

CS 2731 / ISSP 2230

# Introduction to Natural Language Processing

Session 12: RNNs part 2, encoder-decoder

---

Michael Miller Yoder

February 19, 2024



School of Computing and Information

# Course logistics

- [Homework 2](#) grades will be out this week
- Proposal and literature review is **due this Thu Feb 22, 11:59pm**
  - Instructions are on the [project webpage](#)
  - Submit on Canvas
  - One submission per group
- Homework 3 will be released tomorrow, Tue Feb 20

# How to do a literature review

- Look for NLP papers related to your topic in [ACL Anthology](#), [Semantic Scholar](#), and [Google Scholar](#)
- For each paper, note:
  - What they cite in their related work sections (find those papers, iterate)
  - Data
  - Methods
  - Findings
- For at least 4 papers, organize them into themes of approaches, datasets, findings
- Ok: X paper did this, Y paper did this, Z paper did that
- Good: X and Y papers did this, while Z improved with that
- Best: X and Y papers did this, Z improved, nobody has yet to do...

# Midterm course evaluation (OMETs)

- CS 2731:  
<https://go.blueja.io/3ydJumSMRk23LZULkOgifQ>
- ISSP 2230:  
<https://go.blueja.io/U9GquEwBj0WkSuaDGWNx5g>
- All types of feedback are welcome  
(critical and positive)
- **Completely anonymous, will not affect grades**
- Let me know what's working and what to  
improve on while the course is still running!
- Please be as specific as possible
- Available until **tonight, Mon Feb 19 at 11:59pm**



# Lecture overview: RNNs part 2, encoder-decoder

- RNN language modeling
- LSTMs
- RNNs for other NLP tasks
- Encoder-decoder model
- Attention

# RNNs for language modeling

---

# RNN refresher

With a neighbor, talk about the following questions:

1. What are the 2 inputs to the hidden states in a simple 2-layer RNN?
2. What do RNNs allow us to do in NLP that we can't do with feedforward neural networks?

# An RNN Language Model

## output distribution

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

## hidden states

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

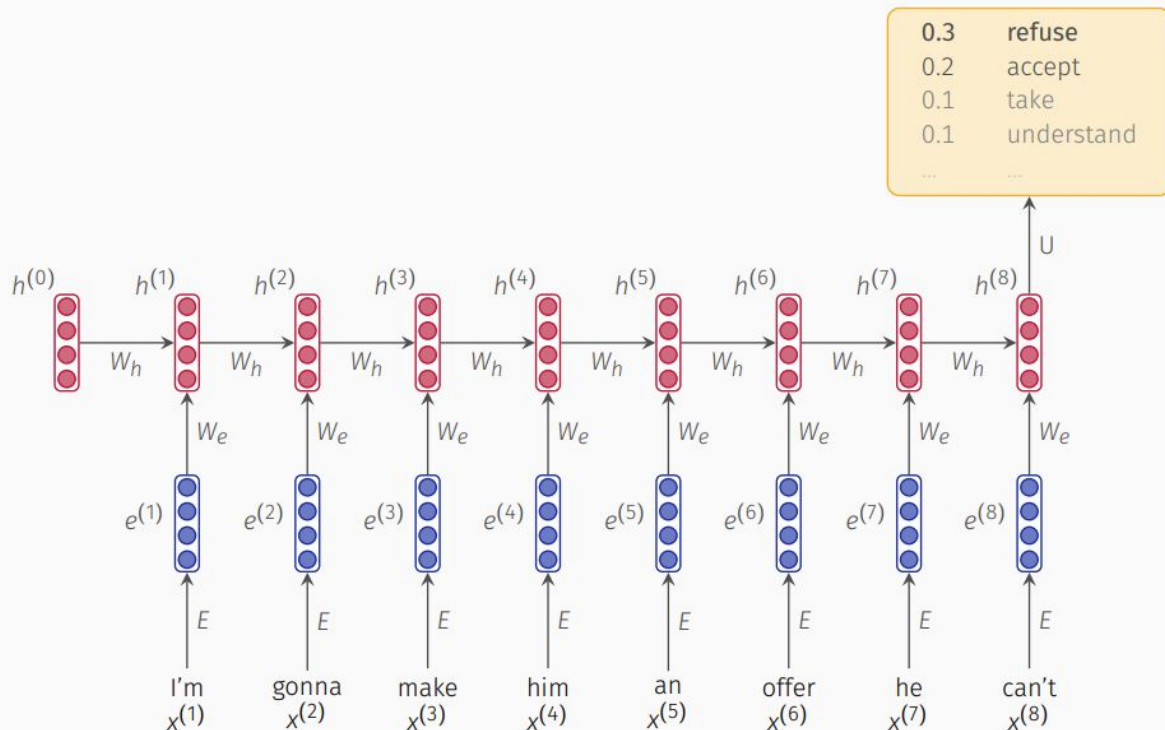
$h^{(0)}$  is the initial hidden state

## word embeddings

$$e^{(t)} = Ex^{(t)}$$

## one-hot vectors

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





# Training an RNN Language Model

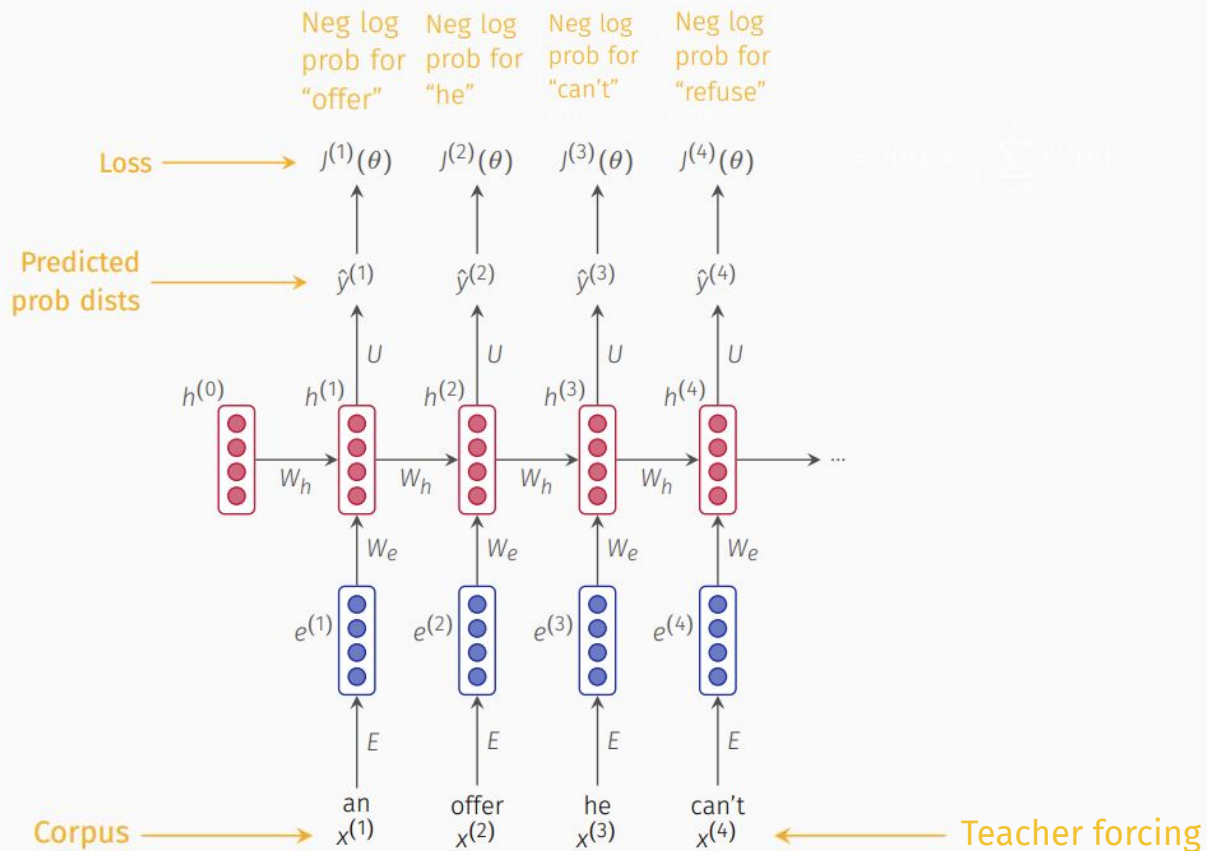
- Get a big corpus of text, which is a sequence of words  $x^{(1)}, \dots, x^{(T)}$
- Feed it into the RNN-LM, computing output distribution  $y^{(t)}$  for every step  $t$ .
- Loss function on step  $t$  is **cross-entropy** between the predicted probability distribution  $\hat{y}^{(t)}$  and the true next word  $y^{(t)}$  (one-hot for  $x^{(t+1)}$ ):

$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = - \sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = - \log \hat{y}_{x_{t+1}}^{(t)}$$

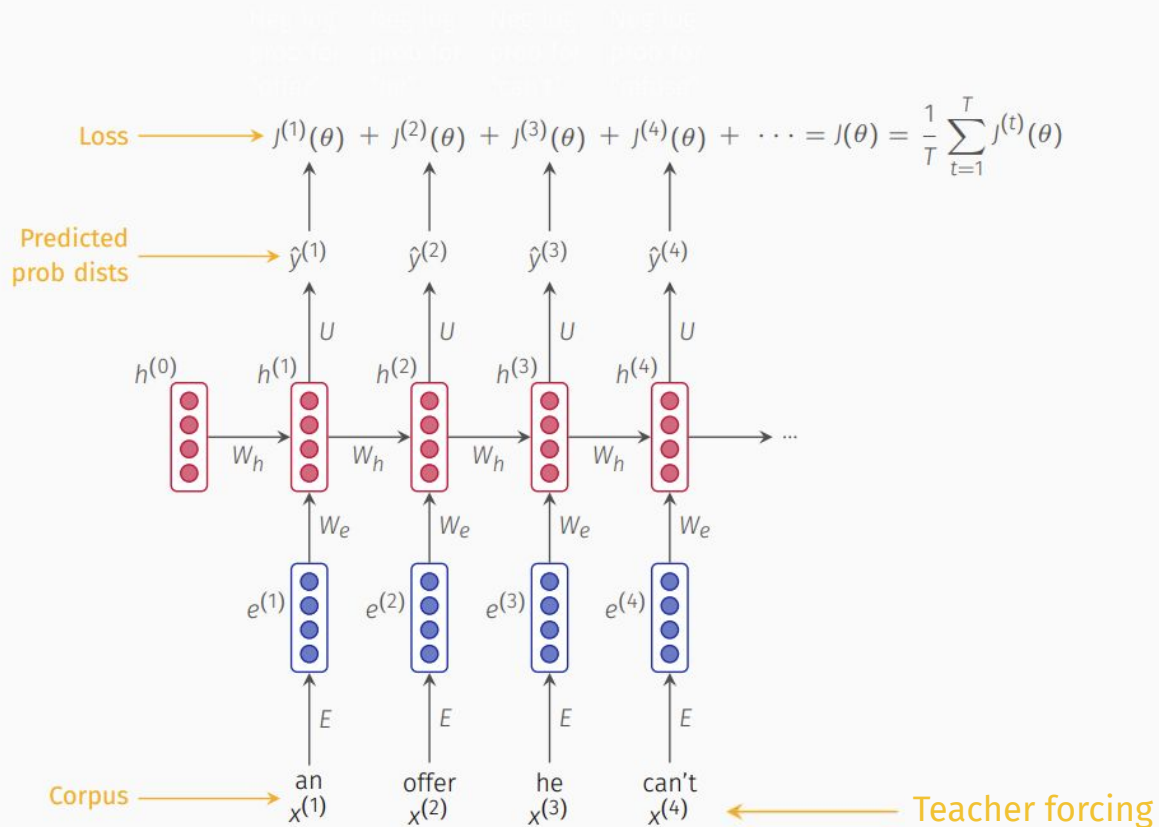
- Average this to get overall loss for the entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T - \hat{y}_{x_{t+1}}^{(t)}$$

# Training an RNN Language Model



# Training an RNN Language Model



# Computing Loss and Gradients in Practice

- In principle, we could compute loss and gradients across the whole corpus  $(x^{(1)}, \dots, x^{(T)})$  but that would be incredibly expensive!

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$$

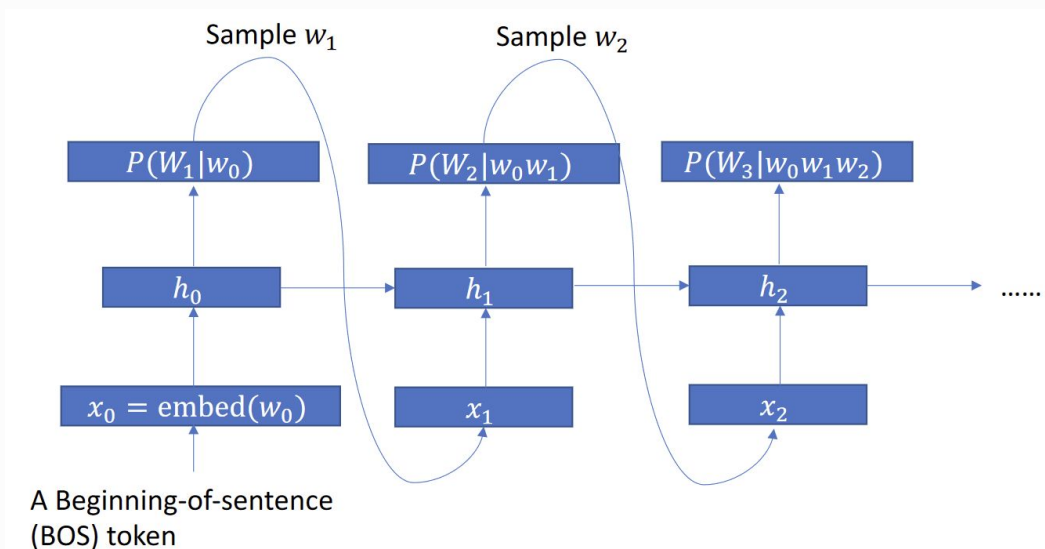
- Instead, we usually treat  $x^{(1)}, \dots, x^{(T)}$  as a document, or even a sentence
- This works much better with **Stochastic Gradient Descent**, which lets us compute loss and gradients for little chunks and update as we go.
- Actually, we do this in batches: compute  $J(\theta)$  for a batch of sentences; update weights; repeat.

## We Will Skip the Details of Backpropagation in RNNs for Now

- The fact that training RNNs involves backpropagation over timesteps, summing as you go, means that it (the **backpropagation through time** algorithm) is a bit more complicated than backpropagation in feedforward neural networks.
- We will skip these details for now, but you will want to learn them if you are doing serious work with RNNs.

# Generation with RNN LMs

- At each time step  $t$ , we sample  $w_t$  from  $P(W_t | \dots)$ , and feed it to the next timestep!
- LM with this kind of generation process is called autoregressive LM



# Advantages of RNN Language Models

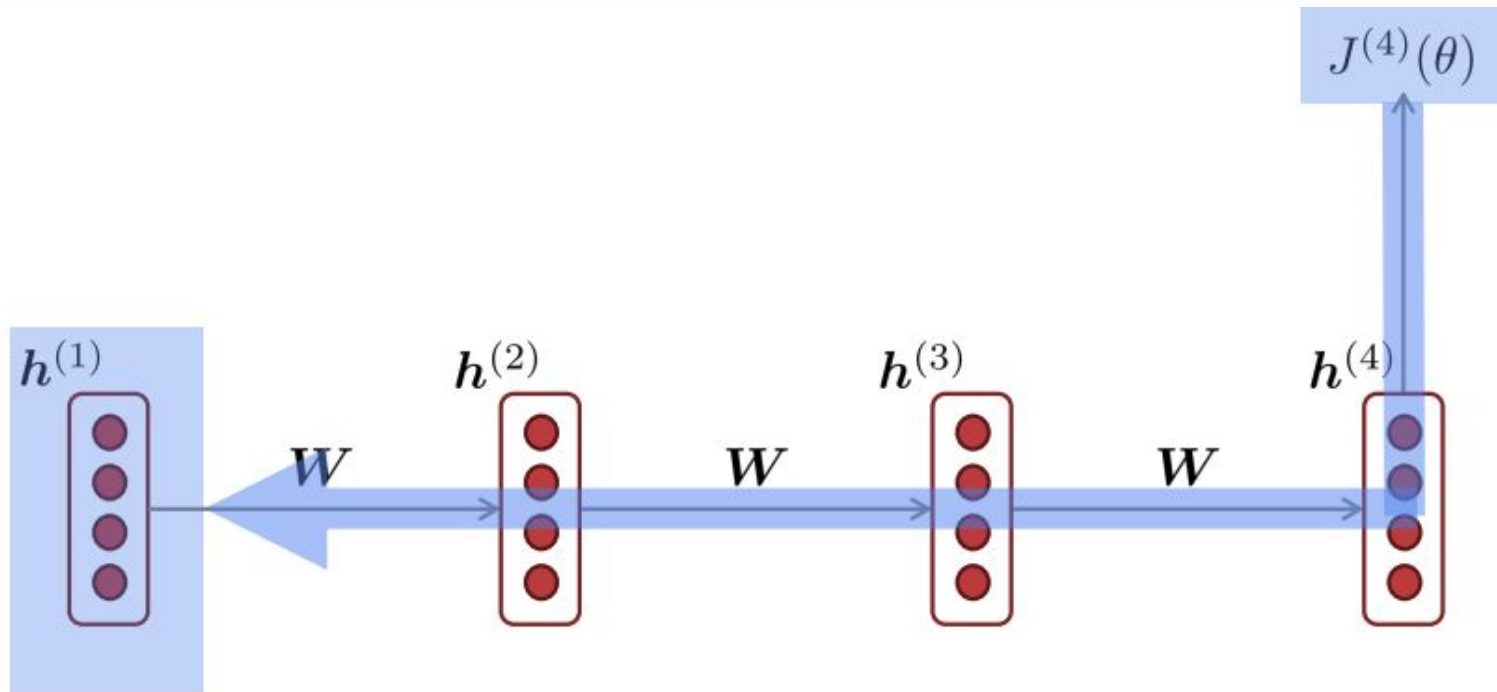
- Input can be of an arbitrary length
- Computation can use information from many steps back (in principle)
- Longer inputs do not mean larger model sizes
- Same weights applied at every time step—**symmetry**

# Disadvantages of RNN LMs

- Recurrent computation is **slow**
  - Computing  $h^{(t)}$  requires computing  $h^{(t-1)}$  which requires computing  $h^{(t-2)}$
  - **Cannot be parallelized**
- In practice, it is difficult to access information from many steps back (cf. the VANISHING GRADIENT PROBLEM)

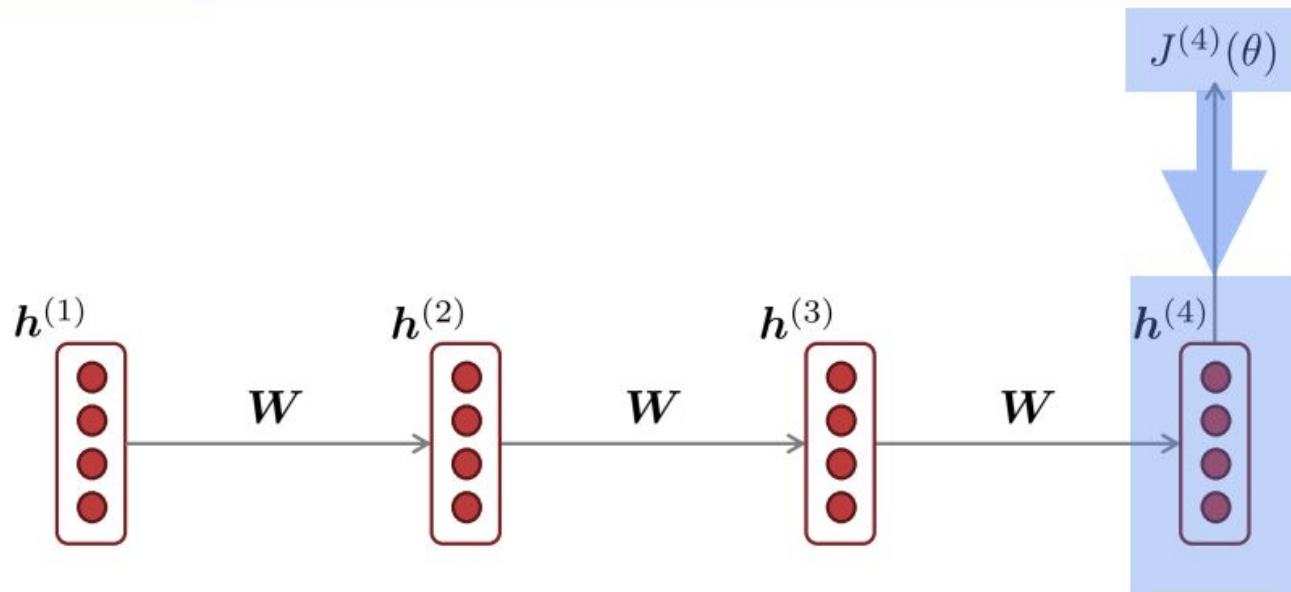


# Vanishing gradient problem



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

# Vanishing gradient problem



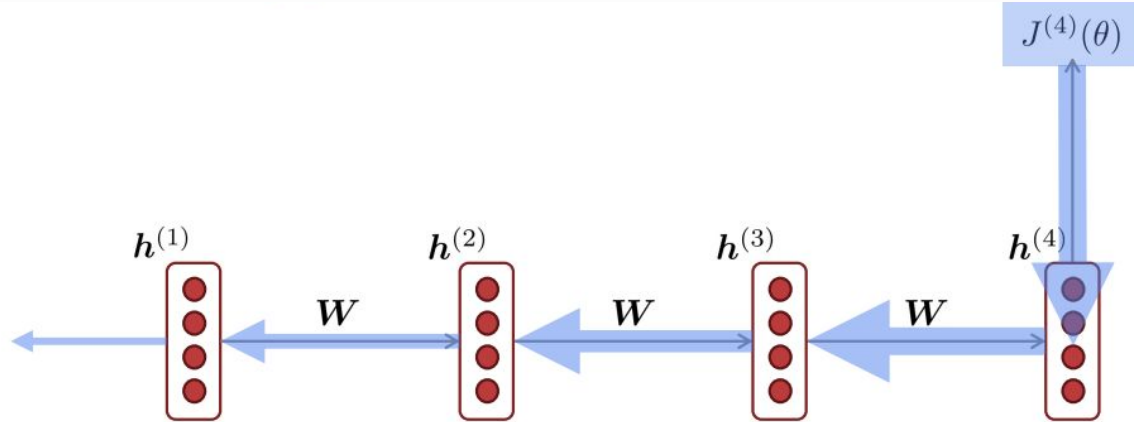
$$\frac{\partial J^{(4)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \times$$

$$\frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \times$$

$$\frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{h}^{(3)}} \times \frac{\partial J^{(4)}}{\partial \mathbf{h}^{(4)}}$$

chain rule!

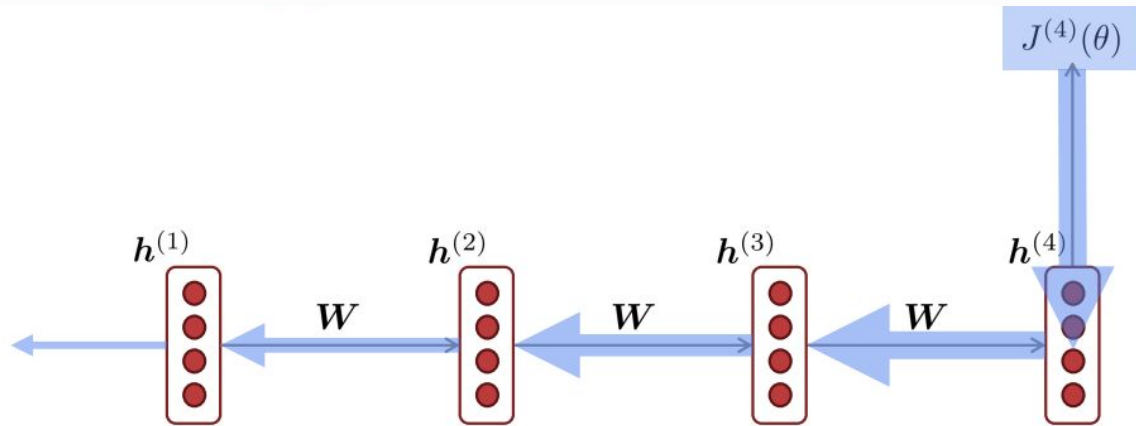
# Vanishing gradient problem



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

What happens if these are small?

# Vanishing gradient problem



$$\frac{\partial J^{(4)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}} \times \frac{\partial \mathbf{h}^{(3)}}{\partial \mathbf{h}^{(2)}} \times \frac{\partial \mathbf{h}^{(4)}}{\partial \mathbf{h}^{(3)}} \times \frac{\partial J^{(4)}}{\partial \mathbf{h}^{(4)}}$$

What happens if these are small?

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

# Vanishing gradient problem

- Gradient signal from far away is lost because it's much smaller than gradient signal from close-by.
- So model weights are basically updated only with respect to near effects, not long-term effects
- **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 be successful, need to model the dependency between “tickets” in the beginning and very end of the paragraph
  - If the gradient is small, can't learn this long-range dependency

# LSTMs

---

# LSTMs Address (but Do not Solve) the Vanishing and Exploding Gradient Problems

- LSTMs: Long Short Term Memory
- Process data sequentially, but keep hidden state through time
- Still subject, at some level, to vanishing gradients, but to a lesser degree than traditional RNNs
- Widely used in language modeling

# LSTMs: real-world success

- In 2013–2015, LSTMs started achieving state-of-the-art results
  - Successful tasks include handwriting recognition, speech recognition, machine translation, parsing, and image captioning, as well as language models
  - LSTMs became the dominant approach for most NLP tasks
- Now (2019–2023), Transformers have become dominant for all tasks
- For example, in WMT (a Machine Translation conference + competition):
  - In WMT 2014, there were 0 neural machine translation systems (!)
  - In WMT 2016, the summary report contains “RNN” 44 times (and these systems won)
  - In WMT 2019: “RNN” 7 times, “Transformer” 105 times

Source: "Findings of the 2016 Conference on Machine Translation (WMT16)", Bojar et al. 2016, <http://www.statmt.org/wmt16/pdf/W16-2301.pdf>

Source: "Findings of the 2018 Conference on Machine Translation (WMT18)", Bojar et al. 2018, <http://www.statmt.org/wmt18/pdf/WMT028.pdf>

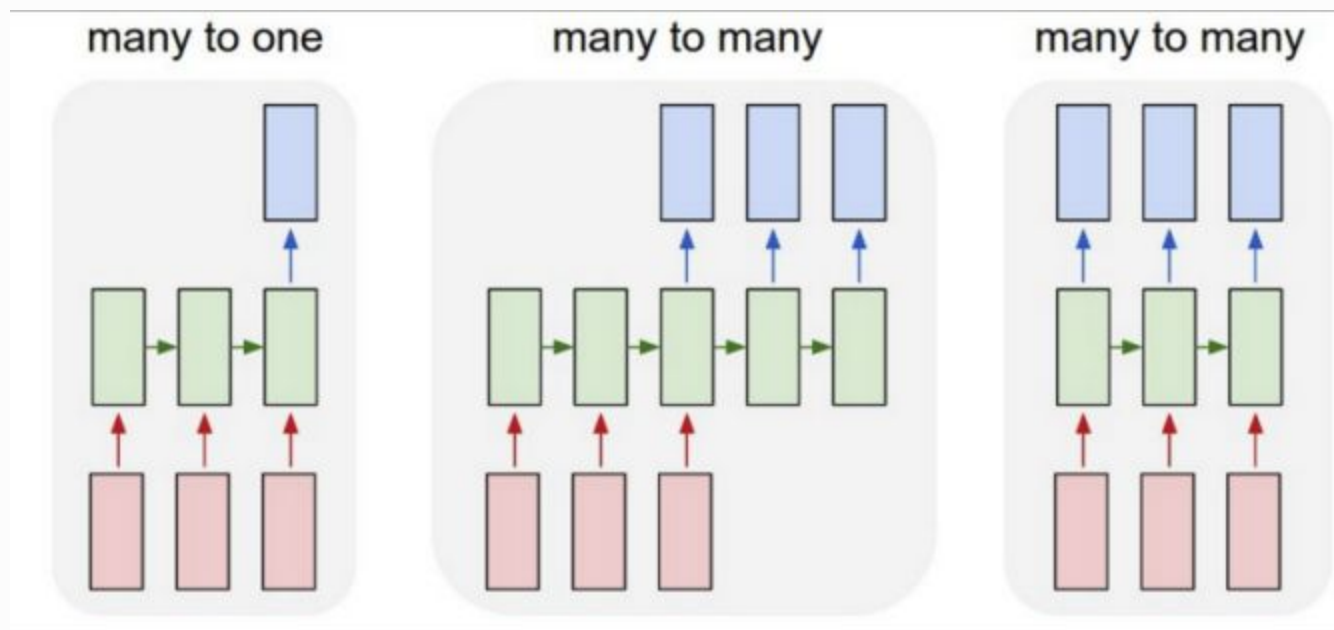
Source: "Findings of the 2019 Conference on Machine Translation (WMT19)", Barrault et al. 2019, <http://www.statmt.org/wmt18/pdf/WMT028.pdf>



# RNNs for other NLP tasks

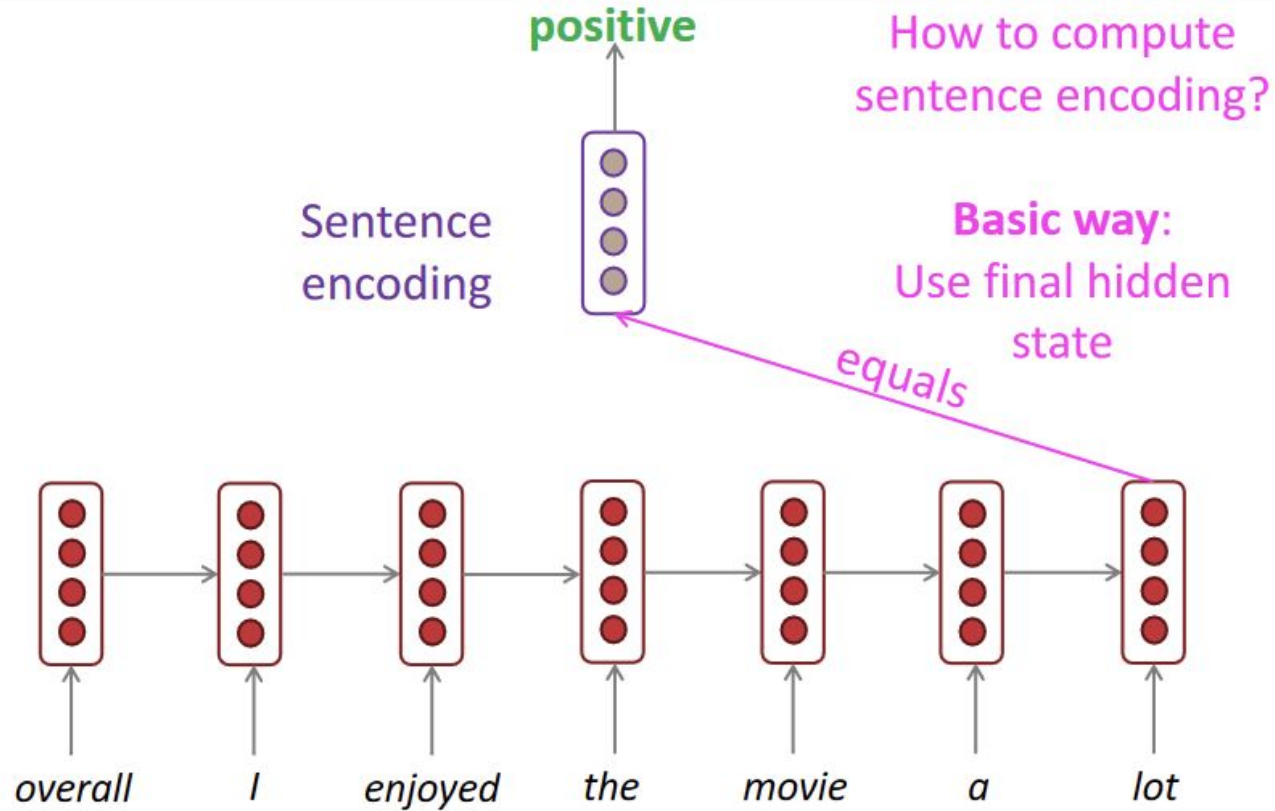
---

# RNNs for tasks other than language modeling

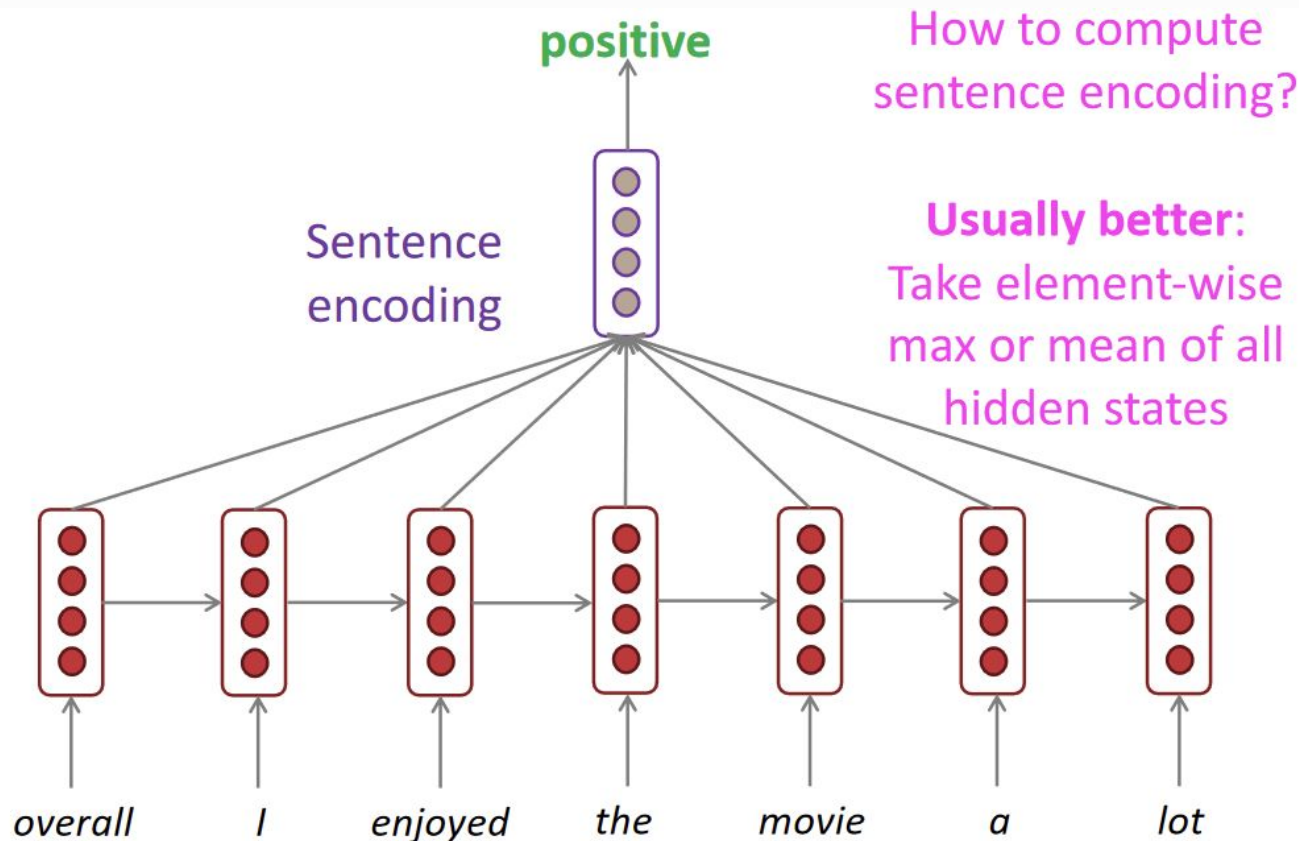


- Text classification (many to one)
- Encoder-decoder, machine translation (many to many)
- Language modeling, sequence labeling (many to many)

# RNNs to encode sentences for text classification



# RNNs to encode sentences for text classification



# Encoder-decoder architecture

---

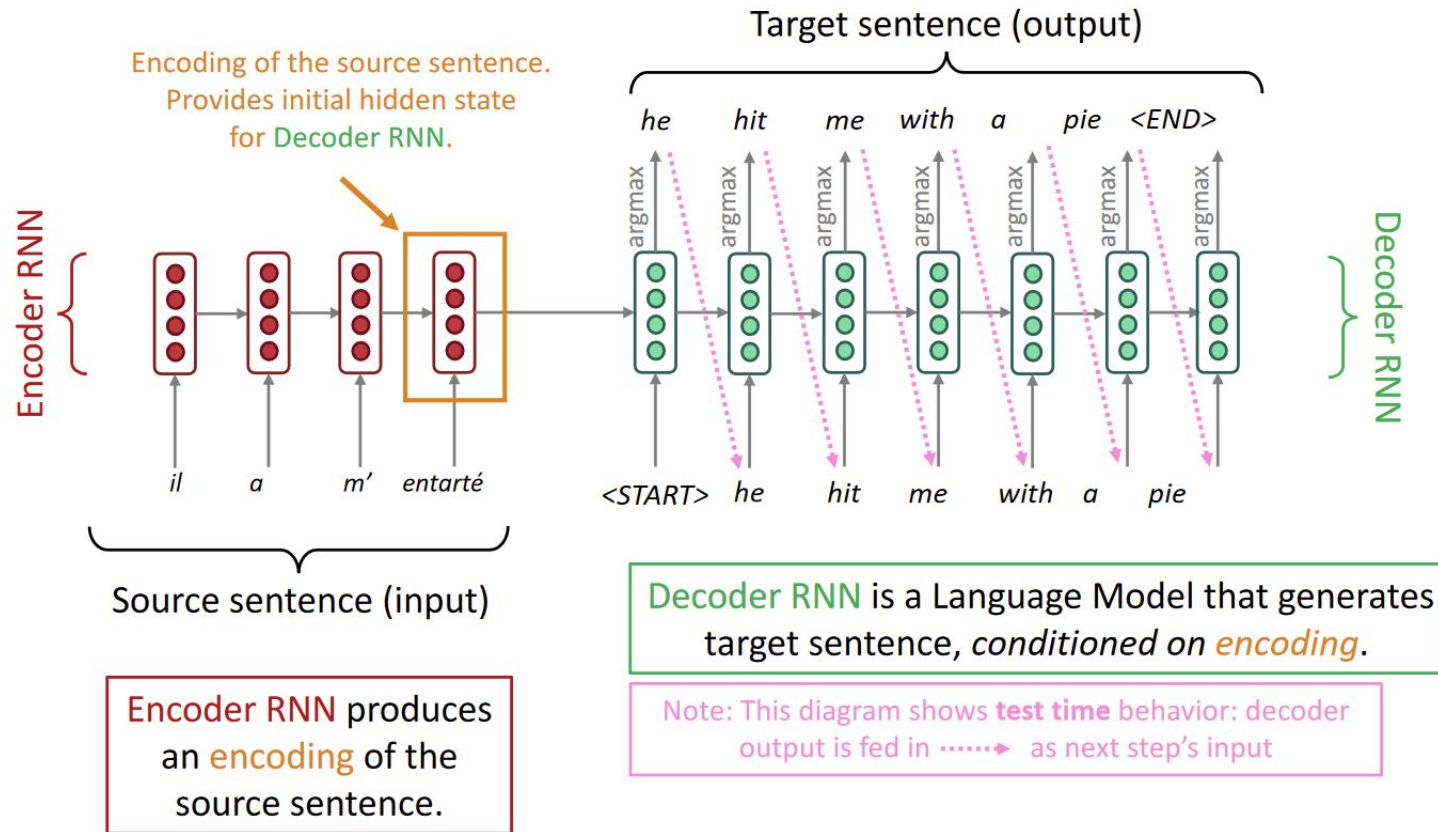
# The many names of encoder-decoder

- Encoder-decoder
- Sequence-to-sequence models
  - seq2seq
- Conditional language modeling

# Parts of an encoder-decoder

- **Encoder:** that accepts an input sequence,  $x^n$ , and generates a corresponding sequence of contextualized representations,  $h^n$ . LSTMs, convolutional neural networks, Transformers can all be encoders
- **Context:**  $c$ , which is a function of  $h^n$ , and conveys the essence of the input to the decoder.
- **Decoder:** which accepts  $c$  as input and generates an arbitrary length sequence of hidden states  $h^m$ , from which a corresponding sequence of output states  $y^m$ , can be obtained

# Encoder-decoder (seq2seq) architecture with RNNs





# Encoder-decoder (seq2seq) is versatile!

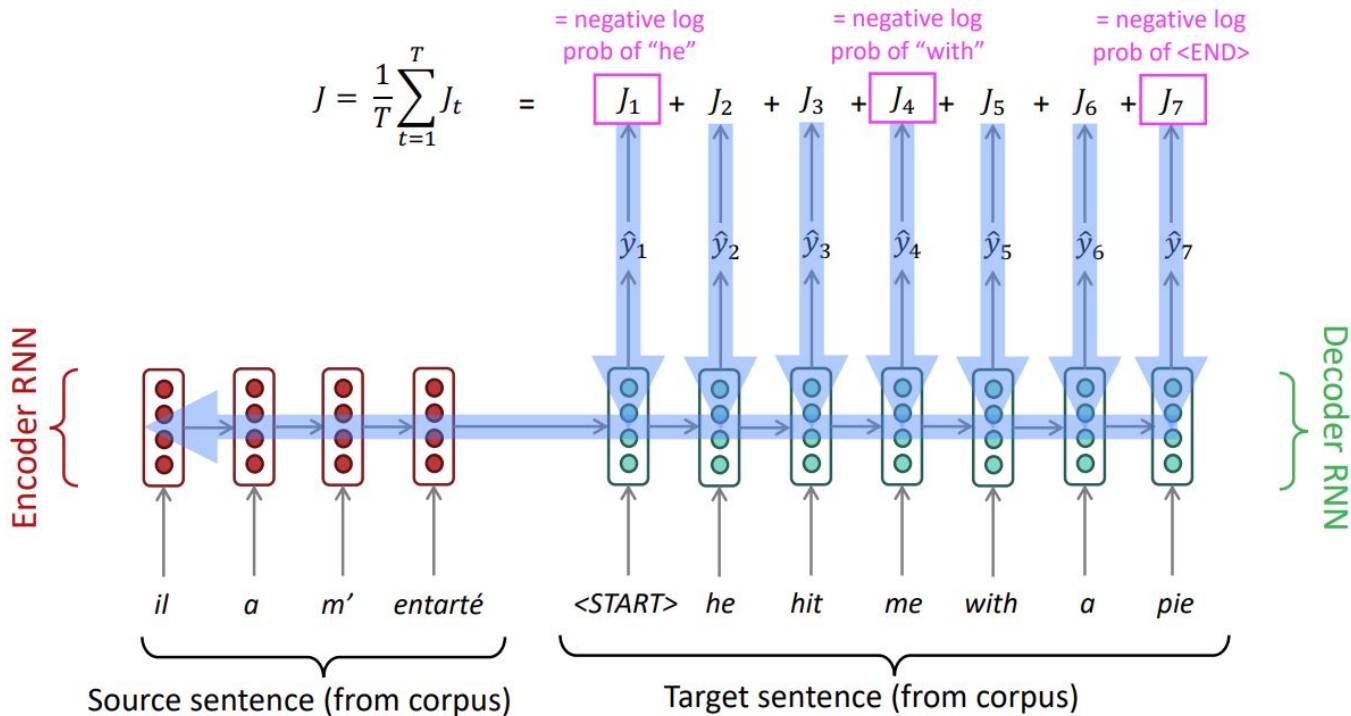
Many NLP tasks can be phrased as sequence-to-sequence:

- Summarization (long text → short text)
- Dialogue (previous utterances → next utterance)
- Parsing (input text → output parse as sequence)
- Code generation (natural language → Python code)

Training corpora needed:

- input <SEPARATOR> output

# Training an encoder-decoder RNN

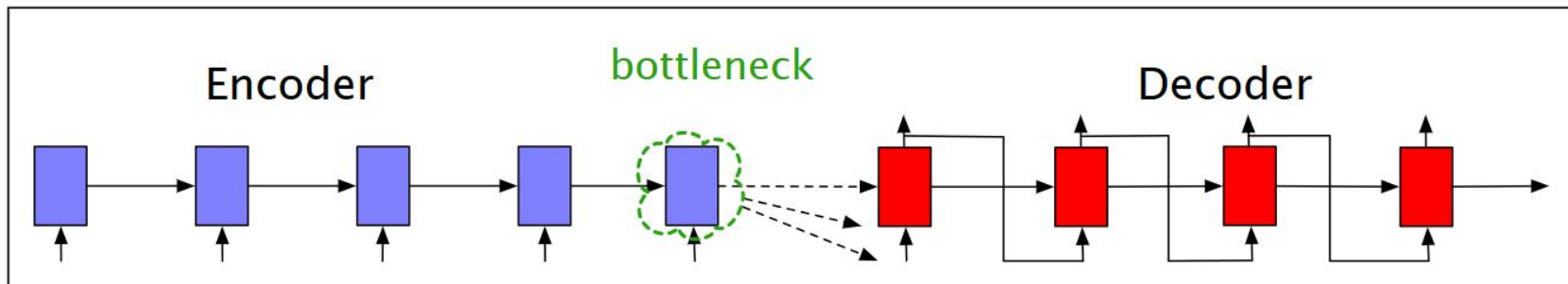


Seq2seq is optimized as a **single system**. Backpropagation operates "*end-to-end*".

# The attention mechanism

---

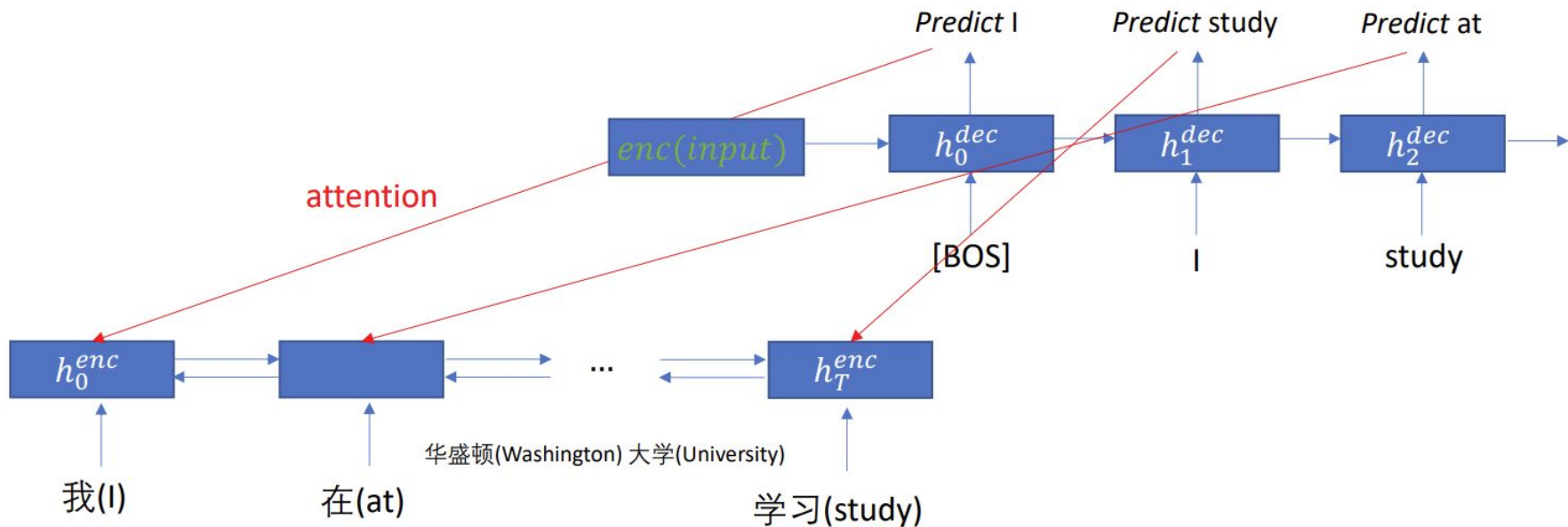
# Attention makes context available beyond the bottleneck



**Figure 9.21** Requiring the context  $c$  to be only the encoder's final hidden state forces all the information from the entire source sentence to pass through this representational bottleneck.

- Bottleneck means that early timesteps in the encoder aren't as accessible
- However, in tasks like MT, we may want to pay attention to different parts of the input in different timesteps.

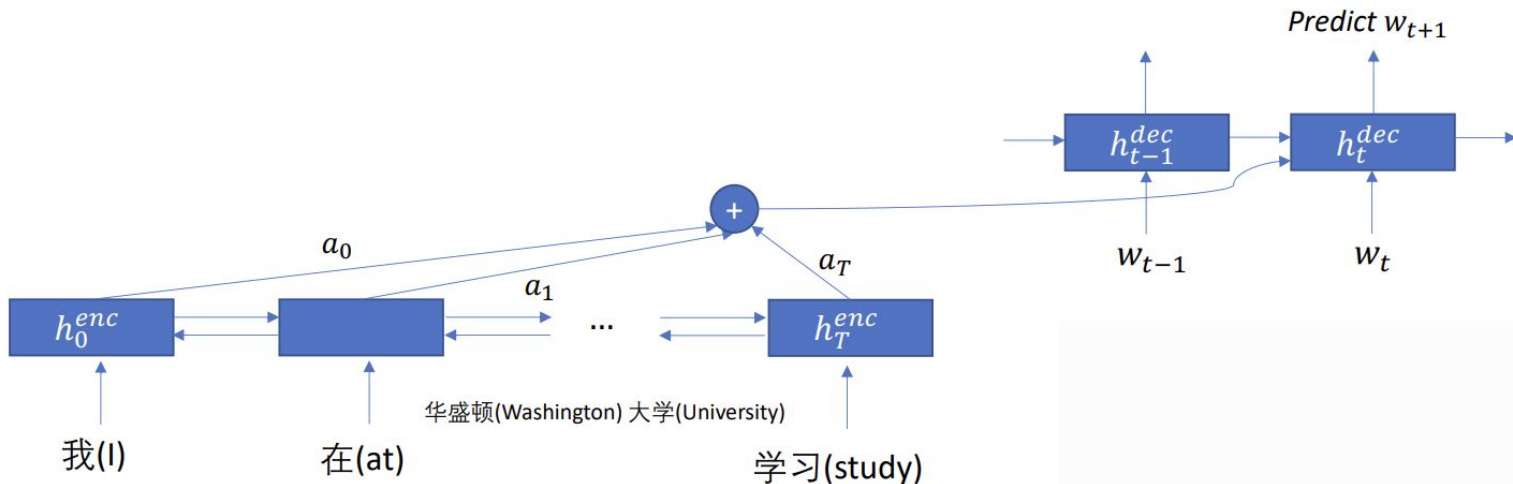
# Attention visualized



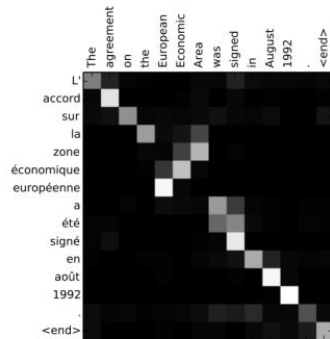
- “I study at University of Washington”
- This alignment is not trivial!
- The attention module is proposed to learn this alignment in an end-to-end fashion.

# The attention mechanism [Bahdanau et al. 2014]

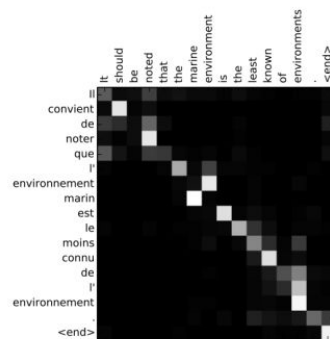
- We now focus on timestep  $t$ .
- For each encoder state  $h_i^{enc}$ , we compute an alignment score  $\hat{a}_i = (h_i^{enc})^T W_a h_{t-1}^{dec}$ .
- Then we get an attention distribution  $a = softmax(\hat{a})$ .
- We can then reweight the encoder states by  $a$  and pass  $\sum_i a_i h_i^{enc}$  to the decoder.
- The parameter  $W_a$  is shared across time steps.



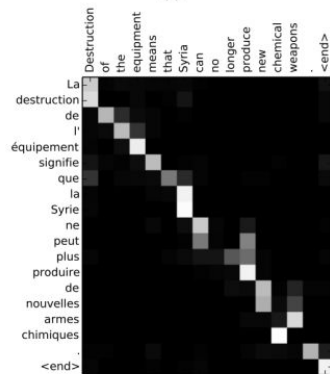
# Attention: learned alignment example



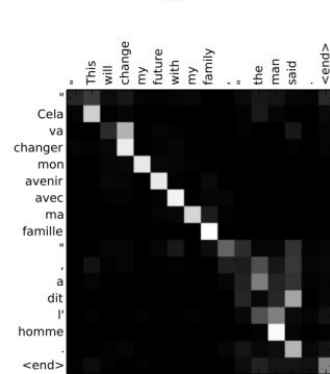
(a)



(b)



(c)



(d)

# Conclusion: RNNs part 2, encoder-decoder

- RNNs are effective neural networks for sequences
  - Can handle sequences of varying length
  - Can “remember” information from earlier timesteps
- RNNs can be used for language modeling and other NLP tasks
- LSTMs are a type of RNNs that handles vanishing gradient problem
- The encoder-decoder framework “encodes” sequential input and then “decodes” sequential output
- The attention mechanism weights encodings of relevant pieces of the input when producing output



*Questions?*