
CS 333:
Natural Language
Processing

Fall 2023

Prof. Carolyn Anderson
Wellesley College

Reminders

- ❖ I'm out of town for a conference most of this week
- ❖ **No help hours Thursday**
- ❖ Cynthia has help hours on Wednesday
- ❖ First Gen in CS lunch this Wednesday
- ❖ CS Colloquium next Wednesday

Recap

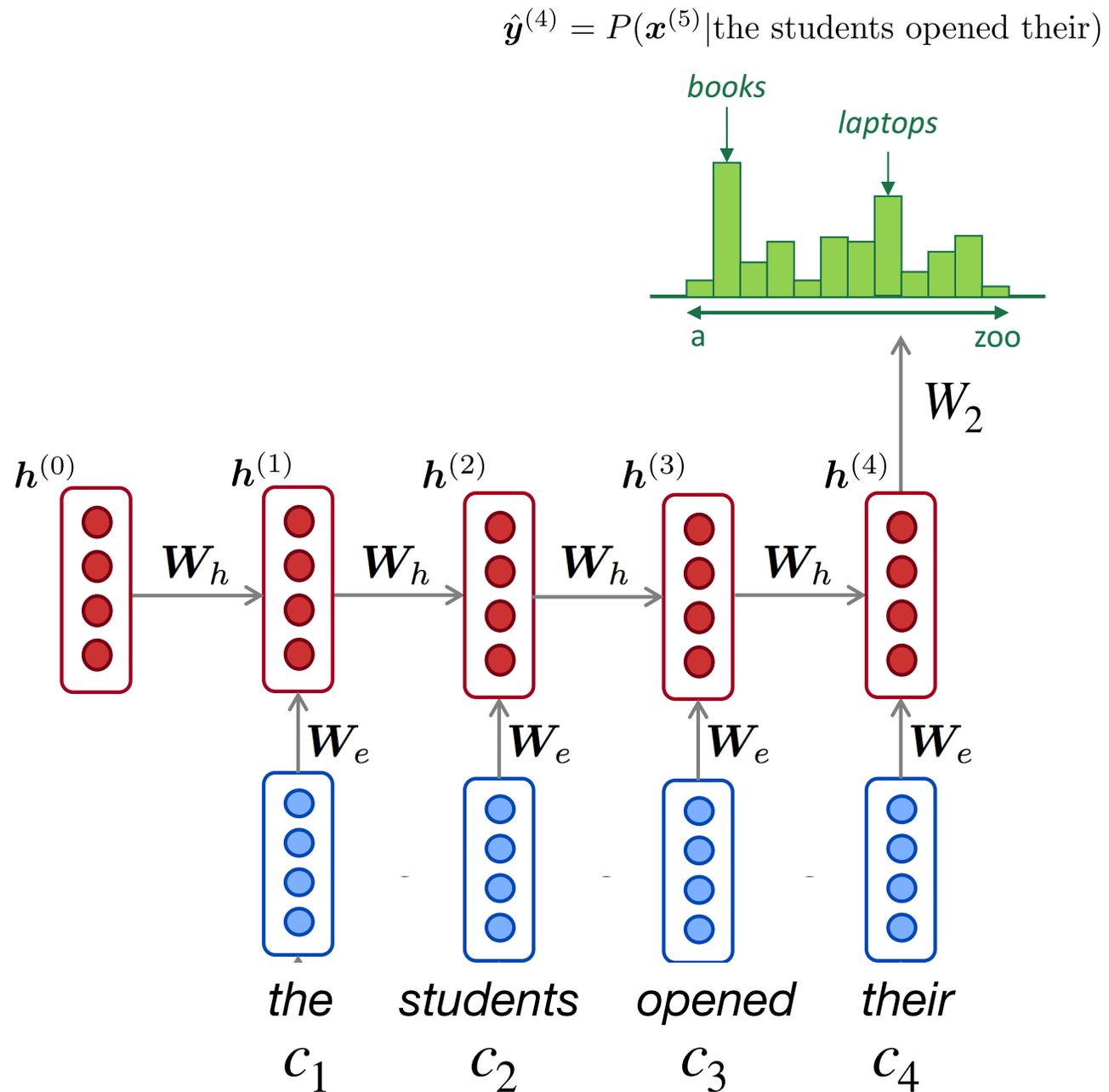
why is this good?

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 \rightarrow representations are shared

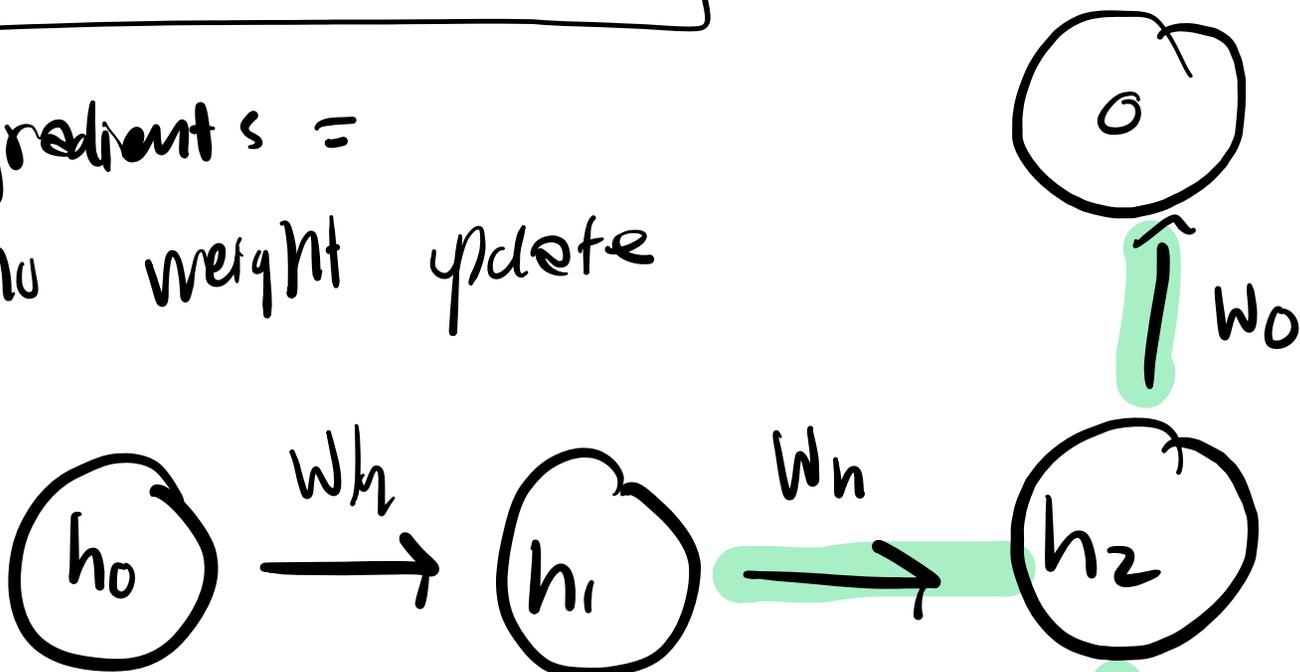
RNN Disadvantages:

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



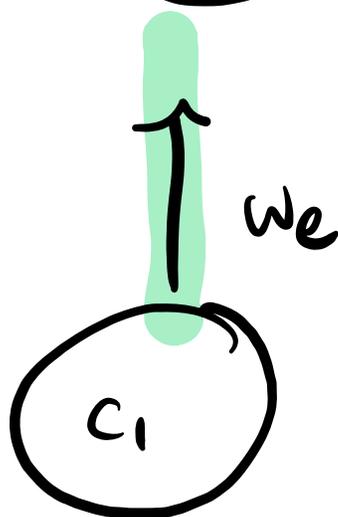
Vanishing Gradients Problem

Zero gradients =
no weight update

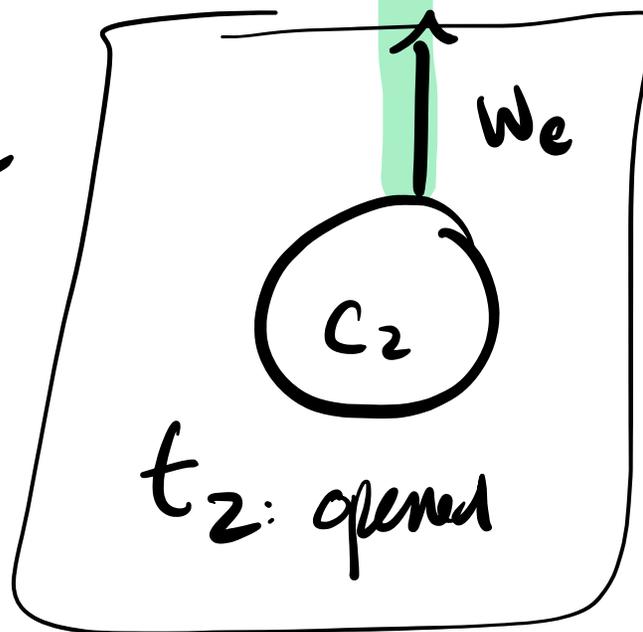


Vanishing
Gradients
Problem

Exploding
Gradients



t_1 : students

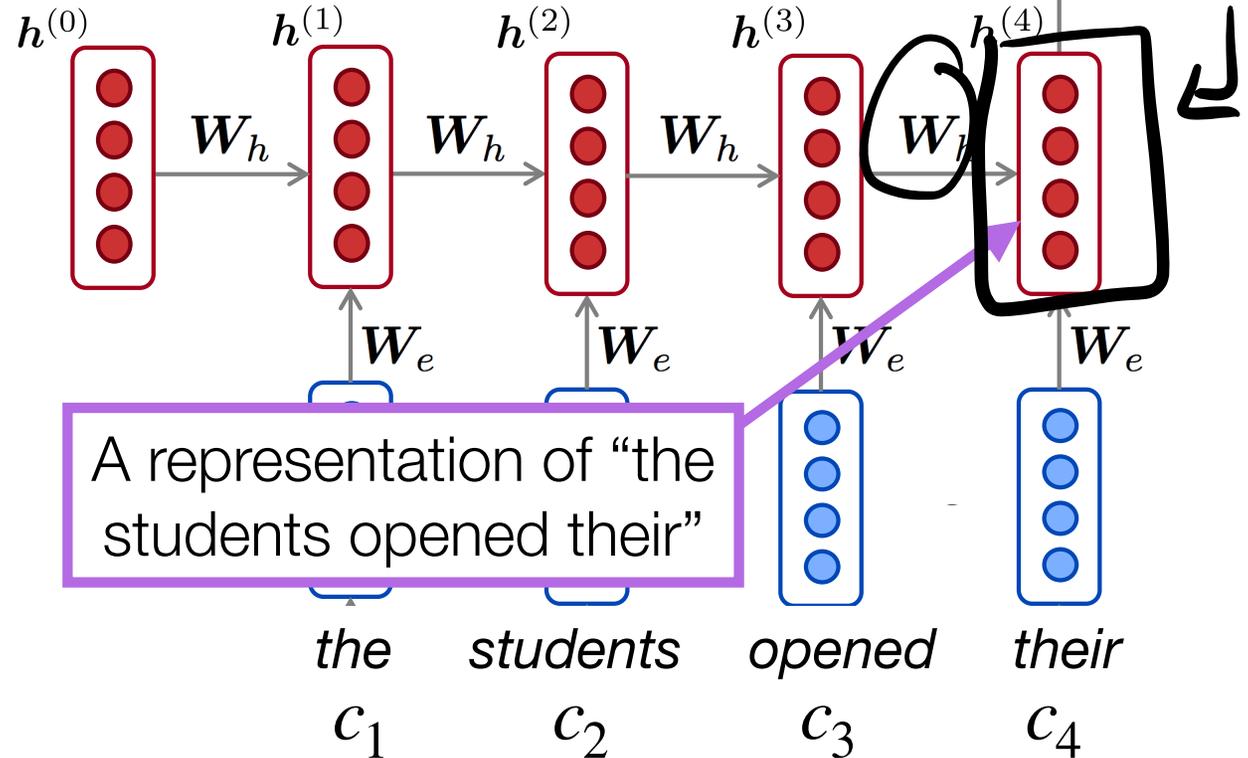


t_2 : opened

RNNs suffer from a **bottleneck** problem

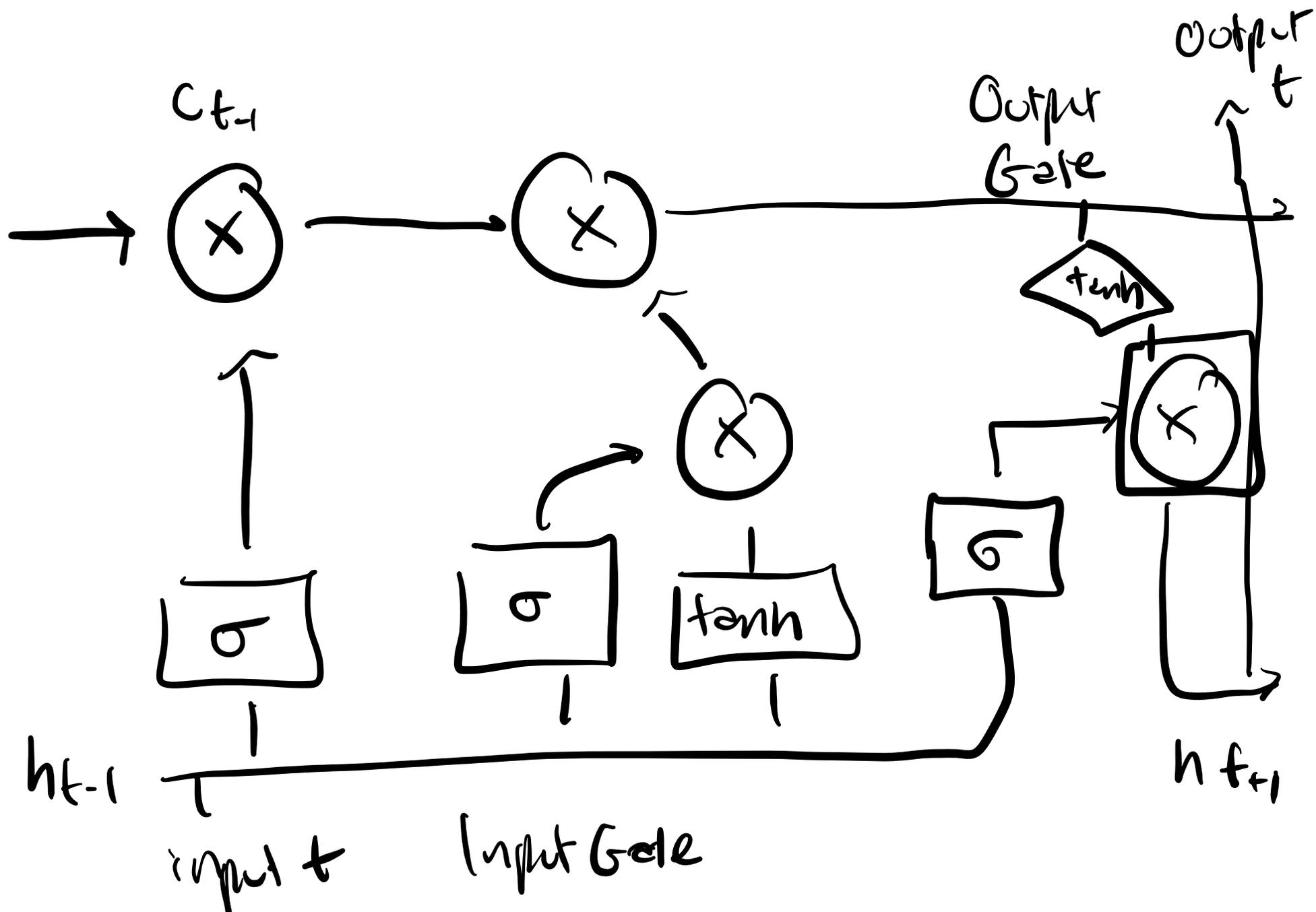
The current hidden representation must encode all of the information about the text observed so far

This becomes difficult especially with longer sequences



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

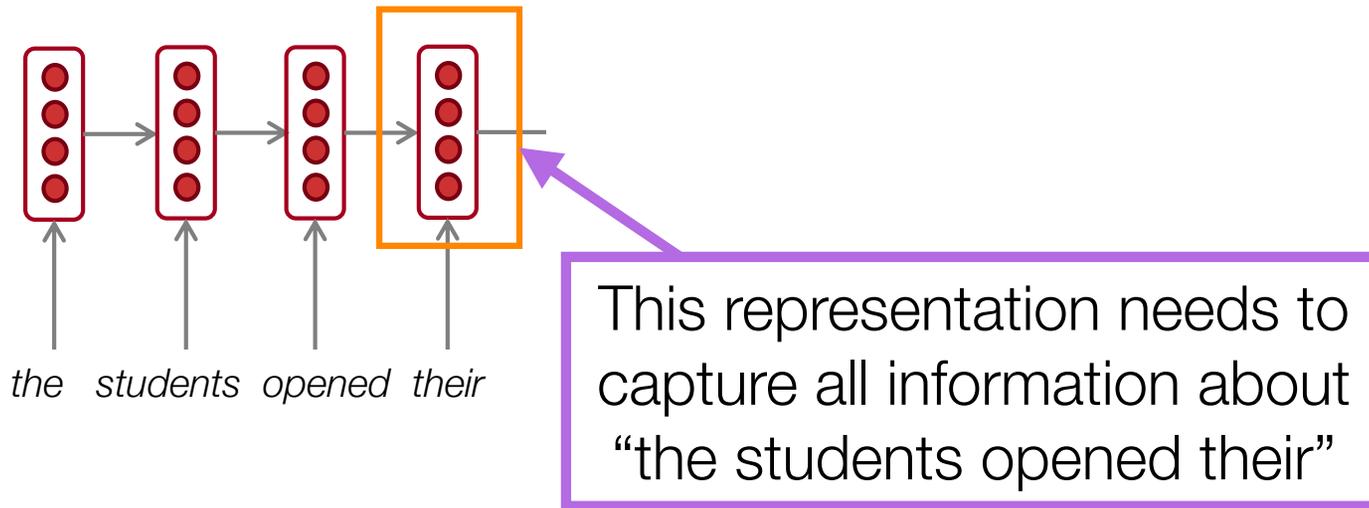
LSTMs



“you can’t cram the meaning
of a whole %&@#&ing
sentence into a single
\$*(&@ing vector!”

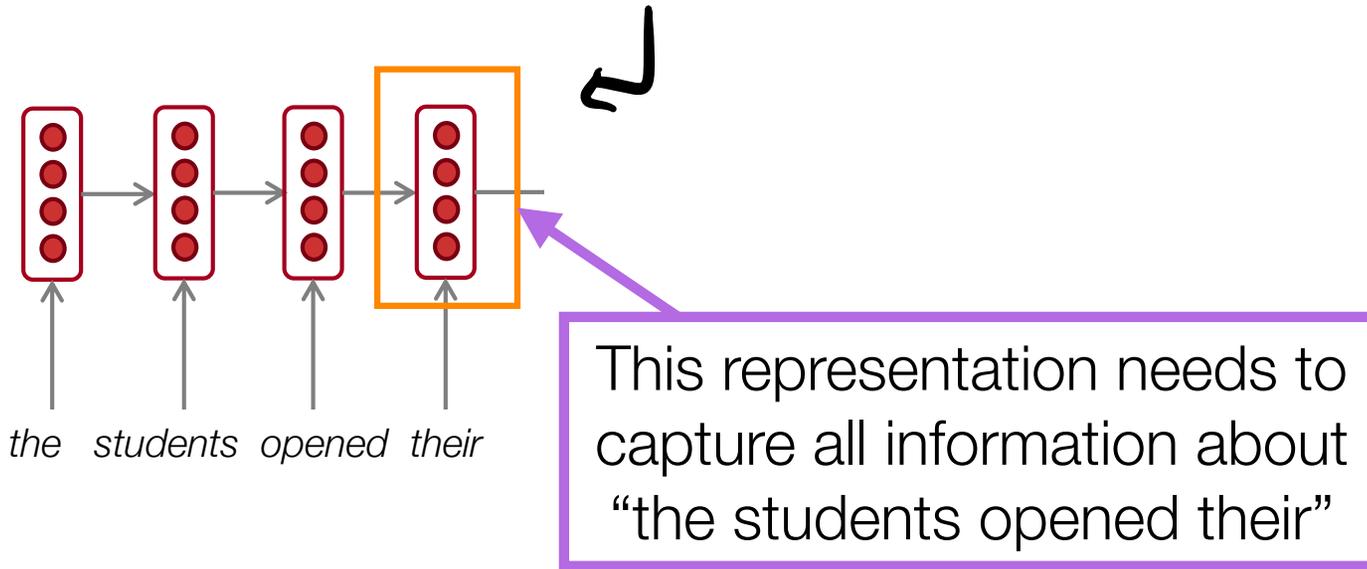
— Ray Mooney (NLP professor at UT Austin)

idea: what if we use multiple vectors?



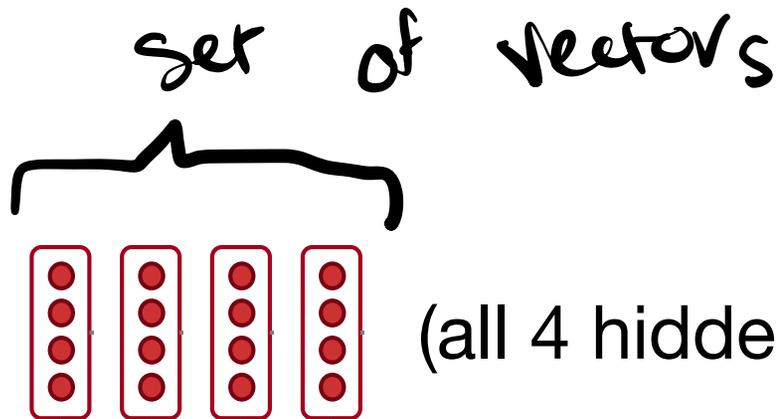
idea: what if we use multiple vectors?

1 vector



Instead of this, let's try:

the students opened their =

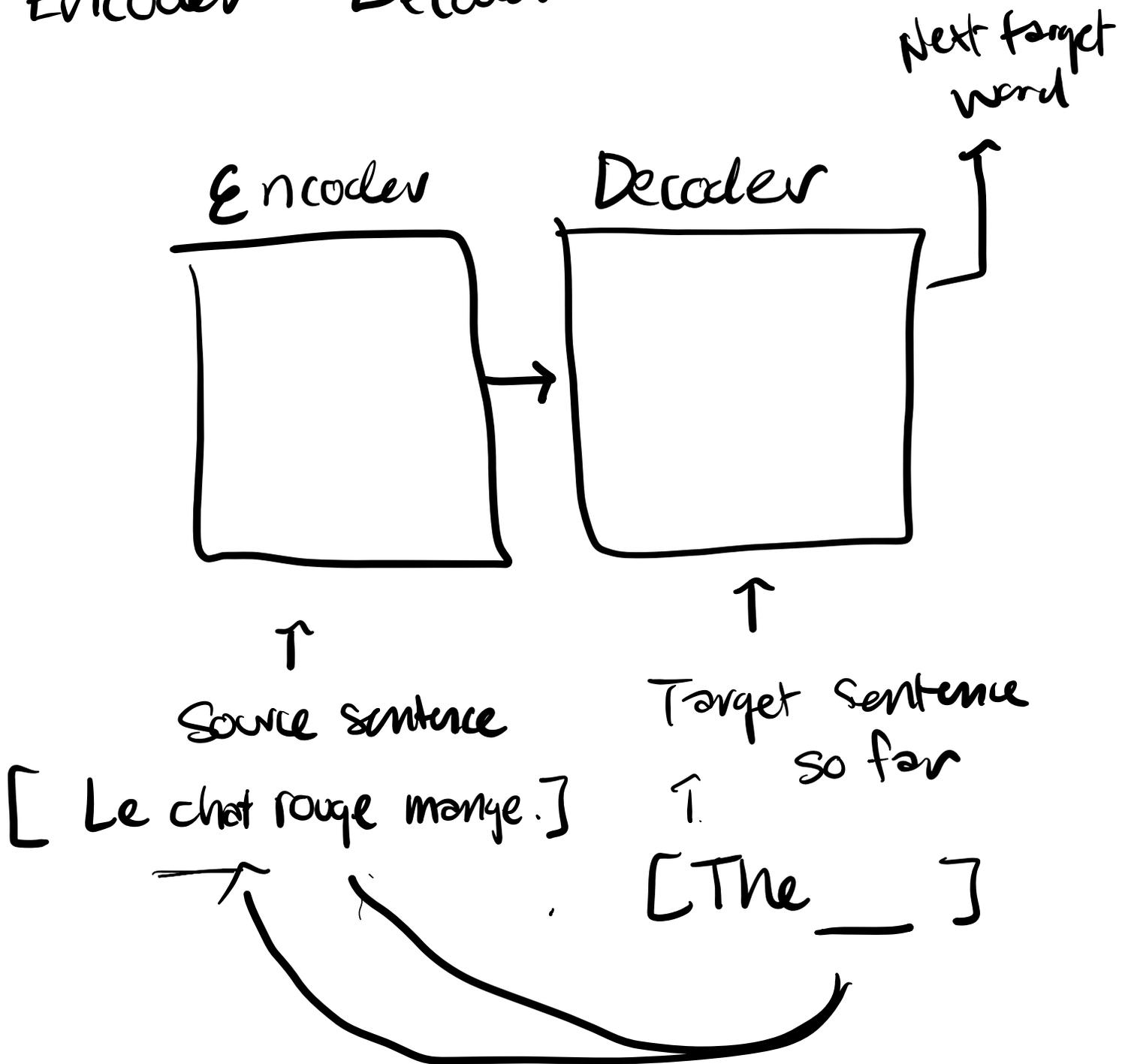


(all 4 hidden states!)

The solution: **attention**

- **Attention mechanisms** (Bahdanau et al., 2015) allow language models to focus on a particular part of the observed context at each time step
 - Originally developed for machine translation, and intuitively similar to *word alignments* between different languages

Encoder - Decoder



Attention

How does it work?

- in general, we have a single *query* vector and multiple *key* vectors. We want to score each query-key pair

in a neural language model, what are the queries and keys?

What from the past is the most useful?

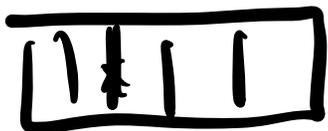
What Is Attention?

Key: o is a task-specific vector
"searching" for important things
from the past context.

Step 1) Measure how similar each x is to k .
dot product

-3.4

$$r_1 = k \cdot x_1$$



x_1 : The

2.4

$$r_2 = k \cdot x_2$$



x_2 : students

-0.8

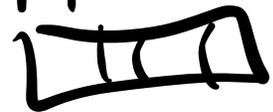
$$r_3 = k \cdot x_3$$



x_3 : quered

-1.2

$$r_4 = k \cdot x_4$$



x_4 : there

What Is Attention?

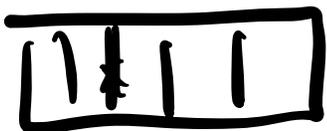
Key: o is a task-specific vector
"searching" for important things
from the past context.

Step 2): Scale scores f_i between $0-1$
(softmax)

0.01

-3.4

$$r_1 = k \cdot x_1$$



x_1 : The

0.7

2.4

$$r_2 = k \cdot x_2$$



x_2 : students

0.24

-0.8

$$r_3 = k \cdot x_3$$

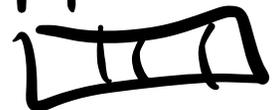


x_3 : quered

0.05

-1.2

$$r_4 = k \cdot x_4$$



x_4 : there

What Is Attention?

Key: o as a task-specific vector

"searching" for important things from the past context.

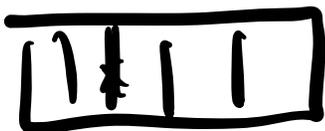
Step 3) Compute a weighted average

Output: a vector of the same dimensions
 o as a word embedding

$$\Sigma = o \cdot x$$

$$o = \begin{matrix} 0.01 \\ -3.4 \end{matrix}$$

$$r_1 = k \cdot x_1$$



x_1 : The

$$o = \begin{matrix} 0.7 \\ 2.4 \end{matrix}$$

$$r_2 = k \cdot x_2$$



x_2 : students

$$o = \begin{matrix} 0.24 \\ -0.8 \end{matrix}$$

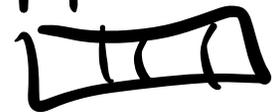
$$r_3 = k \cdot x_3$$



x_3 : quered

$$o = \begin{matrix} 0.05 \\ -1.2 \end{matrix}$$

$$r_4 = k \cdot x_4$$



x_4 : there



Vicki
@vboykis



They don't tell you this in the paper (well they do but you have to read it like 15 times)



Multiplying
a lot of vectors
a lot of times
with scaled softmax



Attention

Why dot product?

- ❖ Dot product provides a measure of similarity between keys and queries.
- ❖ But you might be wondering: *why do we want to pay attention to words that are similar to the current word?*

Why dot product?

- ❖ Dot product provides a measure of similarity between keys and queries.
- ❖ But you might be wondering: *why do we want to pay attention to words that are similar to the current word?*

Consider:

**My brother, a chemist, was late yesterday because he missed the bus.
When he arrived, he was surprised to find that his lab _____**

Why dot product?

- ❖ Dot product provides a measure of similarity between keys and queries.
- ❖ But you might be wondering: *why do we want to pay attention to words that are similar to the current word?*

Consider:

My brother, a chemist, was late yesterday because he missed the bus. When he arrived, he was surprised to find that his lab _____

lab



lab



lab



lab

Lab Assignment

Review available resources on the web:
<http://www.sonomastate.edu/users/ffrahman/sonoma/projects/ca/labview/index.htm>

In-class Lab 1: Introduction to LabVIEW

A- Read <http://www.sonomastate.edu/users/ffrahman/sonoma/projects/ca/labview/index.htm>

B- Follow the steps up to **Profile Tool** Section. In this lab you create a VI to calculate sum and average of several numbers.

C- When you complete the code show it to the instructor.

D- If you have extra time, you can start working on the homework (see below).

Homework:
The homework assignment must be done individually. If you copy the program from another student, both of you will receive **zero** for this assignment.

Watch the video (30 min. only):
<http://www.sonomastate.edu/users/ffrahman/sonoma/projects/ca/labview/snap/default.htm>

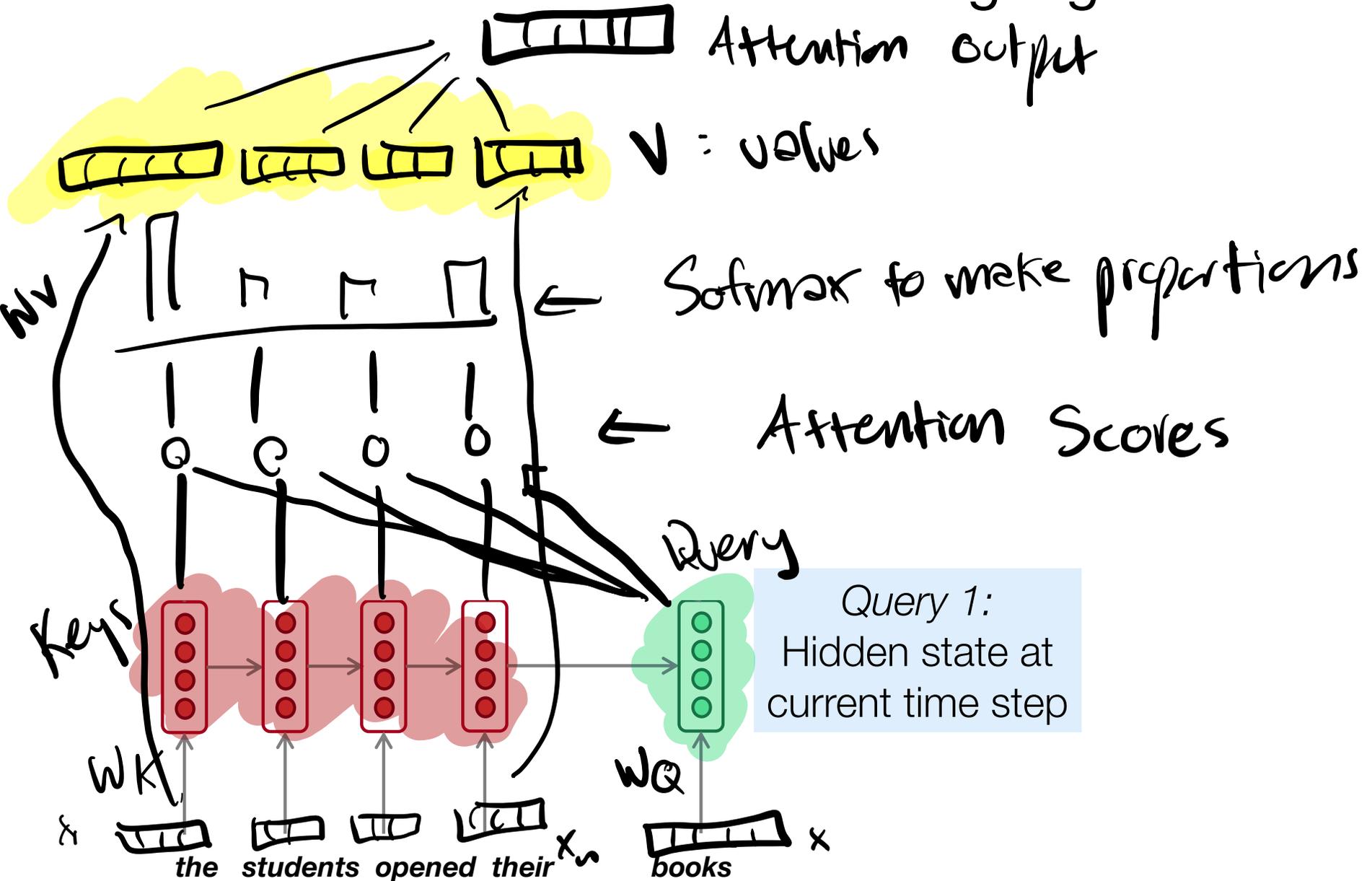
Assignment 1:
Create a simple program that can convert a temperature from the Celsius scale to the Fahrenheit scale: <http://www.sonomastate.edu/users/ffrahman/sonoma/projects/ca/labview/index.htm>. Take a snap shot of the Front panel and Diagram. Place the figures in the table below.

Figure 1. Front Panel VI for Temperature Converter.
Figure 1. Block Diagram for Temperature Converter.

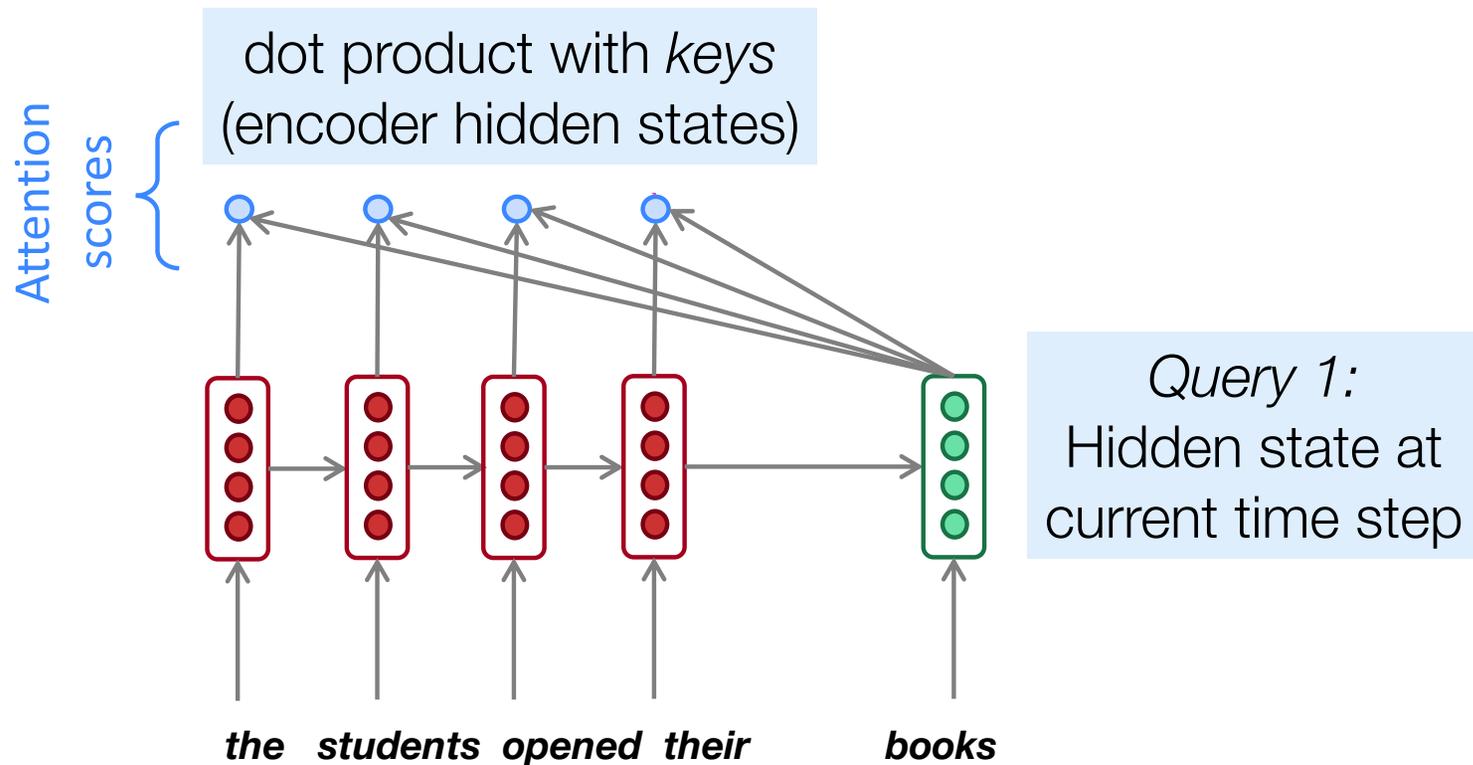
Assignment 2:
Change the code below such that the program generates random numbers between 1-10. Make sure your program works properly. Test it for several values. Take a snap shot of the Front panel and Diagram. Place the figures in the table below.

Figure 1. Front Panel VI for Random Number Generator.
Figure 1. Block Diagram for Random Number Generator.

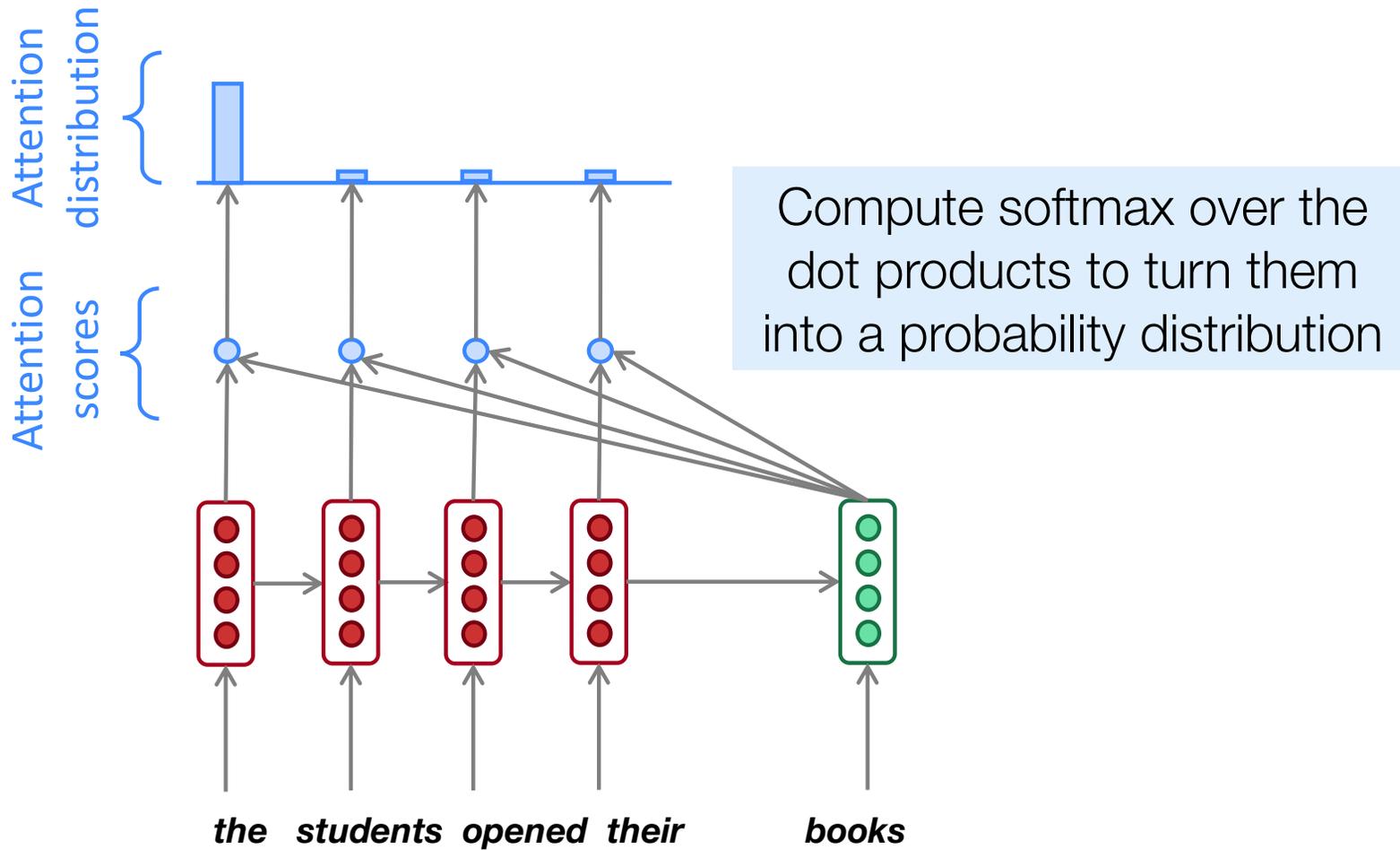
Attention mechanisms in neural language models



Attention mechanisms in neural language models

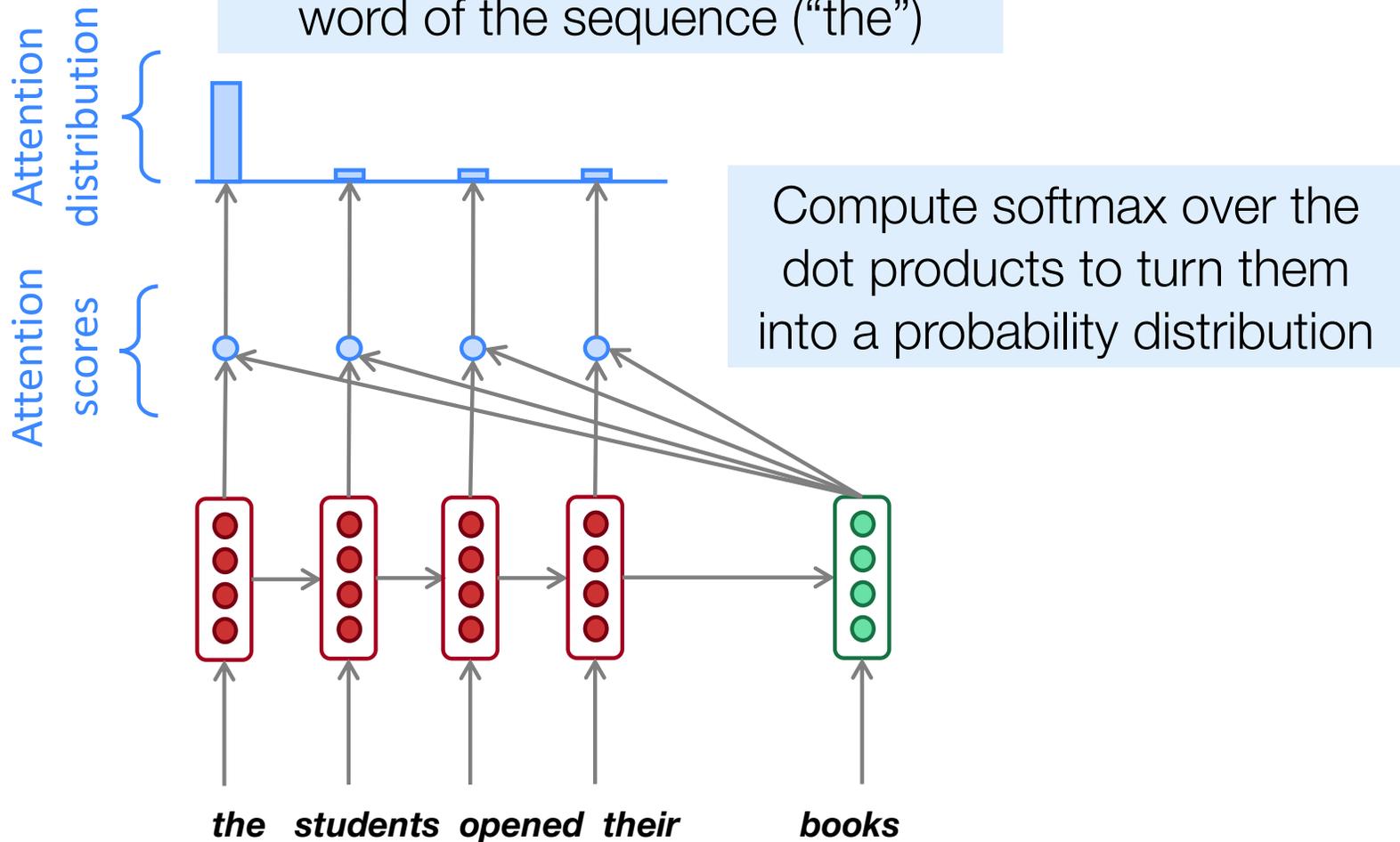


Attention mechanisms in neural language models

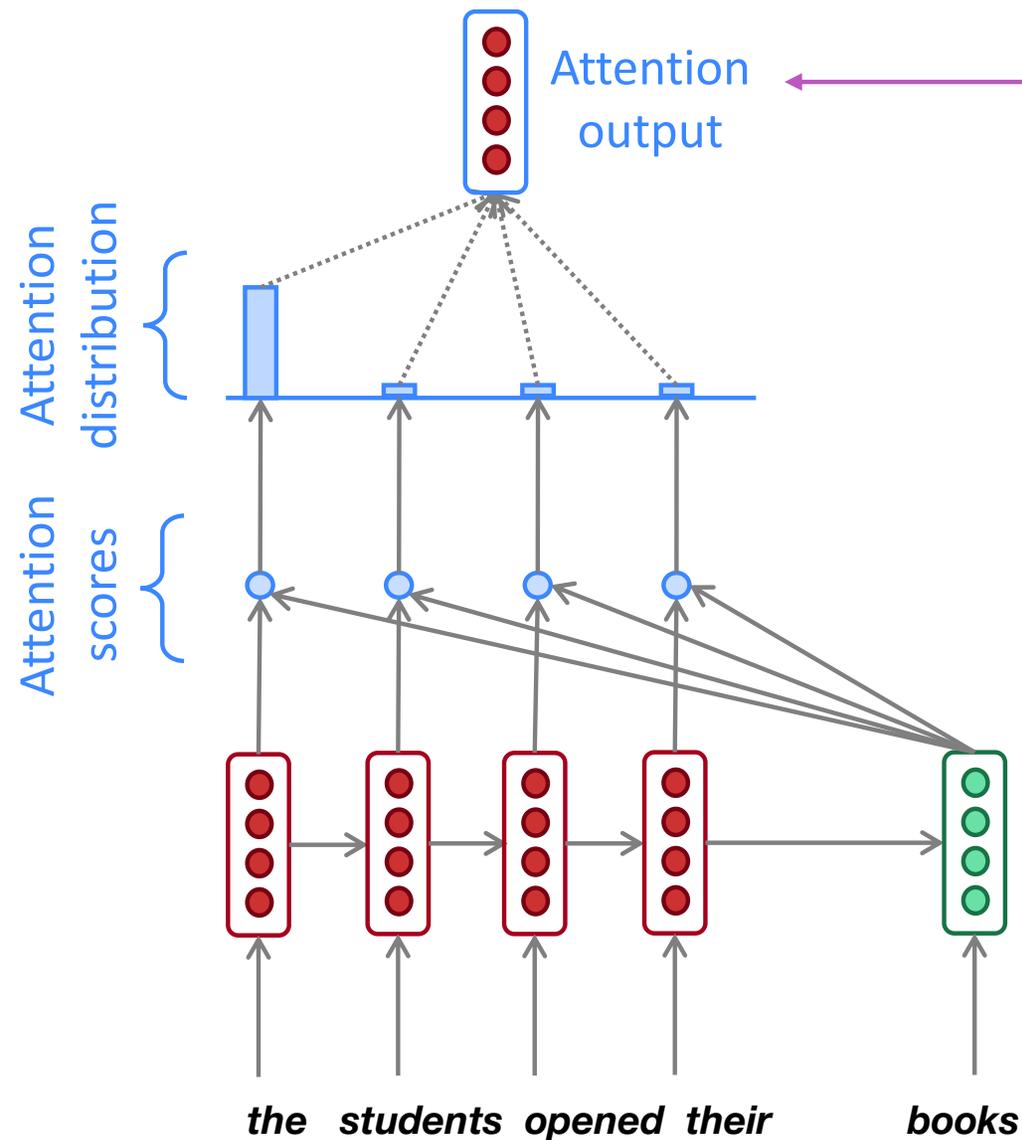


Attention mechanisms in neural language models

At this time step, the attention distribution is focused on the first word of the sequence (“the”)



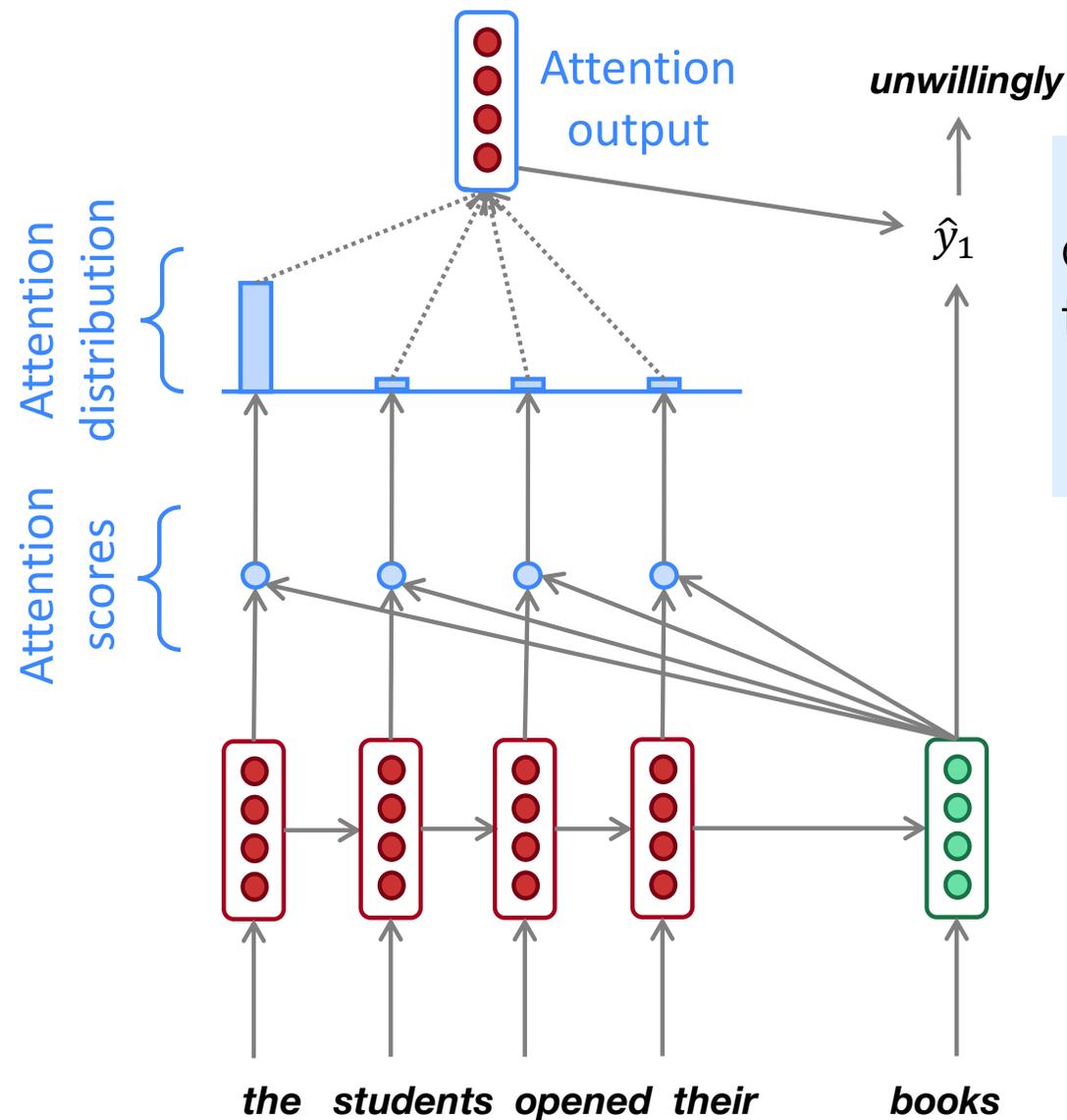
Attention mechanisms in neural language models



We use the attention distribution to compute a weighted average of the hidden states.

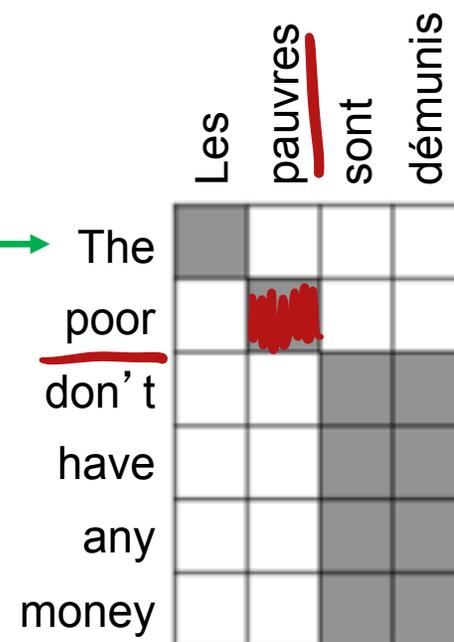
Intuitively, the resulting attention output contains information from hidden states that received high attention scores

Sequence-to-sequence with attention



Concatenate (or otherwise compose) the attention output with the current hidden state, then pass through a softmax layer to predict the next word

- Attention **solves the bottleneck problem**
 - Attention allows decoder to look directly at source; bypass bottleneck
- Attention **helps with vanishing gradient problem**
 - Provides shortcut to faraway states
- Attention provides **some interpretability**
 - By inspecting attention distribution, we can see what the decoder was focusing on
 - We get **alignment for free!**
 - This is cool because we never explicitly trained an alignment system
 - The network just learned alignment by itself



Many variants of attention

- Original formulation: $a(\mathbf{q}, \mathbf{k}) = w_2^T \tanh(W_1[\mathbf{q}; \mathbf{k}])$

- Bilinear product: $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T W \mathbf{k}$

Luong et al., 2015

- Dot product: $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k}$

Luong et al., 2015

- Scaled dot product: $a(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^T \mathbf{k}}{\sqrt{|\mathbf{k}|}}$

Vaswani et al., 2017

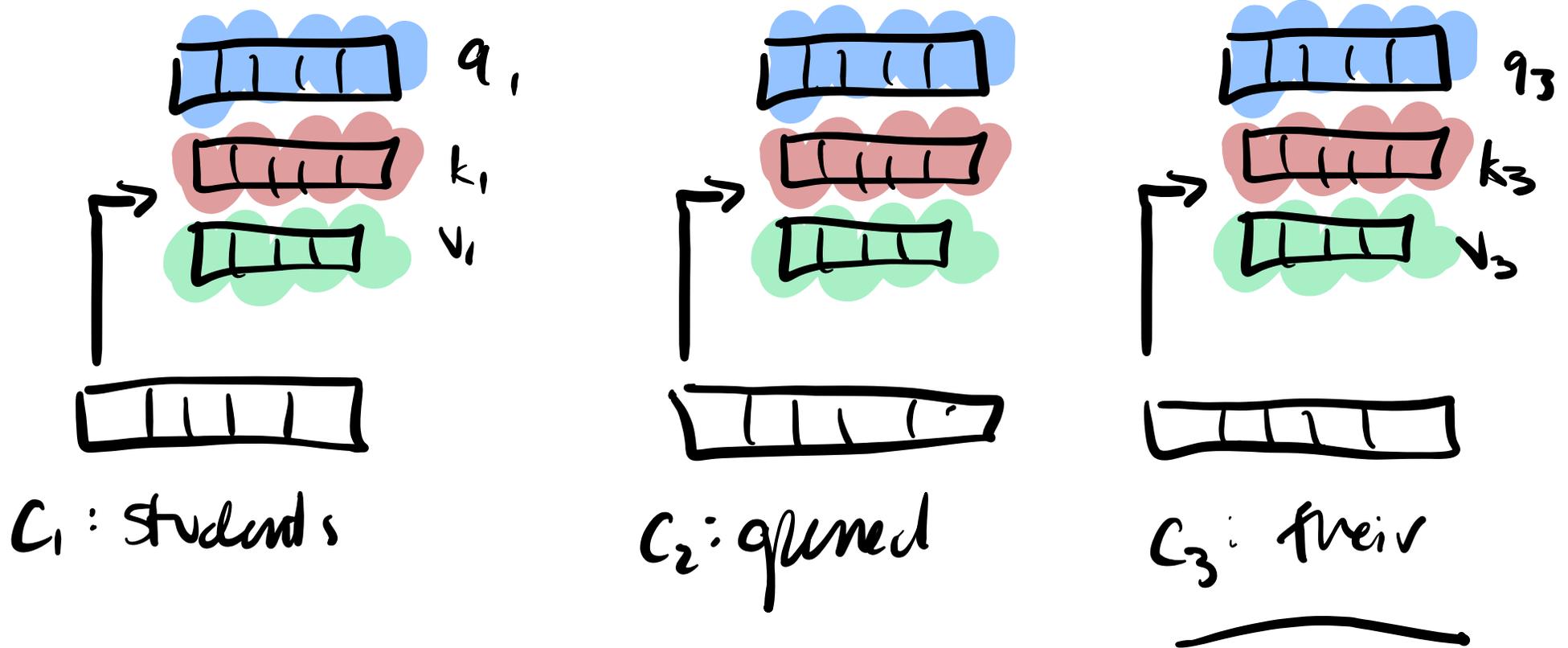
Self-Attention

Self-Attention Motivation: efficiency by parallelization

$$q_i = f(W_q C_i) \quad k_i = f(W_k C_i) \quad v_i = f(W_v C_i)$$

1. Take the dot product between q_3 & every k

$$\langle q_3 k_1 \quad q_3 k_2 \quad \underline{q_3 k_3} \rangle$$

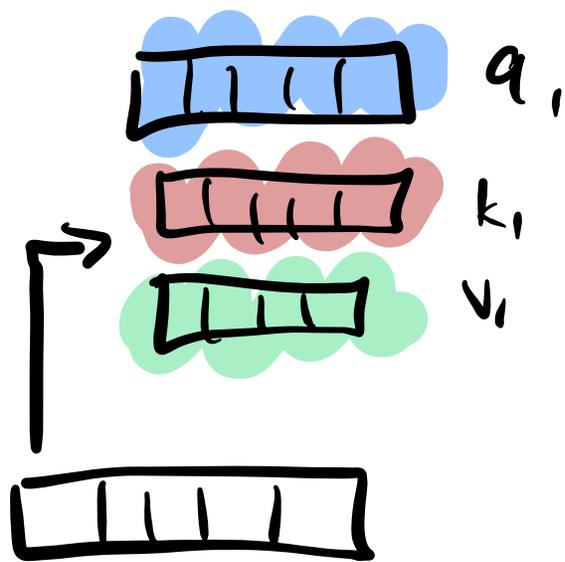


Self-Attention Motivation: efficiency & parallelization

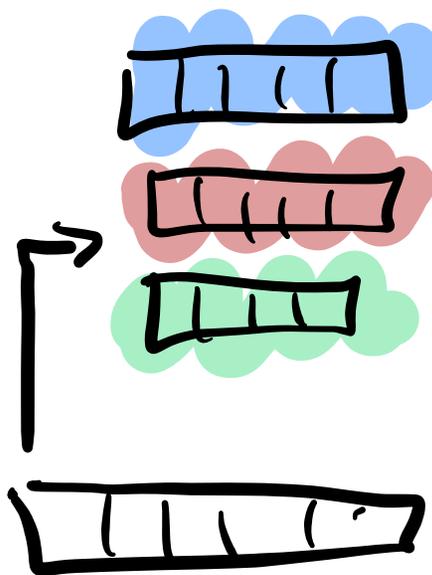
Step 2): Softmax to get a distribution



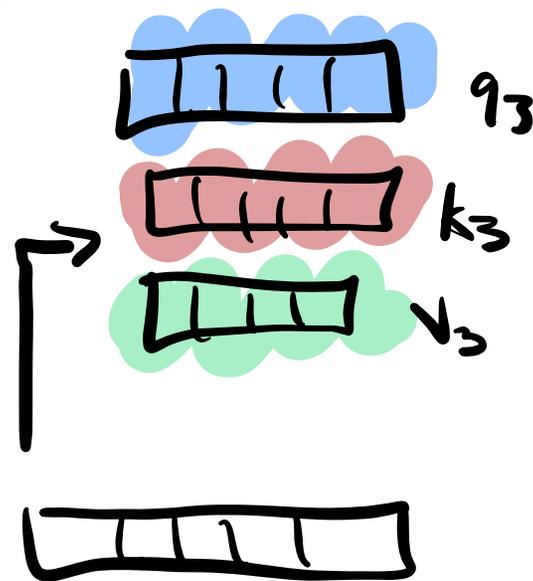
$\langle 0.3 \ 0.5 \ 0.2 \rangle$



C_1 : students



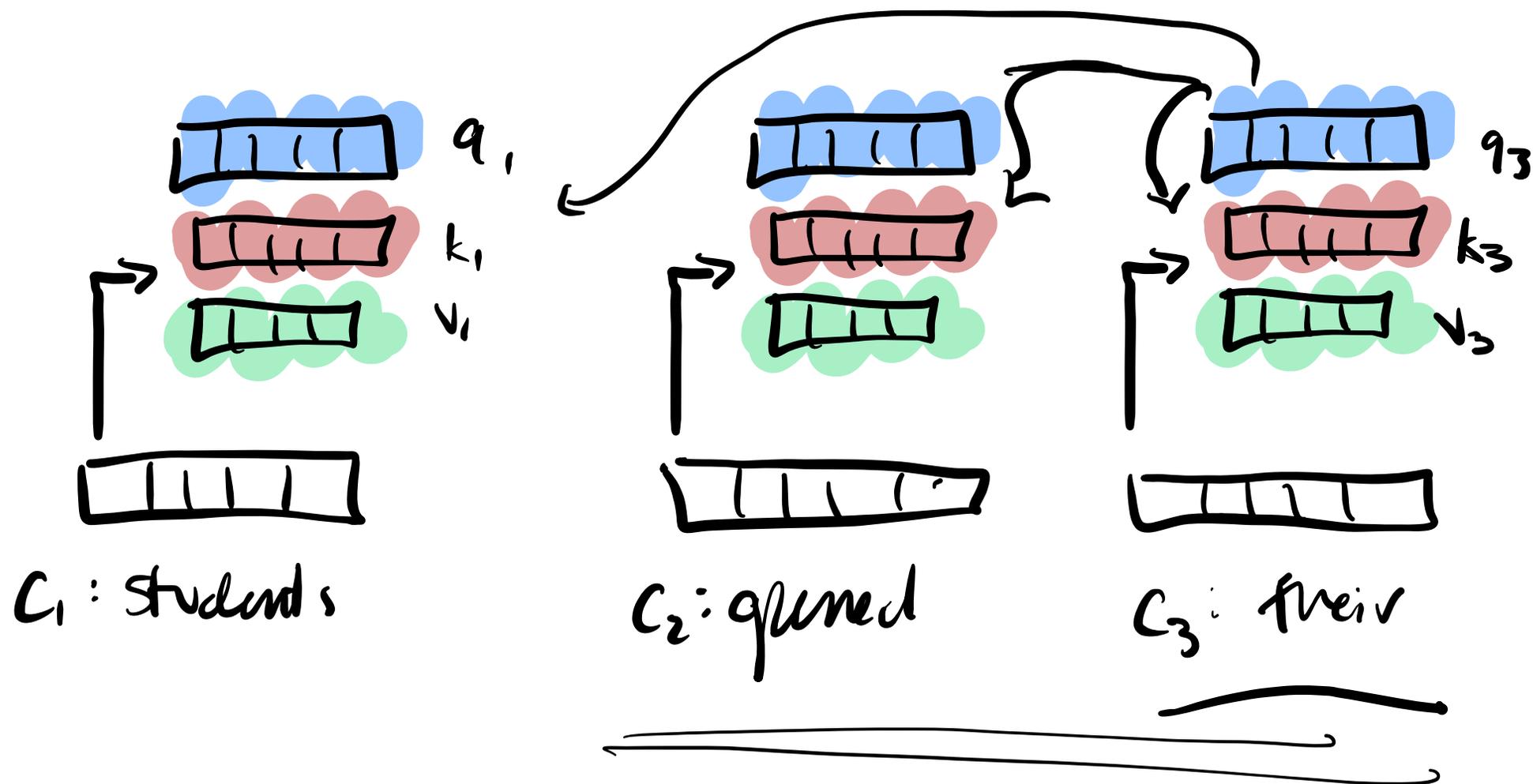
C_2 : opened



C_3 : freir

Self-Attention Motivation: efficiency & parallelization
step 3): Calculate weighted average on values.

$$h_3 = 0.3v_1 + 0.5v_2 + 0.2v_3$$

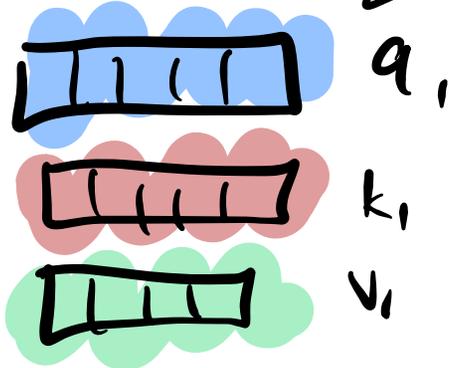
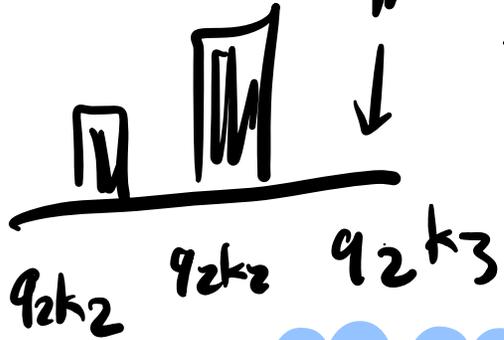


Self-Attention Motivation: efficiency & parallelization

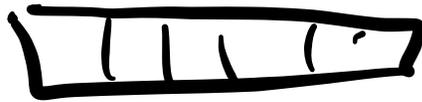
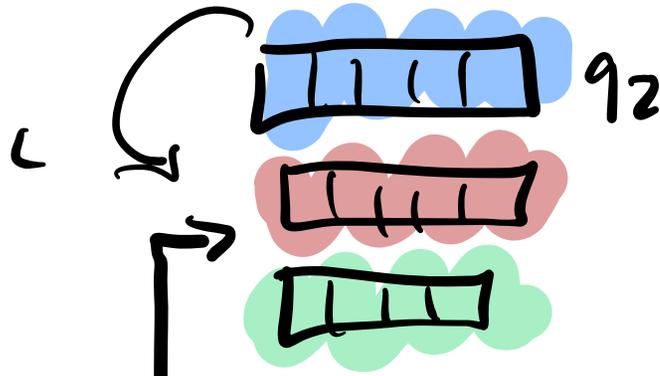
step 3): Calculate weighted average on values.

must be zero!

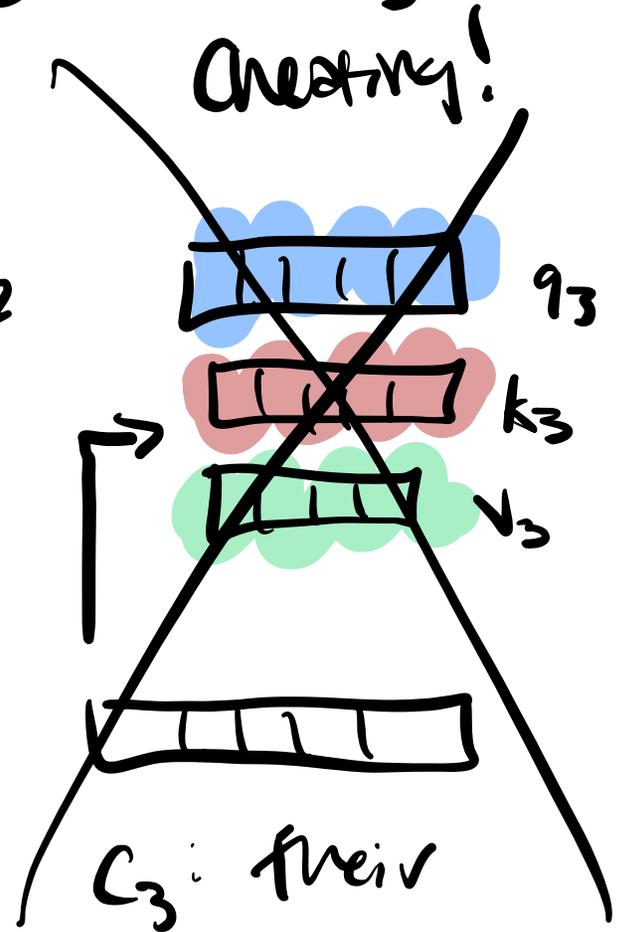
$$h_3 = 0.3v_1 + 0.5v_2 + 0.2v_3$$



C_1 : students



C_2 : opened



C_3 : freier