

MIE524 Data Mining Neural Language Models

Slides by:

Christopher Manning, Jurafsky and Martin (SLP 3rd ed.), Mohit Iyyer

MIE524: Course Topics (Tentative)

Large-scale Machine Learning

Learning Embedding
(NN / AE)

Decision Trees

Ensemble Models
(GBTs)

High-dimensional Data

Locality sensitive hashing

Clustering

Dimensionality reduction

Graph Data

Processing Massive Graphs

PageRank, SimRank

Graph Representation Learning

Applications

Recommender systems

Association Rules

Neural Language Models

Computational Models:

Single Machine

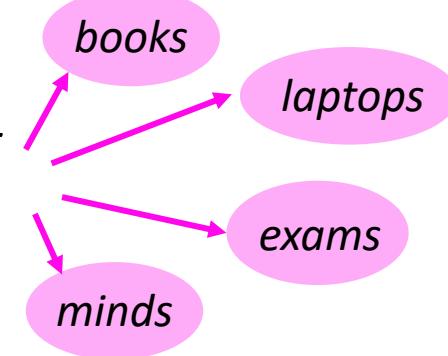
MapReduce/Spark

GPU

2. Language Modeling

- **Language Modeling** is the task of predicting what word comes next

the students opened their _____



- More formally: given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$:

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$$

where $x^{(t+1)}$ can be any word in the vocabulary $V = \{w_1, \dots, w_{|V|}\}$

- A system that does this is called a **Language Model**

Language Modeling

- You can also think of a Language Model as a system that assigns a probability to a piece of text
- For example, if we have some text $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$, then the probability of this text (according to the Language Model) is:

$$\begin{aligned} P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) &= P(\mathbf{x}^{(1)}) \times P(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}) \times \cdots \times P(\mathbf{x}^{(T)} | \mathbf{x}^{(T-1)}, \dots, \mathbf{x}^{(1)}) \\ &= \prod_{t=1}^T P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)}) \end{aligned}$$



This is what our LM provides

Probabilistic Language Modeling

Goal: compute the probability of a sentence or sequence of words:

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

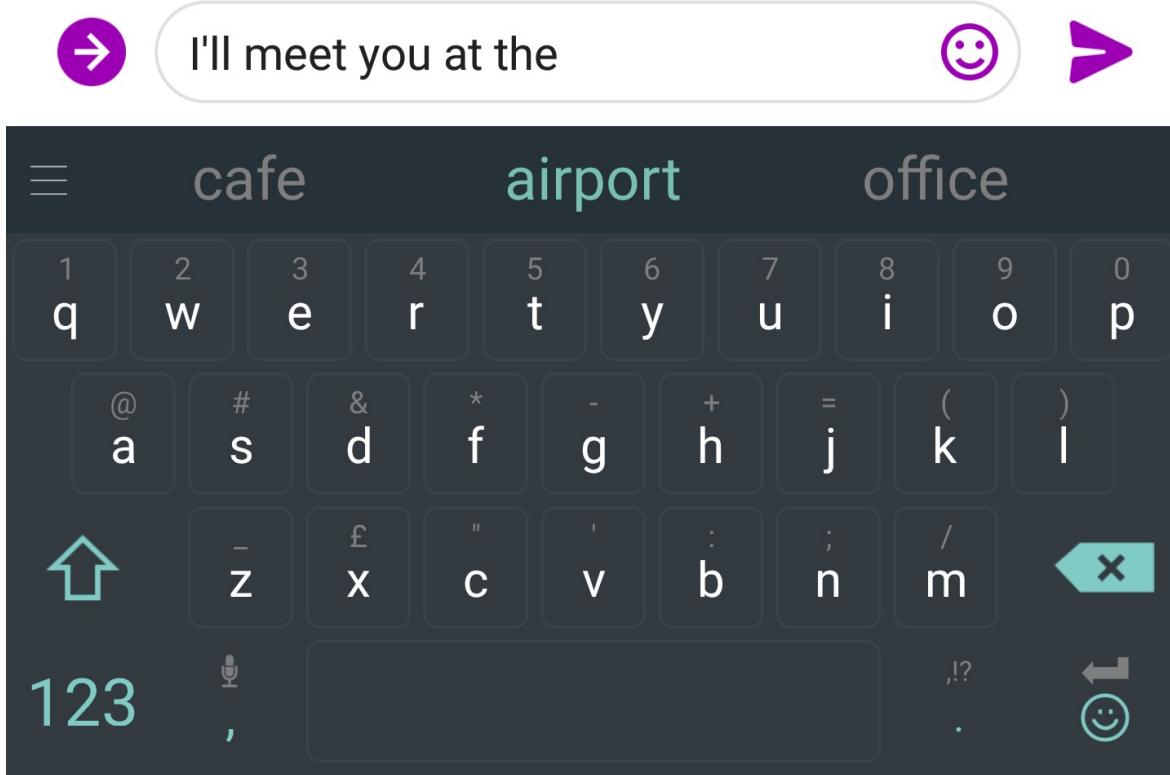
Related task: probability of an upcoming word:

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

A model that computes either of these:

$P(W)$ or $P(w_n | w_1, w_2 \dots w_{n-1})$ is called a **language model**.

You use Language Models every day!



You use Language Models every day!



A screenshot of a Google search interface. At the top, there is a search bar containing the partial query "what is the |". To the right of the search bar is a microphone icon. Below the search bar, a dropdown menu lists ten suggested search queries, each starting with "what is the":

- what is the **weather**
- what is the **meaning of life**
- what is the **dark web**
- what is the **xfl**
- what is the **doomsday clock**
- what is the **weather today**
- what is the **keto diet**
- what is the **american dream**
- what is the **speed of light**
- what is the **bill of rights**

At the bottom of the interface are two buttons: "Google Search" and "I'm Feeling Lucky".

n-gram Language Models

the students opened their _____

- **Question:** How to learn a Language Model?
- **Answer** (pre- Deep Learning): learn an *n-gram Language Model!*
- **Definition:** An *n-gram* is a chunk of n consecutive words.
 - **unigrams:** “the”, “students”, “opened”, “their”
 - **bigrams:** “the students”, “students opened”, “opened their”
 - **trigrams:** “the students opened”, “students opened their”
 - **four-grams:** “the students opened their”
- **Idea:** Collect statistics about how frequent different n-grams are and use these to predict next word.

n-gram Language Models

- First we make a **Markov assumption**: $x^{(t+1)}$ depends only on the preceding $n-1$ words

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)}) = P(x^{(t+1)} | \underbrace{x^{(t)}, \dots, x^{(t-n+2)}}_{n-1 \text{ words}}) \quad (\text{assumption})$$

$$\begin{aligned} \text{prob of a n-gram} &\rightarrow P(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)}) \\ \text{prob of a (n-1)-gram} &\rightarrow P(x^{(t)}, \dots, x^{(t-n+2)}) \end{aligned} \quad (\text{definition of conditional prob})$$

- Question:** How do we get these n -gram and $(n-1)$ -gram probabilities?
- Answer:** By **counting** them in some large corpus of text!

$$\approx \frac{\text{count}(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \dots, x^{(t-n+2)})} \quad (\text{statistical approximation})$$

n-gram Language Models: Example

Suppose we are learning a 4-gram Language Model.

~~as the proctor started the clock, the~~ students opened their _____

discard condition on this

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count(students opened their } \mathbf{w})}{\text{count(students opened their)}}$$

For example, suppose that in the corpus:

- “students opened their” occurred 1000 times
- “students opened their books” occurred 400 times
 - $\rightarrow P(\text{books} | \text{students opened their}) = 0.4$
- “students opened their exams” occurred 100 times
 - $\rightarrow P(\text{exams} | \text{students opened their}) = 0.1$

Should we have discarded
the “proctor” context?

Sparsity Problems with n-gram Language Models

Sparsity Problem 1

Problem: What if “*students opened their w*” never occurred in data? Then w has probability 0!

(Partial) Solution: Add small δ to the count for every $w \in V$. This is called *smoothing*.

$$P(w|\text{students opened their}) = \frac{\text{count(students opened their } w\text{)}}{\text{count(students opened their)}}$$

Sparsity Problem 2

Problem: What if “*students opened their*” never occurred in data? Then we can’t calculate probability for *any w*!

(Partial) Solution: Just condition on “*opened their*” instead. This is called *backoff*.

Note: Increasing n makes sparsity problems worse. Typically, we can’t have n bigger than 5.

Storage Problems with n-gram Language Models

Storage: Need to store count for all n -grams you saw in the corpus.

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count(students opened their } \mathbf{w})}{\text{count(students opened their)}}$$

Increasing n or increasing corpus increases model size!

n-gram Language Models in practice

- You can build a simple trigram Language Model over a 1.7 million word corpus (Reuters) in a few seconds on your laptop*

today the _____

Business and financial news

get probability distribution

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

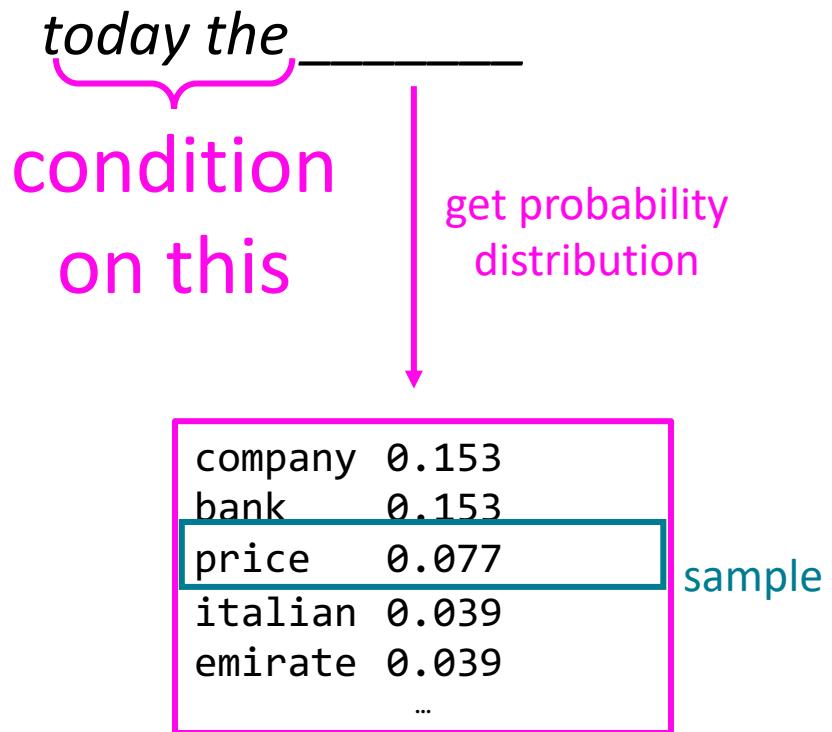
Sparsity problem:
not much granularity
in the probability
distribution

Otherwise, seems reasonable!

* Try for yourself: <https://nlpforhackers.io/language-models/>

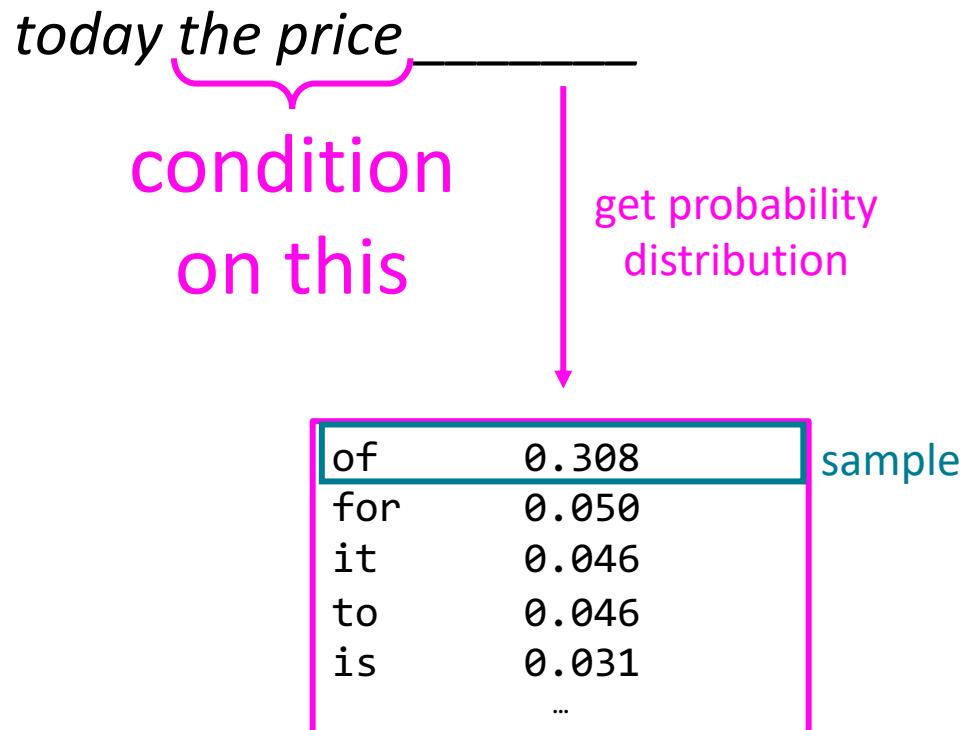
Generating text with a n-gram Language Model

You can also use a Language Model to generate text



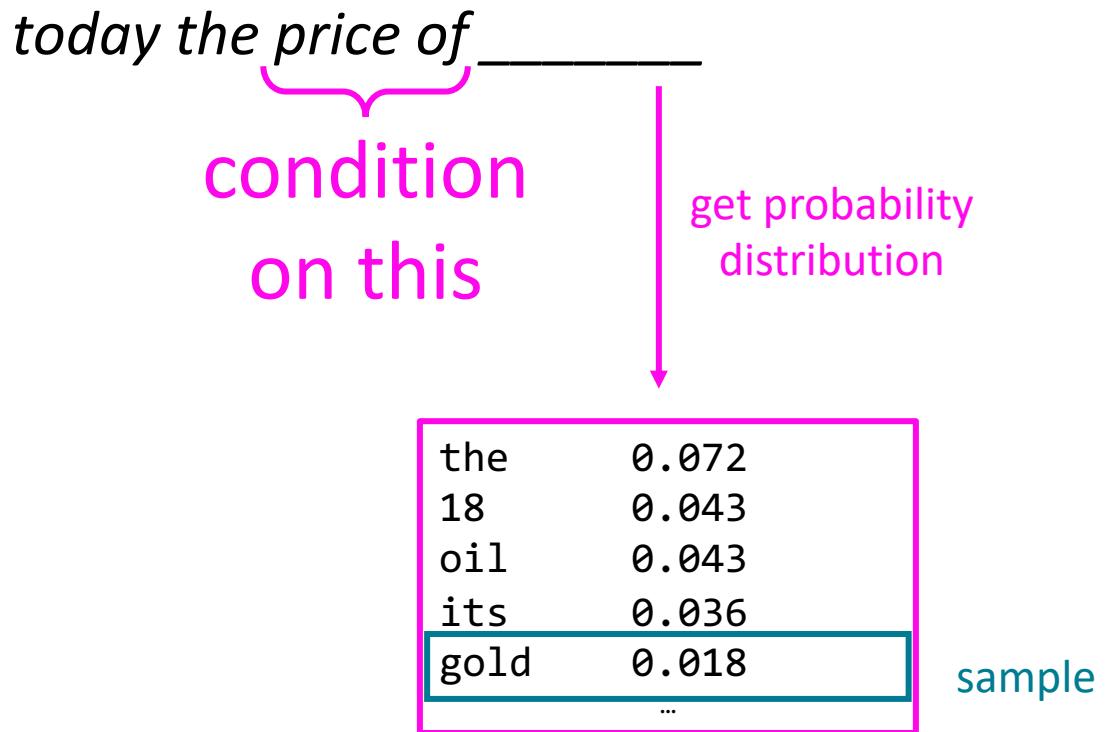
Generating text with a n-gram Language Model

You can also use a Language Model to generate text



Generating text with a n-gram Language Model

You can also use a Language Model to generate text



Generating text with a n-gram Language Model

You can also use a Language Model to generate text

*today the price of gold per ton , while production of shoe
lasts and shoe industry , the bank intervened just after it
considered and rejected an imf demand to rebuild depleted
european stocks , sept 30 end primary 76 cts a share .*

Surprisingly grammatical!

...but **incoherent**. We need to consider more than
three words at a time if we want to model language well.

But increasing n worsens sparsity problem,
and increases model size...

another issue:

- We treat all words / prefixes independently of each other!

students opened their ____

pupils opened their ____

scholars opened their ____

undergraduates opened their ____

students turned the pages of their ____

students attentively perused their ____

...

Shouldn't we *share information* across these semantically-similar prefixes?

one-hot vectors

- n-gram models rely on the “bag-of-words” assumption
- represent each word/n-gram as a vector of zeros with a single 1 identifying its index in the vocabulary

vocabulary
i
hate
love
the
movie
film

movie = $<0, 0, 0, 0, 1, 0>$

film = $<0, 0, 0, 0, 0, 1>$

what are the issues
of representing a
word this way?

all words are equally (dis)similar!

movie = $<0, 0, 0, 0, 1, 0>$

film = $<0, 0, 0, 0, 0, 1>$

dot product is zero!

these vectors are orthogonal

What we want is a representation space in which words, phrases, sentences etc. that are semantically similar also have similar representations!

Enter neural networks!

Students opened their



neural language
model



books

Enter neural networks!

Students opened their

forward pass, or how we compute a prediction of the next word given an existing neural language model



neural language
model

backward pass, or how we train a neural language model on a training dataset using the backpropagation algorithm

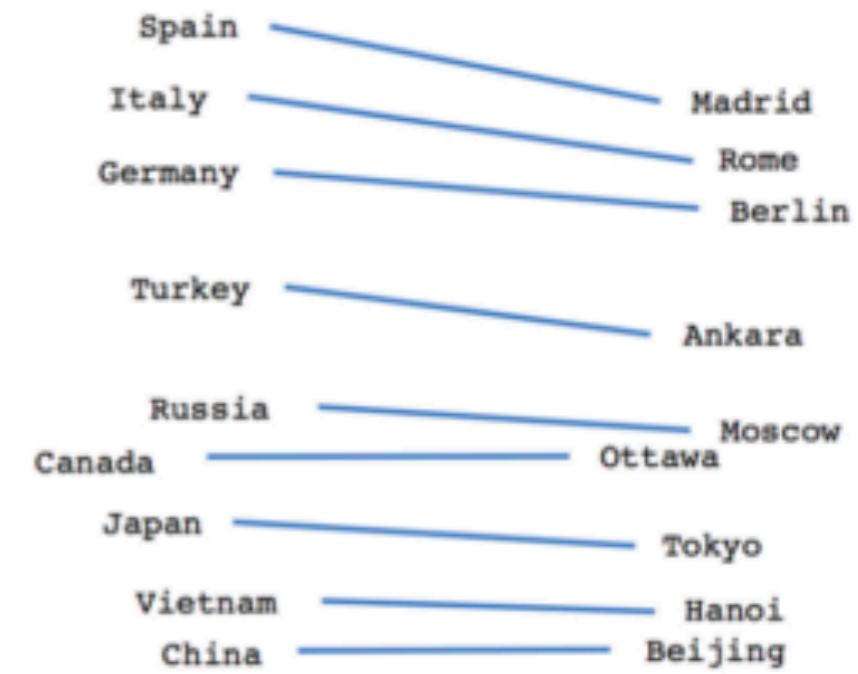
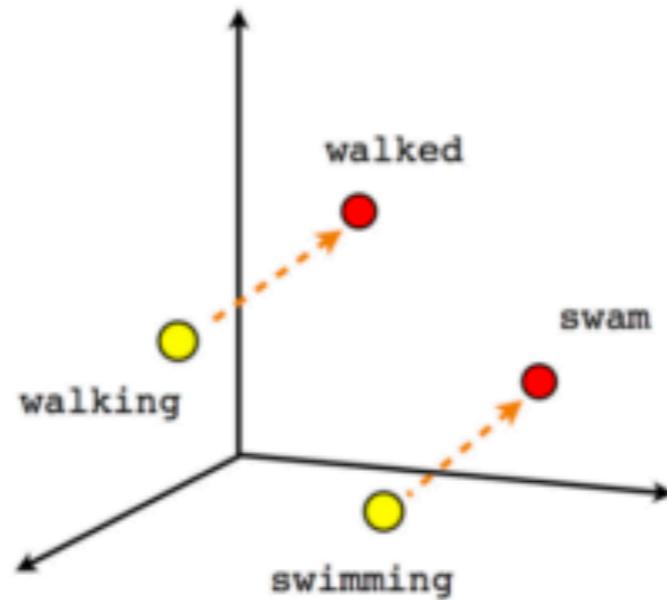
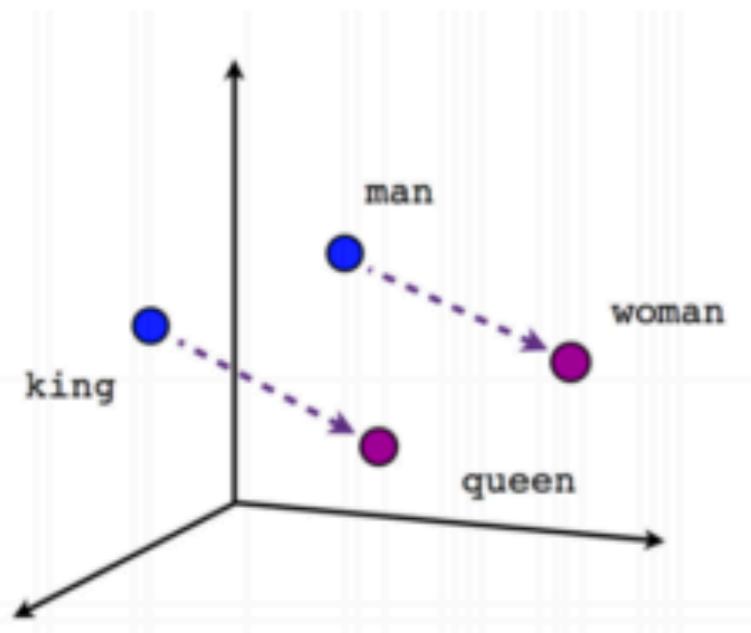


books

words as basic building blocks

- represent words with low-dimensional vectors called **embeddings** (Mikolov et al., NIPS 2013)

$$\text{king} = [0.23, 1.3, -0.3, 0.43]$$



composing embeddings

- neural networks **compose** word embeddings into vectors for phrases, sentences, and documents

neural students opened their
network (  ) = 

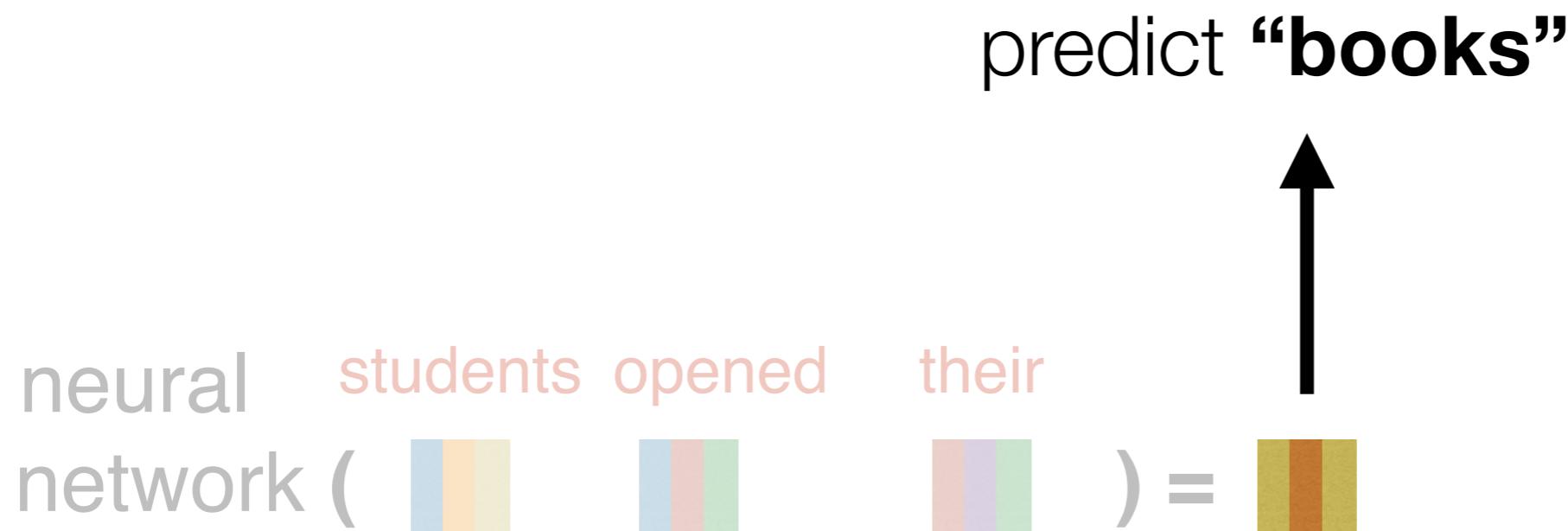
Predict the next word from composed prefix representation

predict “**books**”

↑

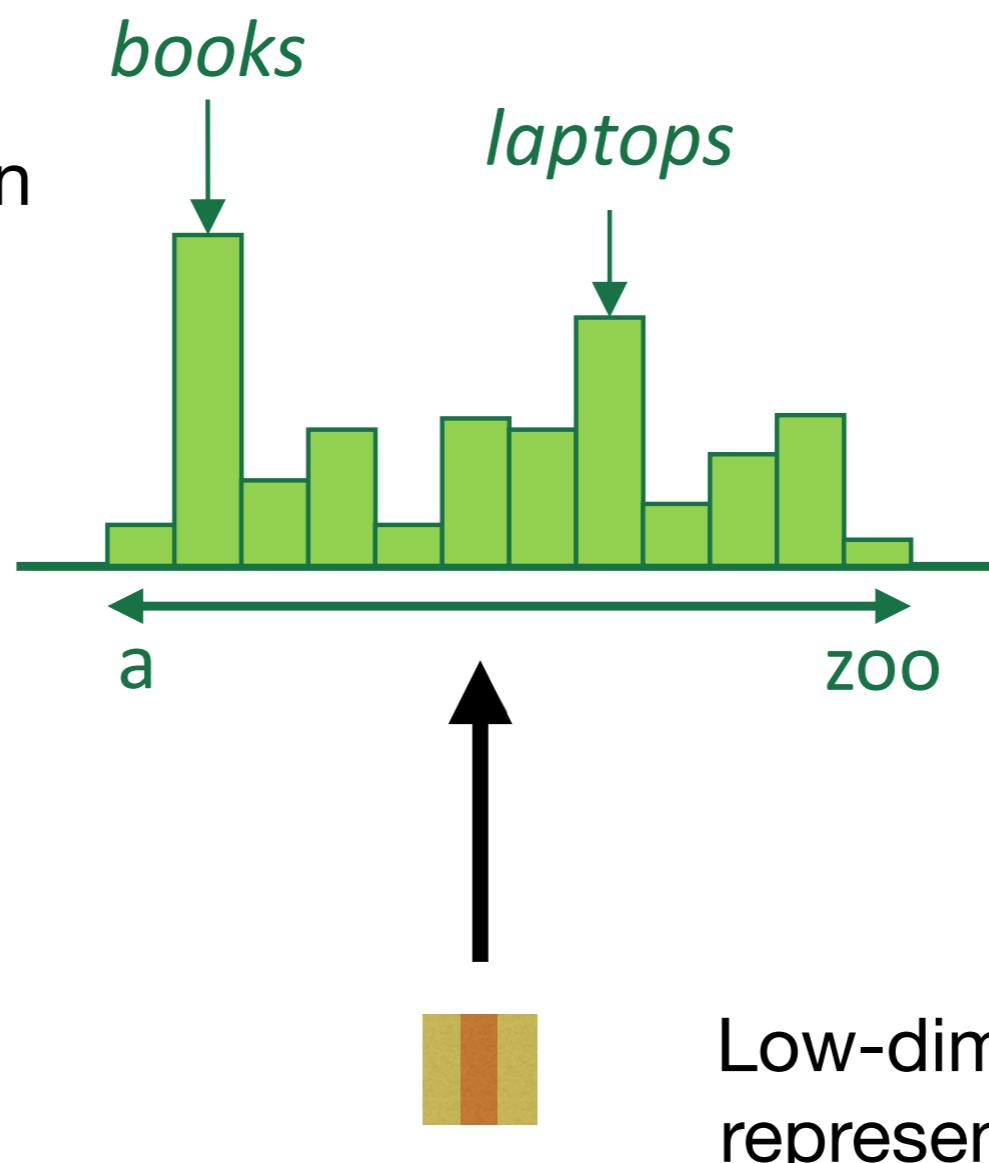
neural students opened their
network ([blue|orange|yellow] [blue|red|green] [brown|purple|green]) = [yellow|orange|yellow]

How does this happen? Let's work our way backwards, starting with the prediction of the next word



$P(w_i | \text{vector for "students opened their"})$

Probability distribution
over the entire
vocabulary



Low-dimensional
representation of
“students opened their”

Let's say our output vocabulary consists of just four words: "books", "houses", "lamps", and "stamps".



Low-dimensional representation of
"students opened their"

Let's say our output vocabulary consists of just four words: "books", "houses", "lamps", and "stamps".

books houses lamps stamps
 $<0.6, 0.2, 0.1, 0.1>$

We want to get a probability distribution over these four words



Low-dimensional representation of "students opened their"

$x = <-2.3, 0.9, 5.4>$



Here's an example 3-d
prefix vector

W is a *weight matrix*. It contains *parameters* that we can *update* to control the final probability distribution of the next word

$$\mathbf{w} = \begin{Bmatrix} 1.2, -0.3, 0.9 \\ 0.2, 0.4, -2.2 \\ 8.9, -1.9, 6.5 \\ 4.5, 2.2, -0.1 \end{Bmatrix}$$

$$\mathbf{x} = <-2.3, 0.9, 5.4>$$



Here's an example 3-d prefix vector

$$\mathbf{w} = \begin{Bmatrix} 1.2, -0.3, 0.9 \\ 0.2, 0.4, -2.2 \\ 8.9, -1.9, 6.5 \\ 4.5, 2.2, -0.1 \end{Bmatrix}$$

W is a *weight matrix*. It contains *parameters* that we can *update* to control the final probability distribution of the next word

first, we'll project our
3-d prefix
representation to 4-d
with a matrix-vector
product

$$\mathbf{x} = <-2.3, 0.9, 5.4>$$



Here's an example 3-d
prefix vector

$$\mathbf{w} = \begin{Bmatrix} 1.2, & -0.3, & 0.9 \\ 0.2, & 0.4, & -2.2 \\ 8.9, & -1.9, & 6.5 \\ 4.5, & 2.2, & -0.1 \end{Bmatrix}$$

$$\mathbf{x} = \langle -2.3, 0.9, 5.4 \rangle$$

intuition: each dimension of \mathbf{x} corresponds to a *feature* of the prefix

intuition: each row
of \mathbf{W} contains
feature weights for a
corresponding word
in the vocabulary

$$\mathbf{w} = \left\{ \begin{array}{l} 1.2, -0.3, 0.9 \\ 0.2, 0.4, -2.2 \\ 8.9, -1.9, 6.5 \\ 4.5, 2.2, -0.1 \end{array} \right\}$$

$$\mathbf{x} = <-2.3, 0.9, 5.4>$$

intuition: each
dimension of \mathbf{x}
corresponds to a
feature of the prefix

intuition: each row of \mathbf{W} contains *feature weights* for a corresponding word in the vocabulary

$$\mathbf{w} = \left\{ \begin{array}{l} 1.2, -0.3, 0.9 \\ 0.2, 0.4, -2.2 \\ 8.9, -1.9, 6.5 \\ 4.5, 2.2, -0.1 \end{array} \right\} \begin{matrix} \text{books} \\ \text{houses} \\ \text{lamps} \\ \text{stamps} \end{matrix}$$

$$\mathbf{x} = <-2.3, 0.9, 5.4>$$

intuition: each dimension of \mathbf{x} corresponds to a *feature* of the prefix

intuition: each row of \mathbf{W} contains *feature weights* for a corresponding word in the vocabulary

CAUTION: we can't easily *interpret* these features.

$$\mathbf{w} = \left\{ \begin{array}{l} 1.2, -0.3, 0.9 \\ 0.2, 0.4, -2.2 \\ 8.9, -1.9, 6.5 \\ 4.5, 2.2, -0.1 \end{array} \right\} \begin{matrix} \text{books} \\ \text{houses} \\ \text{lamps} \\ \text{stamps} \end{matrix}$$

$$\mathbf{x} = <-2.3, 0.9, 5.4>$$

intuition: each dimension of \mathbf{x} corresponds to a *feature* of the prefix

$$\mathbf{Wx} = \langle 1.8, -11.9, 12.9, -8.9 \rangle$$

How did we compute this? It's just the dot product of each row of \mathbf{W} with \mathbf{x} !

$$\mathbf{W} = \left\{ \begin{array}{l} \begin{matrix} 1.2, & -0.3, & 0.9 \\ 0.2, & 0.4, & -2.2 \\ 8.9, & -1.9, & 6.5 \\ 4.5, & 2.2, & -0.1 \end{matrix} \end{array} \right\}$$

$$\mathbf{x} = \langle -2.3, 0.9, 5.4 \rangle$$

$$\mathbf{Wx} = \langle 1.8, -11.9, 12.9, -8.9 \rangle$$

How did we compute this? Just the dot product of each row of \mathbf{W} with \mathbf{x} !

$$\mathbf{W} = \left\{ \begin{array}{l} \begin{matrix} 1.2, & -0.3, & 0.9 \\ 0.2, & 0.4, & -2.2 \\ 8.9, & -1.9, & 6.5 \\ 4.5, & 2.2, & -0.1 \end{matrix} \end{array} \right\}$$
$$\mathbf{x} = \langle -2.3, 0.9, 5.4 \rangle$$
$$1.2 * -2.3 + -0.3 * 0.9 + 0.9 * 5.4$$

Okay, so how do we go from this 4-d vector to a probability distribution?

$$\mathbf{Wx} = \langle 1.8, -11.9, 12.9, -8.9 \rangle$$

We'll use the softmax function!

$$\text{softmax}(x) = \frac{e^x}{\sum_j e^{x_j}}$$

- x is a vector
- x_j is dimension j of x
- each dimension j of the softmaxed output represents the probability of class j

$$\mathbf{Wx} = <1.8, -1.9, 2.9, -0.9>$$

$$\text{softmax}(\mathbf{Wx}) = <0.24, 0.006, 0.73, 0.02>$$

We'll use the softmax function!

$$\text{softmax}(x) = \frac{e^x}{\sum_j e^{x_j}}$$

- x is a vector
- x_j is dimension j of x
- each dimension j of the softmaxed output represents the probability of class j

$$\mathbf{Wx} = <1.8, -1.9, 2.9, -0.9>$$

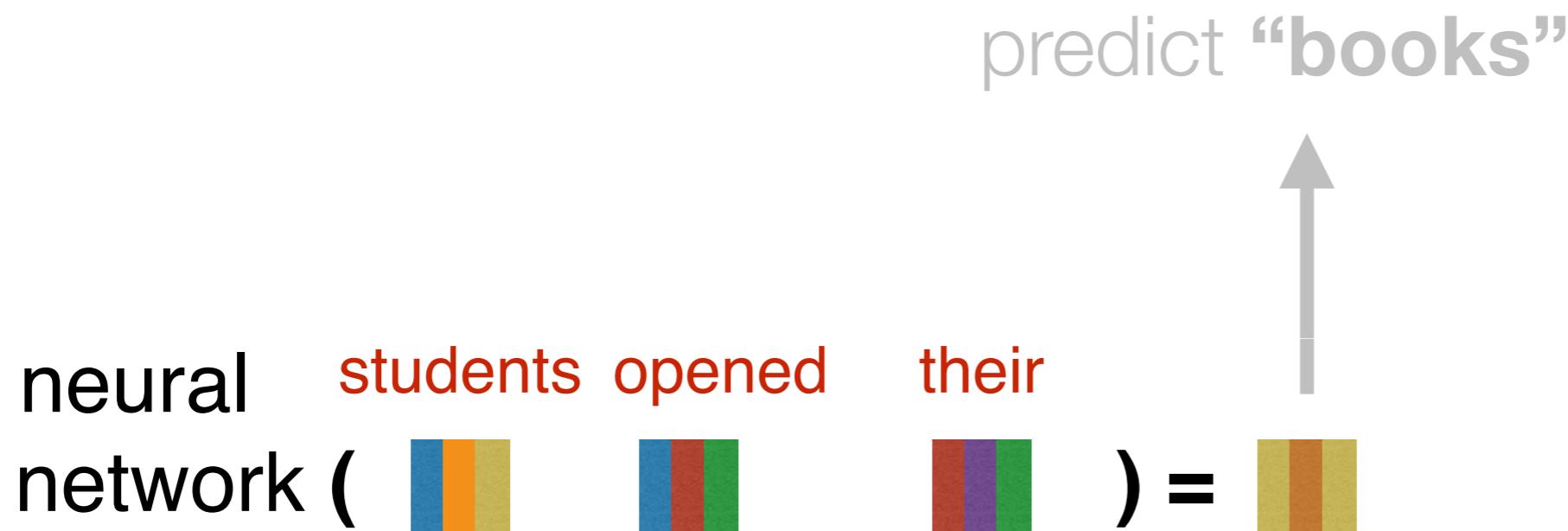
$$\text{softmax}(\mathbf{Wx}) = <0.24, 0.006, 0.73, 0.02>$$

books houses lamps stamps

so to sum up...

- Given a d -dimensional vector representation \mathbf{x} of a prefix, we do the following to predict the next word:
 1. Project it to a V -dimensional vector using a matrix-vector product (a.k.a. a “linear layer”, or a “feedforward layer”), where V is the size of the vocabulary
 2. Apply the softmax function to transform the resulting vector into a probability distribution

Now that we know how to predict “**books**”, let’s focus on how to compute the prefix representation \mathbf{x} in the first place!



Composition functions

input: sequence of word embeddings corresponding to the tokens of a given prefix

output: single vector

Composition functions

input: sequence of word embeddings corresponding to the tokens of a given prefix

output: single vector

- Element-wise functions
 - e.g., just sum up all of the word embeddings!
- Concatenation
- Feed-forward neural networks
- Convolutional neural networks
- Recurrent neural networks
- Transformers (our focus this semester)

Let's look first at *concatenation*, an easy to understand but limited composition function

A fixed-window neural Language Model

as ~~the proctor started the clock~~ the students opened their _____

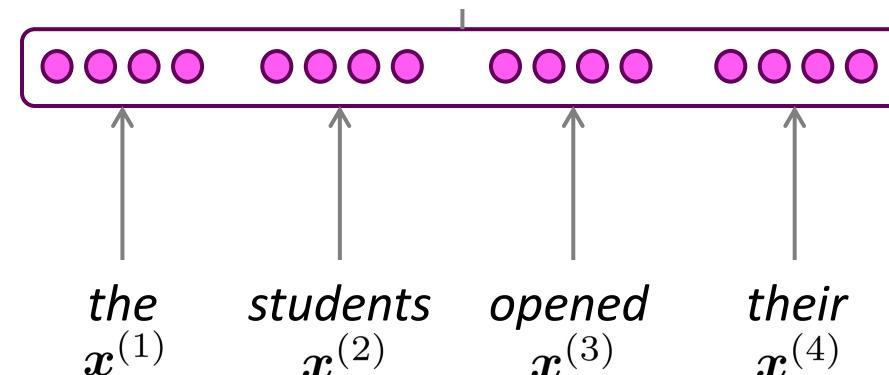
discard

fixed window

A fixed-window neural Language Model

concatenated word embeddings
 $e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$

words / one-hot vectors
 $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$



“Shallow” Encoding

- Simplest encoding approach: **encoder is just an embedding-lookup**

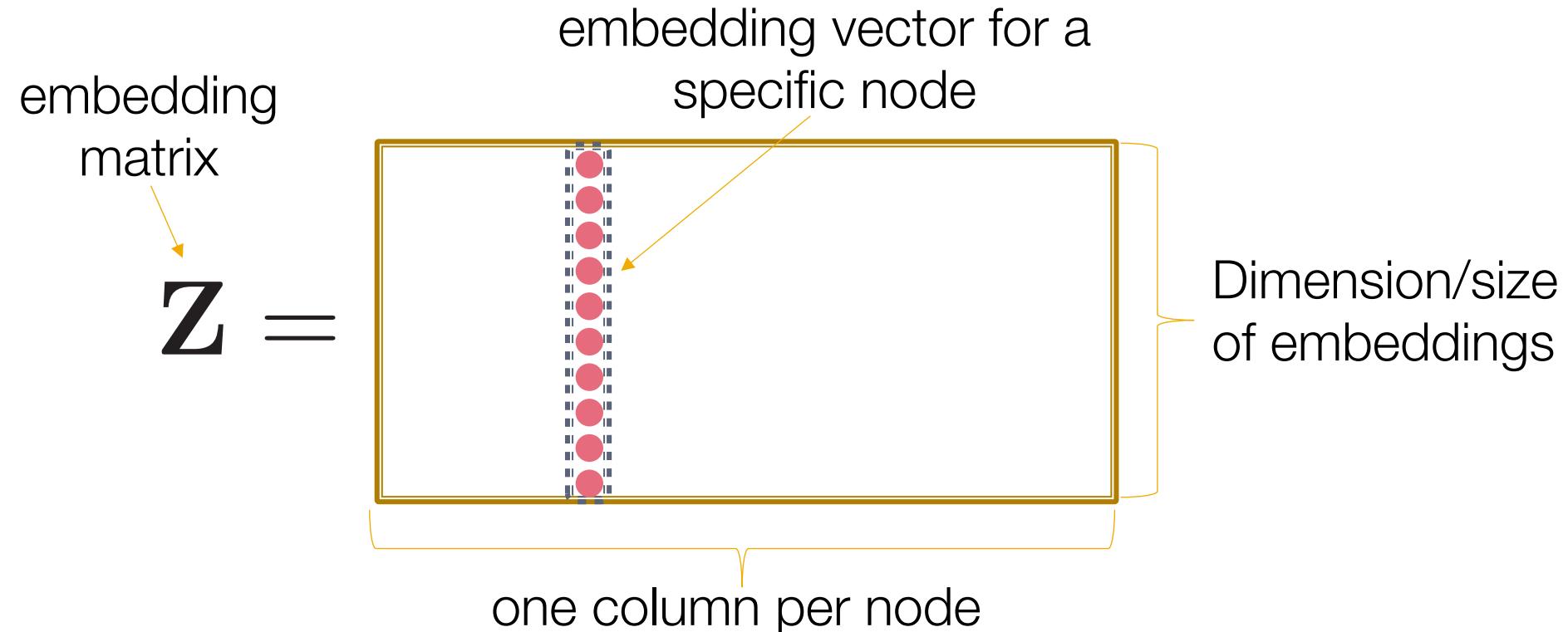
$$\text{ENC}(v) = \mathbf{Z}\mathbf{v}$$

$\mathbf{Z} \in \mathbb{R}^{d \times |\mathcal{V}|}$ Matrix, each column is d -dim node embedding [what we learn!]

$\mathbf{v} \in \mathbb{I}^{|\mathcal{V}|}$ Indicator vector, all zeroes except for a “1” at the position that corresponds to node v

“Shallow” Encoding

- Simplest encoding approach: **encoder is just an embedding-lookup**



A fixed-window neural Language Model

hidden layer

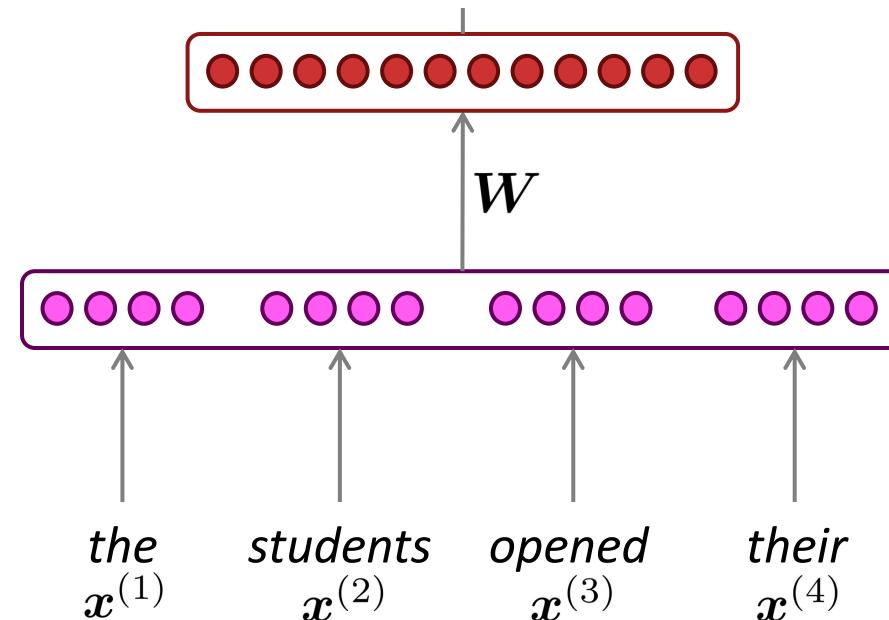
$$h = f(We + b_1)$$

concatenated word embeddings

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

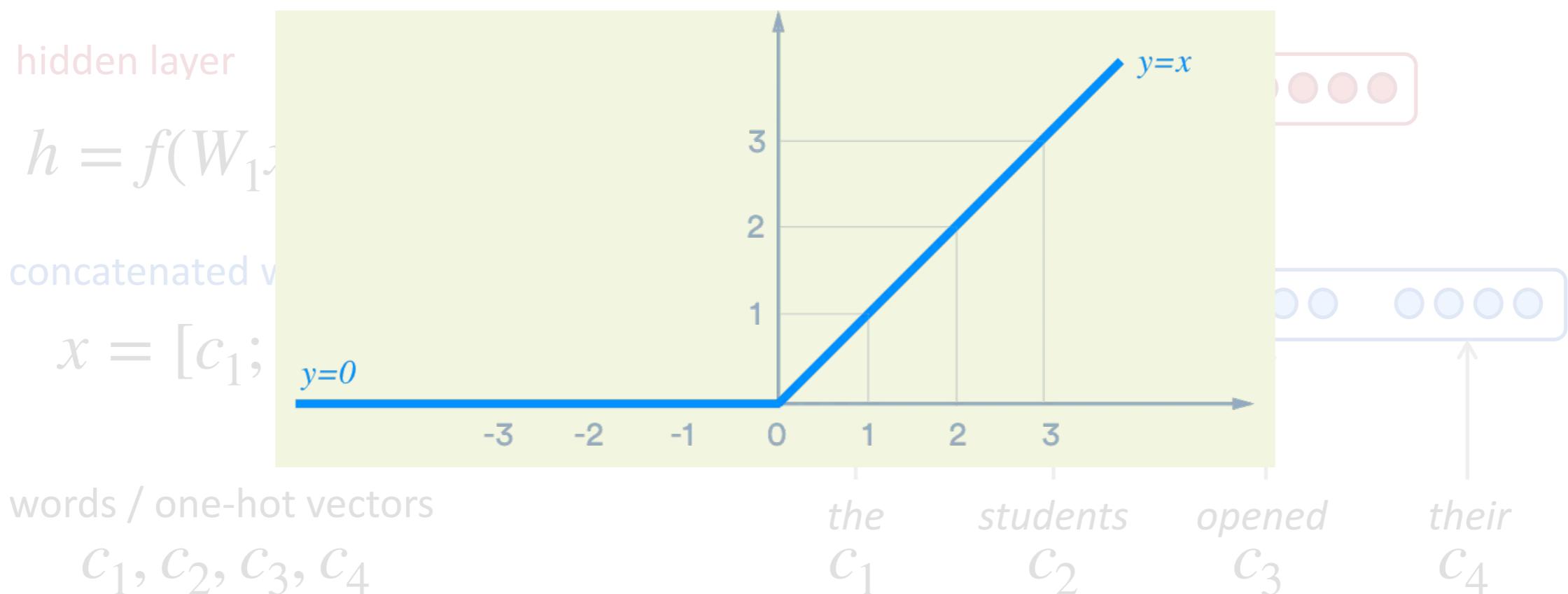
words / one-hot vectors

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



A fixed-window neural Language Model

f is a *nonlinearity*, or an element-wise nonlinear function. The most commonly-used choice today is the rectified linear unit (**ReLU**), which is just $\text{ReLU}(x) = \max(0, x)$. Other choices include **tanh** and **sigmoid**.



A fixed-window neural Language Model

output distribution

$$\hat{y} = \text{softmax}(U\mathbf{h} + \mathbf{b}_2)$$

hidden layer

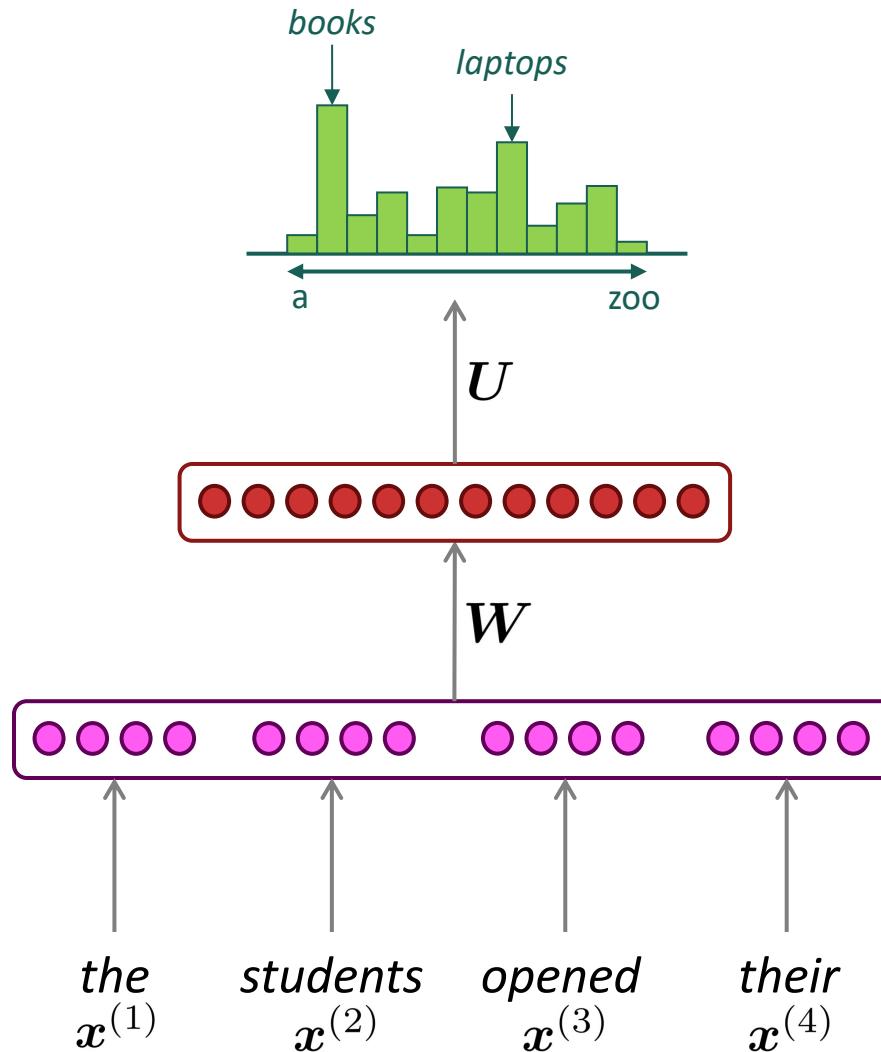
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

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

words / one-hot vectors

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



A fixed-window neural Language Model

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

Improvements over n -gram LM:

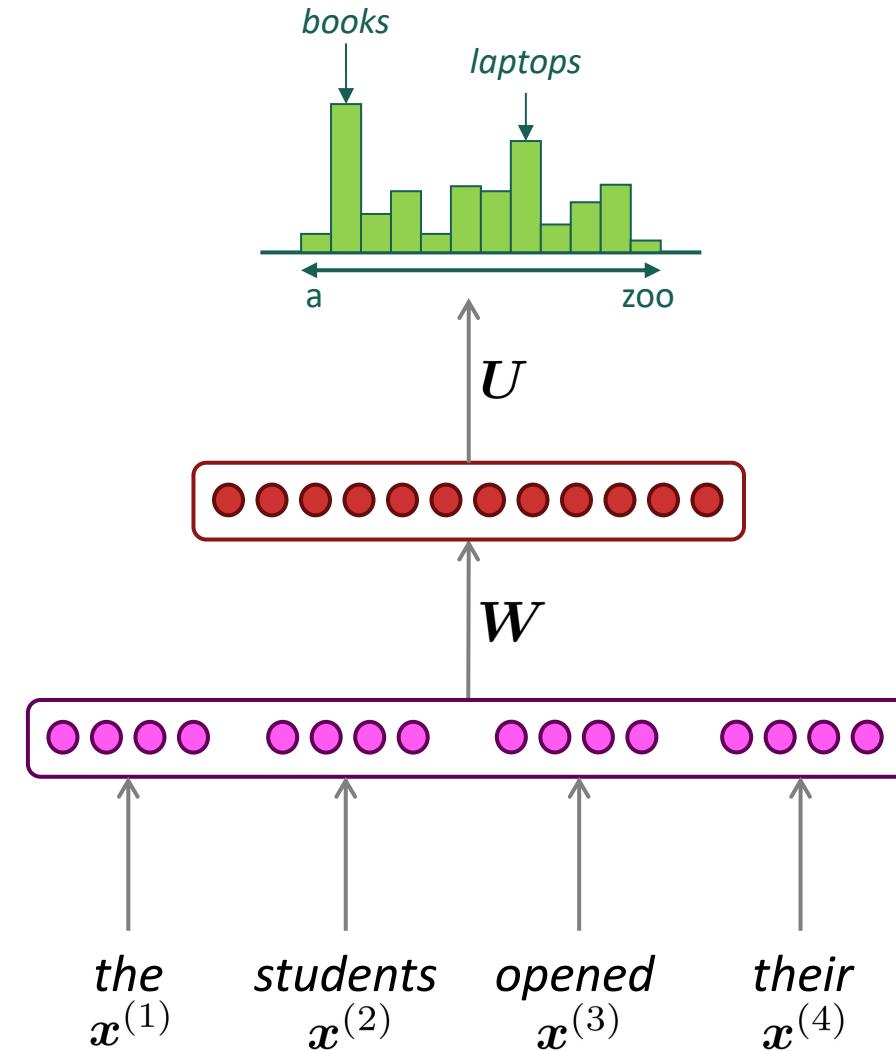
- No sparsity problem
- Don't need to store all observed n -grams

Remaining **problems**:

- Fixed window is **too small**
- Enlarging window enlarges W
- Window can never be large enough!
- $x^{(1)}$ and $x^{(2)}$ are multiplied by completely different weights in W .

No symmetry in how the inputs are processed.

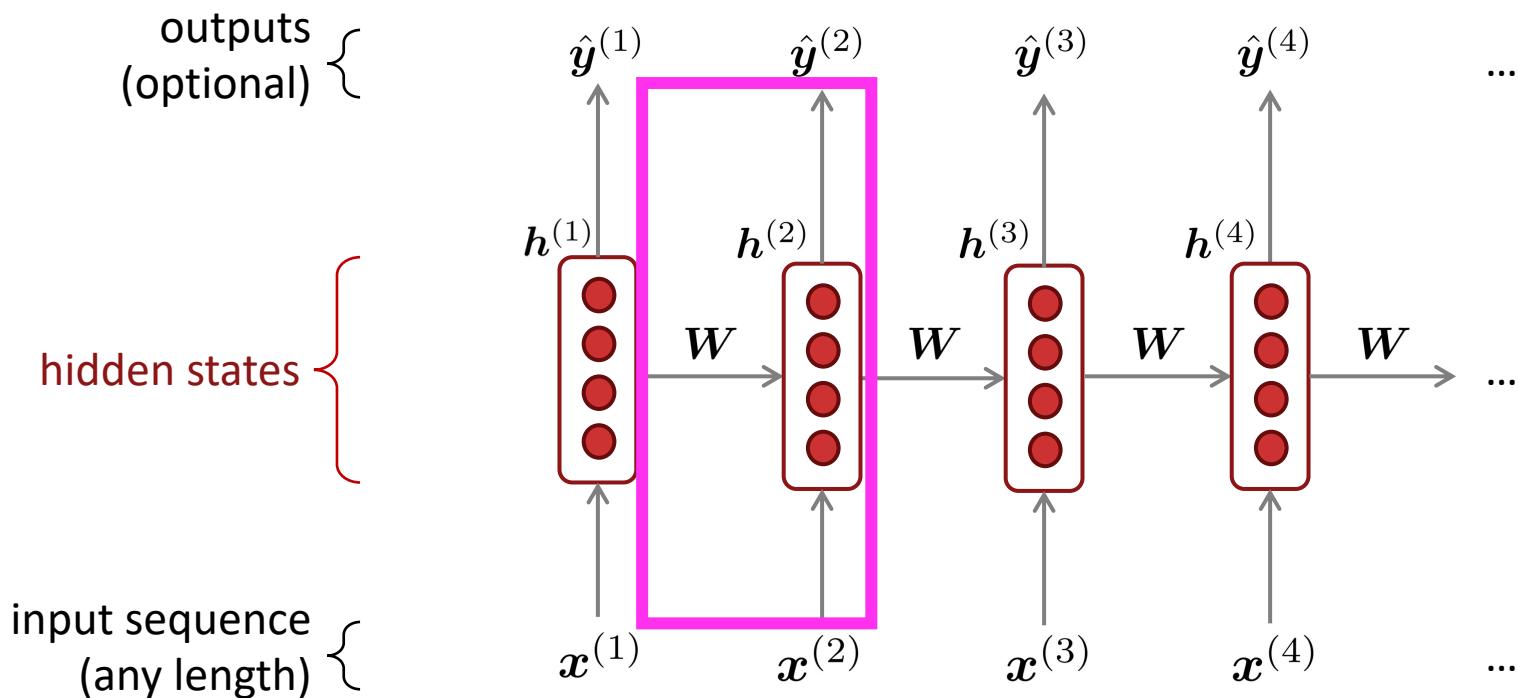
We need a neural architecture
that can process *any length input*



3. Recurrent Neural Networks (RNN)

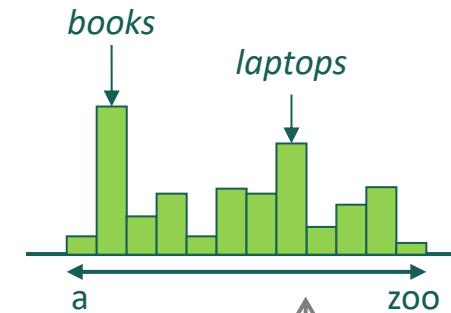
A family of neural architectures

Core idea: Apply the same weights W repeatedly



A Simple RNN Language Model

$$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$



output distribution

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h}^{(t)} + \mathbf{b}_2)$$

hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

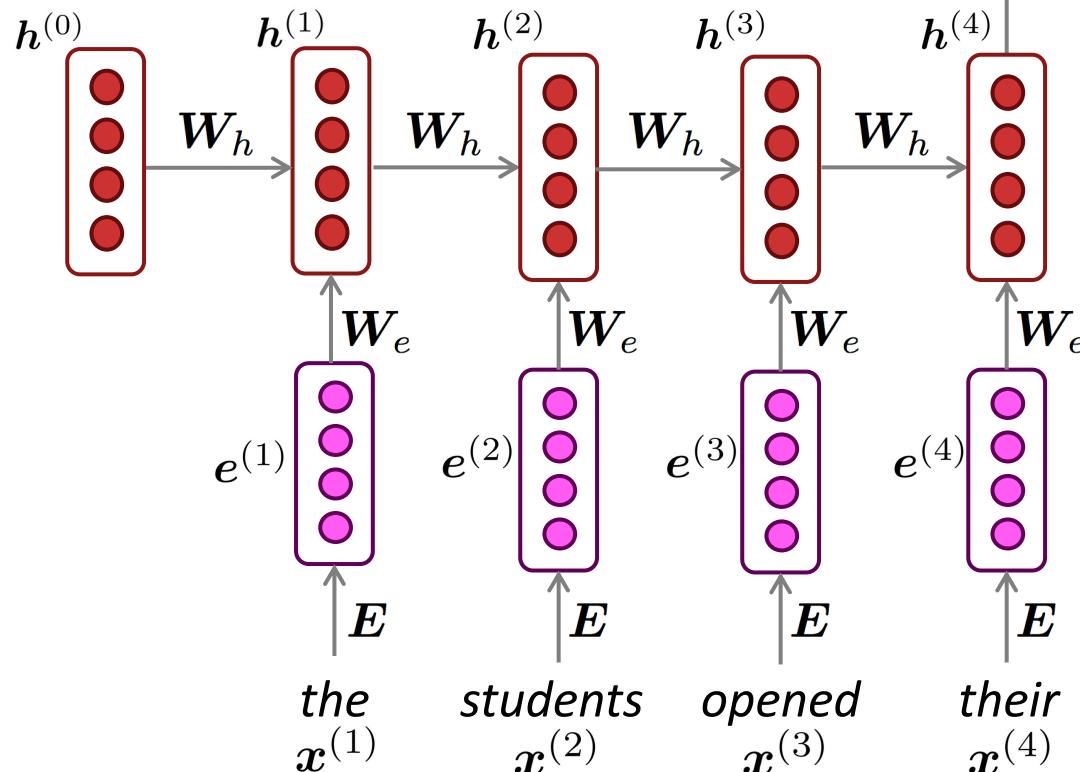
$\mathbf{h}^{(0)}$ is the initial hidden state

word embeddings

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

words / one-hot vectors

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



Note: this input sequence could be much longer now!

RNN Language Models

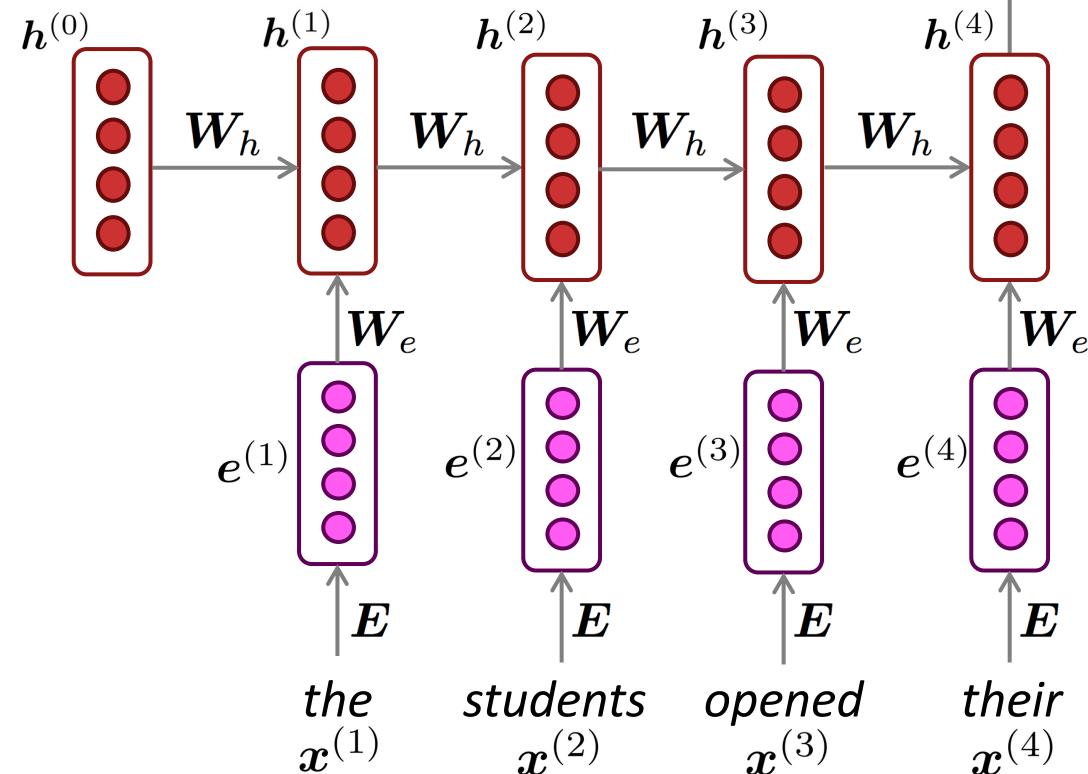
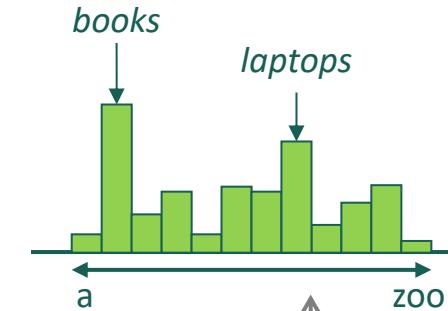
RNN Advantages:

- 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**

$$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$



Training an RNN Language Model

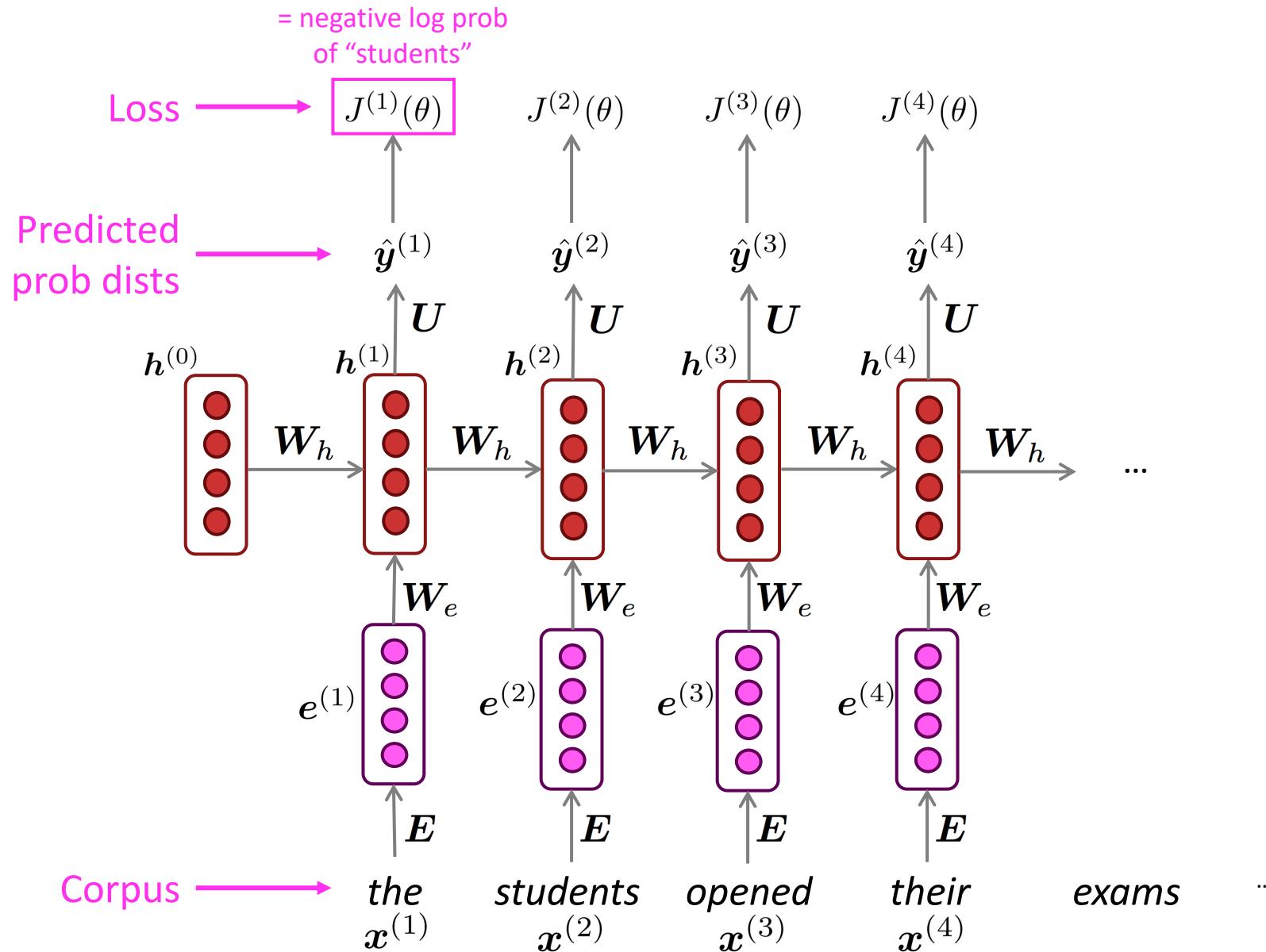
- Get a **big corpus of text** which is a sequence of words $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$
- Feed into RNN-LM; compute output distribution $\hat{\mathbf{y}}^{(t)}$ **for every step t .**
 - i.e., predict probability dist of *every word*, given words so far
- **Loss function** on step t is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

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

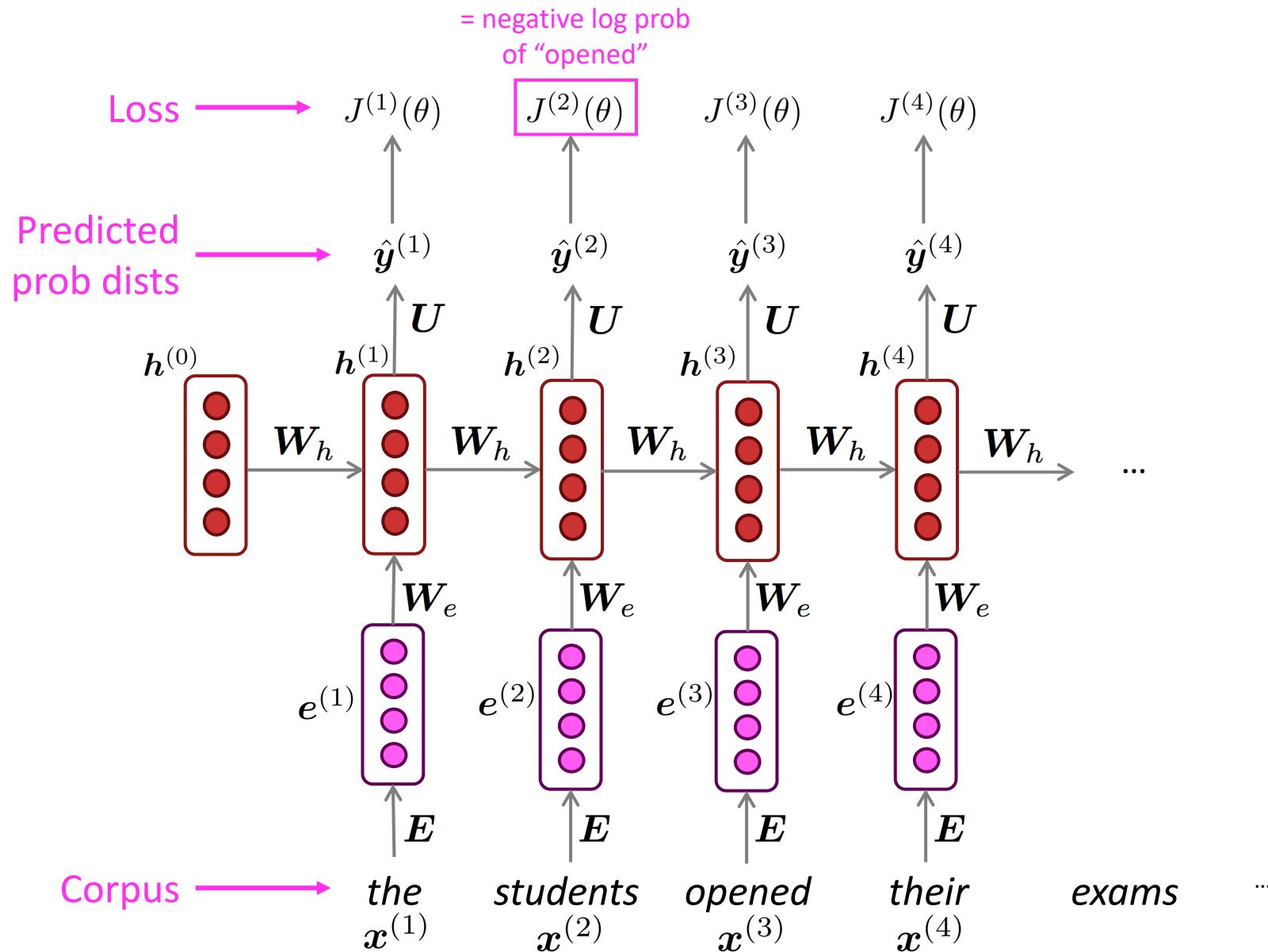
- Average this to get **overall loss** for entire training set:

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

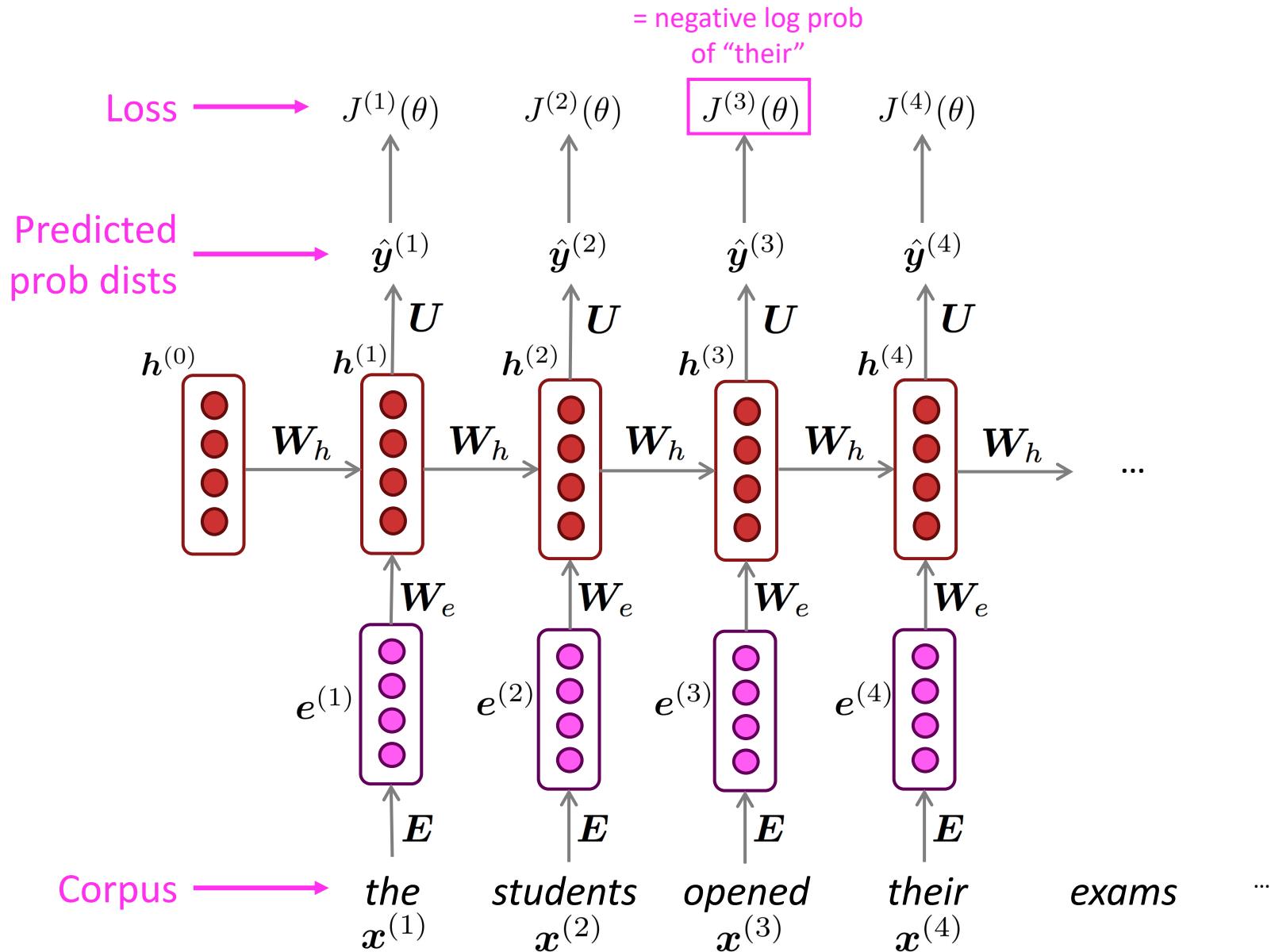
Training an RNN Language Model



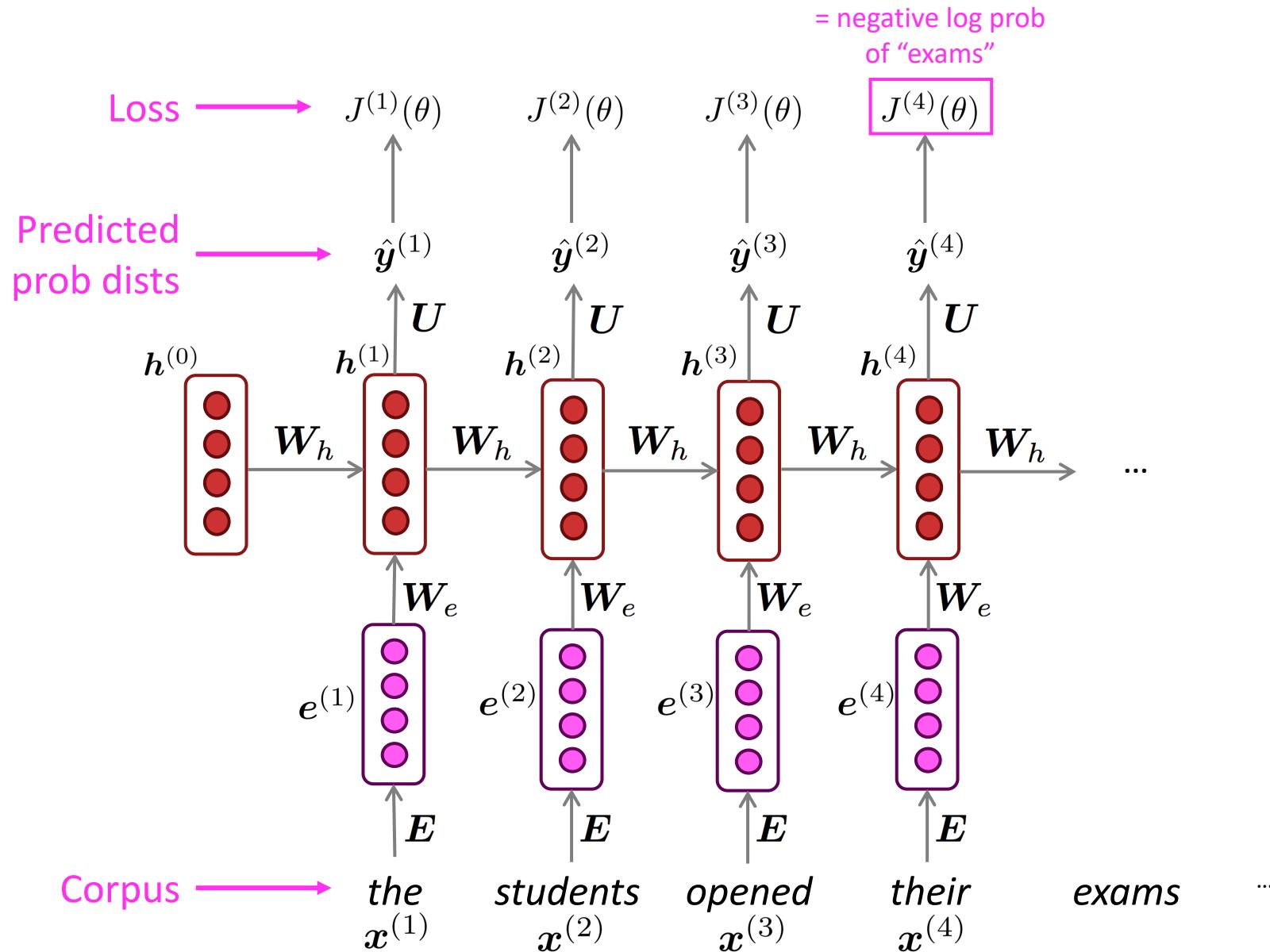
Training an RNN Language Model



Training an RNN Language Model

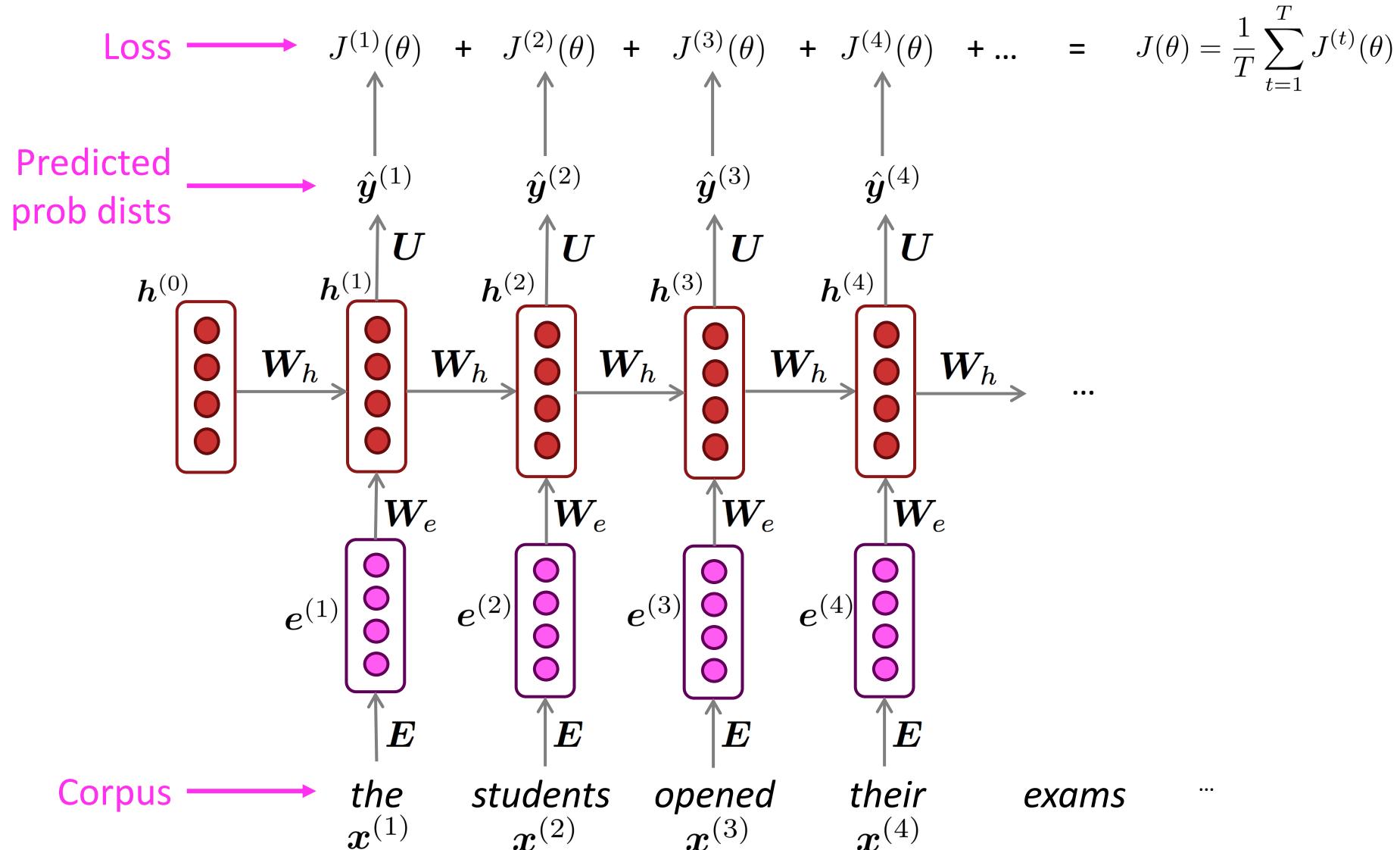


Training an RNN Language Model



Training an RNN Language Model

“Teacher forcing”



Training a RNN Language Model

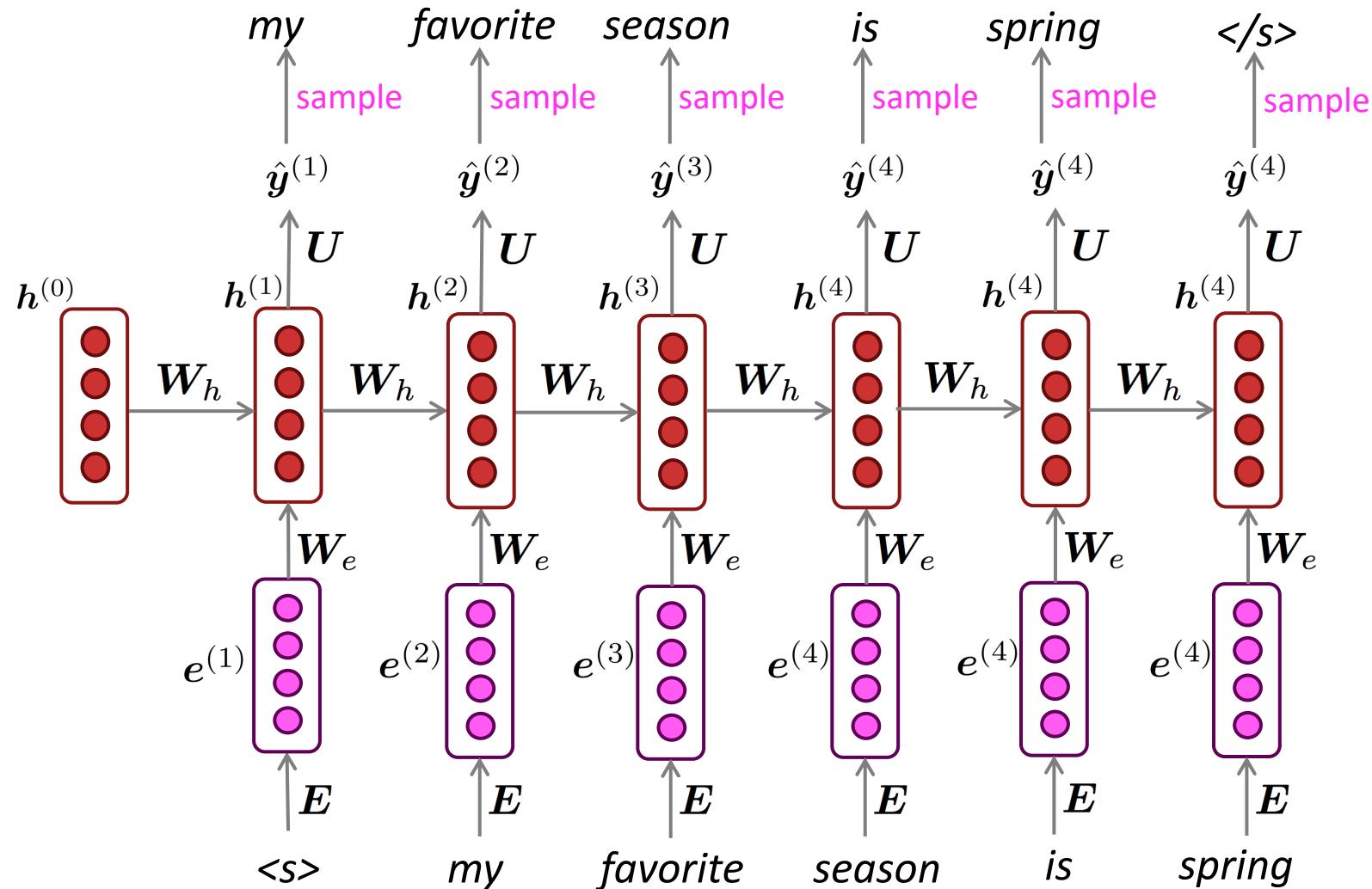
- However: Computing loss and gradients across **entire corpus** $x^{(1)}, \dots, x^{(T)}$ at once is **too expensive** (memory-wise)!

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

- In practice, consider $x^{(1)}, \dots, x^{(T)}$ as a **sentence** (or a **document**)
- Recall: **Stochastic Gradient Descent** allows us to compute loss and gradients for small chunk of data, and update.
- Compute loss $J(\theta)$ for a sentence (actually, a batch of sentences), compute gradients and update weights. Repeat on a new batch of sentences.

Generating with an RNN Language Model (“Generating roll outs”)

Just like an n-gram Language Model, you can use a RNN Language Model to generate text by **repeated sampling**. Sampled output becomes next step’s input.



Generating text with an RNN Language Model

Let's have some fun!

- You can train an RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on **Obama speeches**:



The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.

Source: <https://medium.com/@samim/obama-rnn-machine-generated-political-speeches-c8abd18a2ea0>

Generating text with an RNN Language Model

Let's have some fun!

- You can train an RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on *Harry Potter*:



“Sorry,” Harry shouted, panicking—“I’ll leave those brooms in London, are they?”

“No idea,” said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry’s shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn’t felt it seemed. He reached the teams too.

Source: <https://medium.com/deep-writing/harry-potter-written-by-artificial-intelligence-8a9431803da6>

Generating text with an RNN Language Model

Let's have some fun!

- You can train an RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on **recipes**:

Title: CHOCOLATE RANCH BARBECUE
Categories: Game, Casseroles, Cookies, Cookies
Yield: 6 Servings

2 tb Parmesan cheese -- chopped
1 c Coconut milk
3 Eggs, beaten

Place each pasta over layers of lumps. Shape mixture into the moderate oven and simmer until firm. Serve hot in bodied fresh, mustard, orange and cheese.

Combine the cheese and salt together the dough in a large skillet; add the ingredients and stir in the chocolate and pepper.



Source: <https://gist.github.com/nylki/1efbaa36635956d35bcc>

Evaluating Language Models

- The standard **evaluation metric** for Language Models is **perplexity**.

$$\text{perplexity} = \prod_{t=1}^T \left(\frac{1}{P_{\text{LM}}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} \right)^{1/T}$$



Normalized by
number of words

Inverse probability of corpus, according to Language Model

Lower perplexity is better!

RNNs greatly improved perplexity over what came before

n-gram model →

Increasingly complex RNNs ↓

Model	Perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
LSTM-2048 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
Ours small (LSTM-2048)	43.9
Ours large (2-layer LSTM-2048)	39.8

Perplexity improves
(lower is better) ↓

Source: <https://research.fb.com/building-an-efficient-neural-language-model-over-a-billion-words/>

5. Recap

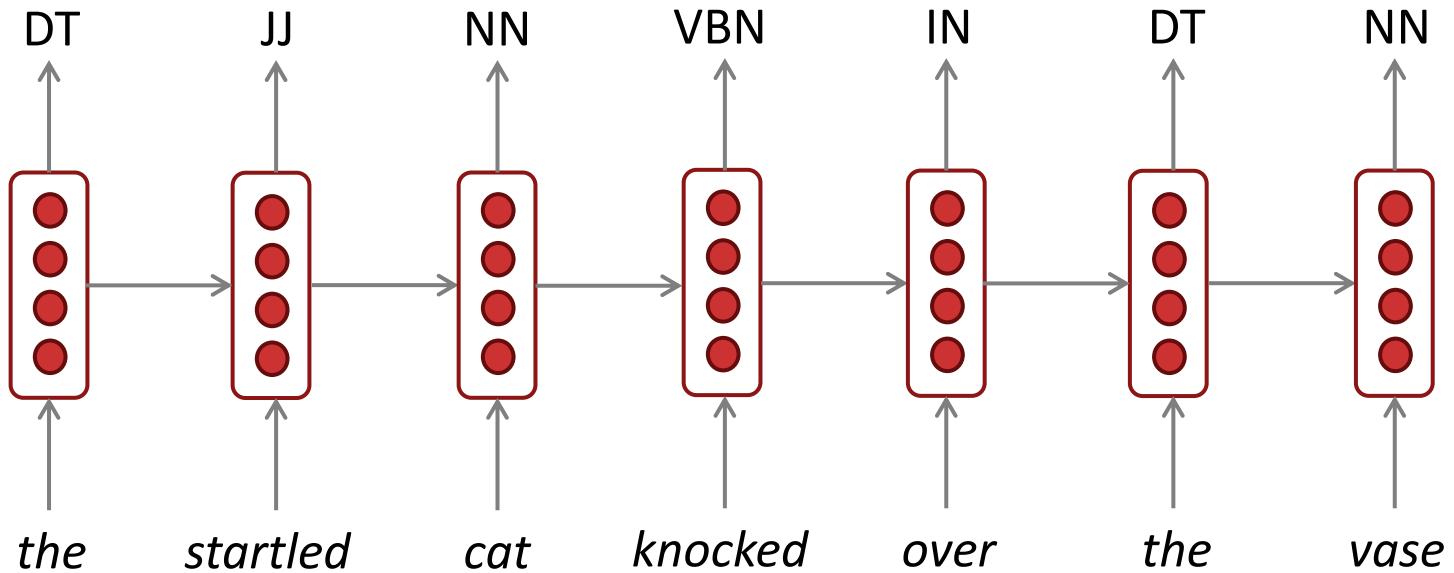
- **Language Model**: A system that predicts the next word
- **Recurrent Neural Network**: A family of neural networks that:
 - Take sequential input of any length
 - Apply the same weights on each step
 - Can optionally produce output on each step
- Recurrent Neural Network \neq Language Model
- We've shown that RNNs are a great way to build a LM (despite some problems)
- RNNs are also useful for much more!

Why should we care about Language Modeling?

- Language Modeling is a **benchmark task** that helps us measure our progress on predicting language use
- Language Modeling is a **subcomponent** of many NLP tasks, especially those involving **generating text** or **estimating the probability of text**:
 - Predictive typing
 - Speech recognition
 - Handwriting recognition
 - Spelling/grammar correction
 - Authorship identification
 - Machine translation
 - Summarization
 - Dialogue
 - etc.
- Everything else in NLP has now been rebuilt upon Language Modeling: **GPT-3 is an LM!**

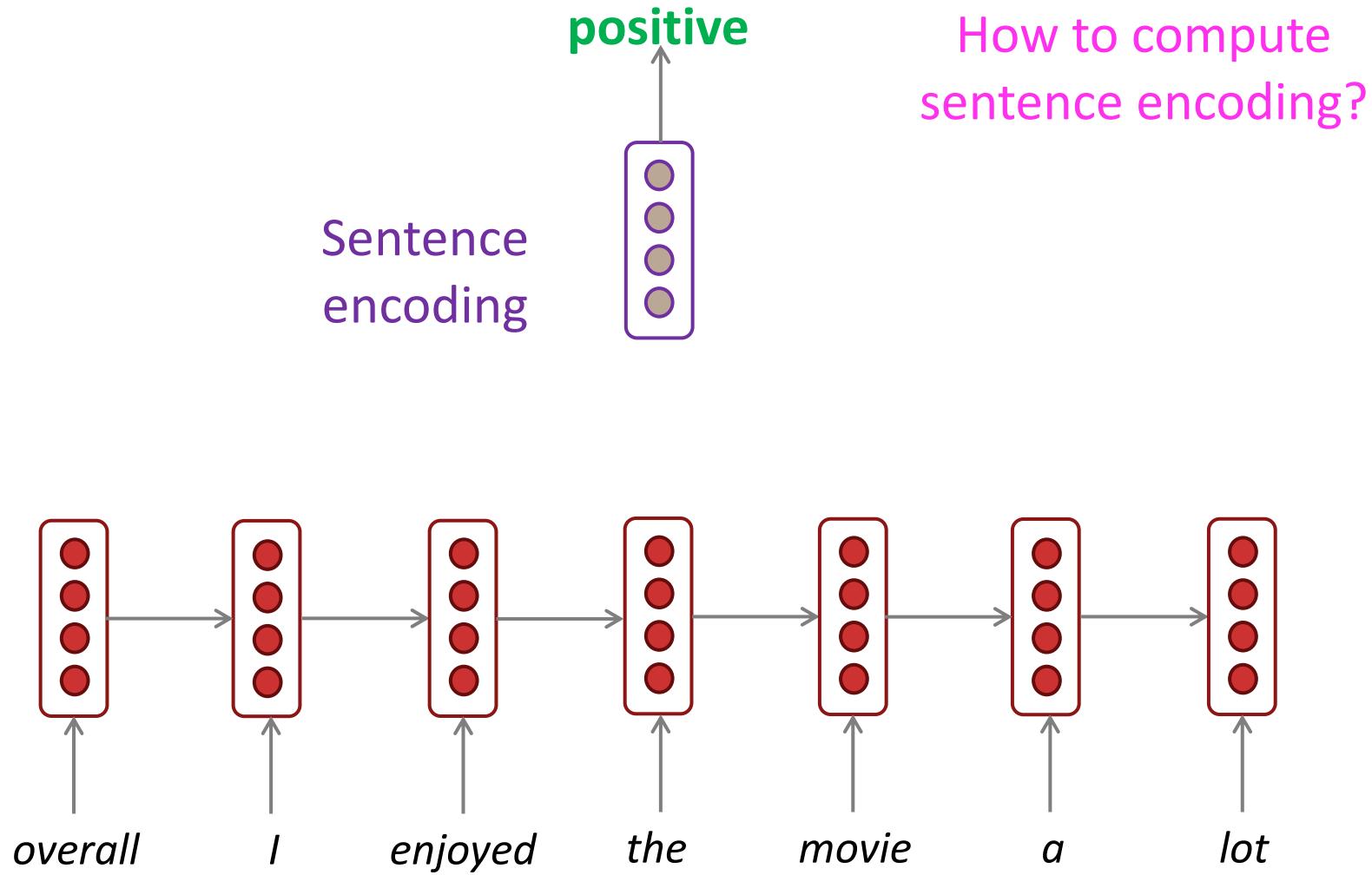
Other RNN uses: RNNs can be used for sequence tagging

e.g., part-of-speech tagging, named entity recognition



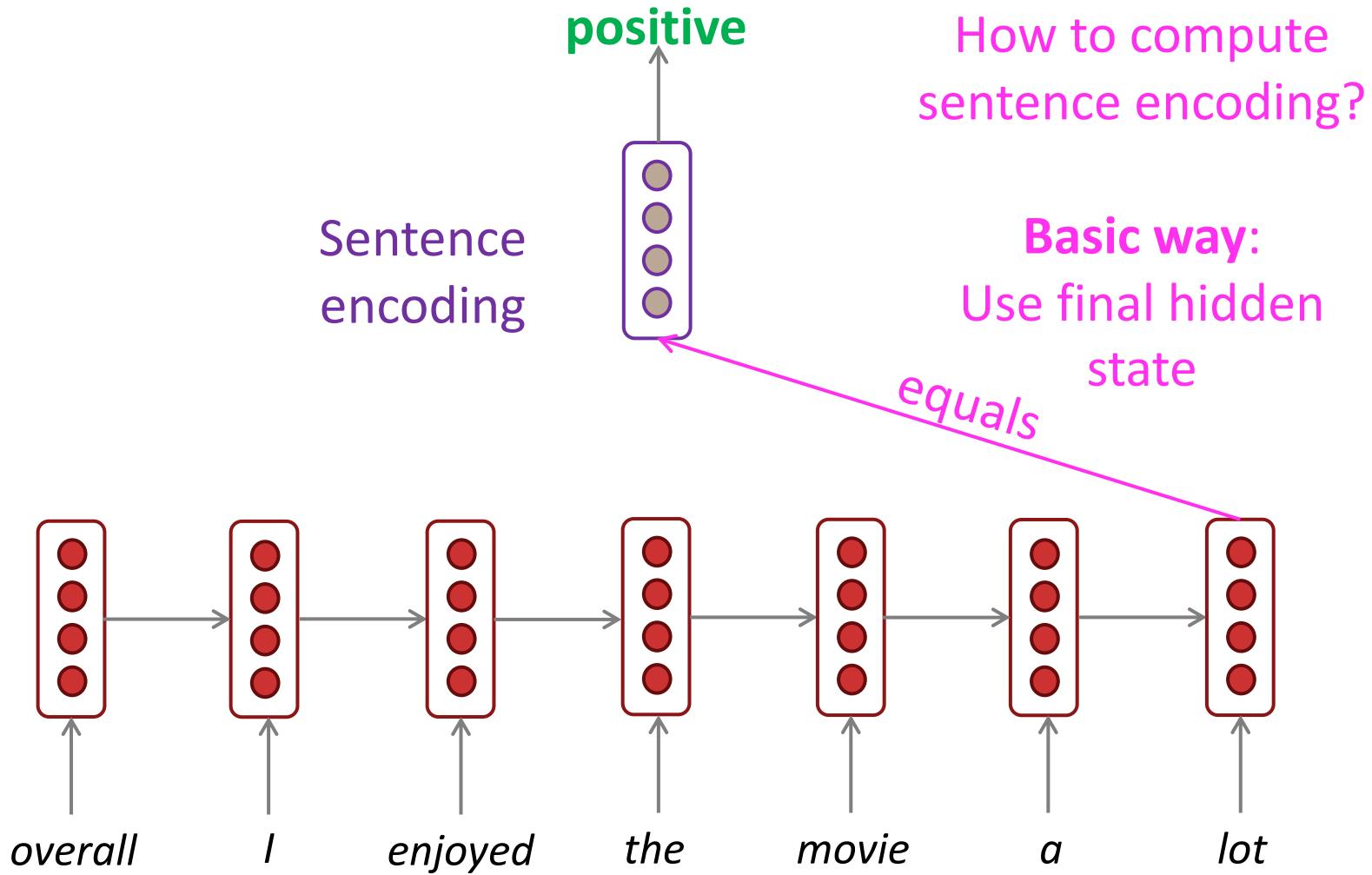
RNNs can be used for sentence classification

e.g., sentiment classification



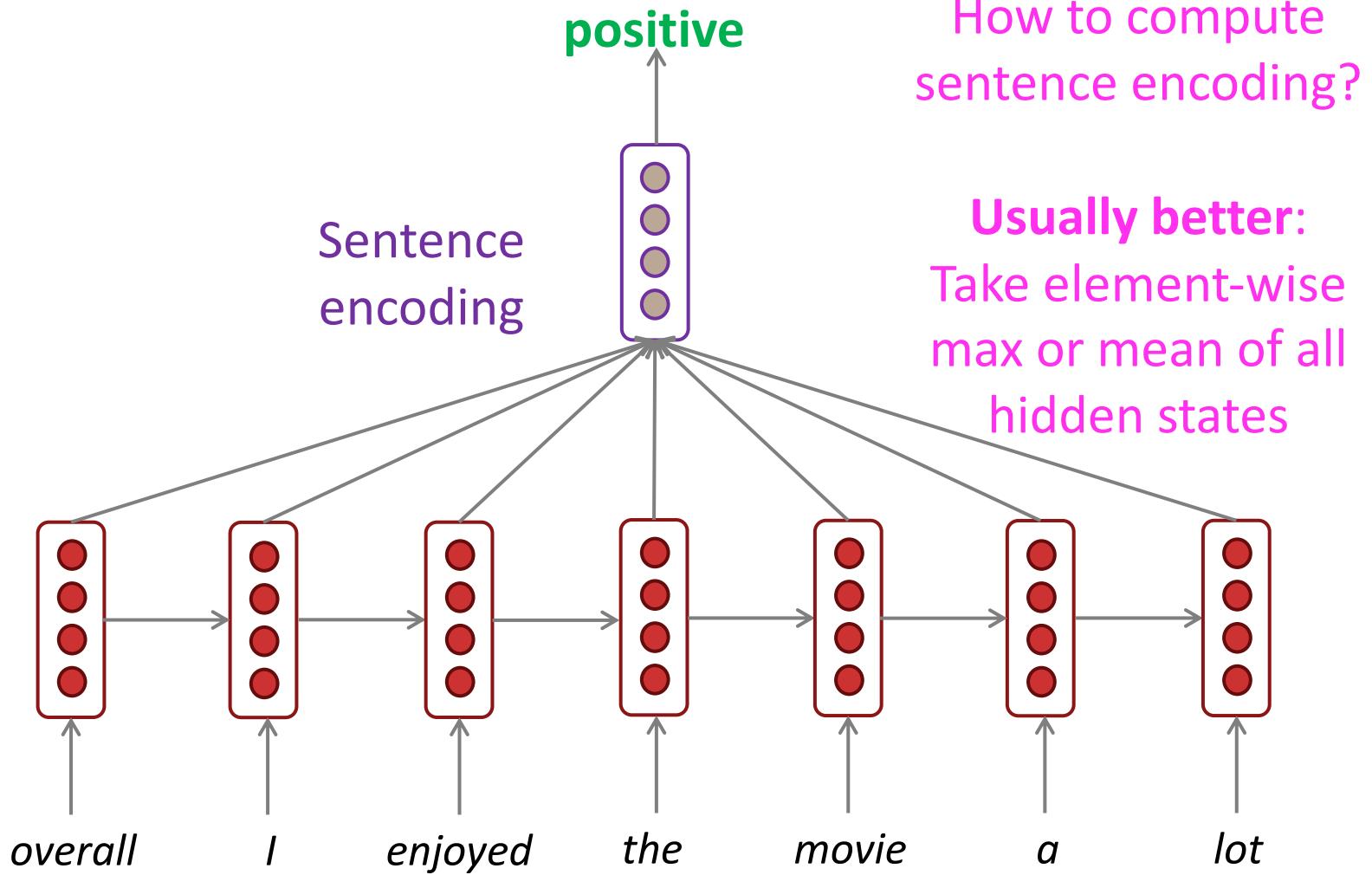
RNNs can be used for sentence classification

e.g., sentiment classification



RNNs can be used for sentence classification

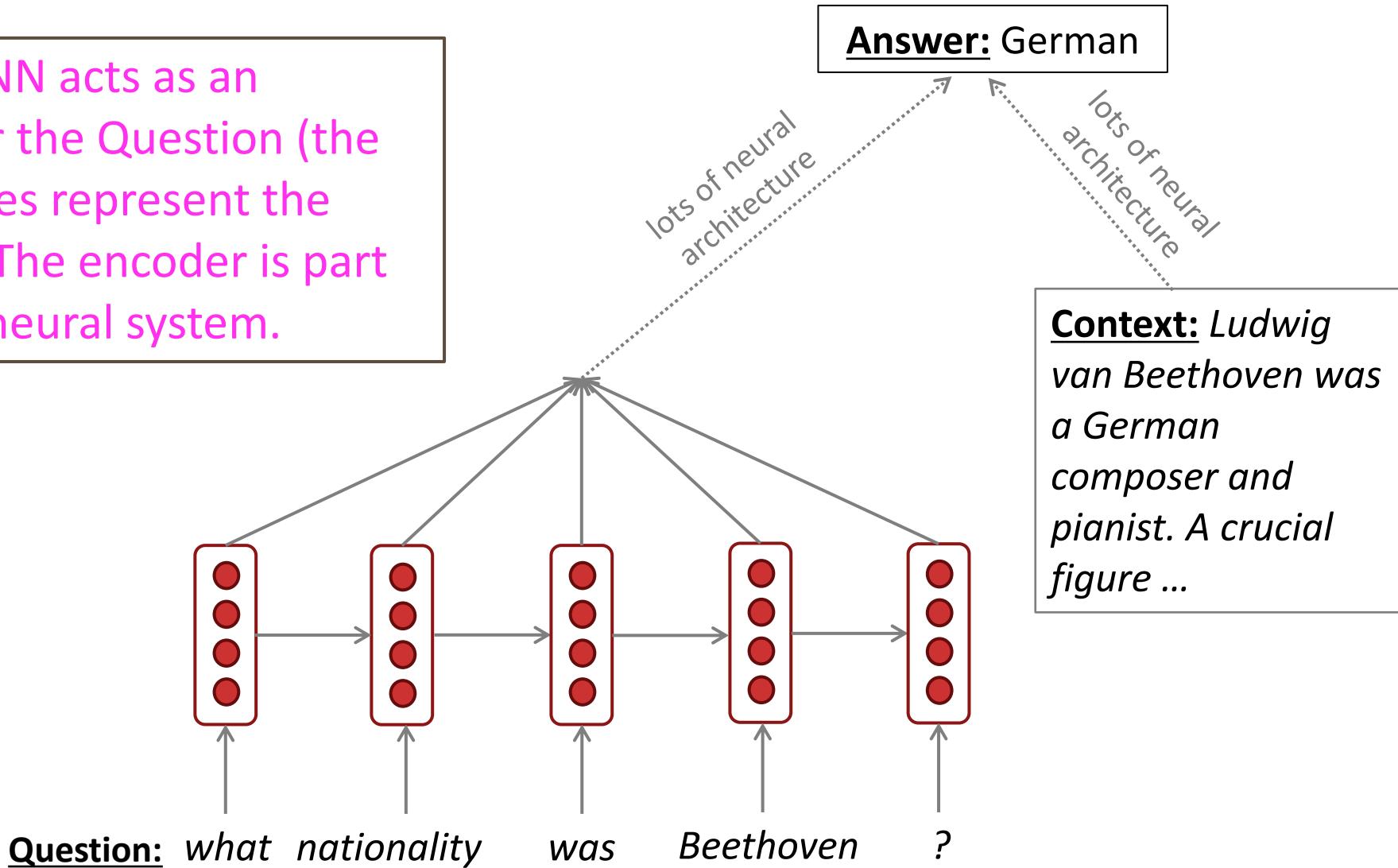
e.g., sentiment classification



RNNs can be used as an encoder module

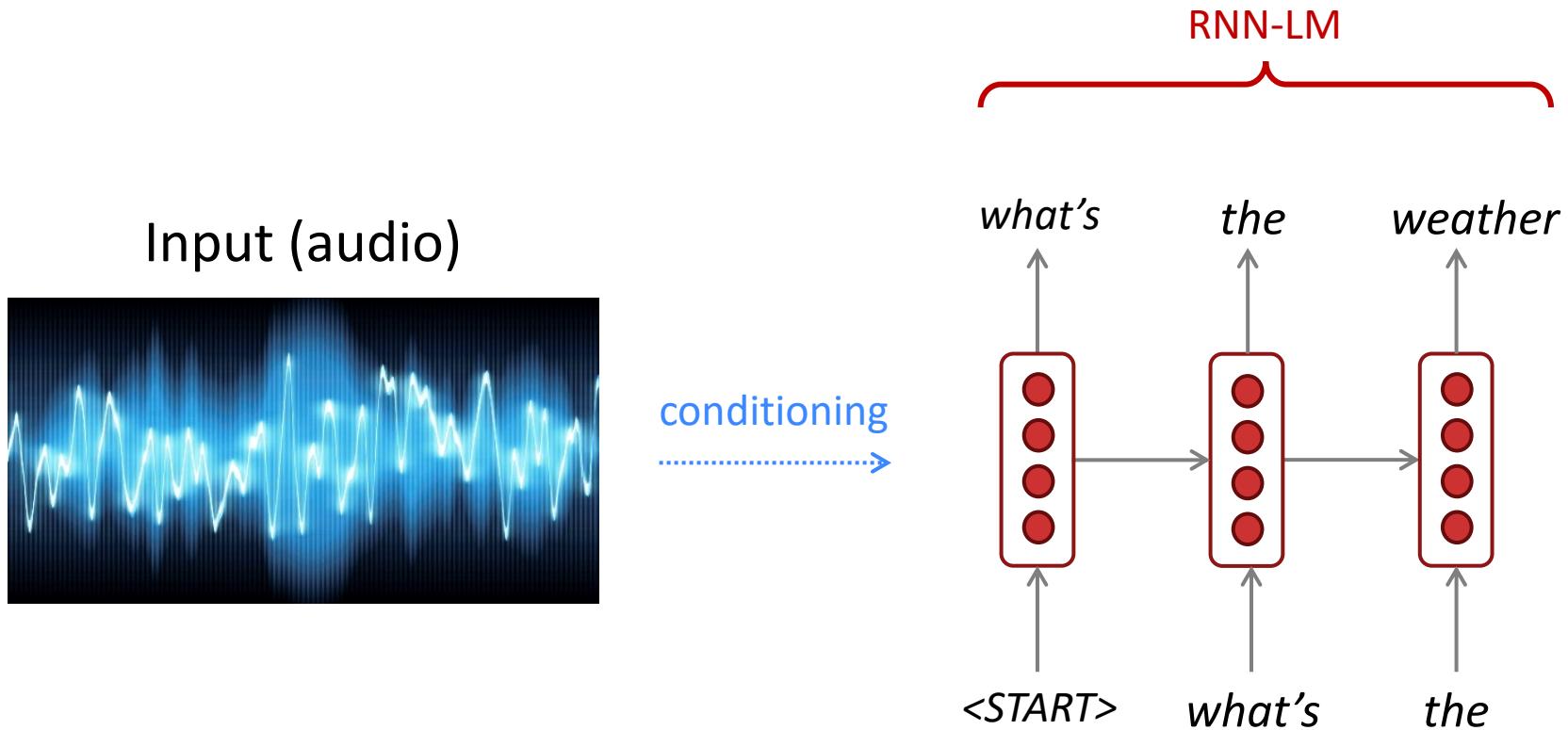
e.g., question answering, machine translation, *many other tasks!*

Here the RNN acts as an **encoder** for the Question (the hidden states represent the Question). The encoder is part of a larger neural system.



RNN-LMs can be used to generate text

e.g., speech recognition, machine translation, summarization



This is an example of a *conditional language model*.
Another popular example for conditional model: Machine Translation