

10-301/601: Introduction to Machine Learning

Lecture 19 – Recurrent Neural Networks

Henry Chai

7/1/24

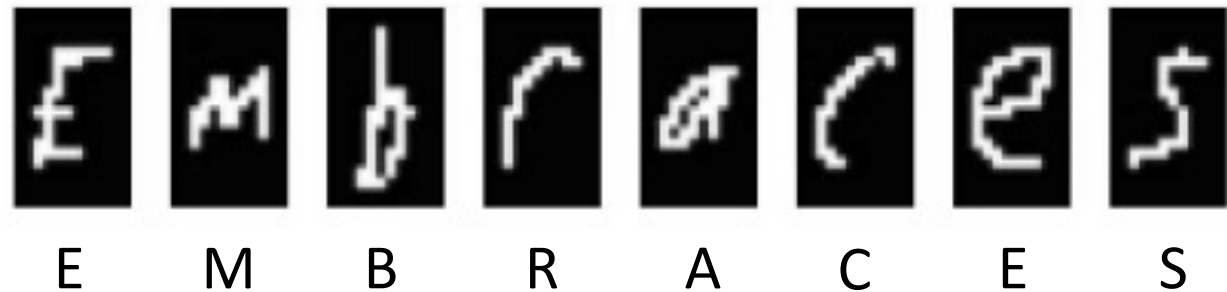
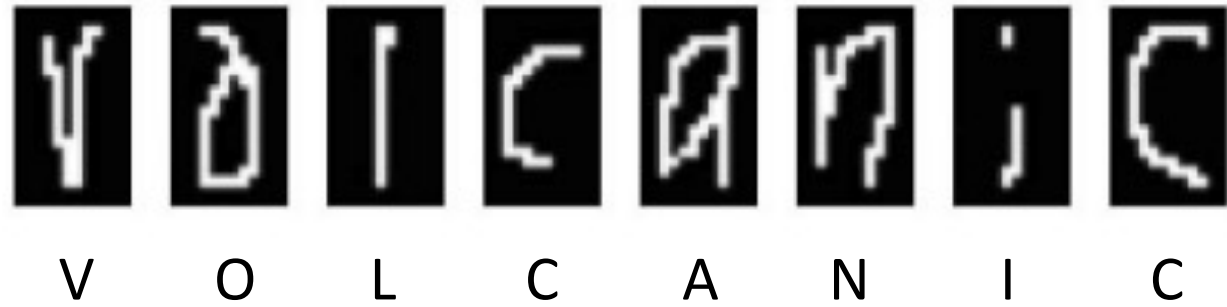
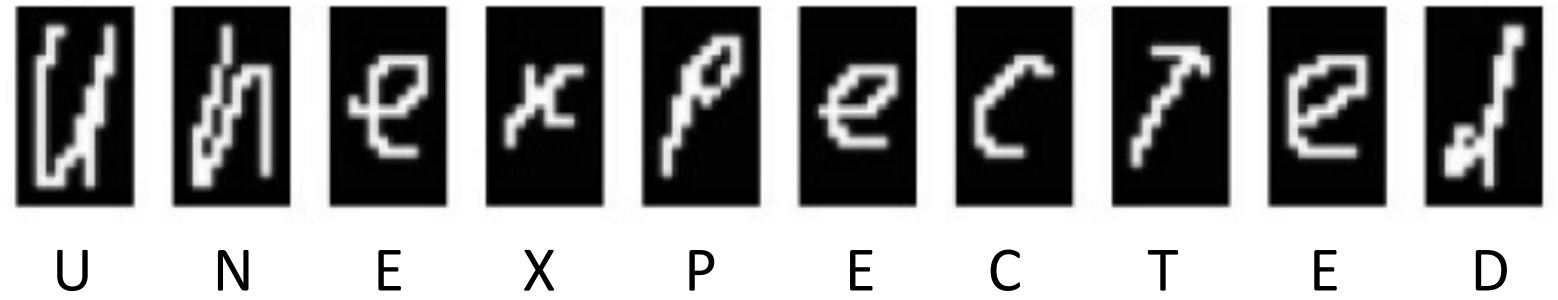
Front Matter

- Announcements
 - HW6 released 6/25, due 7/2 (tomorrow) at 11:59 PM
 - HW7 released 7/2 (tomorrow), due **7/11** at 11:59 PM
 - **Recitation on Wednesday (7/3)**
 - **No class on Thursday (7/4)**
- Recommended Readings
 - Zhang, Lipton, Li & Smola, [Chapters 9 & 10](#)

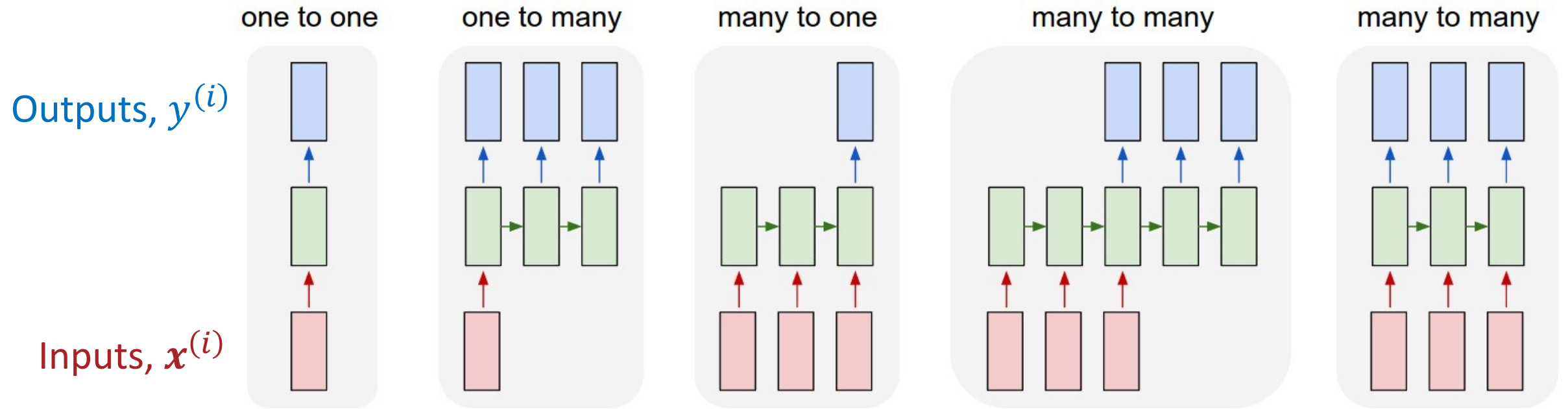
Recurrent Neural Networks

- Neural networks are frequently applied to inputs with some inherent temporal or sequential structure (e.g., text or video) of variable length
- Idea: use the information from previous parts of the input to inform subsequent predictions
- Insight: the hidden layers learn a useful representation (relative to the task)
- Approach: incorporate the output from earlier hidden layers into later ones.

Example: Handwriting Recognition



$$\mathbf{y}^{(i)} = [\mathbf{y}_1^{(i)}, \mathbf{y}_2^{(i)}, \dots, \mathbf{y}_{T_i}^{(i)}]$$



Inputs, $\mathbf{x}^{(i)}$

Outputs, $\mathbf{y}^{(i)}$

$$\mathbf{x}^{(i)} = [\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}]$$

Sequential Data

Which of the following are ways that you could reasonably formulate the handwritten character recognition task?

One-to-one

0%

One-to-many

0%

Many-to-one

0%

Many-to-many (offset)

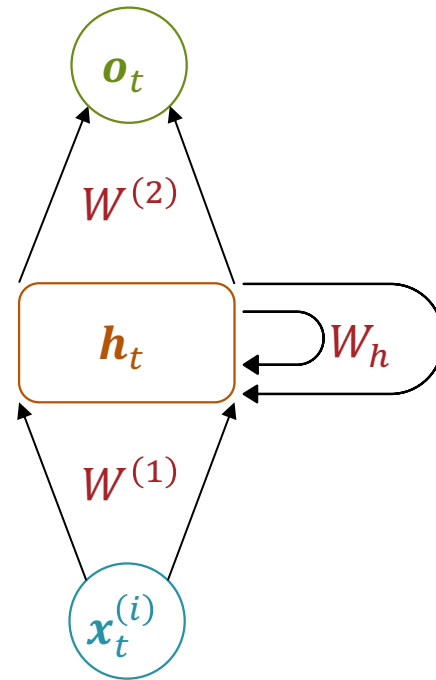
0%

Many-to-many (synced)

0%

Recurrent Neural Networks

$$\mathbf{h}_t = \left[1, \theta \left(W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left(W^{(2)} \mathbf{h}_t \right)$$



- Training dataset consists of (input **sequence**, label **sequence**) pairs, potentially of varying lengths

$$\mathcal{D} = \left\{ \left(\mathbf{x}^{(n)}, \mathbf{y}^{(n)} \right) \right\}_{n=1}^N$$

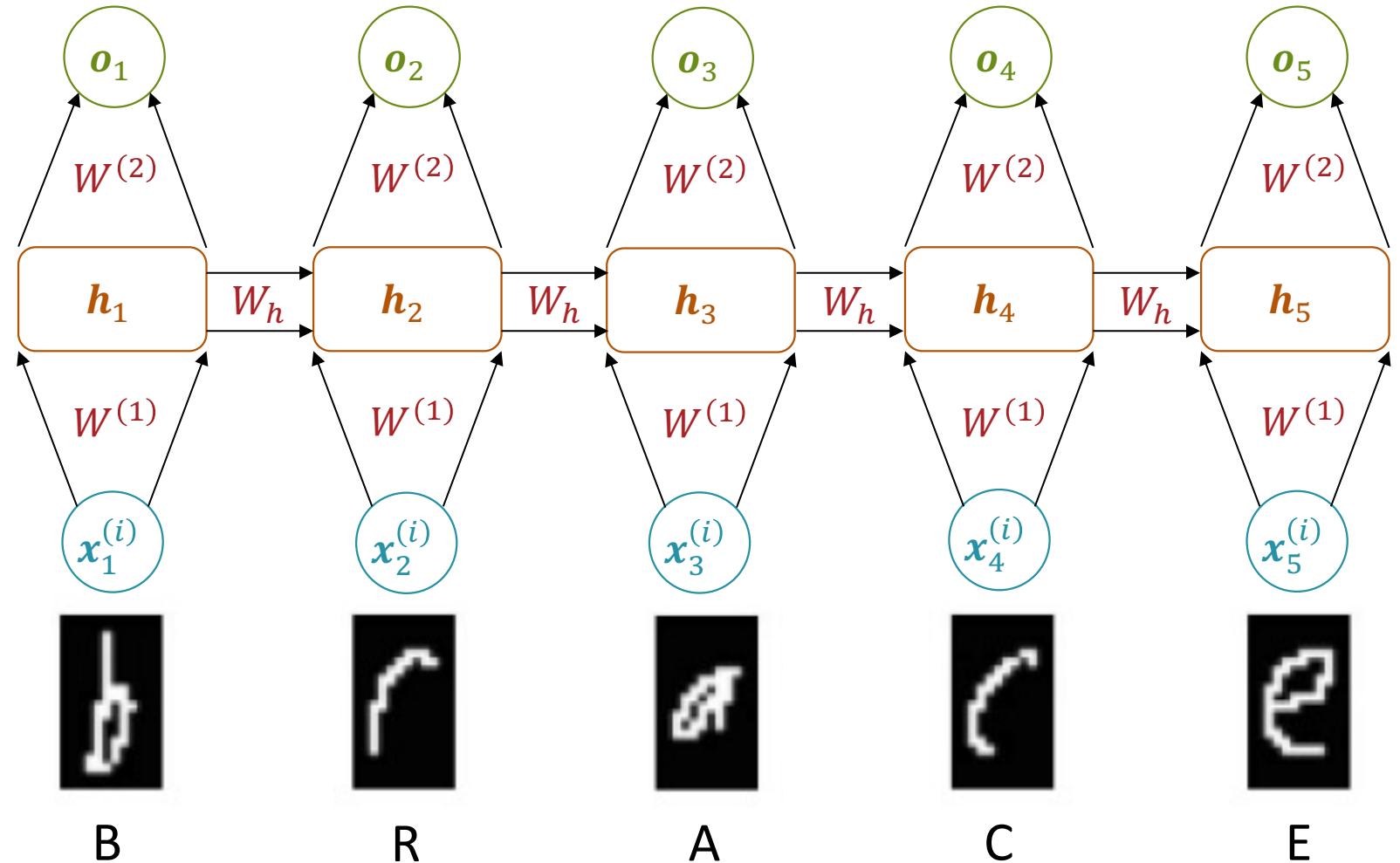
$$\mathbf{x}^{(n)} = \left[\mathbf{x}_1^{(n)}, \dots, \mathbf{x}_{T_n}^{(n)} \right]$$

$$\mathbf{y}^{(n)} = \left[\mathbf{y}_1^{(n)}, \dots, \mathbf{y}_{T_n}^{(n)} \right]$$

- This model requires an initial value for the hidden representation, \mathbf{h}_0 , typically a vector of all zeros

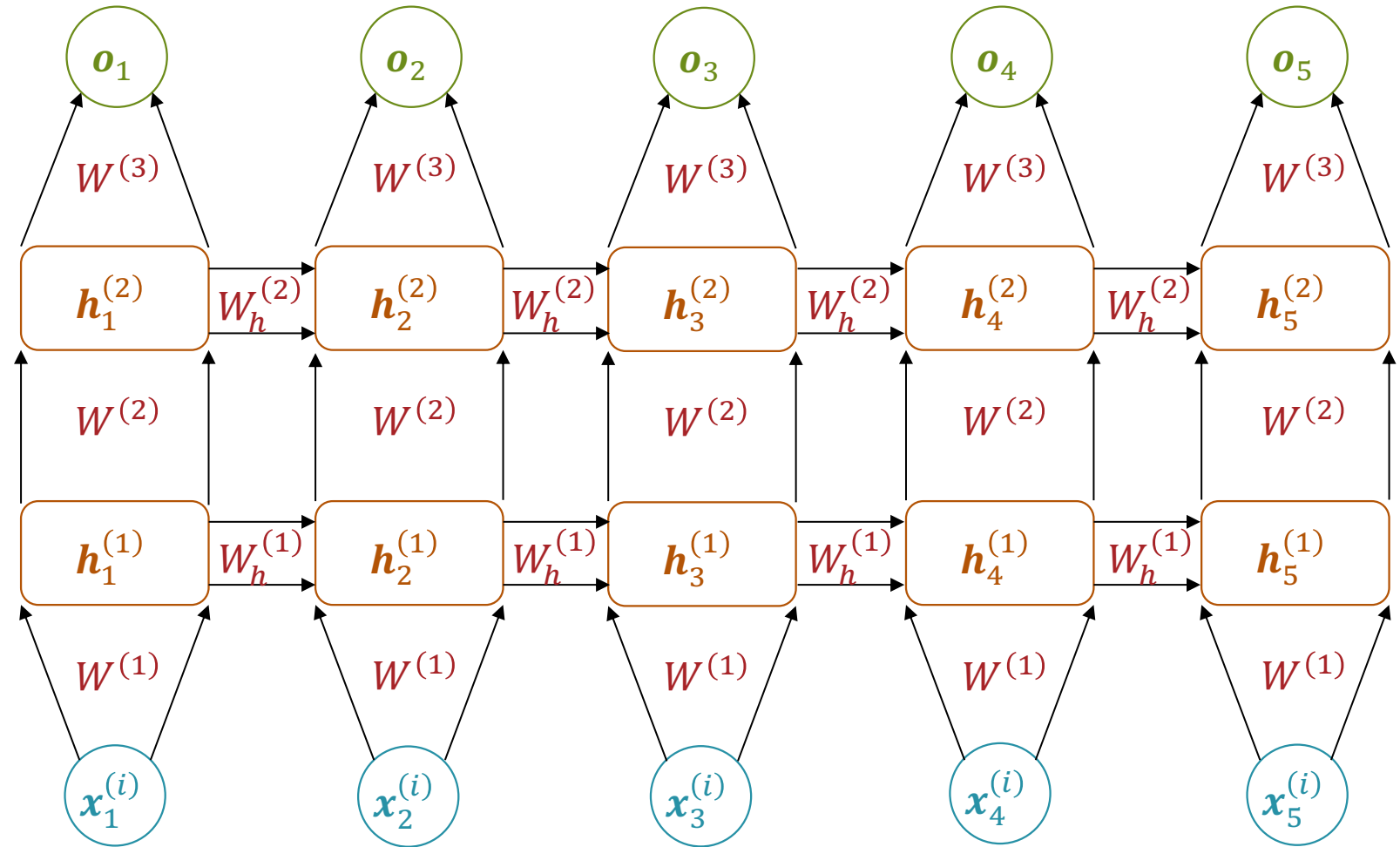
Unrolling Recurrent Neural Networks

$$\mathbf{h}_t = \left[1, \theta \left(W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left(W^{(2)} \mathbf{h}_t \right)$$



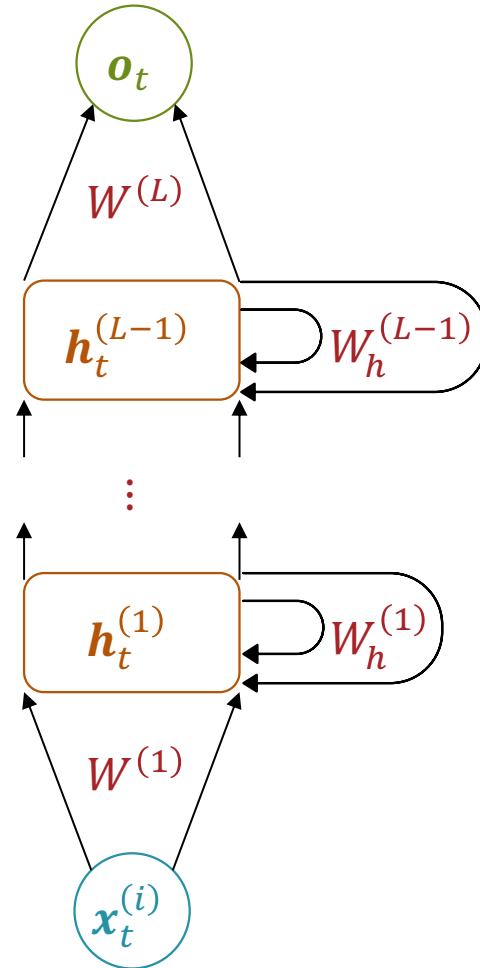
Deep Recurrent Neural Networks

$$\mathbf{h}_t^{(l)} = \left[1, \theta \left(W^{(l)} \mathbf{h}_t^{(l-1)} + W_h^{(l)} \mathbf{h}_{t-1}^{(l)} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left(W^{(L)} \mathbf{h}_t^{(L-1)} \right)$$



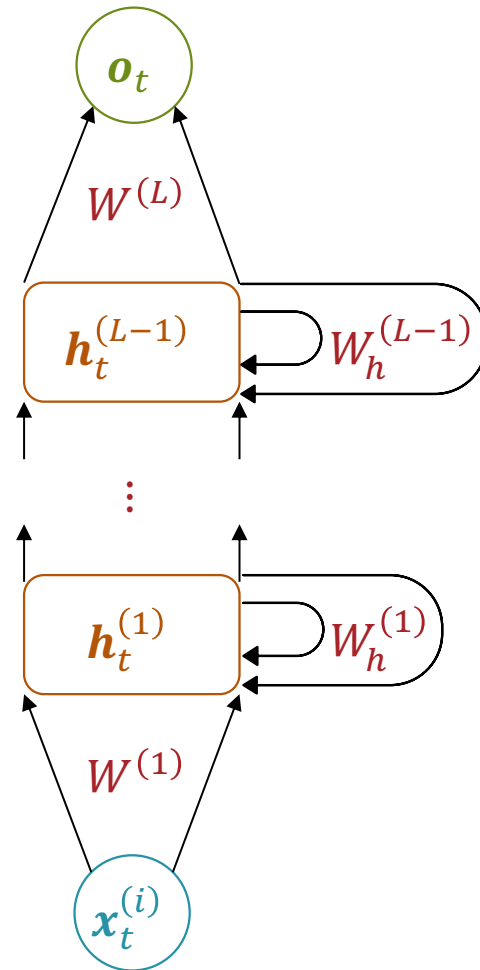
Deep Recurrent Neural Networks

$$\mathbf{h}_t^{(l)} = \left[1, \theta \left(W^{(l)} \mathbf{h}_t^{(l-1)} + W_h^{(l)} \mathbf{h}_{t-1}^{(l)} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left(W^{(L)} \mathbf{h}_t^{(L-1)} \right)$$



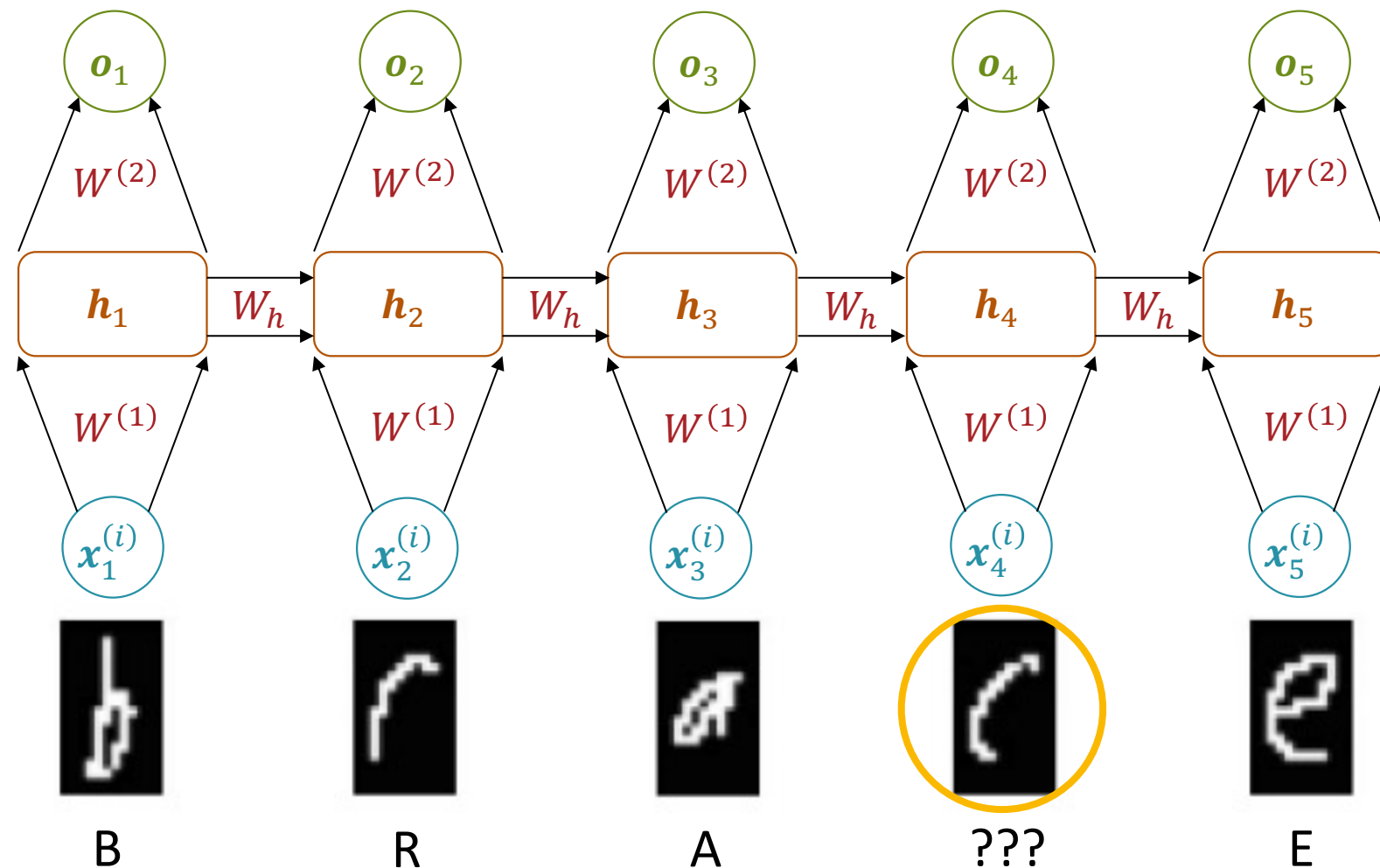
But why do we only pass information forward?
What if later time steps have useful information as well?

$$\mathbf{h}_t^{(l)} = \left[1, \theta \left(W^{(l)} \mathbf{h}_t^{(l-1)} + W_h^{(l)} \mathbf{h}_{t-1}^{(l)} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left(W^{(L)} \mathbf{h}_t^{(L-1)} \right)$$



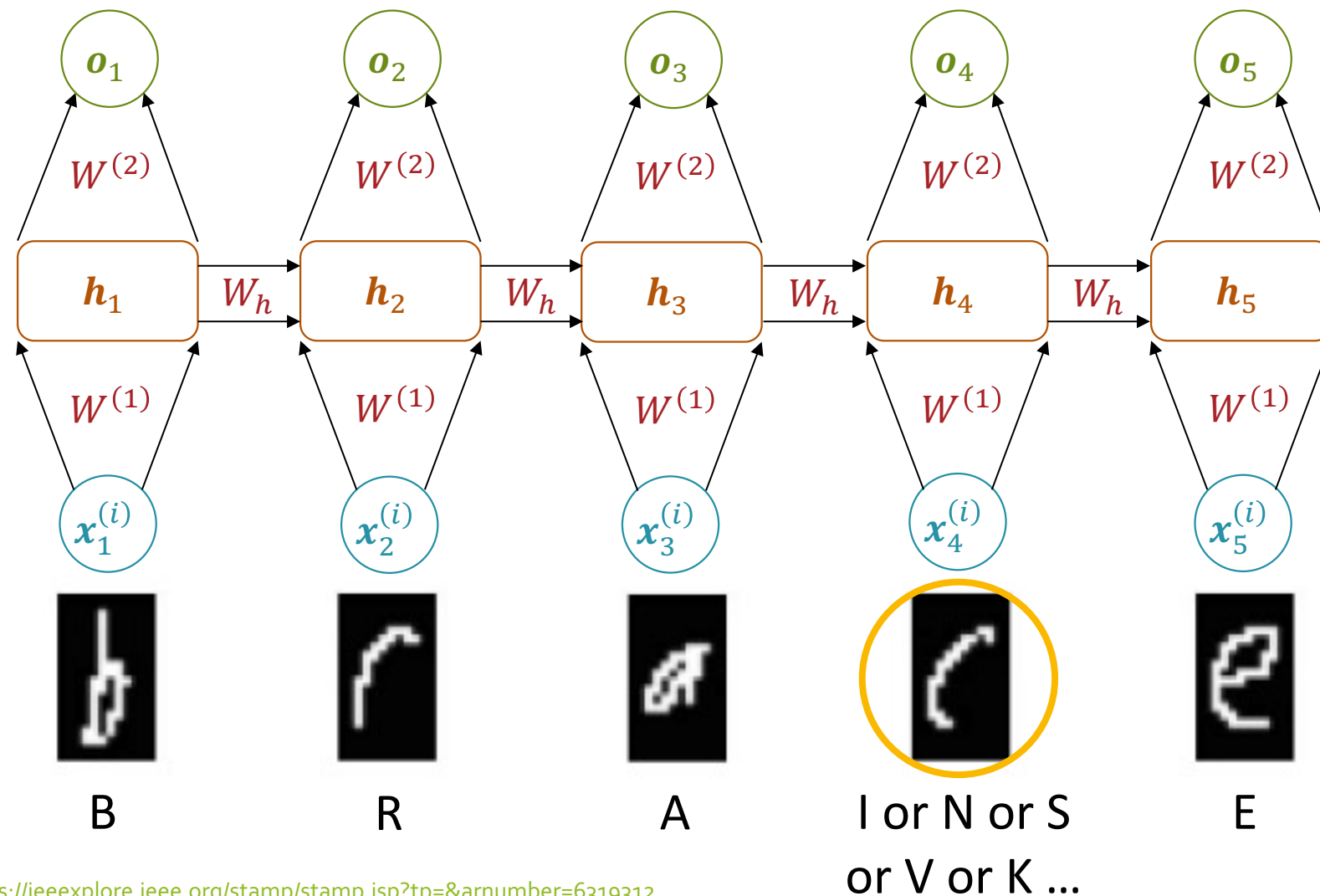
But why do we only pass information forward?
What if later time steps have useful information as well?

$$\mathbf{h}_t = \left[1, \theta \left(W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left(W^{(2)} \mathbf{h}_t \right)$$



But why do we only pass information forward?
 What if later time steps have useful information as well?

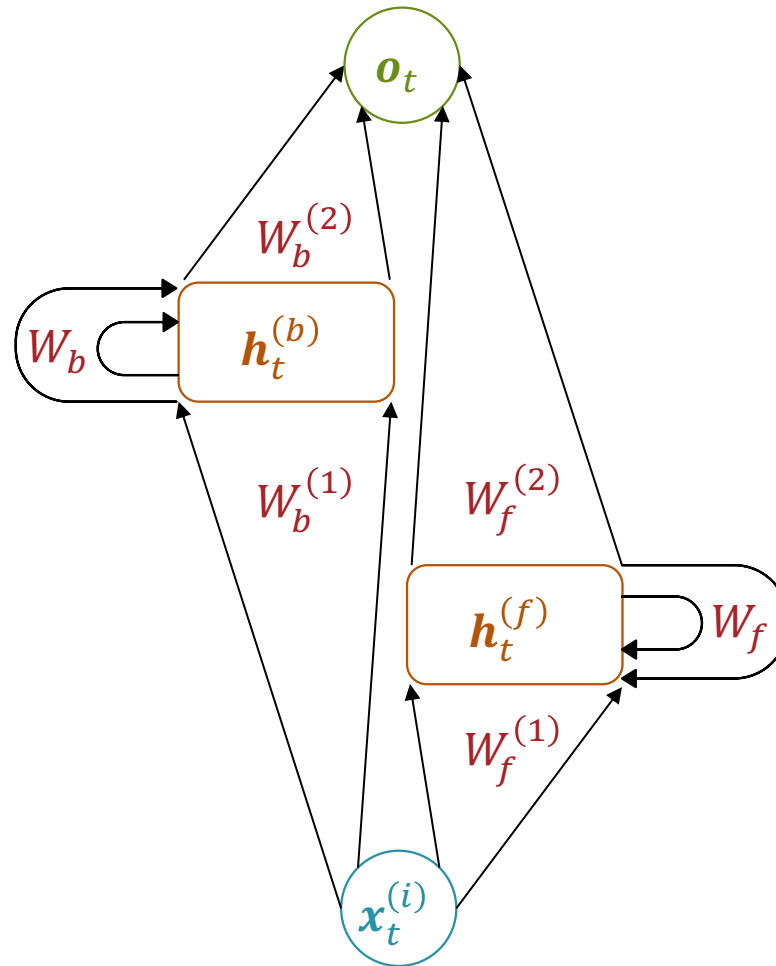
$$\mathbf{h}_t = \left[1, \theta \left(W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left(W^{(2)} \mathbf{h}_t \right)$$



Bidirectional Recurrent Neural Networks

$$\mathbf{h}_t^{(f)} = \left[1, \theta \left(W_f^{(1)} \mathbf{x}_t^{(i)} + W_f \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{h}_t^{(b)} = \left[1, \theta \left(W_b^{(1)} \mathbf{x}_t^{(i)} + W_b \mathbf{h}_{t+1} \right) \right]^T$$

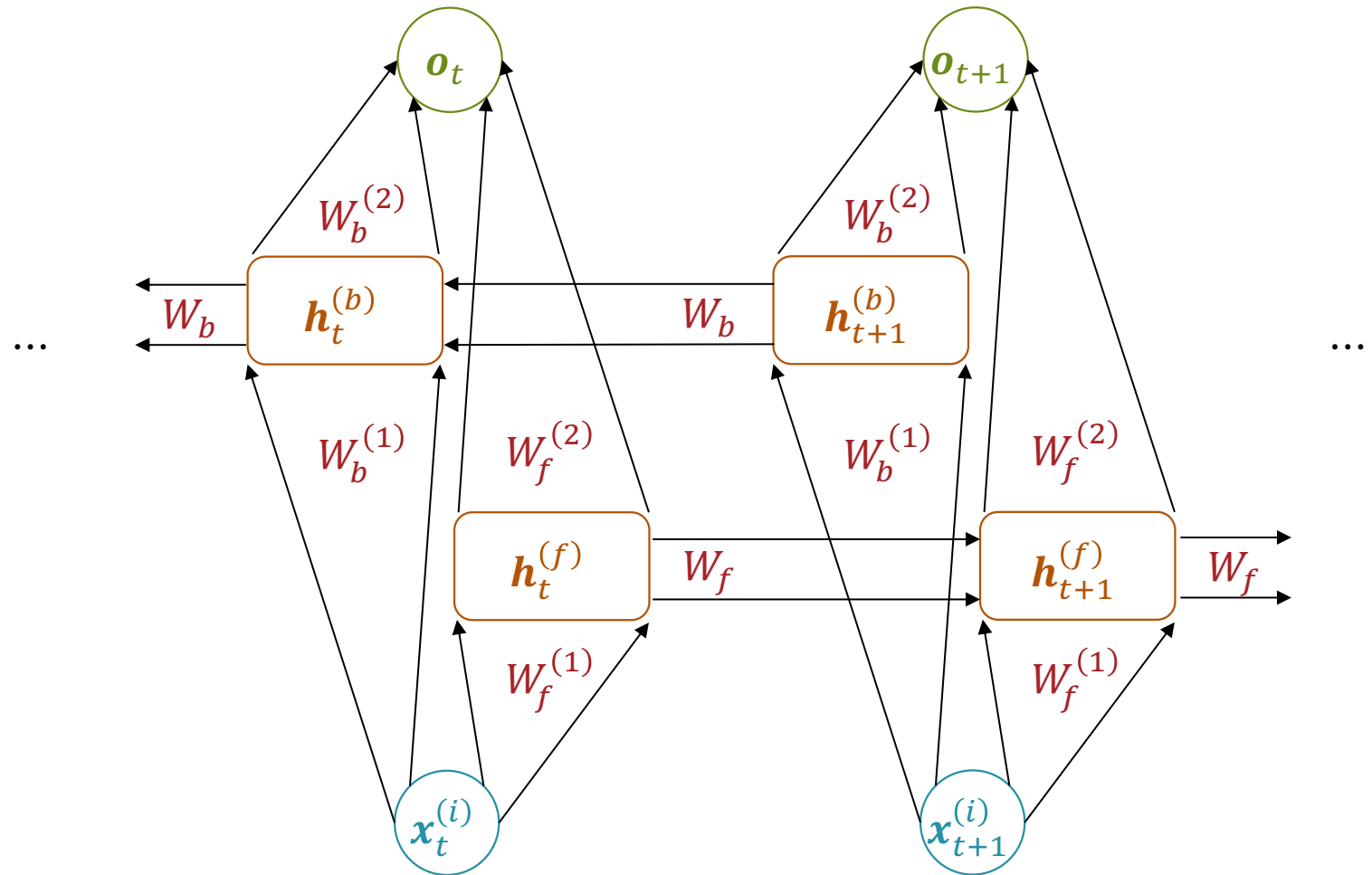
$$\mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left(W_f^{(2)} \mathbf{h}_t^{(f)} + W_b^{(2)} \mathbf{h}_t^{(b)} \right)$$



Unrolling Bidirectional Recurrent Neural Networks

$$o_t = \hat{y}_t^{(i)} = \theta \left(W_f^{(2)} h_t^{(f)} + W_b^{(2)} h_t^{(b)} \right)$$

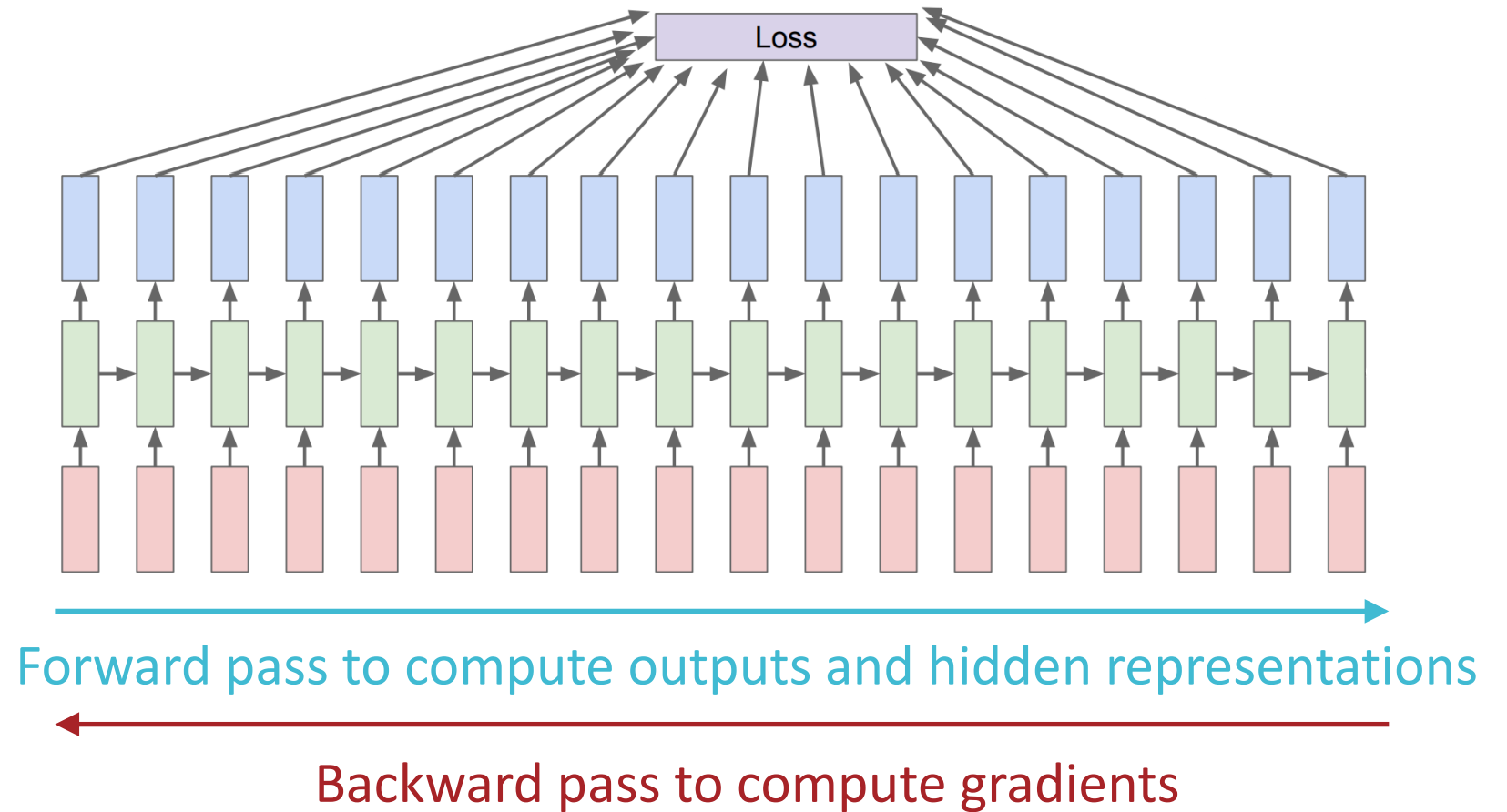
$$h_t^{(f)} = \left[1, \theta \left(W_f^{(1)} x_t^{(i)} + W_f h_{t-1}^{(f)} \right) \right]^T \text{ and } h_t^{(b)} = \left[1, \theta \left(W_b^{(1)} x_t^{(i)} + W_b h_{t+1}^{(b)} \right) \right]^T$$



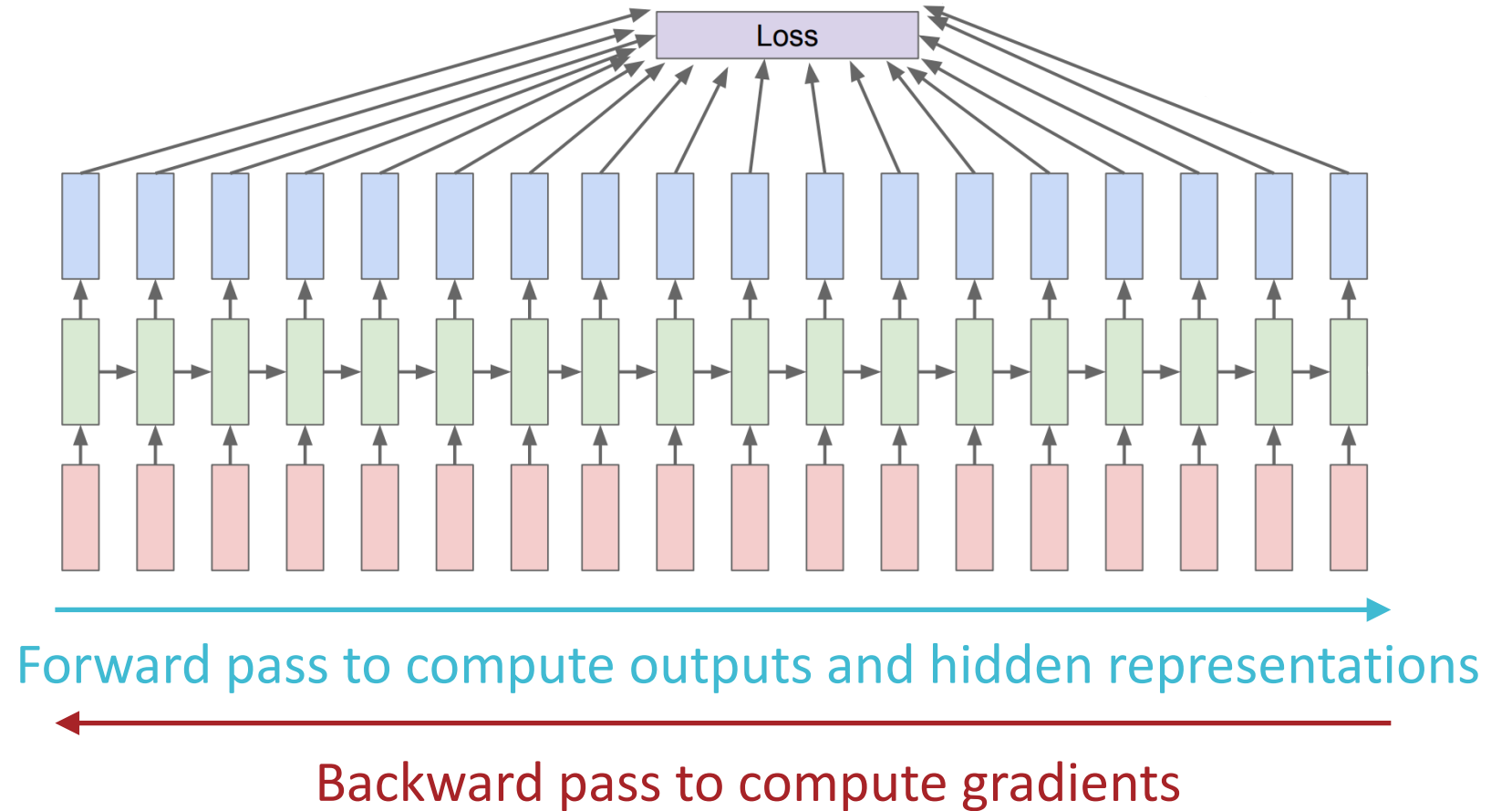
Training RNNs

- A (deep/bidirectional) RNN simply represents a (somewhat complicated) computation graph
 - Weights are shared between different timesteps, significantly reducing the number of parameters to be learned!
- Can be trained using (stochastic) gradient descent/backpropagation → “backpropagation through time”

Training RNNs



Training RNNs: Challenges

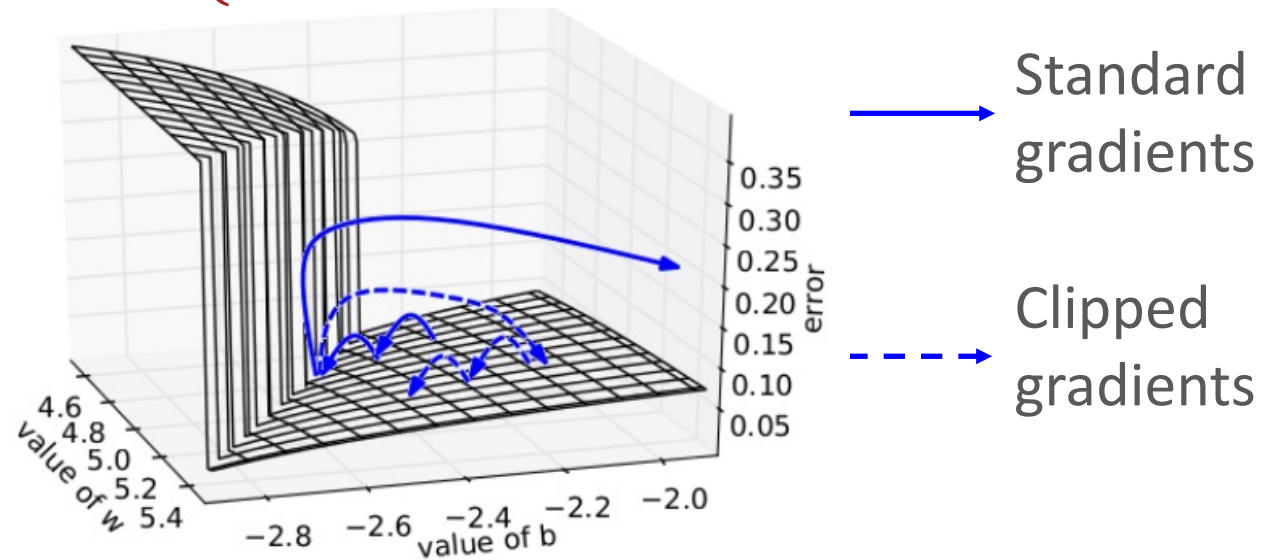


- Issue: as the sequence length grows, the gradient is more likely to explode or vanish

Gradient Clipping (Pascanu et al., 2013)

- Common strategy to deal with exploding gradients: if the magnitude of the gradient ever exceeds some threshold, simply scale it down to the threshold

$$G = \begin{cases} \nabla_W \ell^{(i)} & \text{if } \|\nabla_W \ell^{(i)}\|_2 \leq \tau \\ \left(\frac{\tau}{\|\nabla_W \ell^{(i)}\|_2} \right) \nabla_W \ell^{(i)} & \text{otherwise} \end{cases}$$



Recall: Computing Gradients


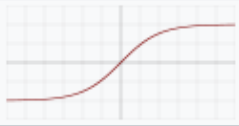

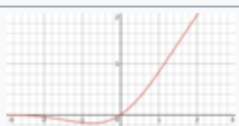



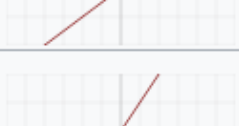
Insight: $s_b^{(l)}$ only affects $\ell^{(i)}$ via $o_b^{(l)}$

$$\text{Chain rule: } \delta_b^{(l)} = \frac{\partial \ell^{(i)}}{\partial o_b^{(l)}} \left(\frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} \right)$$

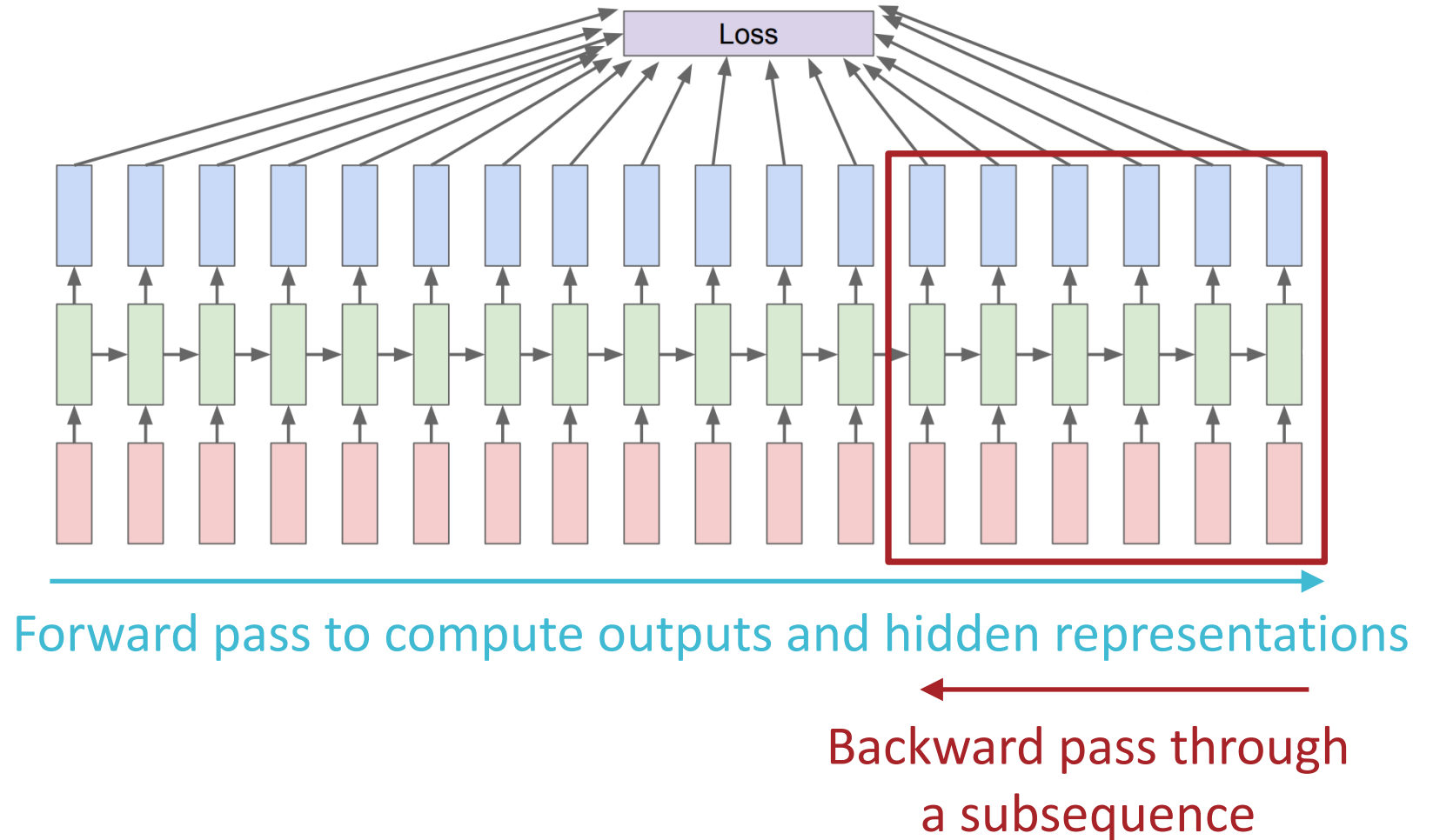
$$\begin{aligned} o_b^{(l)} = \theta \left(s_b^{(l)} \right) &\rightarrow \frac{\partial o_b^{(l)}}{\partial s_b^{(l)}} = \frac{\partial \theta \left(s_b^{(l)} \right)}{\partial s_b^{(l)}} \\ &= 1 - \left(\tanh \left(s_b^{(l)} \right) \right)^2 \leq 1 \end{aligned}$$

when $\theta(\cdot) = \tanh(\cdot)$

Recall: Activation Functions

Logistic, sigmoid, or soft step		$\sigma(x) = \frac{1}{1 + e^{-x}}$
Hyperbolic tangent (\tanh)		$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
Rectified linear unit (ReLU) ^[7]		$\begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max\{0, x\} = x \mathbf{1}_{x>0}$
Gaussian Error Linear Unit (GELU) ^[4]		$\frac{1}{2}x \left(1 + \operatorname{erf} \left(\frac{x}{\sqrt{2}} \right) \right)$ $= x\Phi(x)$
Softplus ^[8]		$\ln(1 + e^x)$
Exponential linear unit (ELU) ^[9]		$\begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ with parameter α
Leaky rectified linear unit (Leaky ReLU) ^[11]		$\begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$
Parametric rectified linear unit (PReLU) ^[12]		$\begin{cases} \alpha x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$ with parameter α

Truncated Backpropagation Through Time



Forward pass to compute outputs and hidden representations

Backward pass through a subsequence

- Idea: limit the number of time steps to backprop through

Language Models

1. Convert raw text into sequence data

$$\mathbf{x}^{(i)} = [\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}]$$

2. Learn or approximate a joint probability distribution over sequences

$$P(\mathbf{x}^{(i)}) = P(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)})$$

3. Sample from the implied conditional distribution to generate new sequences

$$P(\mathbf{x}_{T_i+1} | \mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}) = \frac{P(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}, \mathbf{x}_{T_i+1})}{P(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)})}$$

Tokenization and Embedding

1. Convert raw text into sequence data

$$\mathbf{x}^{(i)} = [\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}]$$

- High-level approach: split raw text into smaller units (“tokens”), then learn a dense, numerical vector representation (“embedding”) for each token

Tokenization

- Example: “Henry is giving a lecture on language models”
- Idea: word-based tokenization
[“henry”, “is”, “giving”, “a”, “lecture”, “on”, “language”,
“models”]

Tokenization

- Example: “Henry is givin’ a lectrue on LMs”
- Idea: word-based tokenization?

[“henry”, “is”, “??”, “a”, “??”, “on”, “???”]

- Can have difficulty trading off between vocabulary size and computational tractability
- Similar words e.g., “model” and “models” can get mapped to completely disparate representations
- Typos or acronyms will likely be out-of-vocabulary (OOV)

Tokenization

- Example: “Henry is givin’ a lectrue on LMs”

- Idea: character-based tokenization:

[“h”, “e”, “n”, “r”, “y”, “i”, “s”, “g”, “i”, “v”, “i”, “n”, “ ’ ”, ...]

- Much smaller vocabularies but a lot of semantic meaning is lost...
- Sequences will be much longer than word-based tokenization, potentially causing computational issues
- Can do well on logographic languages e.g., Kanji 漢字

Tokenization

- Example: “Henry is givin’ a lectrue on LMs”

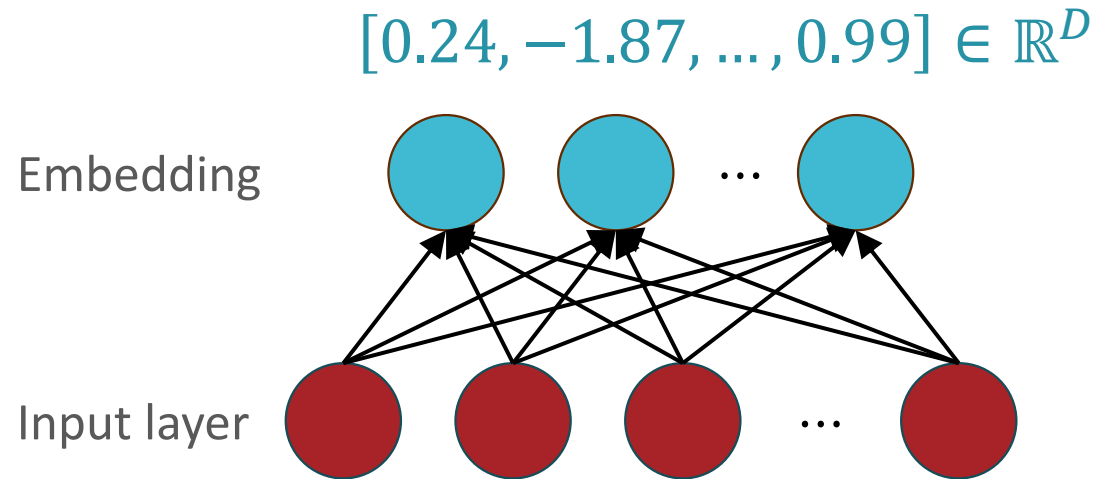
- Common practice: subword tokenization

[“henry”, “is”, “giv”, “##in”, “ ‘ ”, “a”, “lect” “##re”, “on”, “language”, “model”, “#s”]

- Split long or rare words into smaller, semantically meaningful components or *subwords*

Embedding

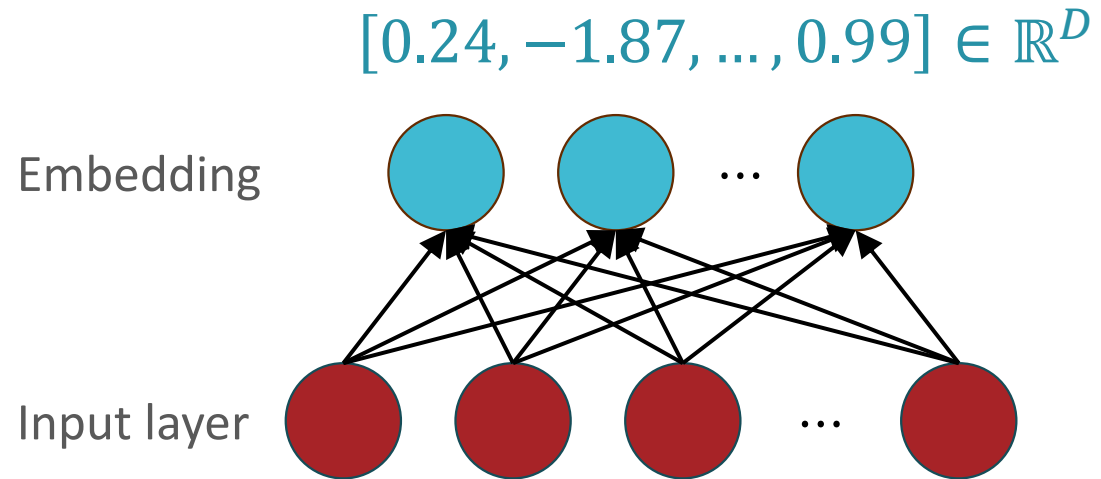
- Given a vocabulary V with $|V|$ tokens, learn an embedding by training a 1-layer, fully-connected feed-forward NN that takes one-hot encoded vectors as input



- Example: “is” $\rightarrow [0, 0, 1, \dots, 0] \in \mathbb{R}^{|V|}$

Okay but how do I go about training this neural network?

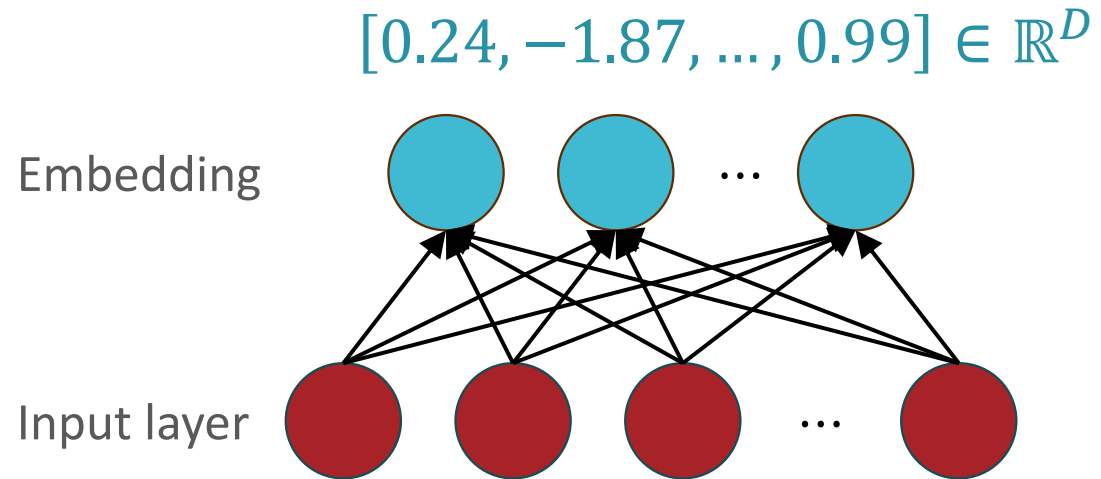
- Given a vocabulary V with $|V|$ tokens, learn an embedding by training a 1-layer, fully-connected feed-forward NN that takes one-hot encoded vectors as input



- Example: "is" $\rightarrow [0, 0, 1, \dots, 0] \in \mathbb{R}^{|V|}$

Embedding

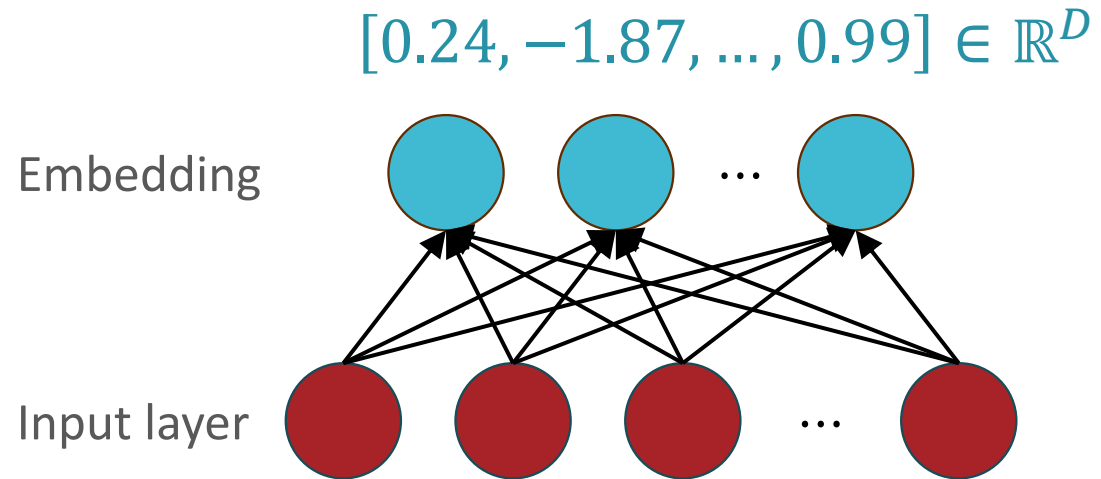
- Given a vocabulary V with $|V|$ tokens, learn an embedding by training a 1-layer, fully-connected feed-forward NN that takes one-hot encoded vectors as input



- Example: “is” $\rightarrow [0, 0, 1, \dots, 0] \in \mathbb{R}^{|V|}$
- Idea: use a pretrained embedding e.g., [word2vec](#) or [GloVe](#)
 - Requires you to use the same vocabulary/tokenization

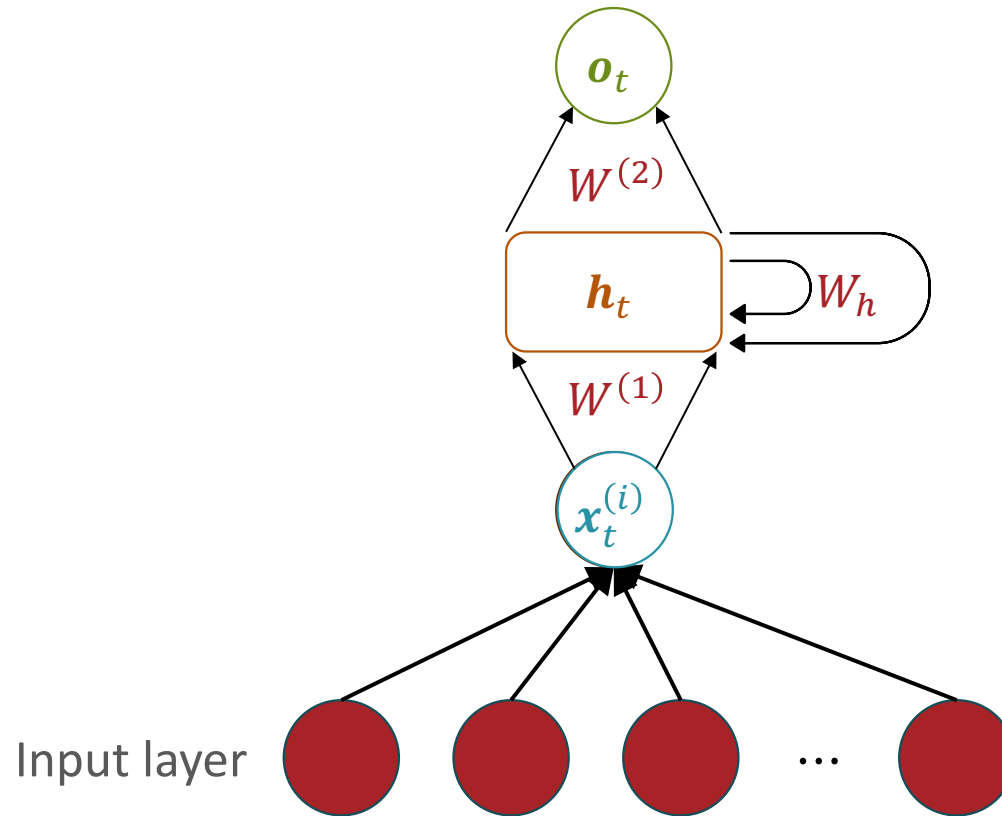
Embedding Layer

- Given a vocabulary V with $|V|$ tokens, learn an embedding by training a 1-layer, fully-connected feed-forward NN that takes one-hot encoded vectors as input



- Example: “is” $\rightarrow [0, 0, 1, \dots, 0] \in \mathbb{R}^{|V|}$
- Common practice: add this 1-layer NN to whatever architecture you’re using and fit it to the task

Embedding Layer



- Example: “is” $\rightarrow [0, 0, 1, \dots, 0] \in \mathbb{R}^{|V|}$
- Common practice: add this 1-layer NN to whatever architecture you’re using and fit it to the task

Key Takeaways

- Recurrent neural networks use contextual information to reason about sequential data
 - Can still be learned using backpropagation → backpropagation through time
 - Susceptible to exploding/vanishing gradients for long training sequences
- Language models fit joint probability distributions to sequences of tokens
 - Tokenization and embedding to generate dense vector representations of texts
 - Can be sampled from to generate text