

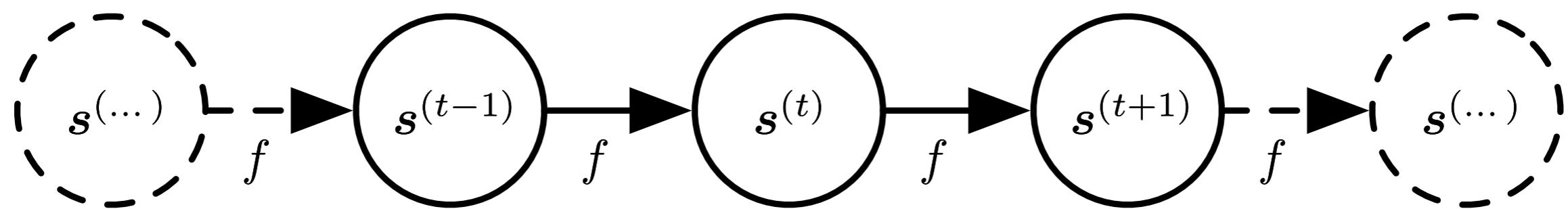
Sequence Modeling: Recurrent and Recursive Nets

Rupam Acharya

Basics

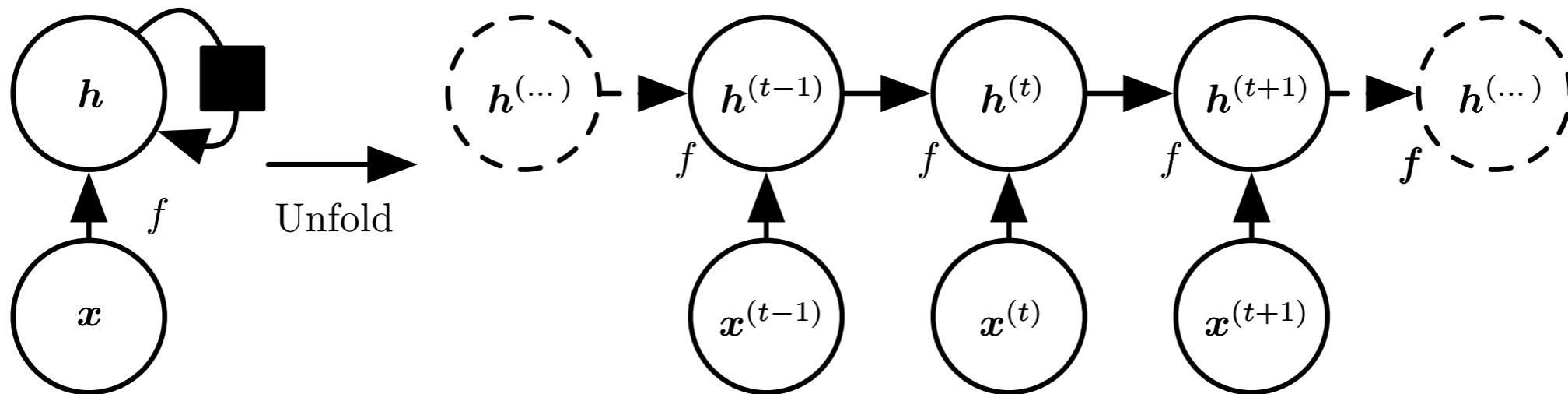
- CNN's are specialized for grid like inputs
- RNN's process sequential inputs
- Can process input of variable size unlike CNN
- Parameter Sharing in RNN helps to work with different length inputs

Classical Dynamical Systems



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

Unfolding Computation Graphs



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

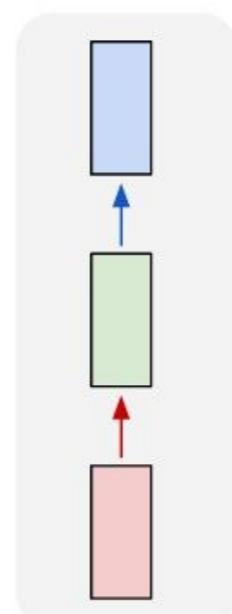
- Predict Future from the past
- Network learns lossy summary of the past sequence

Advantage of Unfolding

- Regardless of sequence length the input sizee is always same
- Same set of parameter is used over the time
- This motivates the RNN

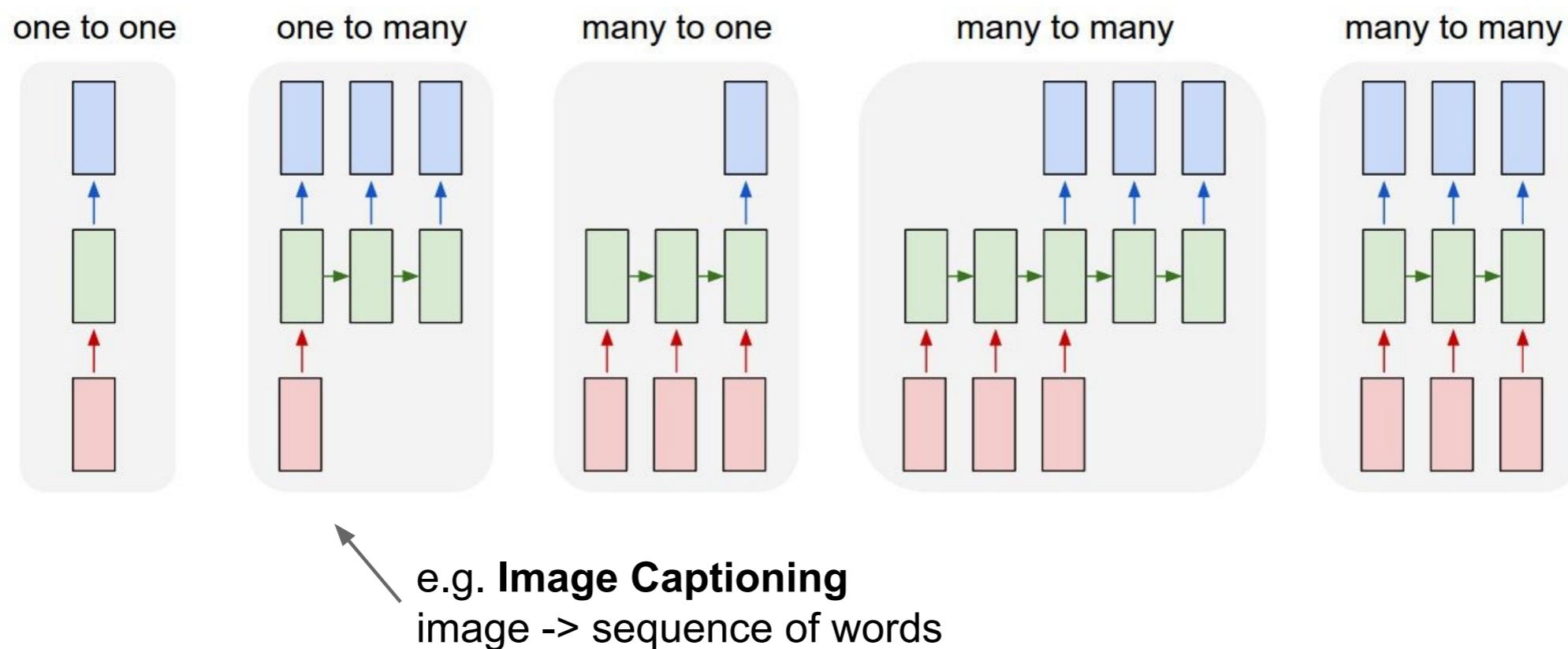
“Vanilla” Neural Network

one to one

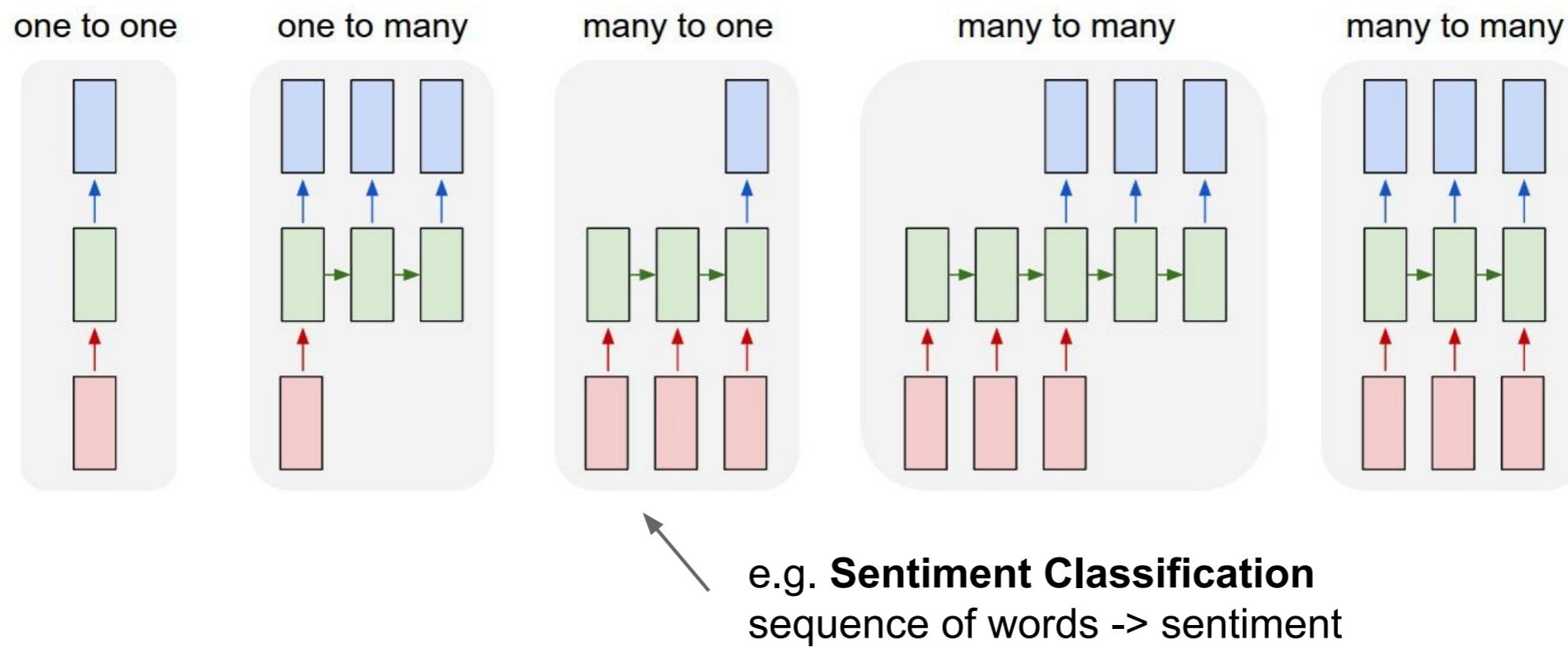


Vanilla Neural Networks

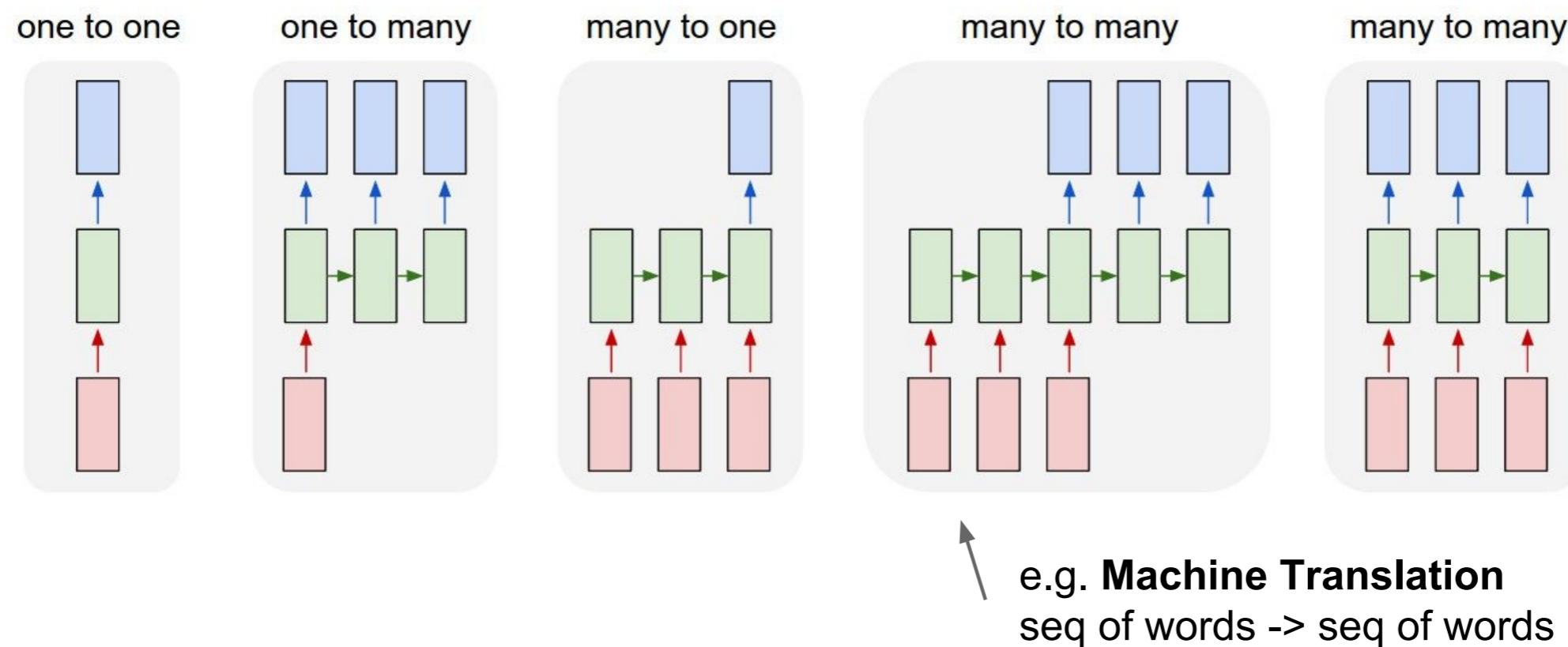
Recurrent Neural Networks: Process Sequences



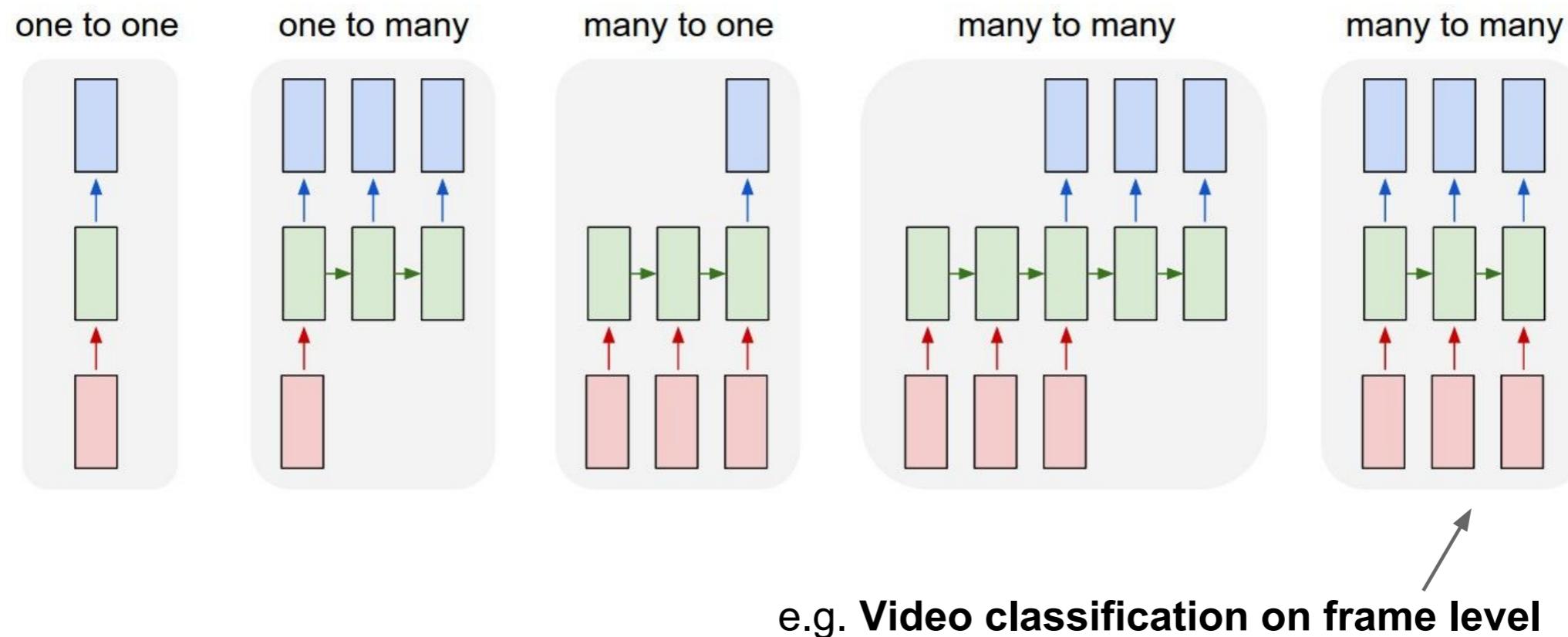
Recurrent Neural Networks: Process Sequences



Recurrent Neural Networks: Process Sequences



Recurrent Neural Networks: Process Sequences



Different Design Pattern

- Recurrent networks that produce an output at each time step and have recurrent connections between hidden units, illustrated in figure 10.3.
- Recurrent networks that 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, illustrated in figure 10.4
- Recurrent networks with recurrent connections between hidden units, that read an entire sequence and then produce a single output, illustrated in figure 10.5.

Recurrent Hidden Units

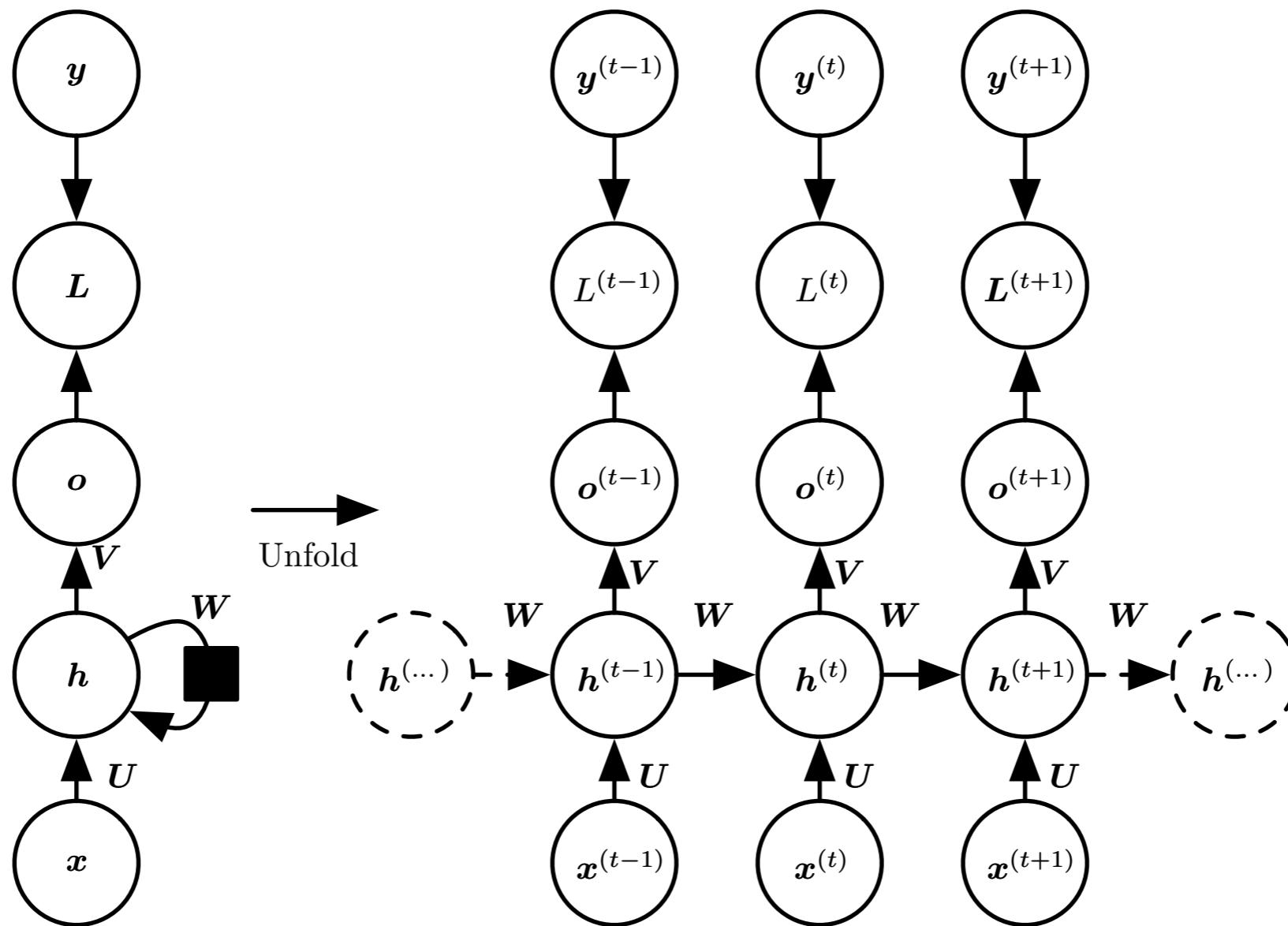


Figure 10.3

Parameter Setting

$t = 1$ to $t = \tau$, we apply the following update equations:

$$\mathbf{a}^{(t)} = \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} \quad (10.8)$$

$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)}) \quad (10.9)$$

$$\mathbf{o}^{(t)} = \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)} \quad (10.10)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)}) \quad (10.11)$$

$$L\left(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}\right) \quad (10.12)$$

$$= \sum_t L^{(t)} \quad (10.13)$$

$$= - \sum_t \log p_{\text{model}}\left(y^{(t)} \mid \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\}\right), \quad (10.14)$$

- Runtime is $O(\tau)$ and cannot be reduced

Recurrence through only the Output

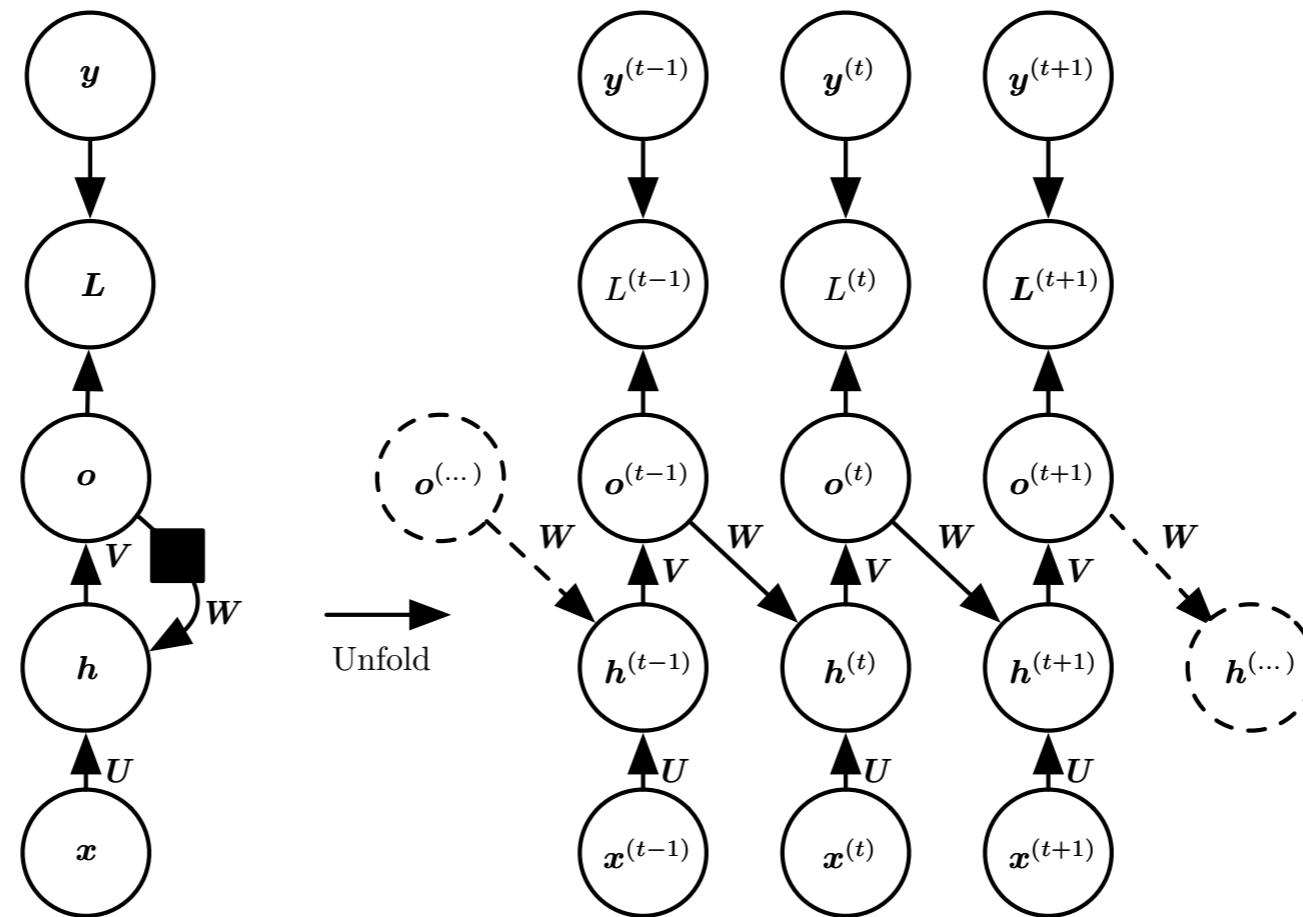


Figure 10.4

- Feedback connection from output to hidden
- Not as strong as hidden to hidden connection

Output to Hidden Connection

- Cannot simulate Universal Turing Machine
- Requires output to capture past information which is not the case
- Advantage is training can be parallelized
- Can be trained using **teacher forcing**

Teacher Forcing

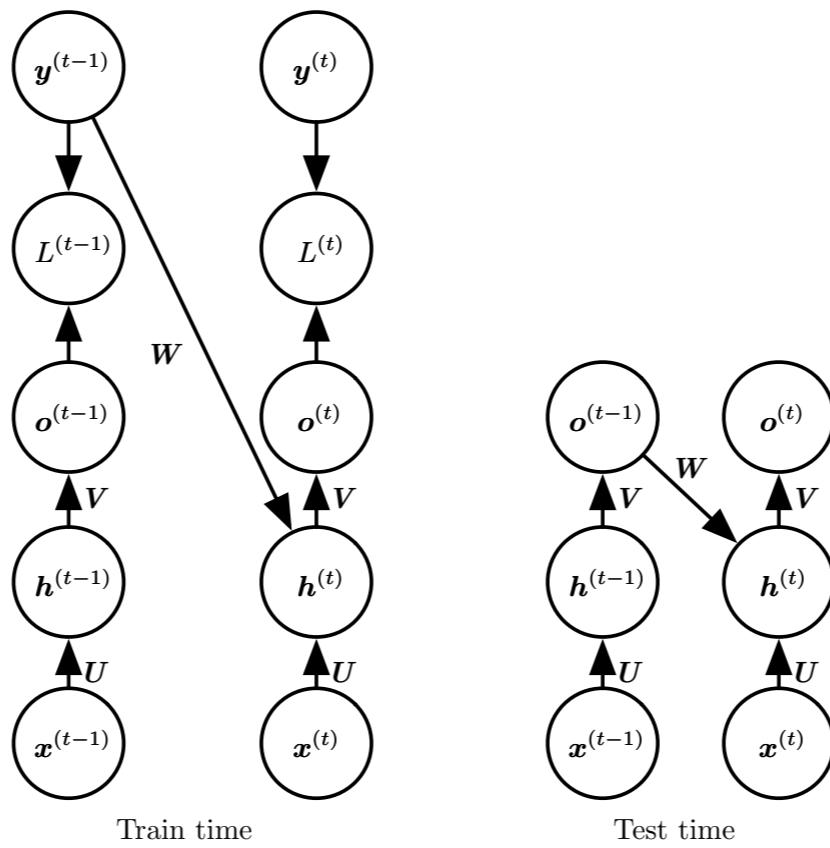


Figure 10.6: Illustration of teacher forcing. Teacher forcing is a training technique that is applicable to RNNs that have connections from their output to their hidden states at the next time step. (*Left*) At train time, we feed the *correct output* $\mathbf{y}^{(t)}$ drawn from the train set as input to $\mathbf{h}^{(t+1)}$. (*Right*) When the model is deployed, the true output is generally not known. In this case, we approximate the correct output $\mathbf{y}^{(t)}$ with the model's output $\mathbf{o}^{(t)}$, and feed the output back into the model.

$$\begin{aligned} & \log p(\mathbf{y}^{(1)}, \mathbf{y}^{(2)} \mid \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \\ &= \log p(\mathbf{y}^{(2)} \mid \mathbf{y}^{(1)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) + \log p(\mathbf{y}^{(1)} \mid \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \end{aligned}$$

Sequence Input, Single Output

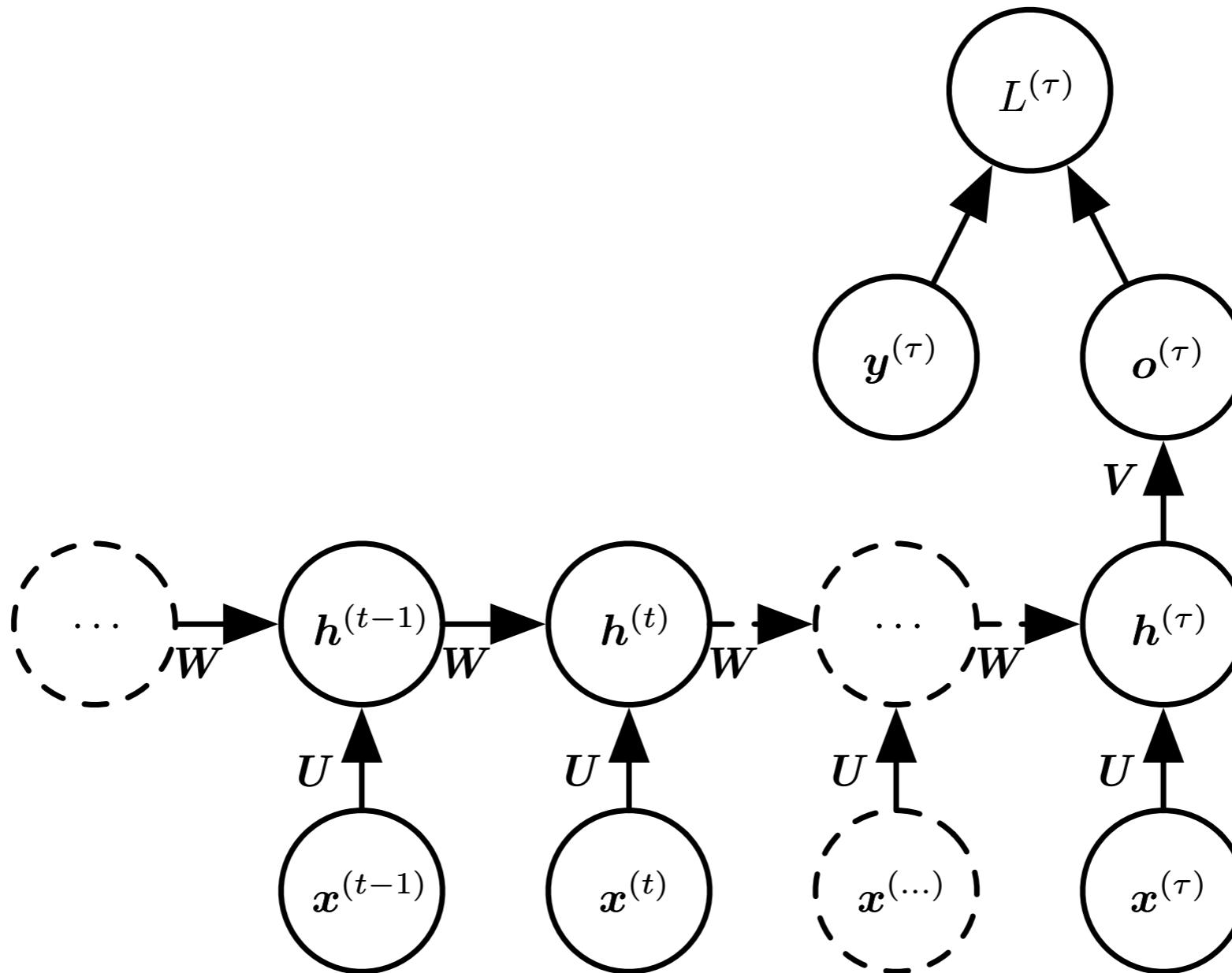


Figure 10.5

Computing Gradient in RNN

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

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

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

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

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

- Need to Compute

Computing Gradient in RNN

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

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

$$\nabla_{\boldsymbol{h}^{(\tau)}} L = \mathbf{V}^\top \nabla_{\boldsymbol{o}^{(\tau)}} L.$$

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

- This is what we have

Computing Gradient in RNN

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

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

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

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

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

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

- This is what we get

Fully Connected Graphical Model

- Chain Rule of Conditional Probability

$$P(\mathbb{Y}) = P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}) = \prod_{t=1}^{\tau} P(\mathbf{y}^{(t)} \mid \mathbf{y}^{(t-1)}, \mathbf{y}^{(t-2)}, \dots, \mathbf{y}^{(1)})$$

- The negative log-likelihood according to this model is

$$L = \sum_t L^{(t)}$$

where

$$L^{(t)} = -\log P(\mathbf{y}^{(t)} = y^{(t)} \mid y^{(t-1)}, y^{(t-2)}, \dots, y^{(1)}).$$

Fully Connected Graphical Model

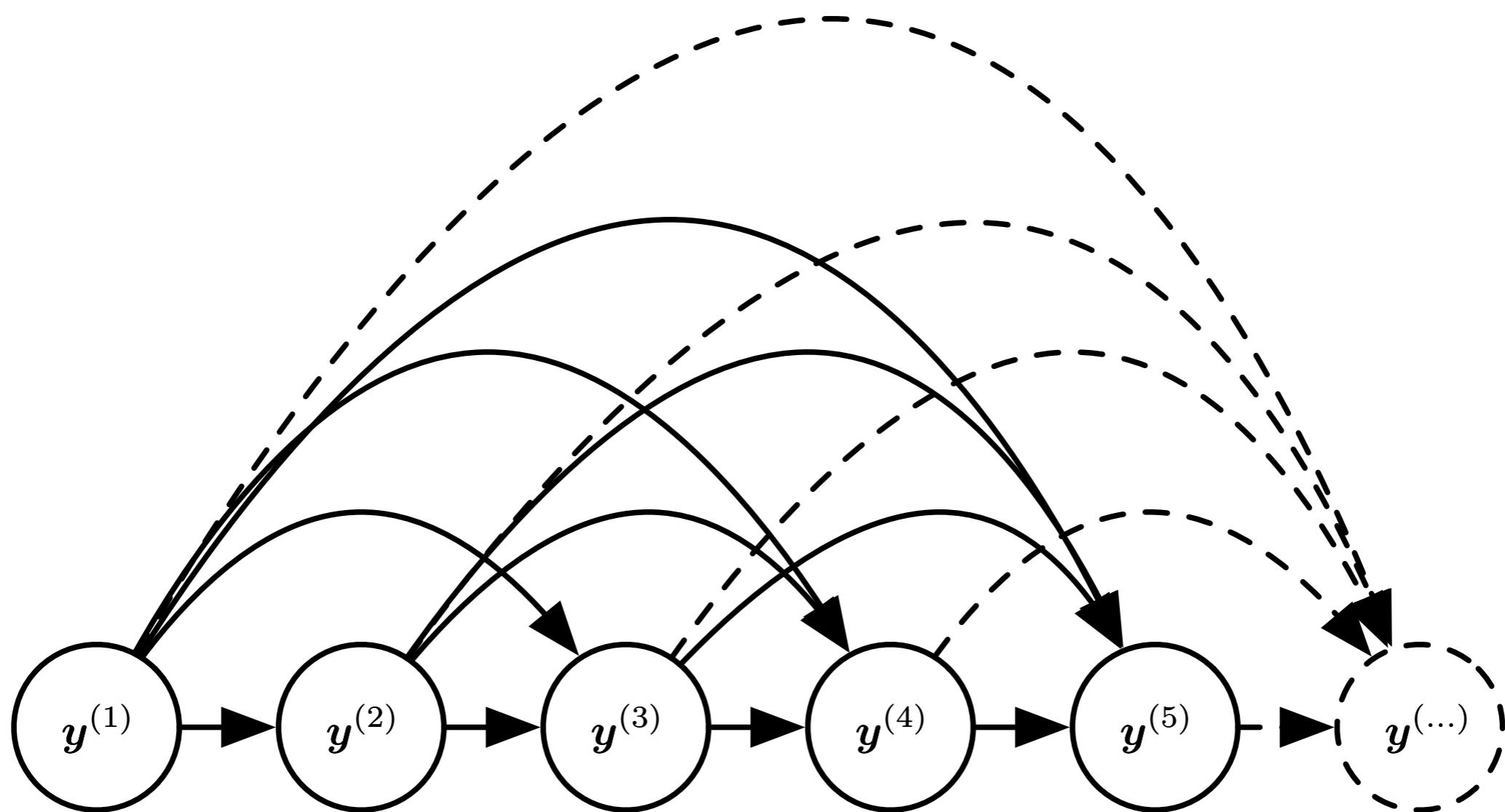
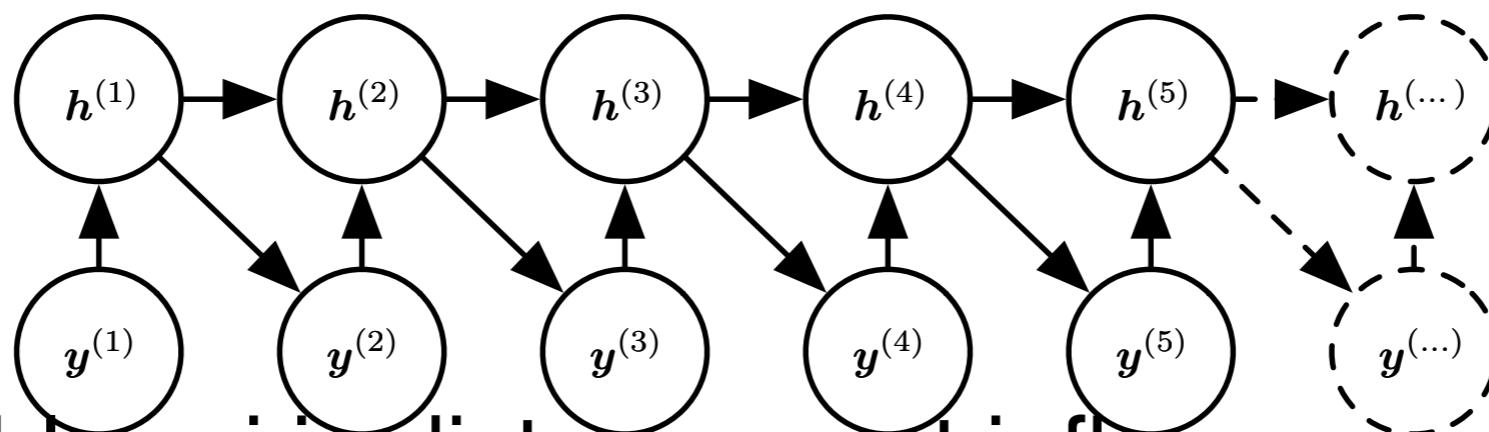


Figure 10.7

(Goodfellow 2016)

RNN Graphical Model



- A variable y_i in distance past influence y_t via h
- The model can be efficiently parameterized by using same conditional probability in each time step
- When the variables are observed the joint distribution can be evaluated efficiently
- Hard to predict missing values and optimizing the parameters

Sampling from Model

- Simply sample from the conditional distribution
- But length of the sequence is not known!!
- If output is a symbol from a vocabulary then add an eos symbol to stop the sampling process
- Generate Bernouli output to determine whether to continue generation or not

Vector to Sequence

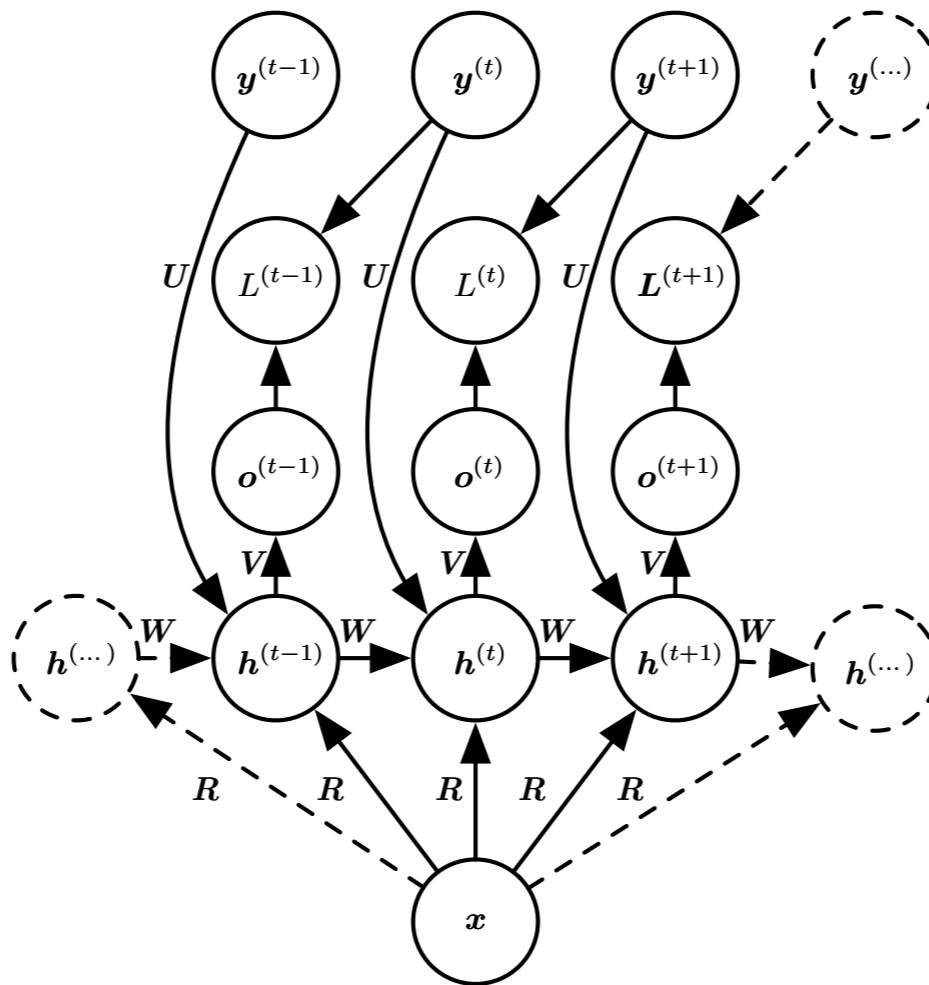


Figure 10.9: An RNN that maps a fixed-length vector \mathbf{x} into a distribution over sequences \mathbf{Y} . This RNN is appropriate for tasks such as image captioning, where a single image is used as input to a model that then produces a sequence of words describing the image. Each element $\mathbf{y}^{(t)}$ of the observed output sequence serves both as input (for the current time step) and, during training, as target (for the previous time step).

Hidden and Output Recurrence

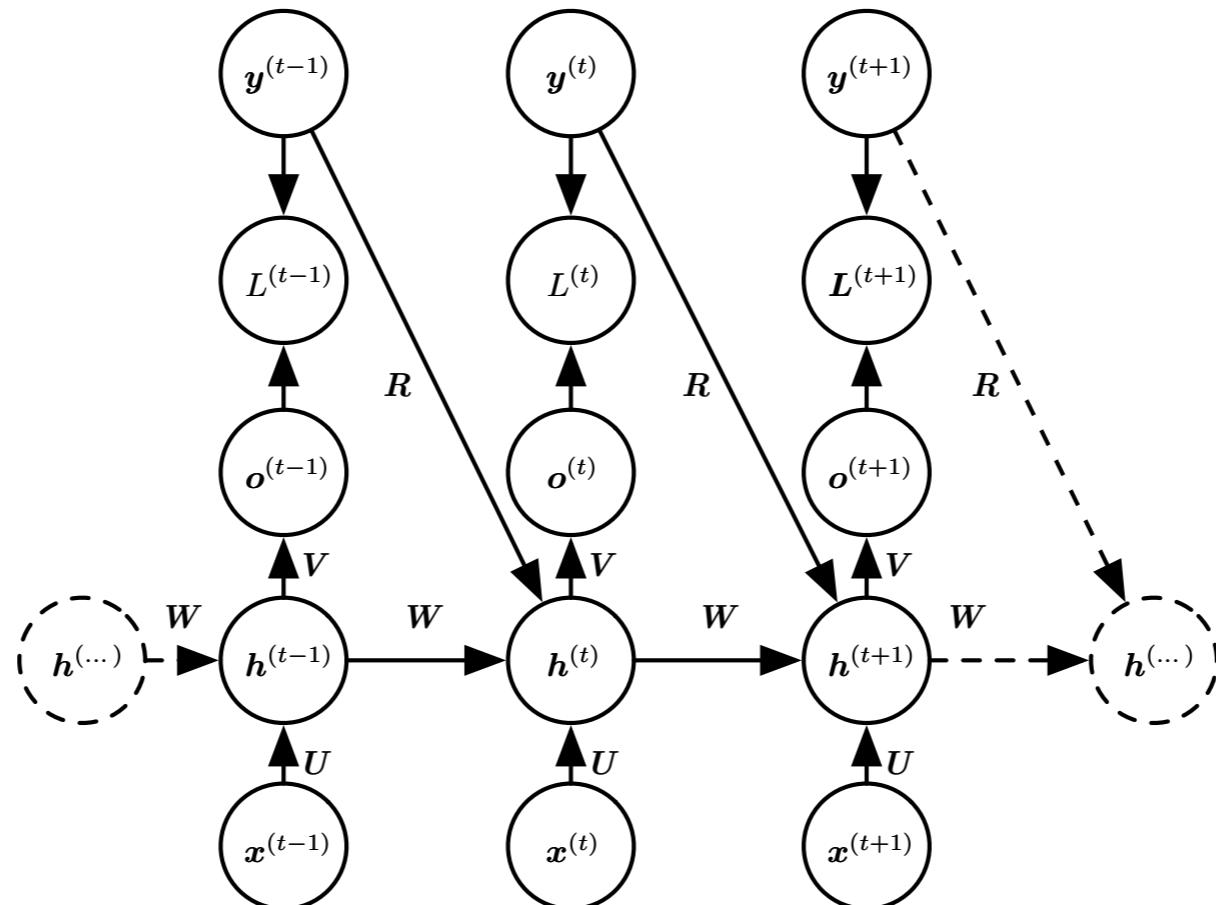


Figure 10.10: A conditional recurrent neural network mapping a variable-length sequence of \mathbf{x} values into a distribution over sequences of \mathbf{y} values of the same length. Compared to figure 10.3, this RNN contains connections from the previous output to the current state. These connections allow this RNN to model an arbitrary distribution over sequences of \mathbf{y} given sequences of \mathbf{x} of the same length. The RNN of figure 10.3 is only able to represent distributions in which the \mathbf{y} values are conditionally independent from each other given the \mathbf{x} values.