

---

CS 333:  
Natural Language  
Processing

---

Fall 2023

Prof. Carolyn Anderson  
Wellesley College

## Computer Science Colloquium Series | Fall 2023

# Supporting Responsible AI Practices in Public Sector Contexts

Anna is a third year PhD student at Carnegie Mellon's Human-Computer Interaction Institute. Her research focuses on improving the design, evaluation, and governance of AI technologies used to inform complex, consequential decisions in real-world organizations. In addition to her research, she will share prior experiences forming collaborations with public sector agencies, doing research internships with industry groups, travelling to conferences, and mentoring undergraduate students. The session will end with an open Q/A discussion on applying to and doing a PhD in Computer Science / Human-Computer Interaction and other topics.

Accessibility and Disability Resources:  
[accessibility@wellesley.edu](mailto:accessibility@wellesley.edu)



***Anna Kawakami '21***

Nov 2nd, 12:45-2:00 | SCI H401

Lunch will be served

Questions??? [eni.mustafara@wellesley.edu](mailto:eni.mustafara@wellesley.edu)

# November 2nd

# Neural Language Models

# language model review

- 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** or **LM**

# n-gram models

$$p(w_j | \text{students opened their}) = \frac{\text{count}(\text{students opened their } w_j)}{\text{count}(\text{students opened their})}$$

# Problems with n-gram Language Models

## Sparsity Problem 1

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

$$p(w_j | \text{students opened their}) = \frac{\text{count}(\text{students opened their } w_j)}{\text{count}(\text{students opened their})}$$

# Problems with n-gram Language Models

## Sparsity Problem 1

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

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

$$p(w_j | \text{students opened their}) = \frac{\text{count}(\text{students opened their } w_j)}{\text{count}(\text{students opened their})}$$

# Problems with n-gram Language Models

**Storage:** Need to store count for all possible  $n$ -grams. So model size is  $O(\exp(n))$ .

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

Increasing  $n$  makes model size huge!

# another issue:

students open their eyes

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

students opened their \_\_\_\_

pupils opened their \_\_\_\_

scholars opened their \_\_\_\_

undergraduates opened their \_\_\_\_

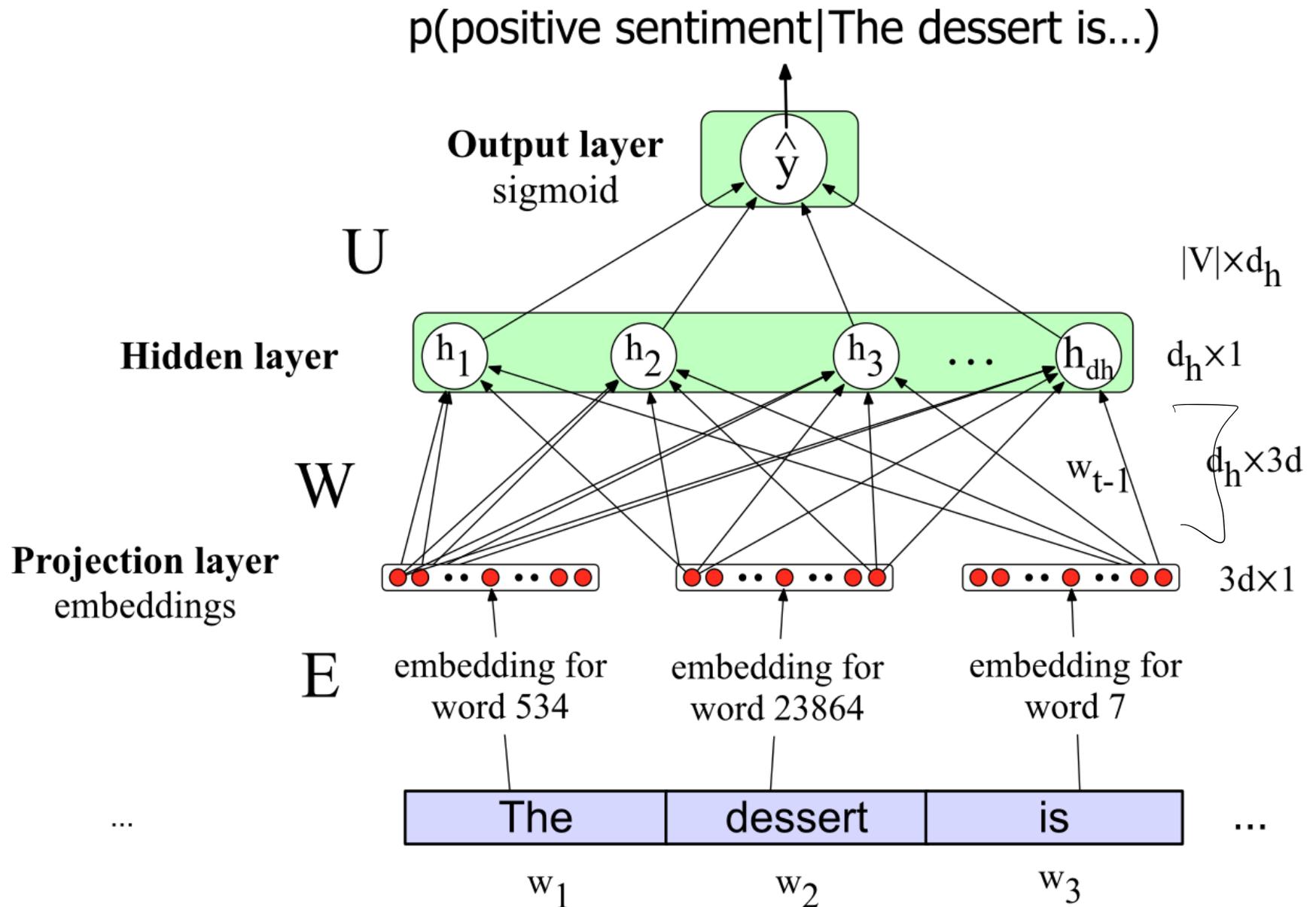
students turned the pages of their \_\_\_\_

students attentively perused their \_\_\_\_

...

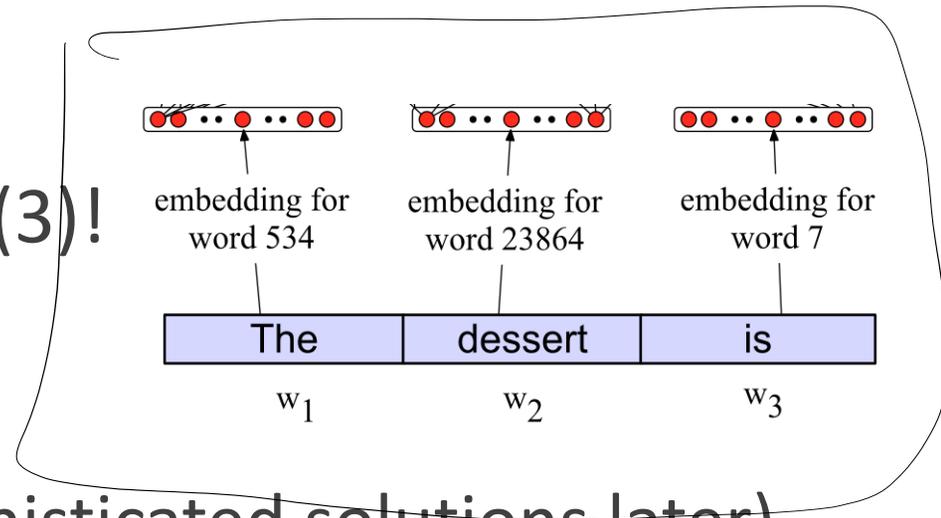
Shouldn't we *share information* across these semantically-similar prefixes?

# Neural Net Classification with embeddings as input features!



# Issue: texts come in different sizes

This assumes a fixed size length (3)!



$$\text{len(embedding)} \times \text{MAX-LEN}$$

Some simple solutions (more sophisticated solutions later)

1. Make the input the length of the longest review
  - If shorter then pad with zero embeddings
  - Truncate if you get longer reviews at test time

Markov assumption
2. Create a single "sentence embedding" (the same dimensionality as a word) to represent all the words
  - Take the mean of all the word embeddings
  - Take the element-wise max of all the word embeddings
  - For each dimension, pick the max value from all words

len(embedding)  
Bag-of-words

# composing embeddings

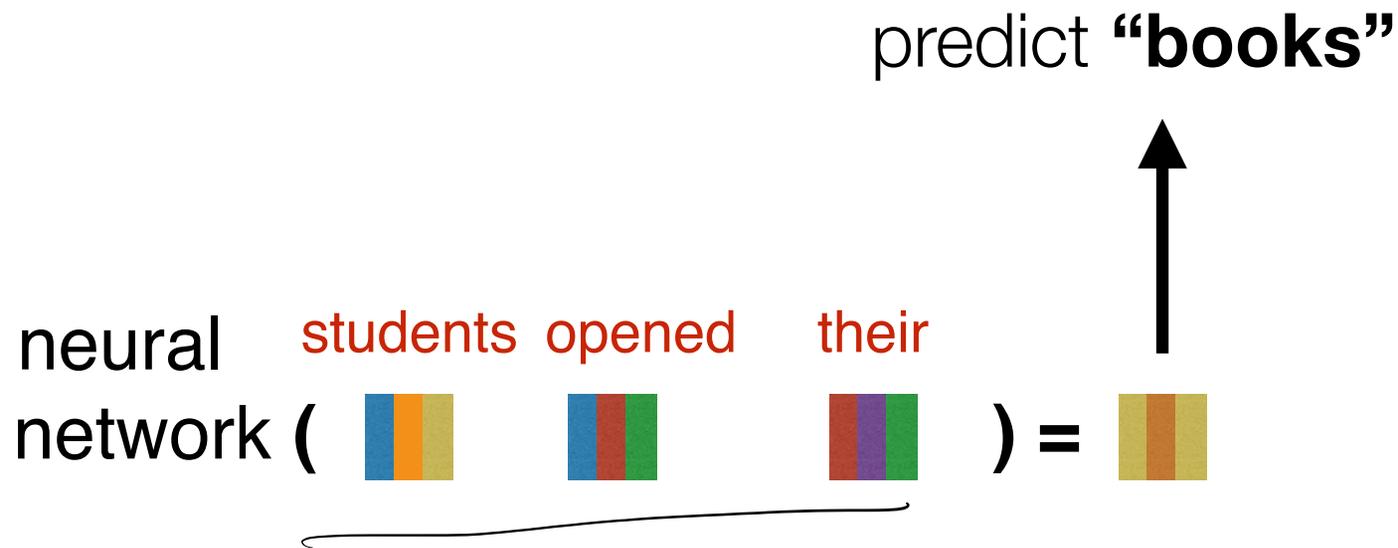
- neural networks **compose** word embeddings into vectors for phrases, sentences, and documents

1. Compose by concatenating:  $\text{len}(E) \times n$

2. Compose by averaging  $\text{len}(E)$



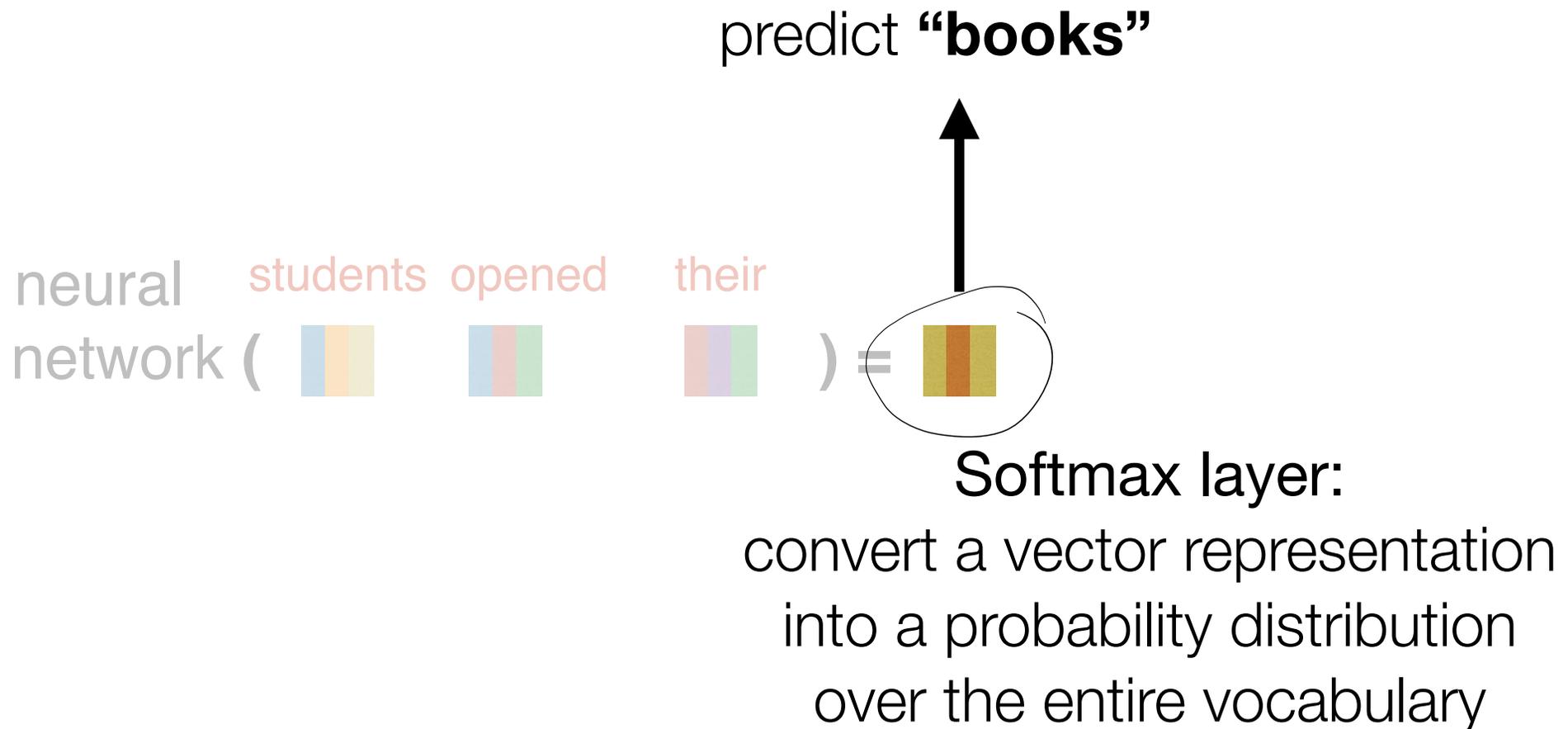
# Predict the next word from composed prefix representation



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

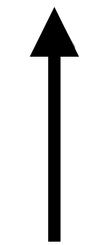
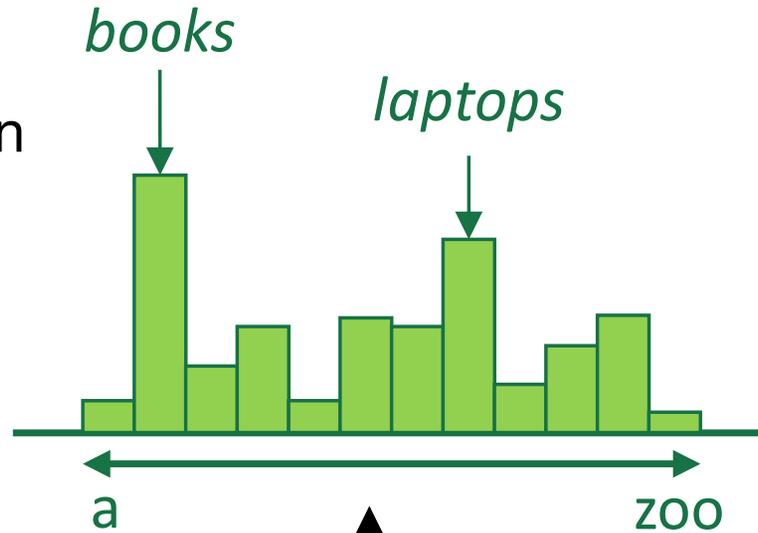


How does this happen? Let's work our way backwards, starting with the prediction of the next word



# Vocabulary

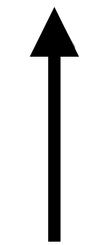
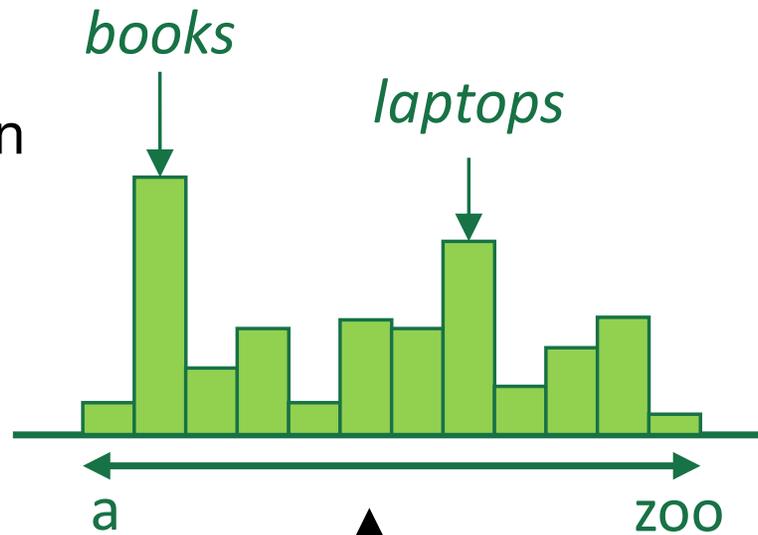
Probability distribution  
over the entire  
vocabulary



Low-dimensional  
representation of  
“students opened their”

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

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".

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

We want to get a probability distribution over these four words

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

✓ SOFTMAX

books houses lamps stamps  
 $\langle 0.2, 0.05, 0.7, 0.05 \rangle$

$$z = W \cdot X = \langle 1.8, -11.9, 12.9, -8.9 \rangle$$

$W \cdot X$  (dot product)

$$W, X = \langle 1.2 \quad -0.3 \quad 0.9 \rangle$$
$$\langle -2.3, 0.9, 5.4 \rangle$$

$$= 1.2 * -2.3 + -0.3 * 0.9$$
$$+ 0.9 * 5.4$$
$$= 1.8$$

$W =$

$$\begin{Bmatrix} 1.2 & -0.3 & 0.9 \\ 0.2 & 0.4 & -2.2 \\ 8.9 & -1.9 & 65 \\ 4.5 & 2.2 & -0.1 \end{Bmatrix}$$

← books  
← houses  
← lamps  
← stamps

$$\langle -2.3, 0.9, 5.4 \rangle$$

| | | | |

$X:$

not human-interpretable!

start with a small vector representation of the sentence prefix

"students opened their"

start with a small  
vector representation  
of the sentence prefix



Low-dimensional  
representation of  
“students opened their”

just like in regression,  
we will learn a set of  
weights



Low-dimensional  
representation of  
“students opened their”

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

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



Here's an example 3-d  
prefix vector

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

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



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

Here's an example 3-d  
prefix vector

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

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

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

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

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

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

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

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

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

intuition: each  
dimension of **x**  
corresponds to a  
*feature* of the prefix

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

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

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

CAUTION: we can't easily *interpret* these features! For example, the second dimension of  $\mathbf{x}$  likely does not correspond to any linguistic property

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

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

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



now we compute the output for this layer by taking the dot product between  $\mathbf{x}$  and  $\mathbf{W}$

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

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

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

$$\begin{aligned} & 1.2 * -2.3 \\ & + -0.3 * 0.9 \\ & + 0.9 * 5.4 \end{aligned}$$

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

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

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

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

**RESULT:**  $\langle 0.6, 0.2, 0.1, 0.1 \rangle$  |

books  
houses  
lamps  
stamps

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

**RESULT:**  $\langle 0.6, 0.2, 0.1, 0.1 \rangle$  |

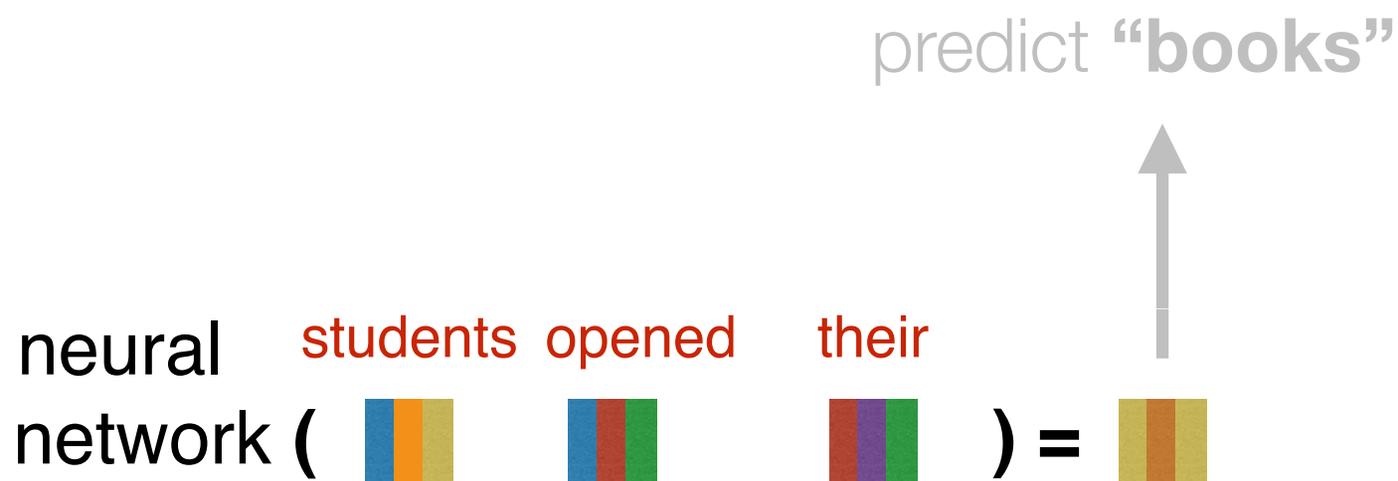
books  
houses  
lamps  
stamps

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

**So far, this is just multi-class regression on word embeddings!**

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

- 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

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

# A fixed-window neural Language Model

$$\hat{y} = \text{softmax}(W_3 h_2)$$

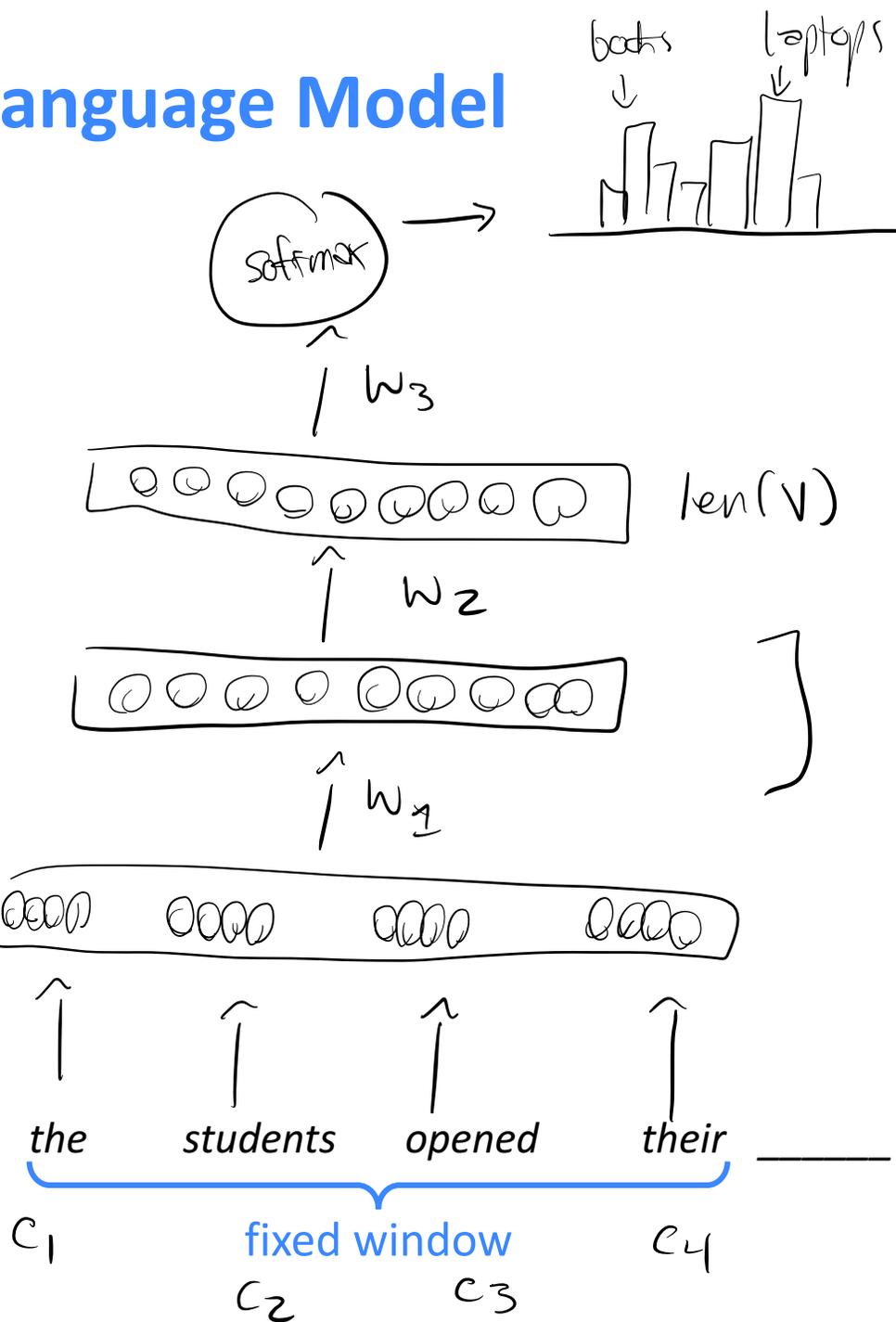
$$h_2 = f(W_2 h_1)$$

$$h_1 = f(W_1 x)$$

f: ReLU  
tanh  
sigmoid/  
softmax

$$X = [c_1; c_2; c_3; c_4]$$

~~as the proctor started the clock~~  
discard



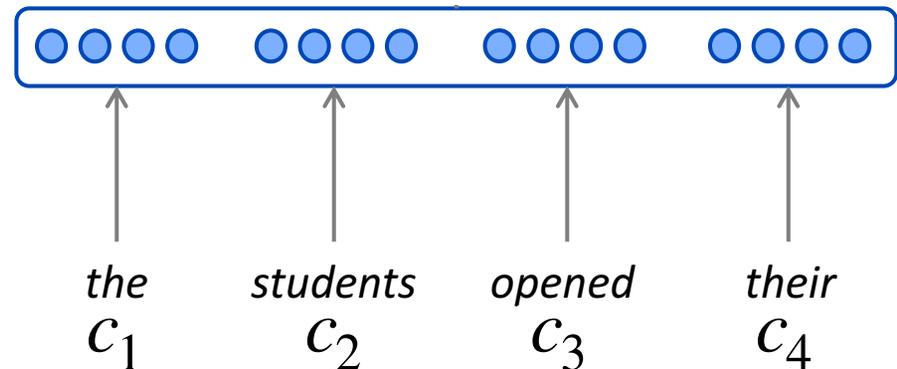
# A fixed-window neural Language Model

concatenated word embeddings

$$x = [c_1; c_2; c_3; c_4]$$

words / one-hot vectors

$$c_1, c_2, c_3, c_4$$



# A fixed-window neural Language Model

hidden layer

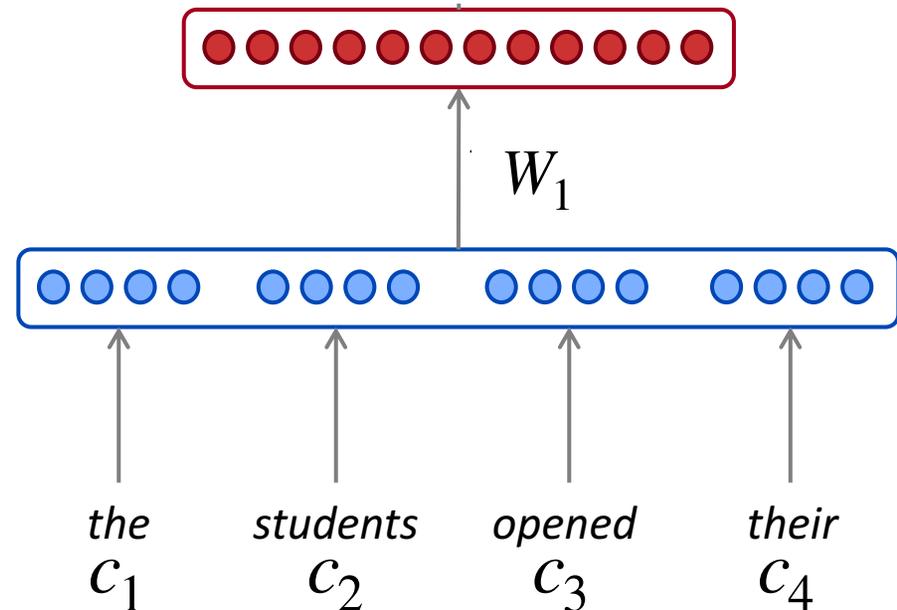
$$h = f(W_1 x)$$

concatenated word embeddings

$$x = [c_1; c_2; c_3; c_4]$$

words / one-hot vectors

$$c_1, c_2, c_3, c_4$$



# 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**.

hidden layer

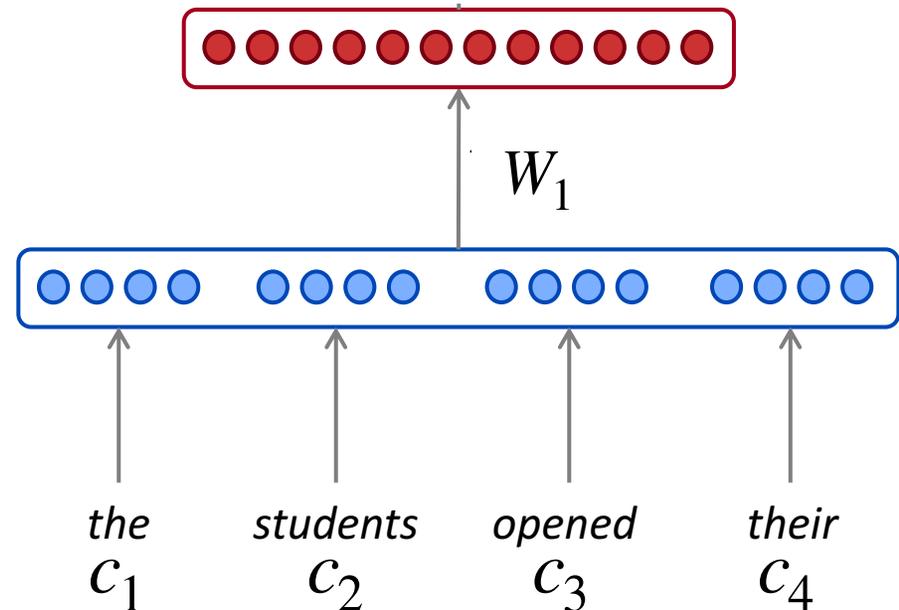
$$h = f(W_1 x)$$

concatenated word embeddings

$$x = [c_1; c_2; c_3; c_4]$$

words / one-hot vectors

$$c_1, c_2, c_3, c_4$$



# A fixed-window neural Language Model

output distribution

$$\hat{y} = \text{softmax}(W_2 h)$$

hidden layer

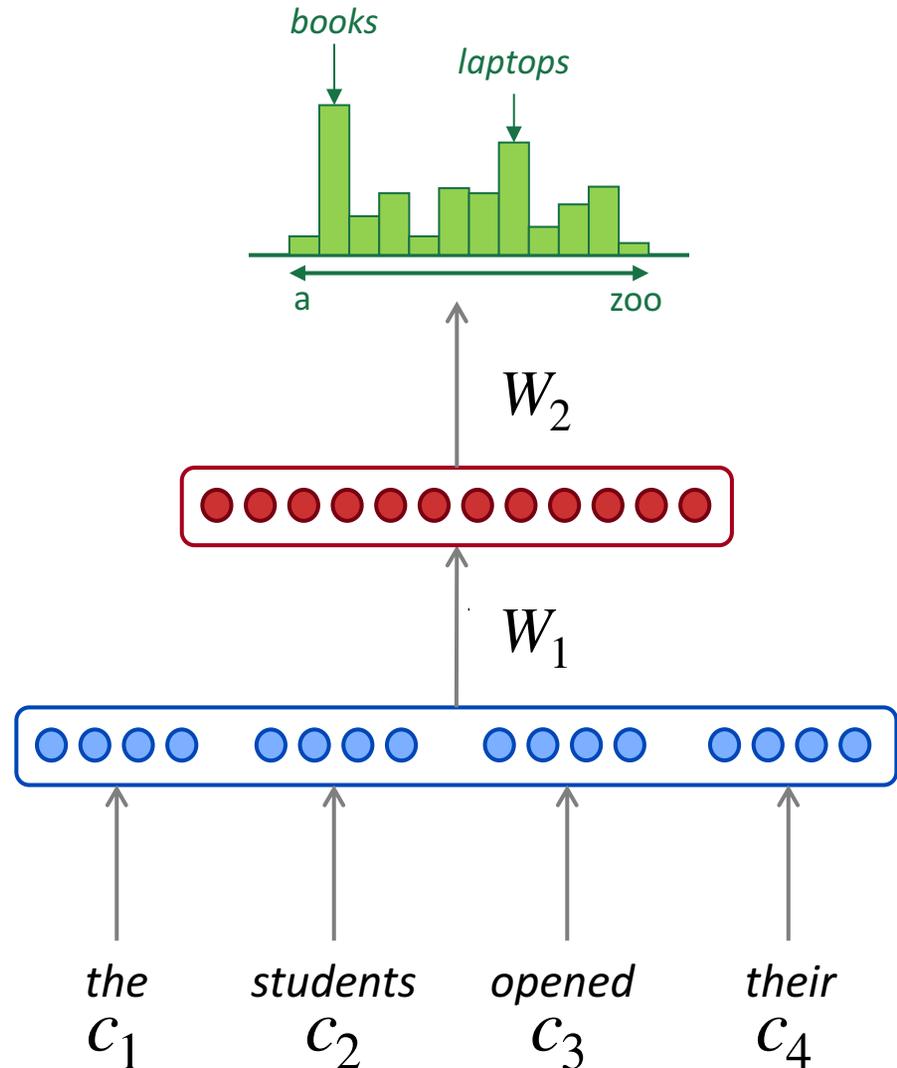
$$h = f(W_1 x)$$

concatenated word embeddings

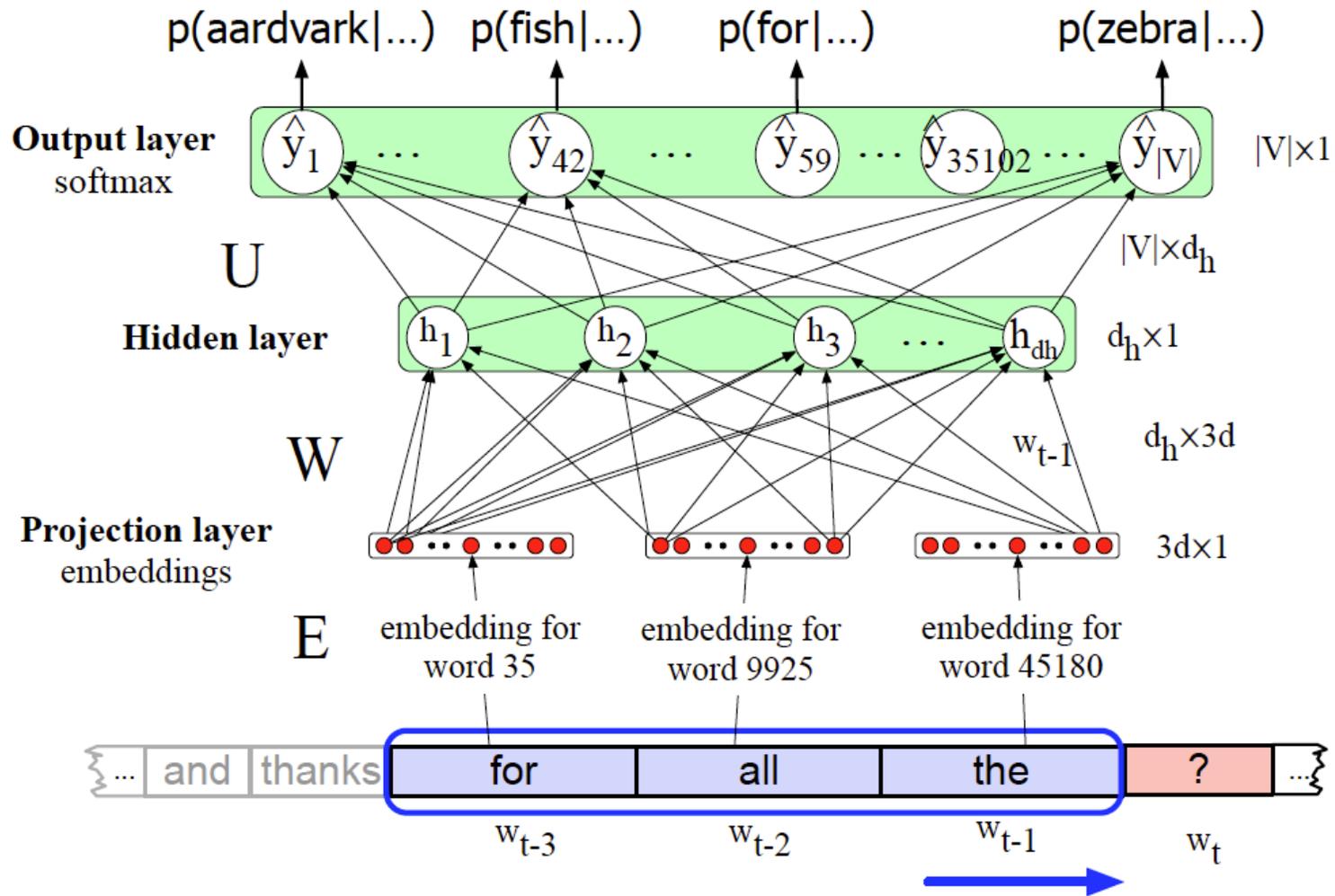
$$x = [c_1; c_2; c_3; c_4]$$

words / one-hot vectors

$$c_1, c_2, c_3, c_4$$

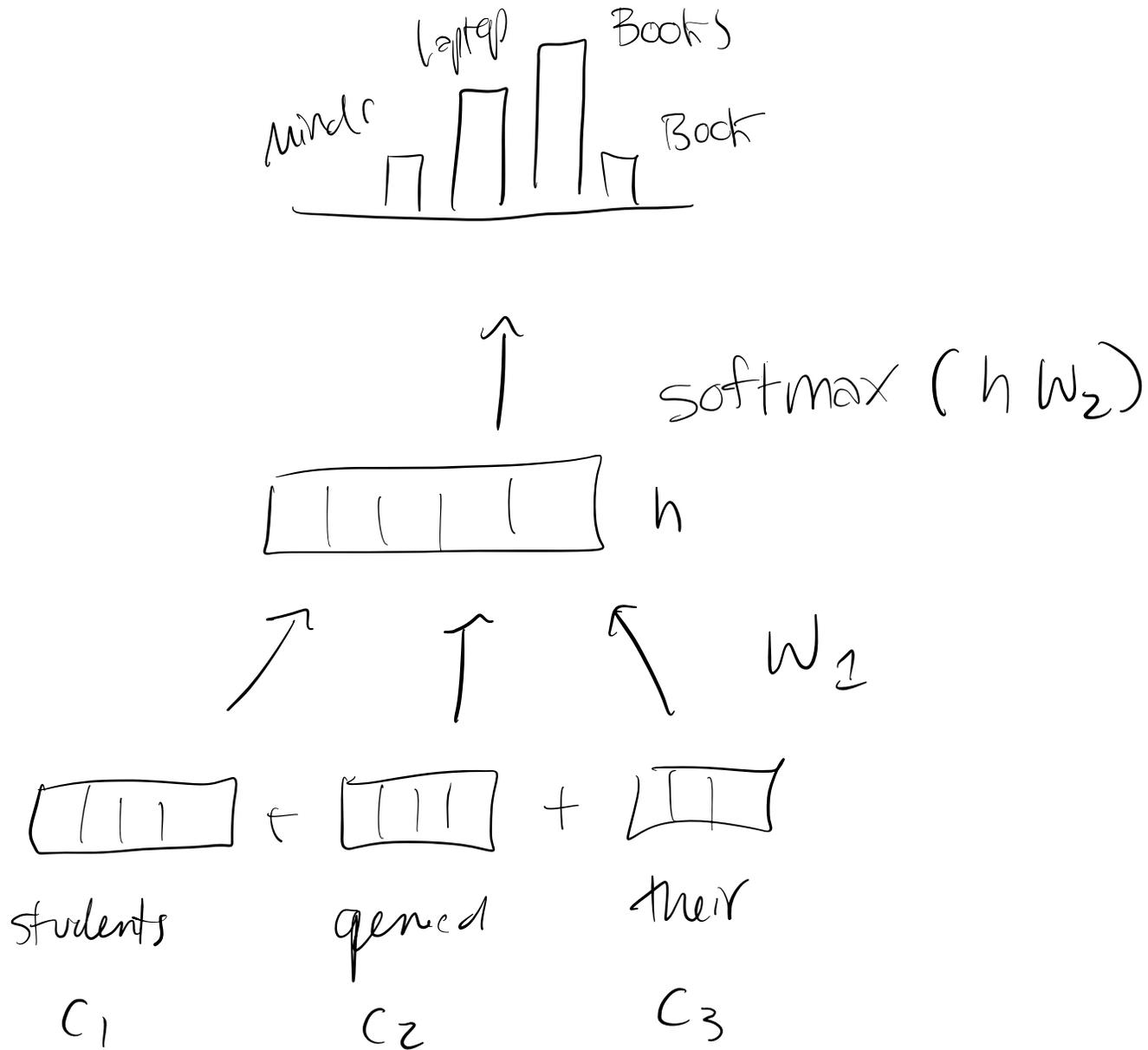


# Neural Language Model



# Training a Fixed-Length Neural Language Model

Goal: given "students open their", predict "books"



## Key Question: what are the parameters?

$W_2$ : pre-softmax weights

$W_1$ : input weights

$C_1, C_2, C_3$ : word embeddings

1. randomly initialize

2. learn weight by updating during training