

# DSCI 565: RECURRENT NEURAL NETWORKS

*This content is protected and may not  
be shared, uploaded, or distributed.*

Ke-Thia Yao

Lecture 12: 2025 October 6

# Recurrent Neural Networks

2

## MLP without Hidden State

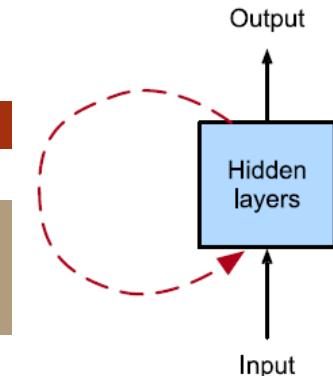
- MLP hidden layer with
  - Input:  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  $n$  batch size,  $d$  inputs
  - Output:  $\mathbf{H} \in \mathbb{R}^{n \times h}$ ,  $h$  hidden units
  - Parameter:  $\mathbf{W}_{\text{xh}} \in \mathbb{R}^{d \times h}$
  - Parameter:  $\mathbf{b}_h \in \mathbb{R}^{1 \times h}$

$$\mathbf{H} = \phi(\mathbf{X} \mathbf{W}_{\text{xh}} + \mathbf{b}_h)$$

## RNN with Hidden State

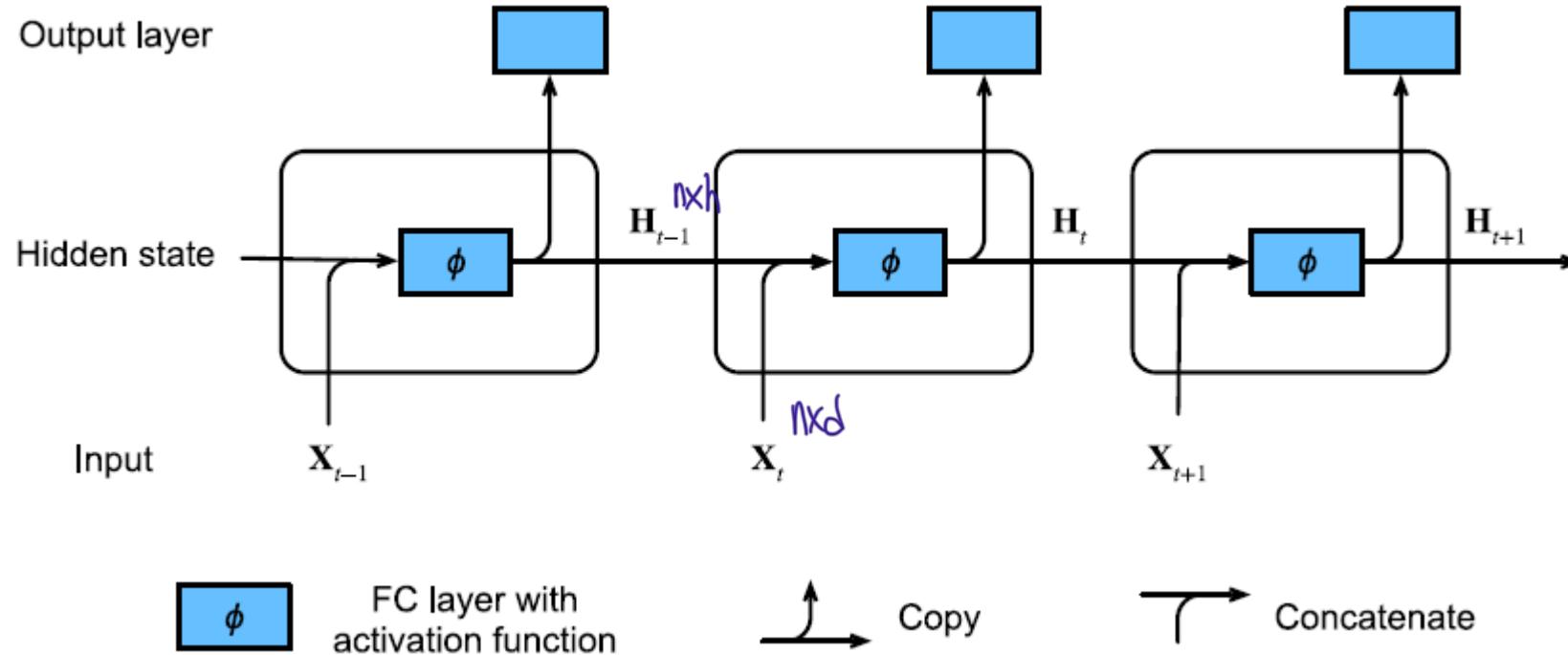
- RNN hidden layer with
  - Input:  $\mathbf{X}_t \in \mathbb{R}^{n \times d}$ ,  $n$  batch size,  $d$  inputs
  - Output:  $\mathbf{H}_t \in \mathbb{R}^{n \times h}$ ,  $h$  hidden units
  - Parameter:  $\mathbf{W}_{\text{xh}} \in \mathbb{R}^{d \times h}$
  - Parameter:  $\mathbf{W}_{\text{hh}} \in \mathbb{R}^{h \times h}$
  - Parameter:  $\mathbf{b}_h \in \mathbb{R}^{1 \times h}$

$$\mathbf{H}_t = \phi(\mathbf{X}_t \mathbf{W}_{\text{xh}} + \mathbf{H}_{t-1} \mathbf{W}_{\text{hh}} + \mathbf{b}_h)$$



# RNN with a Hidden Layer

3



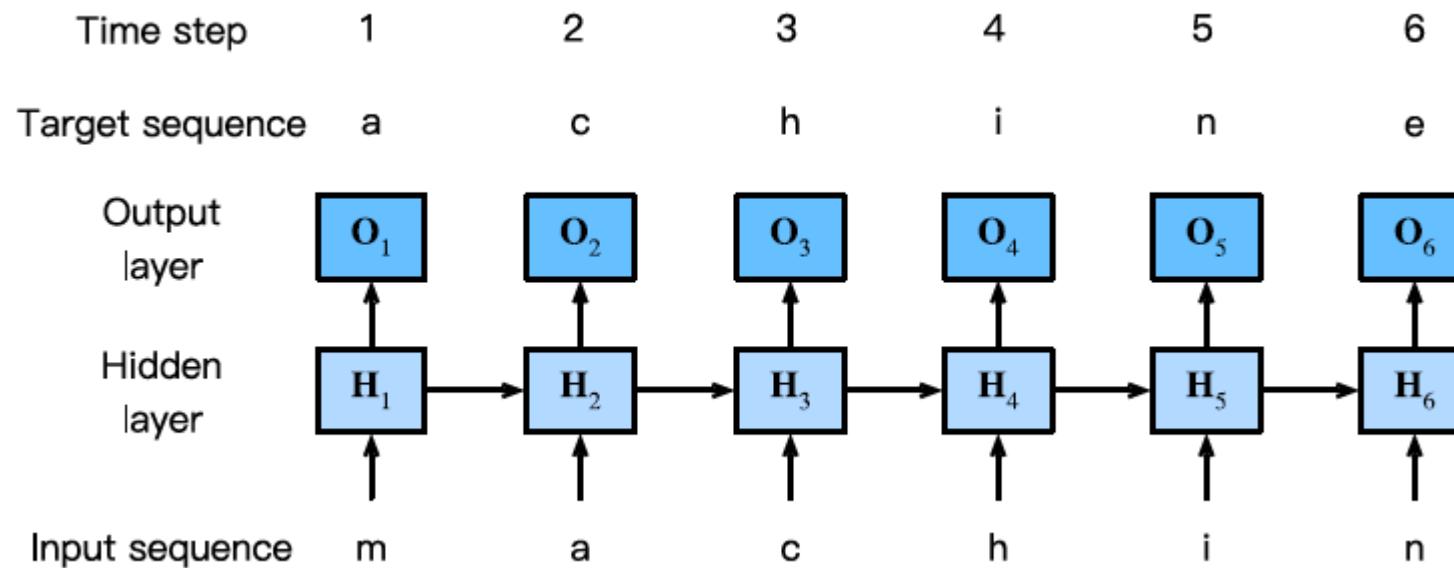
Can simplify by concatenation:

$$\text{concat}(\mathbf{X}_t, \mathbf{H}_t) \in \mathbb{R}^{n \times (d+h)}; \text{concat}(\mathbf{W}_{\text{vh}}, \mathbf{W}_{\text{hh}}) \in \mathbb{R}^{(d+h) \times h}$$

# Character-level RNN

4

Character Level



- Input sequence: machin
- Output sequence: achine
- Each character is encoded as a  $d$ -dimensional vector , .e.g., one-hot encoding

$d$ . dimension

# Gradient Clipping

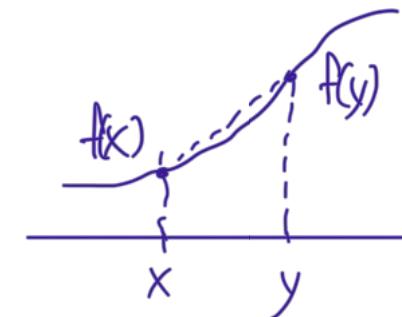
5

- Training RNN uses backpropagation through time of “depth”  $\mathcal{O}(T)$
- This can cause exploding or vanishing gradient problems
- When we update parameters with learning rate  $\eta$ :  $\mathbf{x} \leftarrow \mathbf{x} - \eta \mathbf{g}$
- To limit the gradient within a ball of radius  $\theta$ :

$$\mathbf{g} \leftarrow \min\left(1, \frac{\theta}{\|\mathbf{g}\|}\right) \mathbf{g}$$

- Formally, we say the objective/loss function  $f$  is Lipschitz continuous with constant  $L$ :

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|$$
$$|f(\mathbf{x}) - f(\mathbf{x} - \eta \mathbf{g})| \leq L \eta \|\mathbf{g}\|$$



# Notebook

6

- chapter\_recurrent-neural-networks/rnn-scratch.ipynb
- chapter\_recurrent-neural-networks/rnn-concise.ipynb

# Backpropagation Through Time

7

- A RNN with one-hidden layer be represented at each time step  $t$ :

$$h_t = f(x_t, h_{t-1}, w_h),$$

$$o_t = g(h_t, w_o),$$

- Where  $f, g$  are transformation of the hidden and output layer, respectively
- Notice  $h_t = f(x_t, h_{t-1}, w_h)$  depends on  $x_t, h_{t-1}, w_h$ , and recursively  $h_{t-1} = f(x_{t-1}, h_{t-2}, w_h)$  depends on  $x_{t-1}, h_{t-2}, w_h$
- Computing  $\frac{\partial h_t}{\partial w_h}$  requires chaining through multiple  $f(\dots)$

$$h_t = f(x_t, h_{t-1}, w_h),$$
$$o_t = g(h_t, w_o),$$

# Backpropagation Through Time

8

- The discrepancy between output  $o_t$  and target  $y_t$  is

$$L(x_1, \dots, x_T, y_1, \dots, y_T, w_h, w_o) = \frac{1}{T} \sum_{t=1}^T l(y_t, o_t)$$

- Using the chain rule

$$\begin{aligned}\frac{\partial L}{\partial w_h} &= \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_h} \\ &= \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_o)}{\partial h_t} \boxed{\frac{\partial h_t}{\partial w_h}}\end{aligned}$$

$$h_t = f(x_t, h_{t-1}, w_h),$$

$$o_t = g(h_t, w_o),$$

# Backpropagation Through Time

9

- Chain rule for partial of  $h_t$  wrt  $w_h$

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial w_h}$$

- Notice the above recurrence relation has the form

$$a_t = b_t + c_t a_{t-1}$$

- Then  $a_t = b_t + c_t(b_{t-1} + c_{t-1}a_{t-2}) = b_t + c_t b_{t-1} + c_t c_{t-1} a_{t-2}$
- And

$$a_t = b_t + \sum_{i=1}^{t-1} \left( \prod_{j=i+1}^t c_j \right) b_i$$

$$a_t = b_t + C_t b_{t-1} + C_t C_{t-1} b_{t-2} + \dots$$

# Backpropagation Through Time

10

- Then, substitute

$$a_t = \frac{\partial h_t}{\partial w_h},$$

$$b_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h},$$

$$c_t = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial h_{t-1}},$$

into  $a_t = b_t + c_t a_{t-1}$ .

- We remove the recurrent computation:

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left( \prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}.$$

# Strategies for computing $\frac{\partial h_t}{\partial w_h}$

11

- **Full Computation**
  - Compute for all time steps  $i = 1, \dots, t - 1$
  - Very slow, gradients can explode and sensitive to initial conditions
- **Truncating time steps**
  - Truncate after  $\tau$  steps by terminating sum at  $\partial h_{t-\tau}/\partial w_h$
  - Works well in practice. Biased toward more short-term influences.
- Random truncation
  - Randomly truncate the sequence
  - Empirically does not work much better than regular truncation

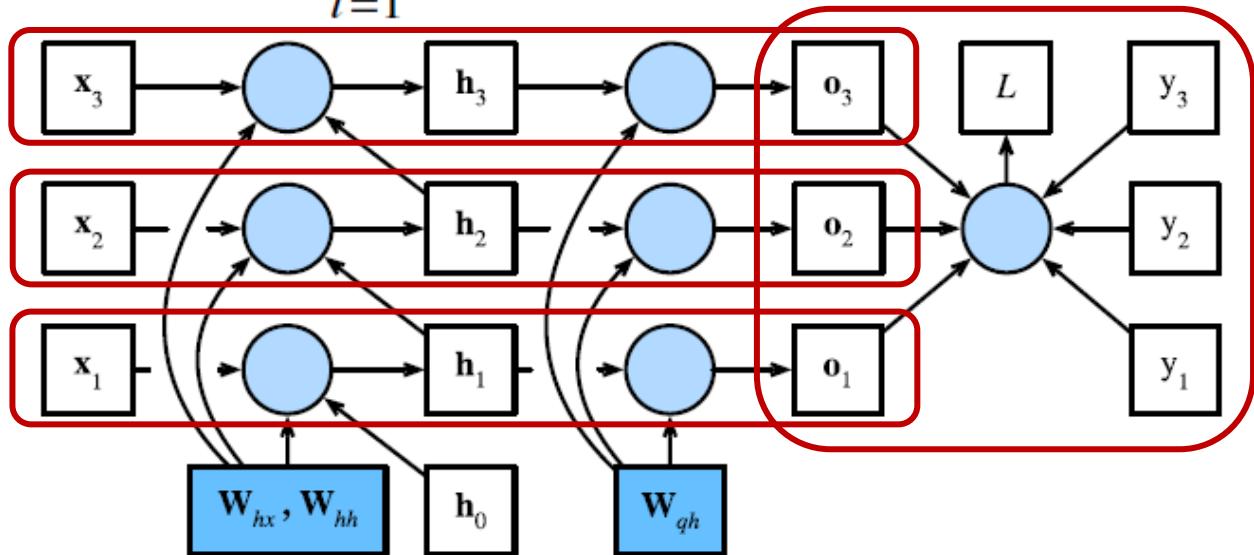
# Forward Pass

12

$$\mathbf{h}_t = \mathbf{W}_{hx} \mathbf{x}_t + \mathbf{W}_{hh} \mathbf{h}_{t-1},$$

$$\mathbf{o}_t = \mathbf{W}_{qh} \mathbf{h}_t,$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{o}_t, y_t) \quad \mathbf{o}_t \in \mathbb{R}^q$$



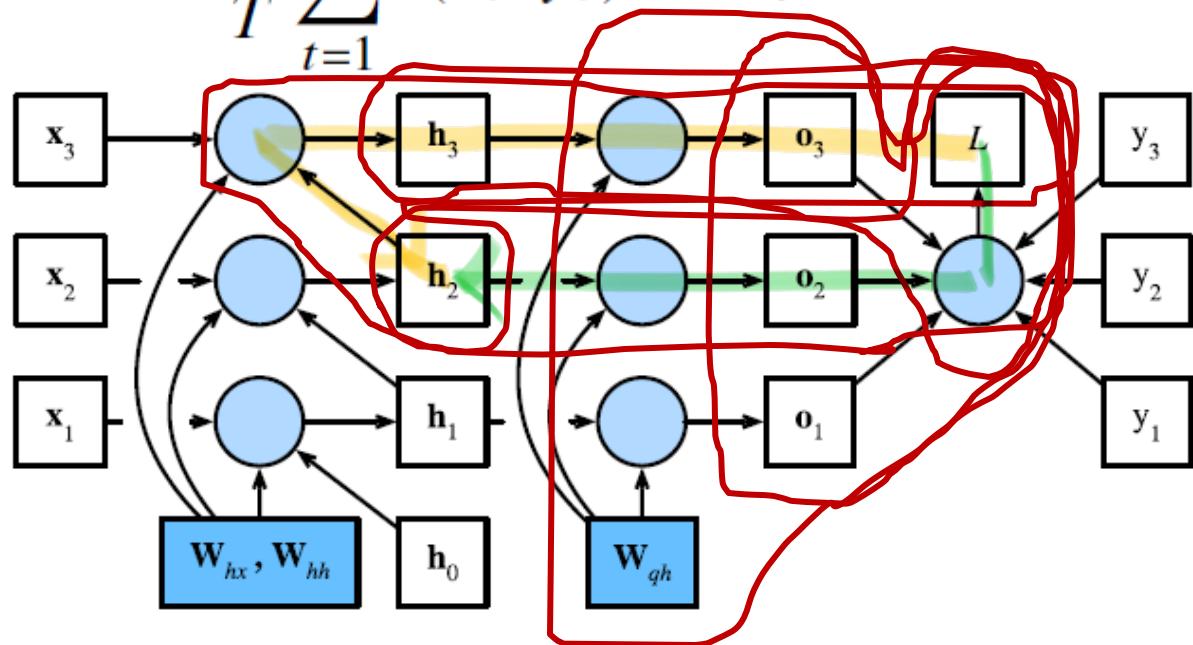
# Backpropagation Through Time in Detail

13

$$\mathbf{h}_t = \mathbf{W}_{hx} \mathbf{x}_t + \mathbf{W}_{hh} \mathbf{h}_{t-1},$$

$$\mathbf{o}_t = \mathbf{W}_{qh} \mathbf{h}_t,$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{o}_t, y_t) \quad \mathbf{o}_t \in \mathbb{R}^q$$



$$\frac{\partial L}{\partial \mathbf{o}_t} = \frac{\partial l(\mathbf{o}_t, y_t)}{T \cdot \partial \mathbf{o}_t} \in \mathbb{R}^q$$

$$\frac{\partial L}{\partial \mathbf{W}_{qh}} = \sum_{t=1}^T \text{prod} \left( \frac{\partial L}{\partial \mathbf{o}_t}, \frac{\partial \mathbf{o}_t}{\partial \mathbf{W}_{qh}} \right) = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{o}_t} \mathbf{h}_t^\top$$

$$\frac{\partial L}{\partial \mathbf{h}_T} = \text{prod} \left( \frac{\partial L}{\partial \mathbf{o}_T}, \frac{\partial \mathbf{o}_T}{\partial \mathbf{h}_T} \right) = \mathbf{W}_{qh}^\top \frac{\partial L}{\partial \mathbf{o}_T}$$

$$\frac{\partial L}{\partial \mathbf{h}_t} = \text{prod} \left( \frac{\partial L}{\partial \mathbf{h}_{t+1}}, \frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t} \right) + \text{prod} \left( \frac{\partial L}{\partial \mathbf{o}_t}, \frac{\partial \mathbf{o}_t}{\partial \mathbf{h}_t} \right)$$

$$= \mathbf{W}_{hh}^\top \frac{\partial L}{\partial \mathbf{h}_{t+1}} + \mathbf{W}_{qh}^\top \frac{\partial L}{\partial \mathbf{o}_t}$$

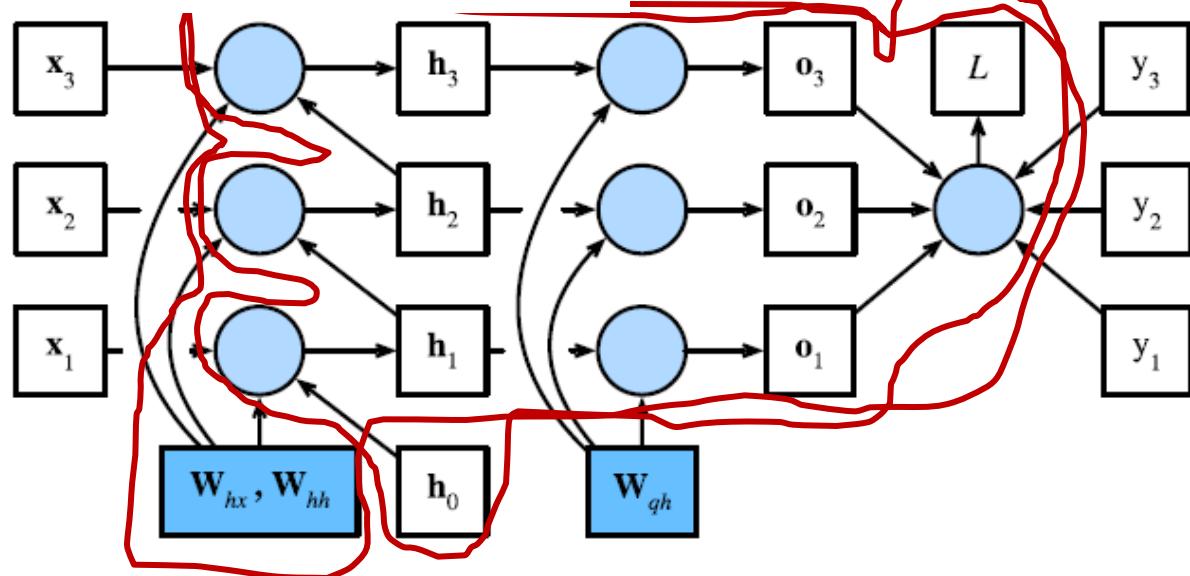
# Backpropagation Through Time in Detail

14

$$\mathbf{h}_t = \mathbf{W}_{hx} \mathbf{x}_t + \mathbf{W}_{hh} \mathbf{h}_{t-1},$$

$$\mathbf{o}_t = \mathbf{W}_{qh} \mathbf{h}_t,$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{o}_t, y_t) \quad \mathbf{o}_t \in \mathbb{R}^q$$



$$\frac{\partial L}{\partial \mathbf{W}_{hx}} = \sum_{t=1}^T \text{prod} \left( \frac{\partial L}{\partial \mathbf{h}_t}, \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_{hx}} \right) = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{h}_t} \mathbf{x}_t^\top,$$

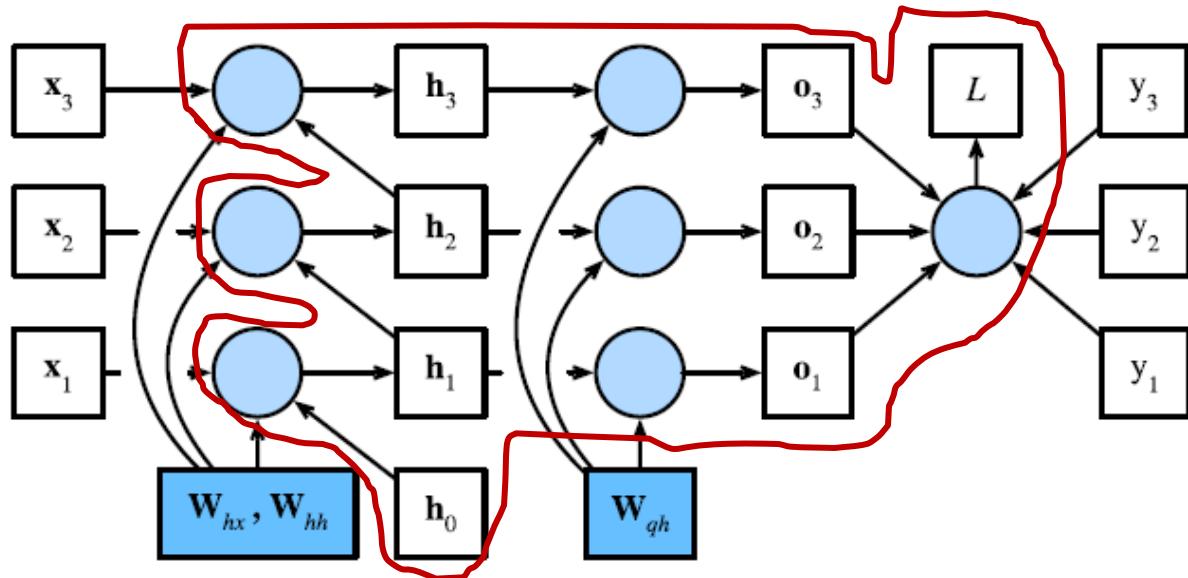
$$\frac{\partial L}{\partial \mathbf{W}_{hh}} = \sum_{t=1}^T \text{prod} \left( \frac{\partial L}{\partial \mathbf{h}_t}, \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_{hh}} \right) = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{h}_t} \mathbf{h}_{t-1}^\top,$$

$$\frac{\partial L}{\partial \mathbf{h}_t} = \sum_{i=t}^T (\mathbf{W}_{hh}^\top)^{T-i} \mathbf{W}_{qh}^\top \frac{\partial L}{\partial \mathbf{o}_{T+i}}$$

# Backpropagation Through Time in Detail

15

$$\frac{\partial L}{\partial \mathbf{h}_t} = \sum_{i=t}^T (\mathbf{W}_{hh}^\top)^{T-i} \mathbf{W}_{qh}^\top \frac{\partial L}{\partial \mathbf{o}_{T+i-t}}$$



- For long sequences, the matrix  $\mathbf{W}_{hh}^T$  is raised to large powers
- Eigenvalues less than 1, it becomes very small
- Eigenvalues greater than 1, it becomes very large
- Matrix becomes ill conditioned, contributing to vanishing and exploding gradients