

# Introduction to Language Models

Tanmoy Chakraborty  
Associate Professor, IIT Delhi  
<https://tanmoychak.com/>

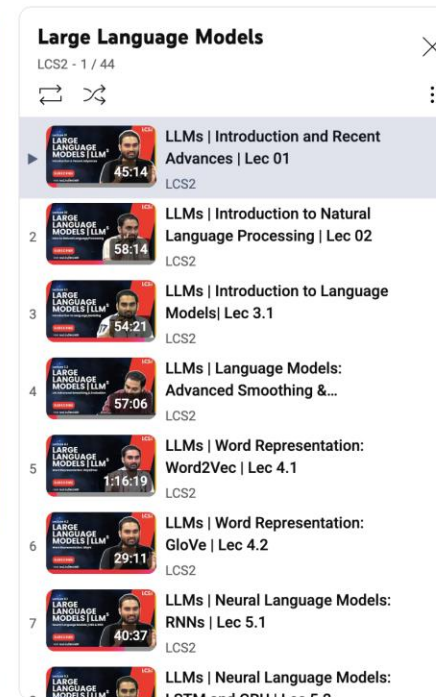
# REMINDER

Last year's (2024) offering



LLMs | Introduction and Recent Advances | Lec 01

<https://lcs2.in/llm2401>

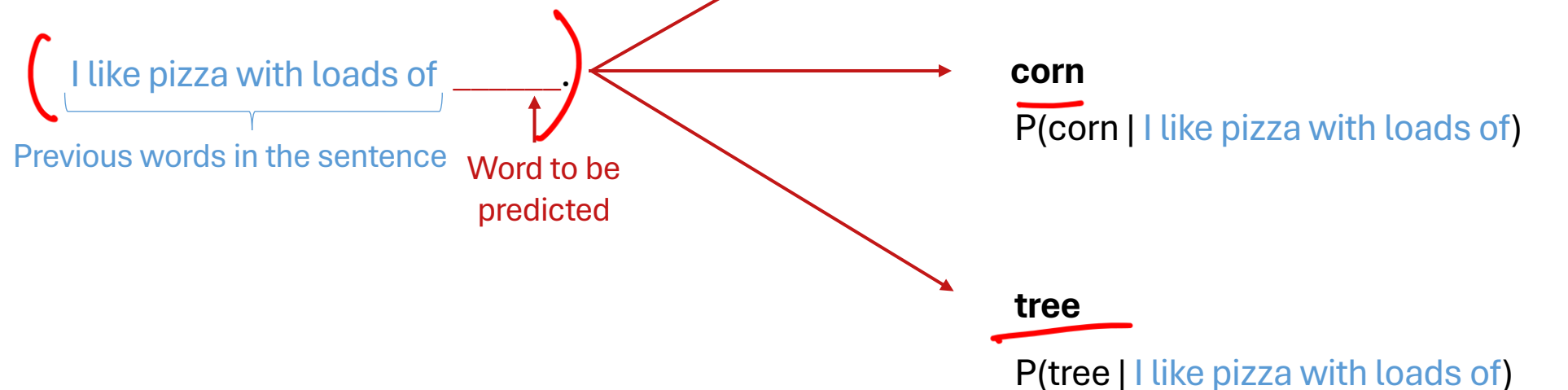


You are advised to study the **first 10 lectures (till Lec 6.1)** of the previous year's course playlist before the **next class on August 4**. Otherwise, you will not be able to follow. Here's the link to the playlist:



# Next Word Prediction

Guess the next word in the sequence...



$P(\text{cheese} \mid \text{I like pizza with loads of}) > P(\text{corn} \mid \text{I like pizza with loads of}) \gg P(\text{tree} \mid \text{I like pizza with loads of})$



# Probabilistic Language Models

- **Goal:** Calculate the probability of a sentence or sequence consisting of  $n$  words

$$P(W) = P(w_1, w_2, w_3, \dots, w_n)$$

or

- **Related Task:** Calculate the probability of the next word conditioned on the preceding words

$$P(w_6 | w_1, w_2, w_3, w_4, w_5)$$

A model that calculates either of these is referred to as a **Language Model (LM)**.



# Estimate Conditional Probabilities

$$P(\text{going} | \text{I am}) = \frac{\text{Count}(\text{I am going})}{\text{Count}(\text{I am})}$$

$$P(\text{begun} | \text{The monsoon season has}) = \frac{\text{Count}(\text{The monsoon season has begun})}{\text{Count}(\text{The monsoon season has})}$$

- **Problem:** Enough data is not available to get an accurate estimate of the above quantities.
- **Solution:** Markov Assumption

# Markov Assumption

**Every next state depends only the previous k states**

- Chain Rule:

$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i | w_1 w_2 \dots w_{i-1})$$

- Applying Markov Assumption we condition on only the preceding k words:

$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i | w_{i-k} \dots w_{i-1})$$

- Probabilistic Language Models exploit the **Chain Rule of Probability** and **Markov Assumption** to build a probability distribution over sequences of words.



# N-gram Language Models

n-gram  
n=1 unigram  
n=2 bigram  
 $\frac{I \text{ am}}{P(a|I)}$

- Let's consider the following conditional probability:

$P(\text{begun} \mid \text{the monsoon season has})$

- An **N-gram model** considers only the preceding **N - 1 words**.
  - Unigram:  $P(\text{begun})$
  - Bigram:  $P(\text{begun} \mid \text{the})$
  - Trigram:  $P(\text{begun} \mid \text{the monsoon})$

**Relation between Markov model and Language Model:**

An N-gram Language Model  $\equiv$  (N - 1) order Markov Model



# Limitation of N-gram Language Models

Stats → tram LM  
↓  
Neural LM

- An insufficient model of language since they are **not effective in capturing long-range dependencies present in language.**

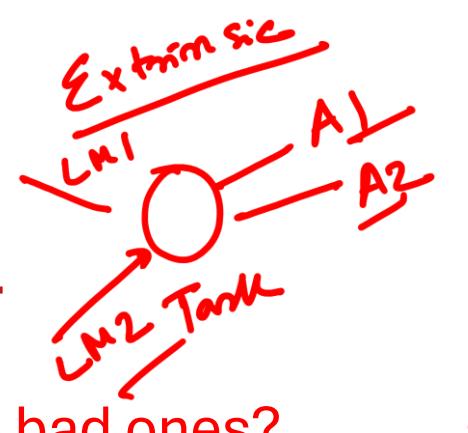
- Example:

✓ The **project**, which he had been working on for months, was finally **approved** by the committee.

The above example highlights the long-distance dependency between “project” and “approved”, where the context provided by earlier words affects the interpretation of later parts of the sentence.



# Evaluation of a Language Model



- Does our language model prefer good sentences to bad ones?
  - Assign higher probability to “real” or “frequently observed” sentences than “ungrammatical” or “rarely observed” sentences
- Terminologies:
  - We optimize the parameters of our model based on data from a **training set**.
  - We assess the model's performance on unseen **test data** that is disjoint from the training data.
  - An evaluation metric provides a measure of the performance of our model on the test set.

# Extrinsic Evaluation

- Measure the effectiveness of a language model by **testing their performance on different downstream NLP tasks**, such as machine translation, text classification, speech recognition.
- Let us consider two different language models: A and B
  - Select a suitable evaluation metric to assess the performance of the language models based on the chosen task.
  - Obtain the evaluation scores for A and B
  - Compare the evaluation scores for A and B



# Intrinsic Evaluation: Perplexity

Thus, for the sentence  $W$ , perplexity is:

$$PP(W) = P(w_1 w_2 \dots w_n)^{-\frac{1}{n}}$$

Applying Chain Rule:

$$PP(W) = \left( \prod \frac{1}{P(w_i | w_1 w_2 \dots w_{i-1})} \right)^{\frac{1}{n}}$$

Applying Markov Assumption ( $n = 2$ ), i.e. **for bigram LM**:

$$PP(W) = \left( \prod \frac{1}{P(w_i | w_{i-1})} \right)^{\frac{1}{n}}$$

Handwritten red annotations showing the derivation of the perplexity formula. It starts with a large expression:  $\sqrt[n]{P(w_1 \dots w_n)}$  with an upward arrow. Below it, the expression is expanded as  $\sqrt[n]{P(w_1) \times P(w_2|w_1) \times \dots \times P(w_n|w_{n-1})}$ .

Minimizing perplexity is the same as maximizing probability.

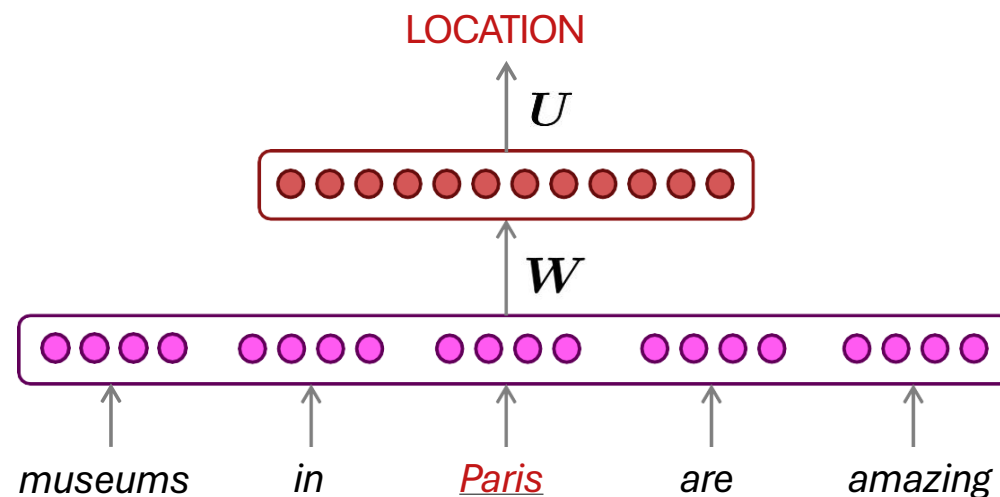


# Neural Language Models

# How to Build a *Neural* Language Model?

- Recall the Language Modeling task:
  - Input:** sequence of words  $x^{(1)}, x^{(2)}, \dots, x^{(t)}$
  - Output:** probability distribution of the next word  $P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$
- How about a window-based neural model?

## Example: NER Task



Word2Vec  
Glove  
FastText



# A Fixed-window Neural Language Model

output distribution

$$\hat{y} = \text{softmax}(U\mathbf{h} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

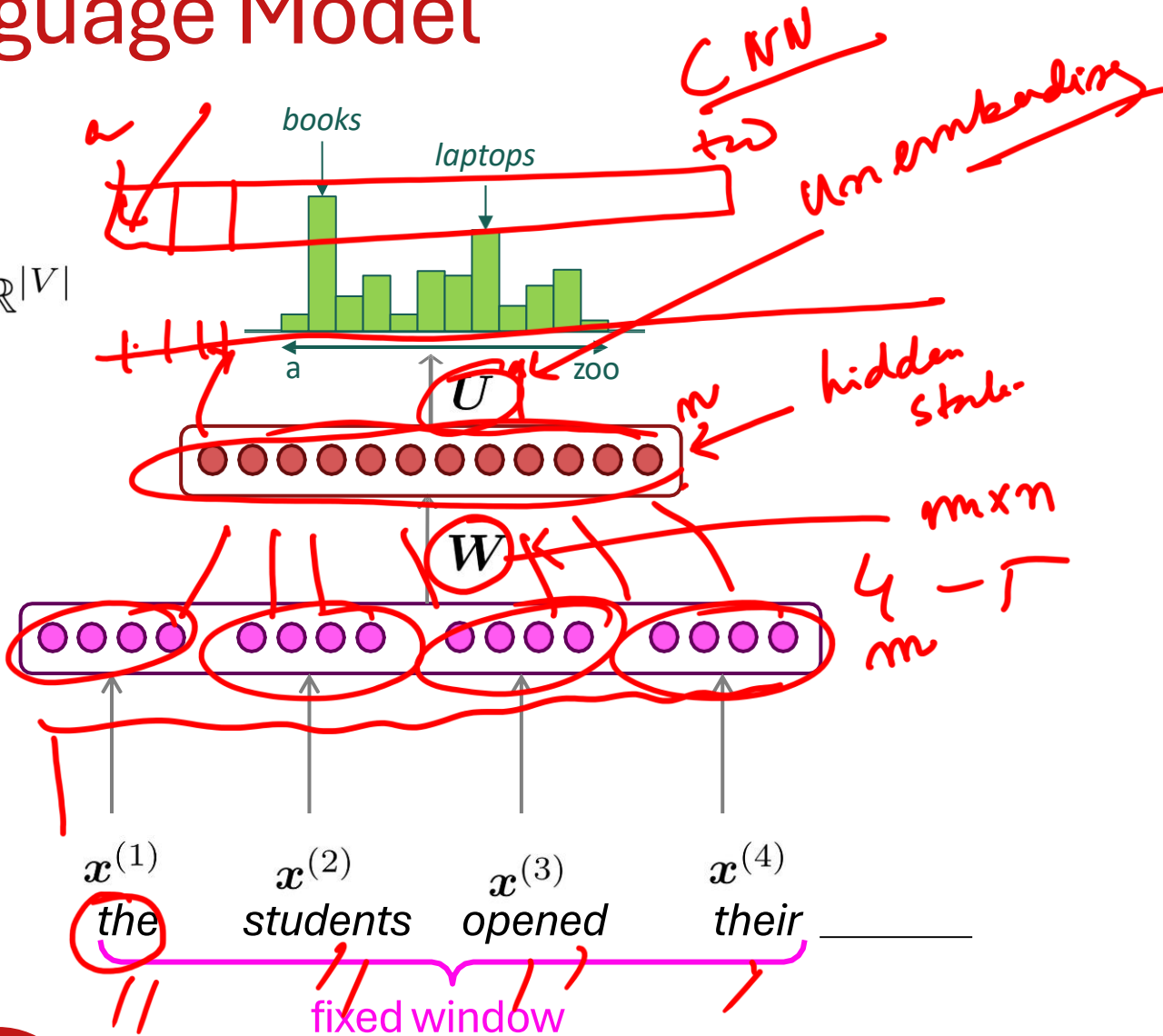
concatenated word embeddings

$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$

~~as the preator started the clock~~



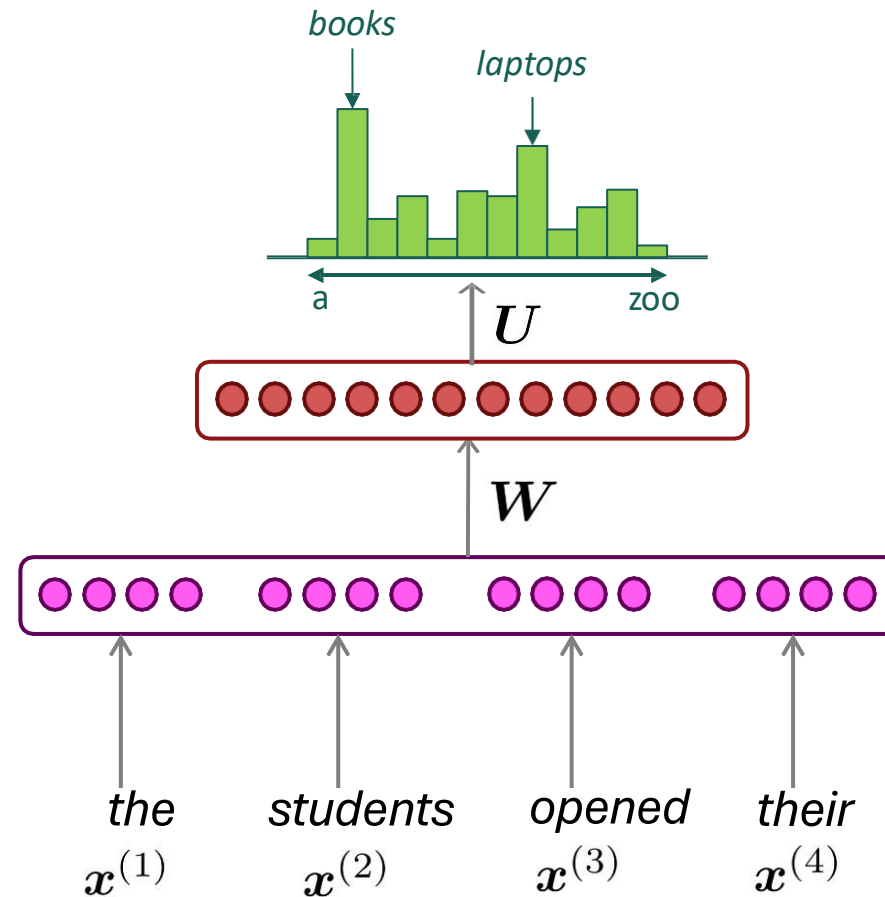
# A Fixed-window Neural Language Model

## Improvements over $n$ -gram LM:

- No sparsity problem
- Don't need to store all observed  $n$ -grams

## Remaining problems:

- Fixed window is **too small** ✓
- Enlarging window enlarges  $W$



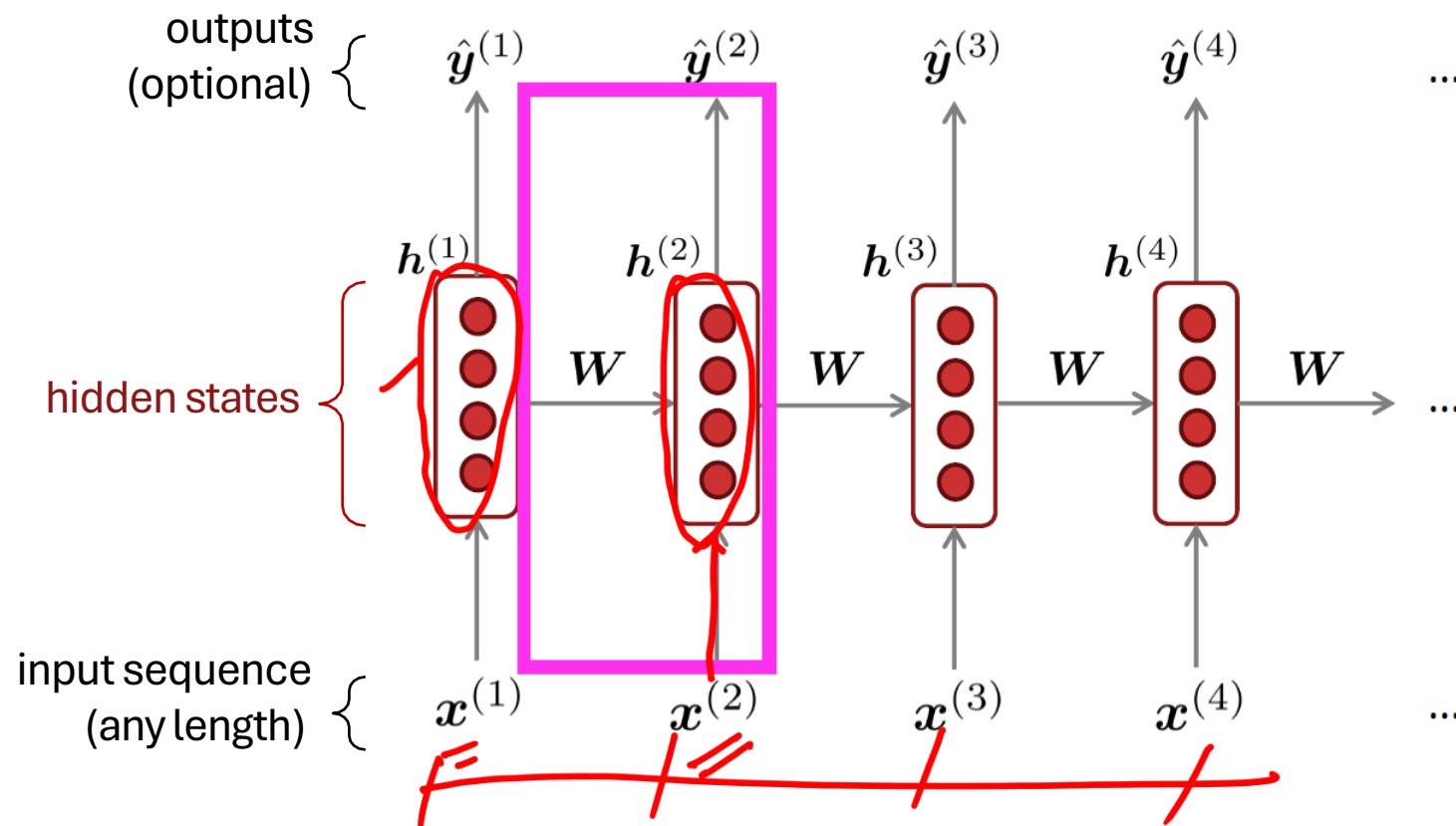
Approximately: Y. Bengio, et al.  
(2000/2003): A Neural Probabilistic  
Language Model

We need a neural  
architecture that can  
process **any length**  
**input**



# Recurrent Neural Networks (RNN)

$$h^2 = f(x^2, h^1)$$



**Core idea:** Apply the same weights  $W$  repeatedly





# A Simple RNN Language Model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(U h^{(t)} + b_2) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$$

$h^{(0)}$  is the initial hidden state

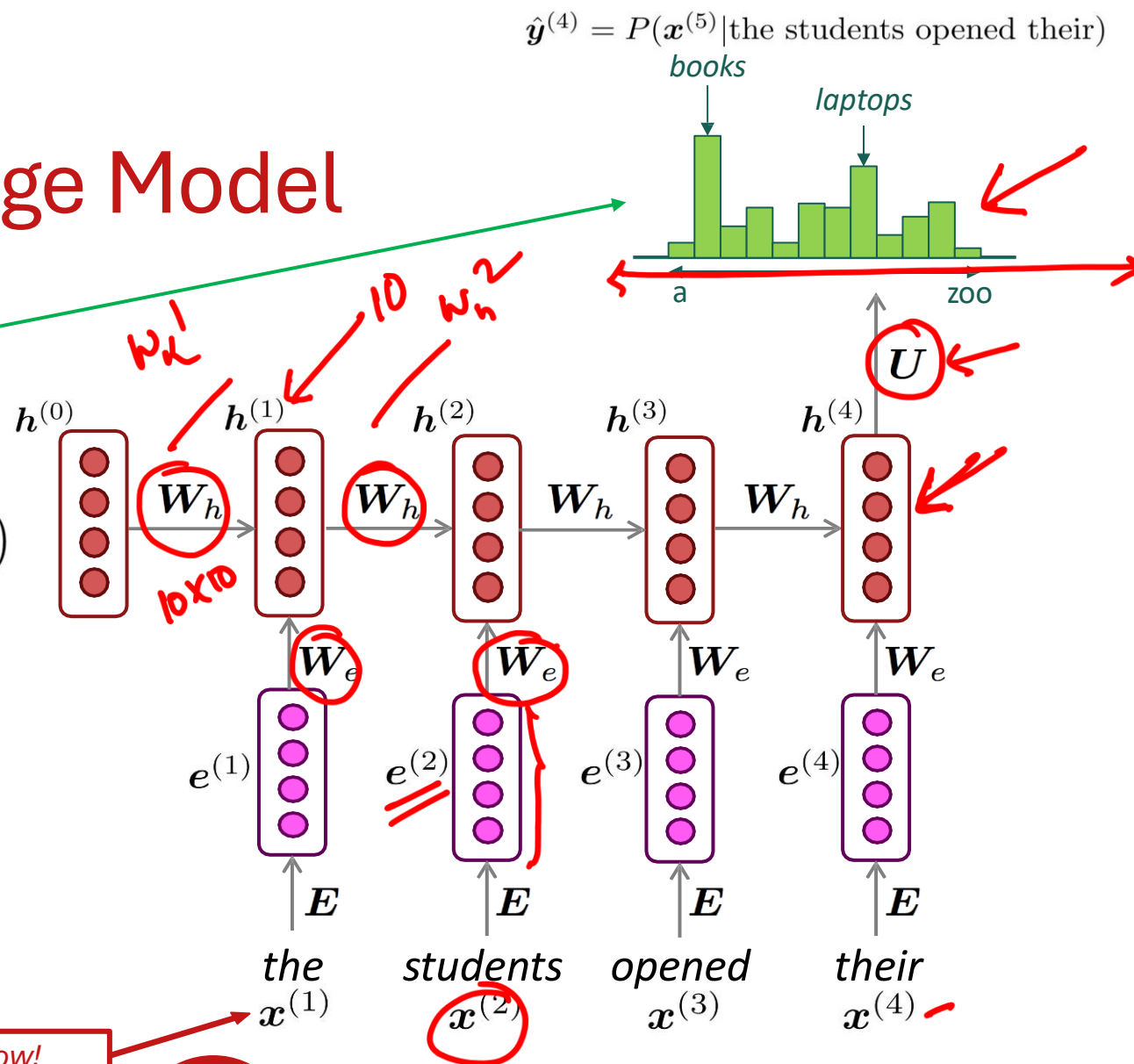
word embeddings

$$e^{(t)} = E x^{(t)}$$

words / one-hot vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

**Note:** this input sequence could be much longer now!



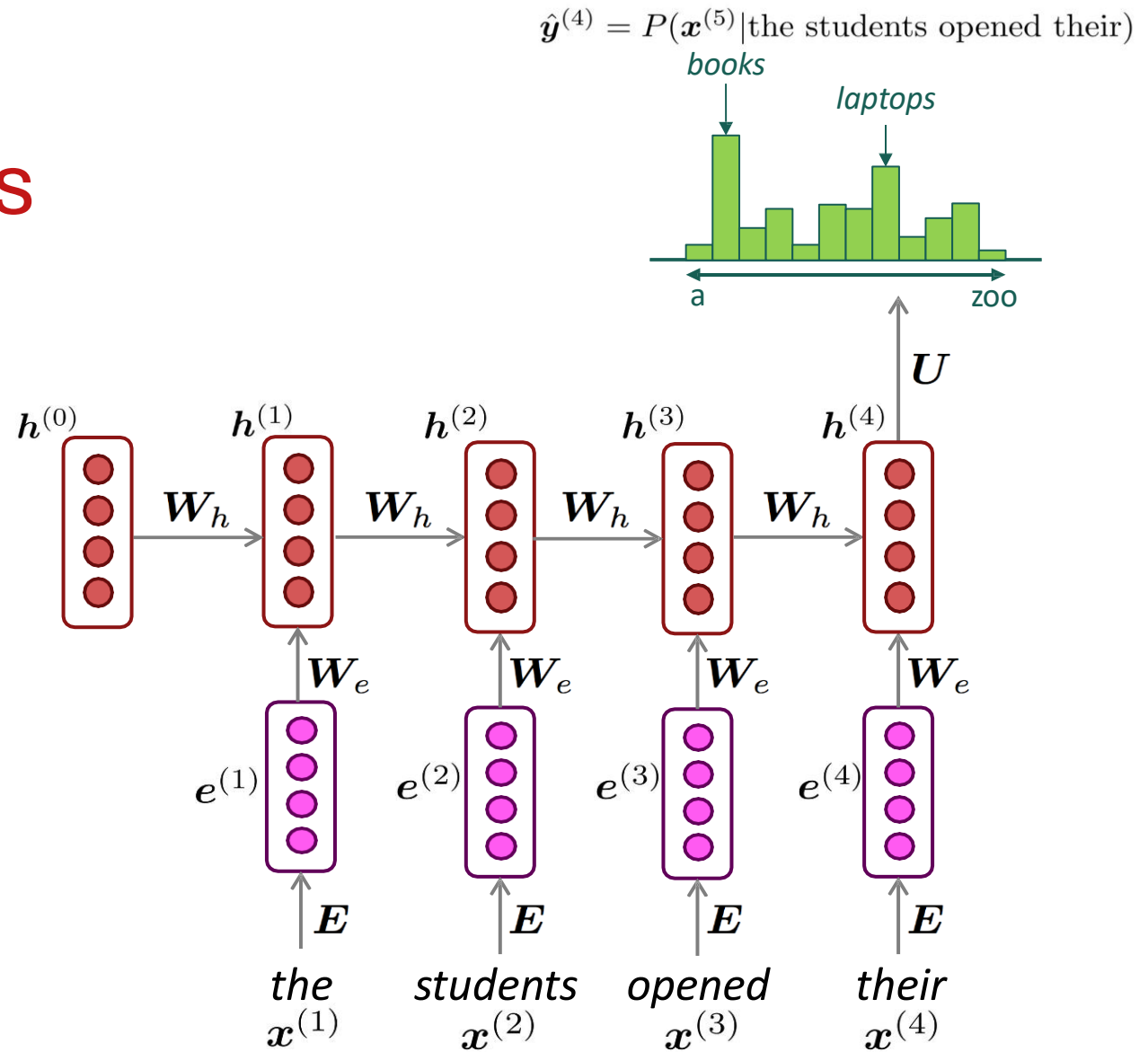
# RNN Language Models

## RNN Advantages:

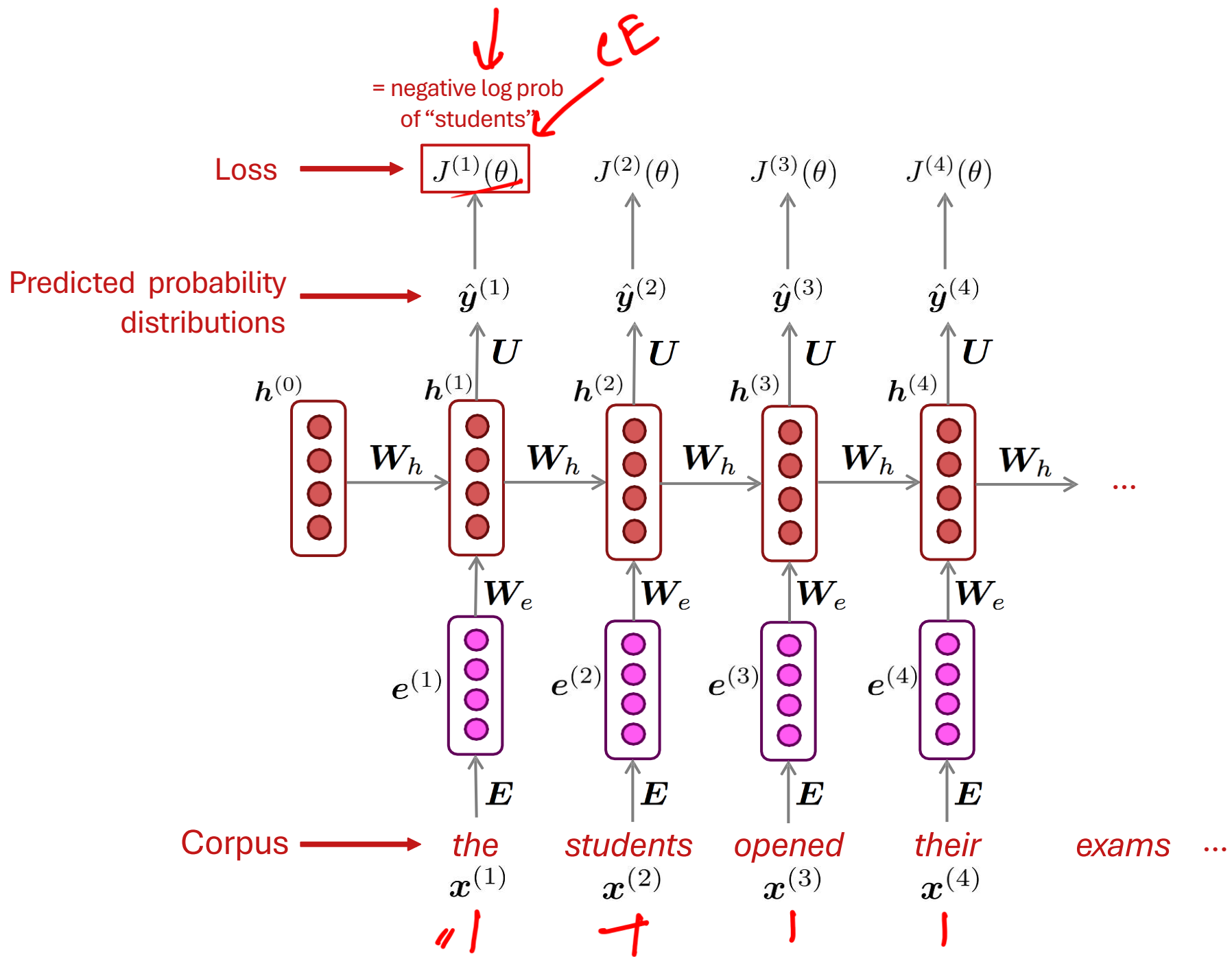
- Can process **any length** input
- Computation for step  $t$  can (in theory) use information from **many steps back**
- **Model size doesn't increase** for longer input context
- Same weights applied on every timestep, so there is **symmetry** in how inputs are processed.

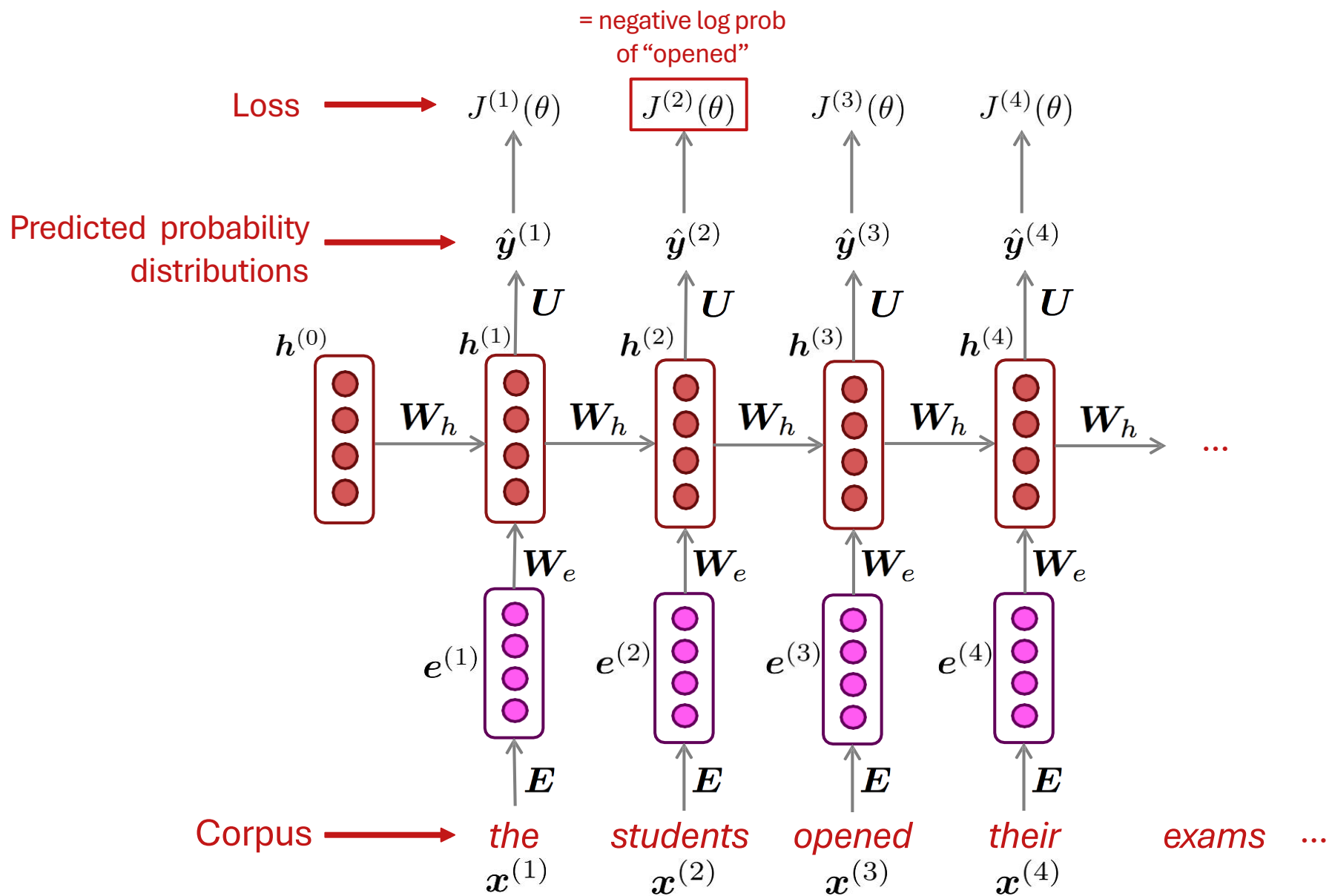
## RNN Disadvantages:

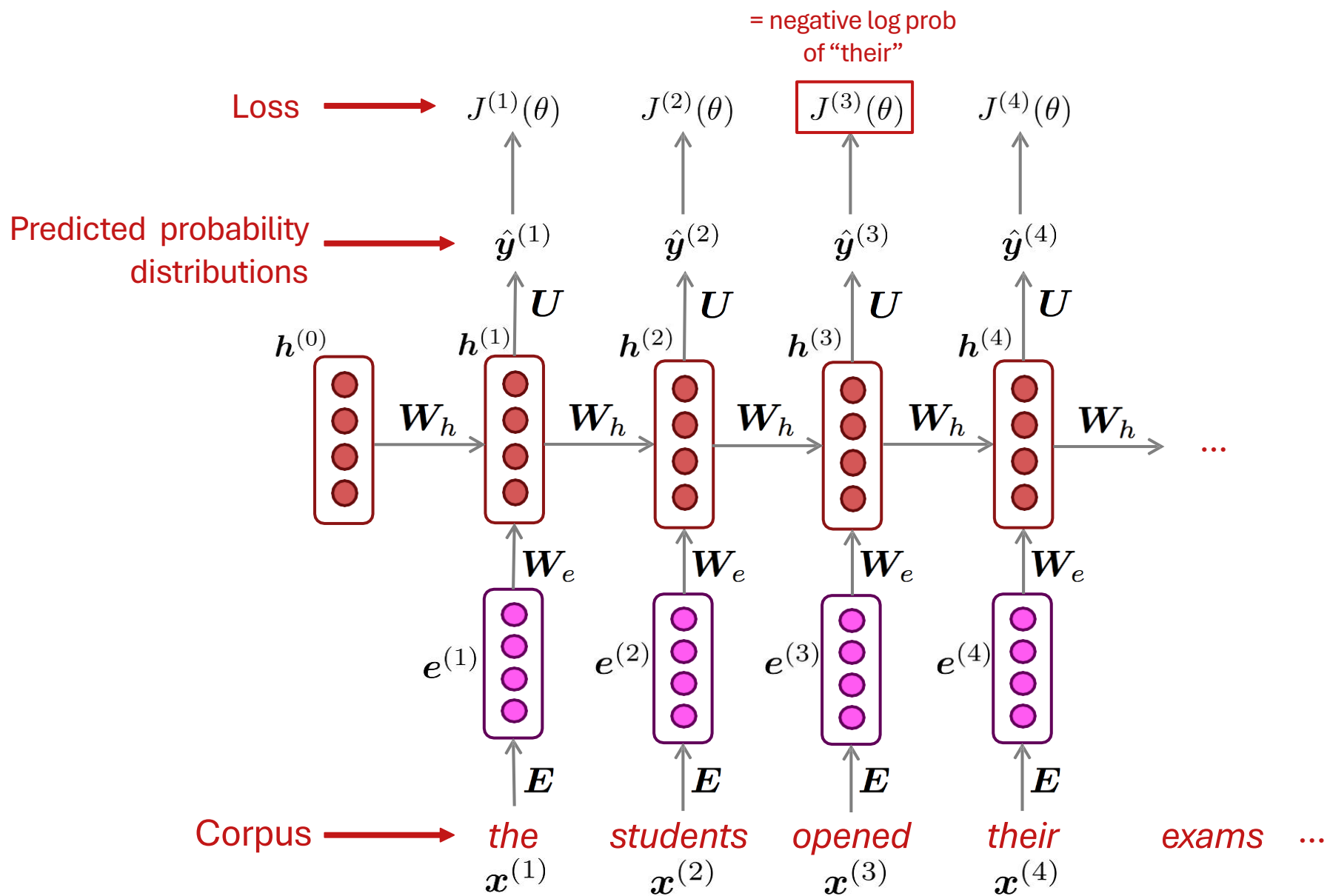
- Recurrent computation is **slow**
- In practice, difficult to access information from **many steps back**

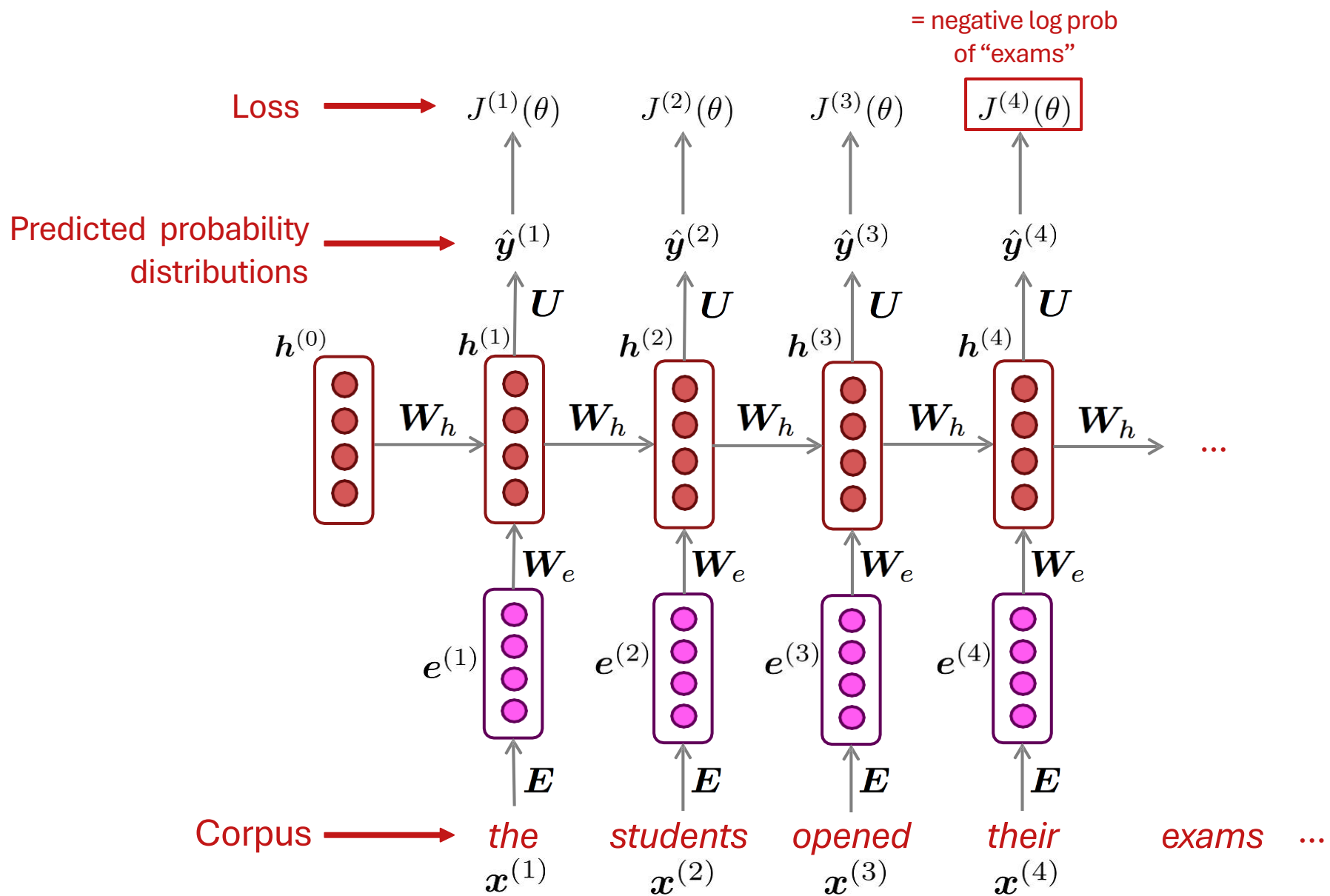


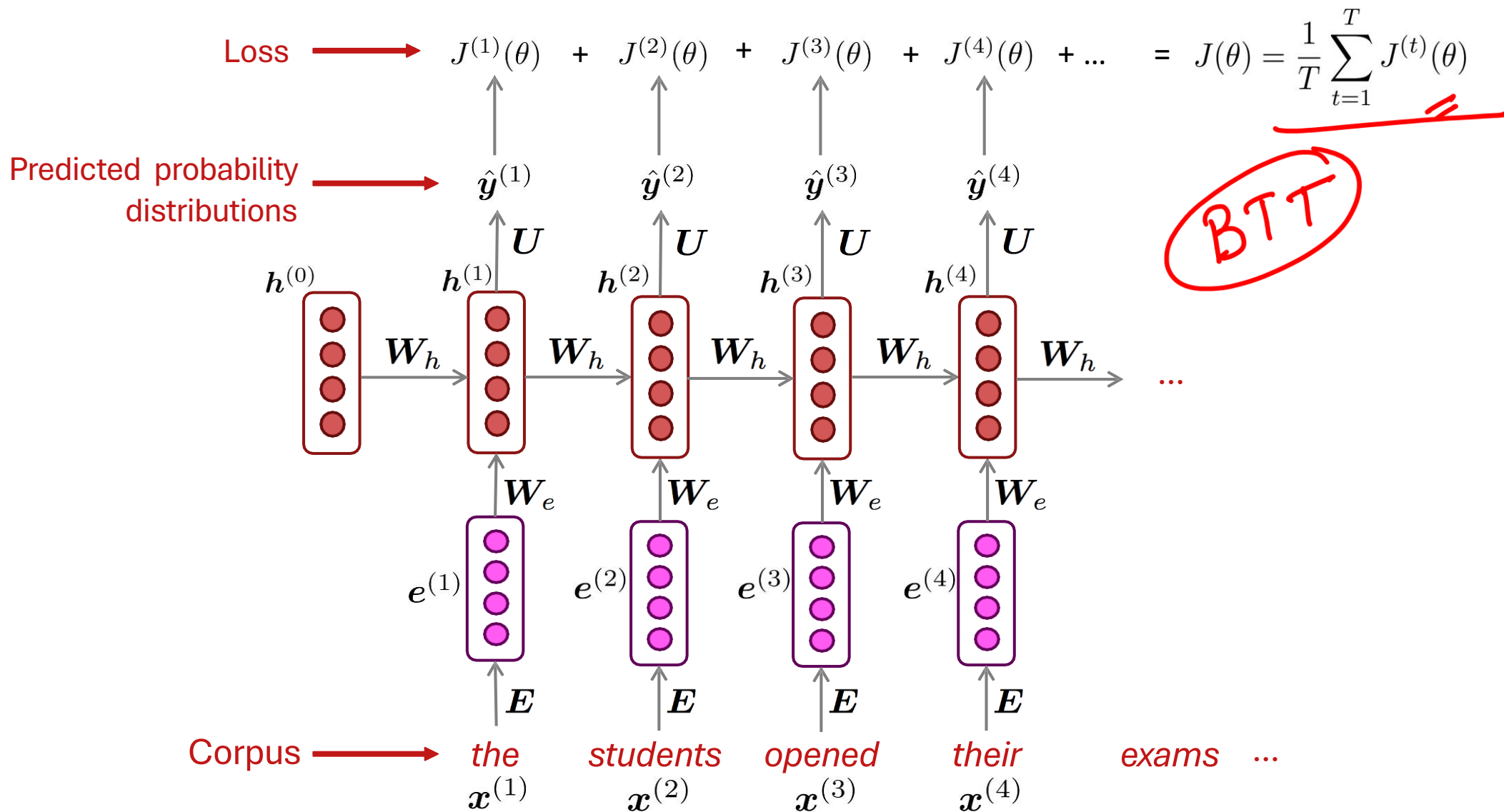
# Training an RNN Language Model







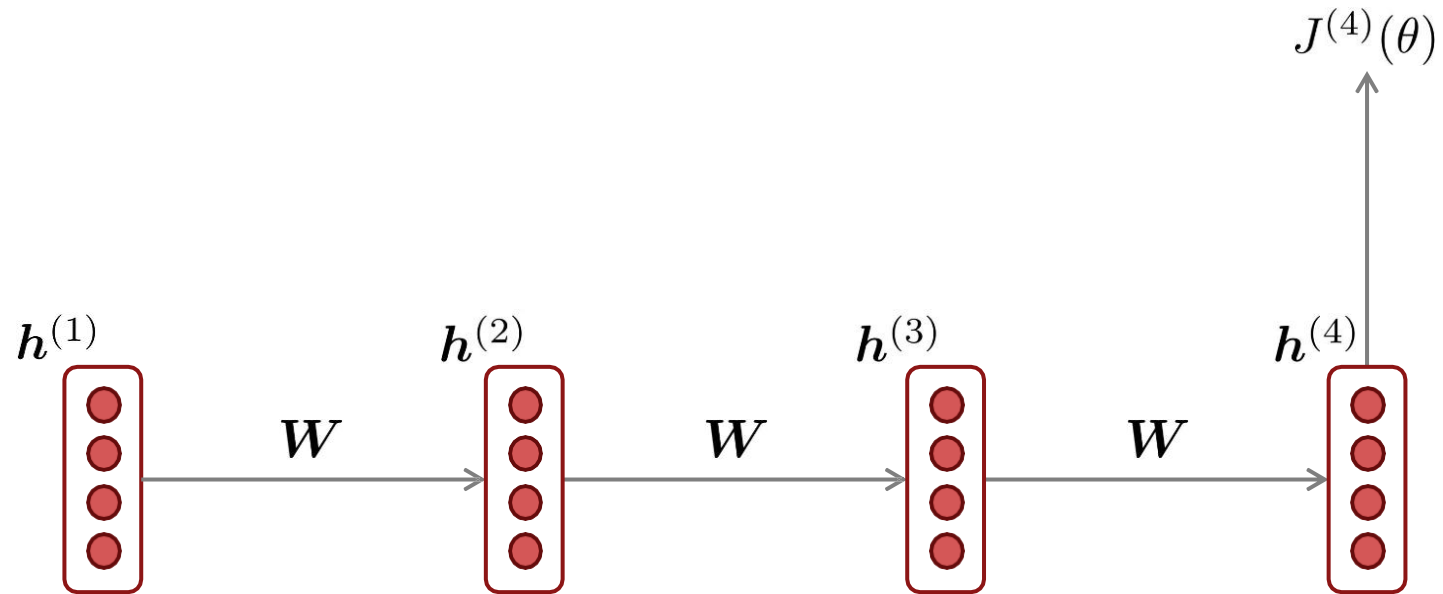




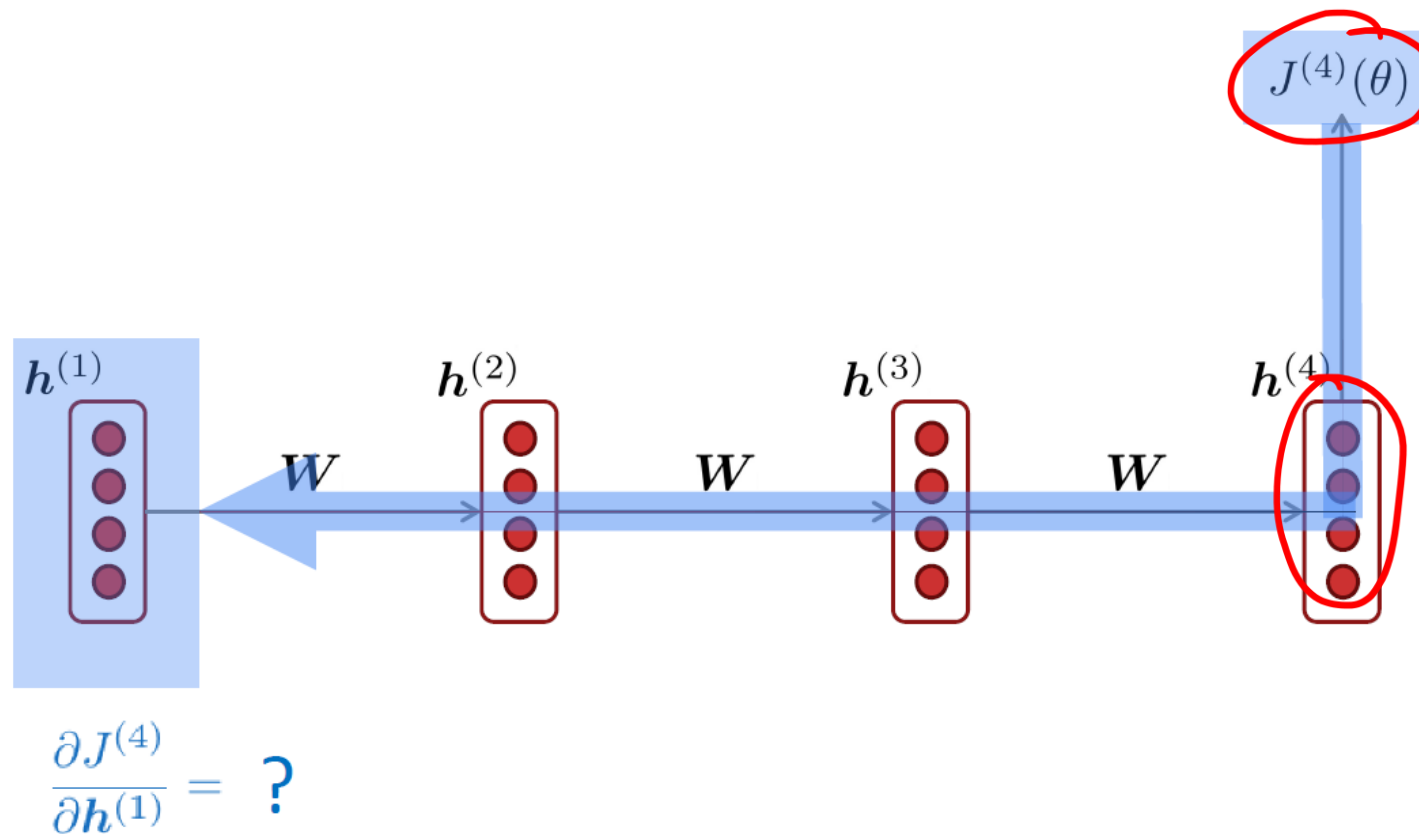


# Problems with RNNs

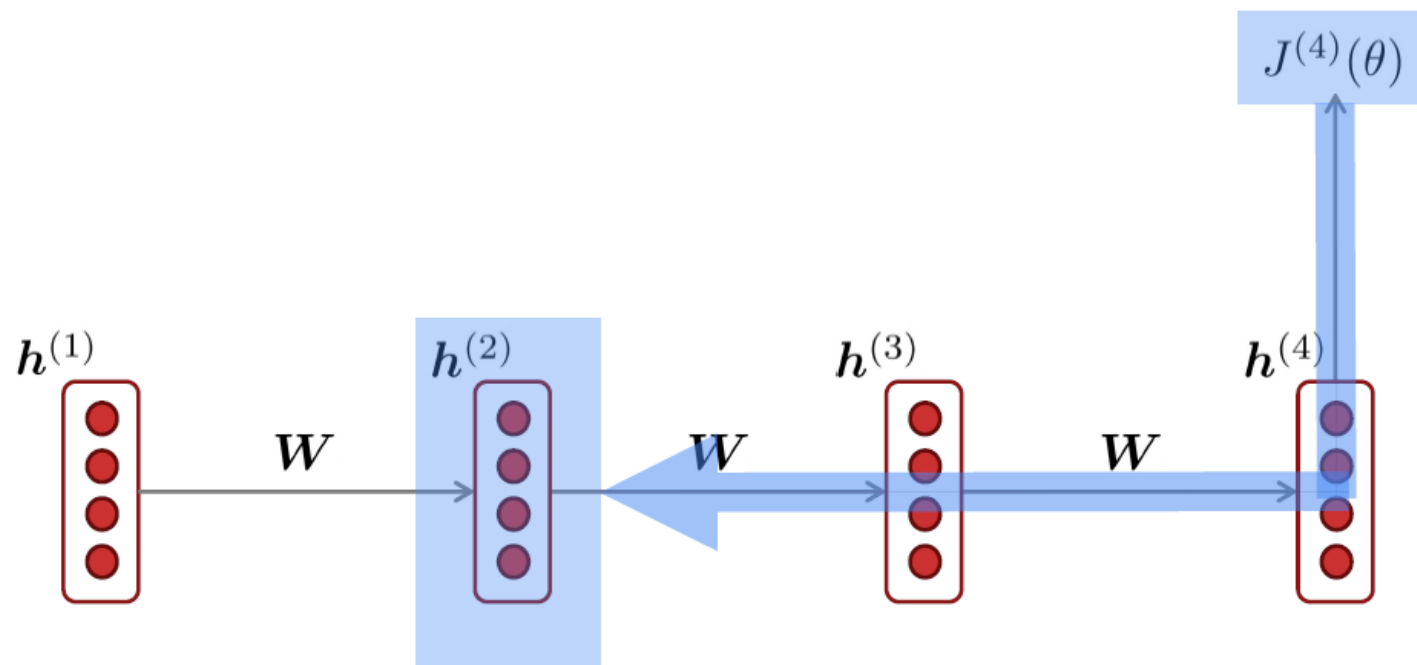
# Vanishing and Exploding Gradients



# Vanishing Gradient Intuition



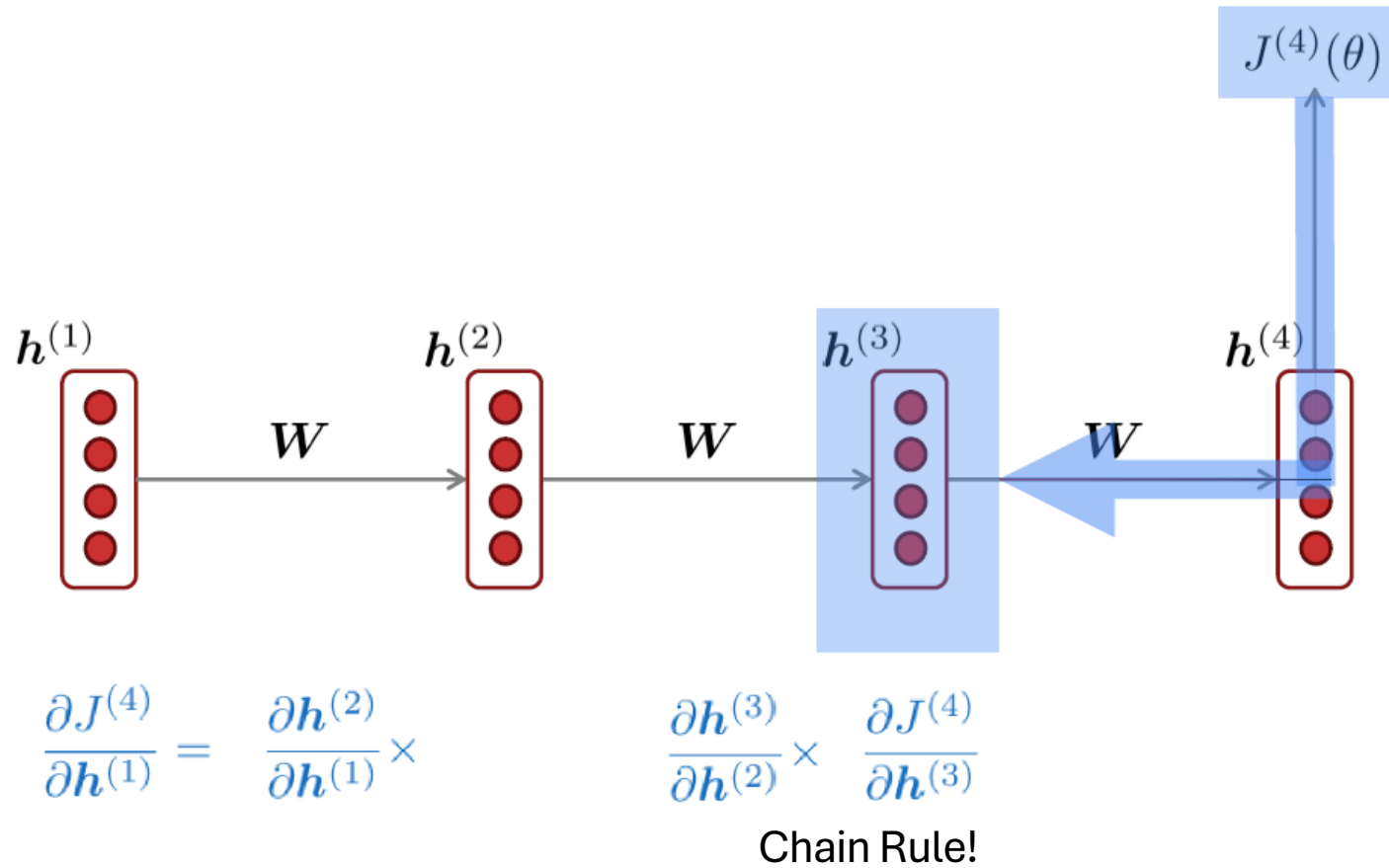
# Vanishing Gradient Intuition



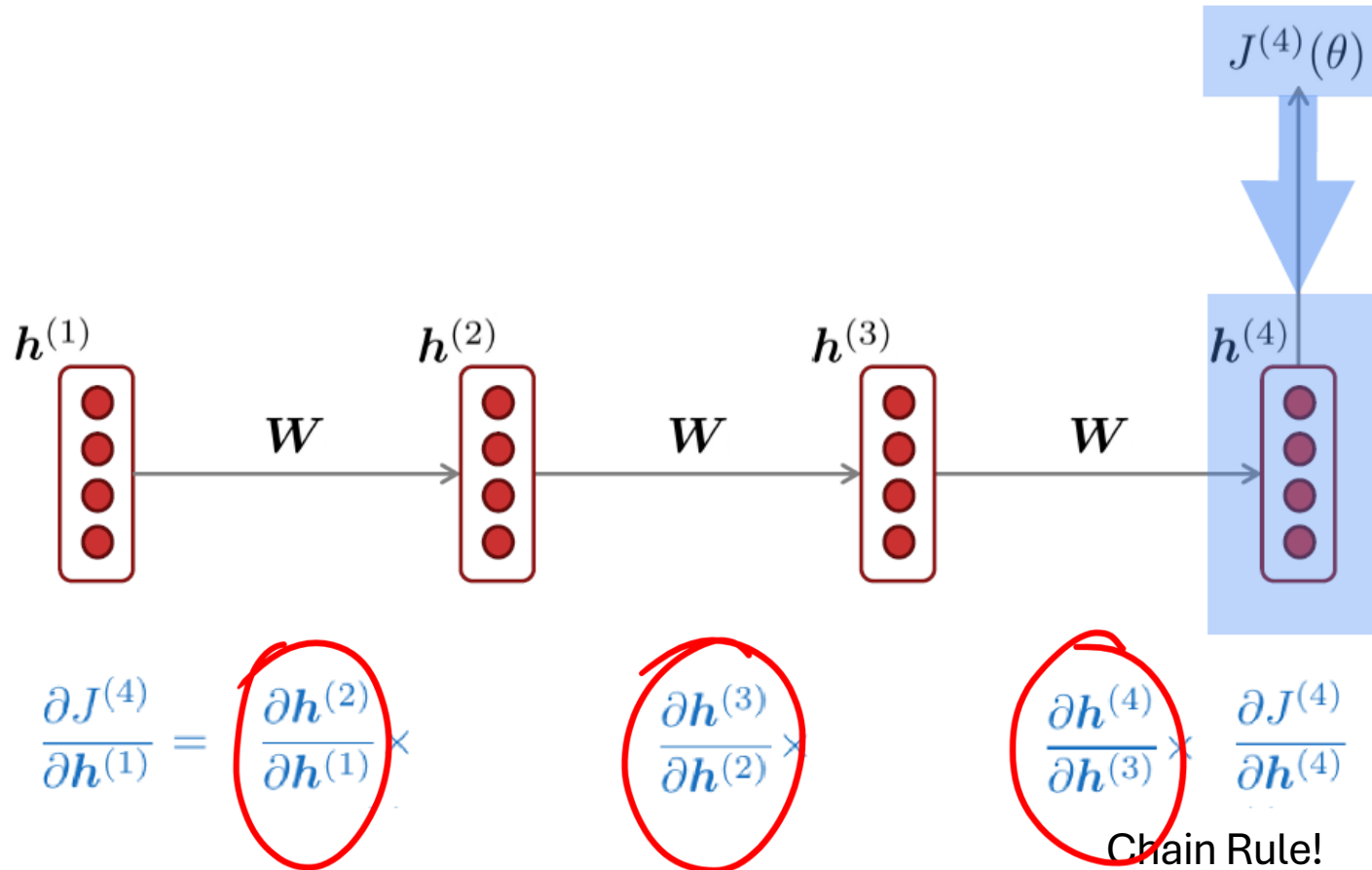
$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial J^{(4)}}{\partial h^{(2)}}$$

Chain Rule!

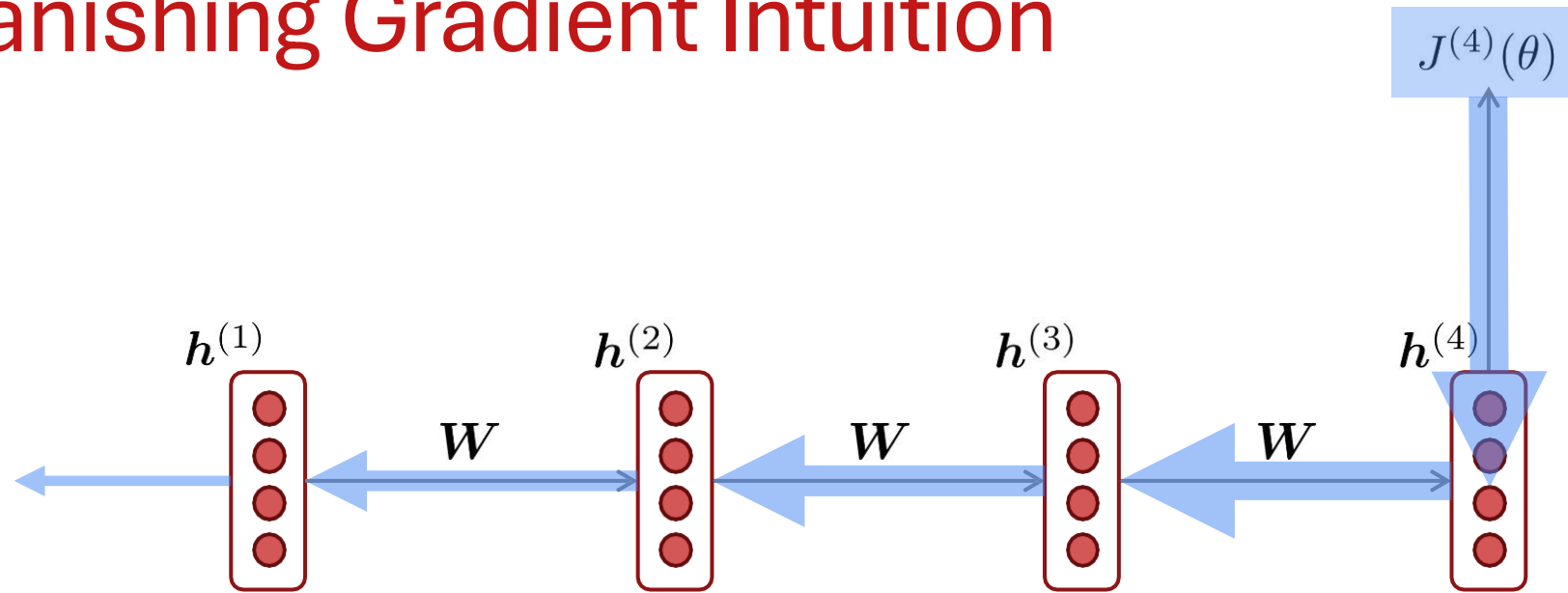
# Vanishing Gradient Intuition



# Vanishing Gradient Intuition



# Vanishing Gradient Intuition



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \left[ \frac{\partial h^{(2)}}{\partial h^{(1)}} \right] \times \left[ \frac{\partial h^{(3)}}{\partial h^{(2)}} \right] \times \left[ \frac{\partial h^{(4)}}{\partial h^{(3)}} \right] \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

What happens if these are small?

**Vanishing gradient problem:** When these are small, the gradient signal gets smaller and smaller as it backpropagates further

# Effect of Vanishing Gradient on RNN-LM

- **LM task:** *When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her \_\_\_\_\_*
- To learn from this training example, the RNN-LM needs to **model the dependency** between “*tickets*” on the 7<sup>th</sup> step and the target word “*tickets*” at the end.
- But if the gradient is small, the model **can’t learn this dependency**
  - So, the model is **unable to predict similar long-distance dependencies** at test time



# LSTMs & GRUs

# LSTM

We have a sequence of inputs  $x^{(t)}$ , and we will compute a sequence of hidden states  $h^{(t)}$  and cell states  $c^{(t)}$ . On 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

$$\mathbf{f}^{(t)} = \sigma(\mathbf{W}_f \mathbf{h}^{(t-1)} + \mathbf{U}_f \mathbf{x}^{(t)} + \mathbf{b}_f)$$

$$\mathbf{i}^{(t)} = \sigma \left( \mathbf{W}_i \mathbf{h}^{(t-1)} + \mathbf{U}_i \mathbf{x}^{(t)} + \mathbf{b}_i \right)$$

$$\mathbf{o}^{(t)} = \sigma \left( \mathbf{W}_o \mathbf{h}^{(t-1)} + \mathbf{U}_o \mathbf{x}^{(t)} + \mathbf{b}_o \right)$$

$$\tilde{\mathbf{c}}^{(t)} = \tanh \left( \mathbf{W}_c \mathbf{h}^{(t-1)} + \mathbf{U}_c \mathbf{x}^{(t)} + \mathbf{b}_c \right)$$

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

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

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

→ All these are vectors of same length  $n$

$$C^{t-1} = \begin{pmatrix} .8 & .1 & .5 \\ 1 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} .8 & 0 & 0 \end{pmatrix}$$

# Gated Recurrent Units (GRUs)

- Proposed by Cho et al. in 2014 as a **simpler alternative to the LSTM**.
- On each timestep  $t$ , we have input  $x^{(t)}$  and hidden state  $h^{(t)}$  (**no cell state**).

**Update gate:** controls what parts of hidden state are updated vs preserved

**Reset gate:** controls what parts of previous hidden state are used to compute new content

**New hidden state content:** reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

**Hidden state:** update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

**How does this solve vanishing gradient?**

Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)

$$\mathbf{u}^{(t)} = \sigma \left( \mathbf{W}_u \mathbf{h}^{(t-1)} + \mathbf{U}_u \mathbf{x}^{(t)} + \mathbf{b}_u \right)$$

$$\mathbf{r}^{(t)} = \sigma \left( \mathbf{W}_r \mathbf{h}^{(t-1)} + \mathbf{U}_r \mathbf{x}^{(t)} + \mathbf{b}_r \right)$$

$$\tilde{\mathbf{h}}^{(t)} = \tanh \left( \mathbf{W}_h (\mathbf{r}^{(t)} \circ \mathbf{h}^{(t-1)}) + \mathbf{U}_h \mathbf{x}^{(t)} + \mathbf{b}_h \right)$$

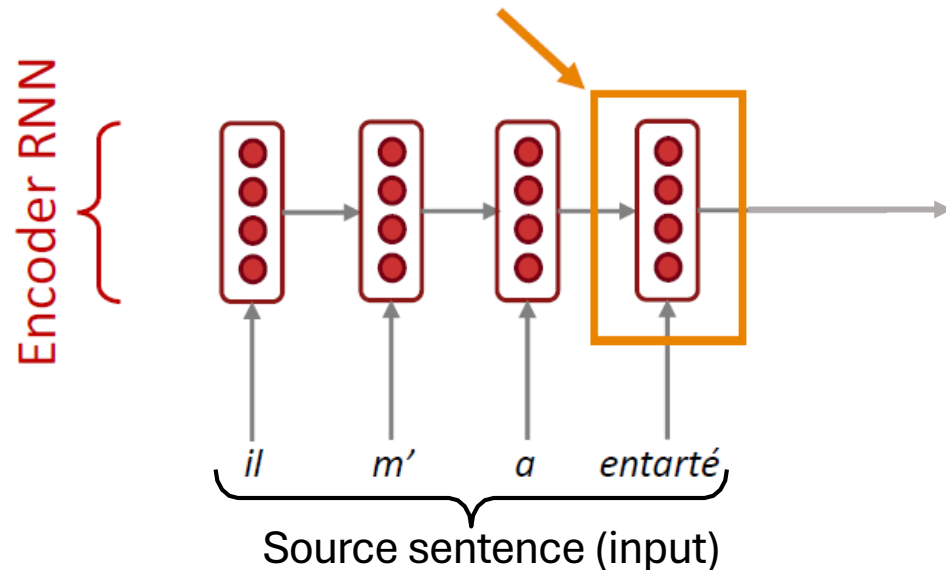
$$\mathbf{h}^{(t)} = (1 - \mathbf{u}^{(t)}) \circ \mathbf{h}^{(t-1)} + \mathbf{u}^{(t)} \circ \tilde{\mathbf{h}}^{(t)}$$

# Sequence-to-Sequence Modeling

# Neural Machine Translation (NMT)

## The Sequence-to-Sequence Model

Encoding of the source sentence.  
Provides initial hidden state  
for Decoder RNN.

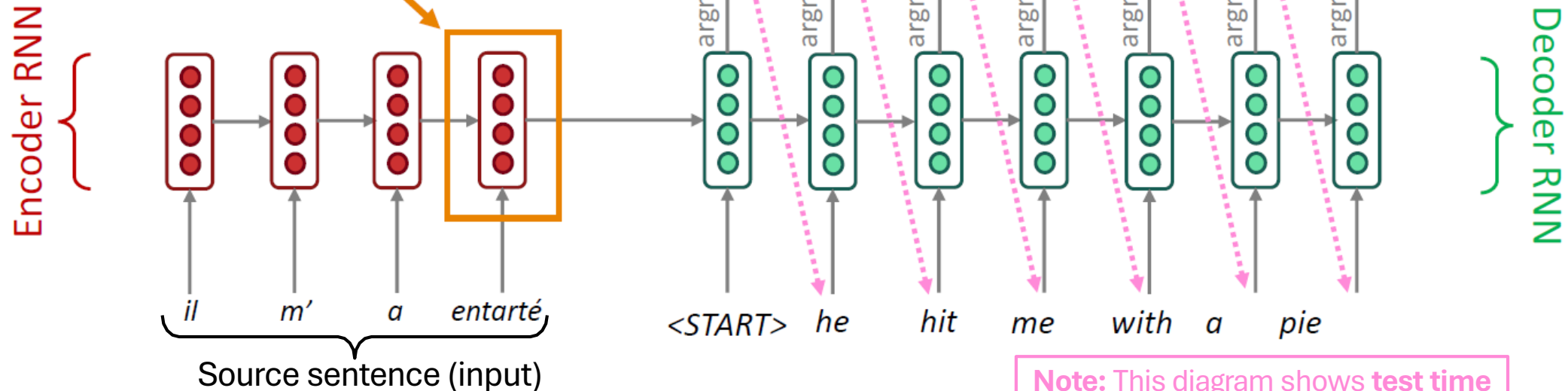


**Encoder RNN produces an encoding of the source sentence.**

# Neural Machine Translation (NMT)

## The Sequence-to-Sequence Model

Encoding of the source sentence.  
Provides initial hidden state  
for Decoder RNN.



**Encoder RNN produces an *encoding* of the source sentence.**

**Decoder RNN is a Language Model that generates target sentence, conditioned on *encoding*.**

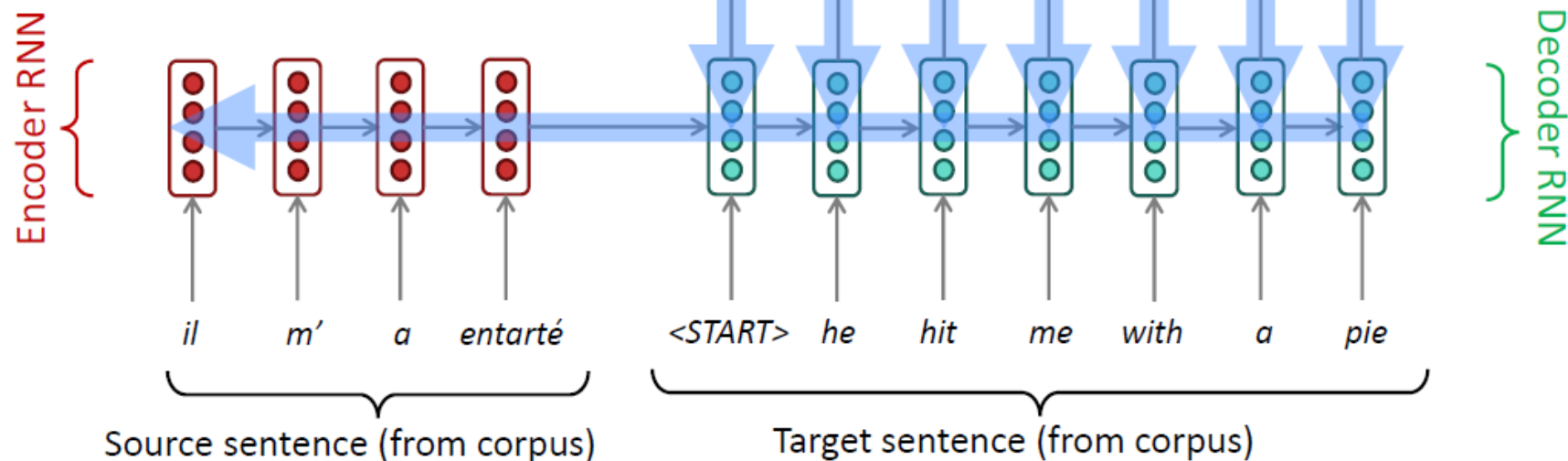
**Note:** This diagram shows **test time** behavior: decoder output is fed in as next step's input

# Training an NMT System

Seq2seq is optimized as a **single system**. Backpropagation operates “*end-to-end*”.

$$J = \frac{1}{T} \sum_{t=1}^T J_t = \boxed{J_1} + J_2 + J_3 + \boxed{J_4} + J_5 + J_6 + \boxed{J_7}$$

= negative log prob of “he”      = negative log prob of “with”      = negative log prob of <END>



# Issues With RNN

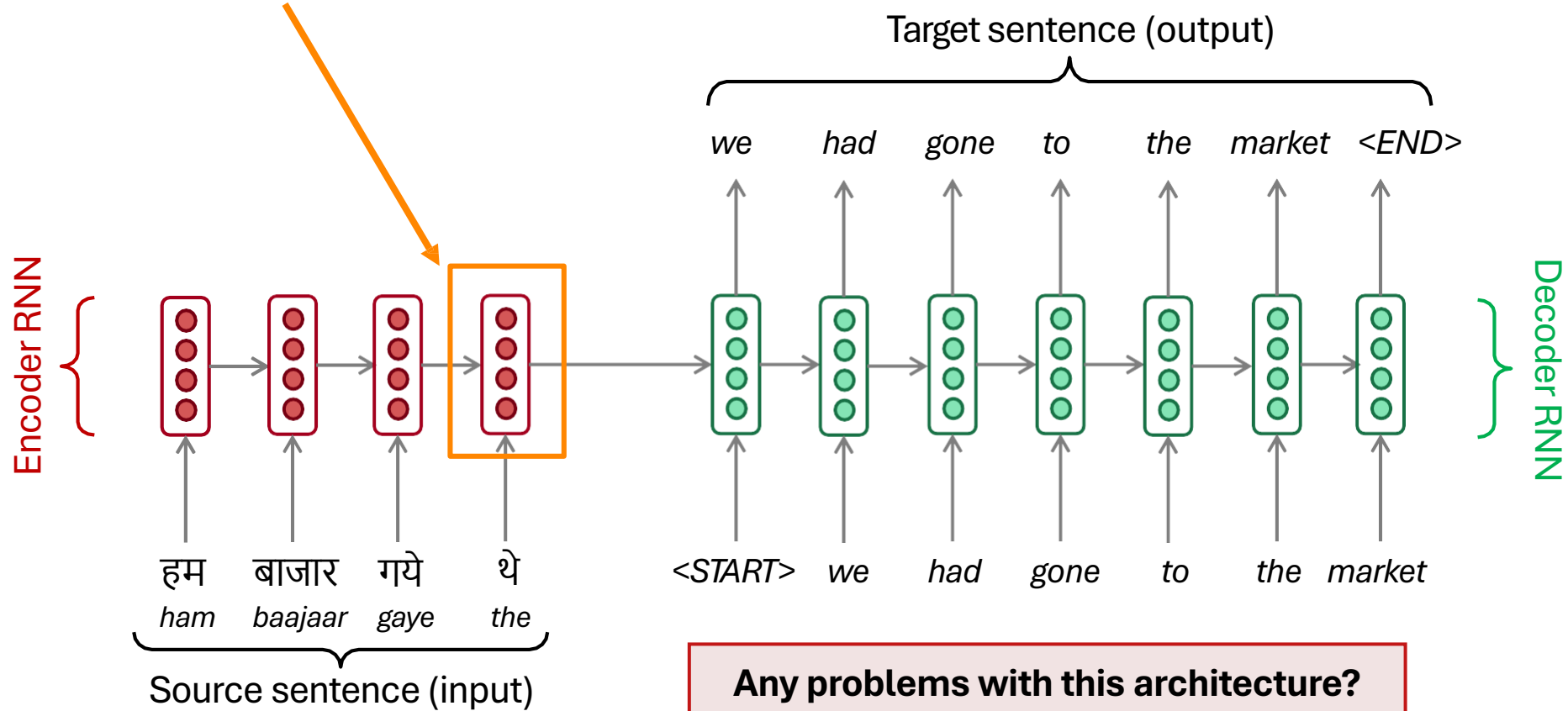
- Linear interaction distance
- Bottleneck problem
- Lack of parallelizability

# ATTENTION



# Sequence-to-Sequence: The Bottleneck Problem

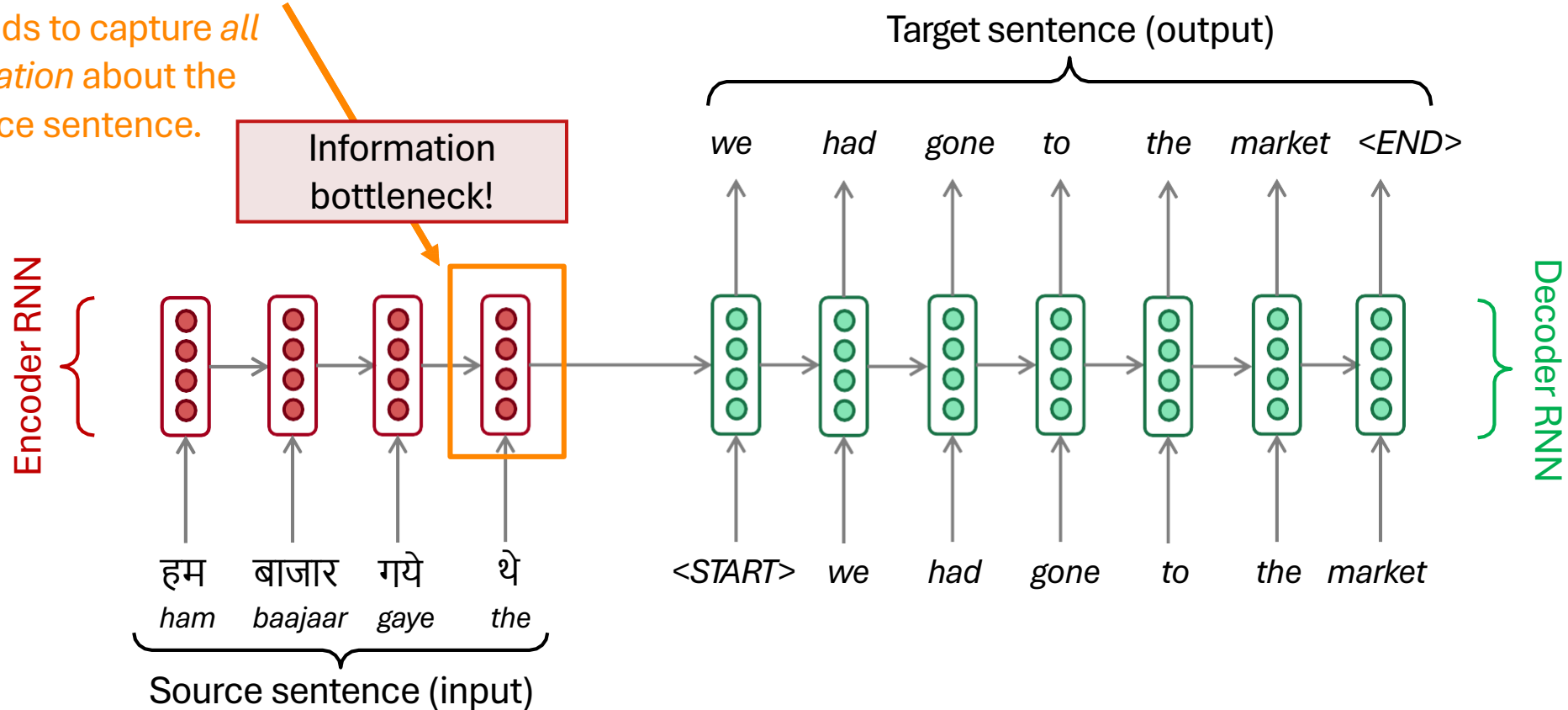
Encoding of the source sentence



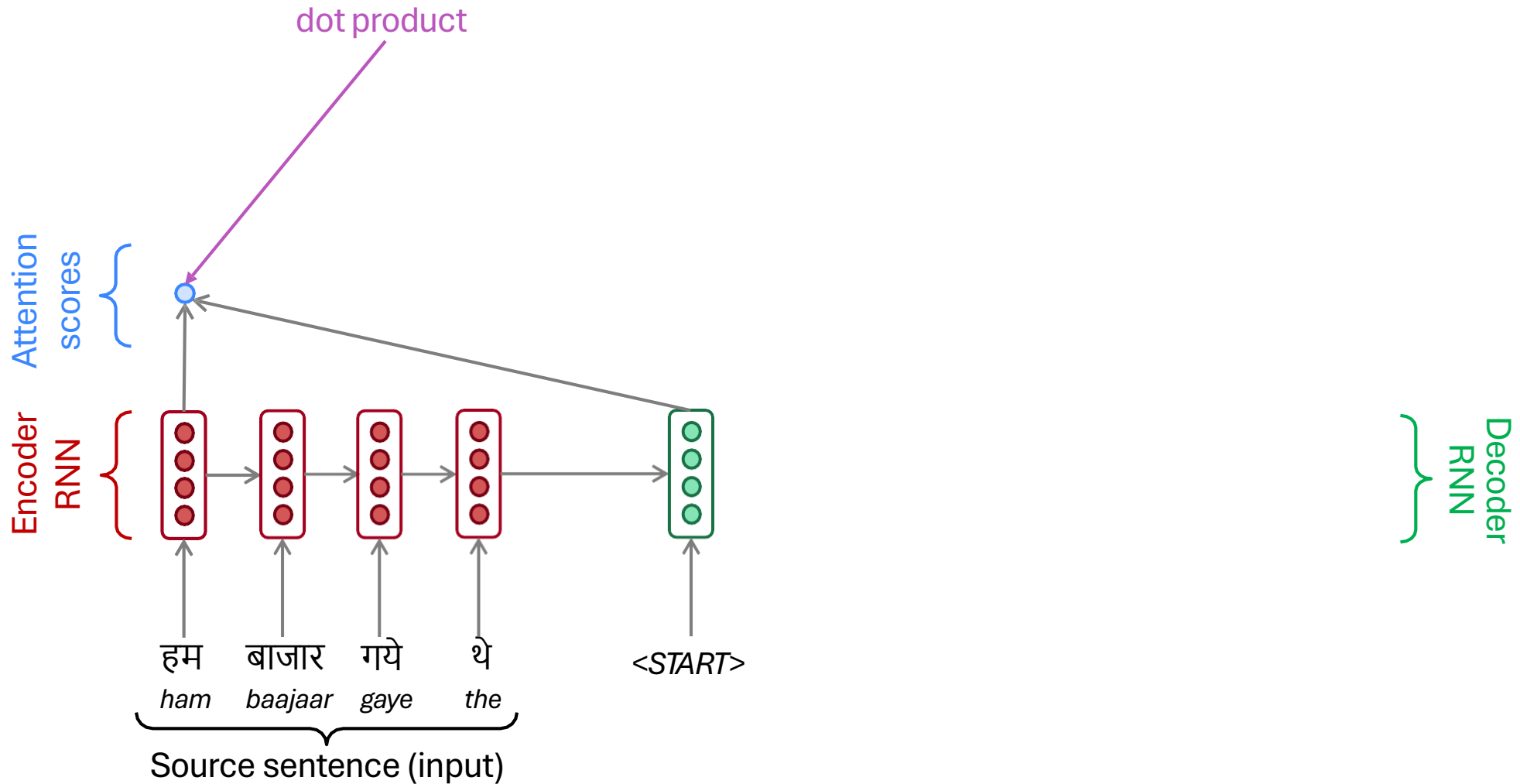
# Sequence-to-Sequence: The Bottleneck Problem

Encoding of the source sentence

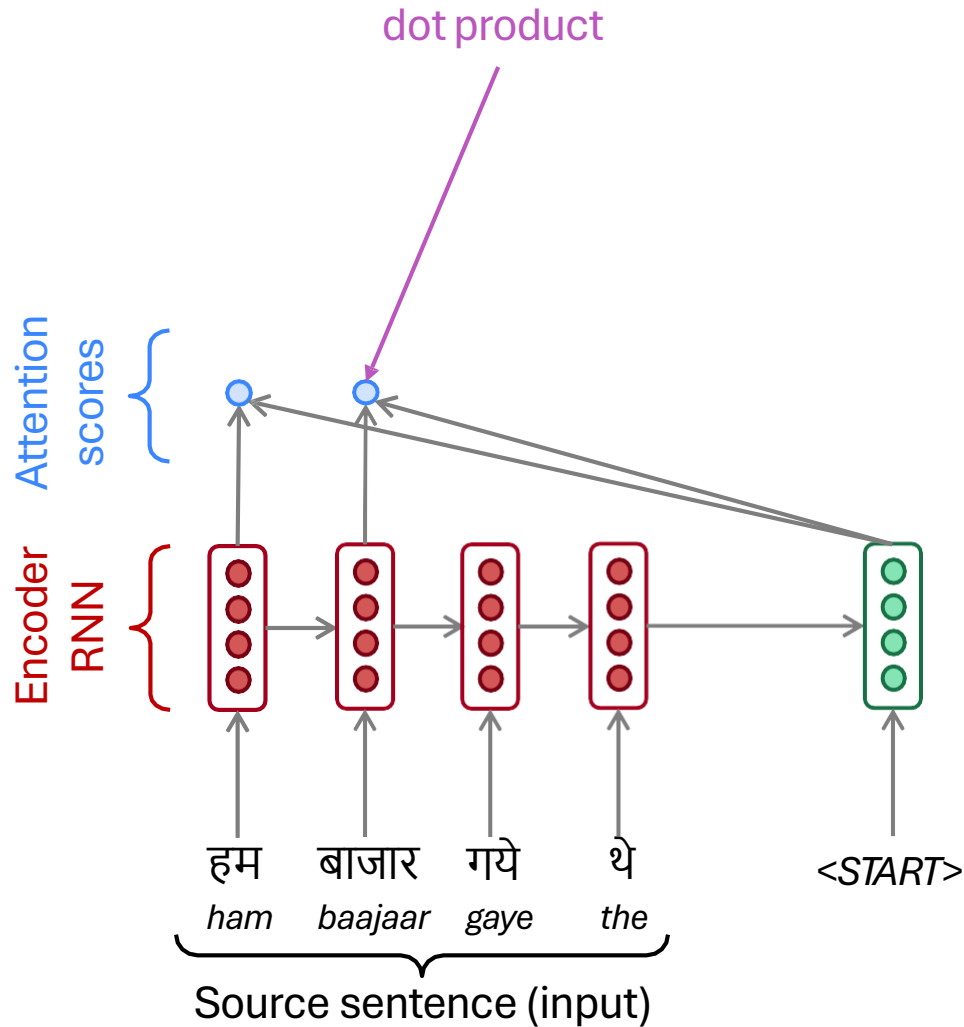
This needs to capture *all* information about the source sentence.



# Sequence-to-Sequence With Attention

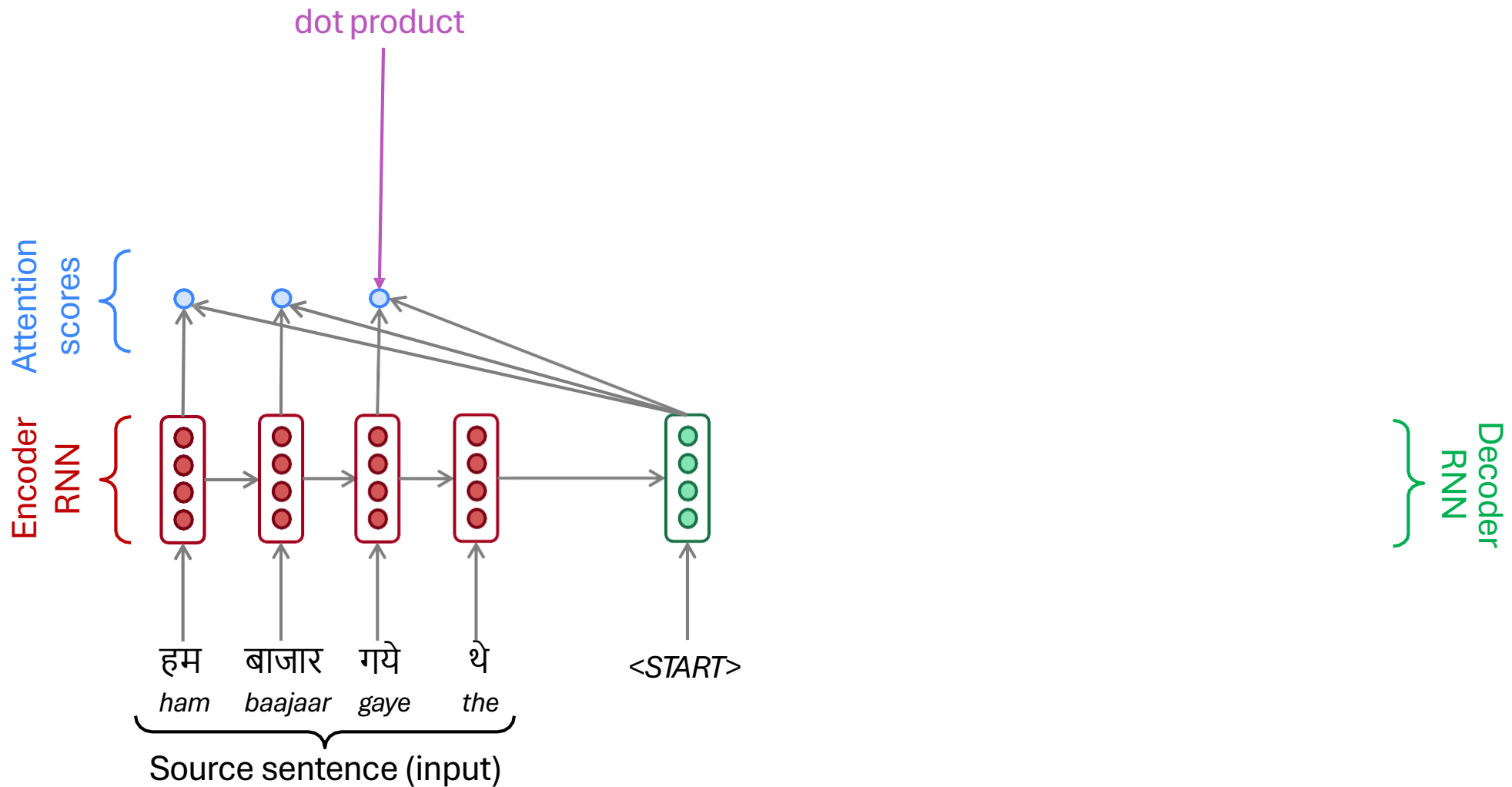


# Sequence-to-Sequence With Attention

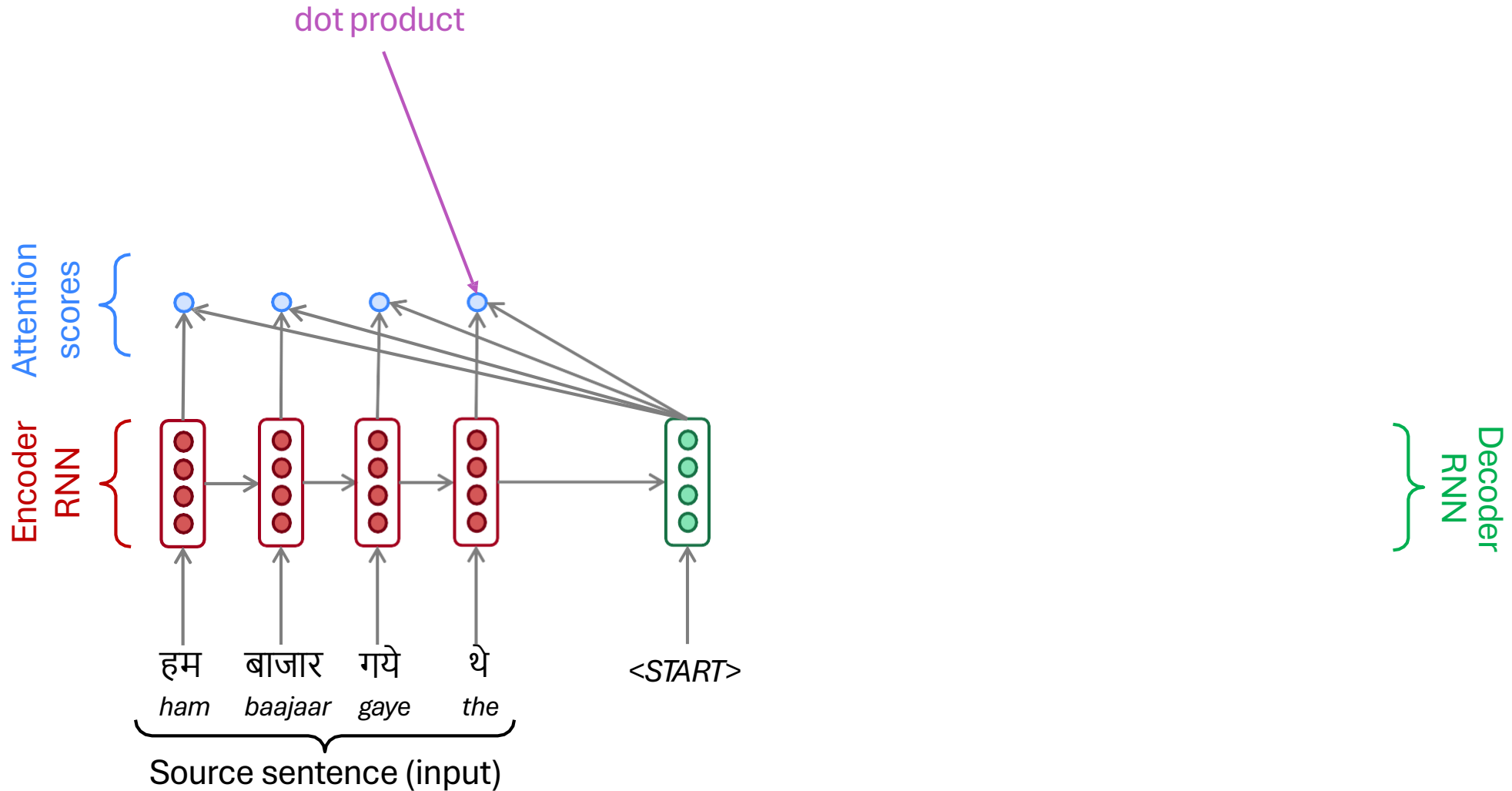


Decoder  
RNN

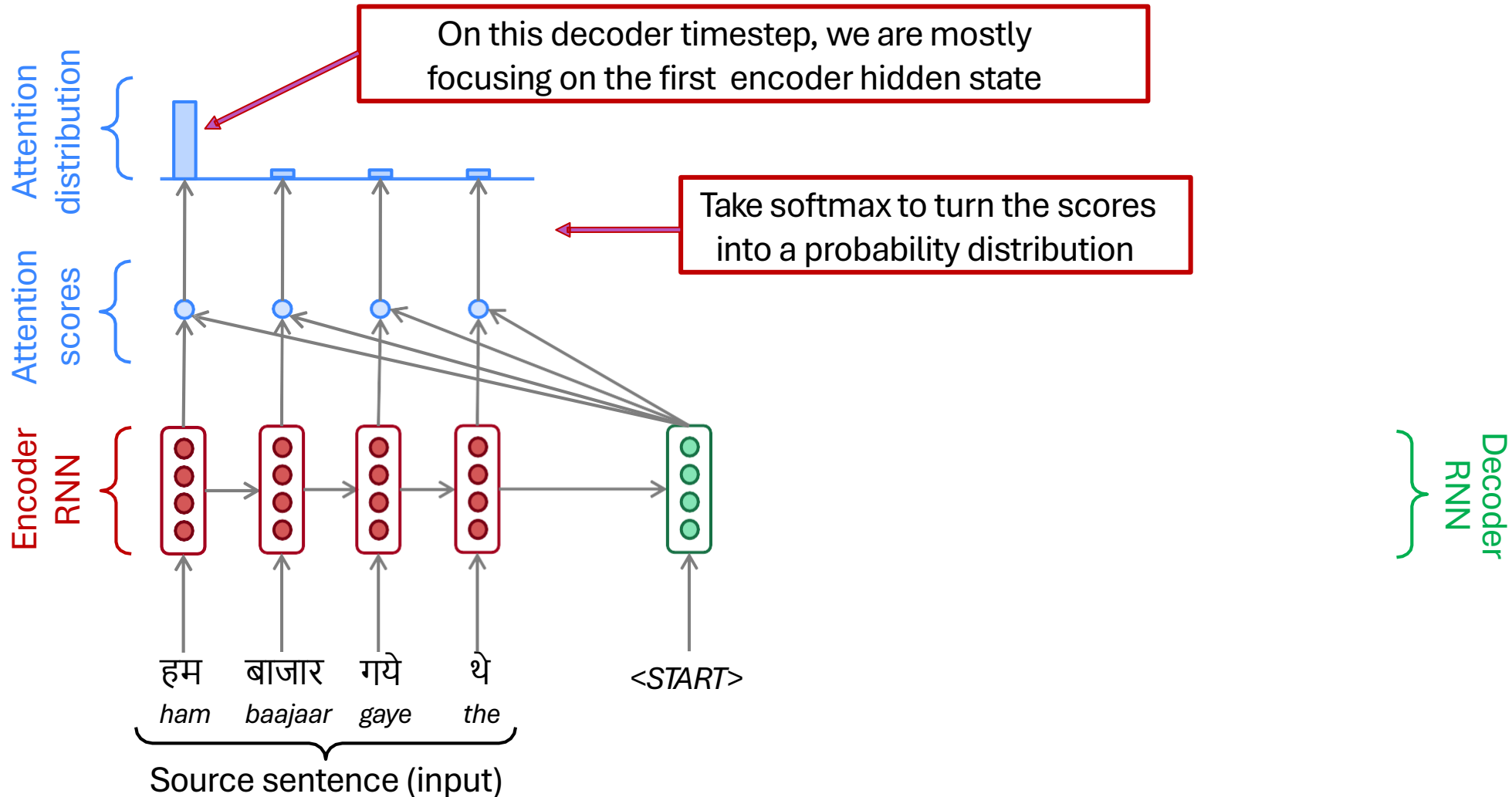
# Sequence-to-Sequence With Attention



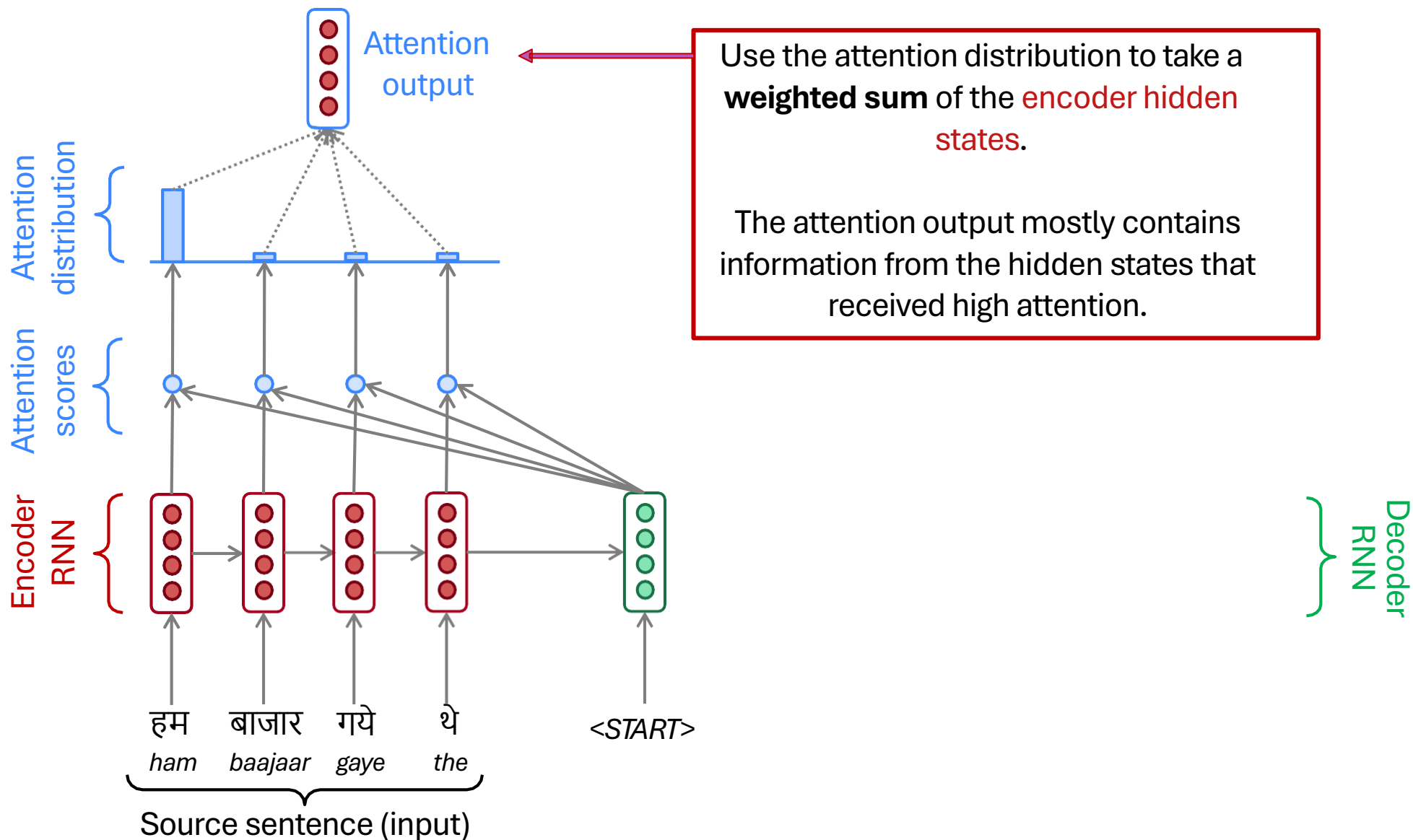
# Sequence-to-Sequence With Attention



# Sequence-to-Sequence With Attention

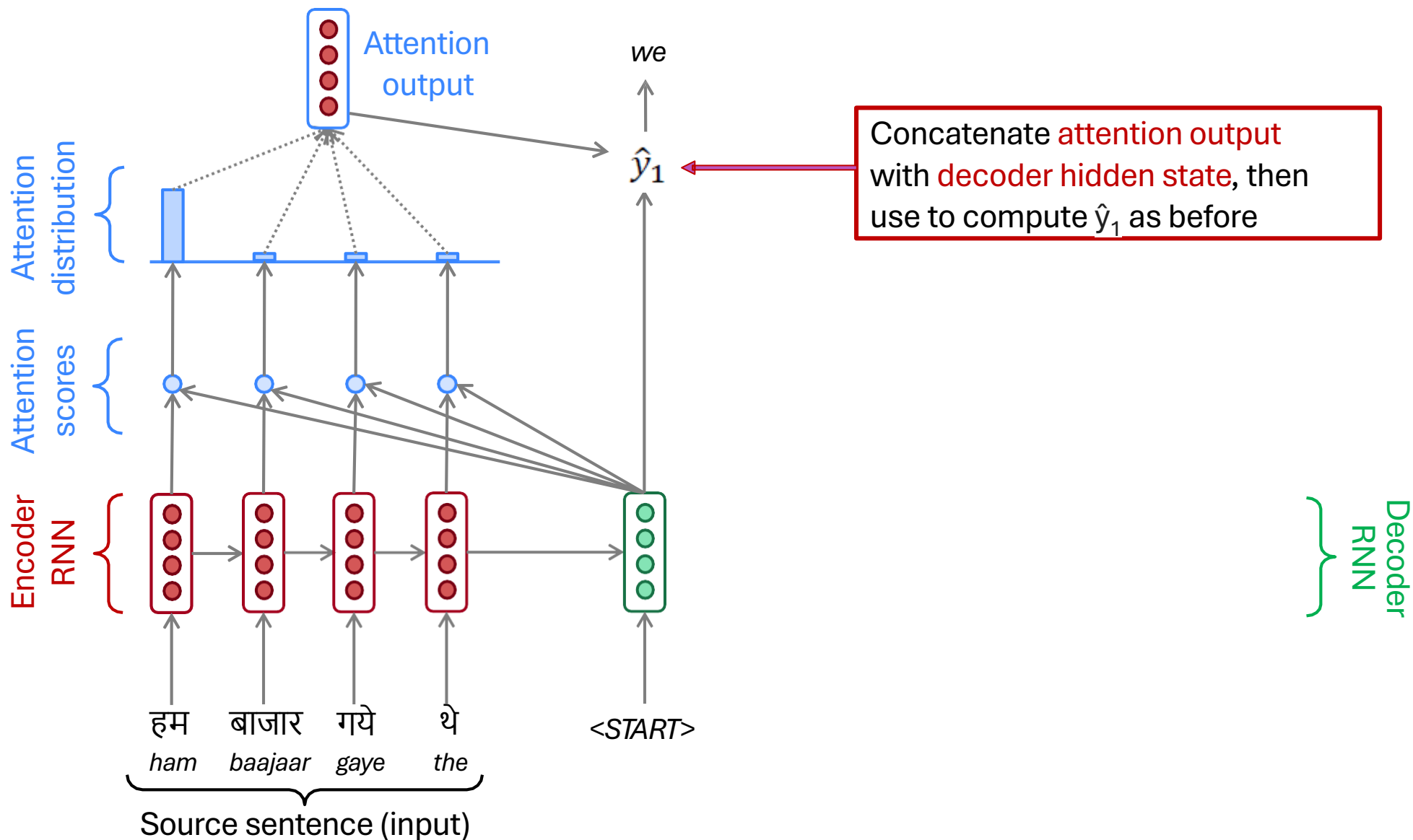


# Sequence-to-Sequence With Attention

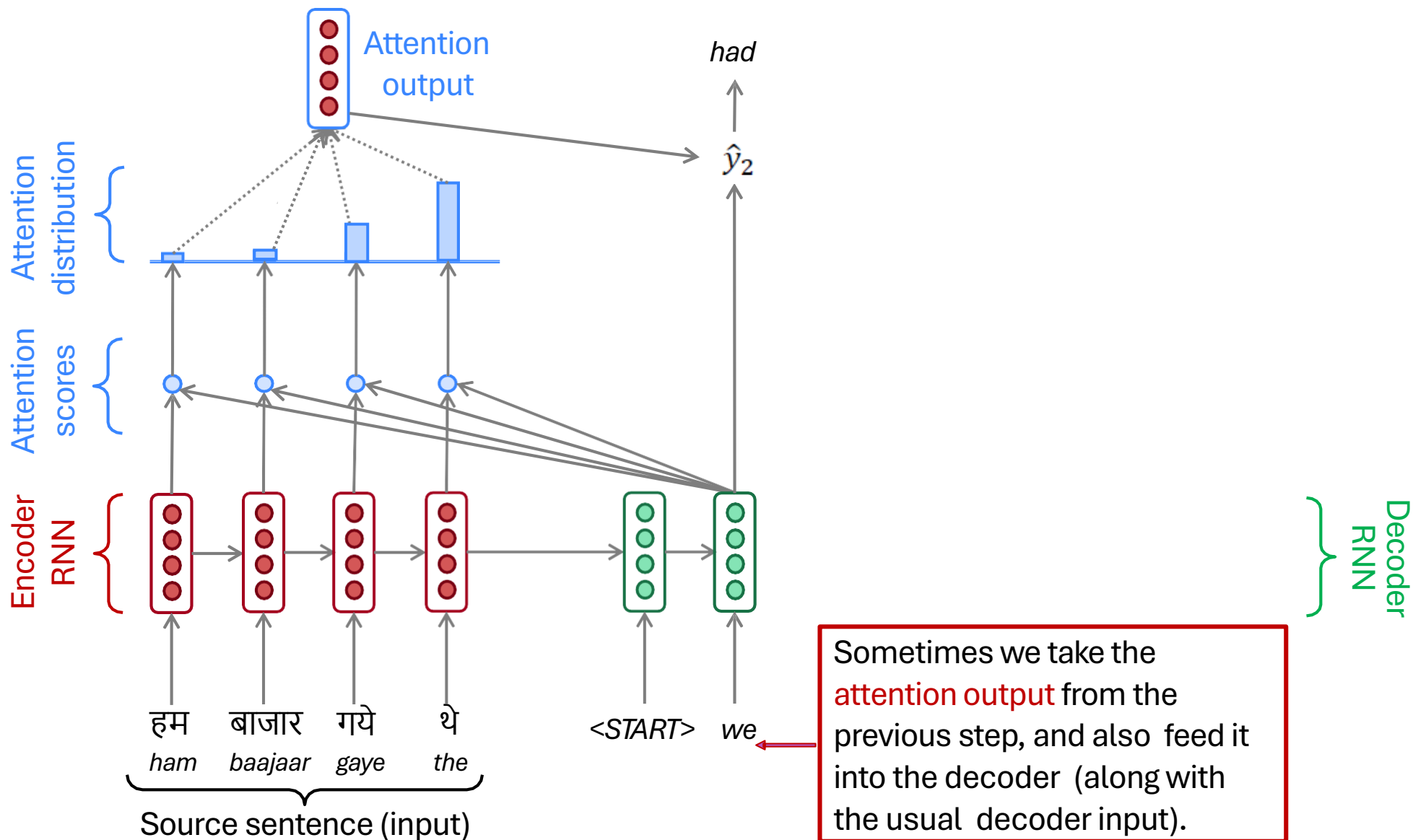




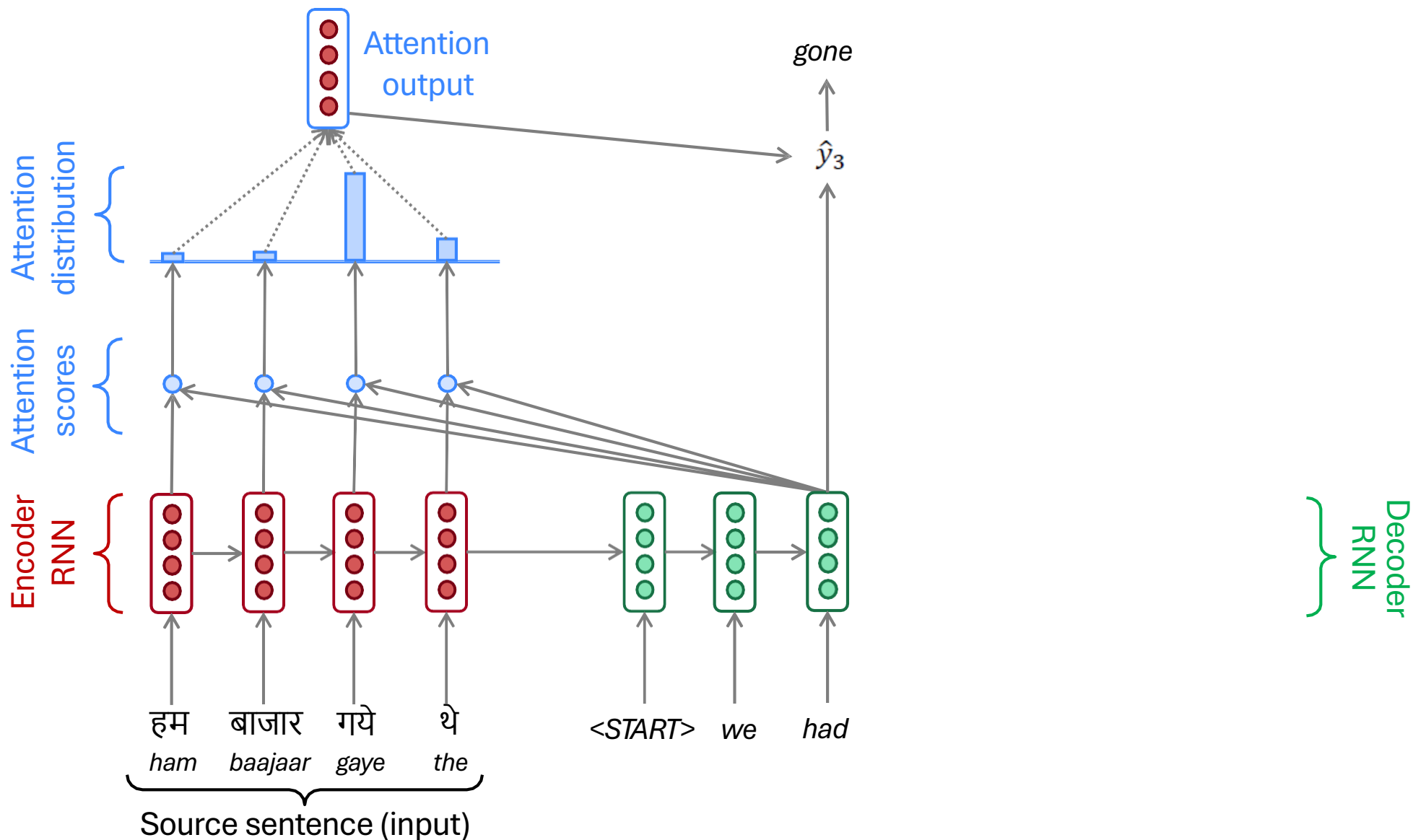
# Sequence-to-Sequence With Attention



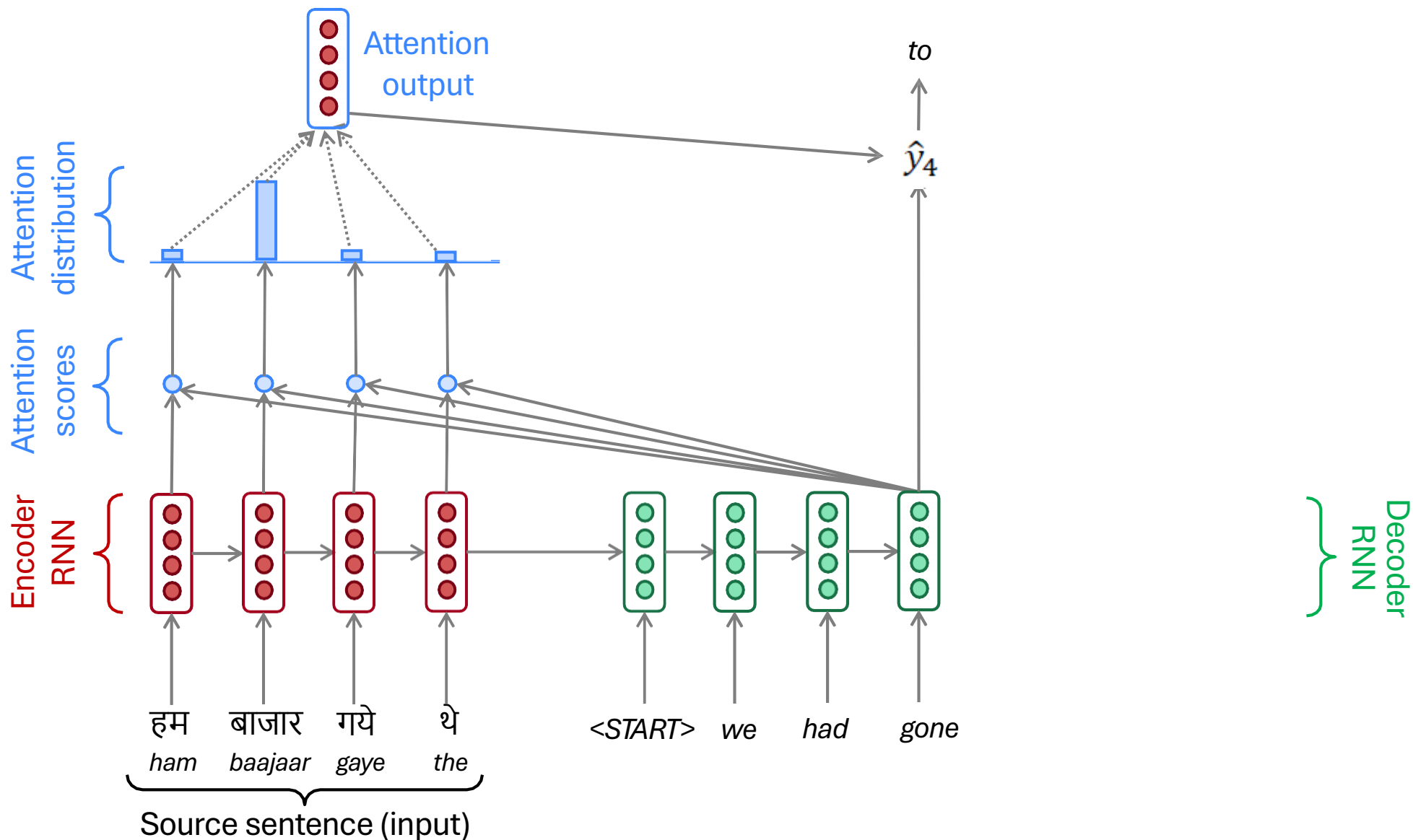
# Sequence-to-Sequence With Attention



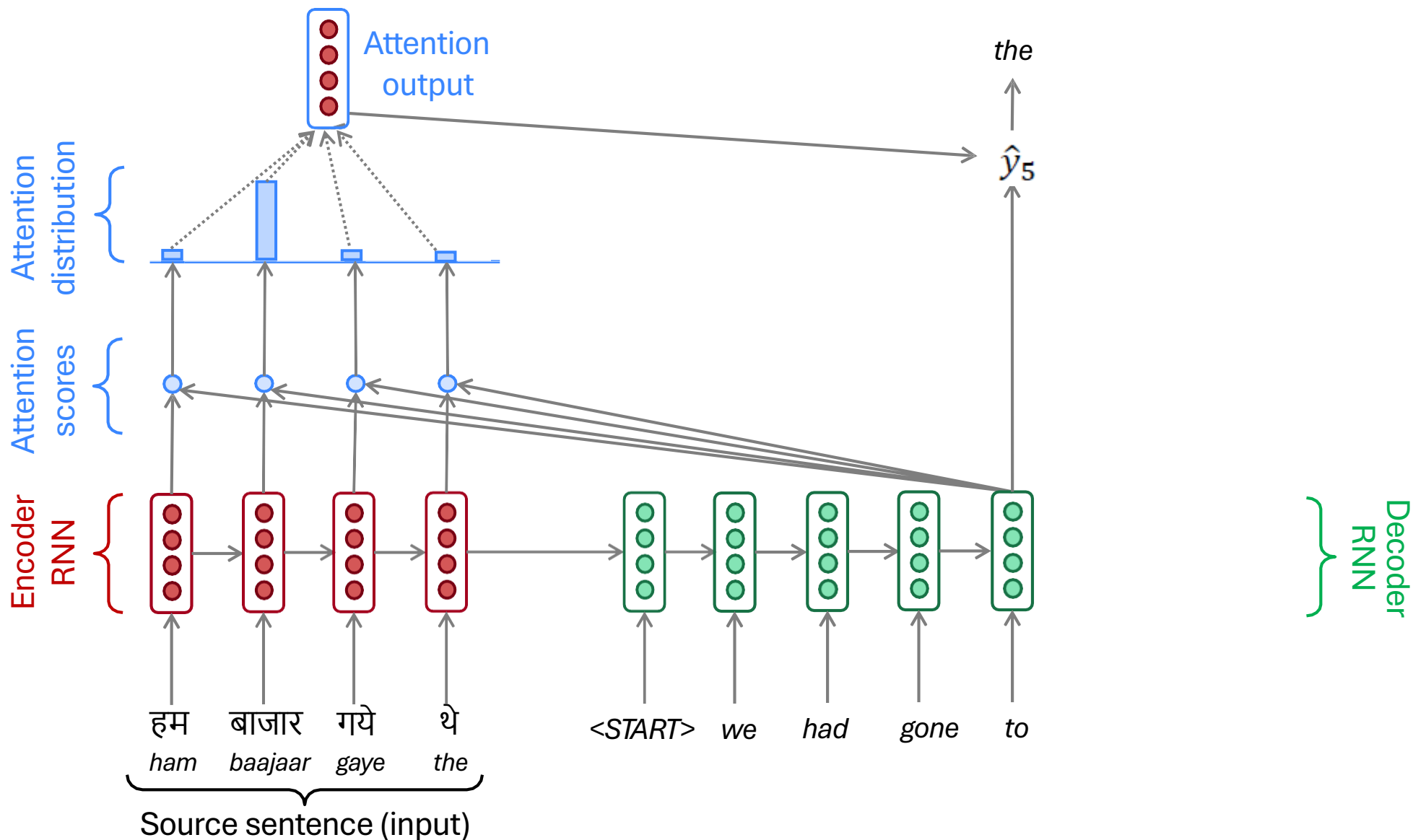
# Sequence-to-Sequence With Attention



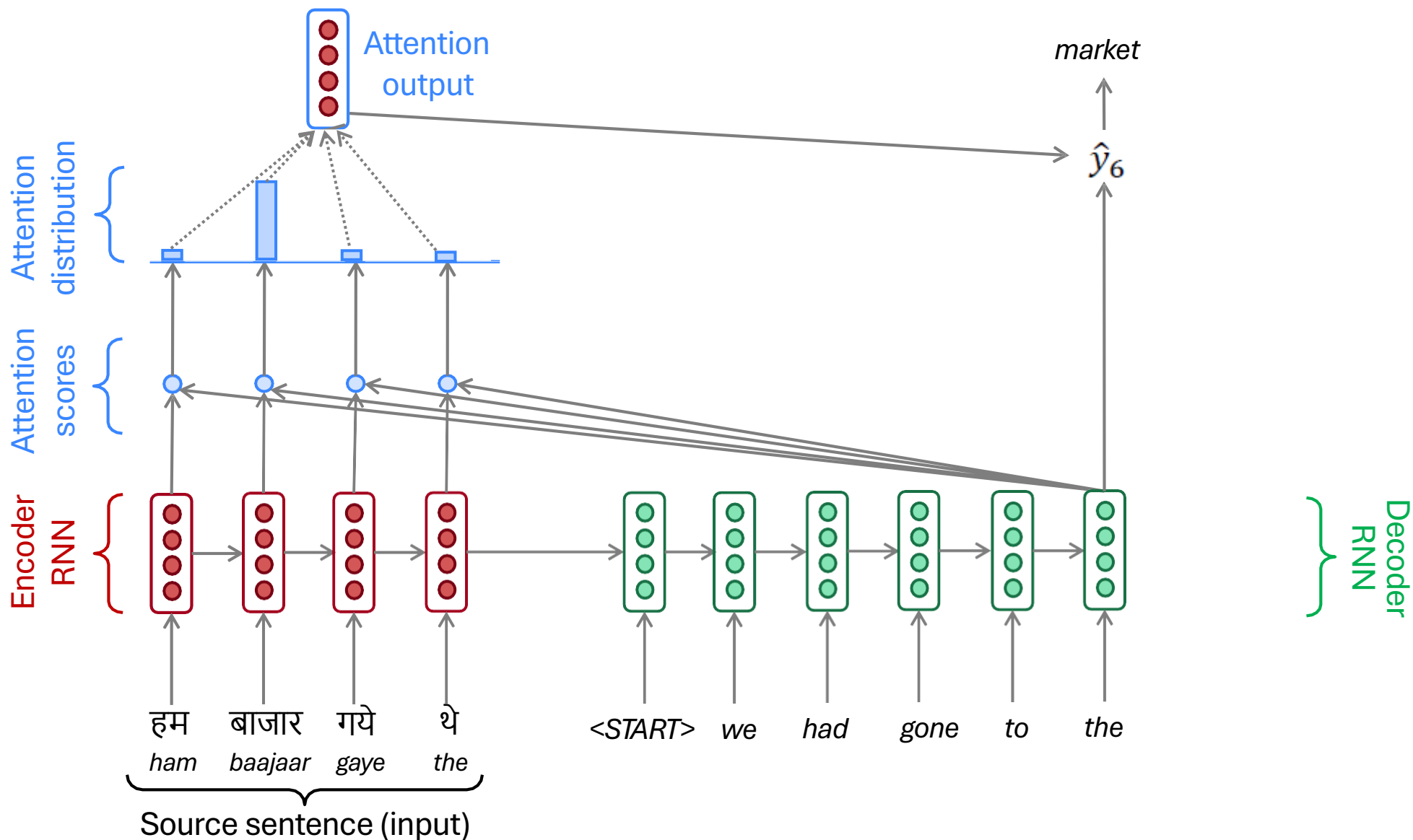
# Sequence-to-Sequence With Attention



# Sequence-to-Sequence With Attention



# Sequence-to-Sequence With Attention



# Attention is Great

- Attention significantly **improves NMT performance**
  - It's very useful to allow decoder to focus on certain parts of the source
- Attention **solves the 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 (soft) **alignment for free!**
  - This is cool because we never explicitly trained an alignment system
  - The network just learned alignment by itself

	he	hit	me	with	a	pie
il						
a						
m'						
entarté						

# Attention is a *General* Deep Learning Technique

- We've seen that attention is a great way to improve the sequence-to-sequence model for Machine Translation.
- However: You can use attention in *many architectures* (not just seq2seq) and *many tasks* (not just MT)
- *More general definition of attention*:
  - Given a set of vector *values*, and a vector *query*, *attention* is a technique to compute a weighted sum of the values, dependent on the query.
- We sometimes say that the *query attends to the values*.
- For example, in the seq2seq + attention model, each decoder hidden state (query) *attends to* all the encoder hidden states (values).
- **Intuition**:
  - The weighted sum is a *selective summary* of the information contained in the values, where the query determines which values to focus on.
  - Attention is a way to obtain a *fixed-size representation of an arbitrary set of representations* (the values), dependent on some other representation (the query).



# Variants of Attention

- Original formulation:  $a(\mathbf{q}, \mathbf{k}) = w_2^T \tanh(W_1[\mathbf{q}; \mathbf{k}])$
- Bilinear product:  $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T W \mathbf{k}$  Luong et al., 2015
- Dot product:  $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k}$  Luong et al., 2015
- Scaled dot product:  $a(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^T \mathbf{k}}{\sqrt{|\mathbf{k}|}}$  Vaswani et al., 2017

**More information:**

“Deep Learning for NLP Best Practices”, Ruder, 2017. <http://ruder.io/deep-learning-nlp-best-practices/index.html#attention>

“Massive Exploration of Neural Machine Translation Architectures”, Britz et al, 2017, <https://arxiv.org/pdf/1703.03906.pdf>

# Self-Attention