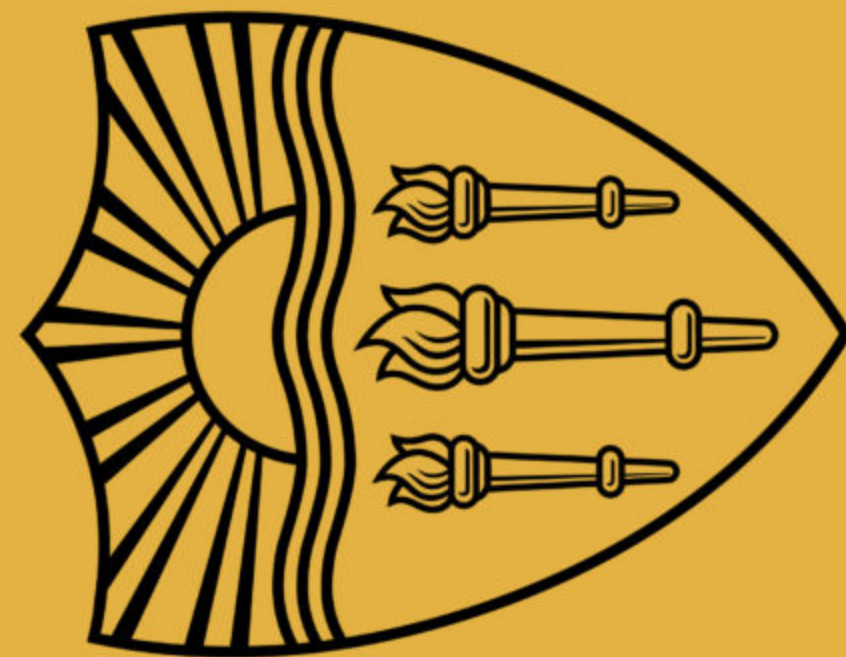


C
S
D



Lecture 10: Sequence-to-Sequence

Instructor: Swabha Swayamdipta
USC CSCI 499 LMs in NLP
Feb 21, 2024 Spring



Slides mostly adapted from Dan Jurafsky, Chris Manning, Mohit Iyer and Richard Socher

Logistics / Announcements

- Today: Graded Quiz 2
- Project Proposals Graded and Feedback Shared
- Next Week: Project Discussions - A flipped classroom format where I answer your questions on the project feedback.
- Next Class:
 - Quiz 3 on Mon 2/26
 - Also HW2 Due

Early Neural Language Models

| | | | |
|---------|--------------------------------------|--------|----------------------------------|
| Jan 29: | Logistic Regression | SLIDES | |
| Jan 31: | Logistic Regression II | SLIDES | HW1 Due |
| Feb 5: | Word Embeddings I | SLIDES | |
| Feb 7: | Word Embeddings II | SLIDES | HW2 Released PROPOSAL DUE |
| Feb 12: | Feedforward Neural Nets and Backprop | SLIDES | |
| Feb 14: | Recurrent Neural Network LMs | SLIDES | |

Modern Neural Language Models

| | | | |
|--------------------|--|----------------------------|--|
| Feb 19: | No Class | PRESIDENT'S DAY | |
| Feb 21: | Sequence-To-Sequence and Attention | | |
| Feb 26: | Transformers - Building Blocks I | | HW2 Due |
| Feb 28: | | PROJECT DISCUSSIONS | |
| Mar 4: | Transformers - Building Blocks II | | HW3 Released PROGRESS REPORT DUE |
| Mar 6: | PyTorch Demo for Transformers | | |
| Mar 11: | No Class | SPRING BREAK | |
| Mar 13: | No Class | SPRING BREAK | |

Lecture Outline

- Quiz 2 Answers
- Recap: Recurrent Neural Nets
- Applications of RNNs
- Sequence-to-Sequence Modeling with Encoder-Decoder Networks
- Attention Mechanism

Quiz 2: Answers - Discussed in class / Redacted

Recap: Recurrent Neural Nets

Feedforward LMs

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

Logit: Unnormalized Scores

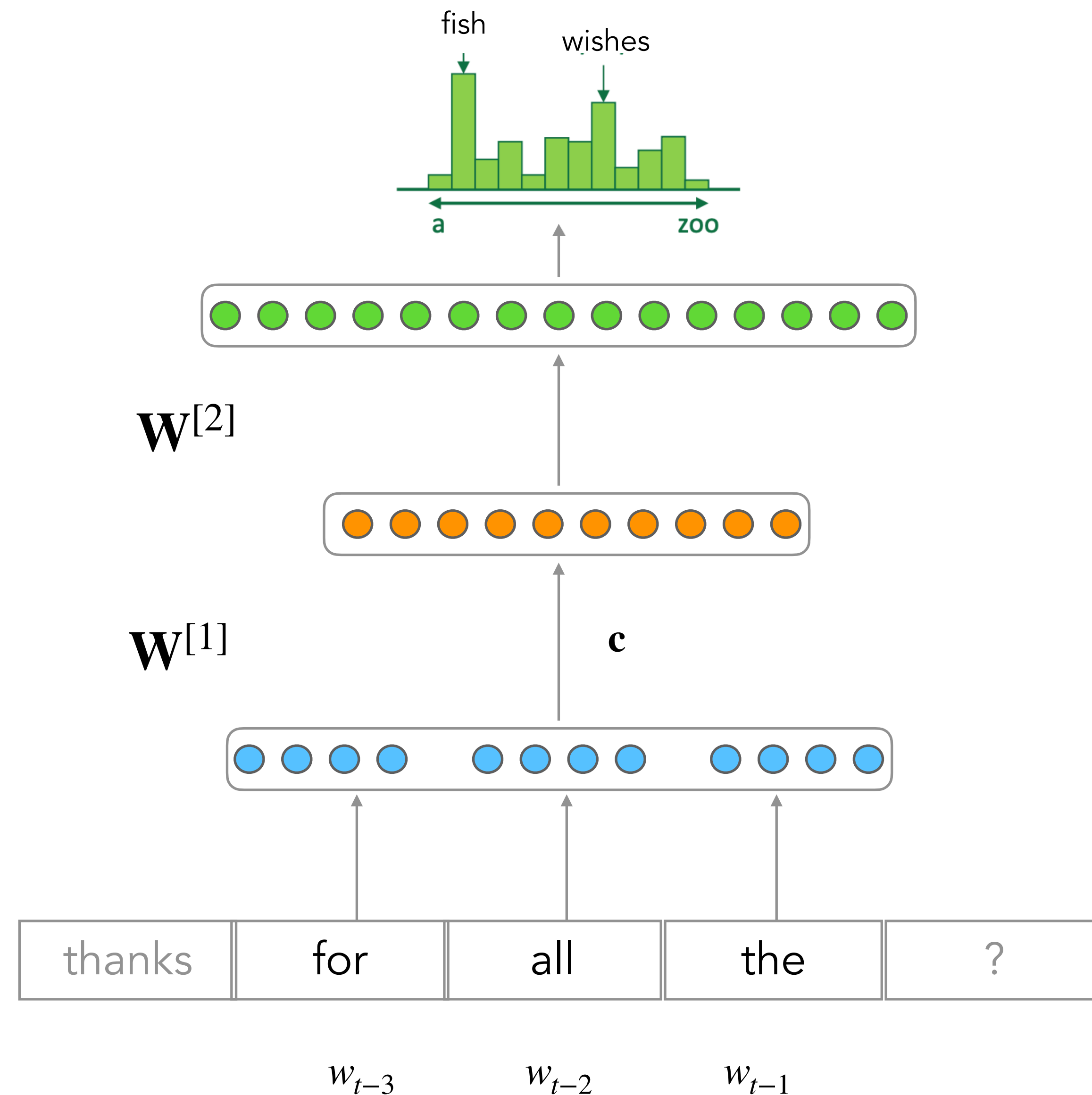
Hidden layer:

$\mathbf{h} = g(\mathbf{W}^{[1]}\mathbf{c})$

Word Embeddings

(concatenated in the window) \mathbf{c}_i

Feedforward networks are also called multilayer perceptrons



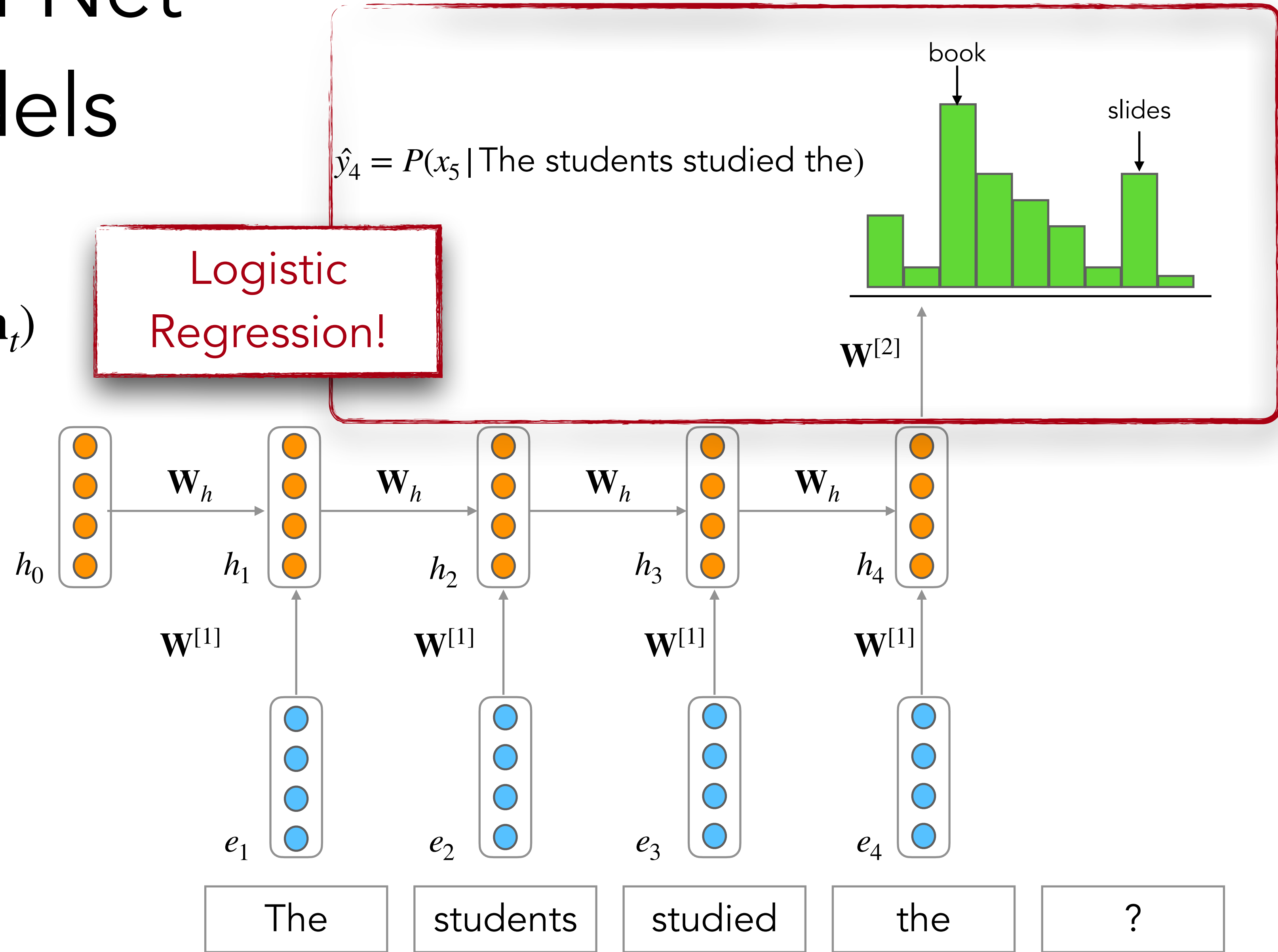
Recurrent Neural Net Language Models

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

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

Initial hidden state: \mathbf{h}_0

Word Embeddings, \mathbf{e}_i



Logistic Regression!

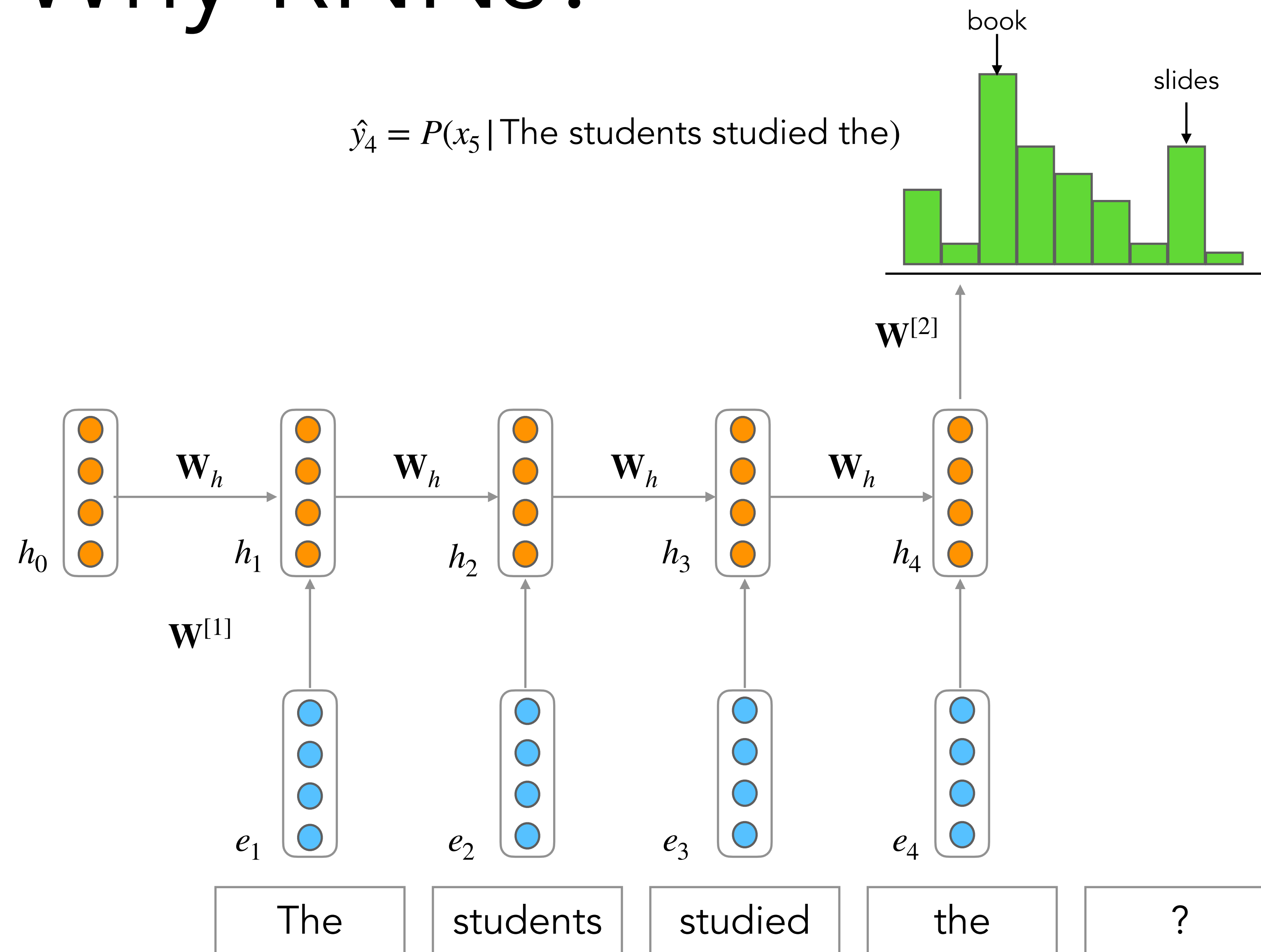
Recurrent Neural Networks

- Recurrent Neural Networks
 - Contain one hidden state \mathbf{h}_t per time step! Serves as a memory of the entire history...
 - Output of each neural unit at time t based both on
 - the current input at t and
 - the hidden state from time $t - 1$
- As the name implies, RNNs have a recursive formulation
 - dependent on its own earlier outputs as an input!
- RNNs thus don't have
 - the limited context problem that n-gram models have, or
 - the fixed context that feedforward language models have,
 - 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

Why RNNs?

RNN Advantages:

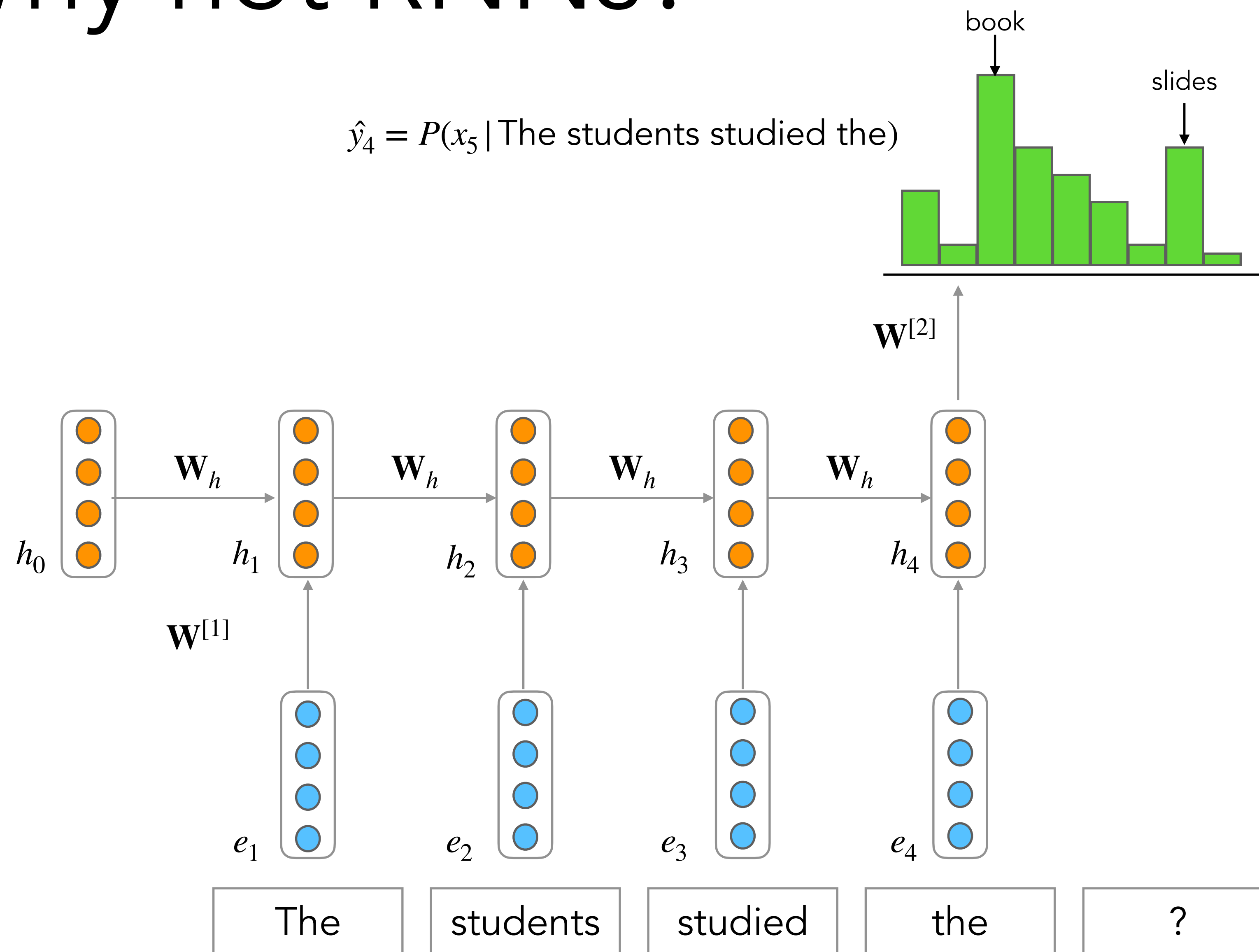
- Can process any length input
- Model size doesn't increase for longer input
- Computation for step t can (in theory) use information from many steps back
- Weights $\mathbf{W}^{[1]}$ are shared (tied) across timesteps \rightarrow Condition the neural network on all previous words



Why not RNNs?

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, \dots, 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 $y_t = x_{t+1}$:

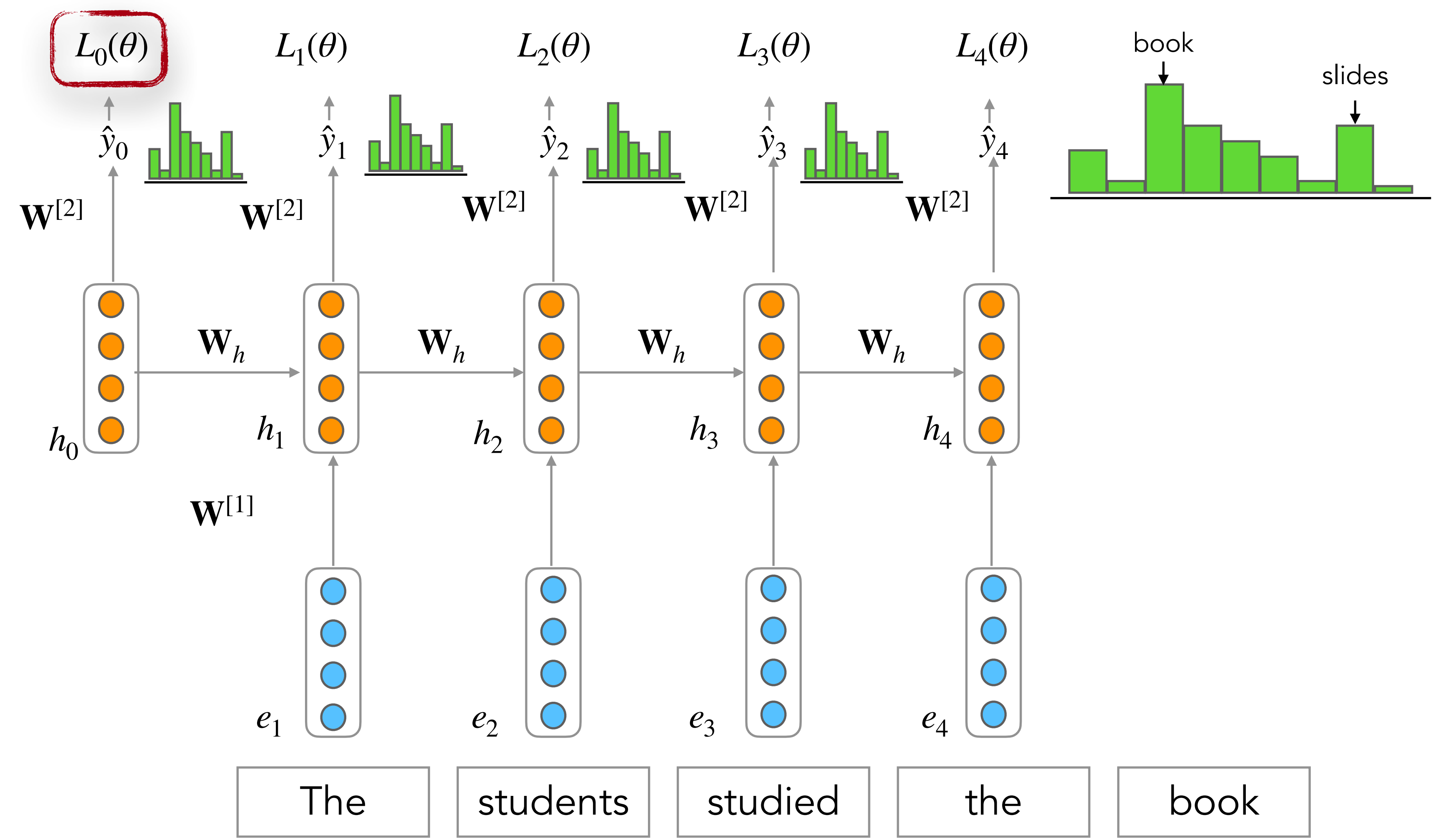
$$L_{CE}(\hat{y}_t, y_t; \theta) = - \sum_{v \in V} \mathbb{I}[y_t = v] \log \hat{y}_t = - \log p_{\theta}(x_{t+1} | x_{\leq t})$$

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

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

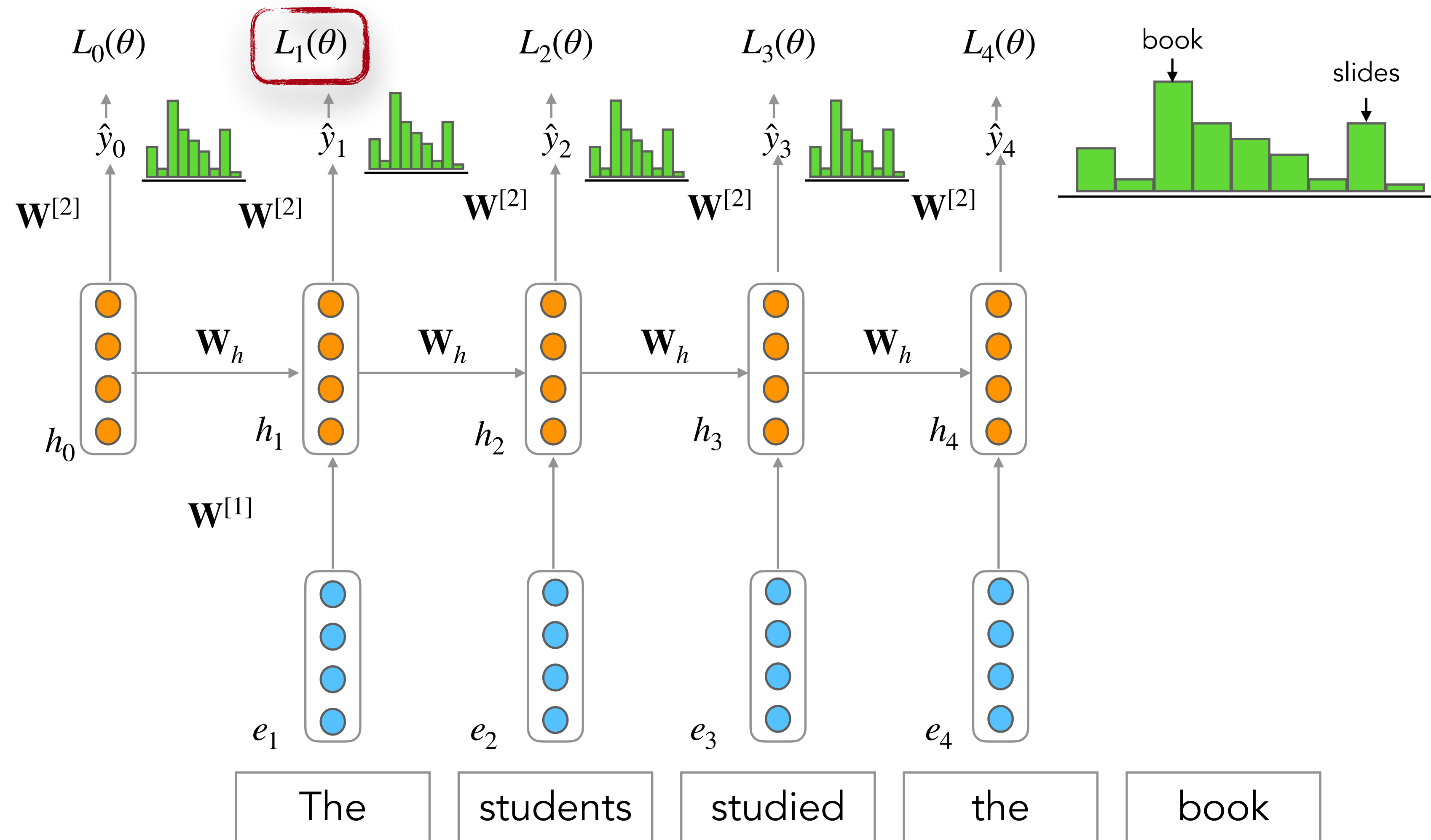
Loss

negative log prob. of "The"

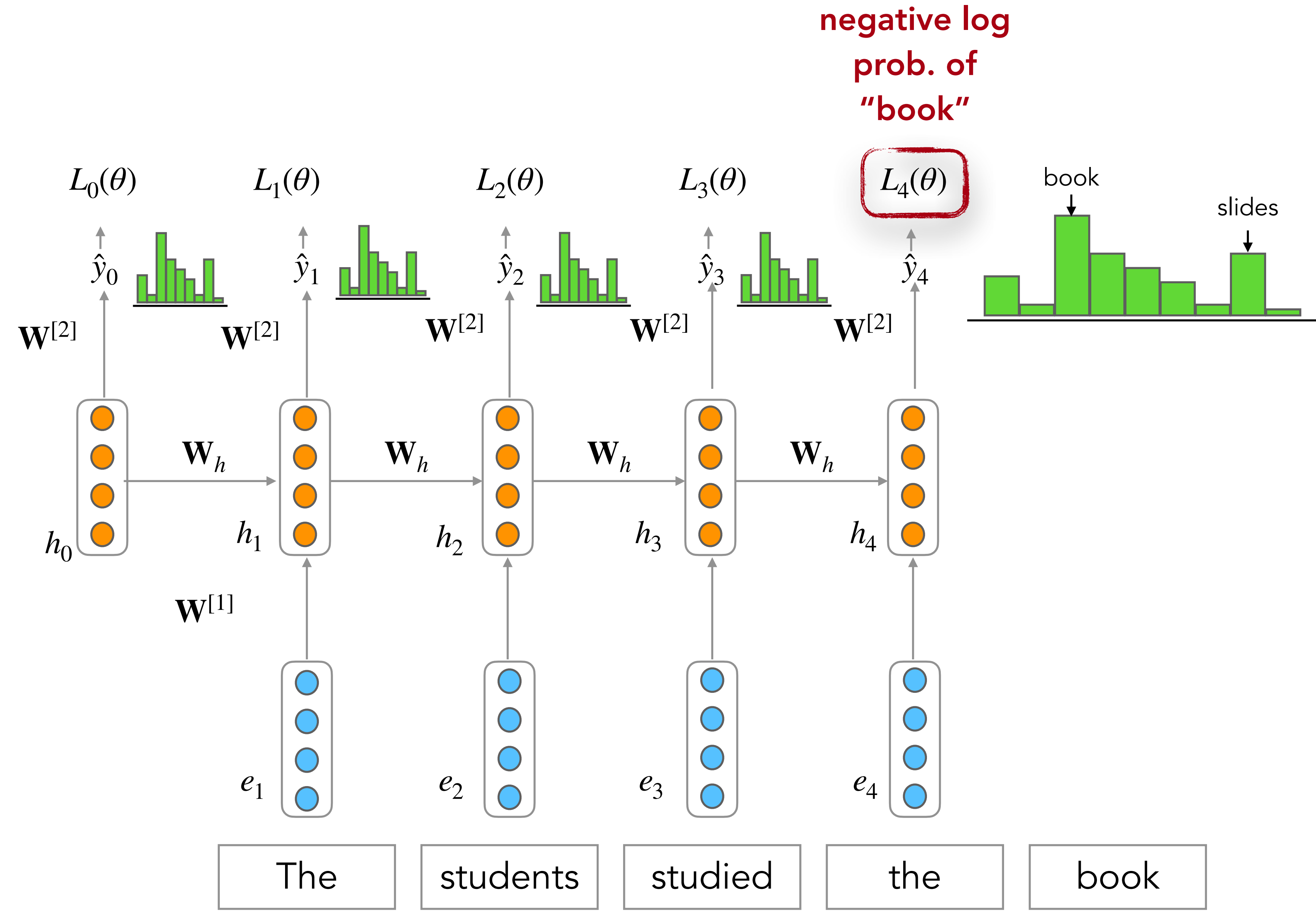


Loss

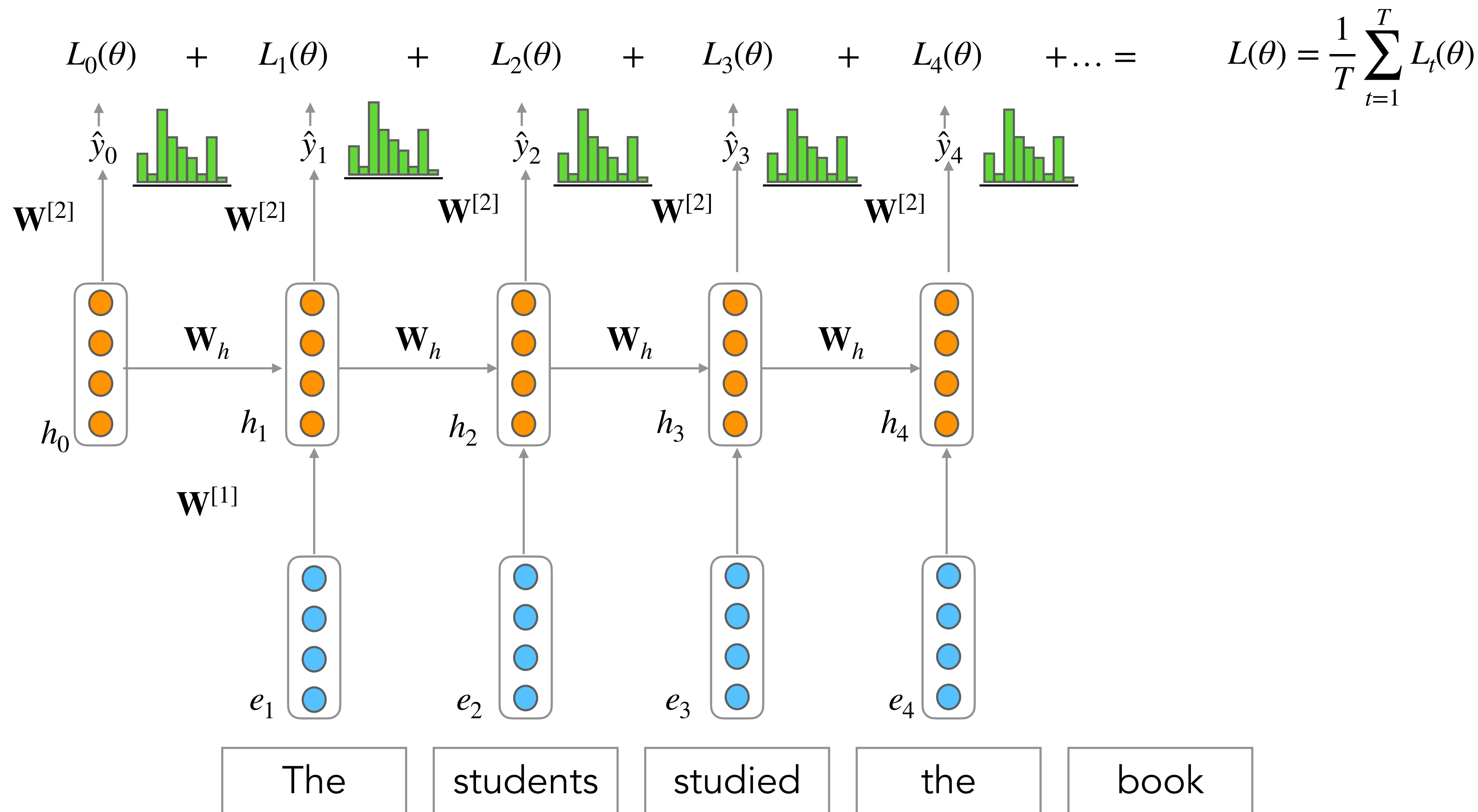
negative log prob. of "students"



Loss

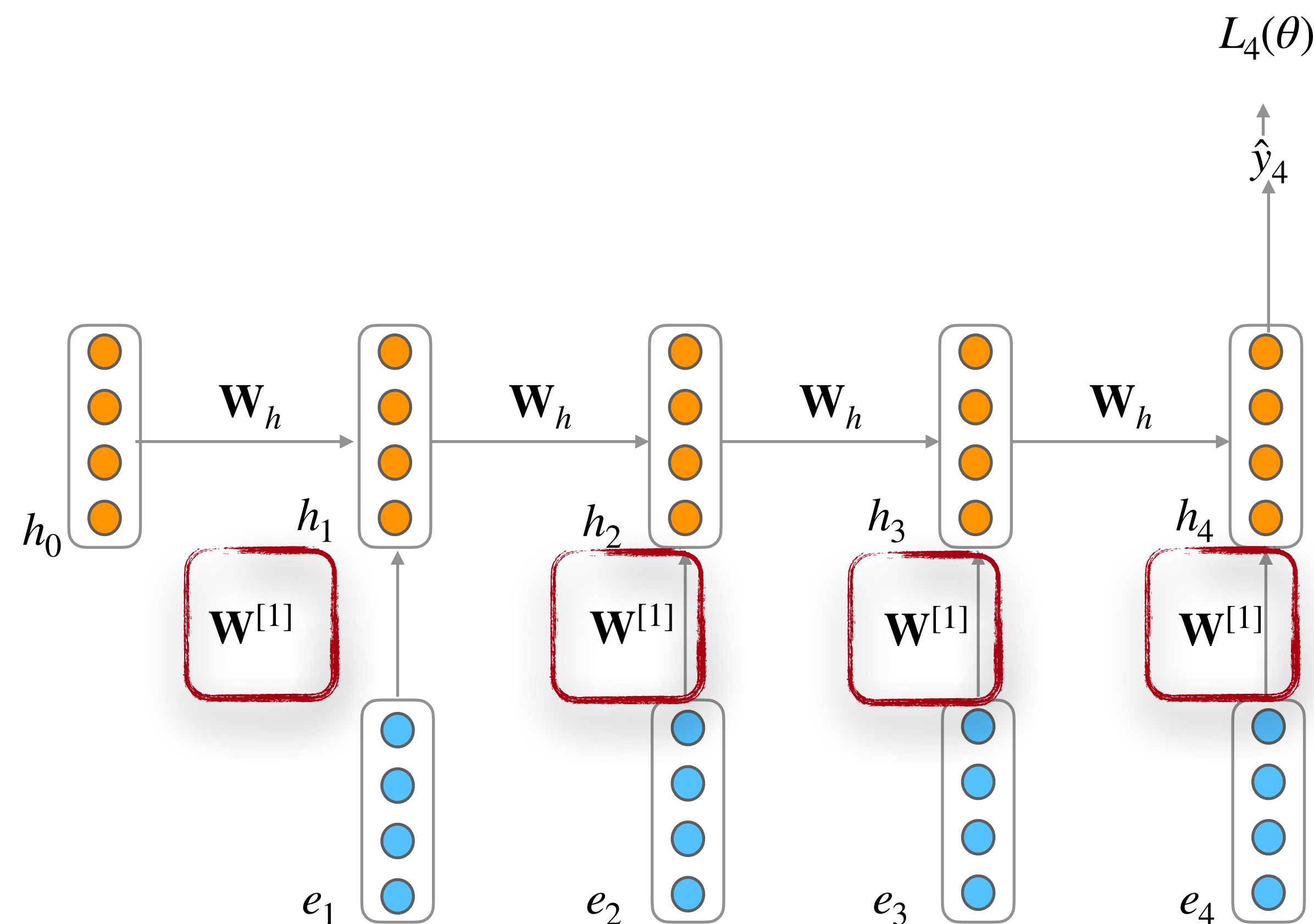


Loss

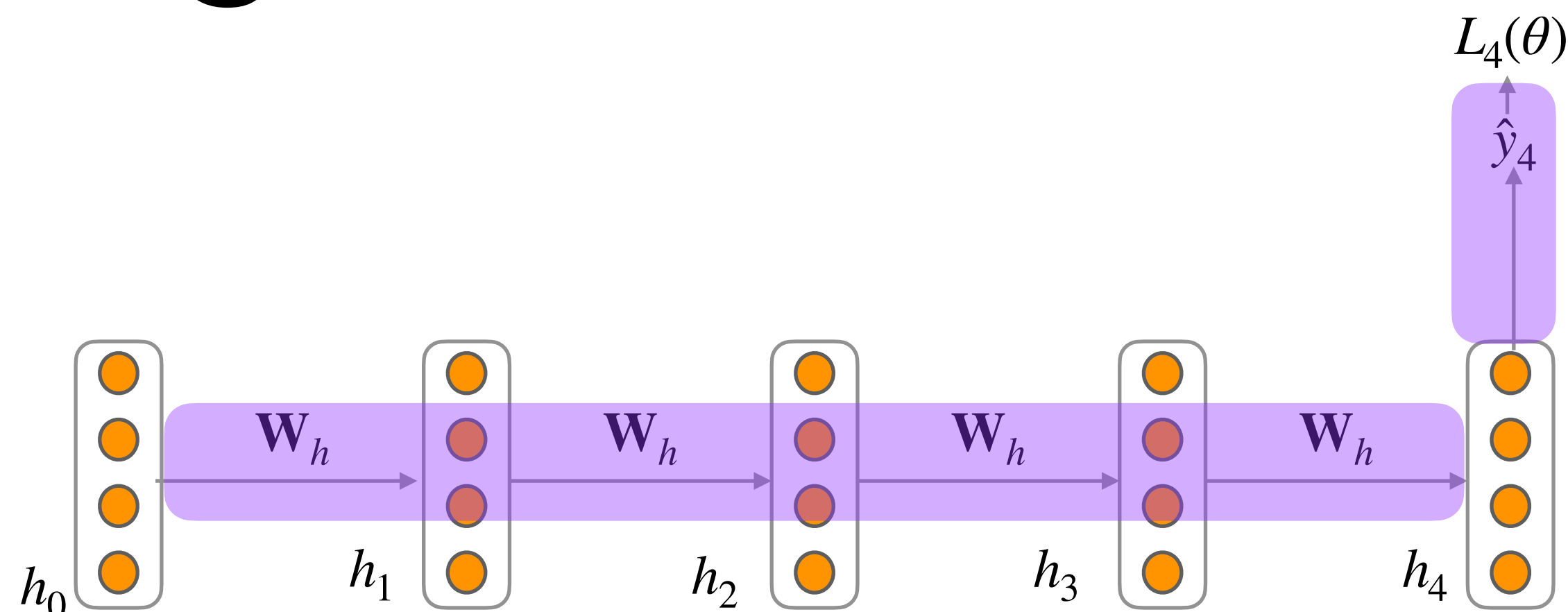


Training RNNs is hard

- Multiply the same matrix at each time step during forward propagation
- Ideally inputs from many time steps ago can modify output y
- This leads to something called the **vanishing gradient problem**



The Vanishing Gradient Problem: Intuition



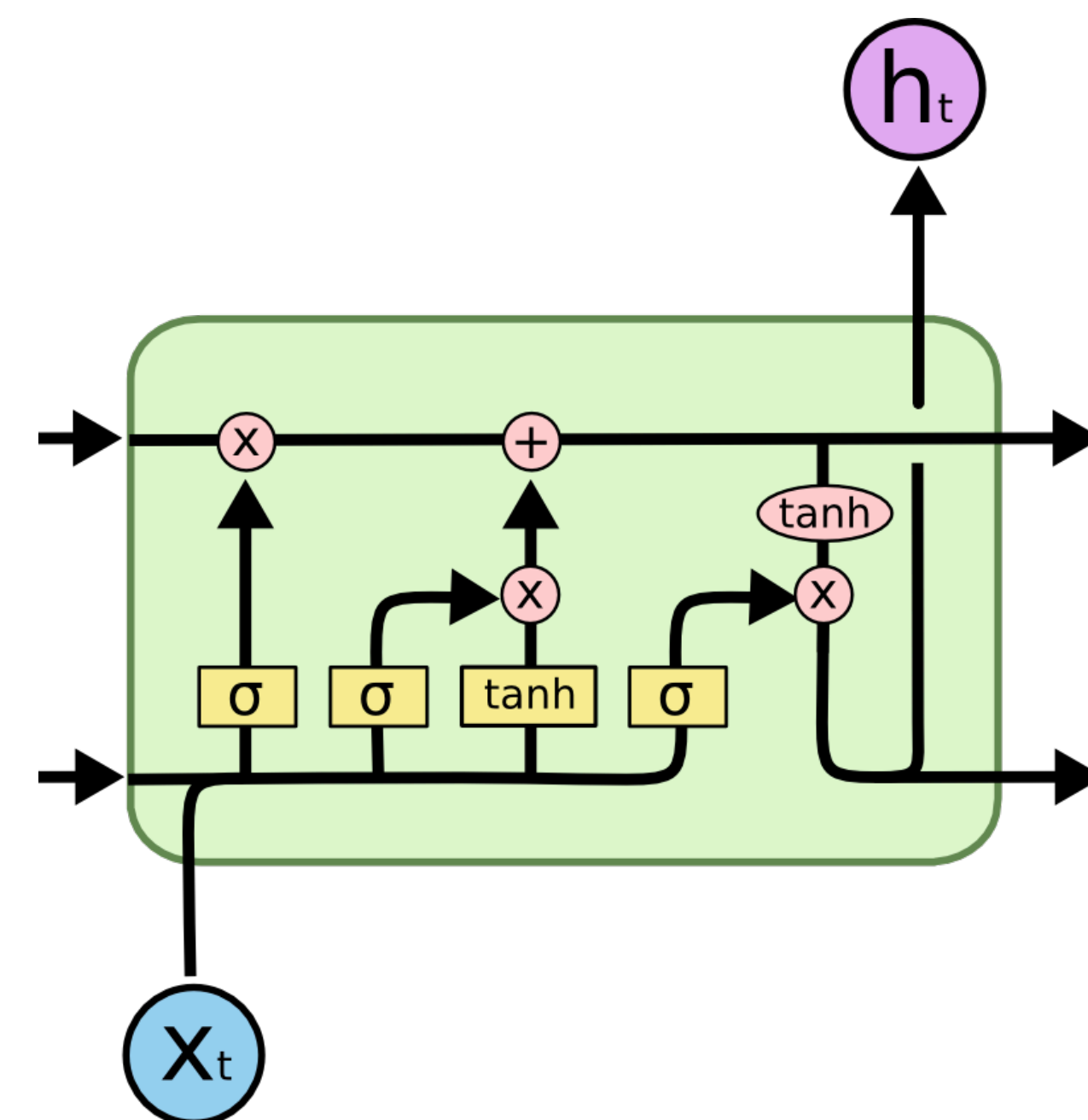
When these gradients are small, the gradient signal gets smaller and smaller as it backpropagates further...

$$\begin{aligned}
 \frac{\partial L_4}{\partial h_0} &= \frac{\partial h_1}{\partial h_0} \times \frac{\partial L_4}{\partial h_1} \\
 &= \frac{\partial h_1}{\partial h_0} \times \frac{\partial h_2}{\partial h_1} \times \frac{\partial L_4}{\partial h_2} \\
 &= \frac{\partial h_1}{\partial h_0} \times \frac{\partial h_2}{\partial h_1} \times \frac{\partial h_3}{\partial h_2} \times \frac{\partial L_4}{\partial h_3} \\
 &= \frac{\partial h_1}{\partial h_0} \times \frac{\partial h_2}{\partial h_1} \times \frac{\partial h_3}{\partial h_2} \times \frac{\partial h_4}{\partial h_3} \times \frac{\partial L_4}{\partial h_4}
 \end{aligned}$$

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

Long Short-Term Memory RNNs (LSTMs)

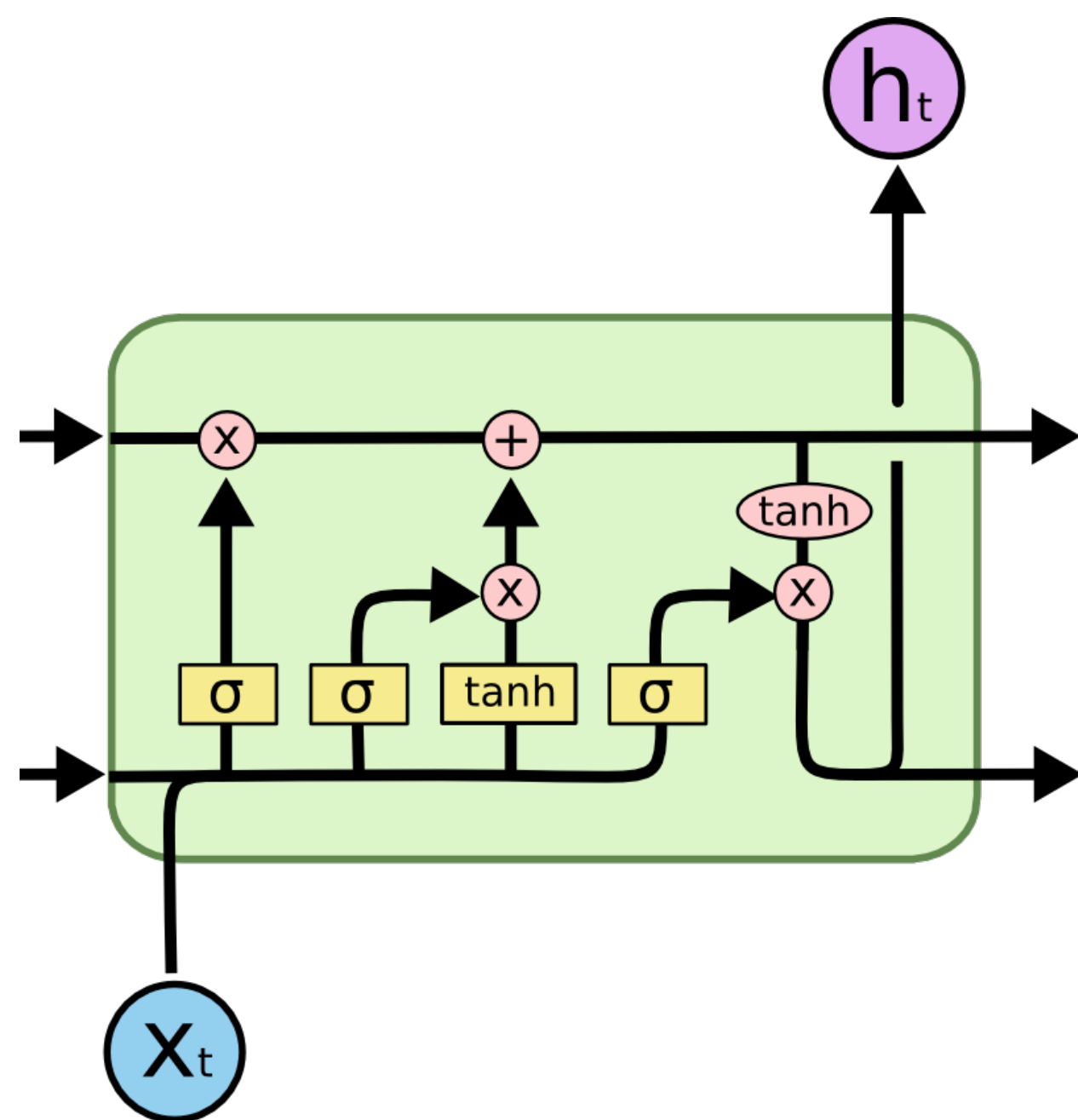
- At time step t , introduces a new cell state $\mathbf{c}_t \in \mathbb{R}^d$
 - In addition to a hidden state $\mathbf{h}_t \in \mathbb{R}^d$
 - The cell stores long-term information (memory)
 - The LSTM can read, erase, and write information from the cell!
 - The cell becomes conceptually rather like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

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

Sigmoid function: all gate values are between 0 and 1

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

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

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

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

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

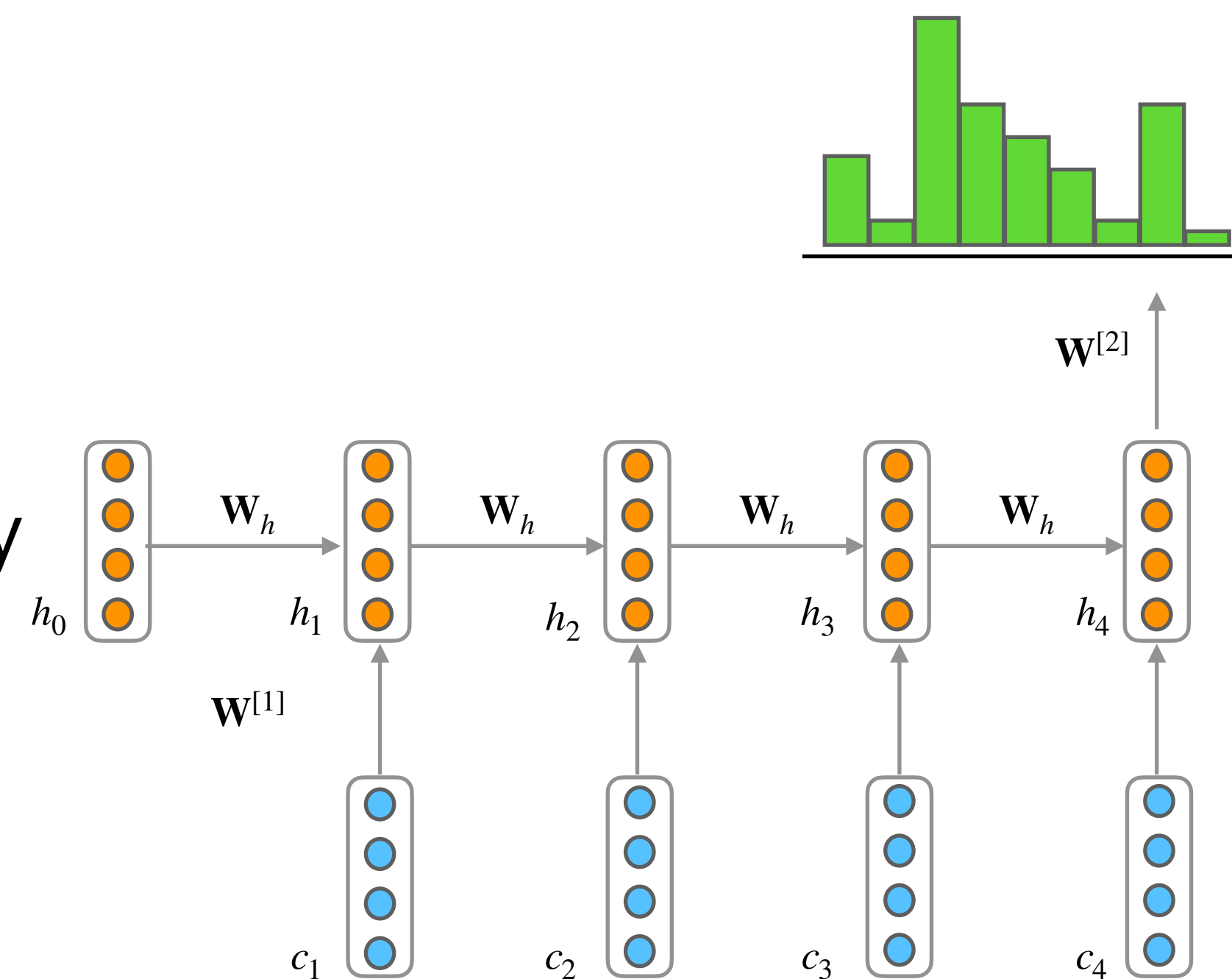
$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

All these are vectors of same length n

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

Summarizing RNNs

- 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
- Training can be expensive and might lead to vanishing gradients
- More advanced architectures: LSTMs (Long Short-Term Memories)



Can be applied to both classification and generation tasks

Applications of RNNs

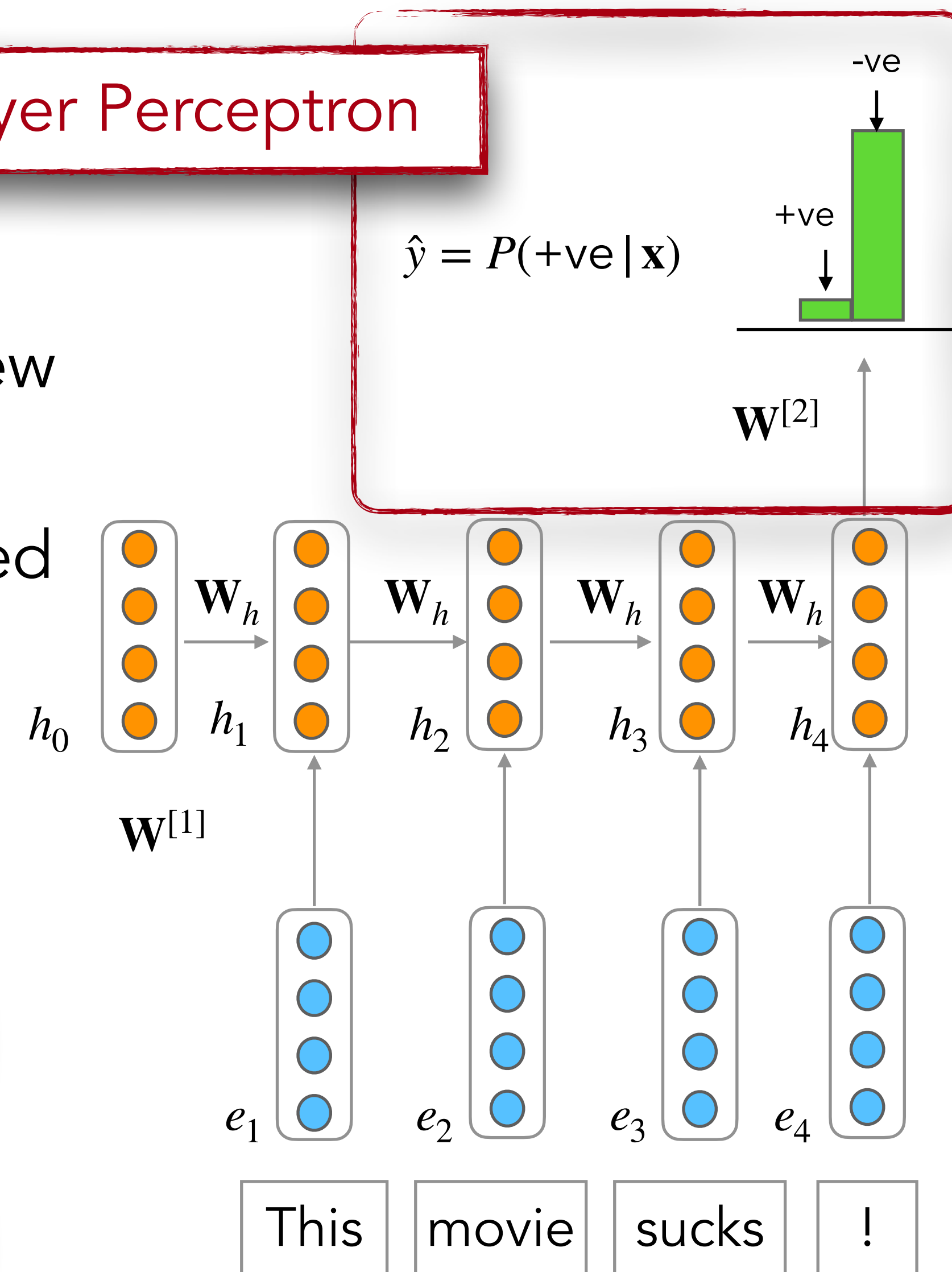
RNNs for Sequence Classification

- \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

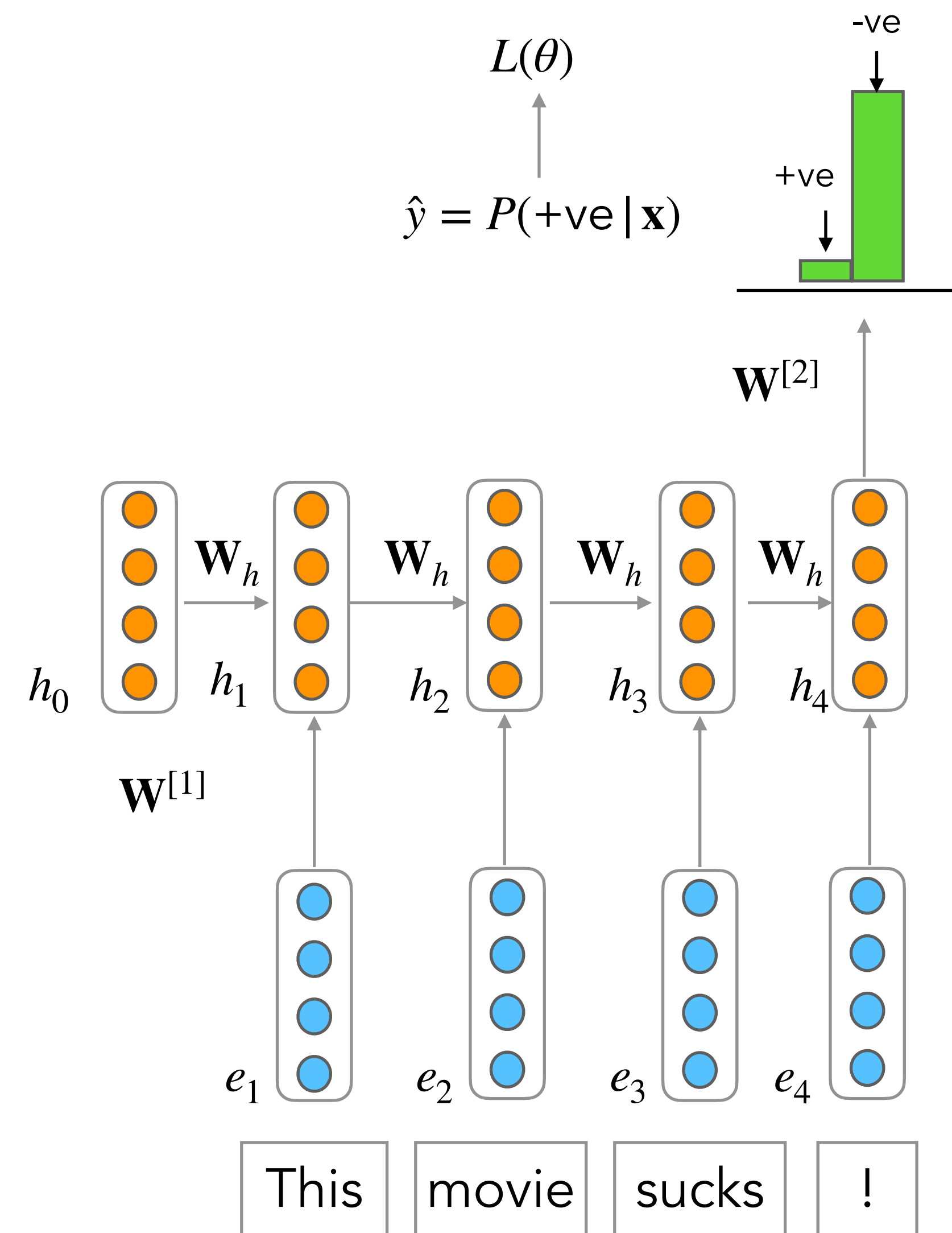
Max pooling

Multilayer Perceptron



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



Generation with RNNLMs

Remember sampling from n-gram LMs?

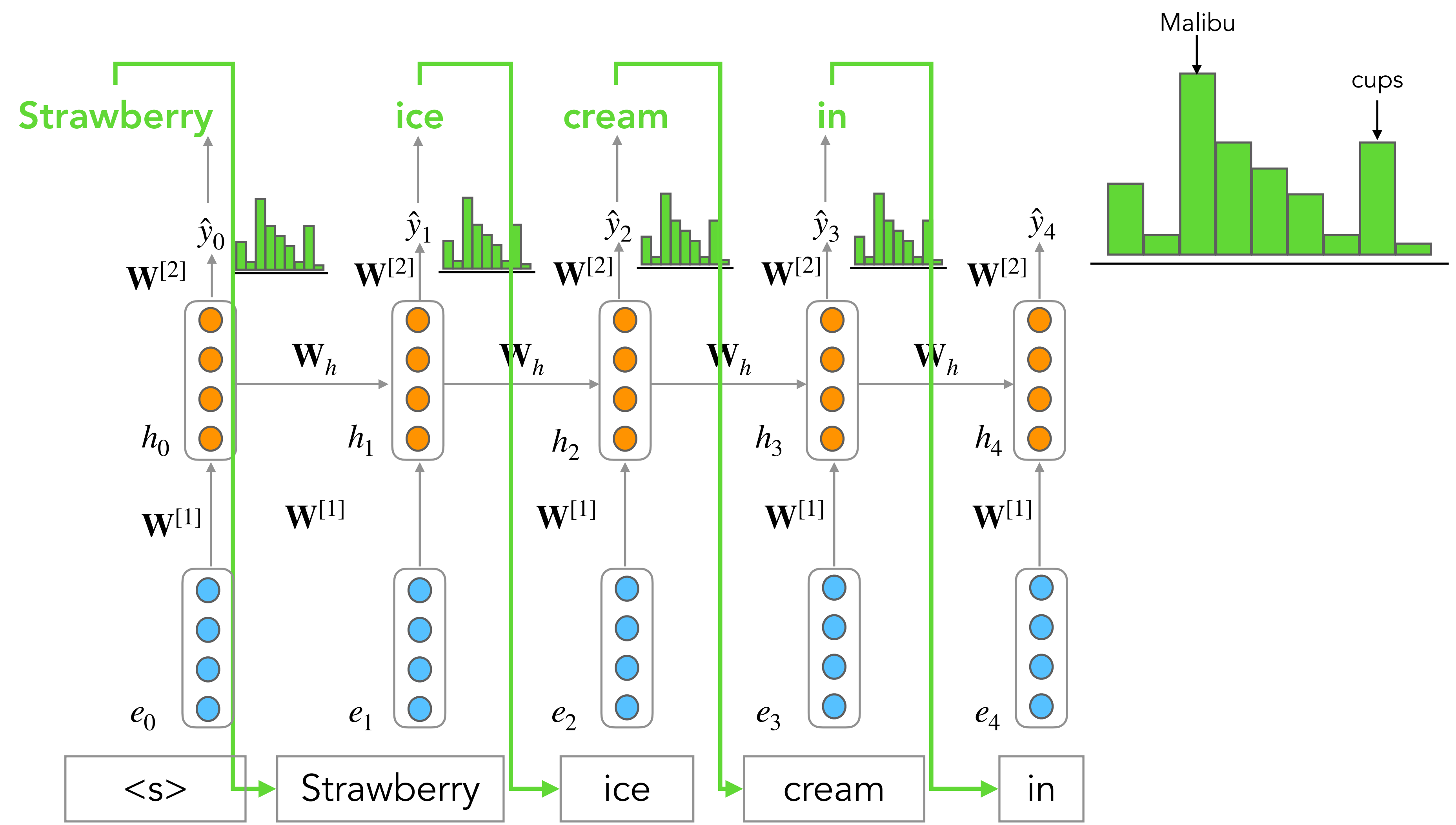
- Similar to sampling from n-gram LMs
- First randomly sample a word to begin a sequence based on its suitability as the start of a sequence
- Then continue to sample words conditioned on our previous choices until
 - we reach a pre-determined length,
 - or an end of sequence token is generated

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

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


$$\hat{y}_4 = P(x_5 | \text{Strawberry ice cream in})$$

arg max \hat{y}_i
 Initial hidden state: \mathbf{h}_0



Generation with RNNLMs

1. Sample a word in the output from the softmax distribution that results from using the beginning of sentence marker, $\langle s \rangle$, 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, $\langle /s \rangle$, 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**

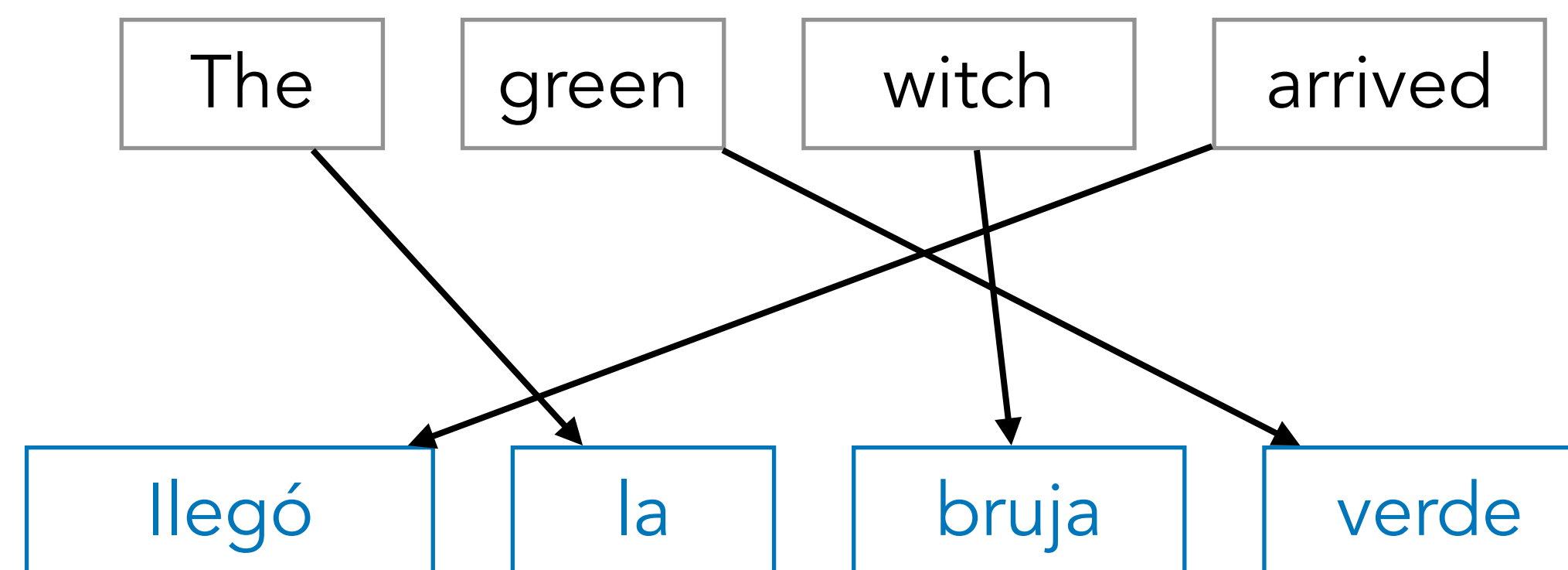
Can do better than $\langle s \rangle$!

Provide rich task-appropriate context!

(Neural) Machine Translation

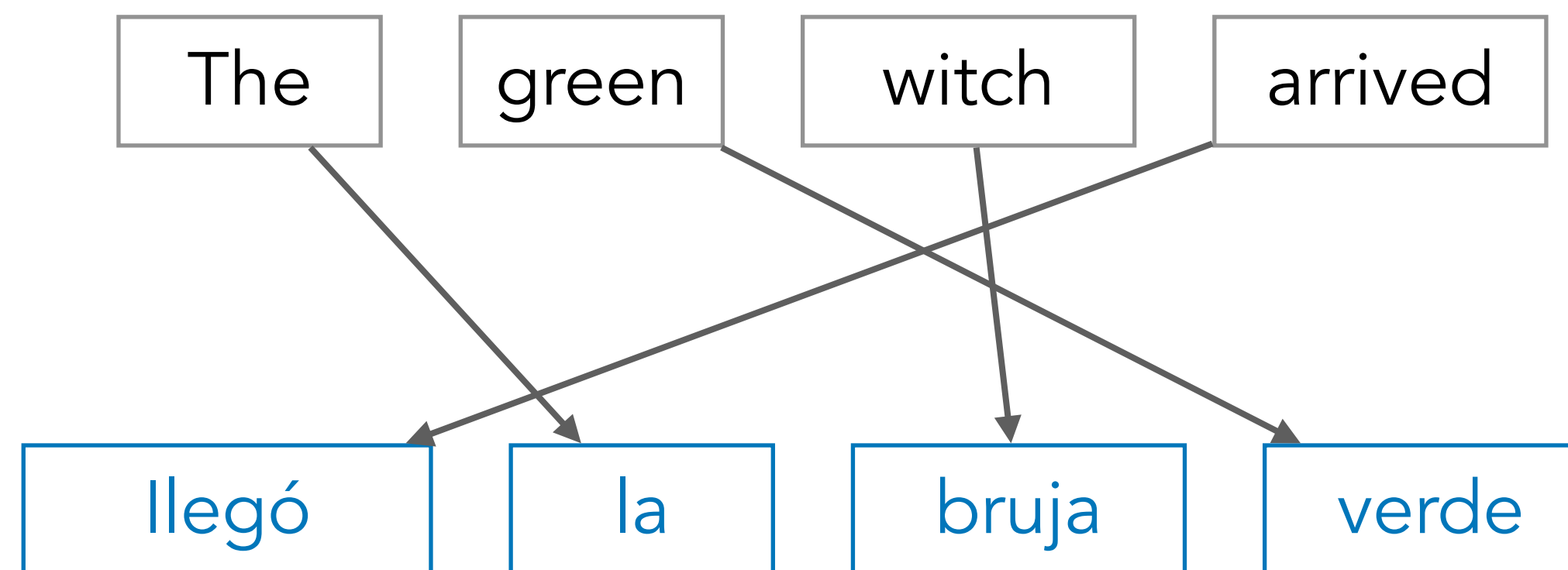
Provide rich task-appropriate context!

- 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)

Sequence-to-Sequence Generation



- 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

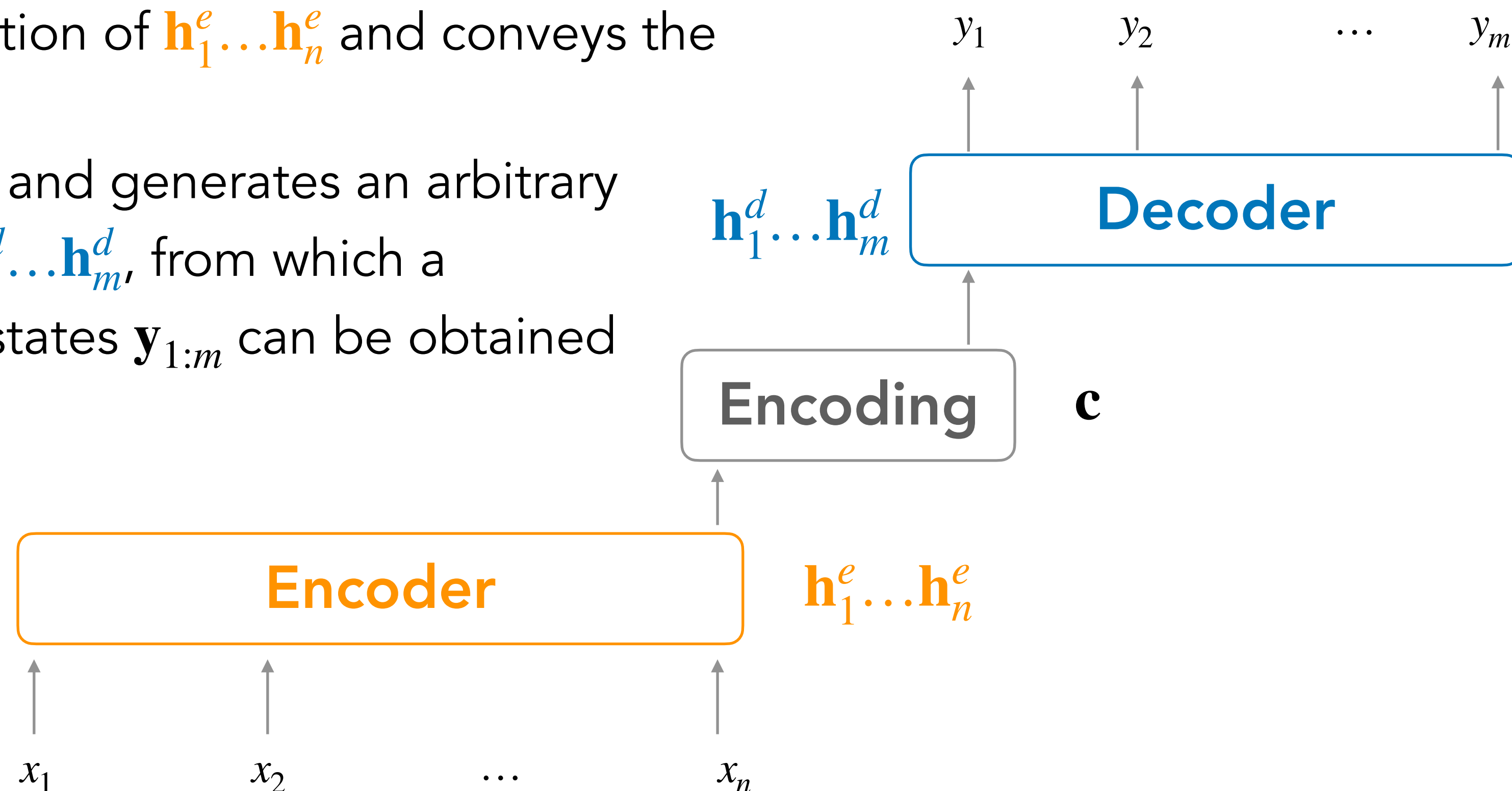
Sequence-to-Sequence Modeling with Encoder-Decoder Networks

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, $\mathbf{h}_1^e \dots \mathbf{h}_n^e$
2. A **encoding** vector, \mathbf{c} which is a function of $\mathbf{h}_1^e \dots \mathbf{h}_n^e$ 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 $\mathbf{h}_1^d \dots \mathbf{h}_m^d$, 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



Produces an encoding of the source sequence

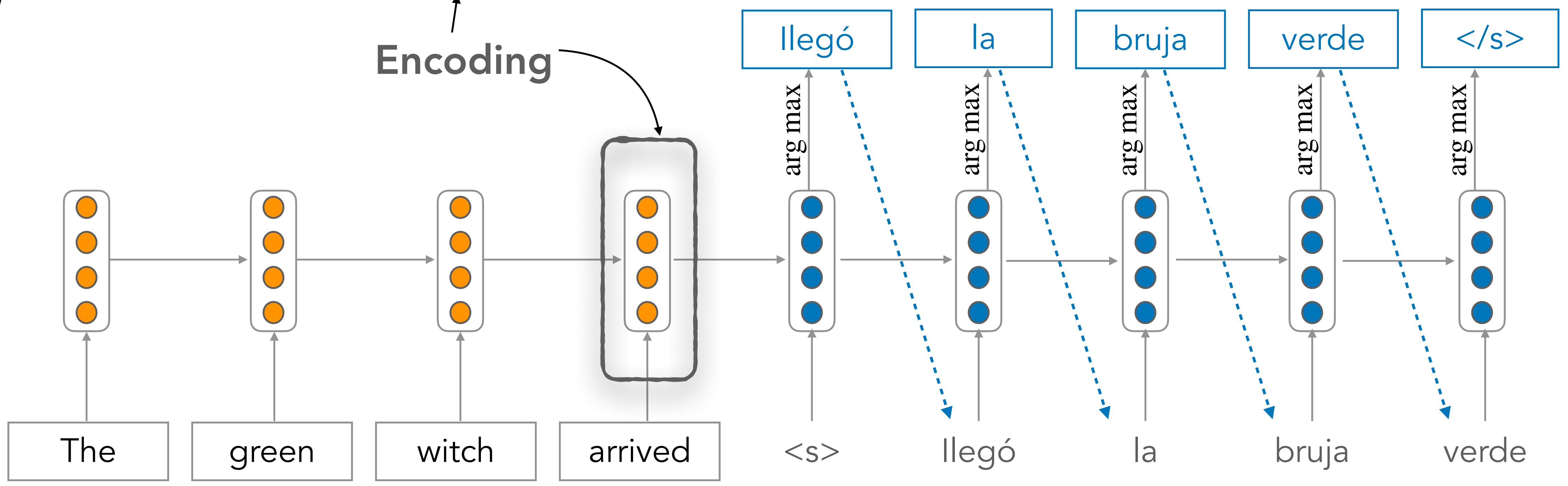
Represents input sequence. Provides initial hidden state for Decoder RNN

Encoding

Target Sentence y

Encoder RNN

Decoder RNN



Source Sentence x

Language Model that produces the target sentence conditioned on the encoding

Encoder RNN

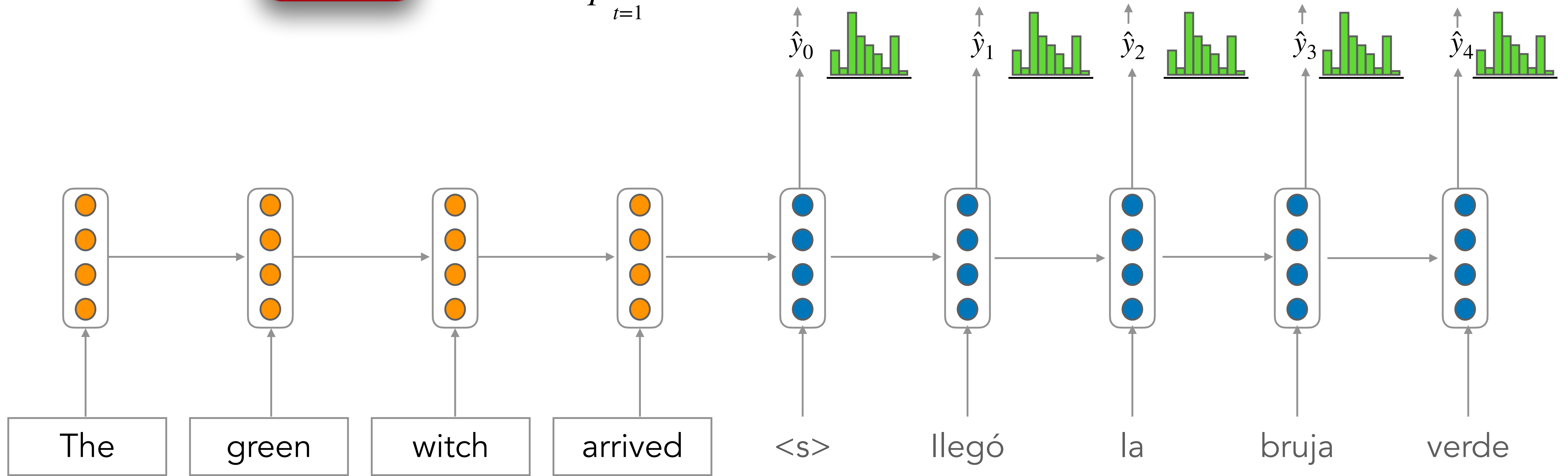
Decoder RNN

Loss

negative log prob. of "llegó"

negative log prob. of "</s>"

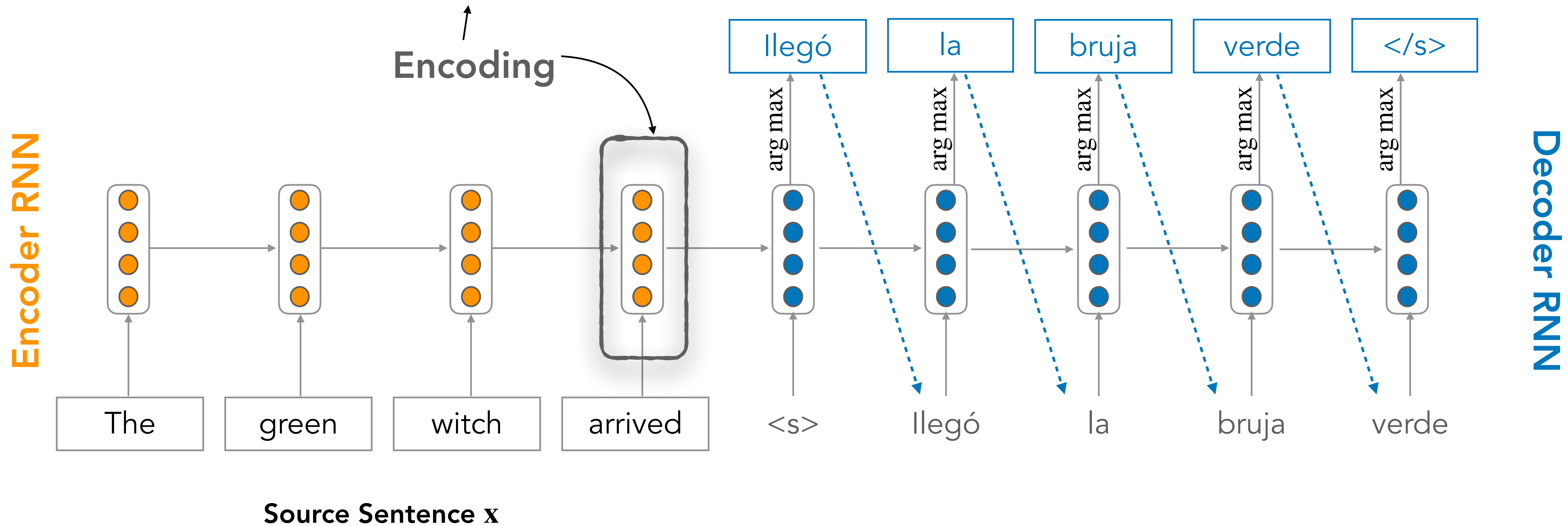
$$L(\theta) = \frac{1}{T} \sum_{t=1}^T L_t(\theta) = L_0(\theta) + L_1(\theta) + L_2(\theta) + L_3(\theta) + L_4(\theta)$$



Source Sentence x

Target Sentence y

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

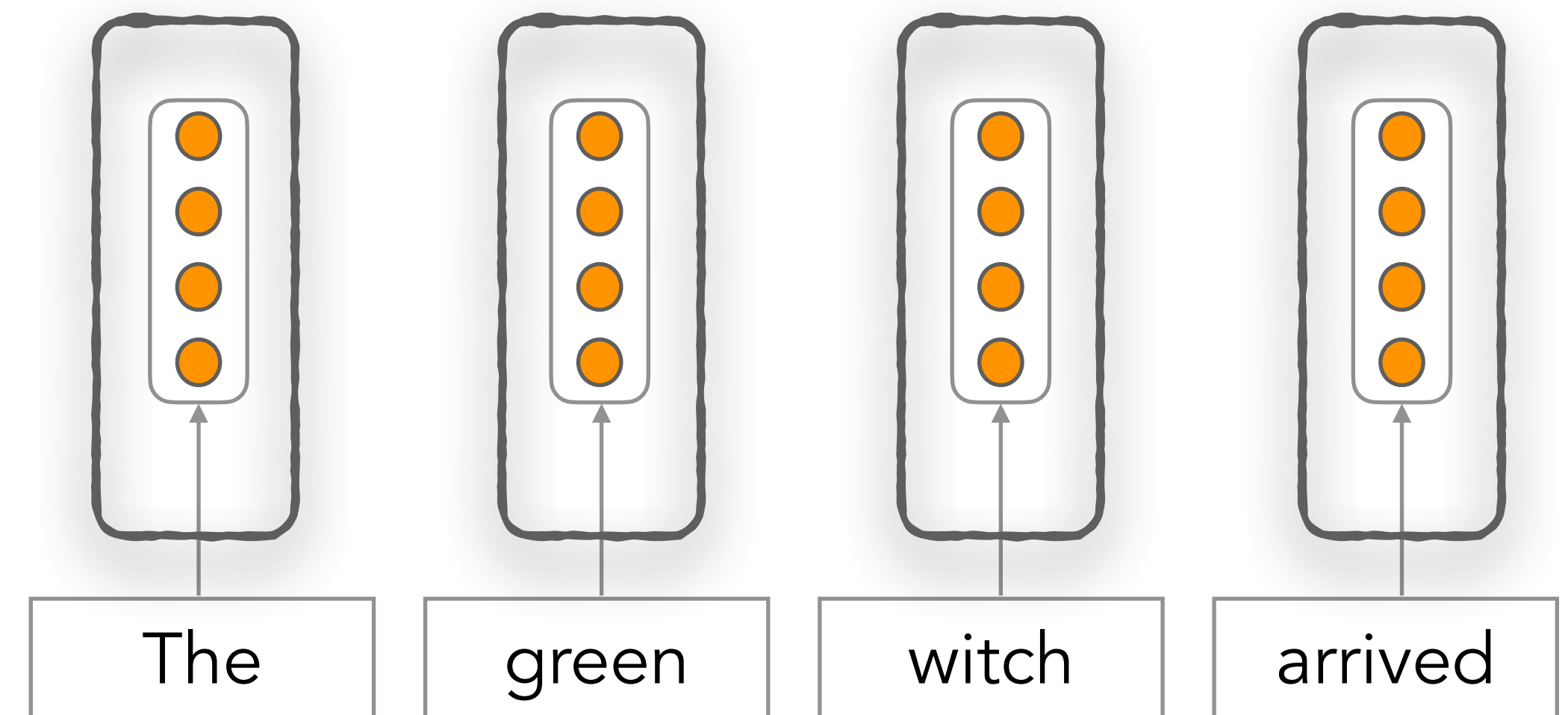
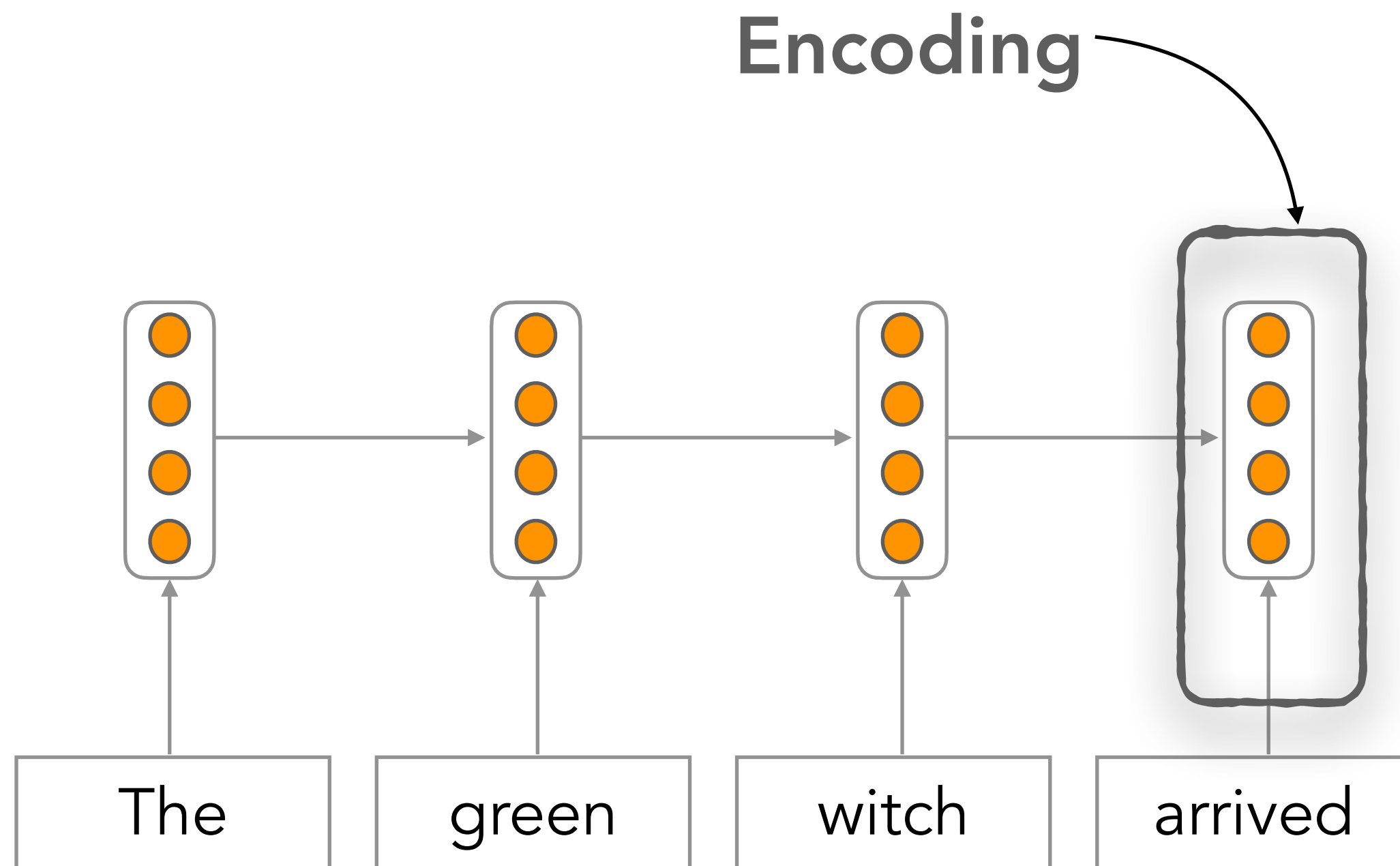


“you can't cram the meaning of a whole sentence into a single vector!”

– Ray Mooney, Professor of Computer Science, UT Austin

Information Bottleneck: One Solution

Encoder RNN



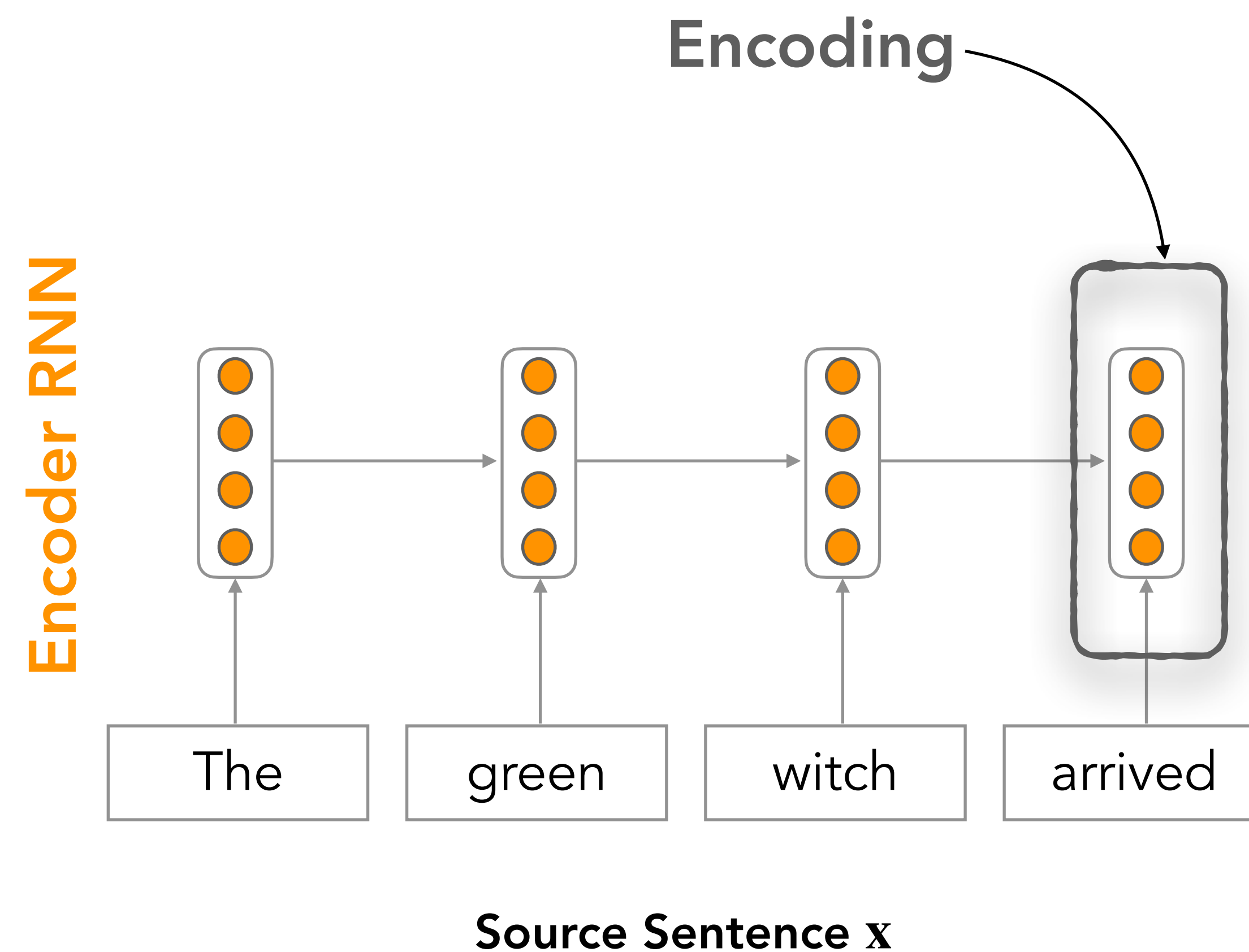
What if we had access to all hidden states?

How to create this?

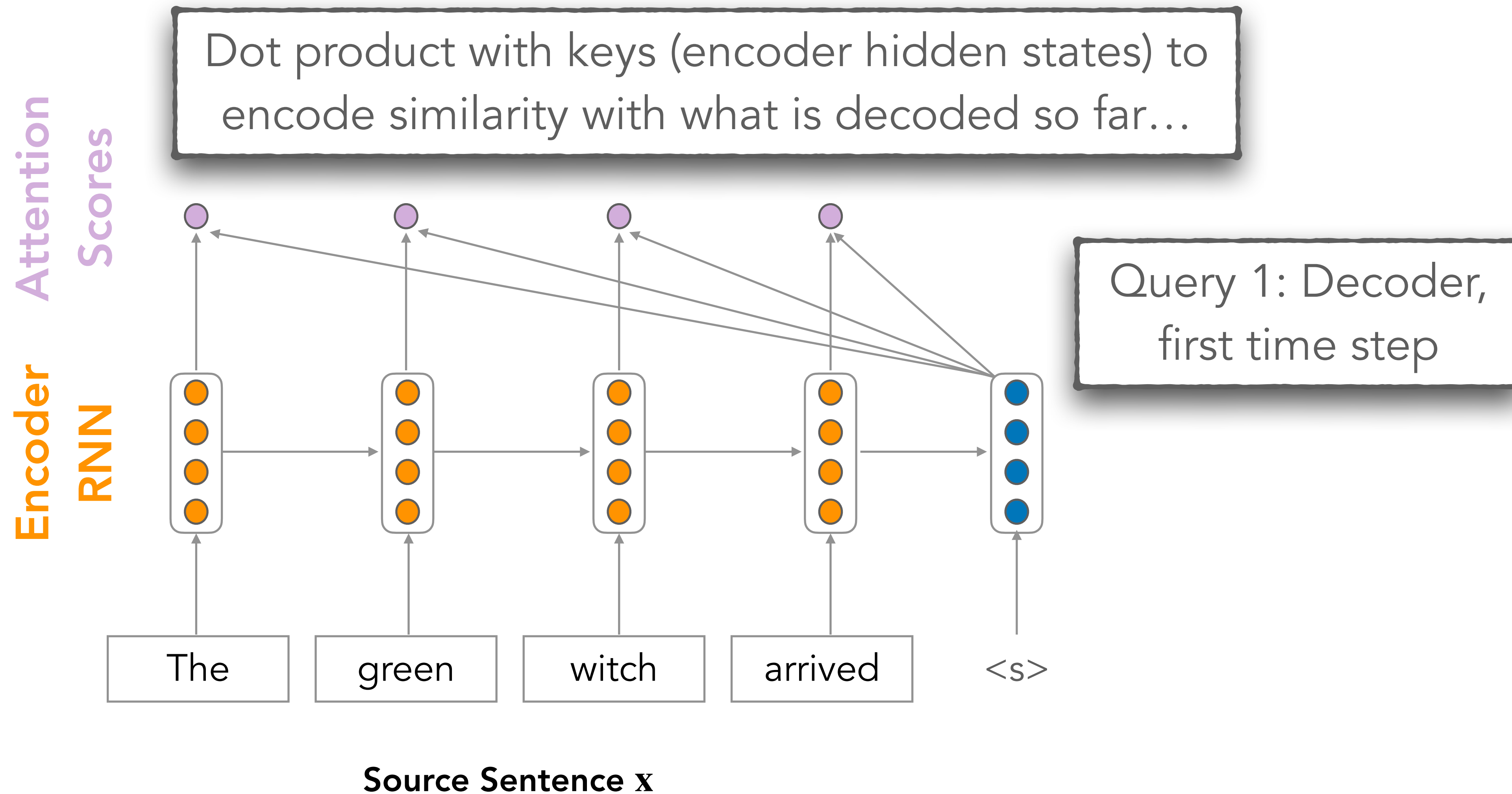
Attention Mechanism

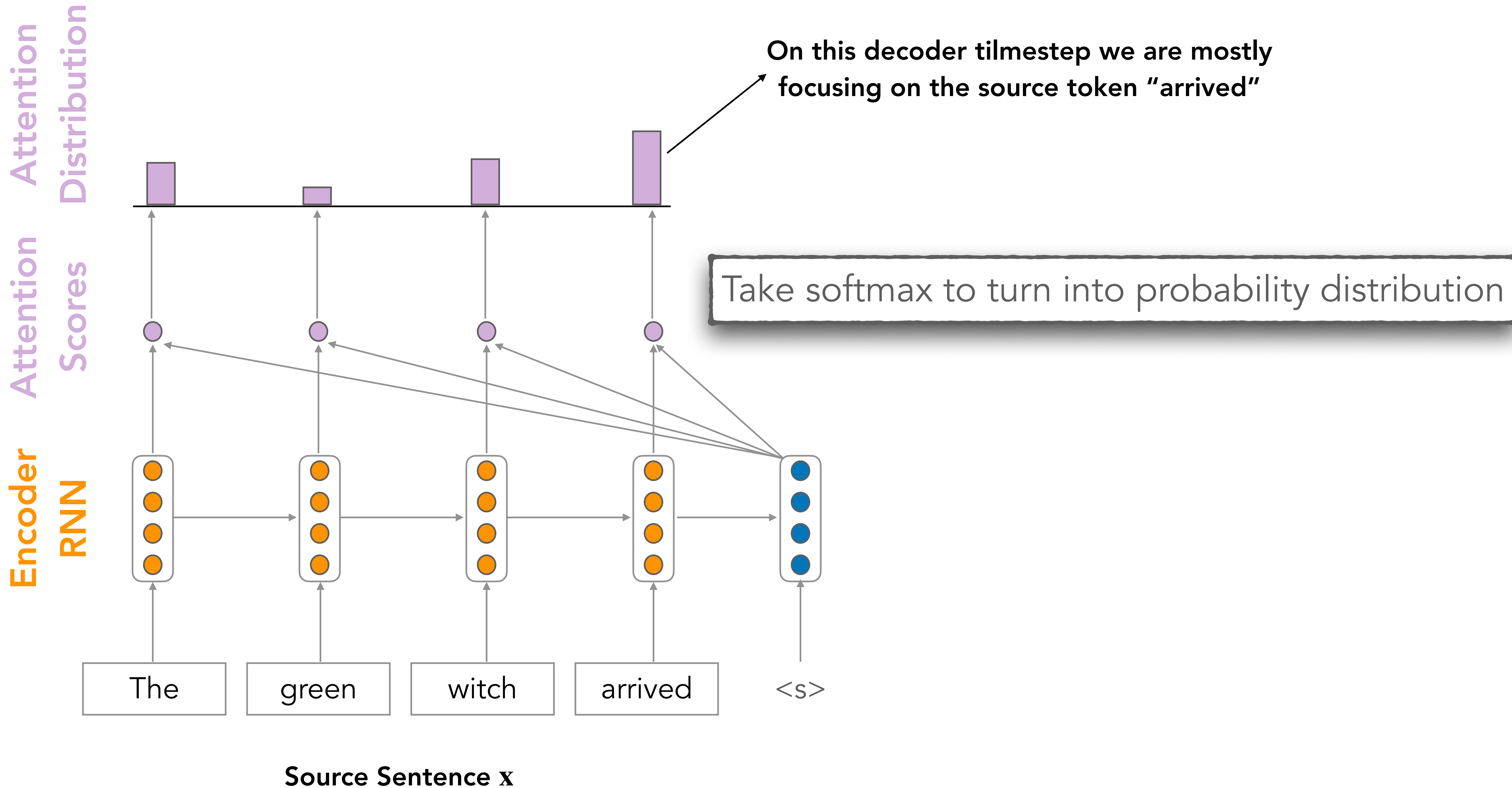
Attention Mechanism

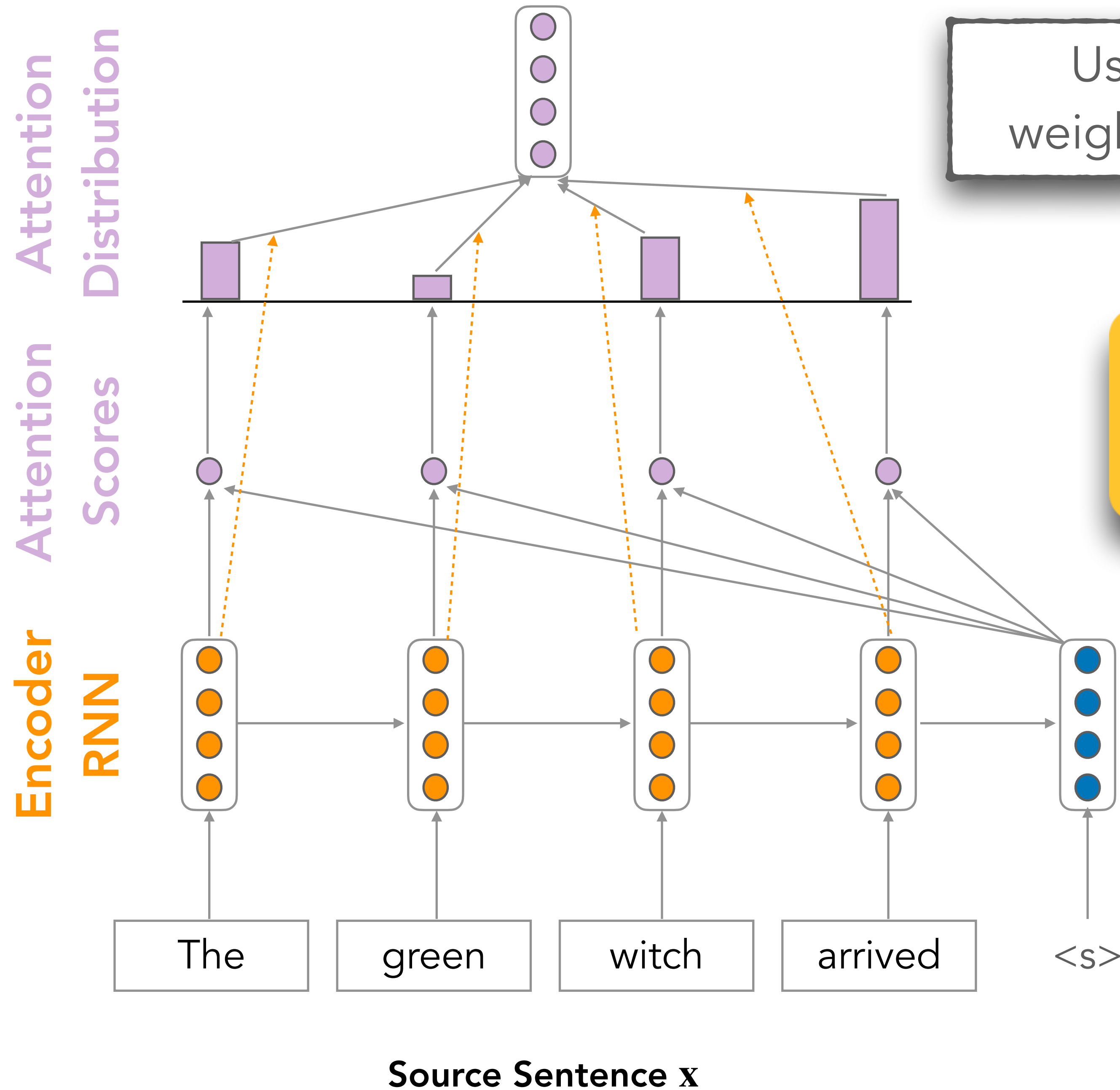
- Attention mechanisms allow the decoder to focus on a particular part of the source sequence at each time step
- Single fixed-length vector \mathbf{c}_t by taking a weighted sum of all the encoder hidden states
 - One per time step *of the decoder!*
- In general, we have a single **query** vector and multiple **key** vectors.
- We want to score each query-key pair



Seq2Seq with Attention

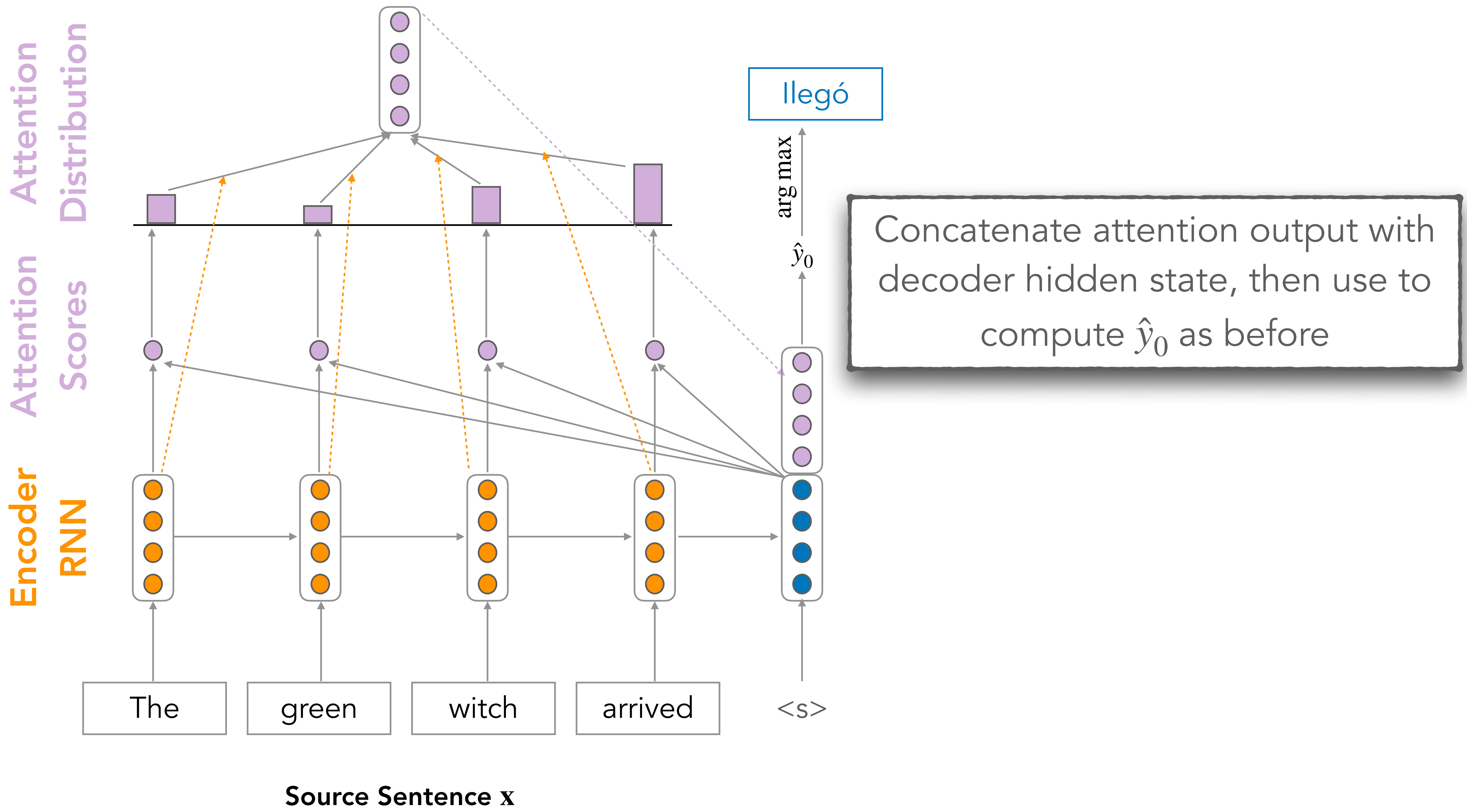


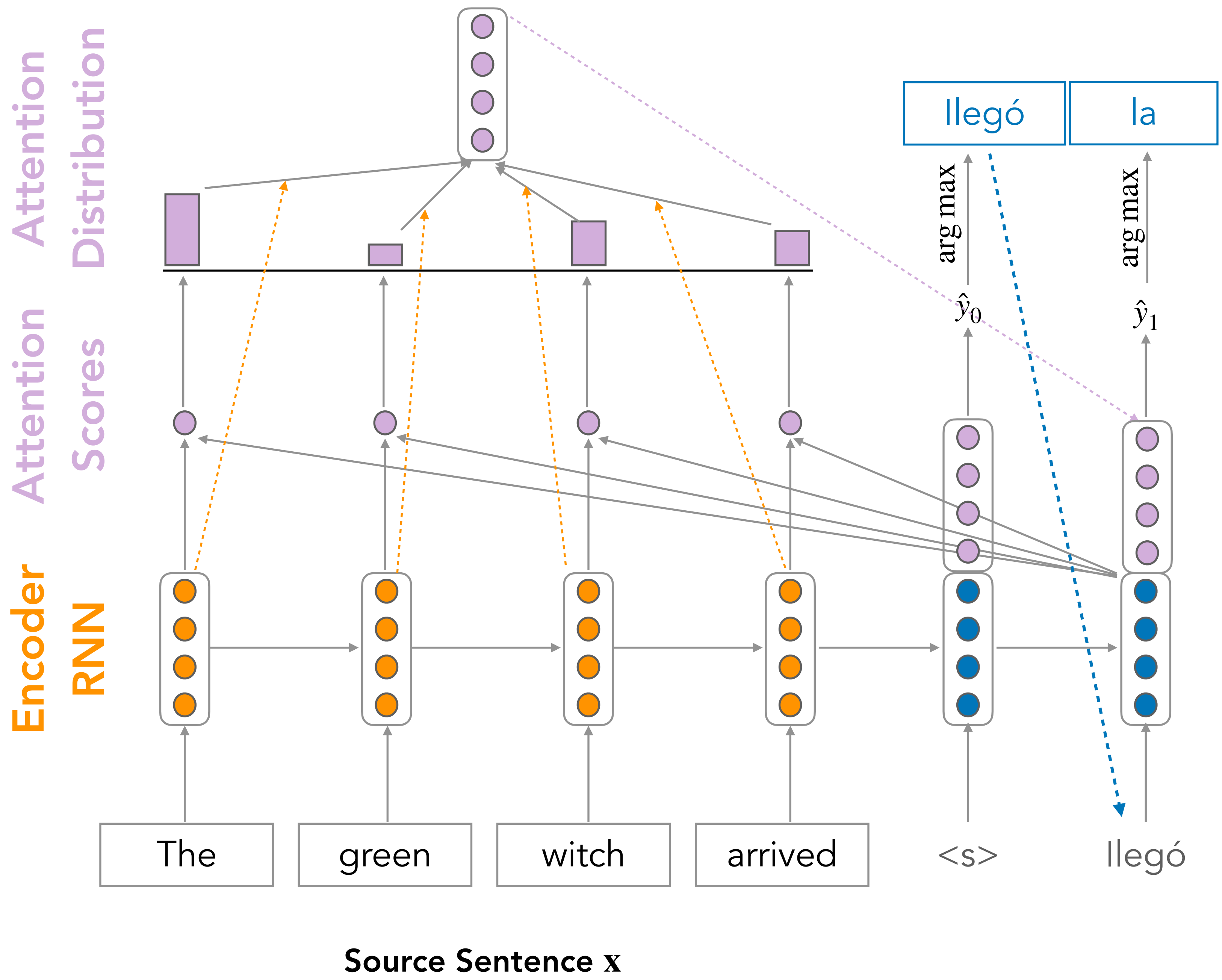




Use the attention distribution to take a weighted sum of the encoder hidden states.

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

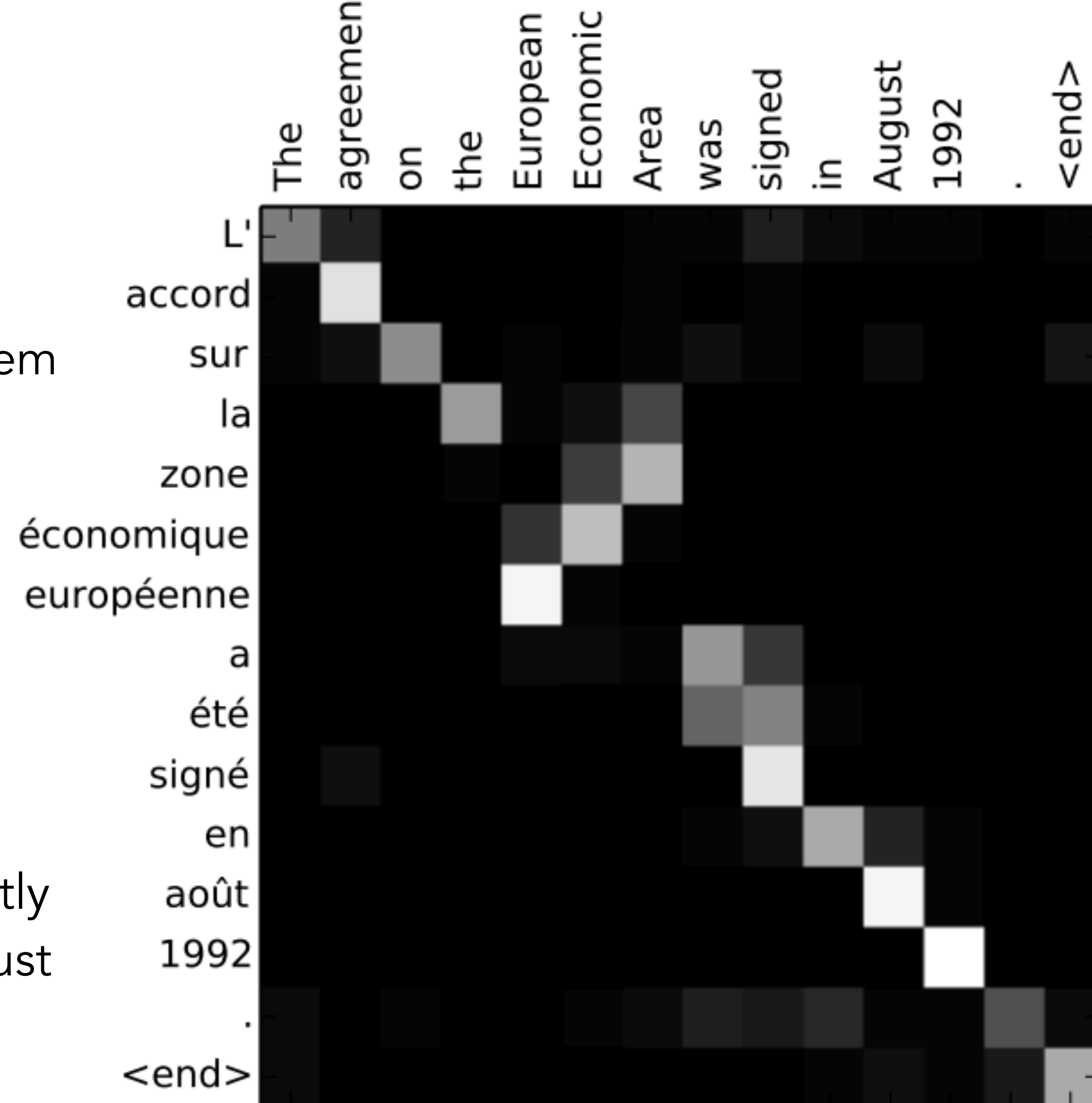




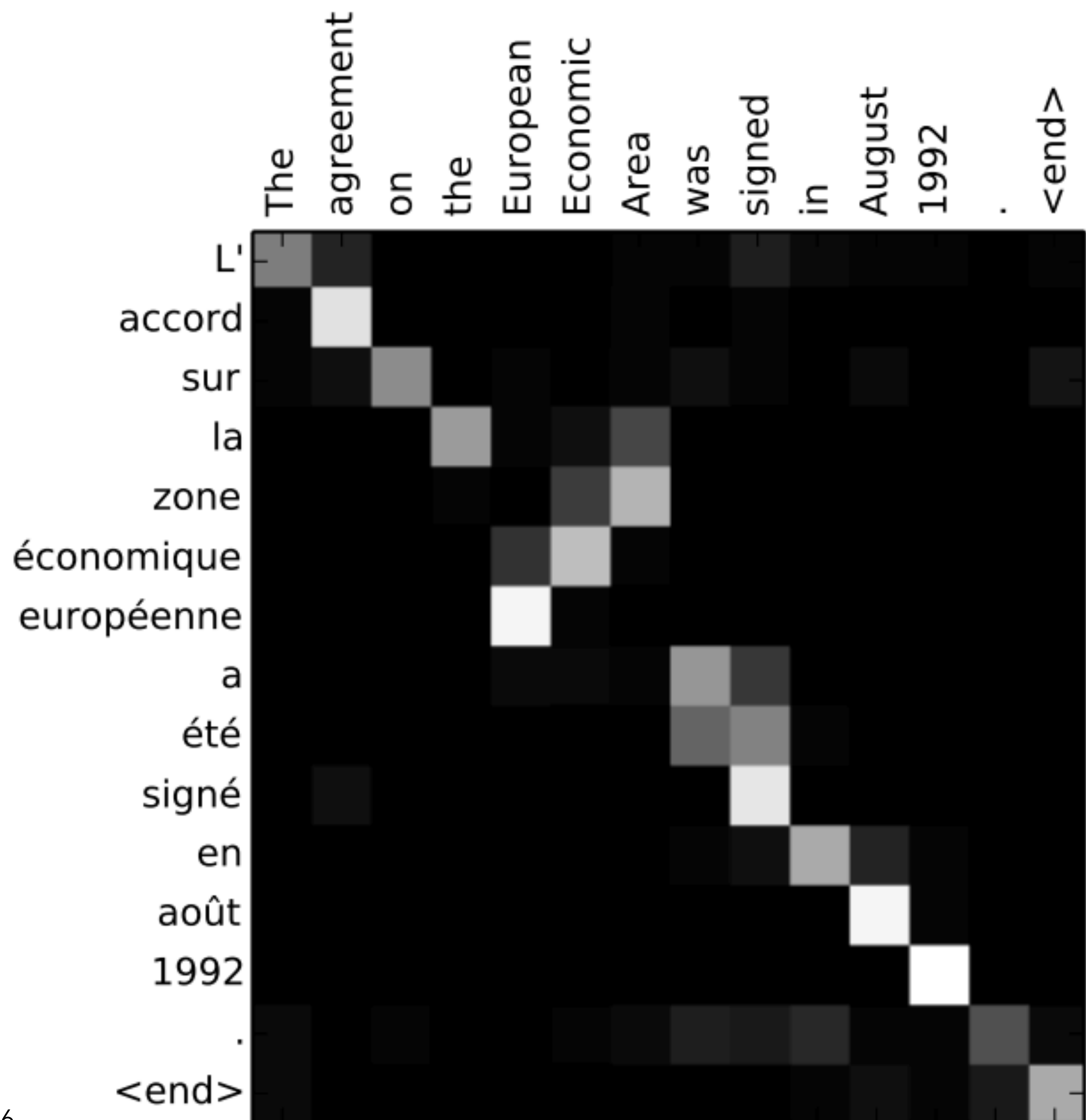
Query 2: Decoder, second time step

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



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 Class: More on attention: self-Attention and Transformers