# Intro: neural networks

## CS 685, Fall 2025

Advanced Natural Language Processing
https://people.cs.umass.edu/~brenocon/cs685_f25/

## Brendan O'Connor

College of Information and Computer Sciences
University of Massachusetts Amherst

- Thanks for HW1!
- Midterm #1: will be in-class on either Tu 10/7 or Th 10/9
- HW2: will be released later this week

# Exercise #1

- Instructions
  - Please work individually without internet on these questions.
  - Write answers on one sheet of paper or word processing doc. Feel free to complete them after lecture.
  - Submit a PDF to Gradescope, "Exercise 1" before Thursday's lecture.
- **Questions (based on mandatory readings)**
  - 1. In last week's Landauer and Dumais reading, they systematically varied their model's main/only hyperparameter, and examined its impact on downstream task performance. What was that hyperparameter?
  - 2. What was the downstream task?
    (Optional: what was the result of their experiment?)
  - 3. In addition to the logistic sigmoid, what two nonlinear functions for NNs are introduced in today's reading (J&M Chapter 6)?

**Target:** levied
**Choices:** (a) imposed
(b) believed
(c) requested
(d) correlated
***Solution:*** *(a) imposed*

- TOEFL (near-)synonym task
  - Acontextual version
- Choose answer with highest cosine to E("levied")
- https://aclweb.org/aclwiki/ TOEFL_Synonym_Questions_(State_of_t he_art)
- Hyperparameter sweep!
  - You should always get a cap shape if you've searched enough
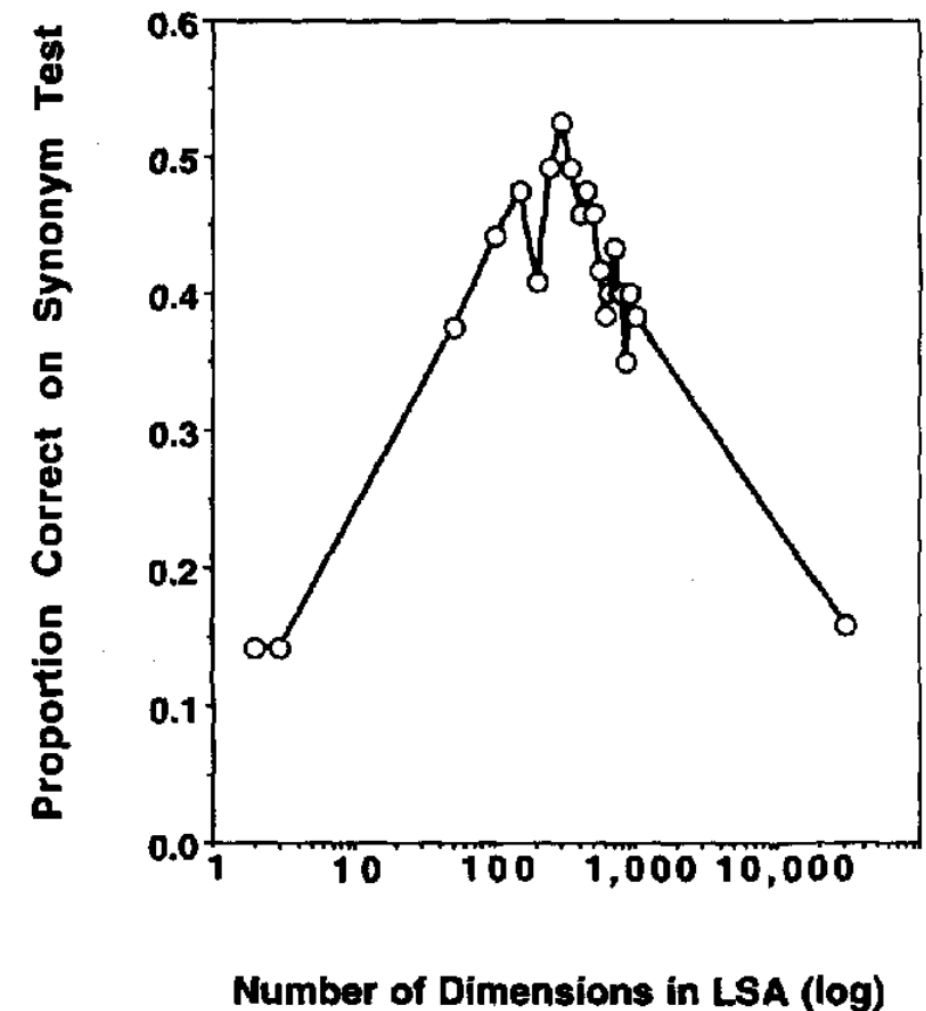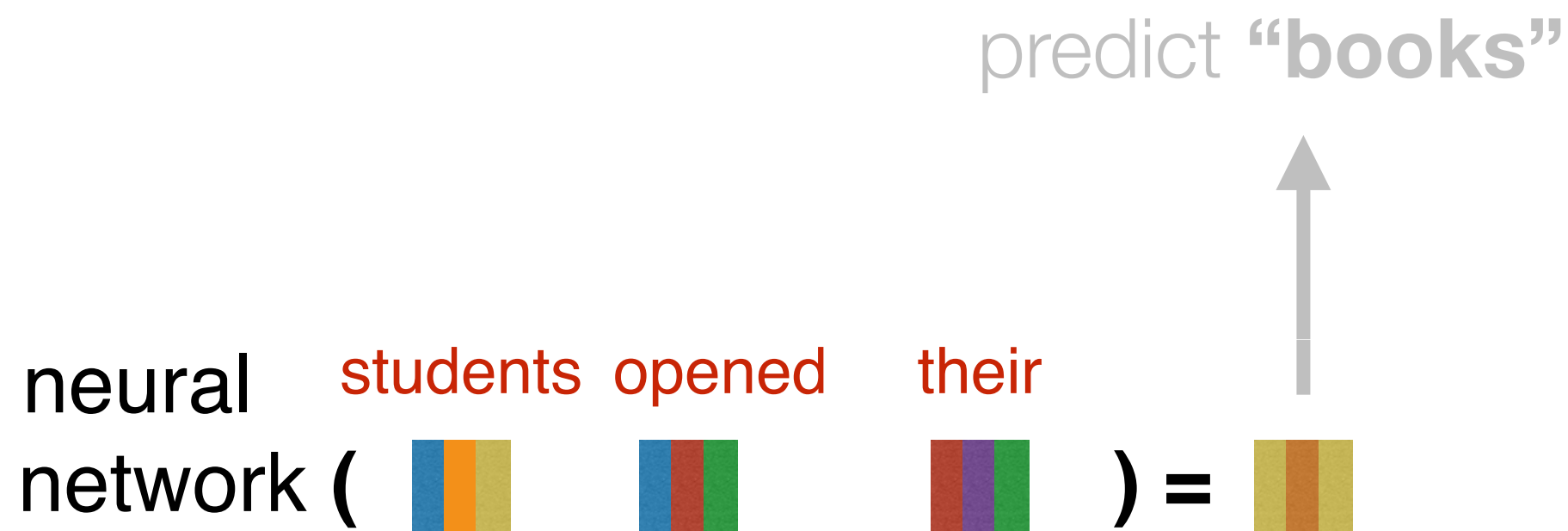


*Figure 3.* The effect of number of dimensions retained in latent-semantic-analysis (LSA)–singular-value-decomposition (SVD) simulations of word-meaning similarities. The dependent measure is the proportion of 80 multiple-choice synonym test items for which the model chose the correct answer. LSA was trained on text samples from 30,473 articles in an electronic file of text for the *Groliers Academic American Encyclopedia*.

# Today: Neural Networks for Language

- Artificial neural network models of language incorporate
    - Word embeddings  (done!)
    - Cool/controversial modeling metaphor  (new?)
    - Non-linearities or feedforward layers  (new)
    - Flexible architectures  (new)

    - Learning from predictive loss minimization (review - Thurs)
    - Flexible architectures enabled by autodiff-based learning (new - Thurs)
        - important to deal with sequential and dependency structure in language
        - main architectures:
            - Recurrent NNs
            - Attention and Transformer NNs

# how do we compute a prefix representation *x* for next-word prediction?

predict **"books"**

neural network **(** students opened their **) =**

# Composition functions

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

*output*: single vector

- Element-wise functions (e.g. sum)
- Concatenation

- Feed-forward neural networks
- Recurrent neural networks
- Convolutional neural networks *[skipping]*

- Transformers (self-attention) neural networks

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

review of Bengio's model (non-linearities
skipped last week)

# A fixed-window neural Language Model

*as    the    proctor    started    the    clock*        *the        students        opened        their*  _____
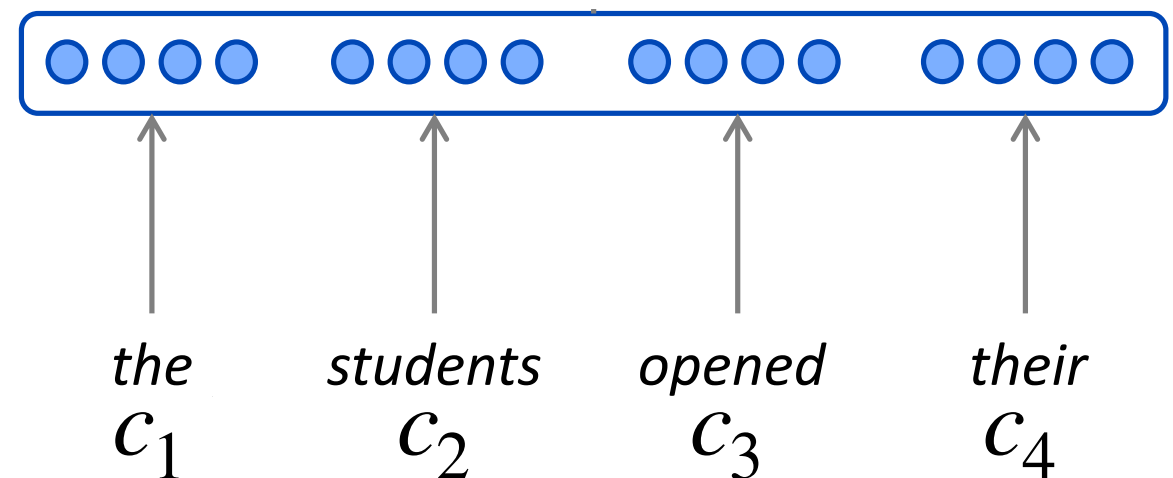
discard

fixed window

# A fixed-window neural Language Model

concatenated word embeddings

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

words / one-hot vectors

$$c_1, c_2, c_3, c_4$$

*the*
$c_1$

*students*
$c_2$

*opened*
$c_3$

*their*
$c_4$

# 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

$$c_1, c_2, c_3, c_4$$

$W_1$

*the*     *students*     *opened*     *their*
$c_1$       $c_2$          $c_3$         $c_4$

# 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 **ReLu**(x) = max(0, x). Other choices include **tanh** and **sigmoid**.

hidden layer

$$h = f(W_1 \ldots$$

concatenated v

$$x = [c_1; \ldots$$

$y=x$

$y=0$

words / one-hot vectors

$c_1, c_2, c_3, c_4$

the
$c_1$

students
$c_2$

opened
$c_3$

their
$c_4$

# A fixed-window neural Language Model
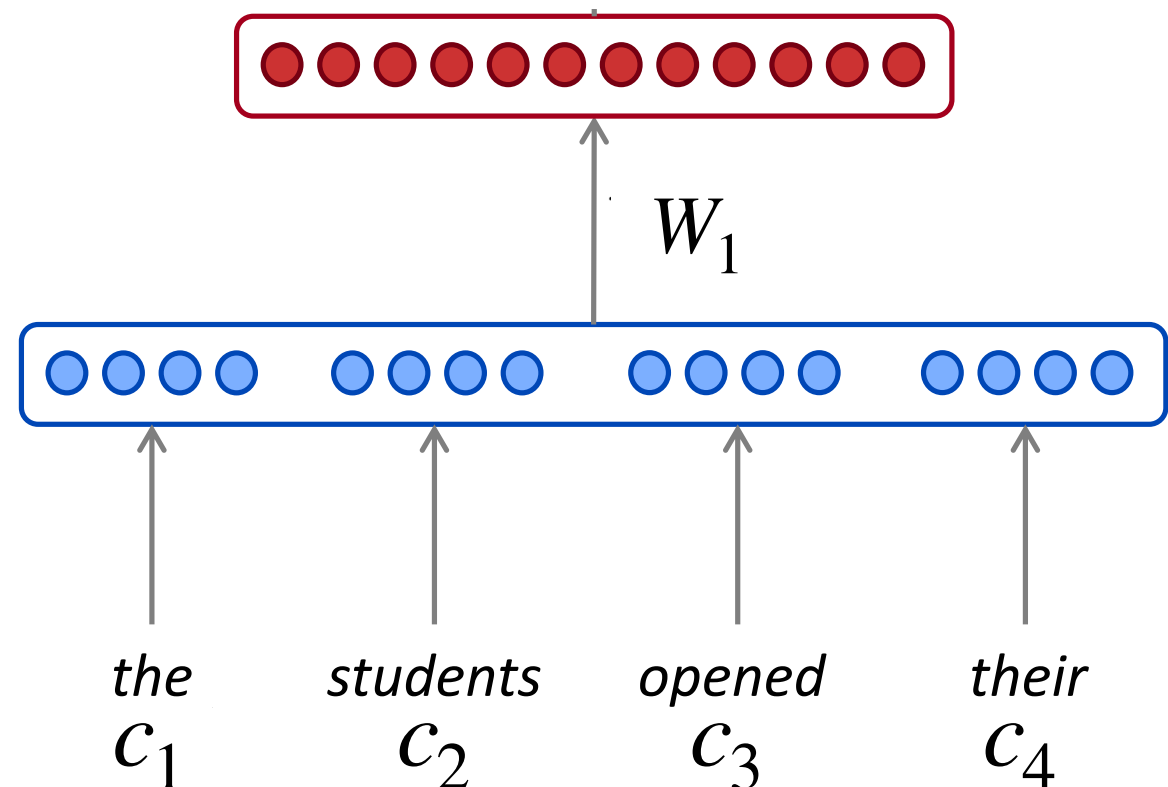
output distribution

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

hidden layer

$$h = f(W_1 x)$$

concatenated word embeddings

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

words / one-hot vectors

$$c_1, c_2, c_3, c_4$$

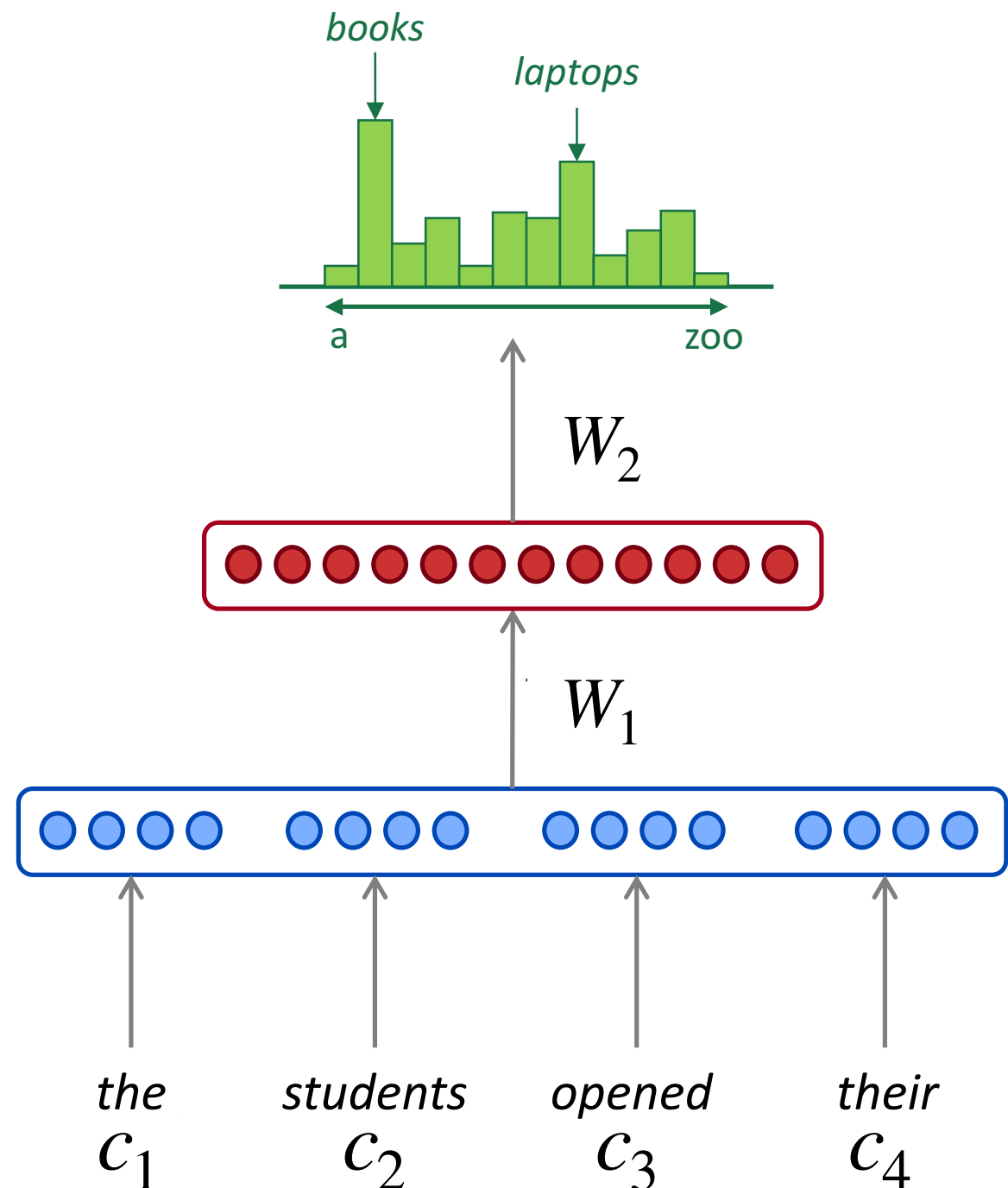books

laptops

a          zoo

$W_2$

$W_1$

*the*       *students*       *opened*       *their*
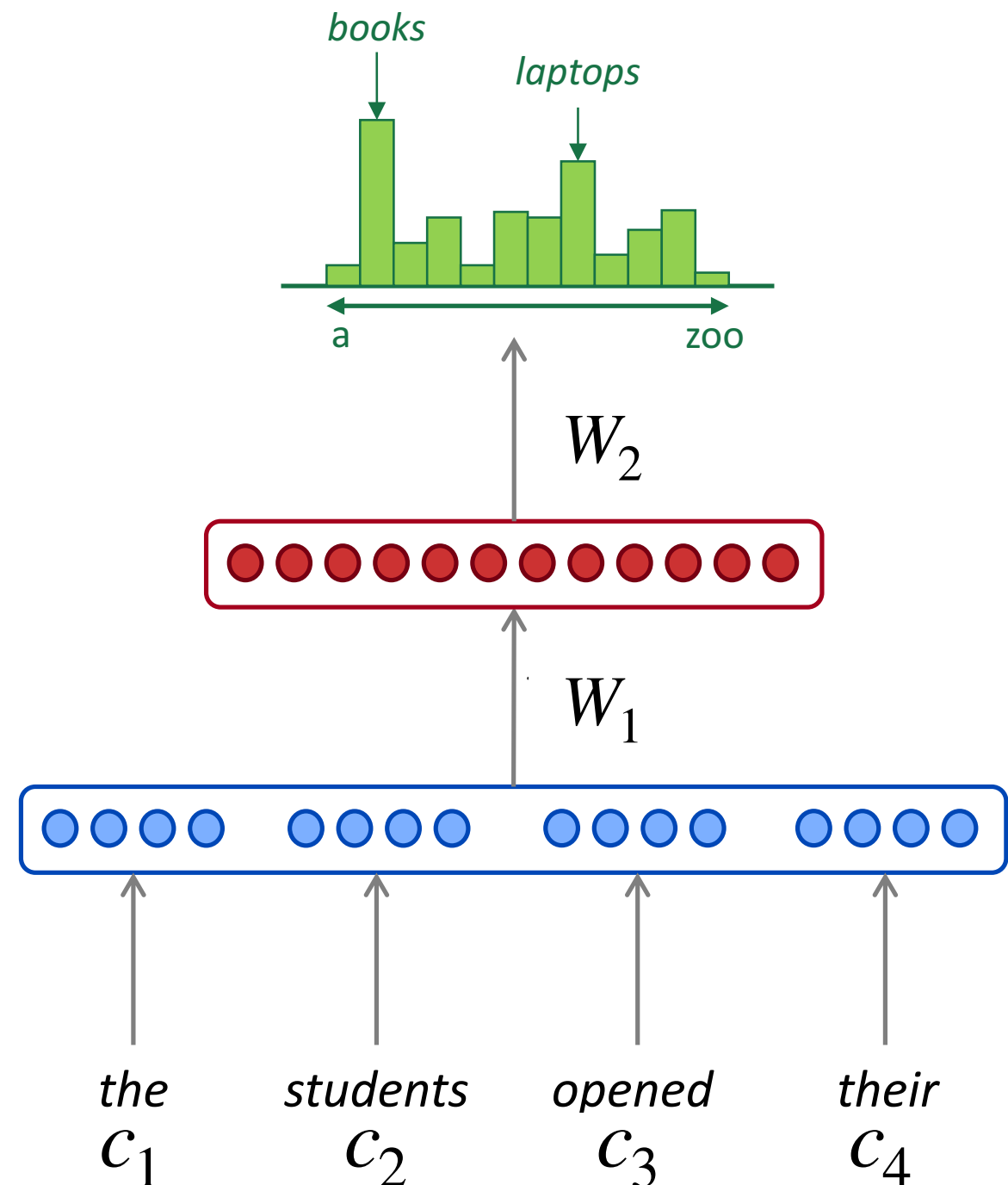$c_1$          $c_2$            $c_3$          $c_4$

**how does this compare to a normal n-gram model?**

**Improvements** over *n*-gram LM:
- No sparsity problem

Remaining **problems**:
- Fixed window is too small
- Enlarging window enlarges $W$
- Window can never be large enough!
- Each $c_i$ uses different rows of $W$. We don't share weights across the window.

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

books

laptops

a          zoo

$$W_2$$

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

$$W_1$$

$$\hat{y} = \mathrm{softmax}(Uh + b_2)$$

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

the          students          opened          their
$C_1$          $C_2$          $C_3$          $C_4$

$$|x^{(1)}, \ldots, \langle x^{(2)}, x^{(3)}, x^{(4)}, x^{(4)}, x^{(4)}$$

# NN non-linearities

- Why use an element-wise "squashing" function?

https://playground.tensorflow.org/

# Feedforward (MLP) Neural Network

If we feed a vector of inputs through a bunch of logistic regression functions, then we get a vector of outputs ...



Layer $L_1$

*But we don't have to decide ahead of time what variables these logistic regressions are trying to predict!*

# Feedforward (MLP) Neural Network

… which we can feed into another logistic regression function



*It is the loss function that will direct what the intermediate hidden variables should be, so as to do a good job at predicting the targets for the next layer, etc.*

# Feedforward (MLP) Neural Network

Before we know it, we have a multilayer neural network….

# Question

- Say we have a many-layered feedforward network
  - $y = g(A\ g(B\ g(C\ x)))$
- Can we use the identity function $g(x)=x$ for the elementwise non-linearity layer?

- Visual demo of the potential usefulness NN non-linearities (for non-language data, at least)

  https://playground.tensorflow.org/

# Recurrent Neural Networks!

# A RNN Language Model

word embeddings

$c_1, c_2, c_3, c_4$



the $E \in \mathbb{R}^{d \times |V|}$ students $E \in \mathbb{R}^{d \times |V|}$ opened $\in \mathbb{R}^{d \times |V|}$ their $\in \mathbb{R}^{d \times |V|}$

$c_1 \qquad c_2 \qquad c_3 \qquad c_4$

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

# A RNN Language Model

**hidden states**

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

h^(0) is initial hidden state!

$h^{(0)}$

**word embeddings**

$$c_1, c_2, c_3, c_4$$

the $\boldsymbol{E} \in$ students $\mathbb{R}^{d \times |V|}$ opened $\mathbb{R}^{d \times |V|}$ their $\in \mathbb{R}^{d \times |V|}$

$$c_1 \qquad c_2 \qquad c_3 \qquad c_4$$

$$\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)},$$

# A RNN Language Model

hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

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

word embeddings

$$c_1, c_2, c_3, c_4$$



$h^{(0)}$  $h^{(1)}$

$W_h$

$W_e$

the students opened their

$E \in \mathbb{R}^{d \times |V|}$  $E \in \mathbb{R}^{d \times |V|}$  $E \in \mathbb{R}^{d \times |V|}$  $E \in \mathbb{R}^{d \times |V|}$

$c_1$      $c_2$      $c_3$      $c_4$

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

# A RNN Language Model

hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

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

word embeddings

$$c_1, c_2, c_3, c_4$$



$h^{(0)}$  $h^{(1)}$  $h^{(2)}$

$W_h$  $W_h$

$W_e$  $W_e$

the $E \in \mathbb{R}^{d \times |V|}$ students $E \in \mathbb{R}^{d \times |V|}$ opened $\mathbb{R}^{d \times |V|}$ their $\in \mathbb{R}^{d \times |V|}$

$c_1$  $c_2$  $c_3$  $c_4$

$x^{(1)}, \ldots, x^{(2)}, x^{(3)},$
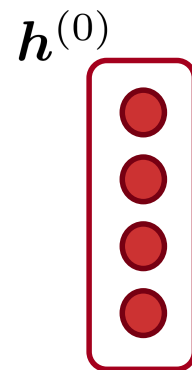
# A RNN Language Model

hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

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

word embeddings

$$c_1, c_2, c_3, c_4$$



$h^{(0)}$    $h^{(1)}$    $h^{(2)}$    $h^{(3)}$

$W_h$   $W_h$   $W_h$

$W_e$   $W_e$   $W_e$

the students opened their

$c_1$    $c_2$    $c_3$    $c_4$

$E \in \mathbb{R}^{d \times |V|}$

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

# A RNN Language Model

hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

h(0) is initial hidden state!

word embeddings

$$c_1, c_2, c_3, c_4$$
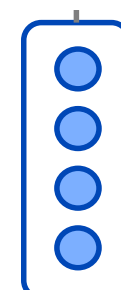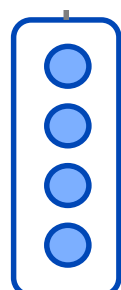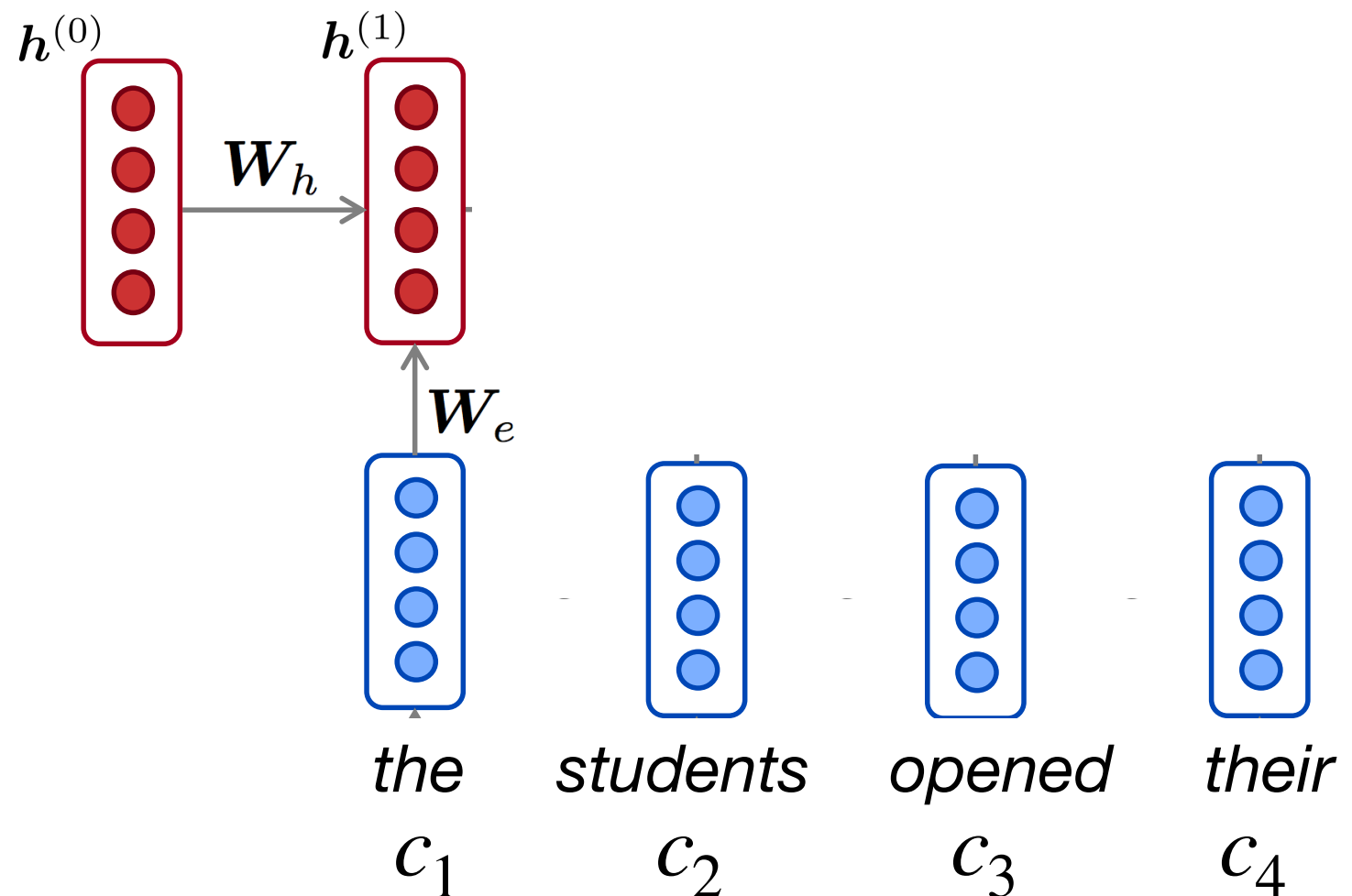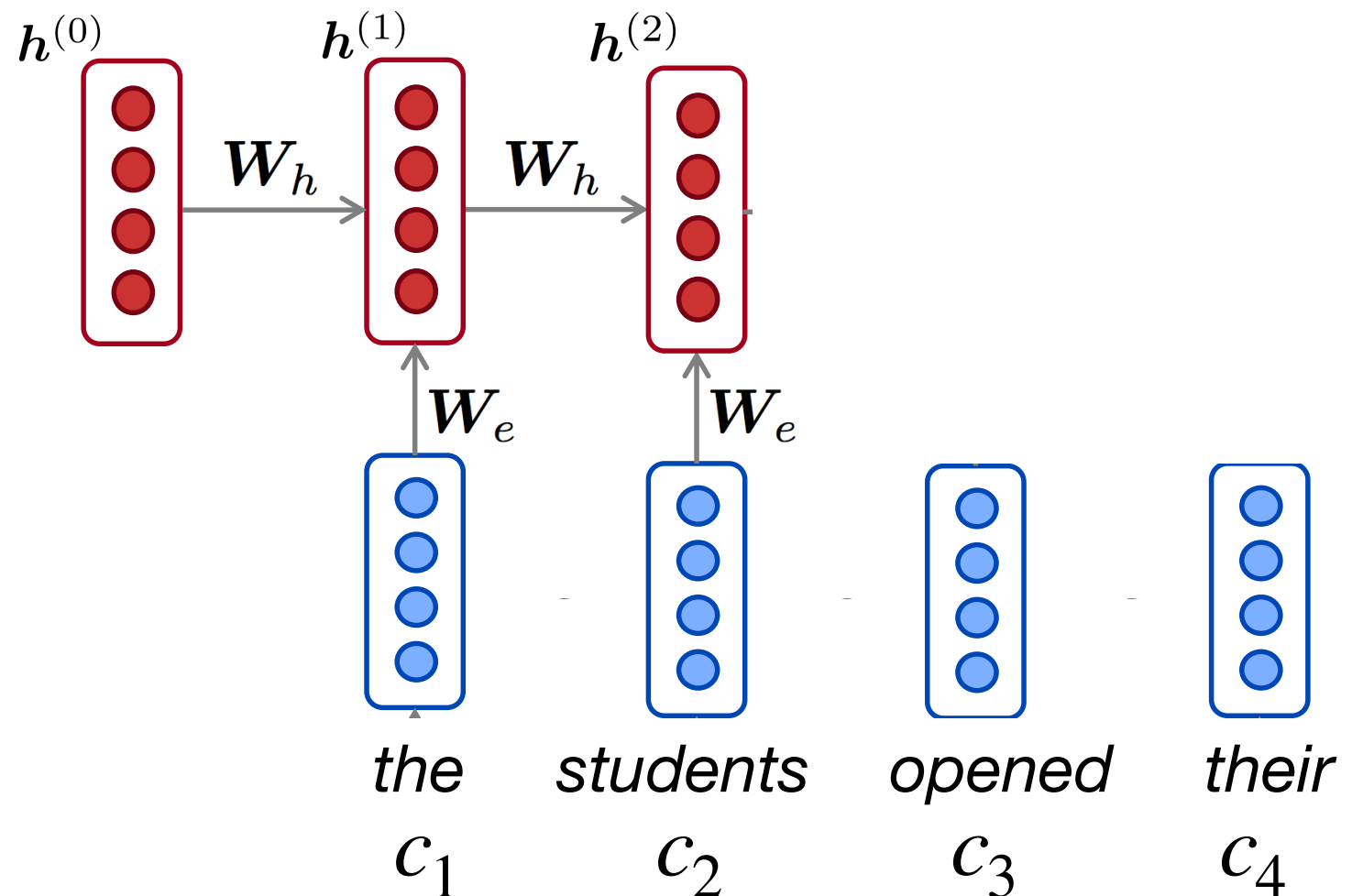
# A RNN Language Model

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

books

laptops

output distribution

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

a                    zoo

$W_2$

$\boldsymbol{h}^{(0)}$    $\boldsymbol{h}^{(1)}$    $\boldsymbol{h}^{(2)}$    $\boldsymbol{h}^{(3)}$    $\boldsymbol{h}^{(4)}$

hidden states

$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t)$$

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

$\boldsymbol{W}_h$    $\boldsymbol{W}_h$    $\boldsymbol{W}_h$    $\boldsymbol{W}_h$

$\boldsymbol{W}_e$    $\boldsymbol{W}_e$    $\boldsymbol{W}_e$    $\boldsymbol{W}_e$

word embeddings

$$c_1, c_2, c_3, c_4$$

the$\boldsymbol{E} \in$ students$\mathbb{R}^{d \times |V|}$ opened$\mathbb{R}^{d \times |V|}$ their$\in \mathbb{R}^{d \times |V|}$

$c_1$                $c_2$                $c_3$                $c_4$

$$\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)},$$

## why is this good?

RNN **Advantages**:
- Can process any length input
- Model size doesn't increase for longer input
- Computation for step $t$ can (in theory) use information from many steps back
- Weights are shared across timesteps → representations are shared

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

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

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

laptops

$|V|$

a      zoo

$W_o$

$\boldsymbol{U} \in \mathbb{R}^{|V| \times D_h}$

$\boldsymbol{h}^{(0)}$    $\boldsymbol{h}^{(1)}$    $\boldsymbol{h}^{(2)}$    $\boldsymbol{h}^{(3)}$

$\boldsymbol{h}^{(0)}$

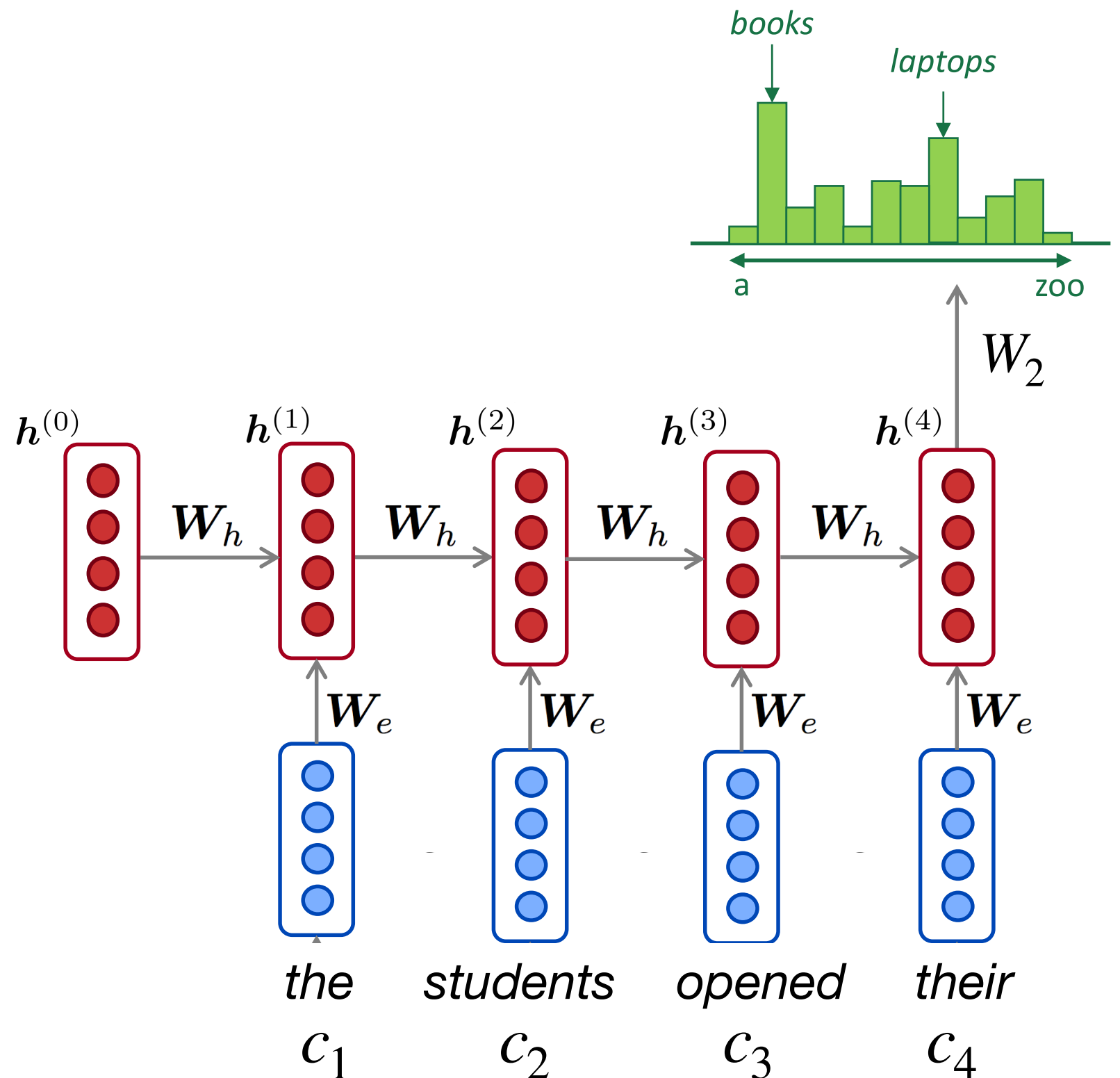$\boldsymbol{h}^{(1)}, \boldsymbol{h}^{(1)}, \boldsymbol{h}\,\boldsymbol{h}^{(1)}, \boldsymbol{h}\,\boldsymbol{h}^{(1)}, \boldsymbol{h}^{(2)}, \boldsymbol{h}^{(3)}, \boldsymbol{h}^{(4)}$

$\boldsymbol{W}_h \in \mathbb{R}^{D_h} \boldsymbol{W}_h \in \mathbb{R}^{D_h} \boldsymbol{W}_h \in \mathbb{R}^{D_h} \boldsymbol{W}_h \in \mathbb{R}^{D_h \times D_h}$

$\boldsymbol{W}_e \in \mathbb{R}^{D_h} \boldsymbol{W}_e \in \mathbb{R}^{D_h} \boldsymbol{W}_e \in \mathbb{R}^{D_h} \boldsymbol{W}_e \in \mathbb{R}^{D_h \times d}$

$\boldsymbol{e}^{(1)}, \ldots, \boldsymbol{e}^{(t)} \in \mathbb{R}^{d}; \boldsymbol{e}^{(1)}, \boldsymbol{e}^{(2)}, \boldsymbol{e}^{(3)}, \boldsymbol{e}^{(4)}$

$\boldsymbol{E} \in \mathbb{R}^{d \times |V|} \boldsymbol{E} \in \mathbb{R}^{d \times |V|} \in \mathbb{R}^{d \times |V|} \in \mathbb{R}^{d \times |V|} \mathbb{R}^{d \times |V|}$

$c_1$    $c_2$    $c_3$    $c_4$

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

29

- stopped here 9/16

# Gradient-based learning

- Goal: learn all model parameters W
- Loss function L(W) based on dataset
- Choose W to minimize L(W) be following the negative gradient of the loss
- Intuition: cross-entropy gradient shifts probability mass to the data

- Next time: gradient learning for arbitrary NN functions; issues and solutions