

Natural Language Processing

CSE 325/425



Sihong Xie

Lecture 13:

- Recurrent neural networks (RNN)

Language model review

Bi-gram $P(w_t|w_{t-1}) = \frac{\text{Count } w_t}{\text{Count } w_{t-1}}$

n -gram $P(w_t|[w_{t-1}, \dots, w_{t-n+1}]) = \frac{\text{Count } [w_t, \dots, w_{t-n+1}]}{\text{Count } [w_{t-1}, \dots, w_{t-n+1}]}$

Two issues

- Data sparsity: the occurrences of many $[w_t, \dots, w_{t-n+1}]$ are zeros.
- Model complexity: number of parameters increases exponentially in n .

The two issues are related: if we want longer range dependencies, we increase n , then both the data sparsity and model complexity become worse.

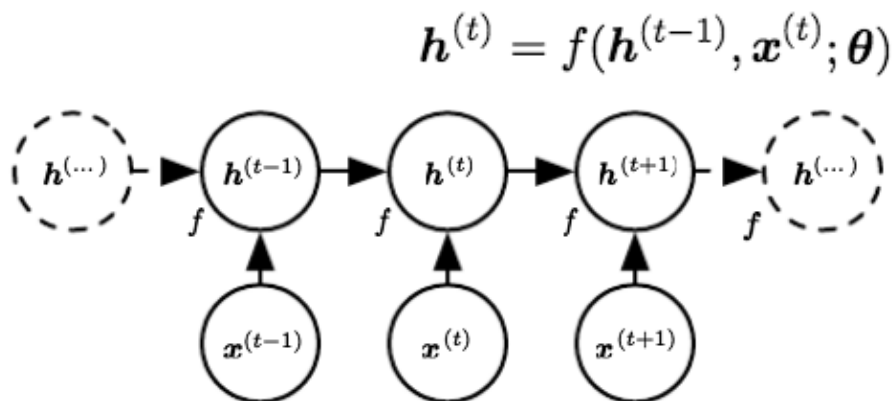
The students walked in the room and asked the ____?____ about the quiz questions.

Address the issues using neural network

Don't store the n -grams, but use a fixed-size model to predict the n -grams.

- model complexity is fixed.
- no data sparsity issue (no n -gram is computed)
- can be generalized to unseen sequences.

Recurrent Neural Networks (RNN)



“Recurrent” because the next \mathbf{h} is compute by the same function that calculates the previous \mathbf{h} .

$\mathbf{h}^{(t)}$ state summarizing what has happened before time step t . (theoretically speaking)

$\boldsymbol{\theta}$ a single model specifying how to transit to the next state (independent of t).

A running example

Recurrent Neural Networks

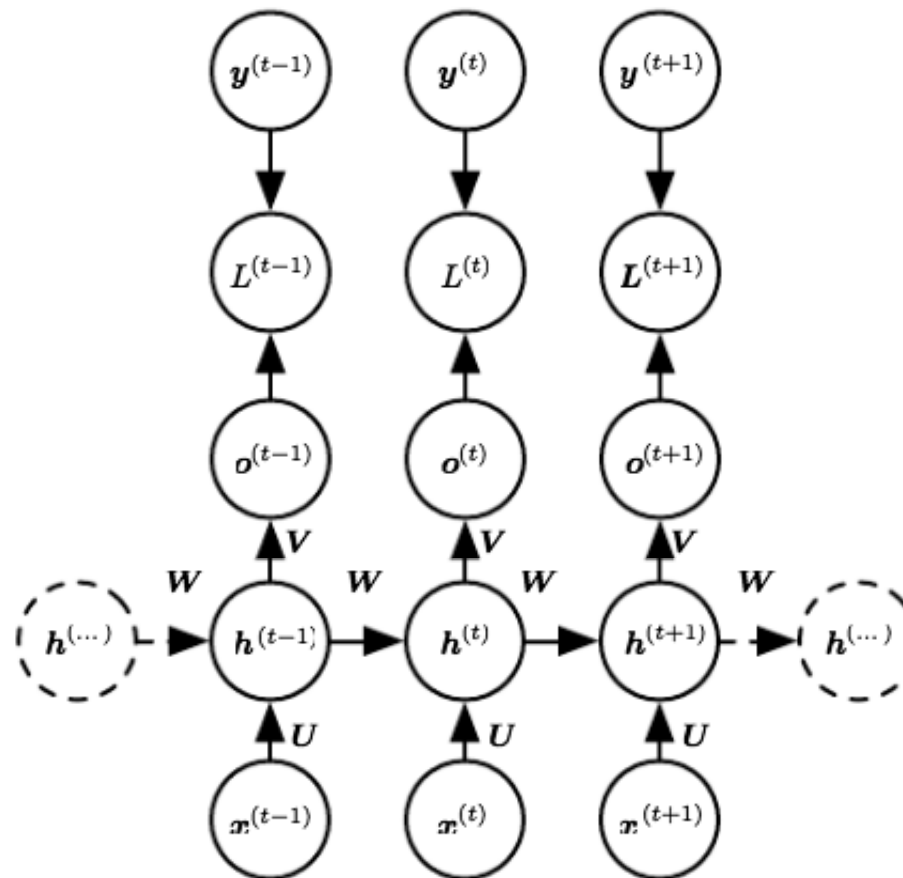
Output sequence
(e.g., POS tags)

Loss function

Output units

Hidden states

Input sequence
(e.g., sentence)



Training data:

$$\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}$$

Trainable parameters:

$$\theta = \{U, W, V, b, c\}$$

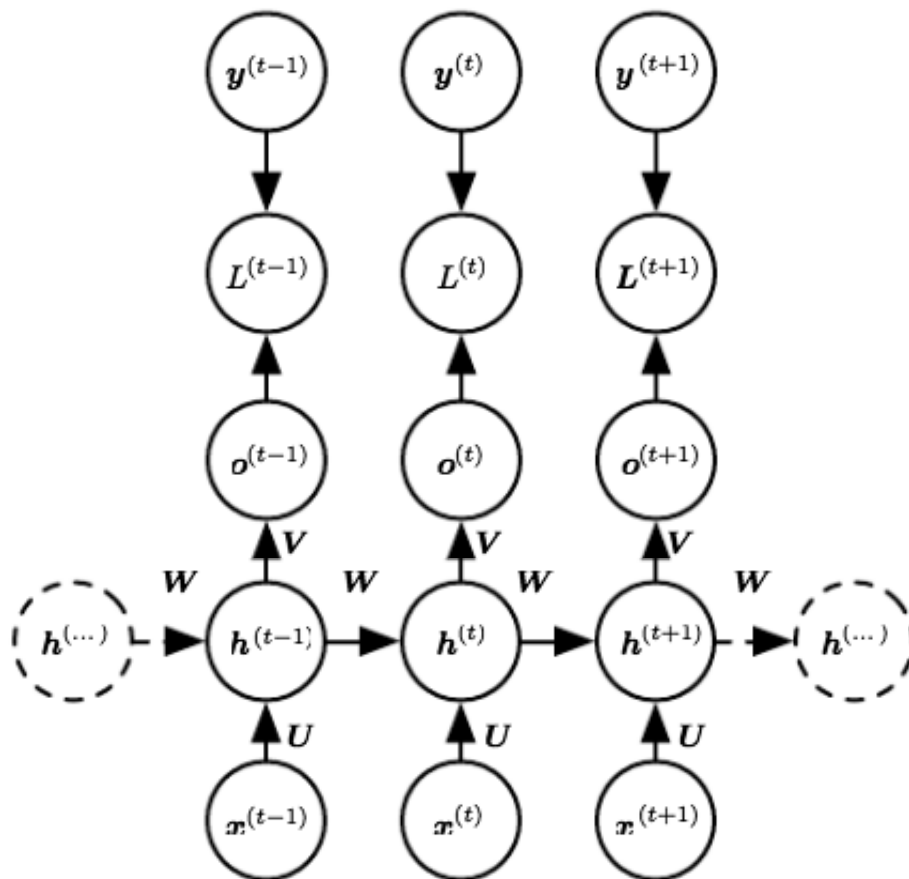
V : maps from h to o

W : maps from h to h

U : maps from x to h

b, c : biases

RNN forward pass



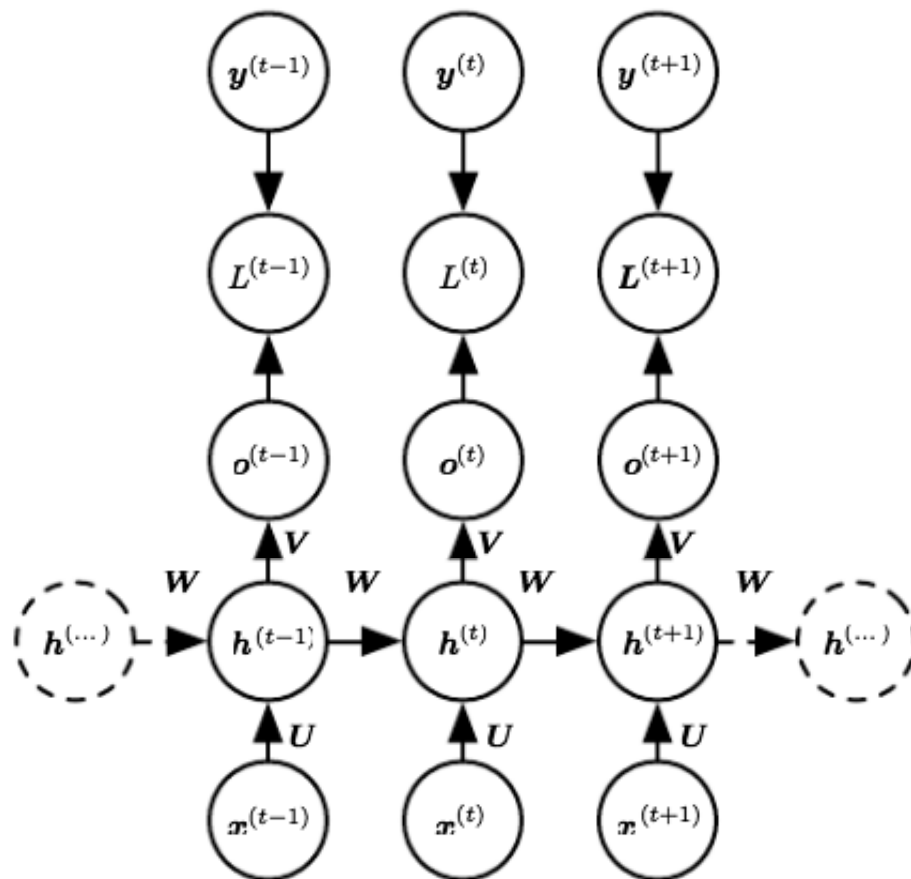
$$\mathbf{a}^{(t)} = \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)},$$

$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)}), \quad \text{activation function}$$

$$\mathbf{o}^{(t)} = \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)},$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)}), \quad \text{probability distribution}$$

RNN forward pass



Negative log likelihood (NLL) loss,
or the “perplexity”

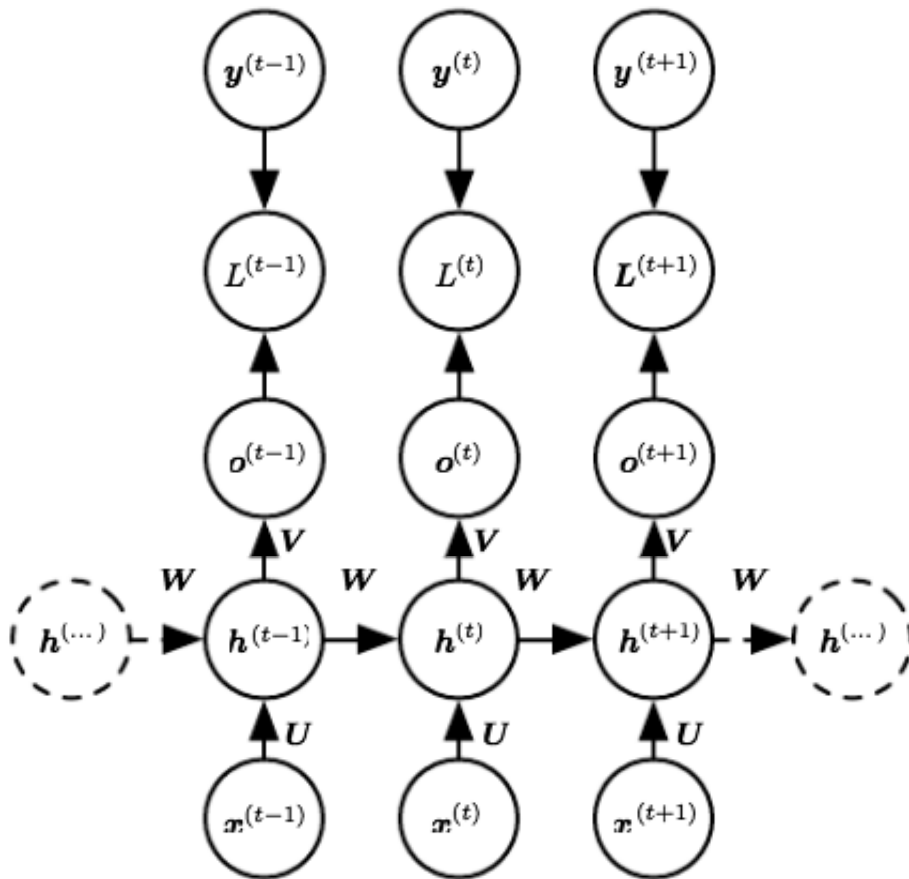
$$\begin{aligned}
 & L\left(\{x^{(1)}, \dots, x^{(\tau)}\}, \{y^{(1)}, \dots, y^{(\tau)}\}\right) \\
 &= \sum_t L^{(t)} \\
 &= - \sum_t \log p_{\text{model}}\left(\hat{y}^{(t)} \mid \{x^{(1)}, \dots, x^{(t)}\}\right)
 \end{aligned}$$

the ground truth label.

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

A running example

RNN back propagation



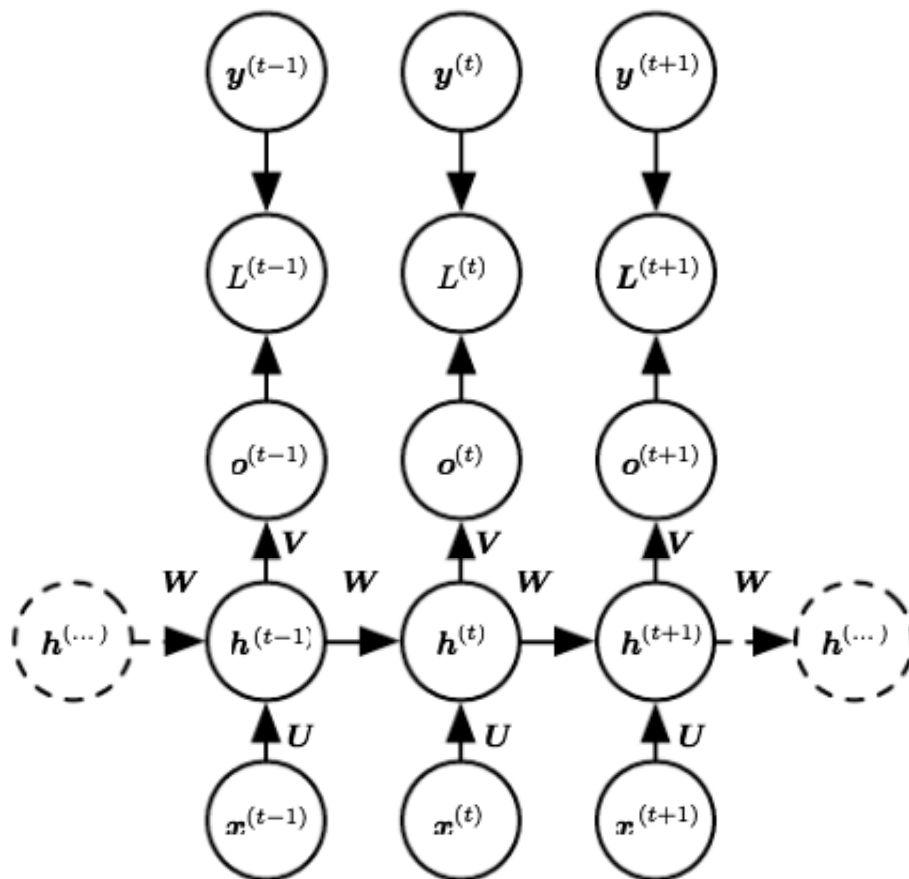
BPTT (Back Propagation through Time)

- Used for gradient descent training;
- A special name for RNN back-propagation;
- Need all information in the forward pass, making BPTT sequential and hard to parallelize.
- $\theta = \{U, W, V, b, c\}$ used in all steps.
- Two derivative rules applied:

$$\nabla_x(f(x) + g(x)) = \nabla_x f(x) + \nabla_x g(x)$$

$$\nabla_x(f(g(x))) = \nabla_q f(g(x)) \times \nabla_x g(x)$$

RNN back propagation



BPTT (Back Propagation through Time)

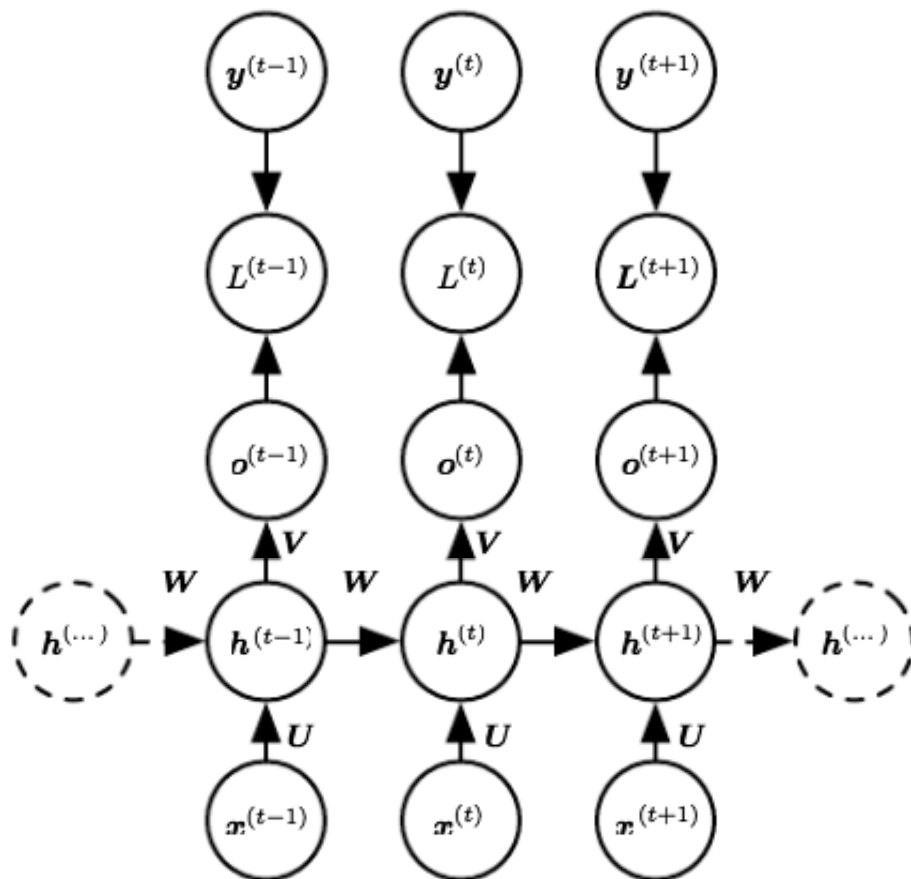
- focus on each step t (the final gradient is the sum of all gradients at each step).

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

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

(an element of $\nabla_{o^{(t)}} L$)

RNN backprop (base case)



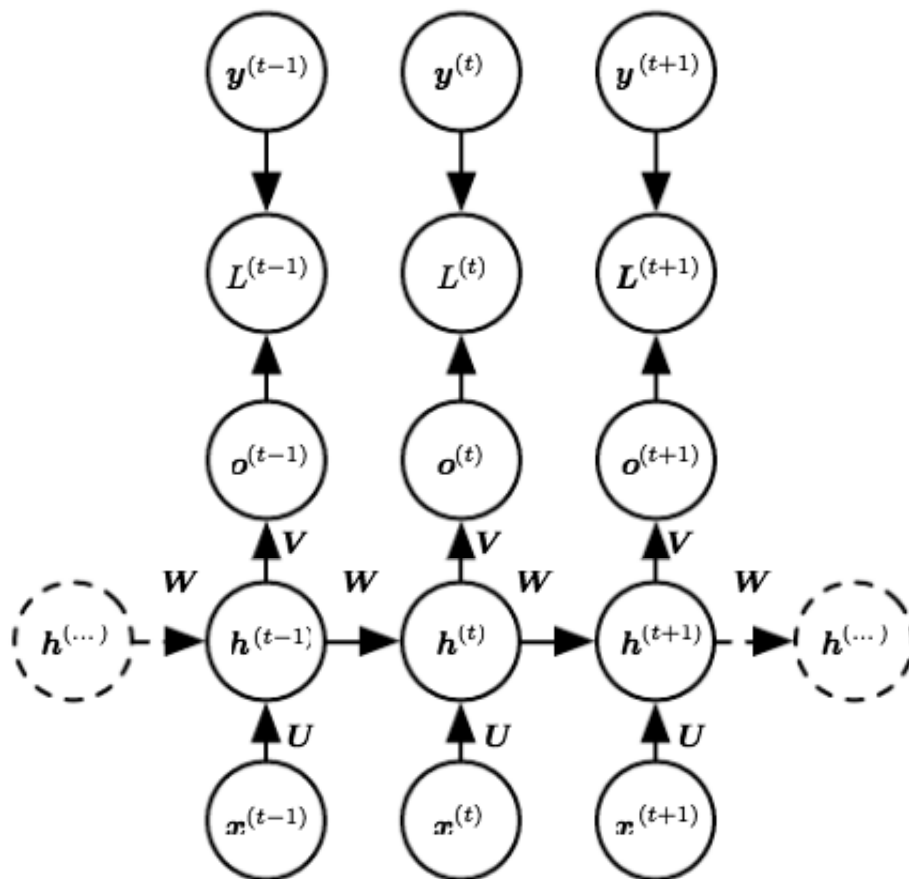
BPTT (Back Propagation through Time)

- Base case: at the final step τ

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

$$\text{since } o^{(\tau)} = Vh^{(\tau)} + c$$

RNN backprop (recurrent)



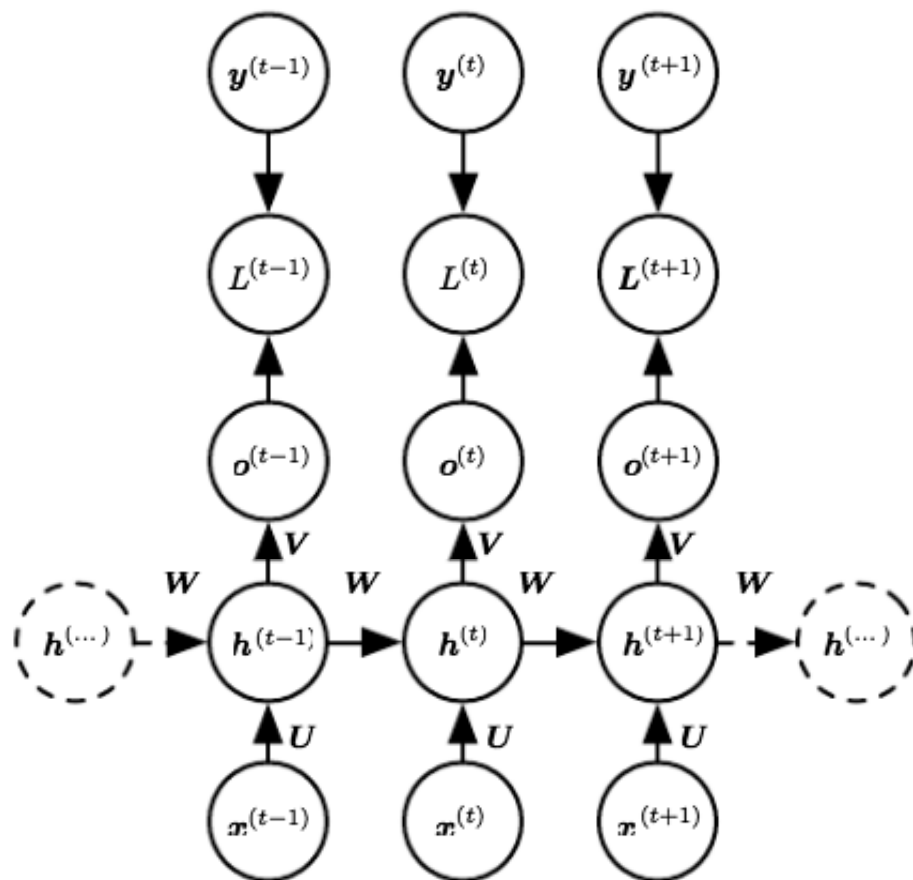
BPTT (Back Propagation through Time)

- Recursively, at any step $1 \leq t < \tau$

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

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

RNN backprop (recurrent)



BPTT (Back Propagation through Time)

- Recursively, at any step $1 \leq t < \tau$

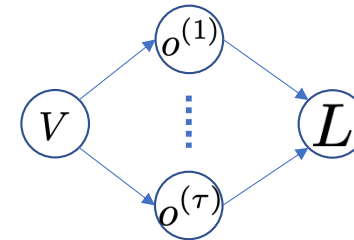
$$\begin{aligned}\nabla_{h^{(t)}} L &= \left(\frac{\partial h^{(t+1)}}{\partial h^{(t)}} \right)^\top (\nabla_{h^{(t+1)}} L) + \left(\frac{\partial o^{(t)}}{\partial h^{(t)}} \right)^\top (\nabla_{o^{(t)}} L) \\ &= \boxed{W^\top \text{diag} \left(1 - \left(h^{(t+1)} \right)^2 \right)} (\nabla_{h^{(t+1)}}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

$$h^{(t+1)} = \tanh(a^{(t+1)}) \quad (\text{element-wise})$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$a^{(t+1)} = W h^{(t)} + U x^{(t+1)} + b$$

RNN backprop (params)



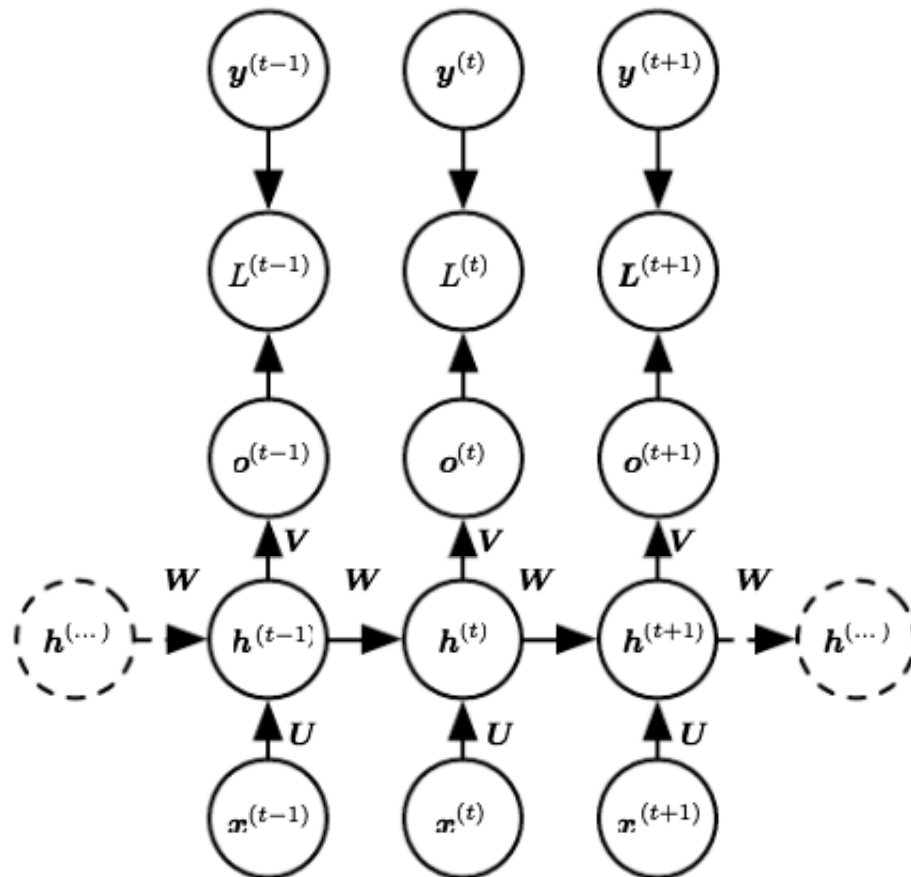
BPTT (Back Propagation through Time)

- at any step $1 \leq t < \tau$

since $\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)} + \mathbf{c}$

$$\nabla_{\mathbf{V}} L^{(t)} = (\nabla_{\mathbf{o}^{(t)}} L^{(t)}) \mathbf{h}^{(t)\top} \quad \nabla_{\mathbf{c}} L^{(t)} = \nabla_{\mathbf{o}^{(t)}} L^{(t)}$$

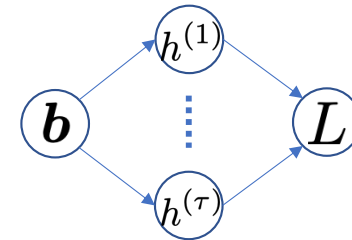
↓ accumulate over time



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

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

RNN backprop (params)

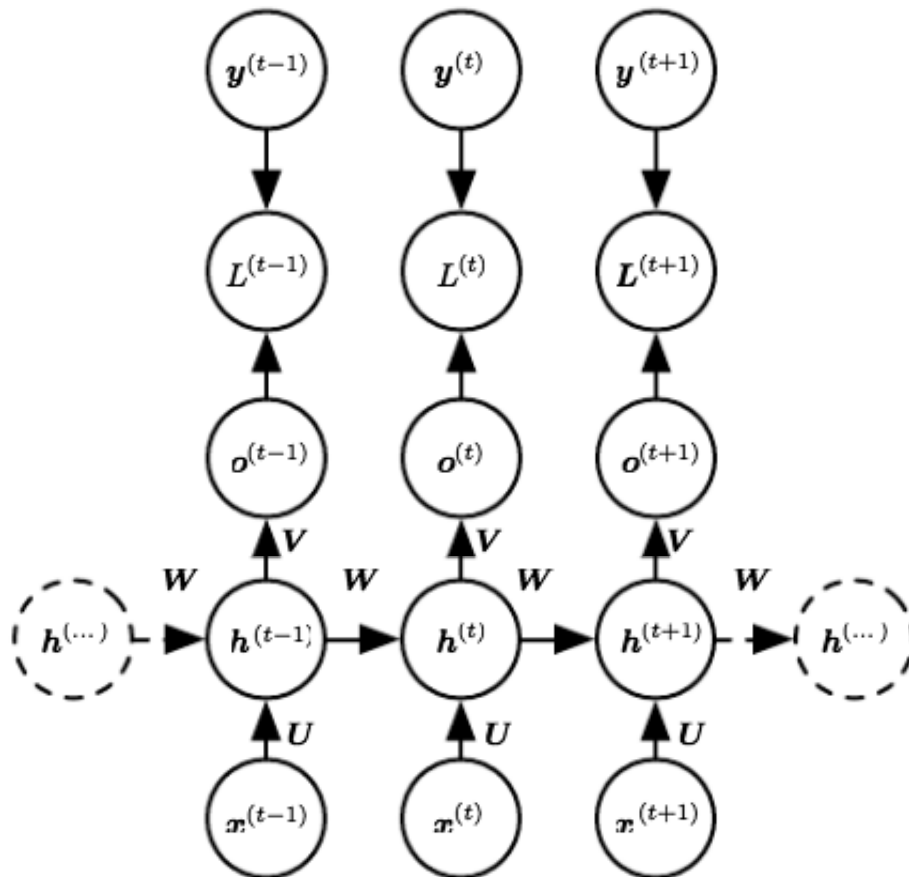


BPTT (Back Propagation through Time)

- at any step $1 \leq t < \tau$

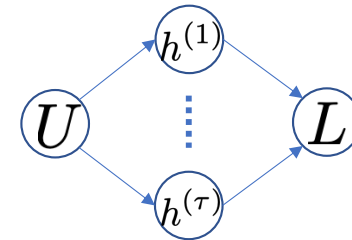
since $\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)})$

$$\mathbf{a}^{(t)} = \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b}$$



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

RNN backprop (params)

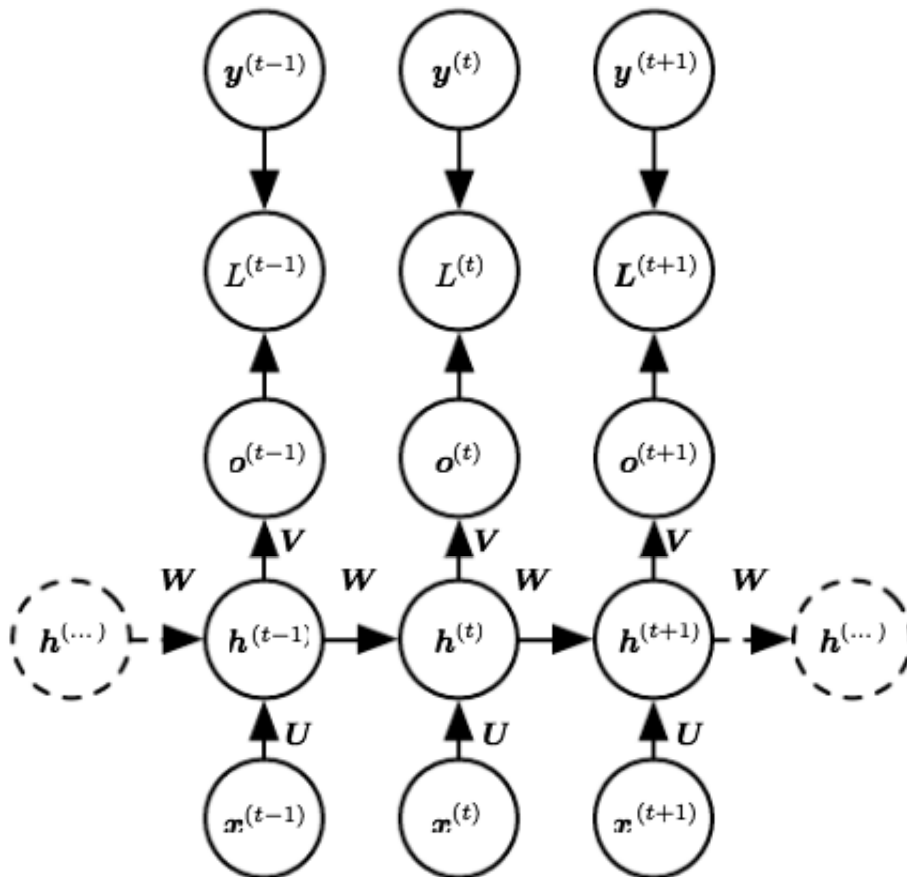


BPTT (Back Propagation through Time)

- at any step $1 \leq t < \tau$

since $\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)})$

$$\mathbf{a}^{(t)} = \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b}$$



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

(We have done $\nabla_{\mathbf{h}^{(t-1)}} L$, and leave $\nabla_{\mathbf{W}} L$ as an exercise.)