# Lecture 10:
# Sequence-to-Sequence Modeling with Encoder-Decoder Architectures

*Instructor: Swabha Swayamdipta*
*USC CSCI 544 Applied NLP*
*Sep 26, Fall 2024*

Some slides adapted from Dan Jurafsky and Chris Manning

# Announcements

- Tonight at 11:59 PM PT: Project Proposal Due
  - Once you propose an idea, you're NOT allowed to completely change it
  - Allowed to make modifications, based on our recommendations
- Next week: Quiz 3
  - Before that: Install Lockdown Browser
  - Cannot take Quiz 3 otherwise
  - Free quiz points!!
- Quiz grades: Please see your total grades and do not be confused by Brightspace options for the correct answer
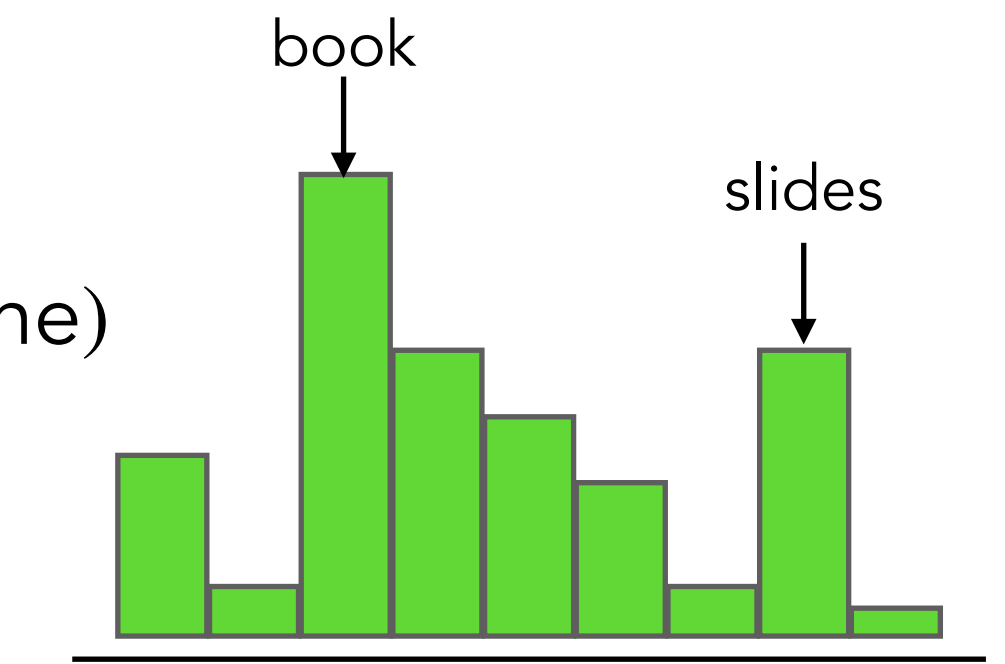
# Lecture Outline

- Announcements
- Recap: Recurrent Neural Nets
- Applications of RNNs
- Seq2Seq Modeling with Encoder-Decoder Networks
- Attention Mechanism
- More on Attention
- Transformers: Self-Attention Networks

# Recap:
# Recurrent Neural Nets

# Recurrent Neural Net Language Models

$$\hat{y}_4 = P(x_5 | \text{The students studied the})$$

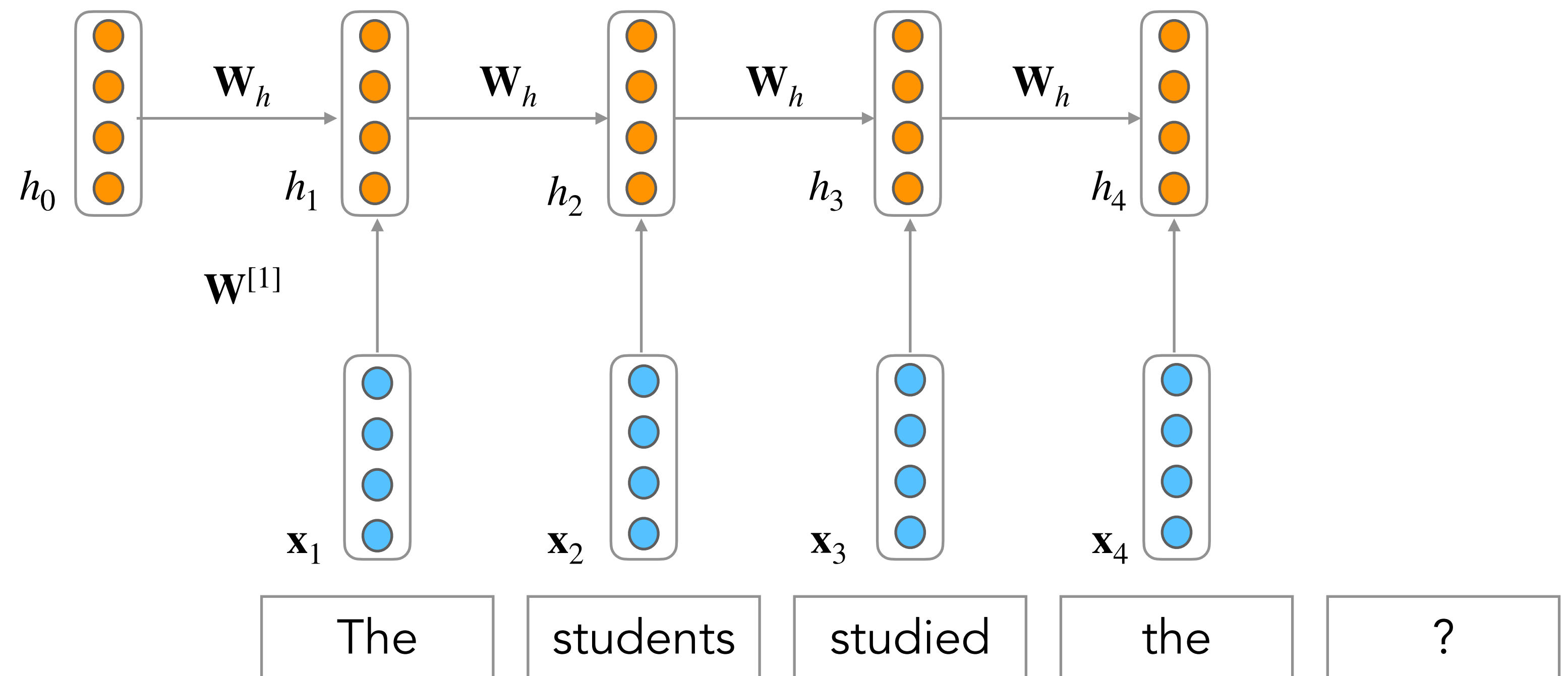Output layer: $\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{W}^{[2]}\mathbf{h}_t)$

Hidden layer:

$$\mathbf{h}_t = g(\mathbf{W}_h\mathbf{h}_{t-1} + \mathbf{W}^{[1]}\mathbf{x}_t)$$

Initial hidden state: $\mathbf{h}_0$

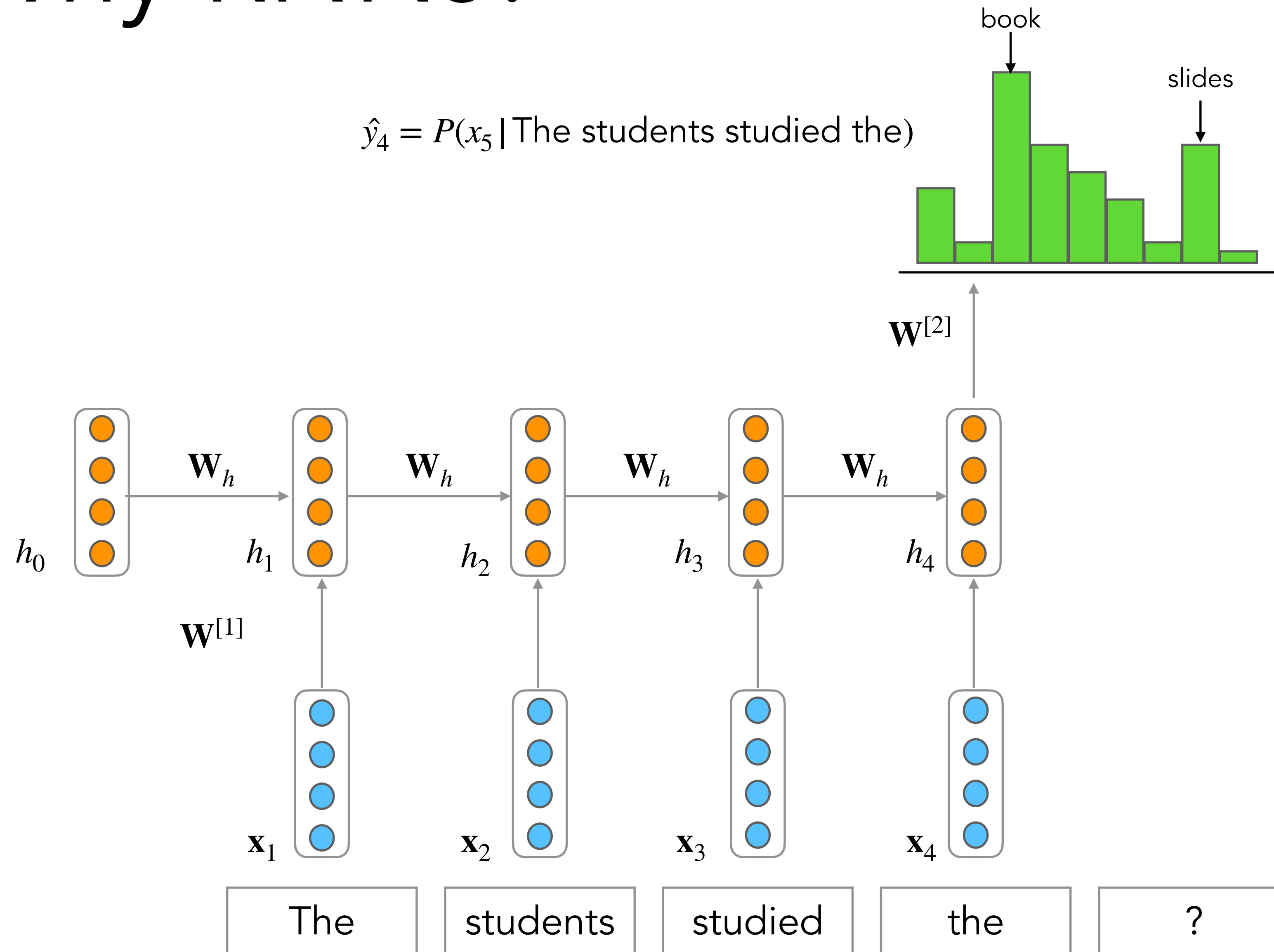Word Embeddings, $\mathbf{x}_i$



book

slides

$\mathbf{W}^{[2]}$

$\mathbf{W}_h$   $\mathbf{W}_h$   $\mathbf{W}_h$   $\mathbf{W}_h$

$h_0$   $h_1$   $h_2$   $h_3$   $h_4$

$\mathbf{W}^{[1]}$

$\mathbf{x}_1$   $\mathbf{x}_2$   $\mathbf{x}_3$   $\mathbf{x}_4$

| The | students | studied | the | ? |

5

# Why RNNs?

RNN Advantages:

- Can process any length input
- Model size doesn't increase for longer input
- Computation for step $t$ can (in theory) use information from many steps back
- Weights $\mathbf{W}_h$ are shared (tied) across timesteps → Condition the neural network on all previous words
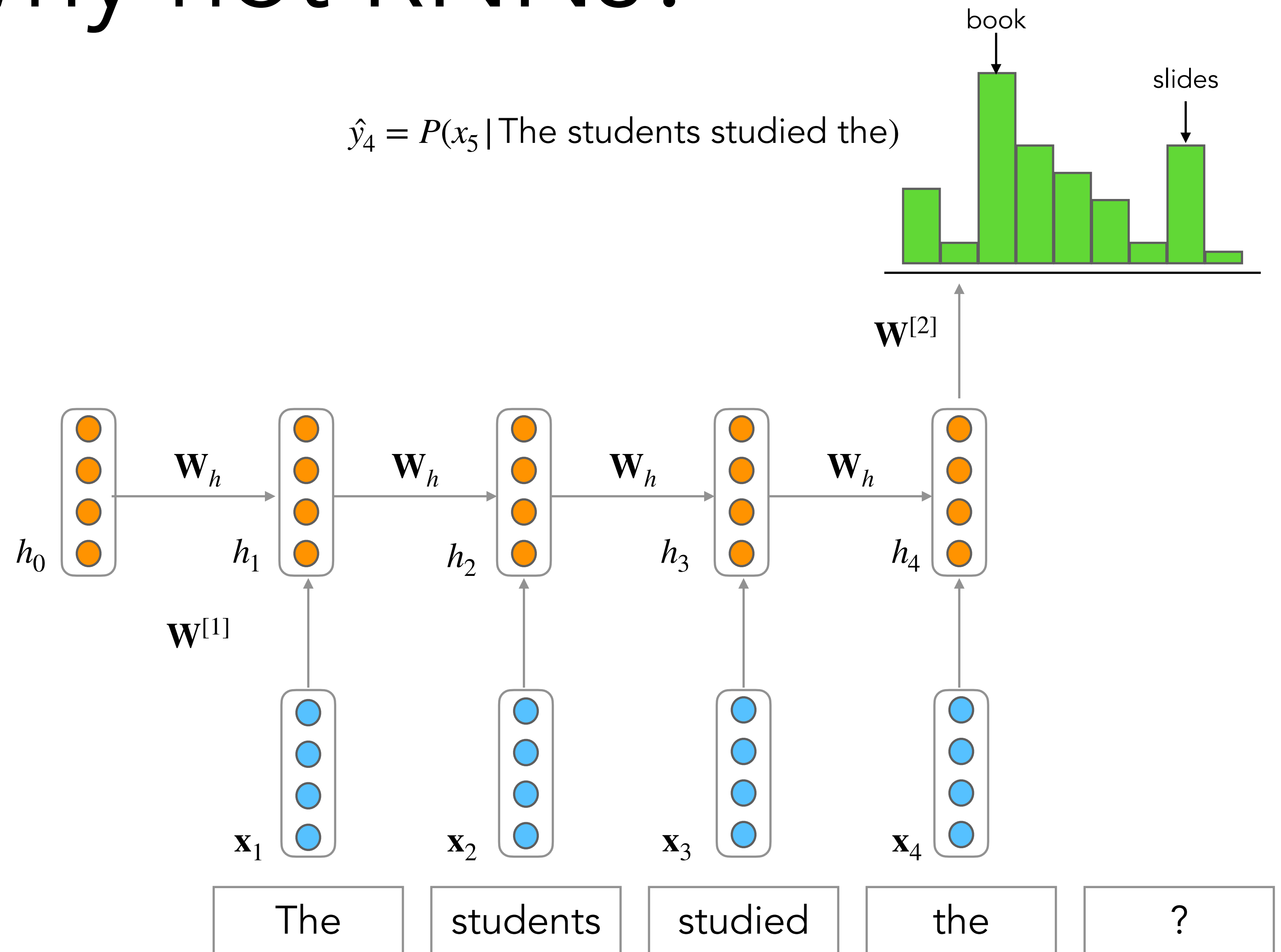- Only need to save previous hidden state in memory (as opposed to an entire window)

$\hat{y}_4 = P(x_5 \,|\, \text{The students studied the})$

book

slides

$\mathbf{W}^{[2]}$

$h_0$    $\mathbf{W}_h$   $h_1$   $\mathbf{W}_h$   $h_2$   $\mathbf{W}_h$   $h_3$   $\mathbf{W}_h$   $h_4$

$\mathbf{W}^{[1]}$

$\mathbf{x}_1$   $\mathbf{x}_2$   $\mathbf{x}_3$   $\mathbf{x}_4$

| The | students | studied | the | ? |

# Why not RNNs?

$$\hat{y}_4 = P(x_5 \mid \text{The students studied the})$$

RNN Disadvantages:
- Recurrent computation is slow
- In practice, difficult to access information from many steps back

# Training Outline

- Get a big corpus of text which is a sequence of words $x_1, x_2, \ldots x_T$
- Feed into RNN-LM; compute output distribution $\hat{y}_t$ for every step $t$
  - i.e. predict probability distribution of every word, given words so far
- Loss function on step $t$ is usual cross-entropy between our predicted probability distribution $\hat{y}_t$, and the true next word $y_t = x_{t+1}$:

$$L_{CE}(\hat{y}_t, y_t; \theta) = -\sum_{v \in V} \mathbb{I}[y_t = v] \log \hat{y}_t = -\log p_\theta(x_{t+1} \mid x_{\leq t})$$

- Average this to get overall loss for entire training set:

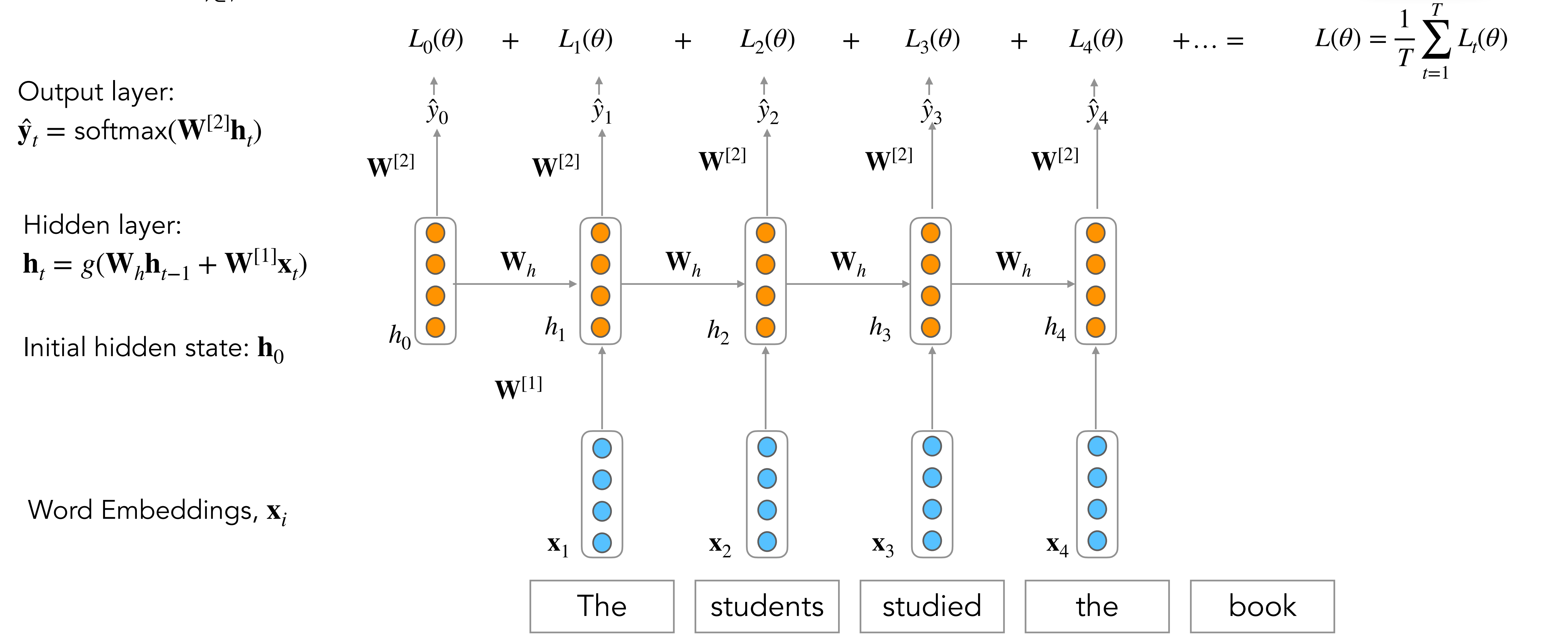$$L(\theta) = \frac{1}{T} \sum_{t=1}^{T} L_{CE}(\hat{y}_t, y_t)$$

USC Viterbi

Cross-Entropy Loss,

$$L_{CE}(\hat{\mathbf{y}}_t, y_t; \theta) = -\sum_{v \in V} \mathbb{I}[y_t = v] \log \hat{y}_t;$$

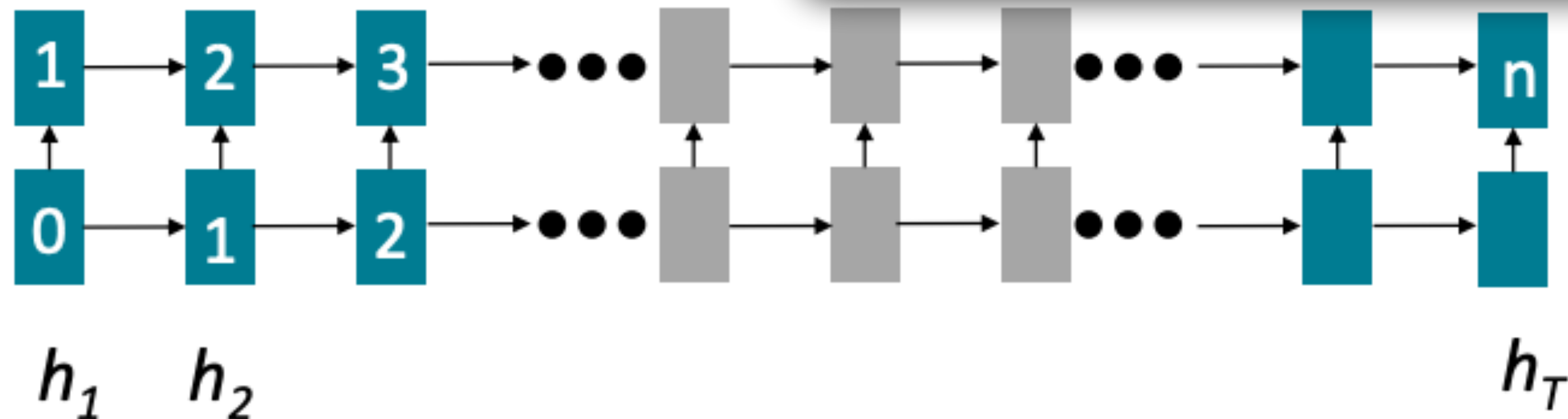$$\theta = [\mathbf{x}; \mathbf{W}^{[1]}; \mathbf{W}_h; \mathbf{W}^{[2]}]$$

Loss

$$L_0(\theta) \quad + \quad L_1(\theta) \quad + \quad L_2(\theta) \quad + \quad L_3(\theta) \quad + \quad L_4(\theta) \quad + \ldots = \qquad L(\theta) = \frac{1}{T}\sum_{t=1}^{T} L_t(\theta)$$

Output layer:

$$\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{W}^{[2]}\mathbf{h}_t)$$

$\hat{y}_0 \qquad \hat{y}_1 \qquad \hat{y}_2 \qquad \hat{y}_3 \qquad \hat{y}_4$

$\mathbf{W}^{[2]} \qquad \mathbf{W}^{[2]} \qquad \mathbf{W}^{[2]} \qquad \mathbf{W}^{[2]} \qquad \mathbf{W}^{[2]}$

Hidden layer:

$$\mathbf{h}_t = g(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}^{[1]}\mathbf{x}_t)$$

$\mathbf{W}_h \qquad \mathbf{W}_h \qquad \mathbf{W}_h \qquad \mathbf{W}_h$

$h_0 \qquad h_1 \qquad h_2 \qquad h_3 \qquad h_4$

Initial hidden state: $\mathbf{h}_0$

$\mathbf{W}^{[1]}$

Word Embeddings, $\mathbf{x}_i$

$\mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3 \qquad \mathbf{x}_4$

| The | students | studied | the | book |

9

# Training RNNs is hard: Parallelizability

- Forward and backward passes have **O(sequence length)** unparallelizable operations!
  - GPUs can perform a bunch of independent computations at once!
  - But future RNN hidden states can't be computed in full before past RNN hidden states have been computed
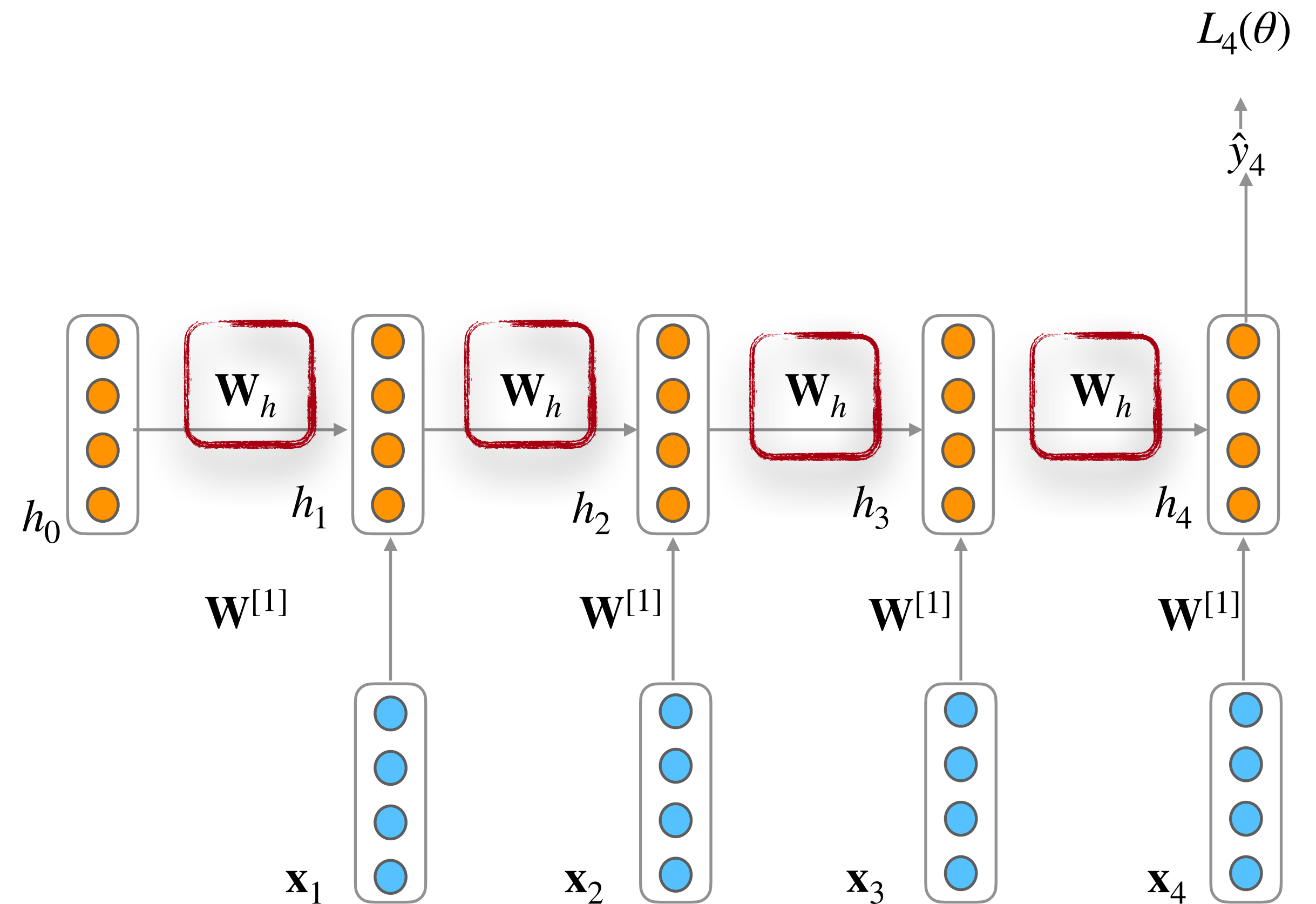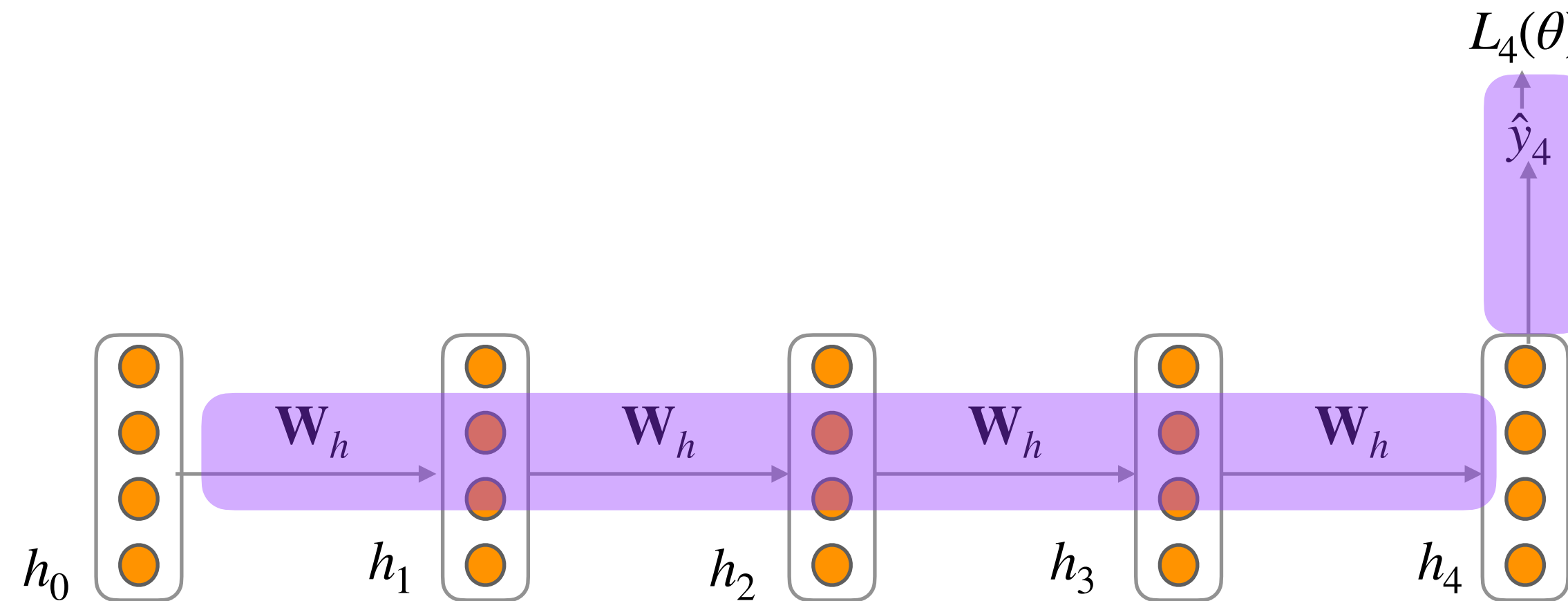
Inhibits training on very large datasets!



Numbers indicate min # of steps before a state can be computed

# Training RNNs is hard: Gradients

- Multiply the same matrix at each time step during forward propagation
  - Advantage: Inputs from many time steps ago can modify output $y$
  - Disadvantage: The **vanishing gradient problem**

# The Vanishing Gradient Problem: Intuition

$L_4(\theta)$

$\hat{y}_4$

$h_0$     $\mathbf{W}_h$     $h_1$     $\mathbf{W}_h$     $h_2$     $\mathbf{W}_h$     $h_3$     $\mathbf{W}_h$     $h_4$

When these gradients are small, the gradient signal gets smaller and smaller as it backpropagates further…

$$\frac{\partial L_4}{\partial h_0} = \qquad \frac{\partial h_1}{\partial h_0} \times \frac{\partial L_4}{\partial h_1}$$

$$= \frac{\partial h_1}{\partial h_0} \times \frac{\partial h_2}{\partial h_1} \times \frac{\partial L_4}{\partial h_2}$$

$$= \frac{\partial h_1}{\partial h_0} \times \frac{\partial h_2}{\partial h_1} \times \frac{\partial h_3}{\partial h_2} \times \frac{\partial L_4}{\partial h_3}$$

$$= \frac{\partial h_1}{\partial h_0} \times \frac{\partial h_2}{\partial h_1} \times \frac{\partial h_3}{\partial h_2} \times \frac{\partial h_4}{\partial h_3} \times \frac{\partial L_4}{\partial h_4}$$

Gradient signal from far away is lost because it's much smaller than gradient signal from close-by

12

# Long Short-Term Memory RNNs (LSTMs)

- At time step $t$, introduces a new cell state $\mathbf{c}_t \in \mathbb{R}^d$
  - In addition to a hidden state $\mathbf{h}_t \in \mathbb{R}^d$
  - The cell stores long-term information (memory)
  - The LSTM can read, erase, and write information from the cell!
    - The cell becomes conceptually rather like RAM in a computer

- The selection of which information is erased/written/read is controlled by three corresponding gates:
  - Input gate $\mathbf{i}_t \in \mathbb{R}^d$, Output gate $\mathbf{o}_t \in \mathbb{R}^d$ and Forget gate $\mathbf{f}_t \in \mathbb{R}^d$
  - Each *element* of the gates can be open (1), closed (0), or somewhere in between
  - The gates are dynamic: their value is computed based on the current context

# LSTMs

Given a sequence of inputs $x_t$ , we will compute a sequence of hidden states $h_t$ and cell states $c_t$

At timestep $t$ :

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

**Sigmoid function:** all gate values are between 0 and 1

$$f^{(t)} = \sigma \left( W_f h^{(t-1)} + U_f x^{(t)} + b_f \right)$$

$$i^{(t)} = \sigma \left( W_i h^{(t-1)} + U_i x^{(t)} + b_i \right)$$

$$o^{(t)} = \sigma \left( W_o h^{(t-1)} + U_o x^{(t)} + b_o \right)$$

**New cell content:** this is the new content to be written to the cell

**Cell state**: erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state**: read ("output") some content from the cell

$$\tilde{c}^{(t)} = \tanh \left( W_c h^{(t-1)} + U_c x^{(t)} + b_c \right)$$

$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

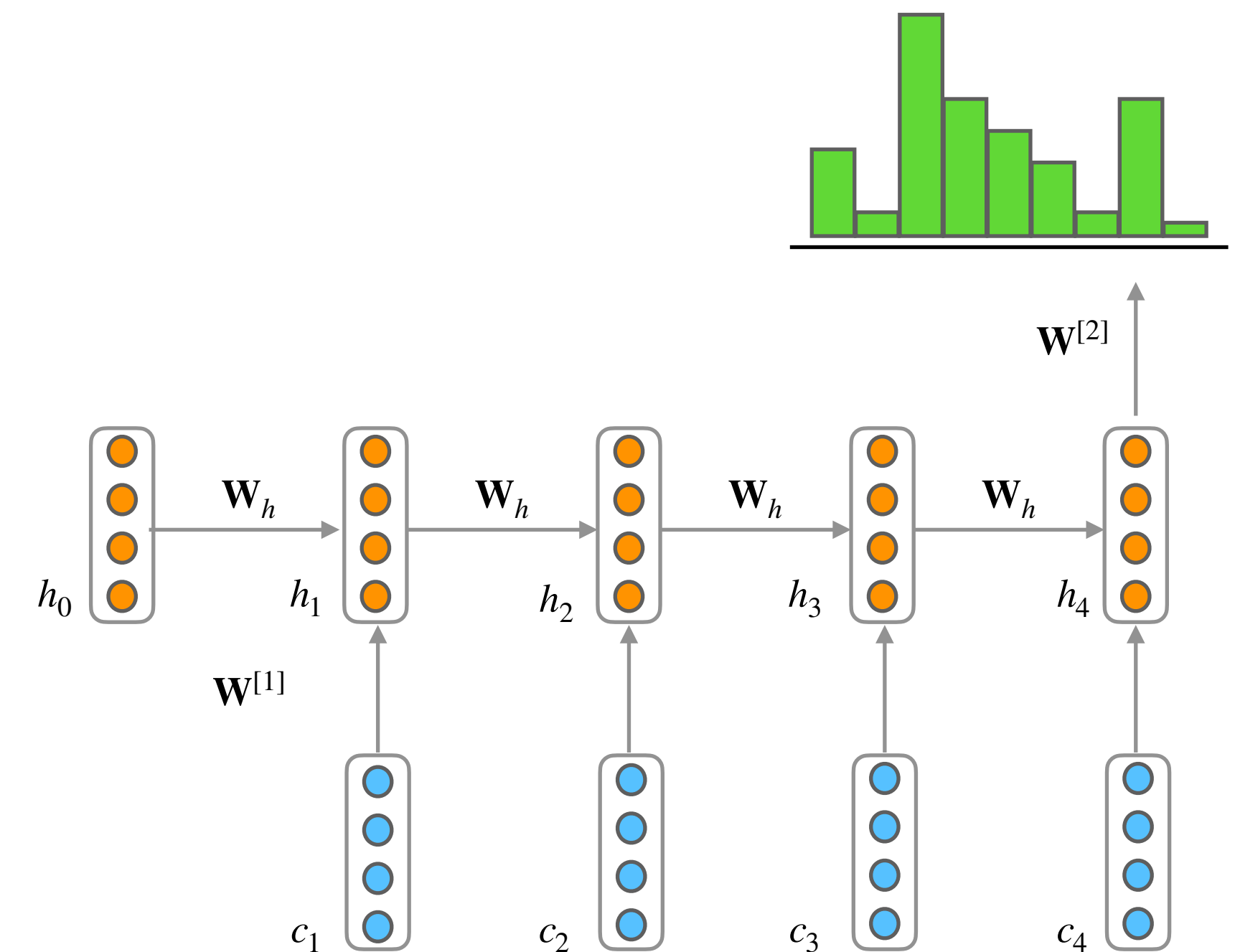$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

All these are vectors of same length $n$

Gates are applied using element-wise (or Hadamard) product: $\odot$

14

# Summarizing RNNs

- Recurrent Neural Networks processes sequences one element at a time
- RNNs do not have
  - the limited context problem of $n$-gram models
  - the fixed context limitation of feedforward LMs
  - since the hidden state can *in principle* represent information about all of the preceding words all the way back to the beginning of the sequence
- But training RNNs is hard
  - Vanishing gradient problem
  - LSTMs address it by incorporating a memory

# Lecture Outline

- Announcements
- Recap: Recurrent Neural Nets
- Applications of RNNs
- Seq2Seq Modeling with Encoder-Decoder Networks
- Attention Mechanism
- More on Attention
- Transformers: Self-Attention Networks
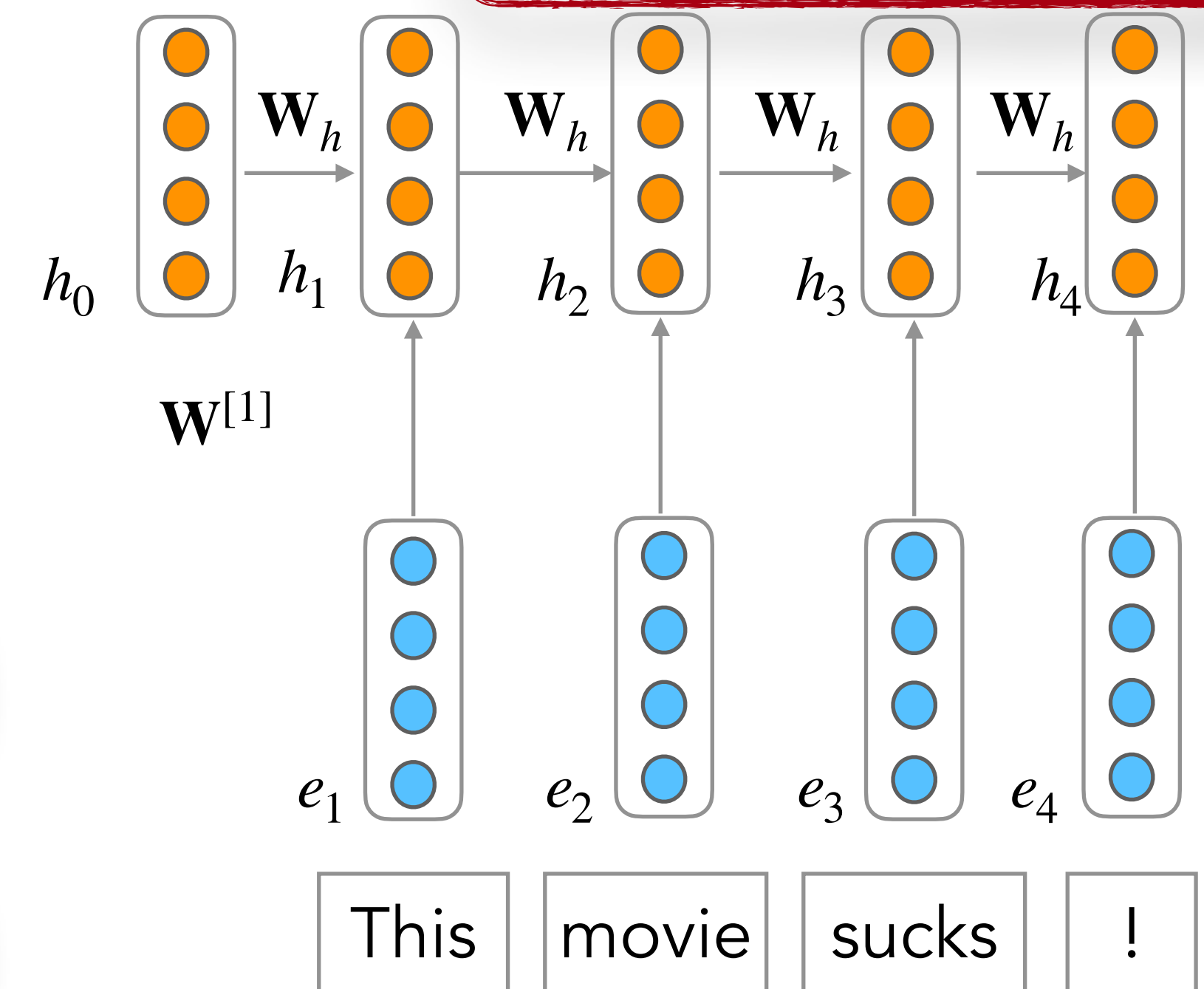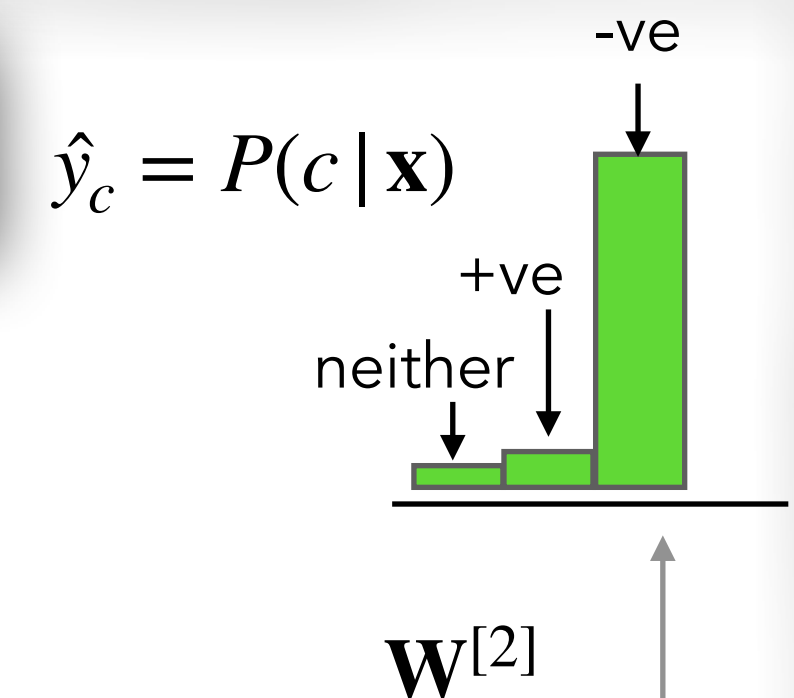
# Applications of RNNs

# RNNs for Sequence Classification

- $\mathbf{x}$ = Entire sequence / document of length $n$
- $y$ = (Multivariate) labels
- Pass $\mathbf{x}$ through the RNN one word at a time generating a new hidden state at each time step
- Hidden state for the last token of the text, $\mathbf{h}_n$ is a compressed representation of the entire sequence
- Pass $\mathbf{h}_n$ to a feedforward network (or multilayer perceptron) that chooses a class via a softmax over the possible classes
- Better sequence representations?
  - could also average all $\mathbf{h}_i$'s or
  - consider the maximum element along each dimension
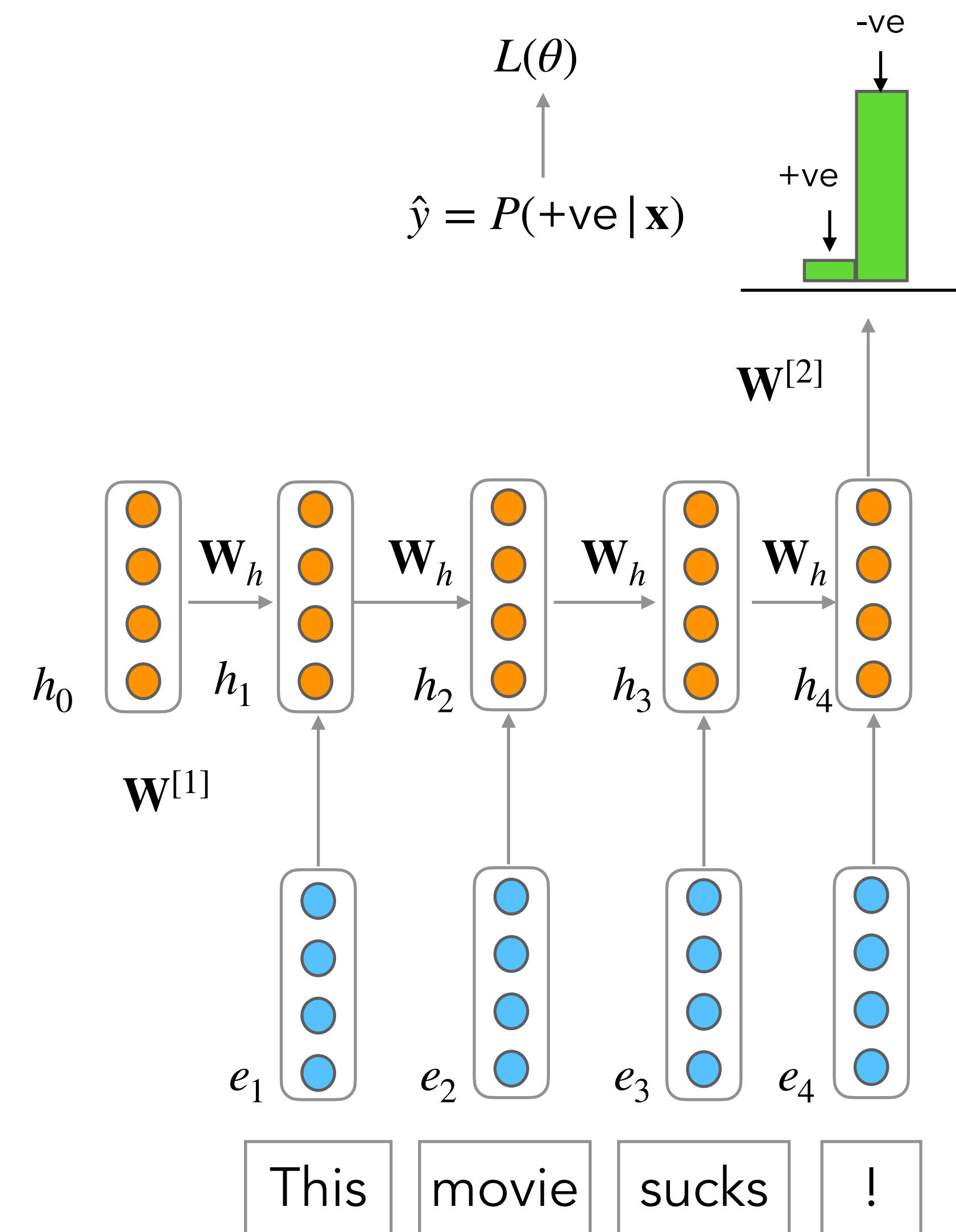
**Mean pooling**

**Max pooling**

**Multilayer Perceptron**

$$\hat{y}_c = P(c \mid \mathbf{x})$$

-ve

+ve

neither

$\mathbf{W}^{[2]}$

$h_0$ $\mathbf{W}_h$ $h_1$ $\mathbf{W}_h$ $h_2$ $\mathbf{W}_h$ $h_3$ $\mathbf{W}_h$ $h_4$

$\mathbf{W}^{[1]}$

$e_1$ $e_2$ $e_3$ $e_4$

This | movie | sucks | !

18

# Training RNNs for Sequence Classification

- Don't need intermediate outputs for the words in the sequence preceding the last element
- Loss function used to train the weights in the network is based entirely on the final text classification task
  - Cross-entropy loss
- Backprop: error signal from the classification is backpropagated all the way through the weights in the feedforward classifier through, to its input, and then through to the three sets of weights in the RNN



19

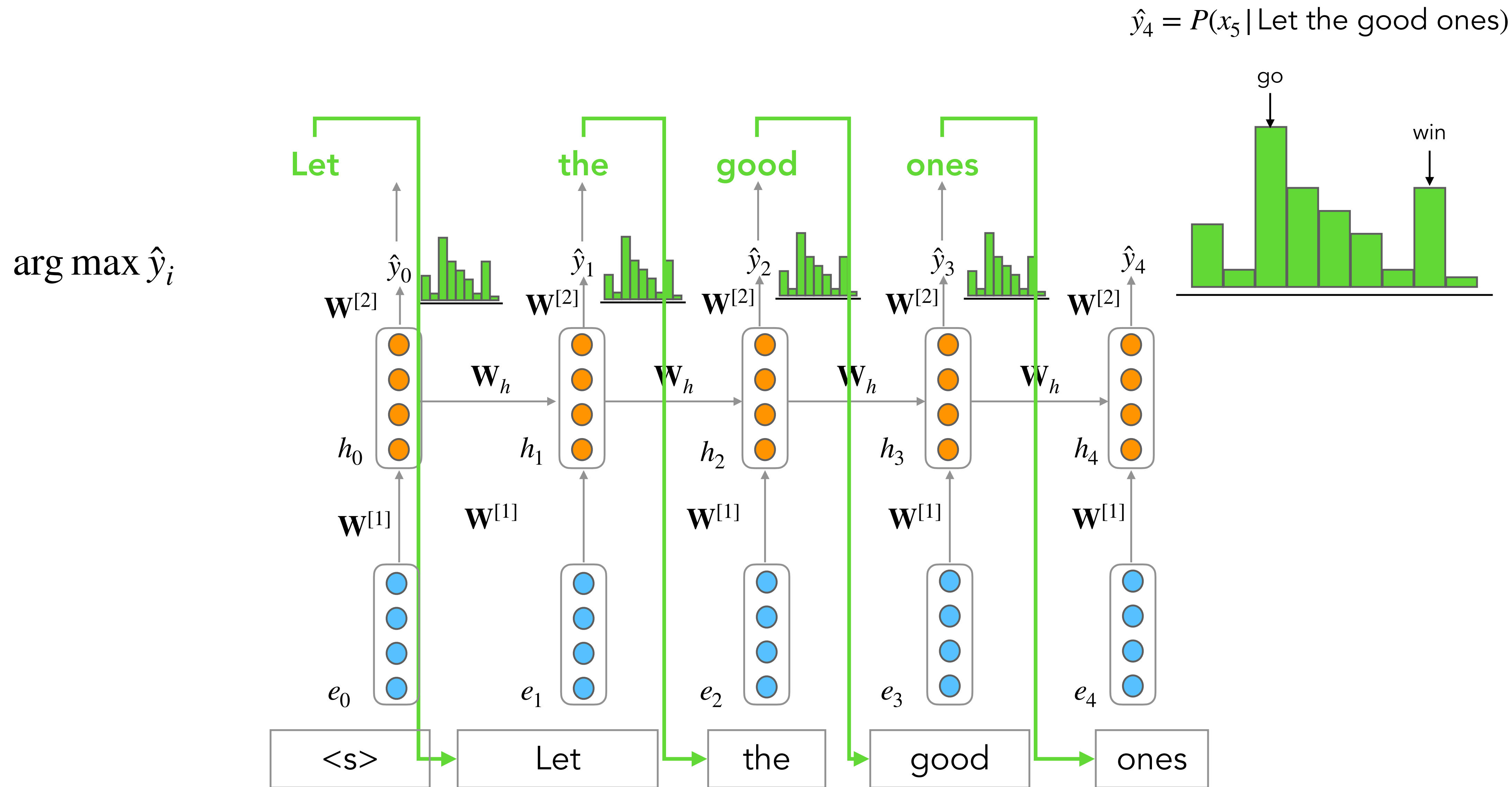# Generation with RNNLMs

Remember sampling from $n$-gram LMs?

- Similar to sampling from $n$-gram LMs
- First randomly sample a word to begin a sequence based on its suitability as the start of a sequence
- Then continue to sample words conditioned on our previous choices until
  - we reach a pre-determined length,
  - or an end of sequence token is generated

1. Choose a random bigram $(<s>, w)$ according to its probability
2. Now choose a random bigram $(w, x)$ according to its probability…and so on until we choose $</s>$

```
<s> I
    I want
      want to
          to eat
              eat Chinese
                  Chinese food
                        food   </s>
I want to eat Chinese food
```

$\hat{y}_4 = P(x_5 \mid \text{Let the good ones})$

$\arg\max \hat{y}_i$

# Generation with RNNLMs

1. Sample a word in the output from the softmax distribution that results from using the beginning of sentence marker, <s>, as the first input.
2. Use the word embedding for that first word as the input to the network at the next time step, and then sample the next word in the same fashion.
3. Continue generating until the end of sentence marker, </s>, is sampled or a fixed length limit is reached.

Repeated sampling of the next word conditioned on previous choices

Autoregressive Generation

# RNNLMs are Autoregressive Models

- Model that predicts a value at time $t$ based on a function of the previous values at times $t - 1, t - 2$, and so on
- Word generated at each time step is conditioned on the word selected by the network from the previous step
- State-of-the-art generation approaches are all autoregressive!
  - Machine translation, question answering, summarization
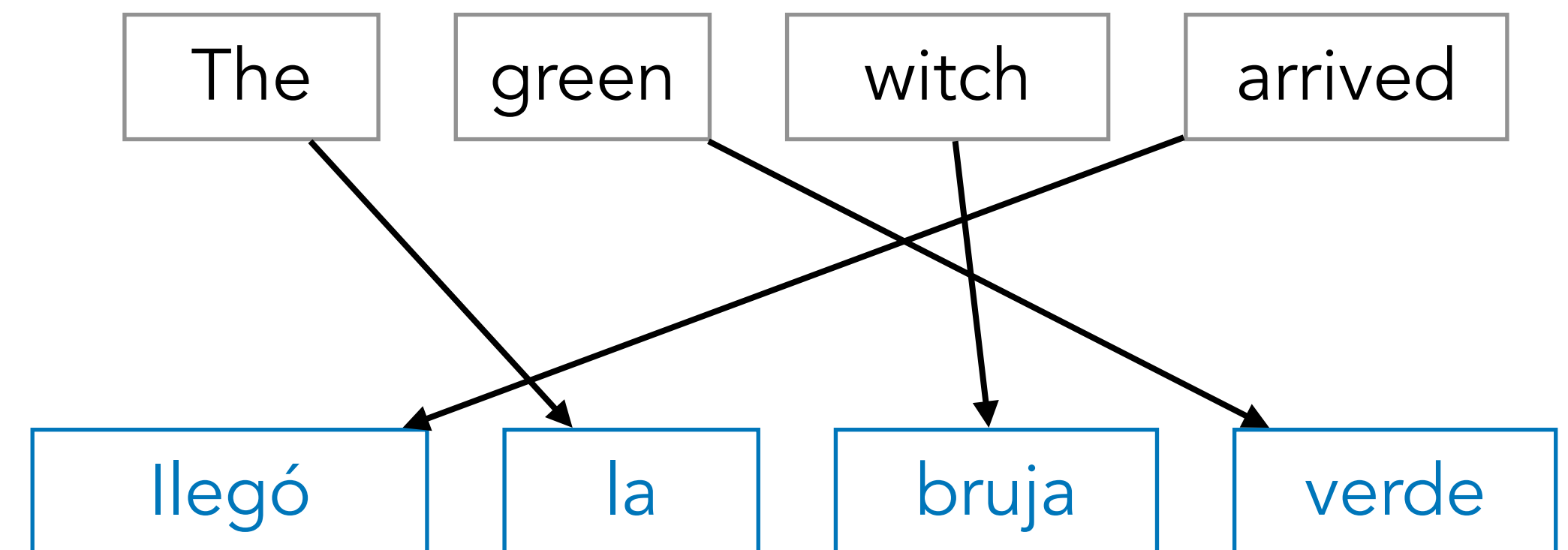- Key technique: prime the generation with the most suitable **context**

Can do better than <s>!

Provide rich task-appropriate context!

23

# (Neural) Machine Translation

Provide rich task-appropriate context!

- Sequence Generation Problem (as opposed to sequence classification)
  - $\mathbf{x}$ = Source sequence of length $n$
  - $\mathbf{y}$ = Target sequence of length $m$
- Different from regular generation from an LM
  - Since we expect the target sequence to serve a specific utility (translate the source)

Sequence-to-Sequence (Seq2seq)

| The | green | witch | arrived |
| --- | --- | --- | --- |

| llegó | la | bruja | verde |
| --- | --- | --- | --- |

# Sequence-to-Sequence Generation

| The | green | witch | arrived |
|-----|-------|-------|---------|

| llegó | la | bruja | verde |
|-------|-----|-------|--------|

- Mapping between a token in the input and a token in the output can be very indirect
  - in some languages the verb appears at the beginning of the sentence; e.g. Arabic, Hawaiian
  - in other languages at the end; e.g. Hindi
  - in other languages between the subject and the object; e.g. English
- Does not necessarily align in a word-word way!

**Need a special architecture to summarize the entire context!**

# Sequence-to-Sequence Models

- Models capable of generating contextually appropriate, arbitrary length, output sequences given an input sequence.
- The key idea underlying these networks is the use of an **encoder network** that takes an input sequence and creates a contextualized representation of it, often called the context.
- This representation is then passed to a **decoder network** which generates a task- specific output sequence.

Encoder-Decoder Networks

# Lecture Outline

- Announcements
- Recap: Recurrent Neural Nets
- Applications of RNNs
- Seq2Seq Modeling with Encoder-Decoder Networks
- Attention Mechanism
- More on Attention
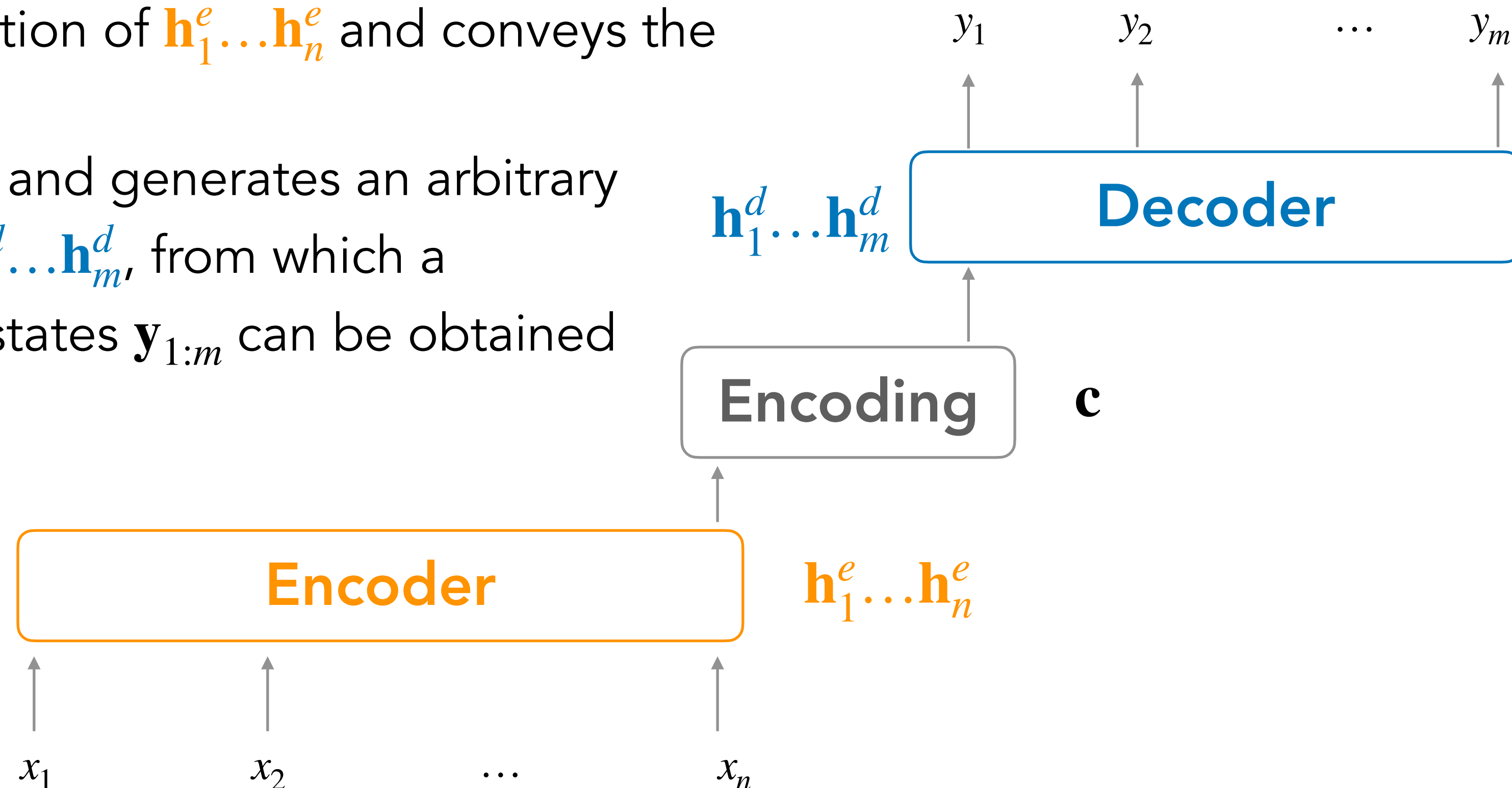- Transformers: Self-Attention Networks

# Sequence-to-Sequence Modeling with Encoder-Decoder Networks

# Encoder-Decoder Networks

Encoder-decoder networks consist of three components:

1. An **encoder** that accepts an input sequence, $\mathbf{x}_{1:n}$ and generates a corresponding sequence of contextualized representations, $\mathbf{h}_1^e \ldots \mathbf{h}_n^e$

2. A **encoding** vector, $\mathbf{c}$ which is a function of $\mathbf{h}_1^e \ldots \mathbf{h}_n^e$ and conveys the essence of the input to the decoder

3. A **decoder** which accepts $\mathbf{c}$ as input and generates an arbitrary length sequence of hidden states $\mathbf{h}_1^d \ldots \mathbf{h}_m^d$, from which a corresponding sequence of output states $\mathbf{y}_{1:m}$ can be obtained

Encoders and decoders can be made of FFNNs, RNNs, or Transformers
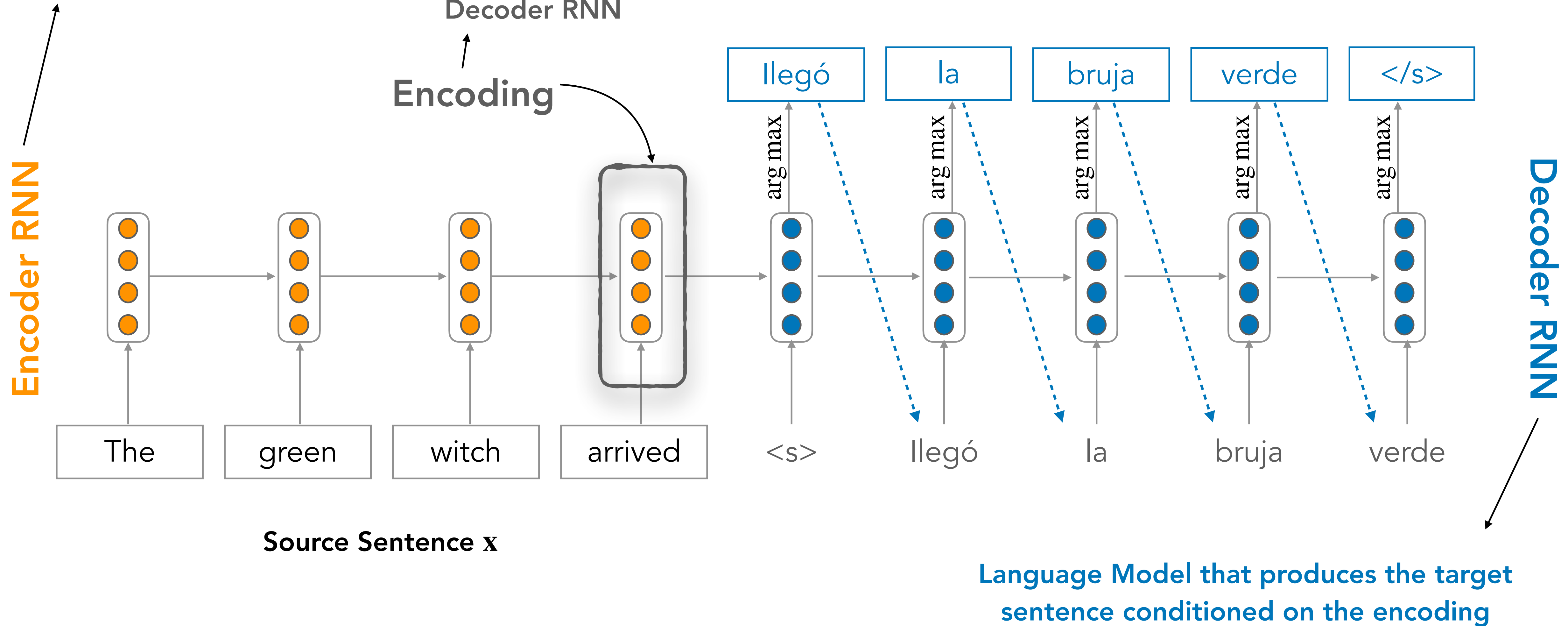
$$y_1 \qquad y_2 \qquad \cdots \qquad y_m$$

$\mathbf{h}_1^d \ldots \mathbf{h}_m^d$ **Decoder**

**Encoding** $\mathbf{c}$

**Encoder** $\mathbf{h}_1^e \ldots \mathbf{h}_n^e$

$$x_1 \qquad x_2 \qquad \cdots \qquad x_n$$

USC Viterbi

**negative log prob. of "llegó"**

**negative log prob. of "</s>"**

**Loss**

$$L(\theta) = \frac{1}{T} \sum_{t=1}^{T} L_t(\theta) = \quad L_0(\theta) \quad + \quad L_1(\theta) \quad + \quad L_2(\theta) \quad + \quad L_3(\theta) \quad + \quad L_4(\theta)$$

$\hat{y}_0 \qquad \hat{y}_1 \qquad \hat{y}_2 \qquad \hat{y}_3 \qquad \hat{y}_4$

**Encoder RNN**

**Decoder RNN**

| The | green | witch | arrived | <s> | llegó | la | bruja | verde |

**Source Sentence x**

**Target Sentence y**

**This needs to capture all information about the source sentence. Information bottleneck!**

**Target Sentence y**

Encoding

**Encoder RNN**

**Decoder RNN**

| llegó | la | bruja | verde | </s> |

arg max — arg max — arg max — arg max — arg max

The green witch arrived

<s> llegó la bruja verde

**Source Sentence x**

"you can't cram the meaning of a whole %&@#&ing sentence into a single $*(&@ing vector!"

*– Ray Mooney, Professor of Computer Science, UT Austin*

# Information Bottleneck: One Solution



**Encoder RNN**

Encoding

The   green   witch   arrived

The   green   witch   arrived

What if we had access to all hidden states?

How to create this?

34

# Lecture Outline

- Announcements
- Recap: Recurrent Neural Nets
- Applications of RNNs
- Seq2Seq Modeling with Encoder-Decoder Networks
- Attention Mechanism
- More on Attention
- Transformers: Self-Attention Networks

# Attention Mechanism

# Attention Mechanism

- Attention mechanisms allow the decoder to focus on a particular part of the source sequence at each time step
- Fixed-length vector $\mathbf{c}_t^{att}$ (attention context vector)
  - Take a weighted sum of all the encoder hidden states
  - One vector per time step *of the decoder*!
  - Weights *attend to* part of the source text relevant for the token the decoder is producing at step $t$
- In general, we have a single **query** vector and multiple **key** vectors.
  - We want to score each query-key pair

**query**

$$\text{Encoding } \mathbf{c}_t^{att} = f(\mathbf{h}_1^e, \ldots, \mathbf{h}_N^e, \mathbf{h}_t^d)$$

**keys**

**Encoder RNN**

| The | green | witch | arrived |

**Source Sentence x**

**Note: Notation different from J&M**

Bahdanau et al., 2015

# Seq2Seq with Attention

**Attention Scores / Attention Logits**

$$\textbf{score}(\mathbf{h}_t^d, \mathbf{h}_j^e) = \mathbf{h}_t^d \cdot \mathbf{h}_j^e$$

Dot product with keys (encoder hidden states) to encode similarity with what is decoded so far…

Query 1: Decoder, first time step

**Encoder RNN**

$\mathbf{h}_1^e$    $\mathbf{h}_2^e$    $\mathbf{h}_3^e$    $\mathbf{h}_4^e$    $\mathbf{h}_0^d$

| The | green | witch | arrived | <s> |

**Source Sentence x**

Dot product attention

38

**Note: Notation different from J&M**

**On this decoder tilmestep we are mostly focusing on the source token "arrived"**

$$\alpha_t = \mathbf{softmax}(\mathbf{score}(\mathbf{h}_t^d, \mathbf{h}_j^e)) \in \Delta^N$$

$$\alpha_{t,j} = \frac{\exp \mathbf{h}_t^d \cdot \mathbf{h}_j^e}{\sum_{n=1}^{N} \exp \mathbf{h}_t^d \cdot \mathbf{h}_n^e}$$

$$\mathbf{score}(\mathbf{h}_t^d, \mathbf{h}_j^e) = \mathbf{h}_t^d \cdot \mathbf{h}_j^e$$

Take softmax to turn into probability distribution

Attention Distribution

Attention Scores

Encoder RNN

$\mathbf{h}_1^e$    $\mathbf{h}_2^e$    $\mathbf{h}_3^e$    $\mathbf{h}_4^e$    $\mathbf{h}_0^d$

The    green    witch    arrived    &lt;s&gt;

**Source Sentence x**

**Note: Notation different from J&M**

**Attention Distribution**

**Attention Scores**

**Encoder RNN**

$$\mathbf{c}_t^{att} = \sum_{i=1}^{N} \alpha_{t,i} \mathbf{h}_i^e \in \mathbb{R}^d$$

Use the attention distribution to take a weighted sum of the encoder hidden states.

$$\alpha_t = \mathbf{softmax}(\mathbf{score}(\mathbf{h}_t^d, \mathbf{h}_j^e)) \in \Delta^N$$

$$\mathbf{score}(\mathbf{h}_t^d, \mathbf{h}_j^e) = \mathbf{h}_t^d \cdot \mathbf{h}_j^e$$

$\mathbf{h}_1^e$ $\mathbf{h}_2^e$ $\mathbf{h}_3^e$ $\mathbf{h}_4^e$ $\mathbf{h}_0^d$

The | green | witch | arrived | <s>

The attention output mostly contains information the hidden states that received high attention.

**Source Sentence x**

**Note: Notation different from J&M**

$$c_t^{att} = \sum_{i=1}^{N} \alpha_{t,i} \mathbf{h}_i^e \in \mathbb{R}^d$$

$$\alpha_t = \text{softmax}(\text{score}(\mathbf{h}_t^d, \mathbf{h}_j^e)) \in \mathbb{R}^N$$

Ilegó

arg max

$\hat{y}_0$

$\mathbf{W}^{[2]}$

$$\text{score}(\mathbf{h}_t^d, \mathbf{h}_j^e) = \mathbf{h}_t^d \cdot \mathbf{h}_j^e$$

$c_0^{att}$

$$\hat{y}_t = \text{softmax}(\mathbf{W}^{[2]}[\mathbf{h}_t^d; c_t^{att}])$$

$$\mathbf{W}^{[2]} \in \mathbb{R}^{V \times 2d}$$

Concatenate attention output with decoder hidden state, then use to compute $\hat{y}_0$ as before

$\mathbf{h}_1^e$  $\mathbf{h}_2^e$  $\mathbf{h}_3^e$  $\mathbf{h}_4^e$  $\mathbf{h}_0^d$

The    green    witch    arrived    <s>

**Source Sentence x**

**Note: Notation different from J&M**

$$\mathbf{c}_t^{att} = \sum_{i=1}^{N} \alpha_{t,i} \mathbf{h}_i^e \in \mathbb{R}^d$$

$$\alpha_t = \text{softmax}(\text{score}(\mathbf{h}_t^d, \mathbf{h}_j^e)) \in \mathbb{R}^N$$

$$\text{score}(\mathbf{h}_t^d, \mathbf{h}_j^e) = \mathbf{h}_t^d \cdot \mathbf{h}_j^e$$
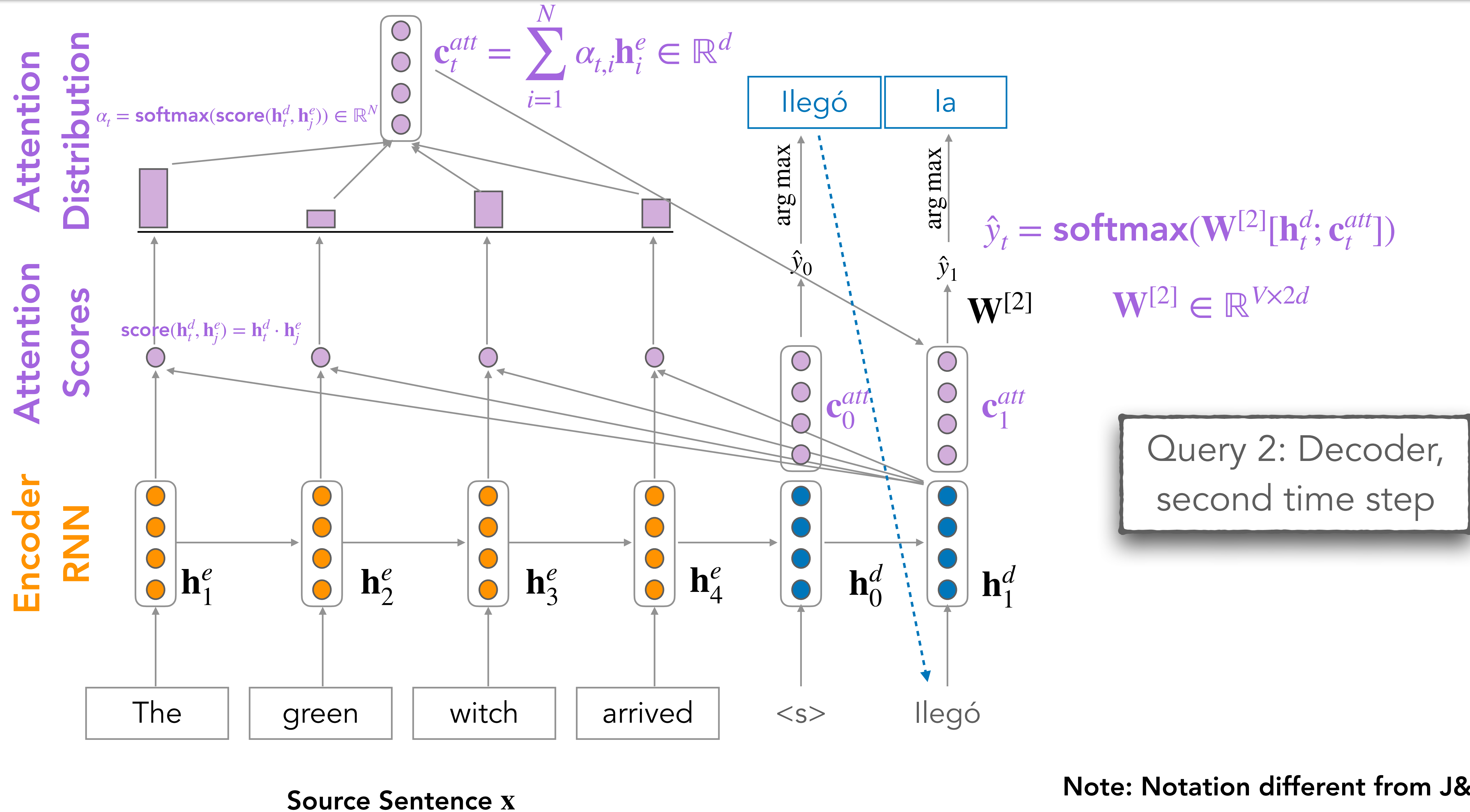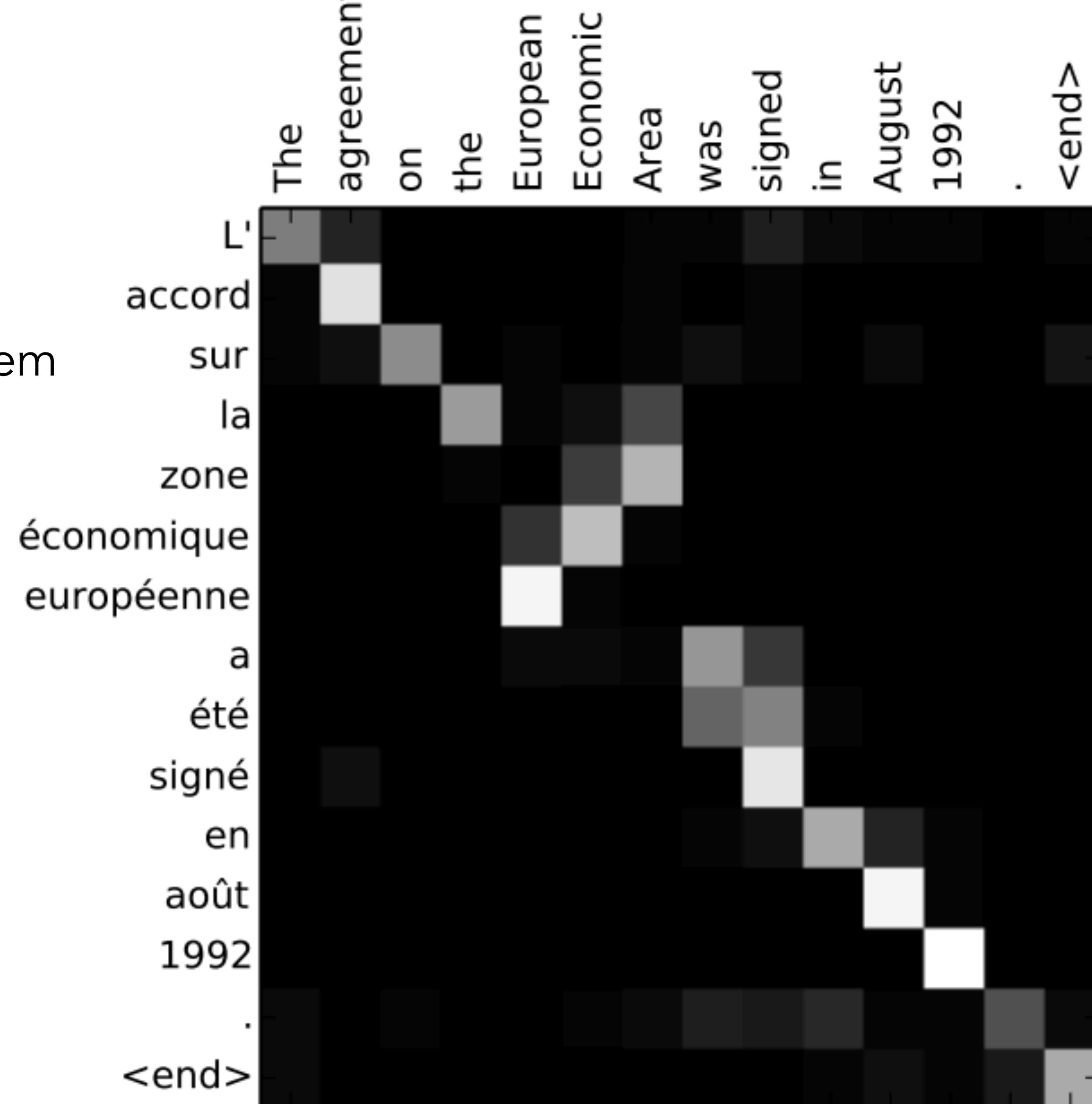
$$\hat{y}_t = \text{softmax}(\mathbf{W}^{[2]}[\mathbf{h}_t^d; \mathbf{c}_t^{att}])$$

$$\mathbf{W}^{[2]} \in \mathbb{R}^{V \times 2d}$$

Attention Distribution

Attention Scores

Encoder RNN

Ilegó    la

arg max    arg max

$\hat{y}_0$    $\hat{y}_1$

$\mathbf{W}^{[2]}$

$\mathbf{c}_0^{att}$    $\mathbf{c}_1^{att}$

$\mathbf{h}_1^e$    $\mathbf{h}_2^e$    $\mathbf{h}_3^e$    $\mathbf{h}_4^e$    $\mathbf{h}_0^d$    $\mathbf{h}_1^d$

The    green    witch    arrived    <s>    Ilegó

Query 2: Decoder, second time step

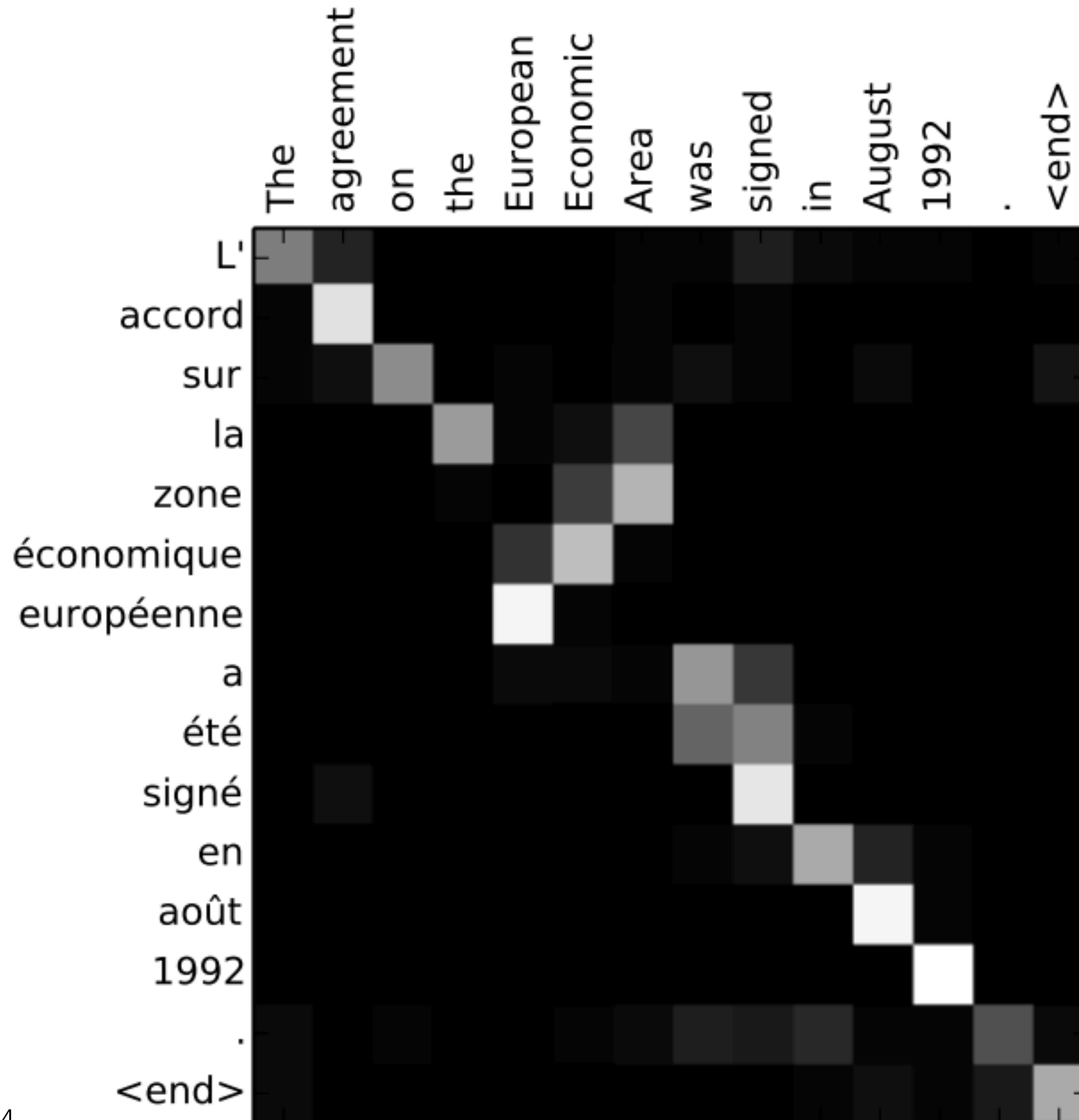**Source Sentence x**

**Note: Notation different from J&M**

# Why Attention?

- Attention significantly **improves** neural machine translation **performance**
  - Very useful to allow decoder to focus on certain parts of the source
- Attention **solves the information bottleneck** problem
  - Attention allows decoder to look directly at source; bypass bottleneck
- Attention **helps with vanishing gradient problem**
  - Provides shortcut to faraway states
- Attention provides some **interpretability**
  - By inspecting attention distribution, we can see what the decoder was focusing on →
  - We get alignment for free! We never explicitly trained an alignment system! The network just learned alignment by itself

43

# Seq2Seq Summary



- Seq2Seq modeling is popular for close-ended generation tasks
    - MT, Summarization, QA
        - Involves an encoder and a decoder
            - Can be any neural architecture!
- Popular Seq2Seq Models using Transformers: BART, T5
- Secret Sauce: Attention
- Next: Self-Attention and Transformers

44