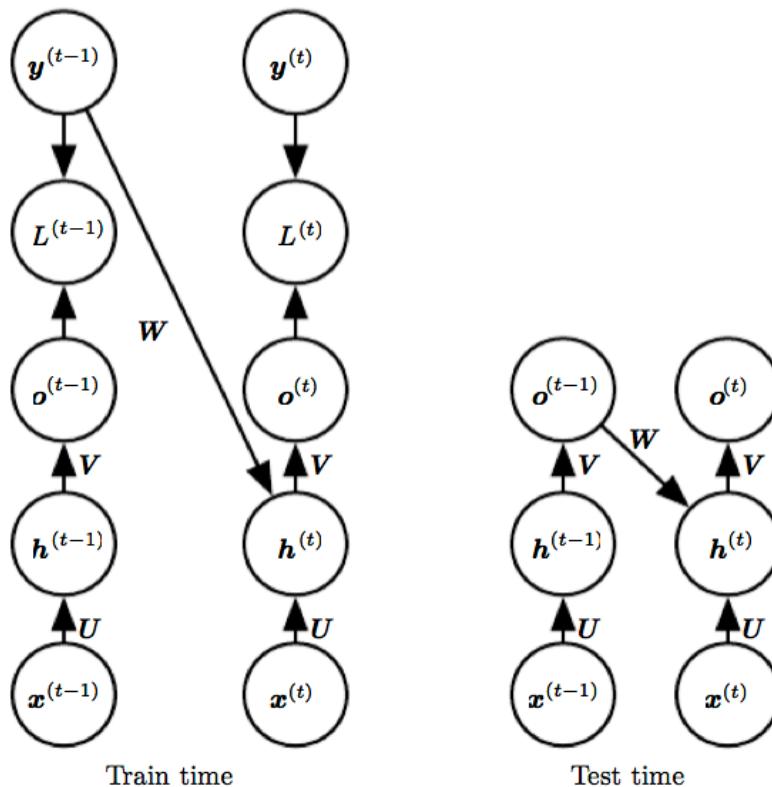
- Single output at the end of the sequence
- Used to summarize a sequence and produce a fixed-size representation
  - To use as input for further processing



(o: output, L: loss, y: target, W: weight)

# Teacher forcing and Networks with output recurrence

- During training the model receives the correct output  $y^{(t)}$  as input at time  $t + 1$



# Teacher forcing and Networks with output recurrence

## ■ Advantage

- Training can be parallelized with the gradient for each step computed in isolation

## ■ Disadvantage

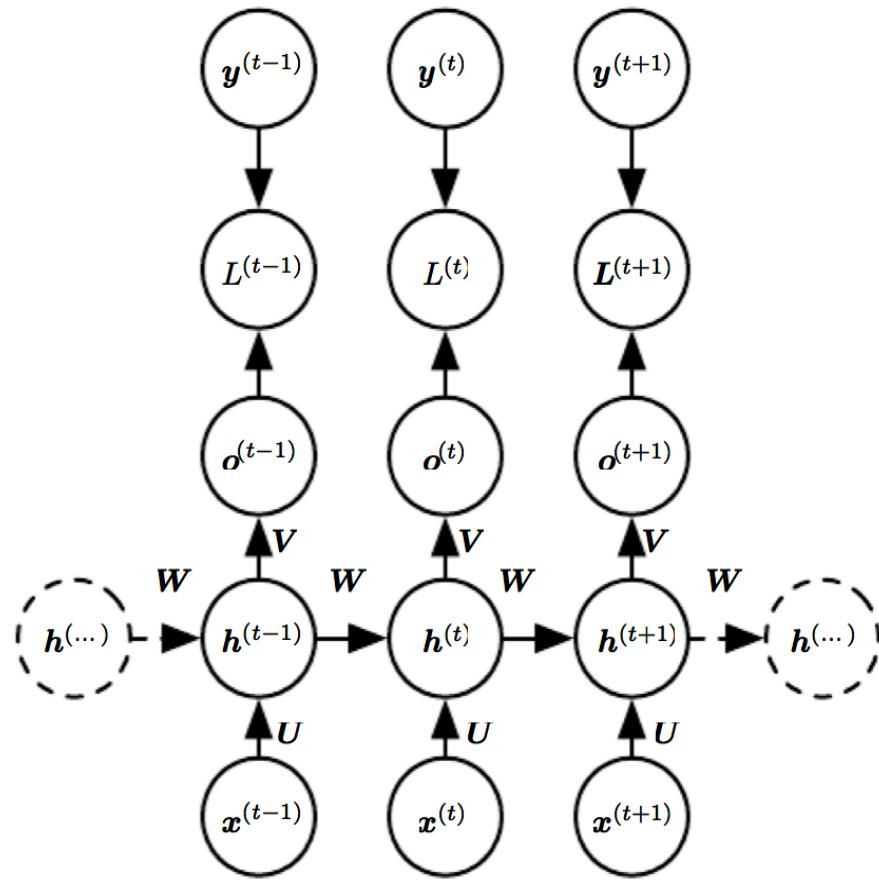
- The kind of inputs that the network sees during training could be quite different from the kind of inputs that it will see at test time

# Back propagation through time (BPTT)

- $\mathbf{h}^{(t)} = \tanh(\mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)})$
- $\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)}$
- $\mathcal{L}(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\})$

$U$ : weight from input to hidden layer  
 $W$ : hidden-to-hidden connection weight  
 $V$ : weight from hidden to output layer  
 $t$ : time index ( $1, 2, \dots, \tau$ )

$\mathbf{x}^{(t)}$ : input layer  
 $\mathbf{h}^{(t)}$ : hidden layer  
 $\mathbf{o}^{(t)}$ : output layer  
 $\mathcal{L}^{(t)}$ : loss function

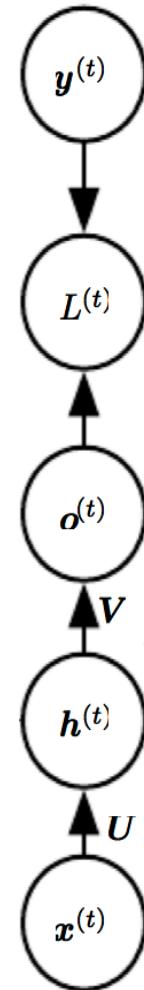


# Back propagation through time (BPTT)

## ■ In feedforward neural network

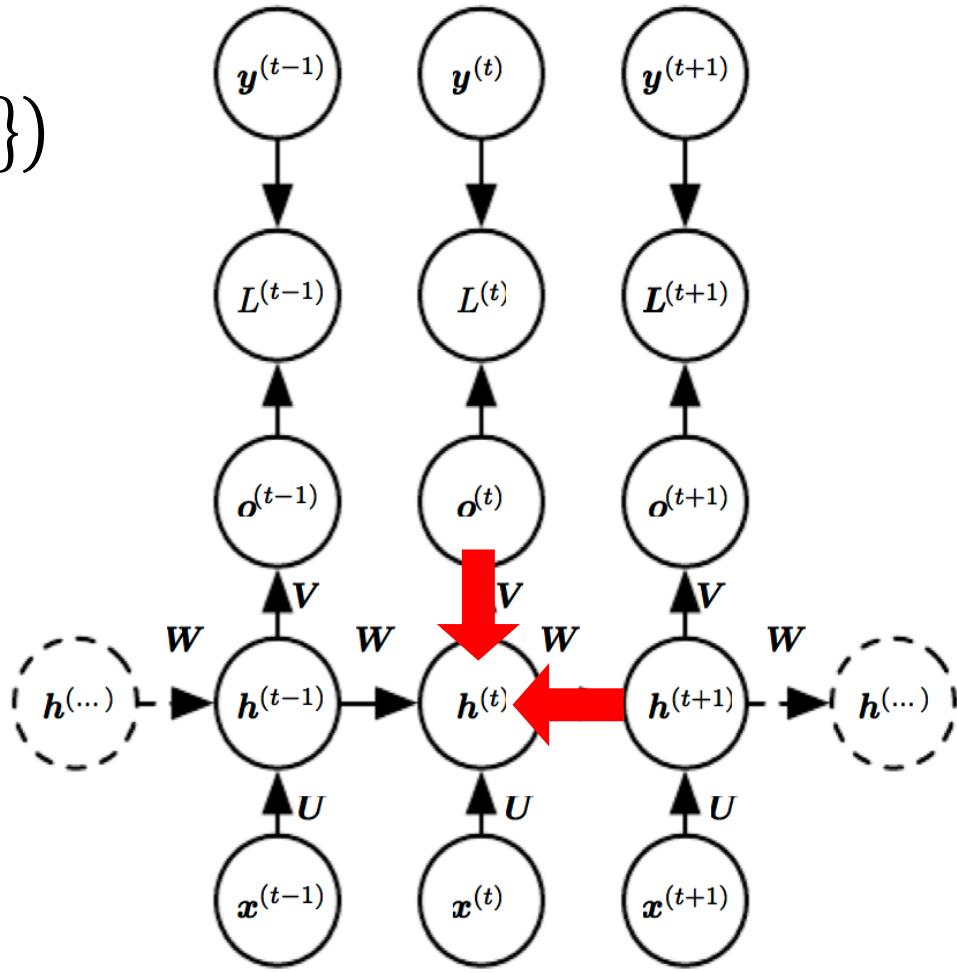
- $\mathbf{h}^{(t)} = \tanh(\mathbf{U}\mathbf{x}^{(t)})$
- $\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)}$
- $L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\})$

- $\nabla_{\mathbf{h}^{(t)}} L = \frac{\partial L}{\partial \mathbf{o}^{(t)}} \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}}$
- $\nabla_{\mathbf{V}} L = \frac{\partial L}{\partial \mathbf{o}^{(t)}} \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{V}}$
- $\nabla_{\mathbf{U}} L = \frac{\partial L}{\partial \mathbf{o}^{(t)}} \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{U}}$



# Back propagation through time (BPTT)

- $\mathbf{h}^{(t)} = \tanh(\mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)})$
- $\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)}$
- $\mathcal{L}(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\})$
- $\nabla_{\mathbf{h}^{(t)}} \mathcal{L}$



# Back propagation through time (BPTT)

- $\nabla_{\mathbf{h}^{(\tau)}} \mathbf{L} = \frac{\partial \mathbf{L}}{\partial \mathbf{h}^{(\tau)}} = \frac{\partial \mathbf{L}}{\partial \mathbf{o}^{(\tau)}} \frac{\partial \mathbf{o}^{(\tau)}}{\partial \mathbf{h}^{(\tau)}} = \mathbf{V}^T \frac{\partial \mathbf{L}}{\partial \mathbf{o}^{(\tau)}} = \mathbf{V}^T \nabla_{\mathbf{o}^{(\tau)}} \mathbf{L}$

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

$$\begin{aligned}\frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} &= \frac{\partial \tanh(\mathbf{a}^{(t+1)})}{\partial \mathbf{h}^{(t)}} = \text{diag} \left( 1 - \tanh^2(\mathbf{a}^{(t+1)}) \right) \frac{\partial \mathbf{a}^{(t+1)}}{\partial \mathbf{h}^{(t)}} = \mathbf{W} \text{diag} \left( 1 - (\mathbf{h}^{(t+1)})^2 \right) \\ &= \left( \mathbf{W} \text{diag} \left( 1 - (\mathbf{h}^{(t+1)})^2 \right) \right)^T \nabla_{\mathbf{h}^{(t)}} \mathbf{L} + \mathbf{V}^T \nabla_{\mathbf{o}^{(t)}} \mathbf{L}\end{aligned}$$

$$= \mathbf{W}^T \text{diag} \left( 1 - (\mathbf{h}^{(t+1)})^2 \right) \nabla_{\mathbf{h}^{(t)}} \mathbf{L} + \mathbf{V}^T \nabla_{\mathbf{o}^{(t)}} \mathbf{L}$$

# Back propagation through time (BPTT)

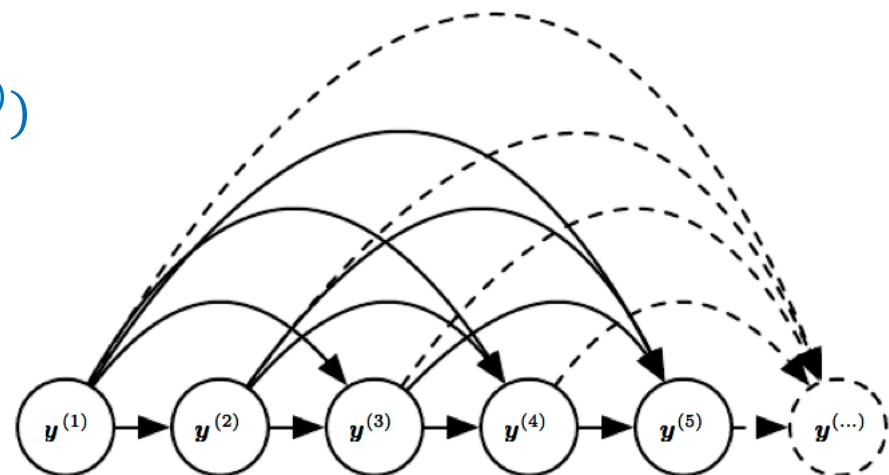
- $\nabla_V L = \sum_t \sum_i \left( \frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_V o_i^{(t)} = \sum_t (\nabla_{o^{(t)}} L) h^{(t)}^T$
- $\nabla_W L = \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{W^{(t)}} h_i^{(t)}$ 
$$= \sum_t \text{diag} \left( 1 - (h^{(t)})^2 \right) (\nabla_{h^{(t)}} L) h^{(t-1)}^T$$
- $\nabla_U L = \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{U^{(t)}} h_i^{(t)}$ 
$$= \sum_t \text{diag} \left( 1 - (h^{(t)})^2 \right) (\nabla_{h^{(t)}} L) x^{(t)}^T$$

# Recurrent Network as Directed Graphical Models

## ■ Fully connected graphical model

- Inefficient parametrization
- ex)  $h^{(t)} = g^{(t)}(x^{(t)}, x^{(t-1)}, \dots, x^{(1)})$

## ■ The edge indicate which variables depend directly on other variables

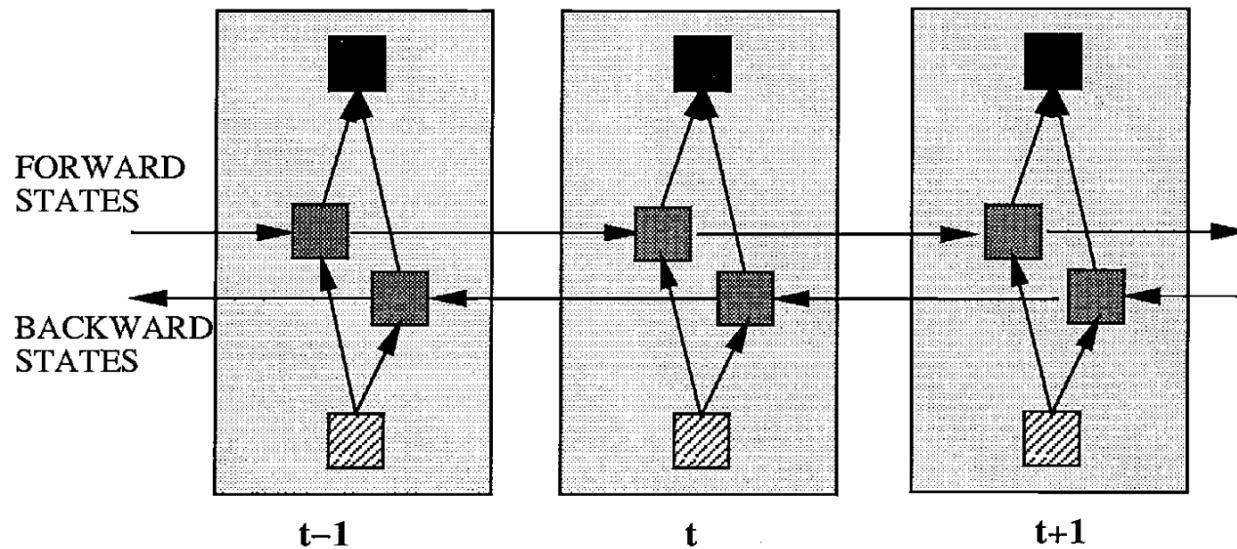


## ■ Efficient parametrization of the joint distribution

$$\text{ex) } h^{(t)} = g^{(t)}(x^{(t)}, x^{(t-1)}, \dots, x^{(1)}) = f(h^{(t-1)}, x^{(t)}; \theta)$$

# Bidirectional RNNs

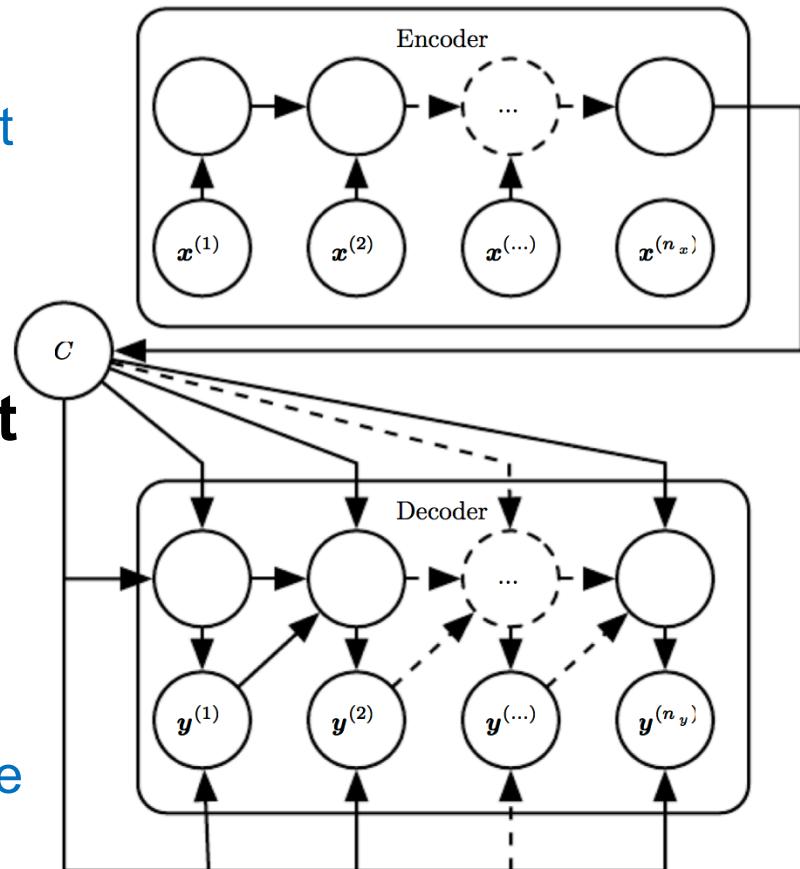
- maintains two hidden layers, one for the forward states and the other for the backward states
- consumes twice as much memory space for its weight and bias parameters
- Output units with a relevant summary of the past in the forward states and the future in the backward states



<http://solarisailab.com/archives/1515>

# Encoder-Decoder Sequence-to-Sequence Architecture

- To map an input sequence to an output sequence which is not necessarily of the same length
- Encoder (or input RNN)
  - process input sequence
  - Emit context C that summarize input
- Decoder (or output RNN)
  - Conditioned on fixed-length vector to generate the output sequence
- The length of input and output sequence can vary from each other
- Limitation
  - Dimension of context C can be too small to summarize a long sequence



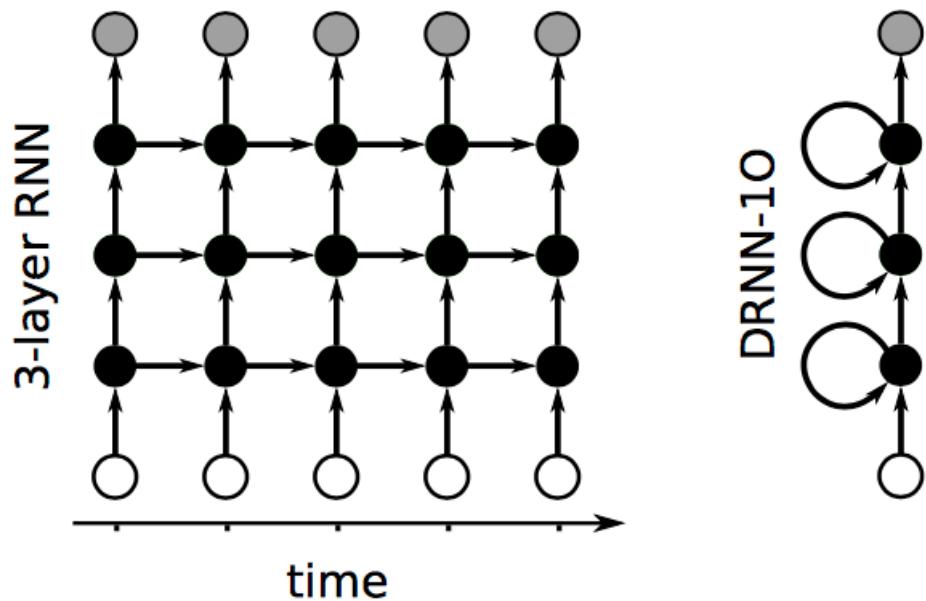
# Deep Recurrent Networks

## ■ 3 parts

- From the input to the hidden state
- From the previous hidden state to the next hidden state
- From the hidden state to the output

## ■ Increase the number of nonlinear steps the gradient at back propagation

- difficult to train the model to capture long-term dependency



Hermans, M., & Schrauwen, B. (2013). Training and analysing deep recurrent neural networks. In *Advances in Neural Information Processing Systems* (pp. 190-198).

# Recursive Neural Networks

- Another generalization of recurrent neural network
- Applied to processing data structures as input to neural nets in natural language processing as well as computer vision
- Advantage
  - For a sequence of same length  $\tau$ , the depth can be reduced from  $\tau$  to  $O(\log \tau)$
- How to best structure the tree

