



CS 2731 / ISSP 2230

# Introduction to Natural Language Processing

Session 11: N-gram language models part 2, RNNs part 1

---

Michael Miller Yoder

February 14, 2024



University of  
Pittsburgh

School of Computing and Information

# Course logistics

- [Homework 2](#) is due **this Thu Oct 3**
  - Text classification
  - Written and programming components
  - Optional Kaggle competition for best LR and NN deception classifiers

# Course logistics

- Project groups have been formed
  - If you want to change groups, let me know
  - Please schedule a group meeting with me this week through my [Bookings link](#) in person (preferred) or on Zoom
    - What questions do you have about completing your project?
  - My office will change on Tuesday to **IS 604B**
- Next project deliverable is the proposal and literature on **Oct 17**
  - Instructions are on the [project webpage](#)
  - It's good to start the literature review early
  - Look for NLP papers in [ACL Anthology](#), [Semantic Scholar](#), and [Google Scholar](#)

# Lecture overview: N-gram language models part 2, RNNs part 1

- Sampling sentences from language models
- The problem of zeros
- Laplace smoothing
- Interpolation and backoff
- Neural language models
- RNN language modeling

# Sampling sentences from language models

---

# The Shannon Visualization Method

- Choose a random bigram ( $\langle s \rangle$ ,  $w$ ) according to its probability
- Now choose a random bigram ( $w$ ,  $x$ ) according to its probability
- And so on until we choose  $\langle /s \rangle$
- Then string the words together

$\langle s \rangle$  I  
I want  
want to  
to eat  
eat Chinese  
Chinese food  
food  $\langle /s \rangle$   
I want to eat Chinese food

# Approximating Shakespeare

1  
gram

–To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have

–Hill he late speaks; or! a more to leg less first you enter

2  
gram

–Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.

–What means, sir. I confess she? then all sorts, he is trim, captain.

3  
gram

–Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.

–This shall forbid it should be branded, if renown made it empty.

4  
gram

–King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;

–It cannot be but so.



# Not Shakespeare (no offense)

1  
gram

Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives

2  
gram

Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her

3  
gram

They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions



# The problem of zeros

---

# The Perils of Overfitting

N-grams only work well for word prediction if the test corpus looks like the training corpus

- In real life, it often doesn't
- We need to train robust models that generalize!
  - One kind of generalization: Zeros!
  - Things that don't ever occur in the training set but occur in the test set

# N-grams in the test set that weren't in the training set

Suppose our bigram LM, trained on Twitter, reads a document by the philosopher Wittgenstein:

*Whereof one cannot speak, thereof one must be silent.*

This contains the bigrams: *whereof one*, *one cannot*, *cannot speak*, *speak [comma]*, *[comma] thereof*, *thereof one*, *one must*, *must be*, *be silent*.

Suppose “whereof one” never occurs in the training corpus (**train**) but whereof occurs 20 times. According to MLE, it's probability is

$$P(\text{one}|\text{whereof}) = \frac{c(\text{whereof, one})}{c(\text{whereof})} = \frac{0}{20} = 0$$

The probability of the sentence is the **product** of the probabilities of the bigrams. What happens if one of the probabilities is zero?

# Two kinds of “zeros”

1. Completely unseen words in the test set
2. Words in unseen contexts in the test set

# Unknown Words

If we know all the words in advanced

- Vocabulary  $V$  is fixed
- Closed vocabulary task

Often we don't know this

- Out Of Vocabulary = OOV words
- Open vocabulary task

Instead: create an unknown word token <UNK>

- Training of <UNK> probabilities
- Create a fixed lexicon  $L$  of size  $V$
- At text normalization phase, any training word not in  $L$  changed to <UNK>
- Now we train its probabilities like a normal word
- At decoding time
- If text input: Use UNK probabilities for any word not in training

# Laplace smoothing

---

# The intuition of smoothing

When we have sparse statistics:

$P(w \mid \text{denied the})$

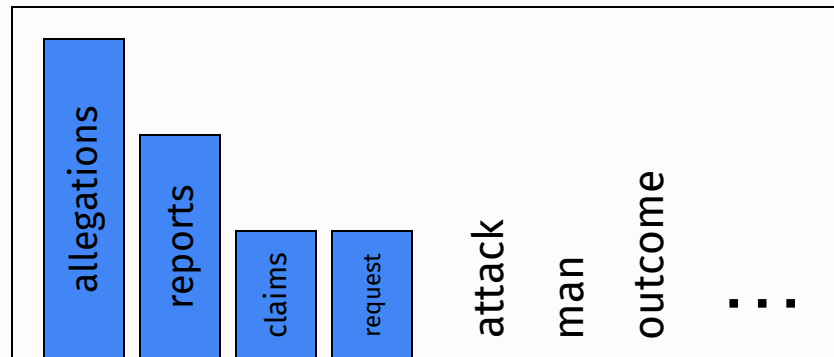
3 allegations

2 reports

1 claims

1 request

7 total



Steal probability mass to generalize better

$P(w \mid \text{denied the})$

2.5 allegations

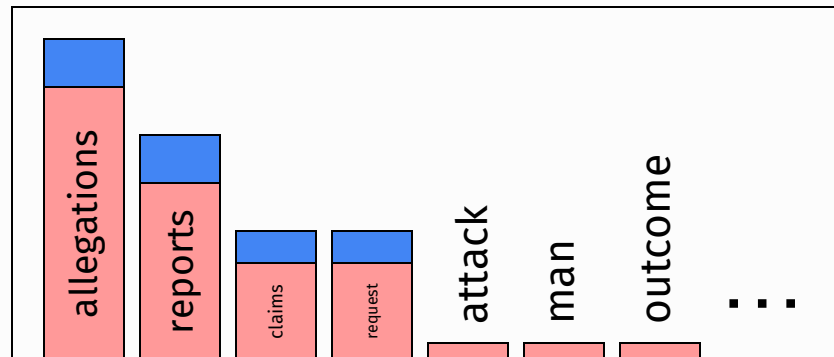
1.5 reports

0.5 claims

0.5 request

2 other

7 total





# Laplace smoothing: Pretending that we saw each word once more

$$\text{MLE estimate } P_{MLE}(w_i|w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$\text{Add-1 estimate } P_{Add-1}(w_i|w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + |V|}$$

Where  $V$  is the vocabulary of the corpus.

# Laplace Smooth Is too Blunt

It shifts too much probability mass away from attested ngrams and onto unattested ngrams, so it isn't used much for ngram LMs any more (there are better methods).

But remember that it does work:

- For text classification
- In domains where the number of zeros isn't as large as with ngrams

# Interpolation and backoff

---

# Backoff and Interpolation Let You Use Less Context

Suppose you have a context you don't know much about (because you have seen few or no relevant ngrams). You can condition your probabilities for these contexts on shorter contexts you know more about.

**Backoff** Use trigram if you have good evidence, otherwise bigram, otherwise unigram.

**Interpolation** Mix unigrams, bigrams, and trigrams together in one (weighted) probability soup.

Interpolation works better; backoff is sometimes cheaper.

# Linear interpolation takes into account different n-grams

The simplest way to do this is to **not** take context into account. The lambdas, in the following formula, are weighting factors:

$$\begin{aligned}\hat{P}(w_n|w_{n-2}w_{n-1}) = & \lambda_1 P(w_n|w_{n-2}w_{n-1}) \\ & + \lambda_2 P(w_n|w_{n-1}) \\ & + \lambda_3 P(w_n)\end{aligned}$$

where

$$\forall i \lambda_i \geq 0 \wedge \sum_i^n \lambda_i = 1$$

That is, the lambdas must sum to one.

# Lambdas Are Tuned Using a Held-Out dev Set



Choose  $\lambda$ s to maximize the probability of held-out data (**dev**):

- Fix the ngram probabilities (on **train**)
- Then search for  $\lambda$ s that give the largest probability to **dev**:

# Web-Scale Ngrams

How to deal with, e.g., Google N-gram corpus Pruning

- Only store N-grams with count > threshold.
- Remove singletons of higher-order n-grams
- Entropy-based pruning

Efficiency

- Efficient data structures like tries
- Bloom filters: approximate language models
- Store words as indexes, not strings
- Use Huffman coding to fit large numbers of words into two bytes
- Quantize probabilities (4-8 bits instead of 8-byte float)

# Stupid Backoff is Stupid but Efficient

No discounting, just use relative frequencies

$$S(w_i | w_{i-k+1}^{i-1}) = \begin{cases} \frac{c(w_{i-k+1}^i)}{c(w_{i-k+1}^{i-1})} & \text{if } c(w_{i-k+1}^i) > 0 \\ 0.4S(w_i | w_{i-k+2}^{i-1}) & \text{otherwise} \end{cases}$$

$$S(w_i) = \frac{c(w_i)}{N}$$



# Neural language models

---

# Why Neural LMs work better than N-gram LMs

Training data:

We've seen: I have to make sure that the cat gets fed.

Never seen: dog gets fed

Test data:

I forgot to make sure that the dog gets \_\_\_

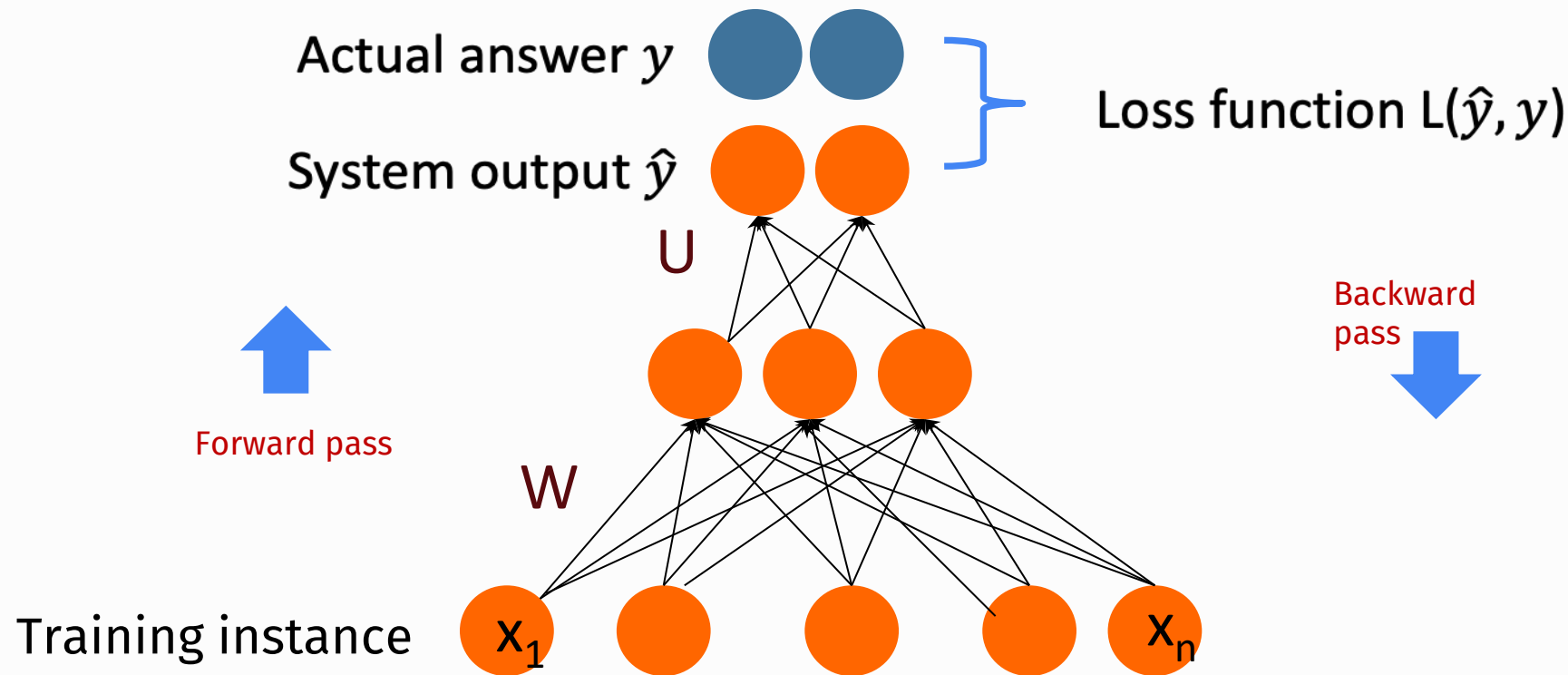
N-gram LM can't predict "fed"!

Neural LM can use similarity of "cat" and "dog" embeddings to generalize and predict "fed" after dog

# Language modeling with recurrent neural networks (RNNs)

---

# Reminder: training a 2-layer network



# There are Two Benefits to Training a Neural Language Model

If you train a neural language model, you get two things:

1. An algorithm that will allow you to predict the next word in a sequence
2. A set of embeddings  $\mathbf{E}$  that can be used to represent words in other tasks (assuming that you did not freeze the embedding layer)

FFNNs take an input of fixed **dimensions**—a fixed number of features, a fixed number of tokens

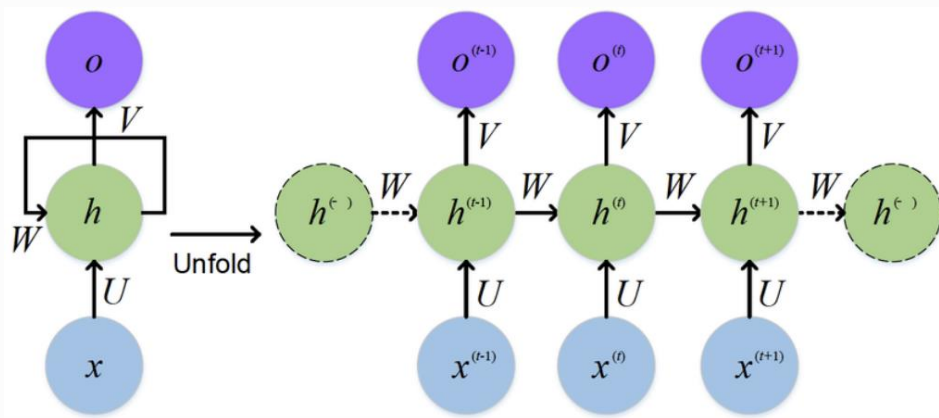
The number tokens in a text—even a sentence—can be **arbitrarily large** (or short)

RNNs help us address this issue



# The architecture of an RNN

- Special kind of multilayer neural network for modeling sequences
- Hidden layers between the input and output receive input not just from the input layer, **but also from the hidden layer at a preceding timestep**
- RNNs can “remember” information from earlier on



# An RNN Language Model

output distribution

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

hidden states

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

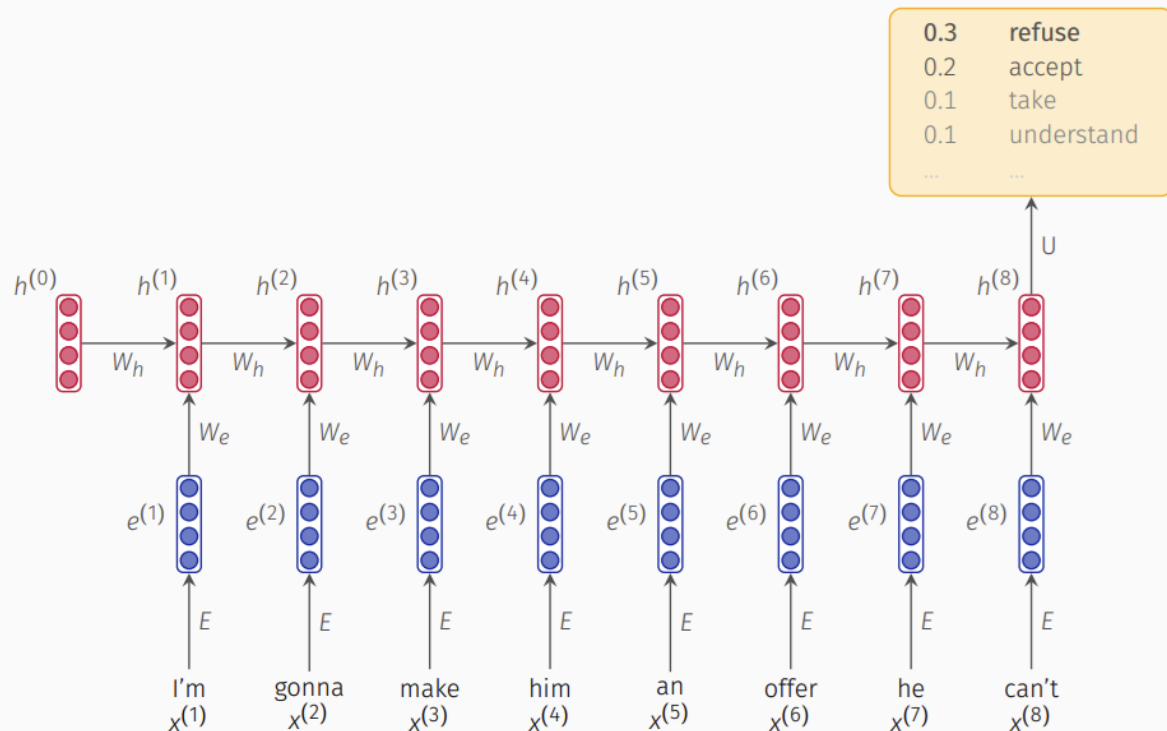
$h^{(0)}$  is the initial hidden state

word embeddings

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

one-hot vectors

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



# Training an RNN Language Model

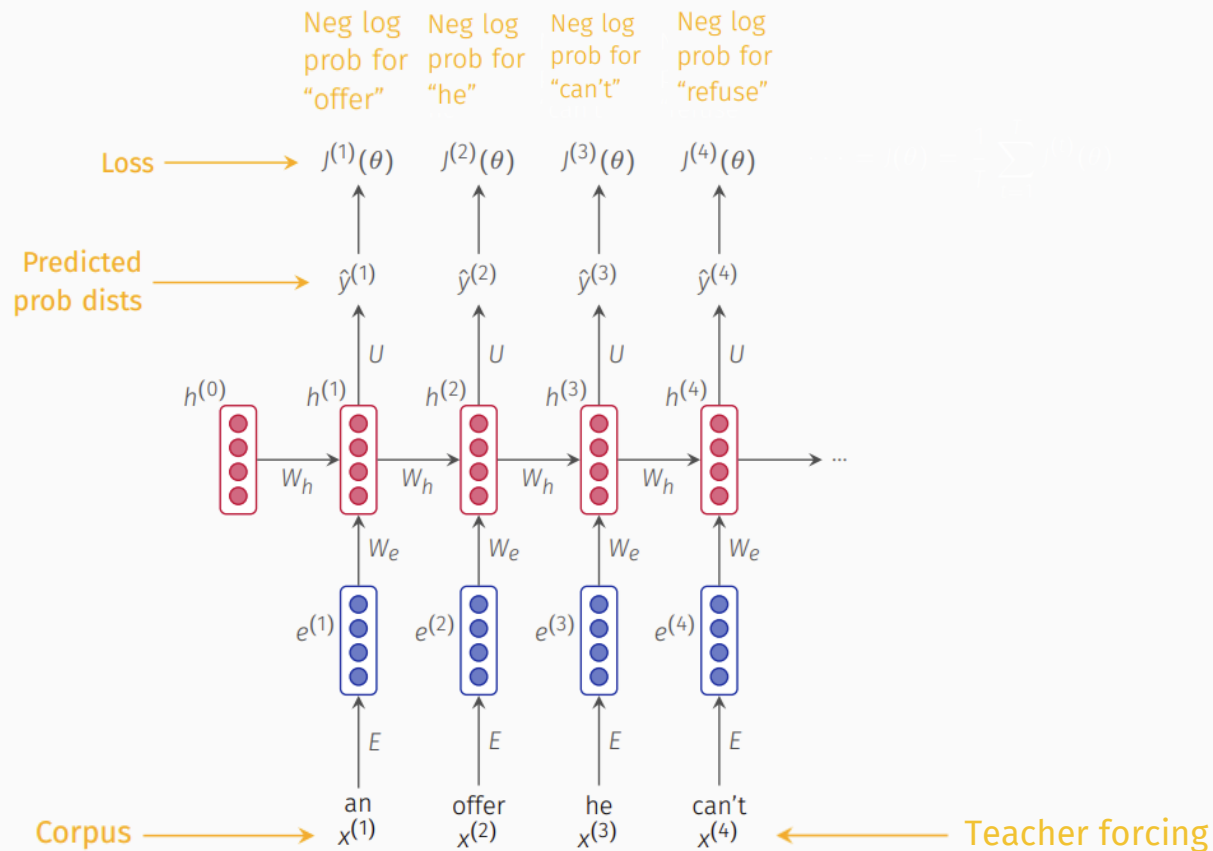
- Get a big corpus of text, which is a sequence of words  $x^{(1)}, \dots, x^{(T)}$
- Feed it into the RNN-LM, computing output distribution  ${}^{(t)}$  for every step  $t$ .
- Loss function on step  $t$  is **cross-entropy** between the predicted probability distribution  $\hat{\mathbf{y}}^{(t)}$ 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

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

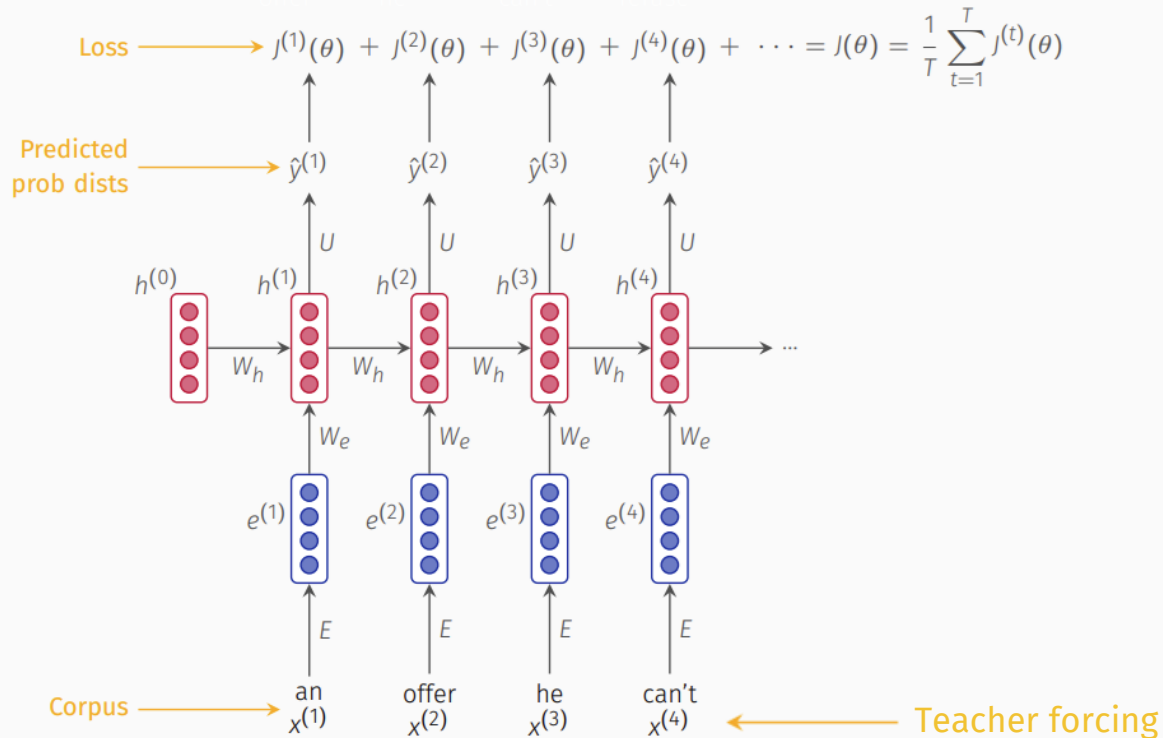
- Average this to get overall loss for the entire training set:

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

# 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!

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

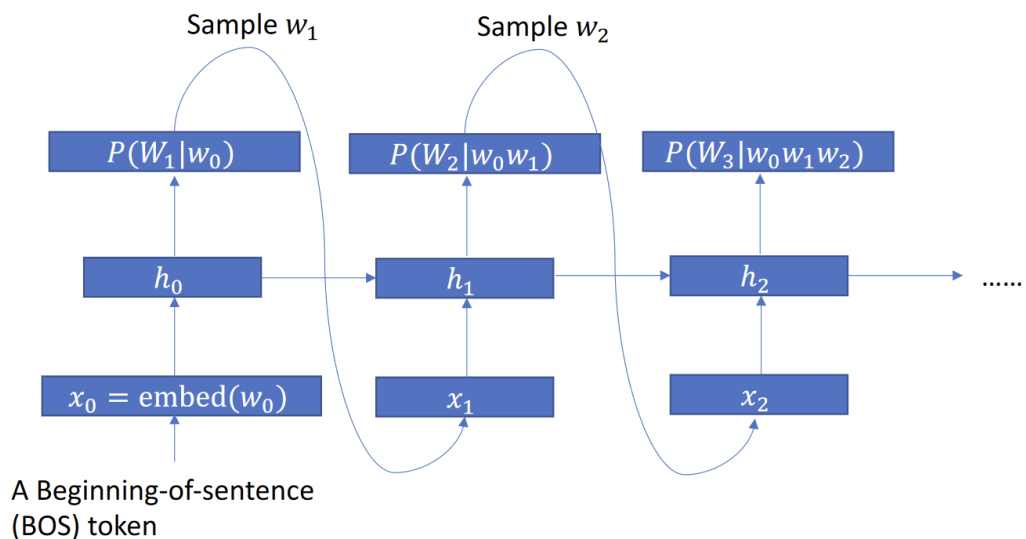
- Instead, we usually treat  $x^{(1)}, \dots, x^{(T)}$  as a document, or even a sentence
- This works much better with **Stochastic Gradient Descent**, which lets us compute loss and gradients for little chunks and update as we go.
- Actually, we do this in batches: compute  $J(\theta)$  for a batch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





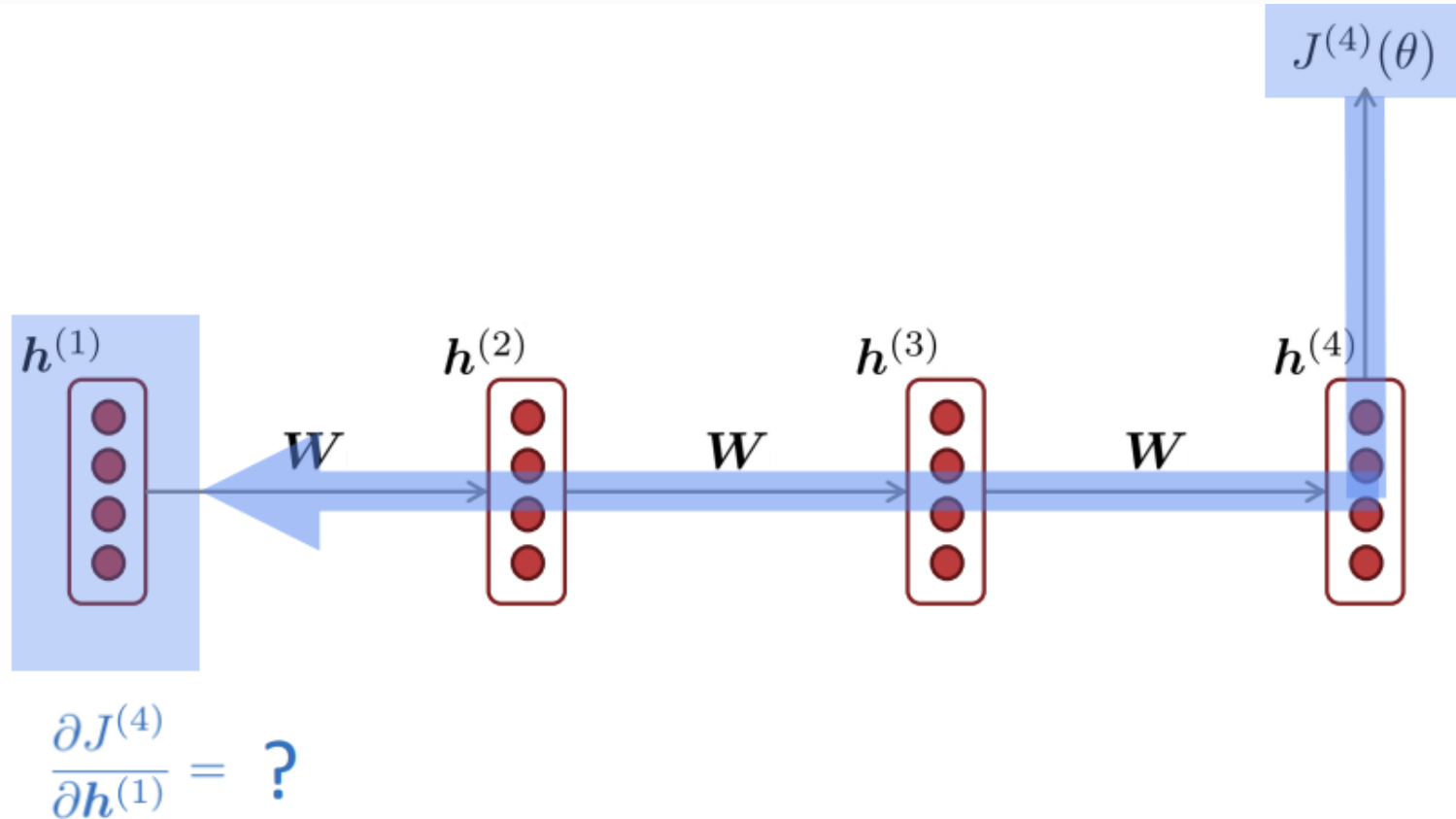
# Advantages of RNN Language Models

- 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**

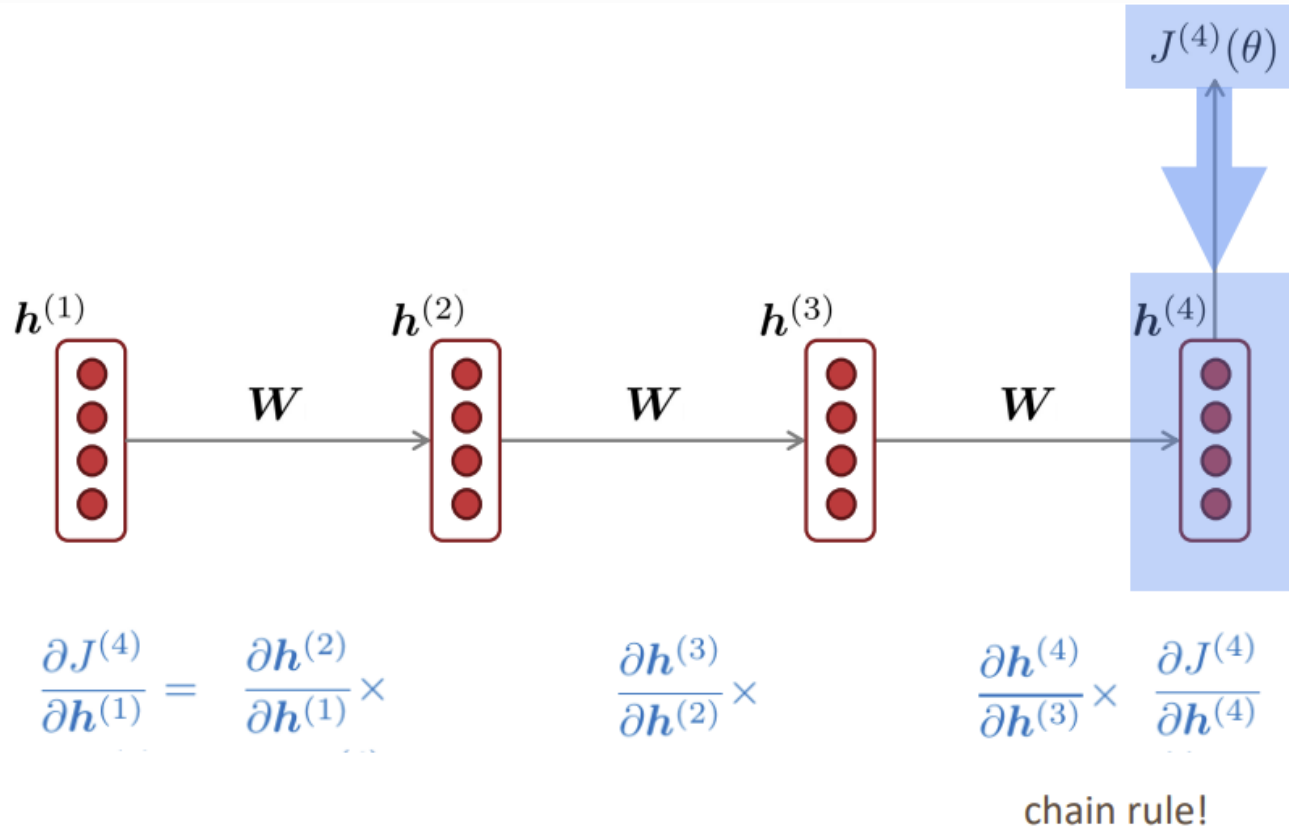
# Disadvantages of RNN LMs

- 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)

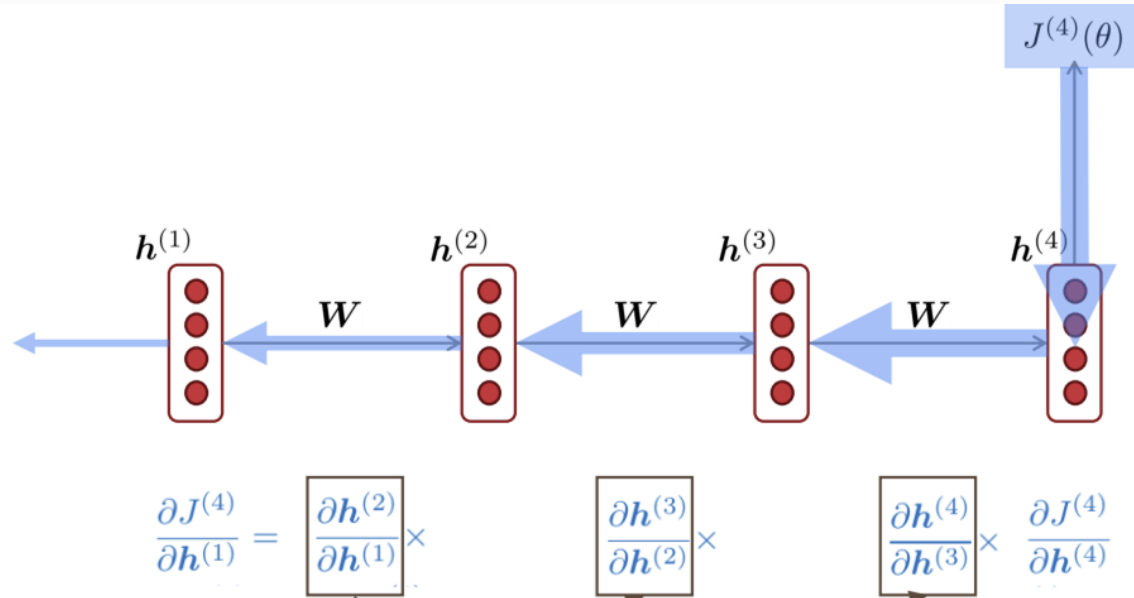
# Vanishing gradient problem



# Vanishing gradient problem

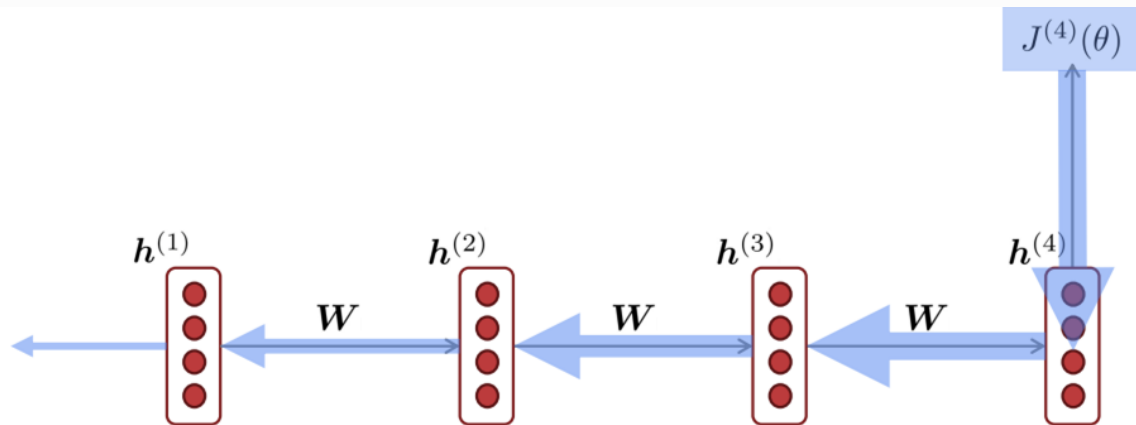


# Vanishing gradient problem



What happens if these are small?

# 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:**  
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