



RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNS,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

## Lecture 9 Recurrent Neural Networks

Shiwei Lan<sup>1</sup>

<sup>1</sup>School of Mathematical and Statistical Sciences  
Arizona State University

STP598 Machine Learning and Deep Learning  
Fall 2021



# Table of Contents

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNS,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

## 1 Introduction

## 2 Recurrent Neural Networks (RNN)

## 3 Variants: Bidirectional RNNs, Deep RNNS, Recursive NNs

## 4 Long-Time Memory: LSTM, GRU

## 5 Software Implementation



# Recurrent Neural Networks

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

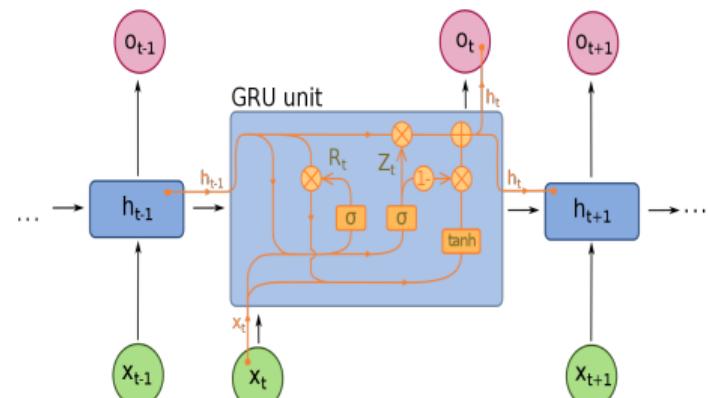
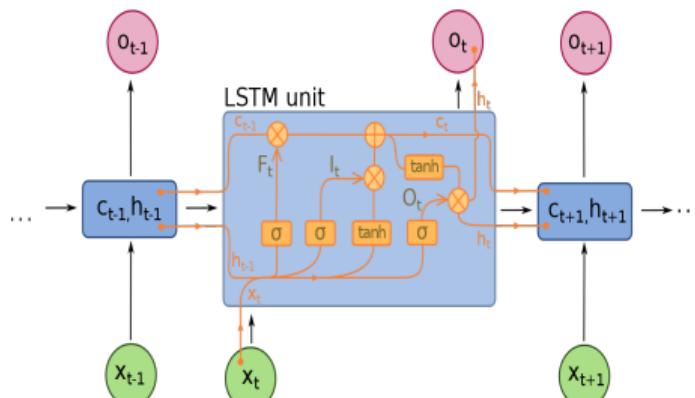
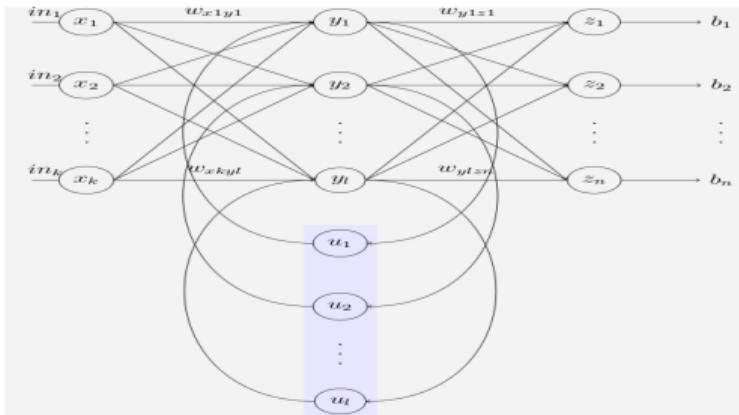
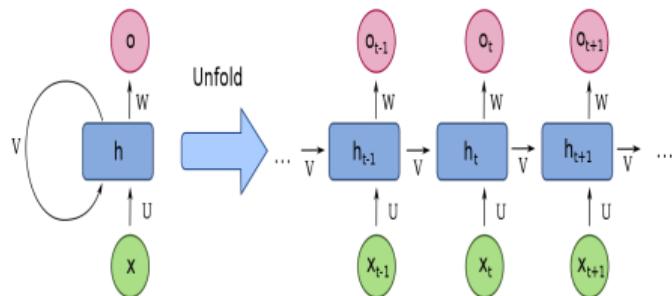
- Recurrent Neural Networks (RNN) are a family of neural networks for processing sequential data, e.g. time-series, a sequence of images, etc.
- CNN is specialized for processing a grid of values (e.g. images); RNN is specialized for processing a sequence of values (e.g. time series).
- The idea going from MLP to RNN dates back to 1980s: *sharing parameters across different parts of a model.*
- What about CNN across a 1-D temporal sequence? Shallow.
- RNNs designed to retain long time memory, e.g. LSTM, GRU, are particularly useful and successful.

# Recurrent Neural Networks

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation



# Table of Contents

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNS,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

## 1 Introduction

## 2 Recurrent Neural Networks (RNN)

## 3 Variants: Bidirectional RNNs, Deep RNNS, Recursive NNs

## 4 Long-Time Memory: LSTM, GRU

## 5 Software Implementation

RNN

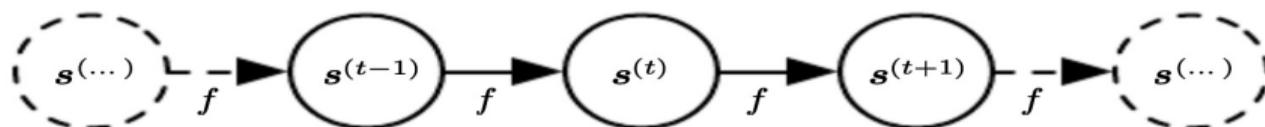
S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

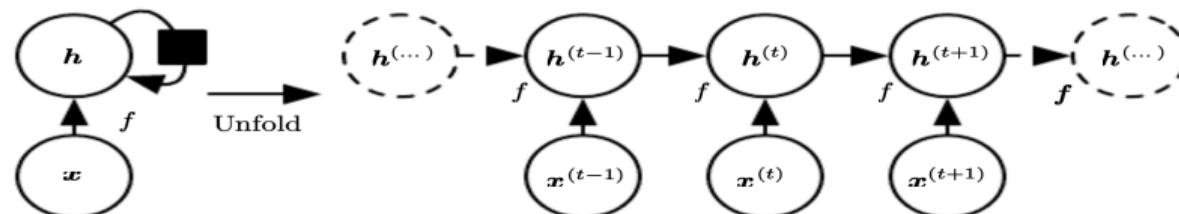
- Consider the classical form of a dynamical system:

$$s^{(t)} = f(s^{(t-1)}; \theta) \quad (1)$$



- Sometimes, the dynamics may be driven by external signal  $x^{(t)}$

$$h^{(t)} = f(h^{(t-1)}, x^{(t)}; \theta) \quad (2)$$





# Unfolding

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

- We can represent the unfolded recurrence after  $t$  steps with a function  $g^{(t)}$ :

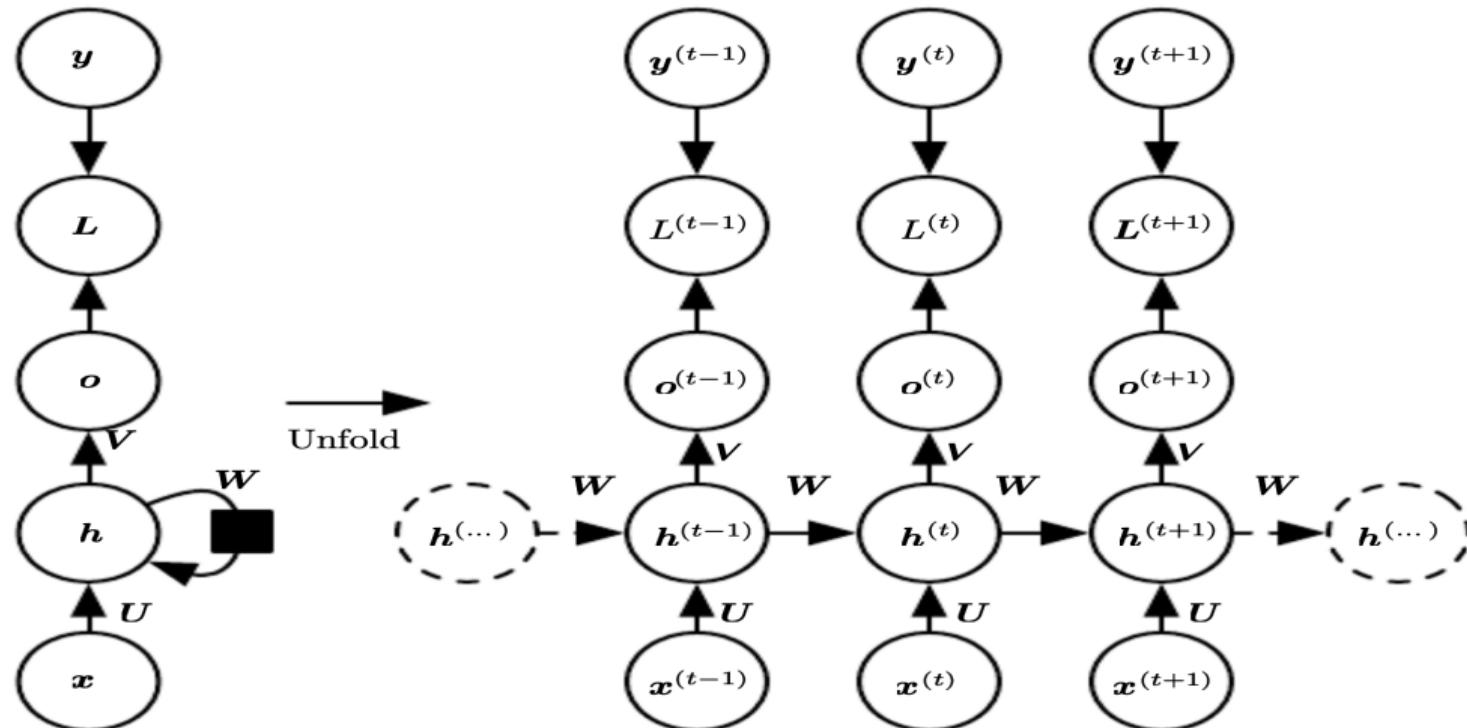
$$h^{(t)} = g^{(t)}(x^{(t)}, x^{(t-1)}, \dots, x^{(1)}) = f(h^{(t-1)}, x^{(t)}; \theta) \quad (3)$$

- The unfolding process thus introduces two major advantages:
  - ① Regardless of the sequence length, the learned model always has the same input size, because it is specified in terms of transition from one state to another state, rather than specified in terms of a variable-length history of states.
  - ② It is possible to use the *same* transition function  $f$  with the same parameters at every time step.

RNN

S.Lan

Introduction

Recurrent  
Neural Networks (RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation



# Recurrent Neural Networks

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

- Assume the output is discrete, as if the RNN is used to predict words or characters.
- Let  $\hat{y}$  be a vector of normalized probabilities over the output.
- Forward propagation begins with a specification of the initial state  $h^{(0)}$ . Then from  $t = 1$  to  $t = \tau$ , we have

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)} \quad (4)$$

$$h^{(t)} = \tanh(a^{(t)}) \quad (5)$$

$$o^{(t)} = c + Vh^{(t)} \quad (6)$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)}) \quad (7)$$

where the parameters are the bias vectors  $b$  and  $c$  along with the weight matrices  $U$ ,  $V$  and  $W$ , respectively for input-to-hidden, hidden-to-output and hidden-to-hidden connections.



# Recurrent Neural Networks

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

- The total loss for a given sequence of  $x$  values paired with a sequence of  $y$  values would be the sum of the losses over all the time steps.
- Let  $L^{(t)}$  be the negative log-likelihood of  $y^{(t)}$  given  $x^{(1)}, \dots, x^{(t)}$ .

$$\begin{aligned} L(\{x^{(1)}, \dots, x^{(\tau)}\}, \{y^{(1)}, \dots, y^{(\tau)}\}) &= \sum_t L^{(t)} \\ &= - \sum_t \log p_{\text{model}}(y^{(t)} | \{x^{(1)}, \dots, x^{(t)}\}) \end{aligned} \tag{8}$$

where  $p_{\text{model}}(y^{(t)} | \{x^{(1)}, \dots, x^{(t)}\})$  is given by reading the entry for  $y^{(t)}$  from the model's output vector  $\hat{y}^{(t)}$ .

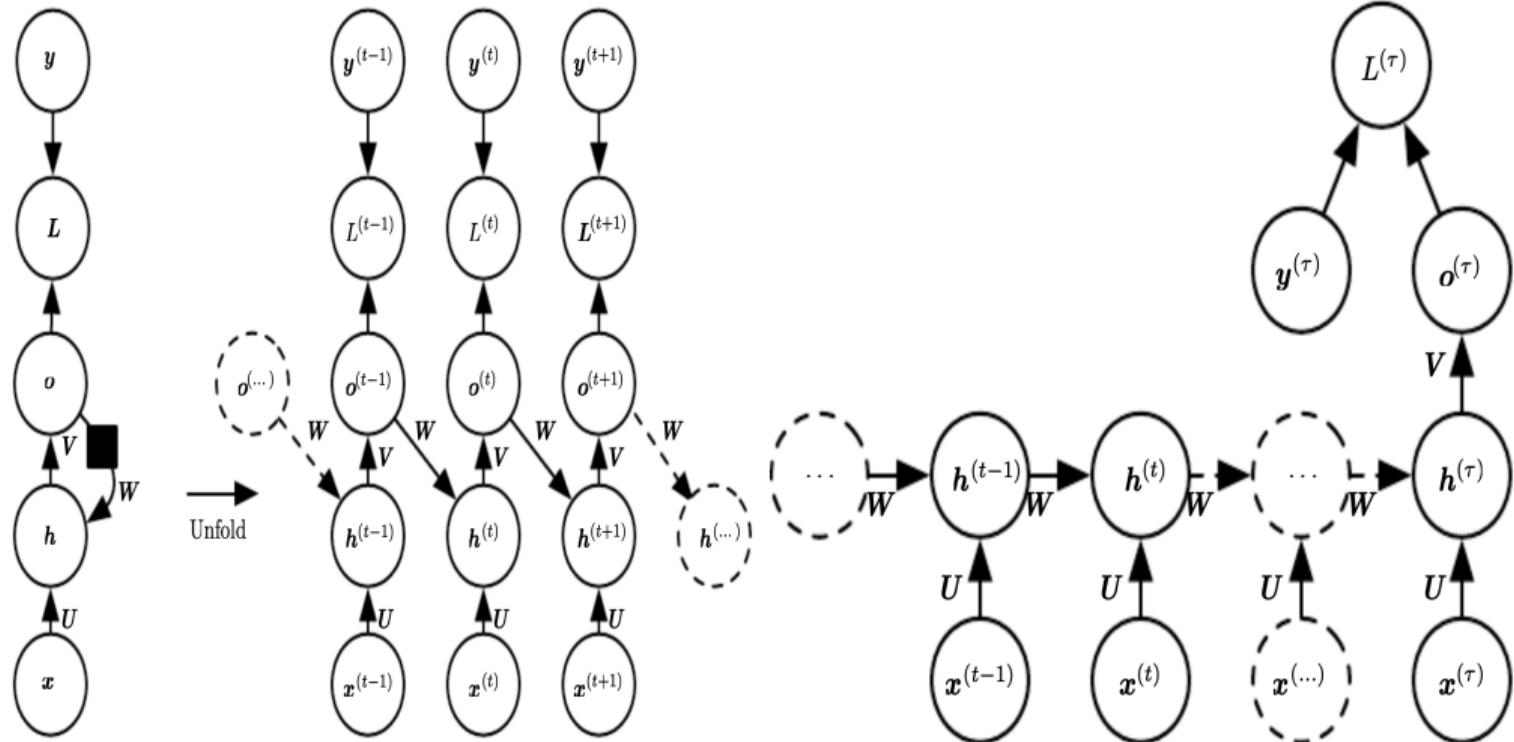
- The runtime is  $\mathcal{O}(\tau)$ ; the memory cost is also  $\mathcal{O}(\tau)$ .
- The back-propagation algorithm applied to the unrolled graph is called **back-propagation through time** (BPTT).

# Alternative RNNs

RNN

S.Lan

Introduction

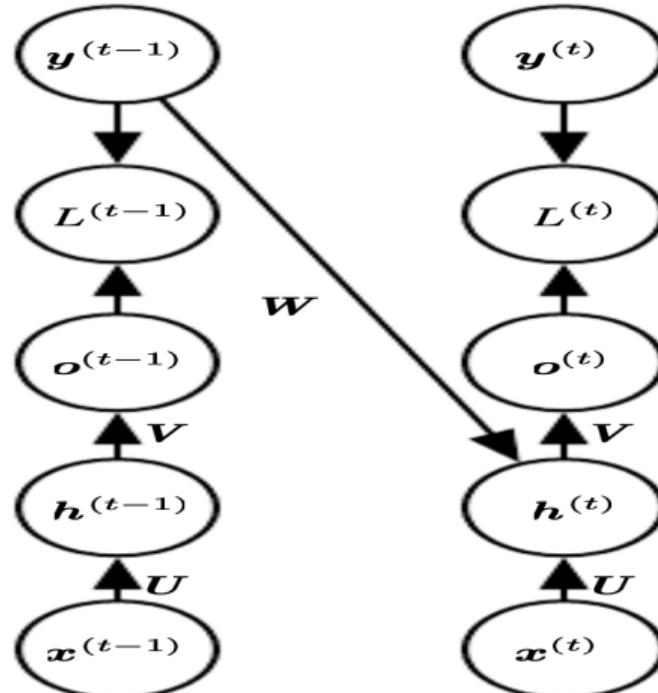
Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

# Teacher Forcing

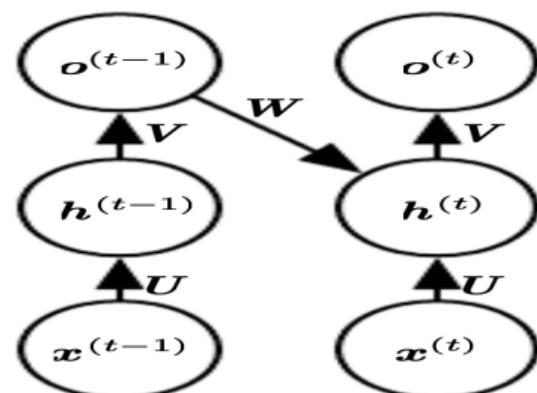
RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

Train time

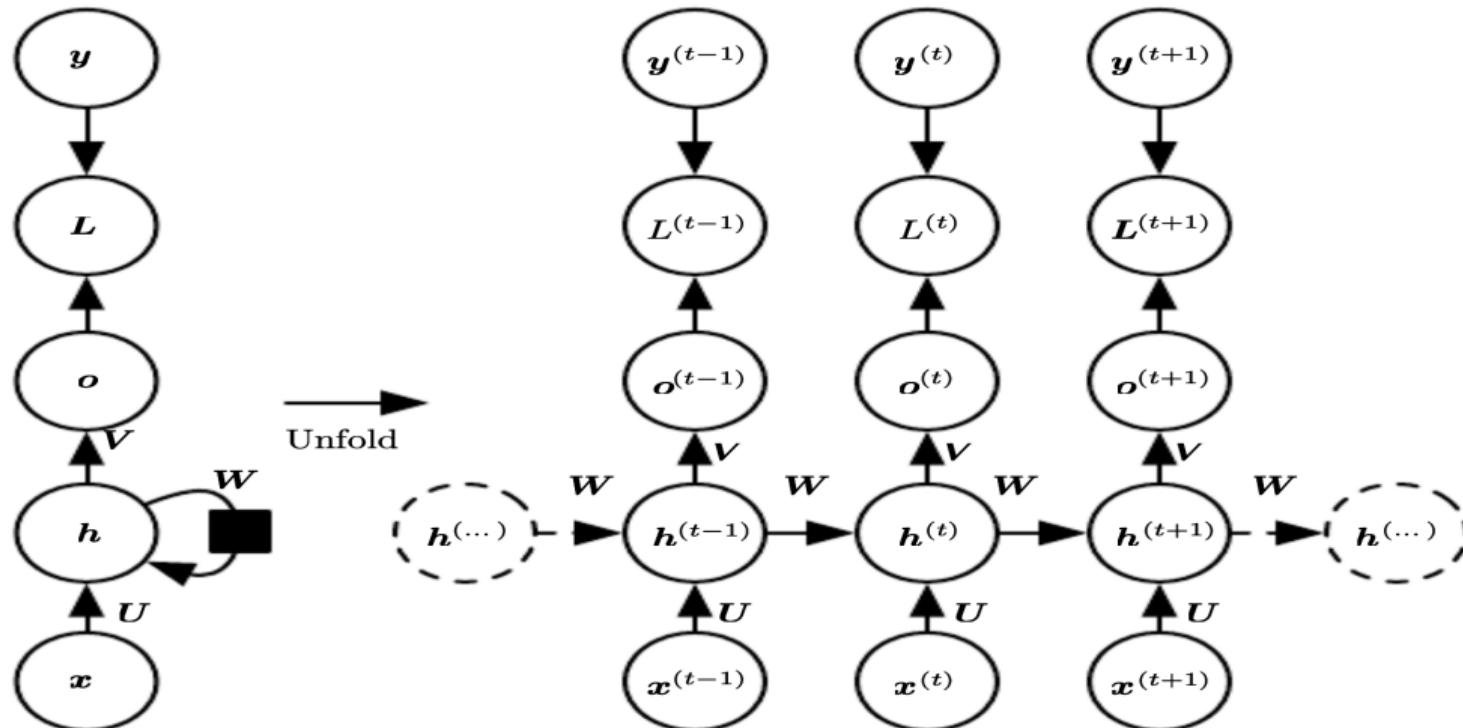


Test time

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation



# BPTT in RNN

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

$$\frac{\partial L}{\partial L^{(t)}} = 1 \quad (9)$$

$$(\nabla_{o^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i,y^{(t)}} \quad (10)$$

$$\nabla_{h^{(\tau)}} L = V^T \nabla_{o^{(t)}} L \quad (11)$$

$$\nabla_{h^{(t)}} L = \left( \frac{\partial h^{(t+1)}}{h^{(t)}} \right)^T (\nabla_{h^{(t+1)}} L) + \left( \frac{\partial o^{(t)}}{h^{(t)}} \right)^T (\nabla_{o^{(t)}} L) \quad (12)$$

$$= W^T (\nabla_{h^{(t+1)}} L) \text{diag} \left( 1 - (h^{(t+1)})^2 \right) + V^T (\nabla_{o^{(t)}} L) \quad (13)$$



# BPTT in RNN

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

$$\nabla_{\mathbf{c}} L = \sum_t \left( \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^\top \nabla_{\mathbf{o}^{(t)}} L = \sum_t \nabla_{\mathbf{o}^{(t)}} L \quad (10.22)$$

$$\nabla_{\mathbf{b}} L = \sum_t \left( \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^\top \nabla_{\mathbf{h}^{(t)}} L = \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) \nabla_{\mathbf{h}^{(t)}} L \quad (10.23)$$

$$\nabla_{\mathbf{v}} L = \sum_t \sum_i \left( \frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\mathbf{v}} o_i^{(t)} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t)^\top} \quad (10.24)$$

$$\begin{aligned} \nabla_{\mathbf{w}} L &= \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{w}^{(t)}} h_i^{(t)} \\ &= \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)^\top} \end{aligned} \quad (10.25) \quad (10.26)$$

$$\nabla_{\mathbf{u}} L = \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{u}^{(t)}} h_i^{(t)} \quad (10.27)$$

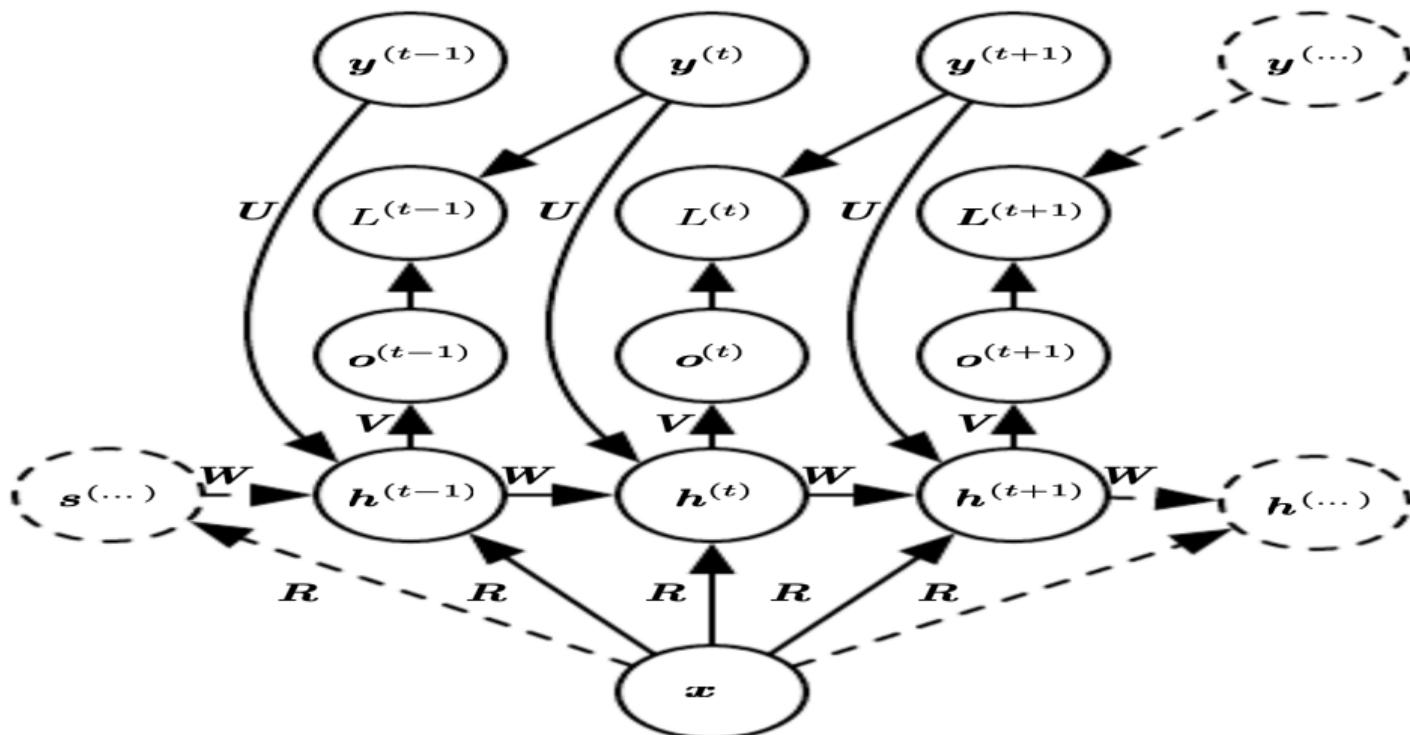
$$= \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{x}^{(t)^\top} \quad (10.28)$$

# What if input is a fixed-size vector?

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation



# Table of Contents

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNS,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

## 1 Introduction

## 2 Recurrent Neural Networks (RNN)

## 3 Variants: Bidirectional RNNs, Deep RNNS, Recursive NNs

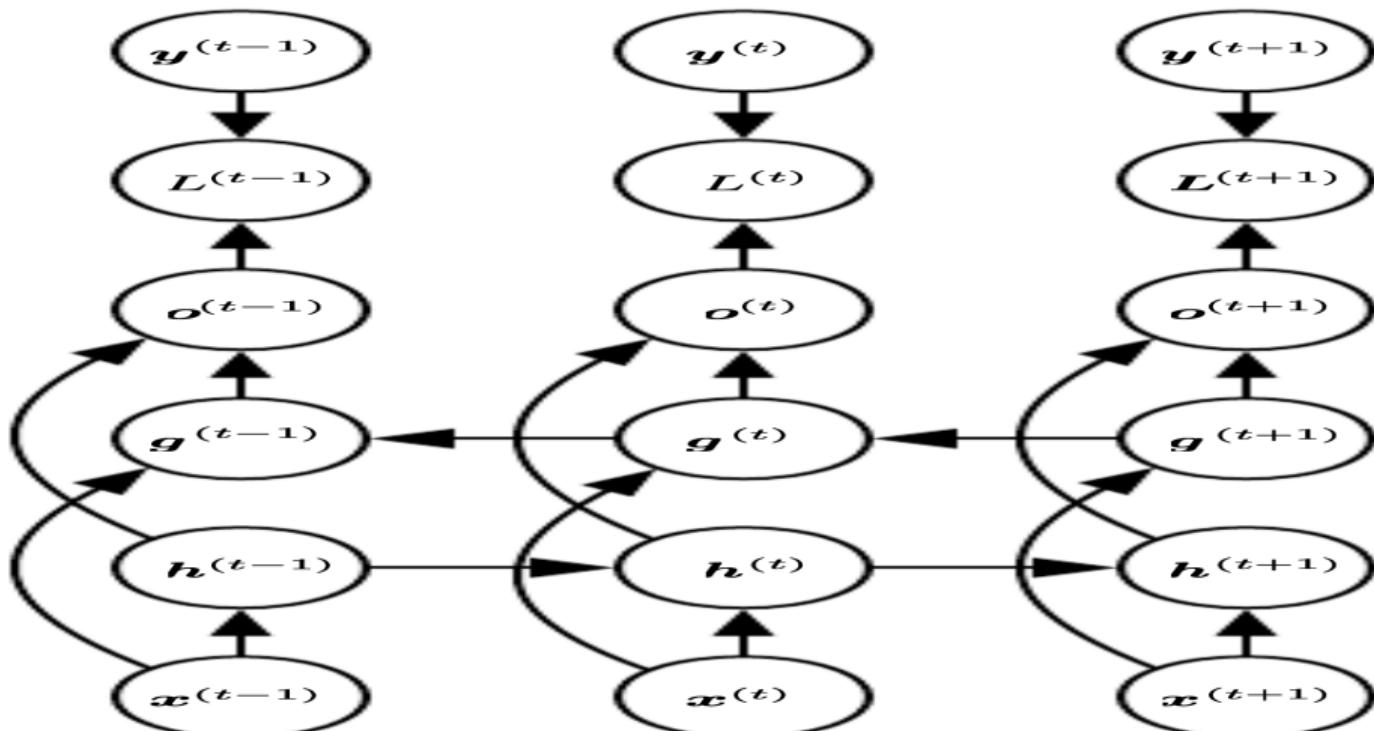
## 4 Long-Time Memory: LSTM, GRU

## 5 Software Implementation

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

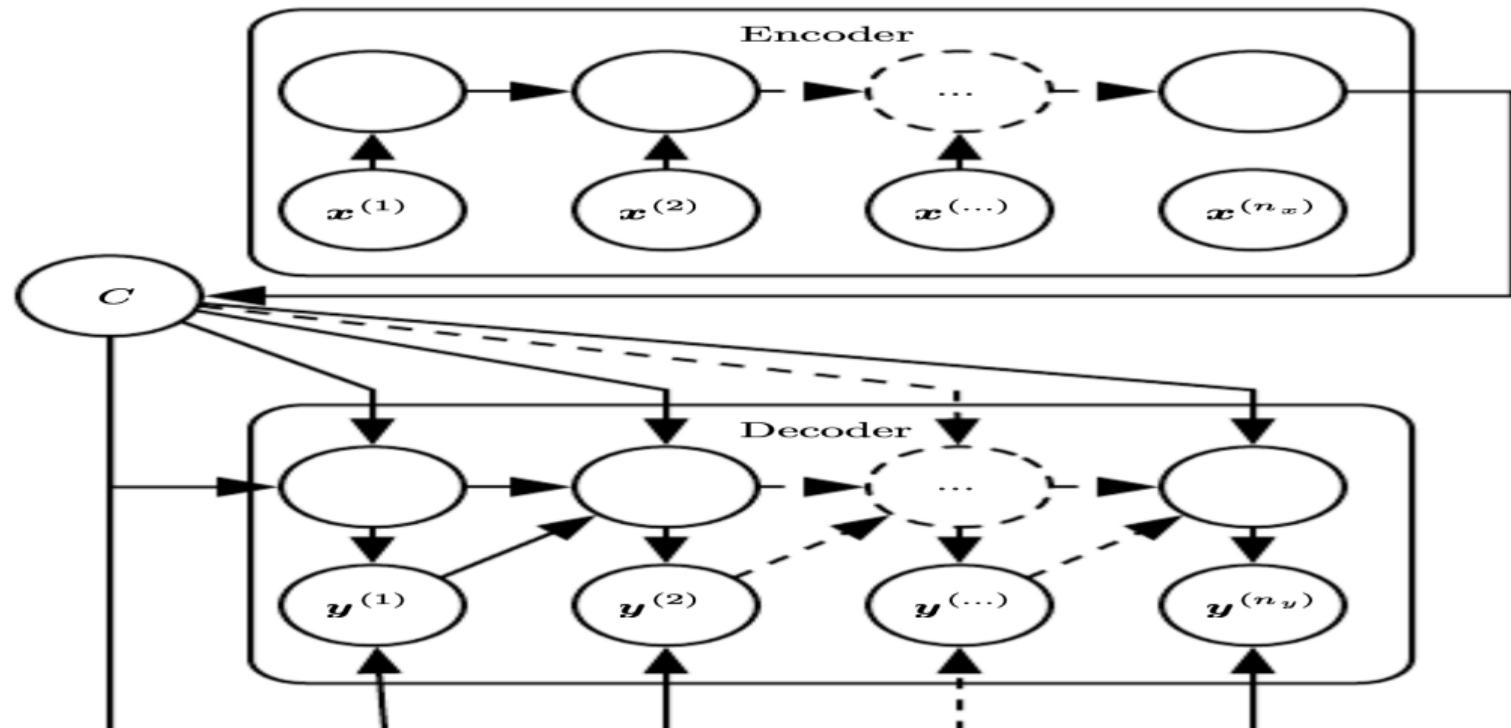
- In many applications we want to output a prediction of  $y^{(t)}$  which may depend on *the whole input sequence*.
- Bidirectional RNNs (Schuster and Paliwal, 1997) have been extremely successful (Graves, 2012) in applications such as handwriting recognition (Graves et al., 2008; Graves and Schmidhuber, 2009), speech recognition (Graves and Schmidhuber, 2005; Graves et al., 2013).
- Bidirectional RNNs compute a representation that depends on *both the past (through  $g^{(t)}$ ) and the future (through  $h^{(t)}$ )* without having to specify a fixed-size window around  $t$ .
- This idea can be naturally extended to 2-dimensional input (images) by having four RNNs in 4 directions: up, down, left, right.

# What if input and output sequences have different length?

RNN

S.Lan

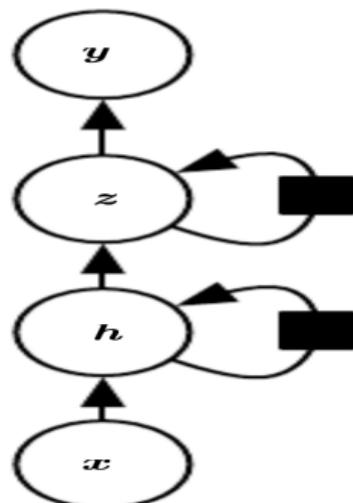
Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

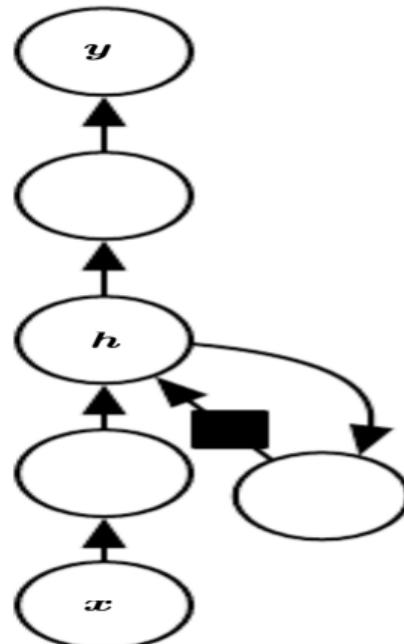
RNN

S.Lan

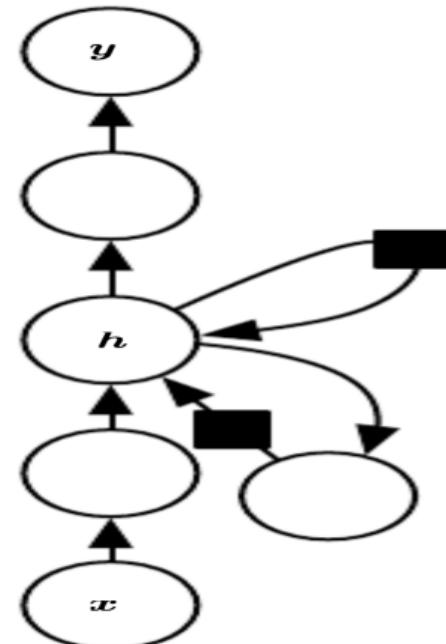
Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

(a)



(b)

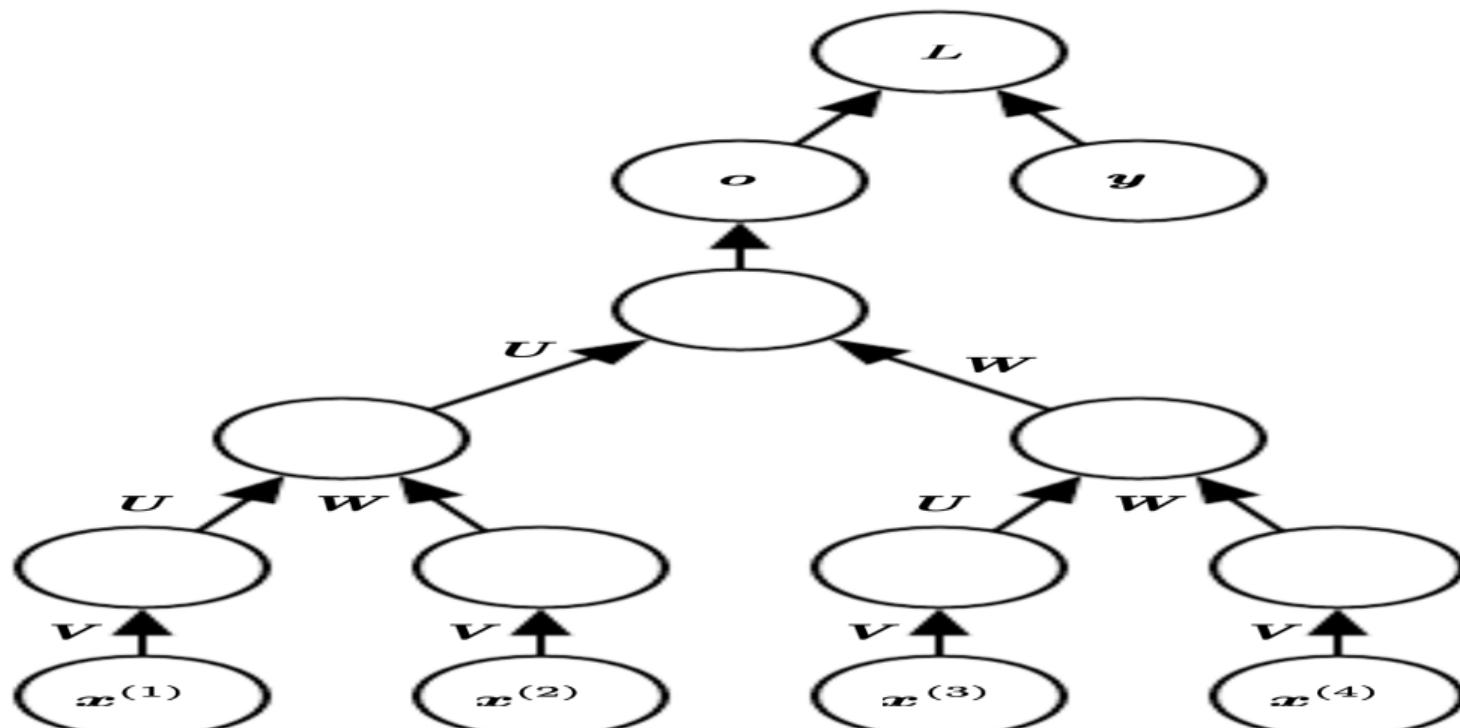


(c)

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation



# Table of Contents

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNS,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

## 1 Introduction

## 2 Recurrent Neural Networks (RNN)

## 3 Variants: Bidirectional RNNs, Deep RNNS, Recursive NNs

## 4 Long-Time Memory: LSTM, GRU

## 5 Software Implementation

# The challenge of long-term memory

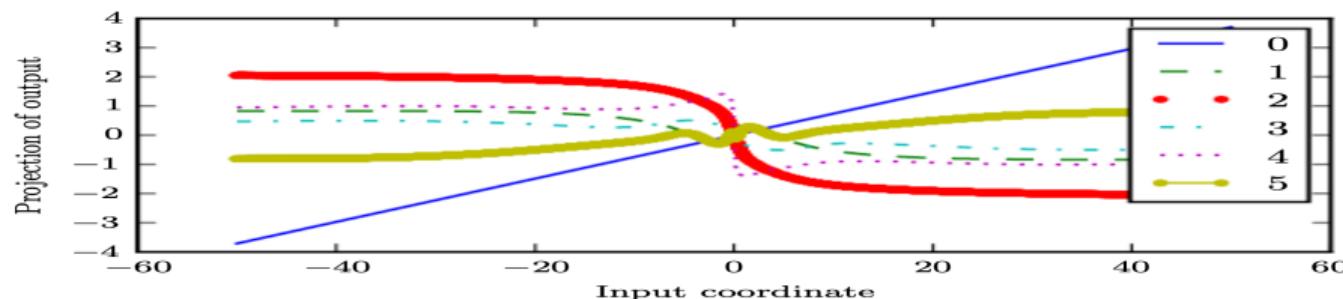
RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

- The basic problem is that gradients propagated over many stages tend to either vanish (most of the time) or explode (rarely, but with much damage to the optimization).
- The difficulty with long-term dependencies arises from the exponentially smaller weights given to long-term interactions (involving the multiplication of many Jacobians) compared to short-term ones.



- A heuristic explanation:

$$h^{(t)} = W^T h^{(t-1)} = Q \Lambda Q^T h^{(t-1)} = \dots = Q \Lambda^t Q^T h^{(0)} \quad (14)$$

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

- **Echo State Networks**(ESNs, Jaeger and Haas, 2004; Jaeger, 2007b): set the recurrent weights and *earn only the output weights*.
- **Leaky Units**(Mozer, 1992; El Hihi and Bengio, 1996): update running average  $\mu^{(t)} = \alpha\mu^{(t-1)} + (1 - \alpha)v^{(t)}$ . Hidden units with linear self-connections behave similarly.
- **Skip/remove connections**(Lin et al., 1996; El Hihi and Bengio, 1996): longer delays, multiple time-scales.

# Long Short-Term Memory (LSTM)

RNN

S.Lan

Introduction

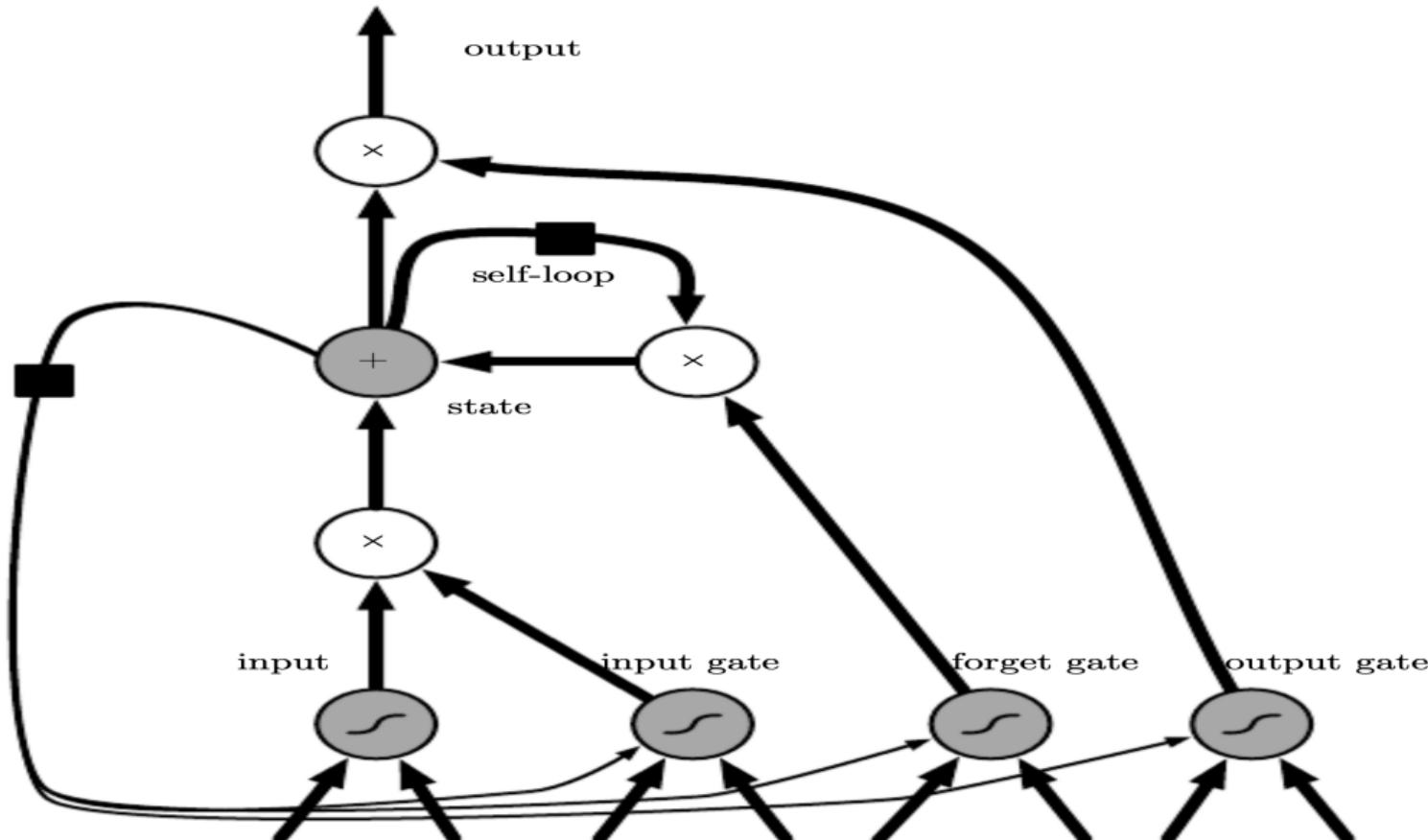
Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

- The clever idea of introducing self-loops to produce paths where the gradient can flow for long durations is a core contribution of the initial long short-term memory (LSTM) model (Hochreiter and Schmidhuber, 1997)
- A crucial addition has been to make the weight on this self-loop conditioned on the context, rather than fixed (Gers et al., 2000).
- By making the weight of this self-loop gated (controlled by another hidden unit), the time scale of integration can be changed dynamically.
- The LSTM has been found extremely successful in many applications, such as unconstrained handwriting recognition (Graves et al., 2009), speech recognition (Graves et al., 2013; Graves and Jaitly, 2014), handwriting generation (Graves, 2013), machine translation (Sutskever et al., 2014), image captioning (Kiros et al., 2014b; Vinyals et al., 2014b; Xu et al., 2015) and parsing (Vinyals et al., 2014a).

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

- The self-loop weight (or the associated time constant) is controlled by a **forget gate** unit  $f_i^{(t)}$  (for time step  $t$  and cell  $i$ )

$$f_i^{(t)} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} \right) \quad (15)$$

- The LSTM cell internal state is updated as follows

$$s_i^{(t)} = f_i^{(t)} s_i^{(t-1)} + g_i^{(t)} \sigma \left( b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)} \right) \quad (16)$$

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

- The **external input gate** unit  $g_i^{(t)}$  is computed similarly

$$g_i^{(t)} = \sigma \left( b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} \right) \quad (17)$$

- The output  $h_i^{(t)}$  of the LSTM cell can also be shut off, via the **output gate**  $q_i^{(t)}$

$$h_i^{(t)} = \tanh(s_i^{(t)}) q_i^{(t)} \quad (18)$$

$$q_i^{(t)} = \sigma \left( b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)} \right) \quad (19)$$

- A simplified version is called gated recurrent units or GRUs (Cho et al., 2014b; Chung et al., 2014, 2015a; Jozefowicz et al., 2015; Chrupala et al., 2015).

# Clipping Gradients

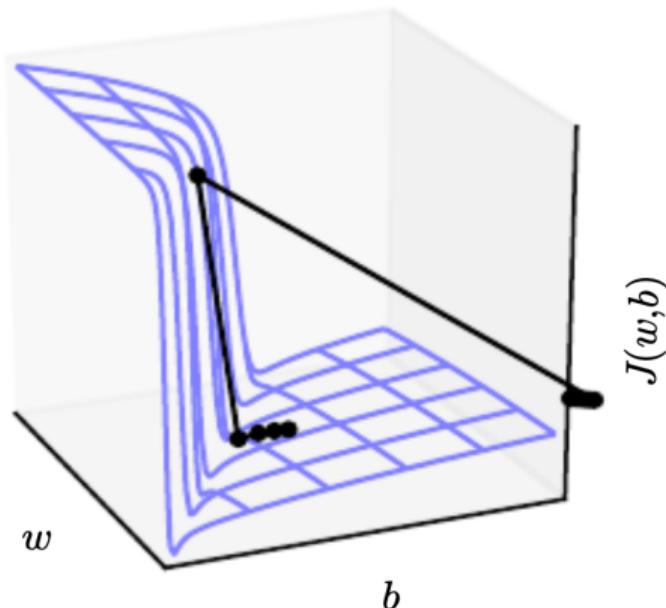
RNN

S.Lan

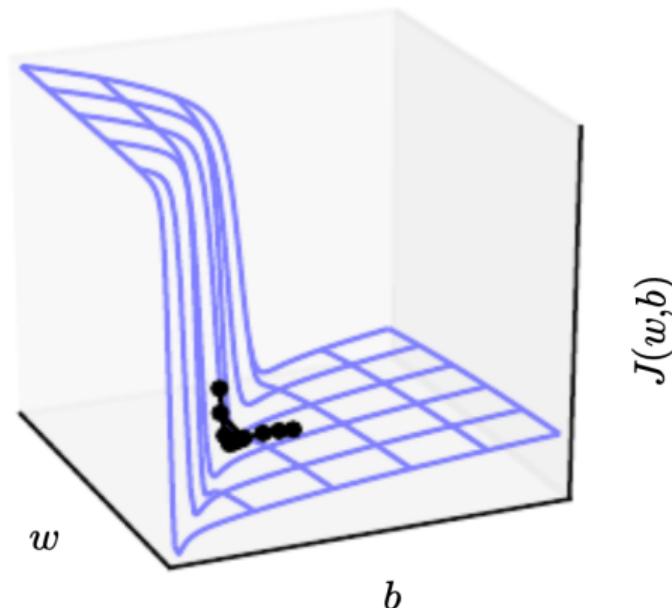
Introduction

Recurrent  
Neural Networks  
(RNN)Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNsLong-Time  
Memory:  
LSTM, GRUSoftware  
Implementation

Without clipping



With clipping





# Table of Contents

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNS,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

## 1 Introduction

## 2 Recurrent Neural Networks (RNN)

## 3 Variants: Bidirectional RNNs, Deep RNNS, Recursive NNs

## 4 Long-Time Memory: LSTM, GRU

## 5 Software Implementation



# TensorFlow

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

- In TensorFlow, we build RNN with `tf.keras.layers.SimpleRNN`, `tf.keras.layers.LSTM` or `tf.keras.layers.GRU`.
- `tf.keras.layers.SimpleRNN`, `tf.keras.layers.LSTM` and `tf.keras.layers.GRU` share some common arguments:
  - `units`: the dimensionality of the output space.
  - `activation`: activation function to use
  - `use_bias`: whether the layer uses a bias vector
  - `go_backwards`: whether to process the input sequence backwards and return the reversed sequence.
  - `unroll`: If True, the network will be unrolled, else a symbolic loop will be used. Unrolling can speed-up a RNN, although it tends to be more memory-intensive. Unrolling is only suitable for short sequences.



# PyTorch

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

- In PyTorch , we build RNN with `torch.nn.RNN`, `torch.nn.LSTM` or `torch.nn.GRU`.
- `torch.nn.RNN` , `torch.nn.LSTM` and `torch.nn.GRU` share some common arguments:
  - `input_size`: the number of expected features in the input  $x$
  - `hidden_size`: the number of features in the hidden state  $h$
  - `num_layers`: number of recurrent layers
  - `nonlinearity` : non-linearity to use, 'tanh' or 'relu'.
  - `bias`: whether to use bias weights  $b_i h$  and  $b_h h$
  - `bidirectional` : whether to becomes a bidirectional RNN



# More Reading

RNN

S.Lan

Introduction

Recurrent  
Neural Networks  
(RNN)

Variants:  
Bidirectional  
RNNs, Deep  
RNNs,  
Recursive NNs

Long-Time  
Memory:  
LSTM, GRU

Software  
Implementation

- TensorFlow

- <https://www.guru99.com/rnn-tutorial.html>
- <https://www.kdnuggets.com/2017/04/build-recurrent-neural-network-tensorflow.html>
- <https://www.tensorflow.org/guide/keras/rnn>

- PyTorch

- [https://www.deeplearningwizard.com/deep\\_learning/practical\\_pytorch/pytorch\\_recurrent\\_neural\\_networks/](https://www.deeplearningwizard.com/deep_learning/practical_pytorch/pytorch_recurrent_neural_networks/)
- <https://blog.floydhub.com/a-beginners-guide-on-recurrent-neural-networks-with-pytorch/>
- [https://pytorch.org/tutorials/intermediate/char\\_rnn\\_classification\\_tutorial.html](https://pytorch.org/tutorials/intermediate/char_rnn_classification_tutorial.html)