
CS 333:
Natural Language
Processing

Fall 2025

Prof. Carolyn Anderson
Wellesley College

Reminders

- ♦ Tuesday -> Monday cycle for final two assignments:
 - HW 7 is due on Monday, 11 / 17
 - HW 8 will be released on Tuesday, 11 / 18 and due on Monday, 11 / 24
- ♦ My next help hours: Thursday 4-5
- ♦ Two CS Colloquia upcoming!

HW 7: Evaluating LLMs

In this assignment, you will practice evaluating LLMs. You'll work with Llama 1B, and test its capabilities on two tasks: Pig Latin decoding, and common-sense reasoning.

The assignment also asks you to do some setup work for HW 8, and to describe your final project topic.

WELLESLEY CS COLLOQUIUM

Professor Lindsey D. Cameron
Wharton School, University of Pennsylvania



Resocializing the Platform: Patchwork Embeddedness and How Workers ReConstitute Digital Platforms

14 NOV 2025 | 3:30 PM | SCI-H105

Snacks will be provided!



cs129@wellesley.edu

Accessibility and Disability:
accessibility@wellesley.edu

Two ML talks for the price of one:

*Co-designing Tools to Measure Student Learning with
Machine Learning and Science Education Research*

Dr. Kaitlin Gili

*Subgroup Validity in Machine Learning for
Echocardiogram Data*

Cynthia Feeney

Thursday Nov. 20 at 12:45-1pm in H-105

HW 6 Curiosity Points

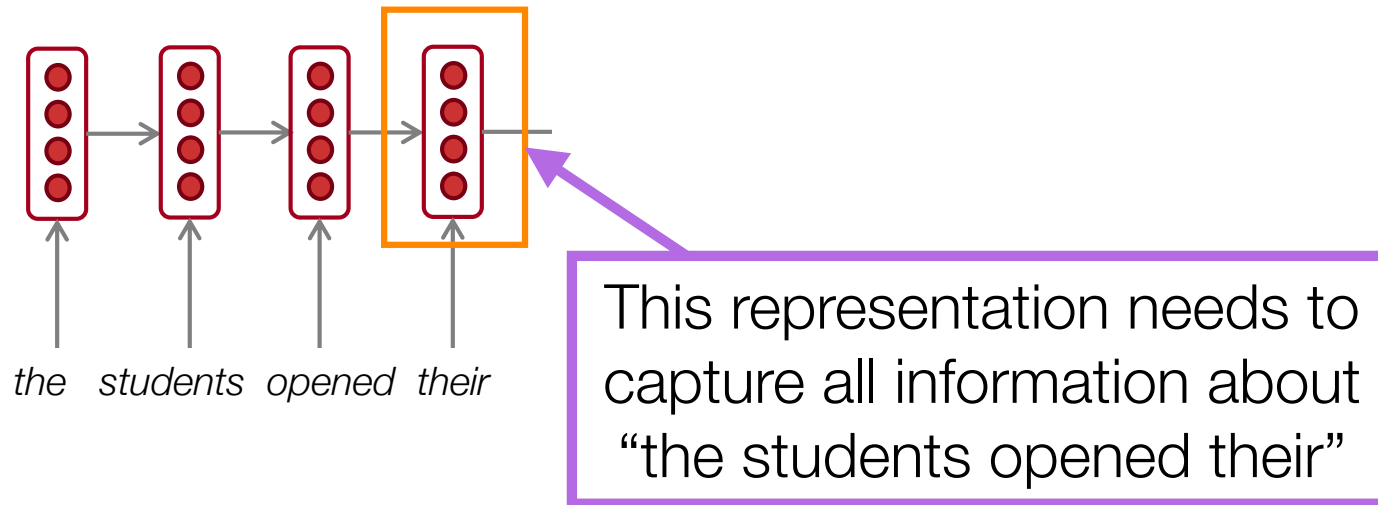
- ♦ Model experimentation:
 - Comparison of non-linear activation functions
 - Experiments with model architecture changes
- ♦ Analysis:
 - Dimensionality reduction and visualization of learned embeddings
 - Visualization of model performance
- ♦ Literature exploration:
 - Read related work on non-linear activation functions, music understanding, neural network models at Google, class imbalance, and BERT models
- ♦ New tasks:
 - Scraping lyrics from new artists
 - Lyric guessing game
 - Song year prediction task

Recap

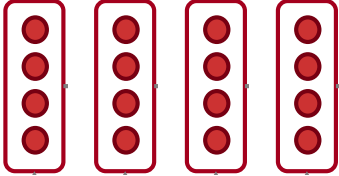
“you can’t cram the meaning
of a whole %&@#&ing
sentence into a single
\$*(&@ing vector!”

— Ray Mooney (NLP professor at UT Austin)

idea: what if we use multiple vectors?



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


What Is Attention?

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


Step 1) Measure the ^{dot product} similarity between query & key

$$r_1 = k_1 \cdot q$$



the

$$r_2 = k_2 \cdot q$$


students

$$r_3 = k_3 \cdot q$$


opened

$$r_4 = k_4 \cdot q$$


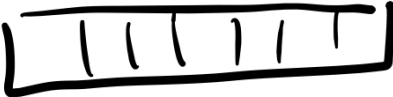
them


Keys: representations of past context

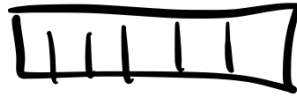
What Is Attention?


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

Step 2: ^{softmax} Scale scores to between 0-1

0.01
-3.4
 $r_1 = k_1 \cdot q$

the

0.7
2.4
 $r_2 = k_2 \cdot q$

students

0.24
-0.8
 $r_3 = k_3 \cdot q$

opened

0.05
-1.2
 $r_4 = k_4 \cdot q$

them

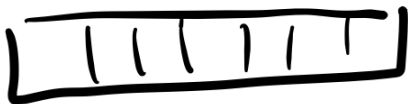
Keys: representations of past context


What Is Attention?

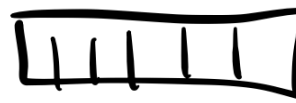
Query: a task-specific vector that is "searching" for important things from the past context.


Step 3: Compute a weighted average of embeddings

Output: a vector of the same dimensions as each word embedding.

0.01
 -3.4
 $r_1 = k_1 \cdot q$

the

0.7
 2.4
 $r_2 = k_2 \cdot q$

students

0.24
 -0.8
 $r_3 = k_3 \cdot q$

opened

0.05
 -1.2
 $r_4 = k_4 \cdot q$

there

Keys: representations of past context



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:

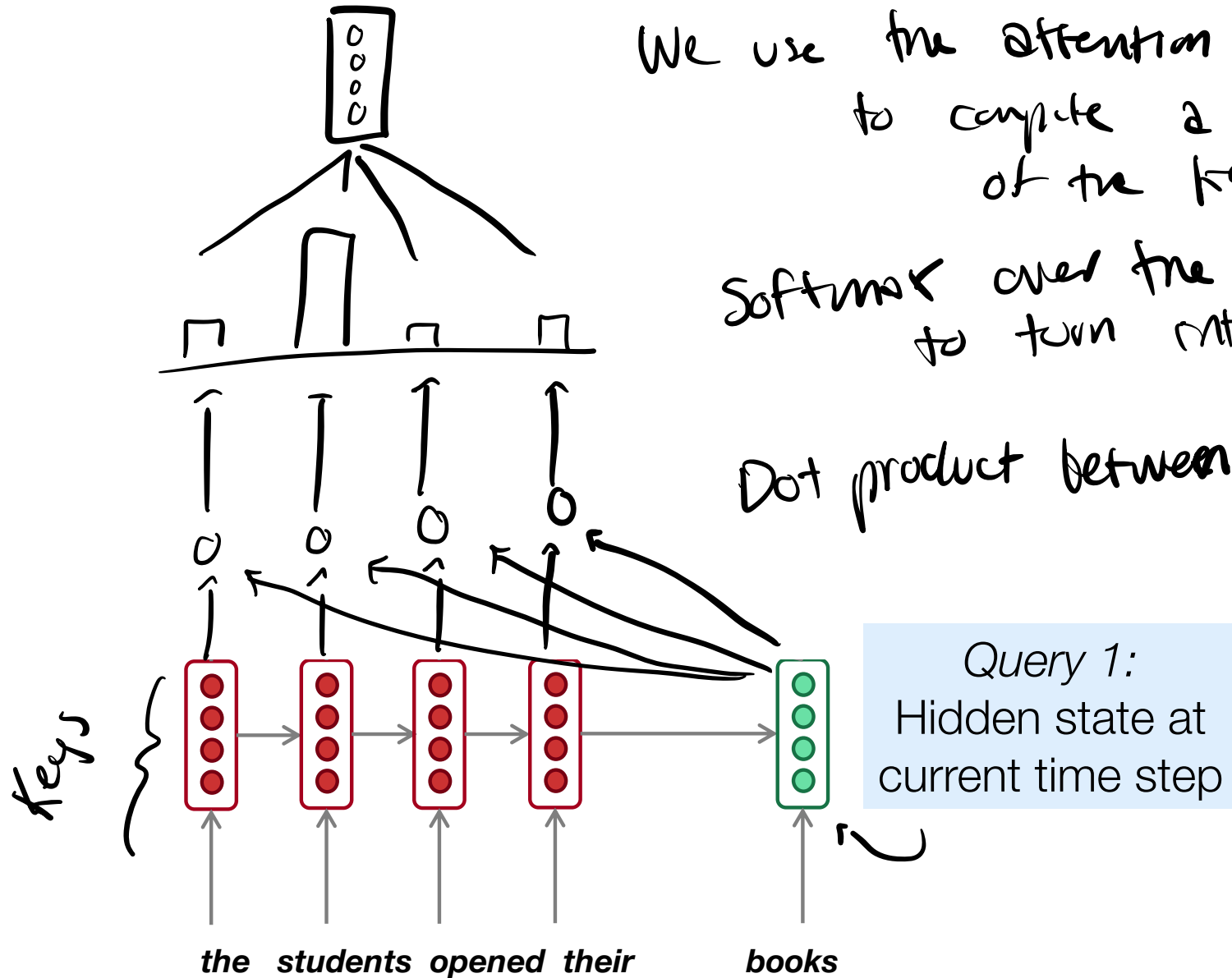
**My brother, a chemist, was late yesterday because he missed the bus.
When he arrived, he was surprised to find that his lab _____**

Attention mechanisms in neural language models

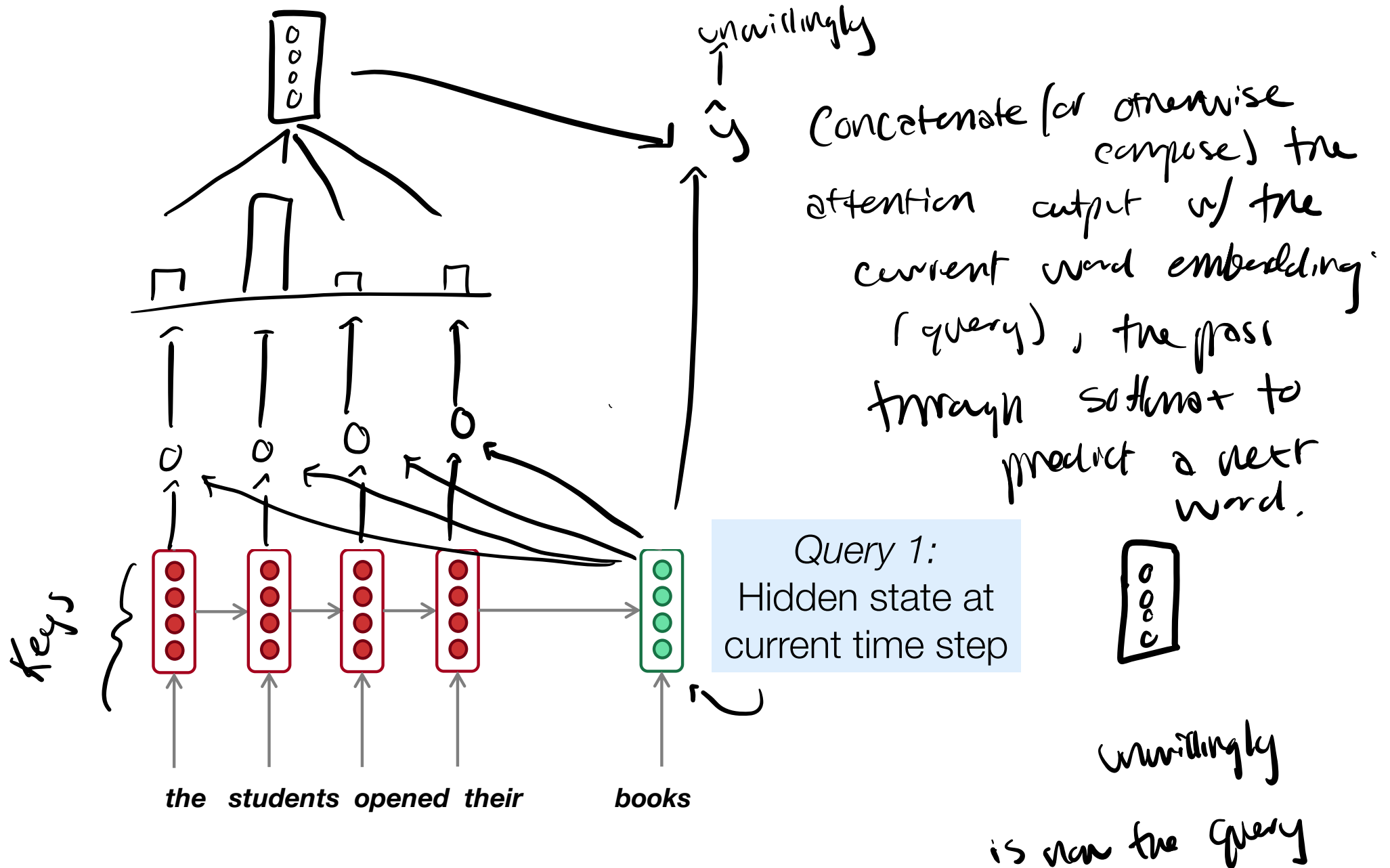
We use the attention distribution to compute a weighted average of the keys

Softmax over the dot products to turn into a probability dist.

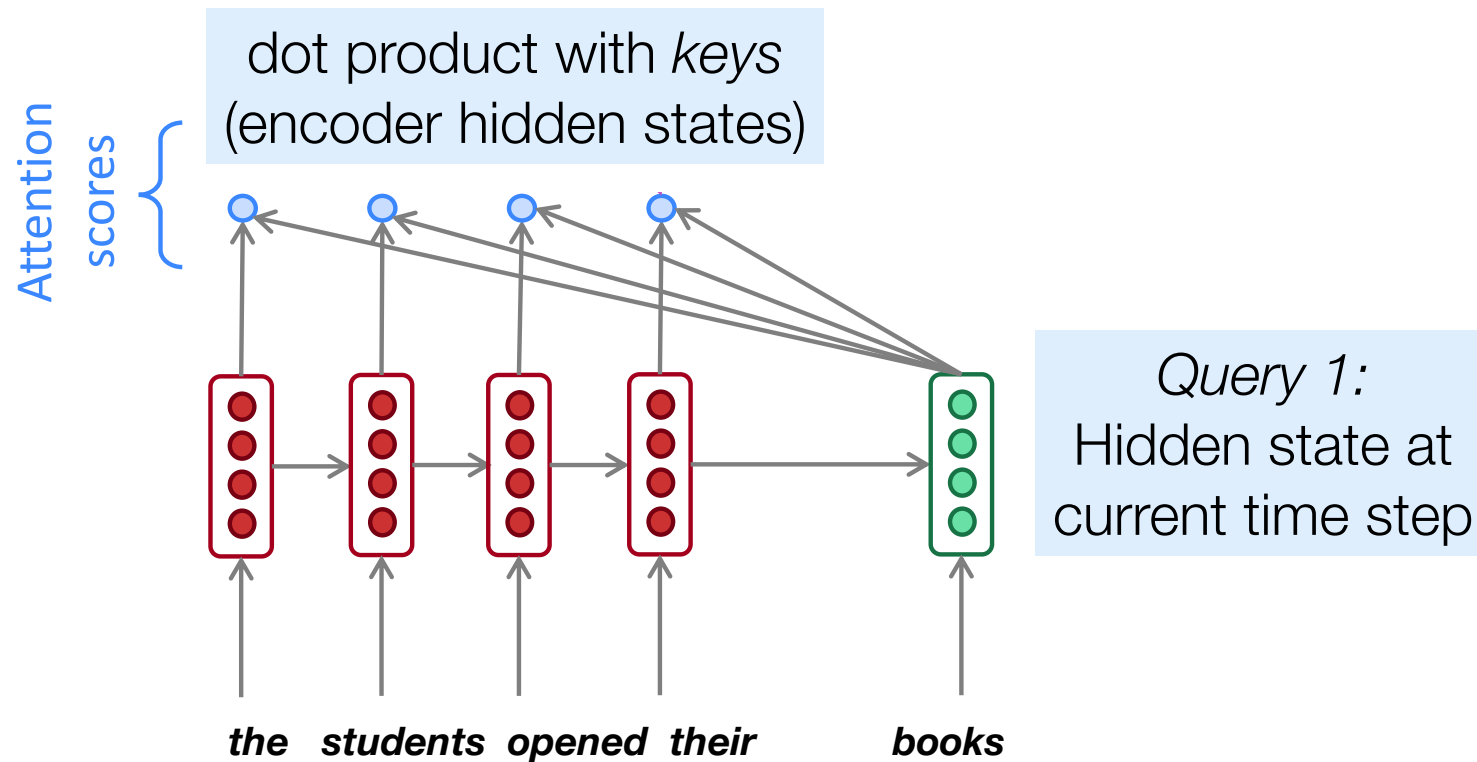
Dot product between $q \cdot k$



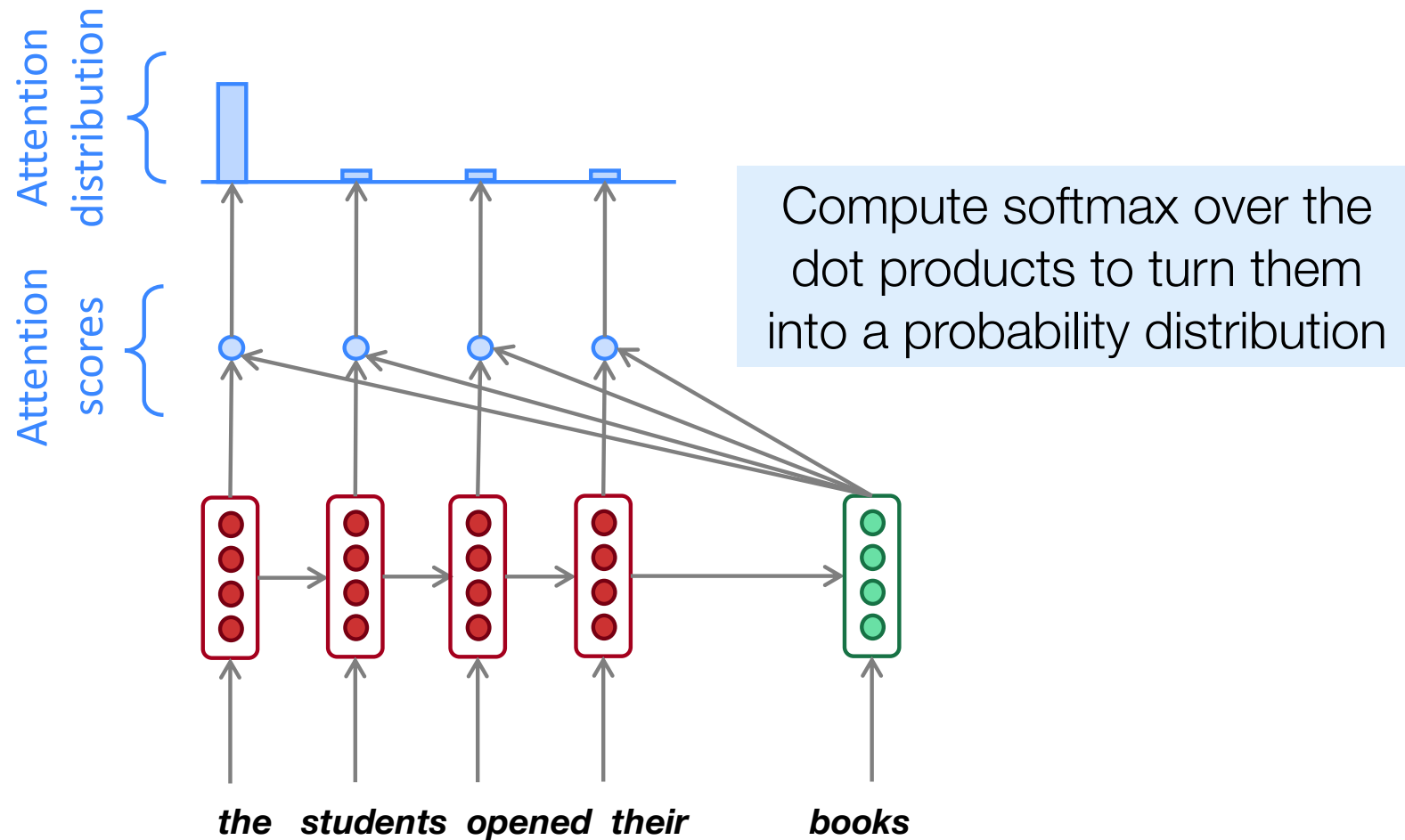
Attention mechanisms in neural language models



Attention mechanisms in neural language models



Attention mechanisms in neural language models

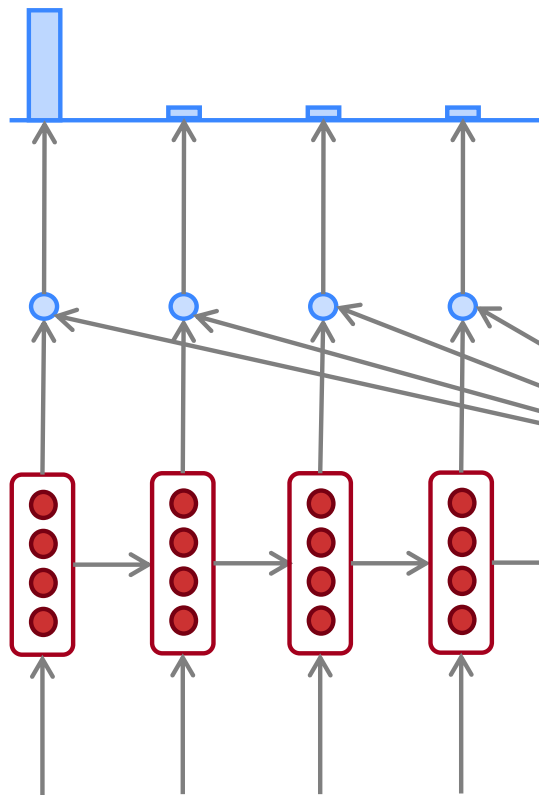


Attention mechanisms in neural language models

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

Attention
distribution

Attention
scores

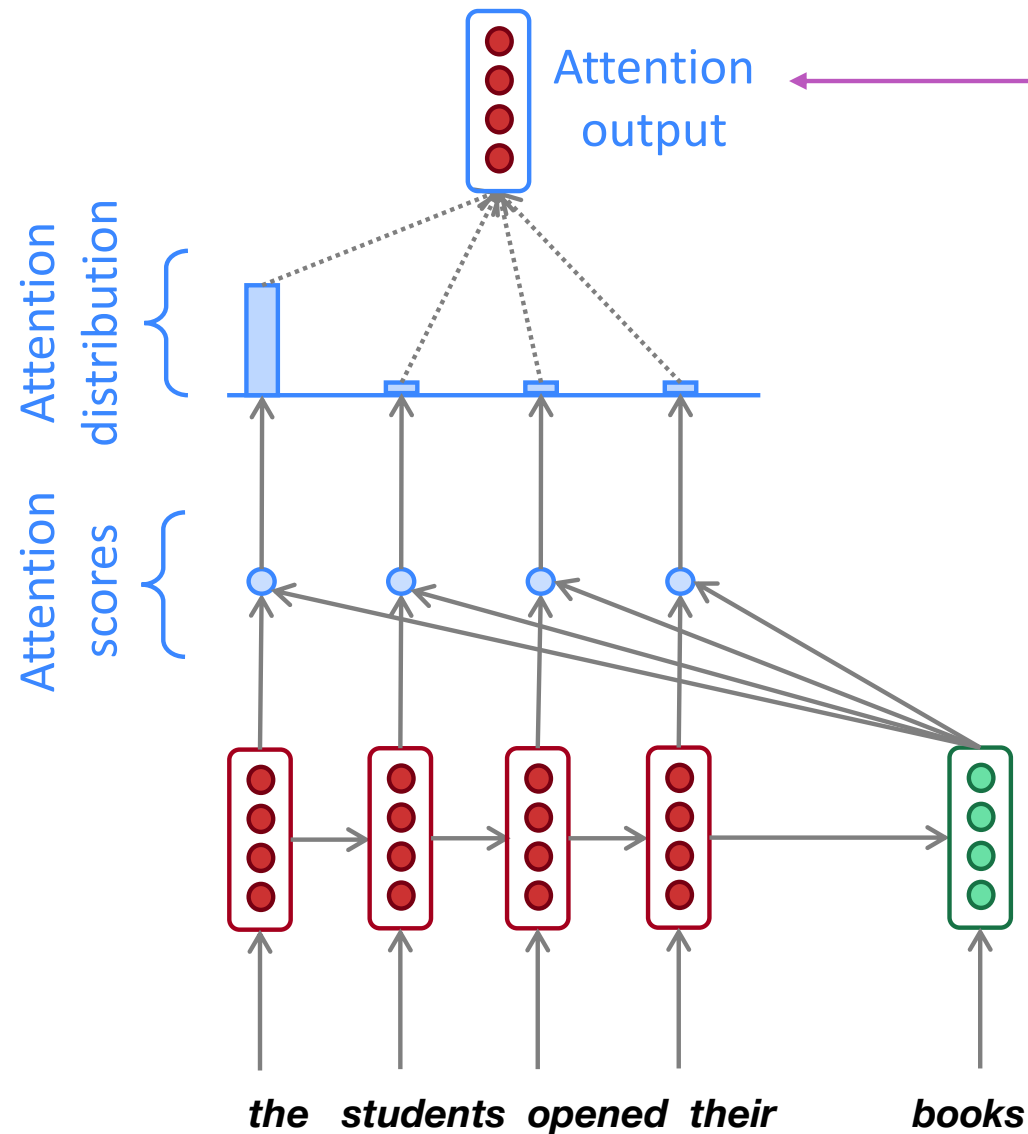


Compute softmax over the dot products to turn them into a probability distribution

the students opened their

books

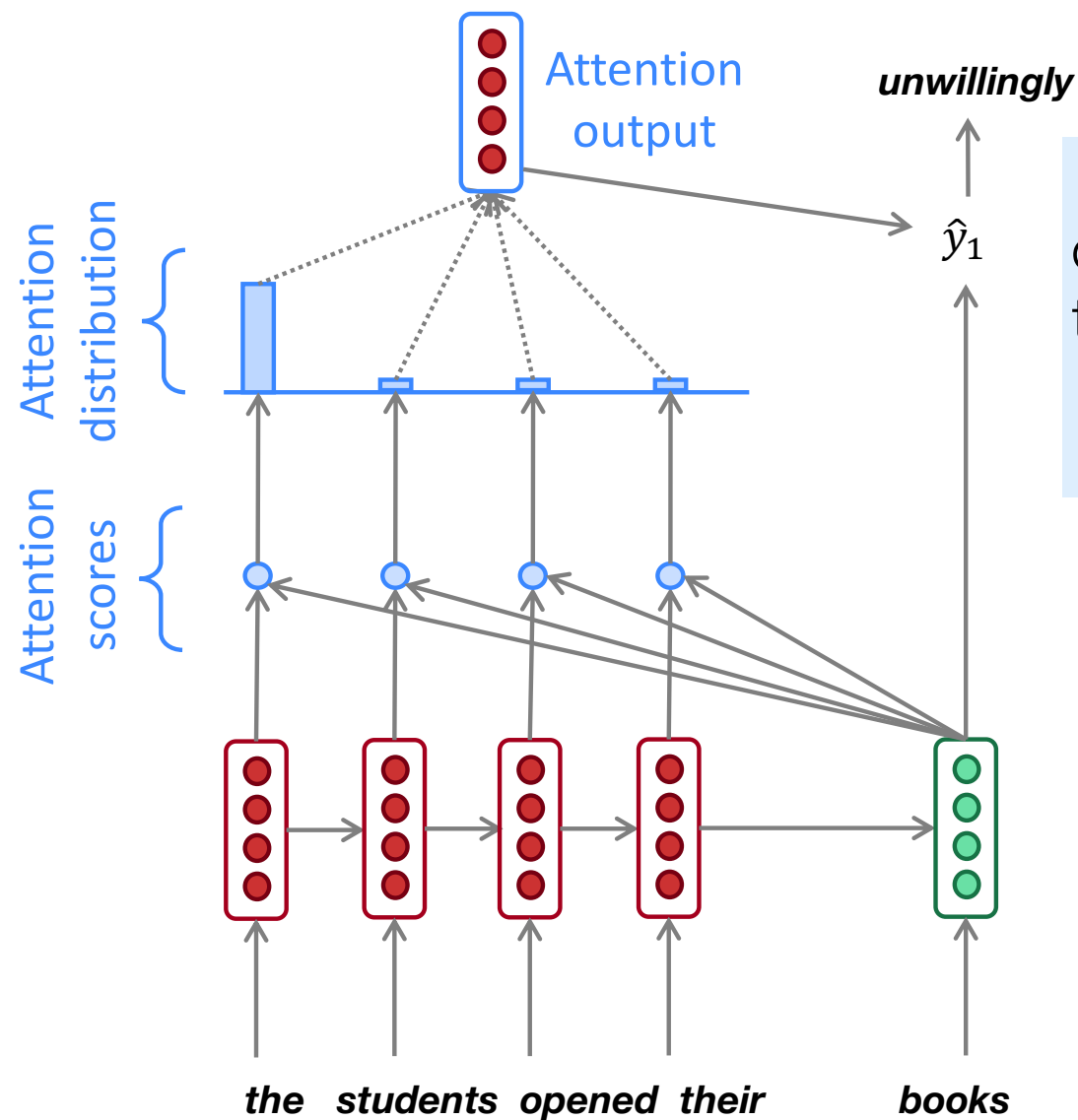
Attention mechanisms in neural language models



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

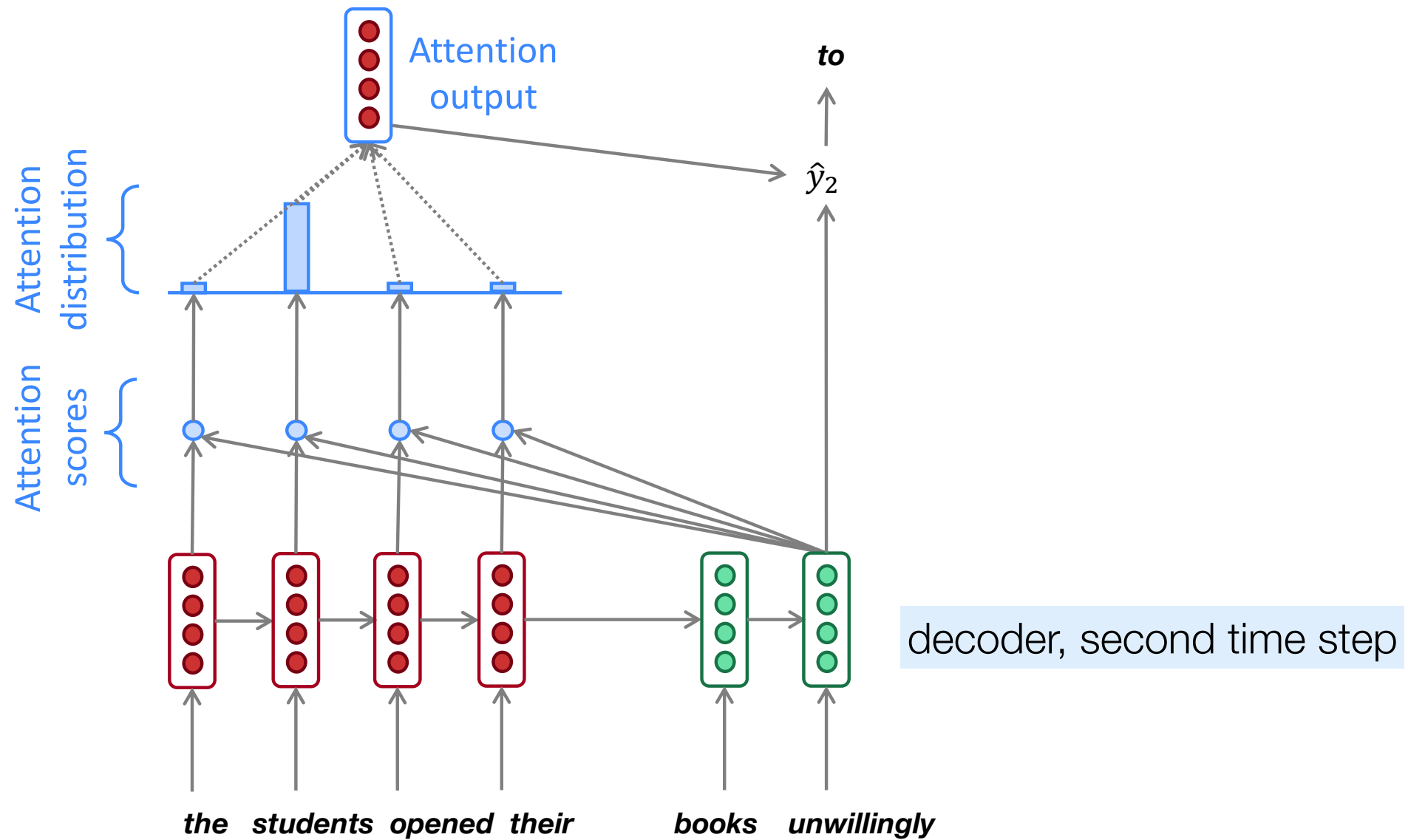
Intuitively, the resulting attention output contains information from hidden states that received high attention scores

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

Sequence-to-sequence with attention



- 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



	Les	pauvres	sont	démunis
The	■			
poor		■		
don't			■	■
have			■	■
any			■	■
money			■	■

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

Parallelizing Self-Attention

