

---

CS 232:  
Artificial Intelligence

Fall 2023

---

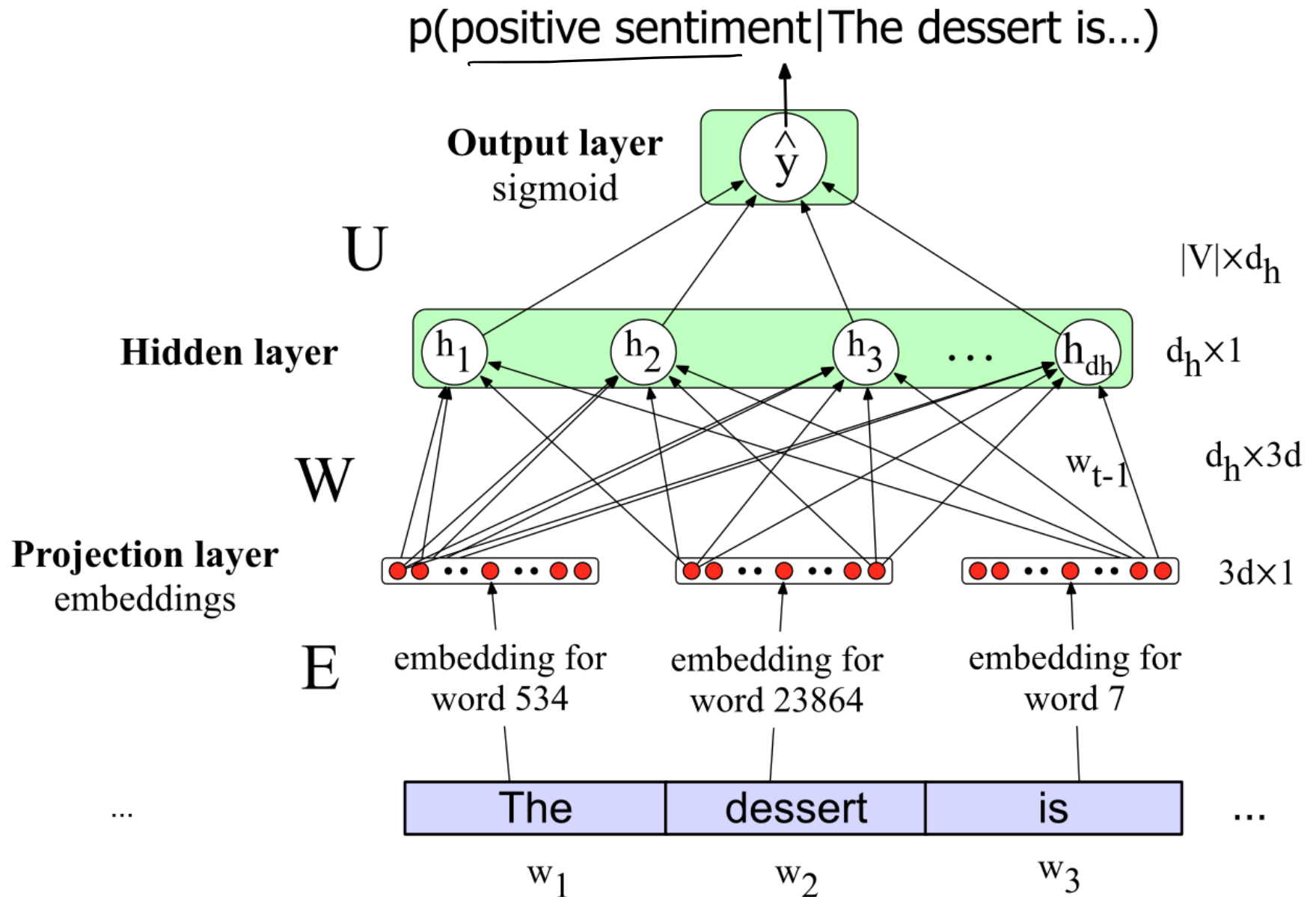
Prof. Carolyn Anderson  
Wellesley College

# Reminder

---

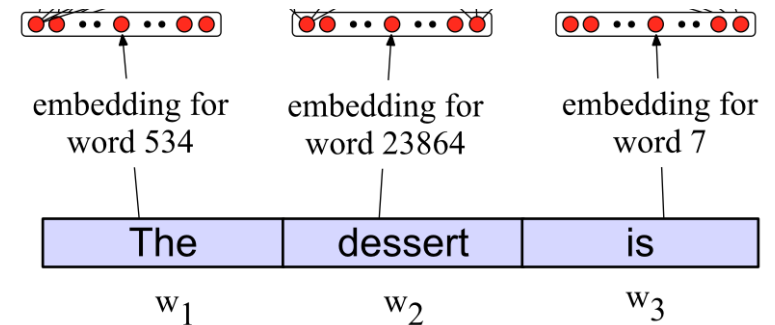
- ❖ I'm out of town for a conference most of the week
- ❖ No help hours on Thursday!
- ❖ Lyra has help hours on Wednesday
- ❖ First Gen in CS lunch this Wednesday
- ❖ CS Colloquium next Wednesday

# Neural Net Classification with embeddings as input features!



# Issue: texts come in different sizes

This assumes a fixed size length (3)!



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
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

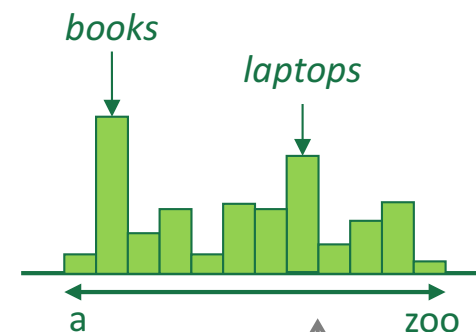
# Recurrent Neural Network

## A RNN Language Model

output distribution

$$\hat{y} = \text{softmax}(W_2 h^{(t)} + b_2)$$

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



hidden states

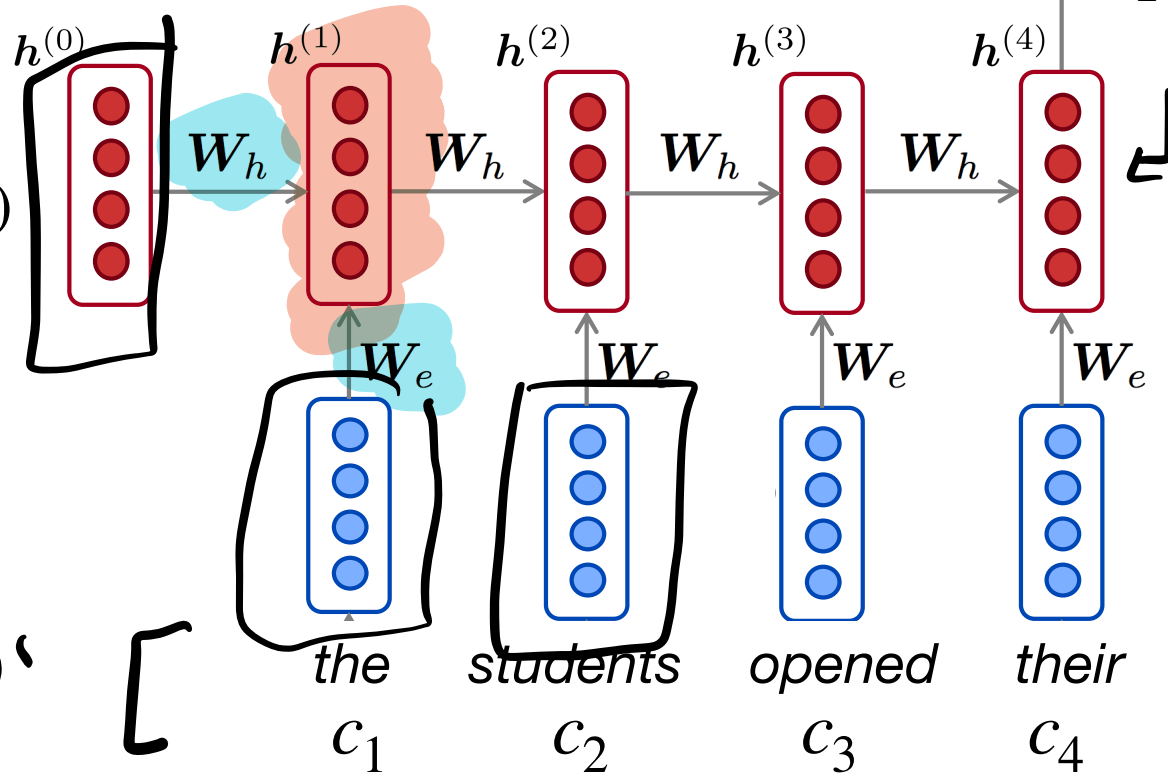
$$h^{(t)} = f(W_h h^{(t-1)} + W_e c_t + b_1)$$

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

word embeddings

$c_1, c_2, c_3, c_4$

timesteps [

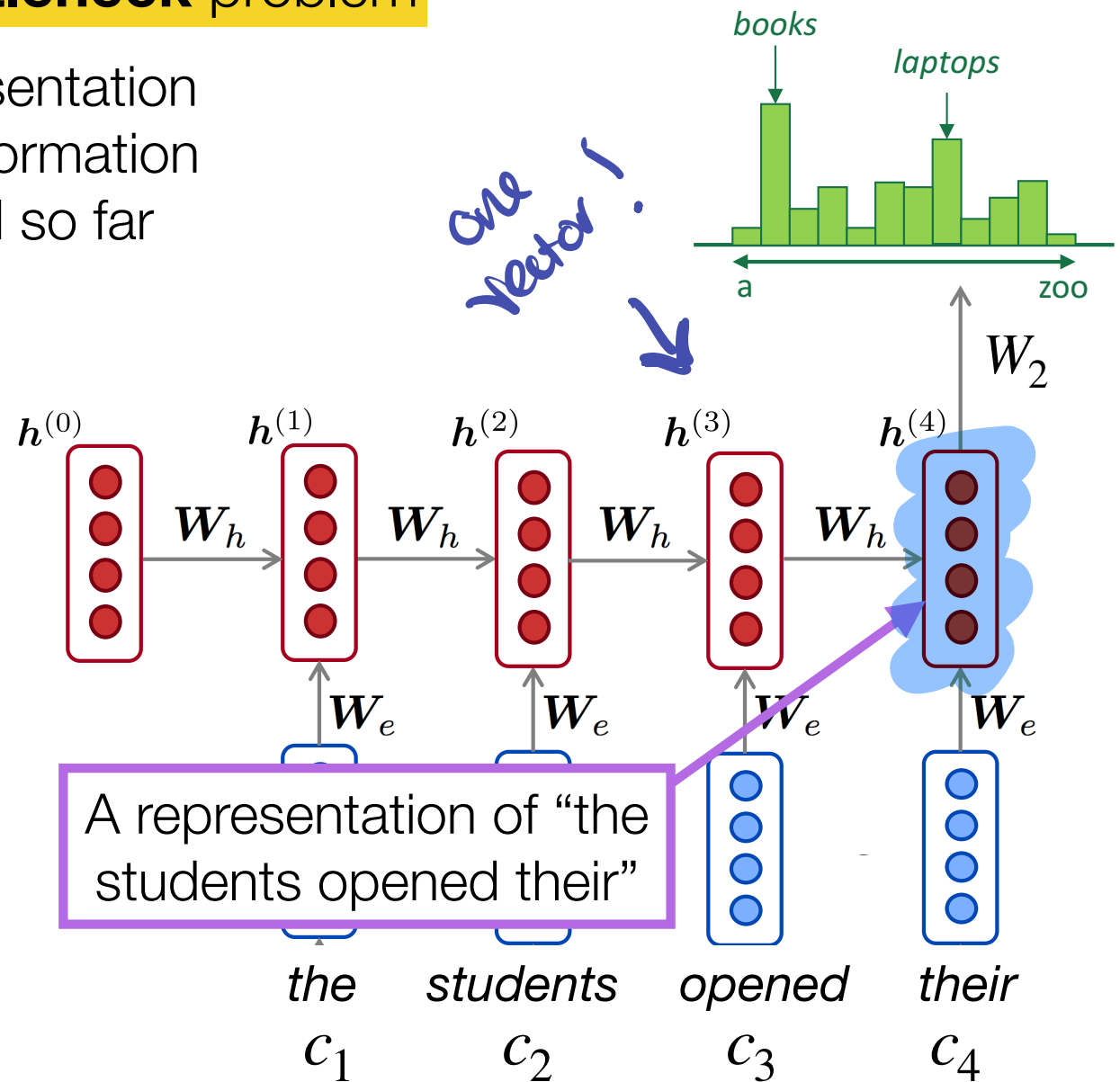


# RNNs suffer from a **bottleneck** problem

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

This becomes difficult especially with longer sequences

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



## why is this good?

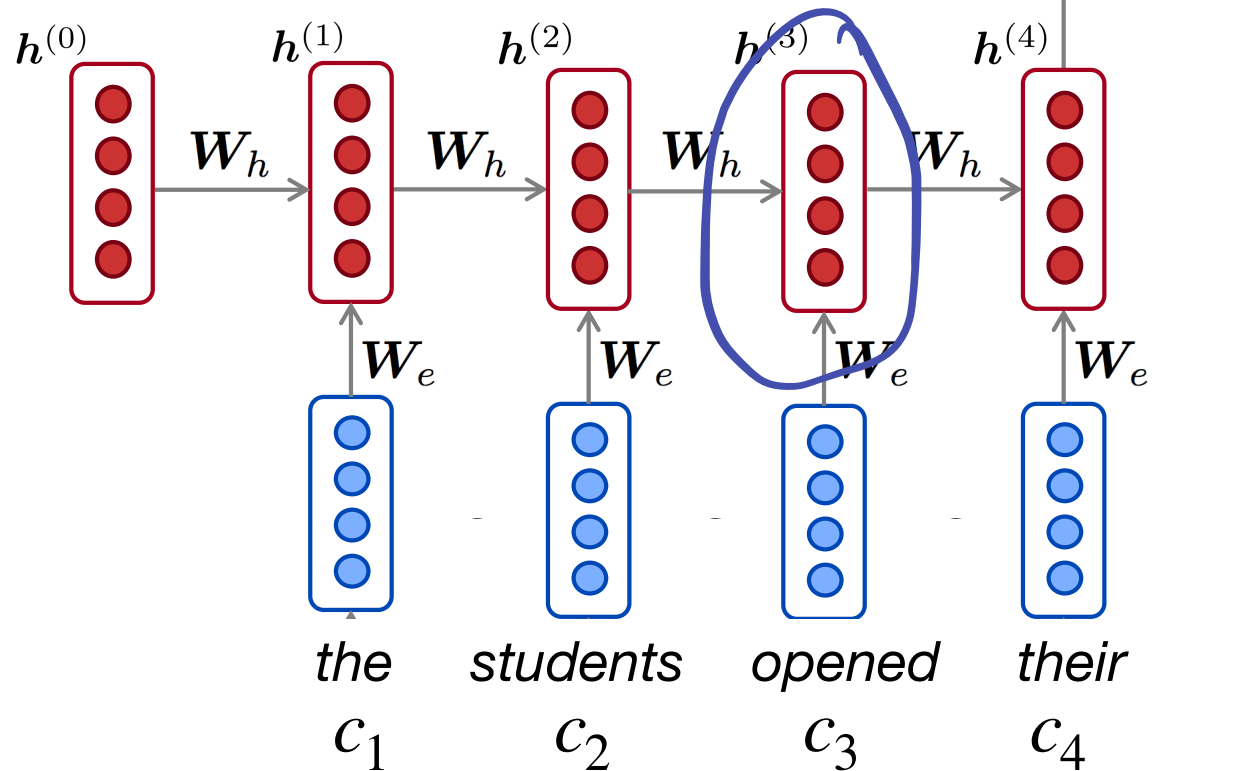
### RNN Advantages:

- 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**

Main Issue:  
access to the  
past depends on  $W_h$



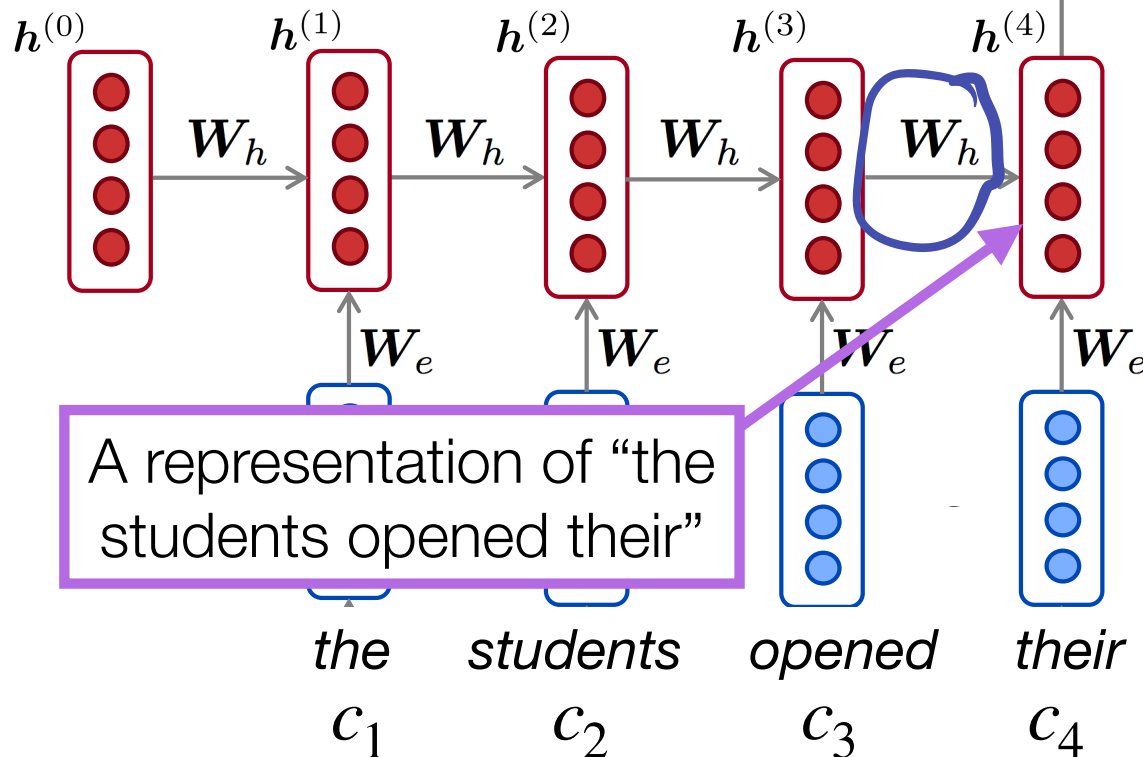
$$\hat{y}^{(4)} = P(x^{(5)} | \text{the students opened their})$$

# RNNs suffer from a **bottleneck** problem

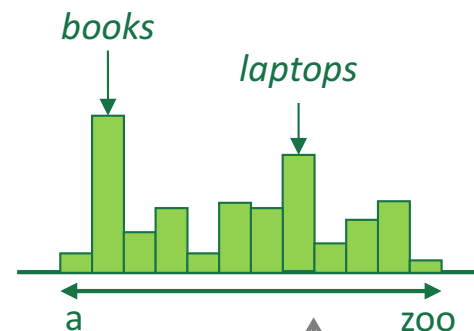
The current hidden representation must encode all of the information about the text observed so far

This becomes difficult especially with longer sequences

$W_h$



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

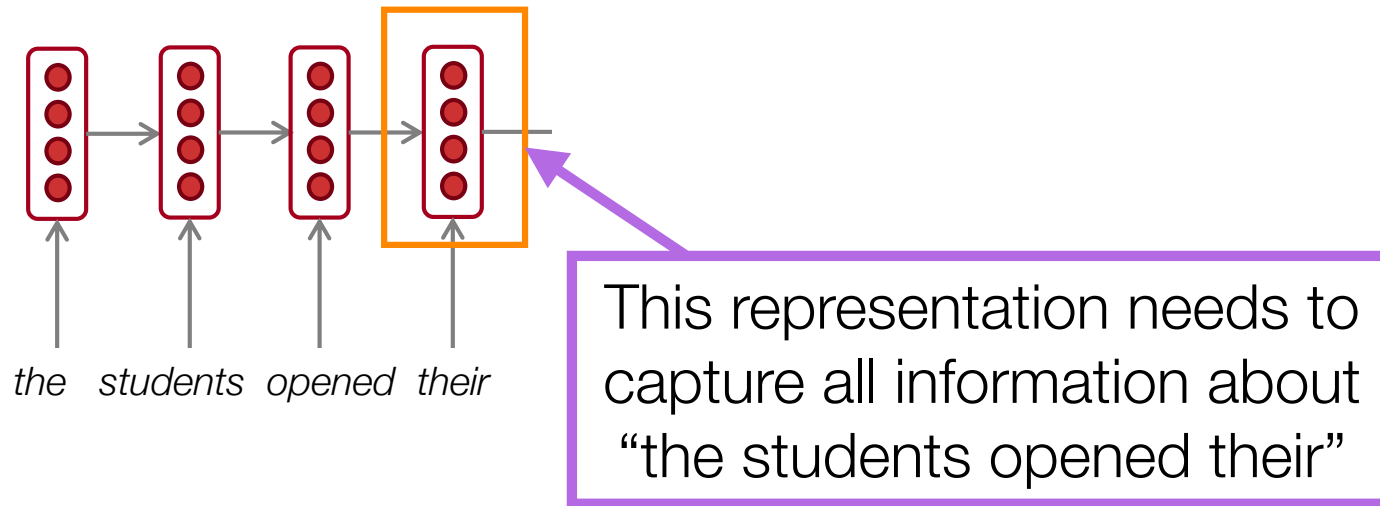




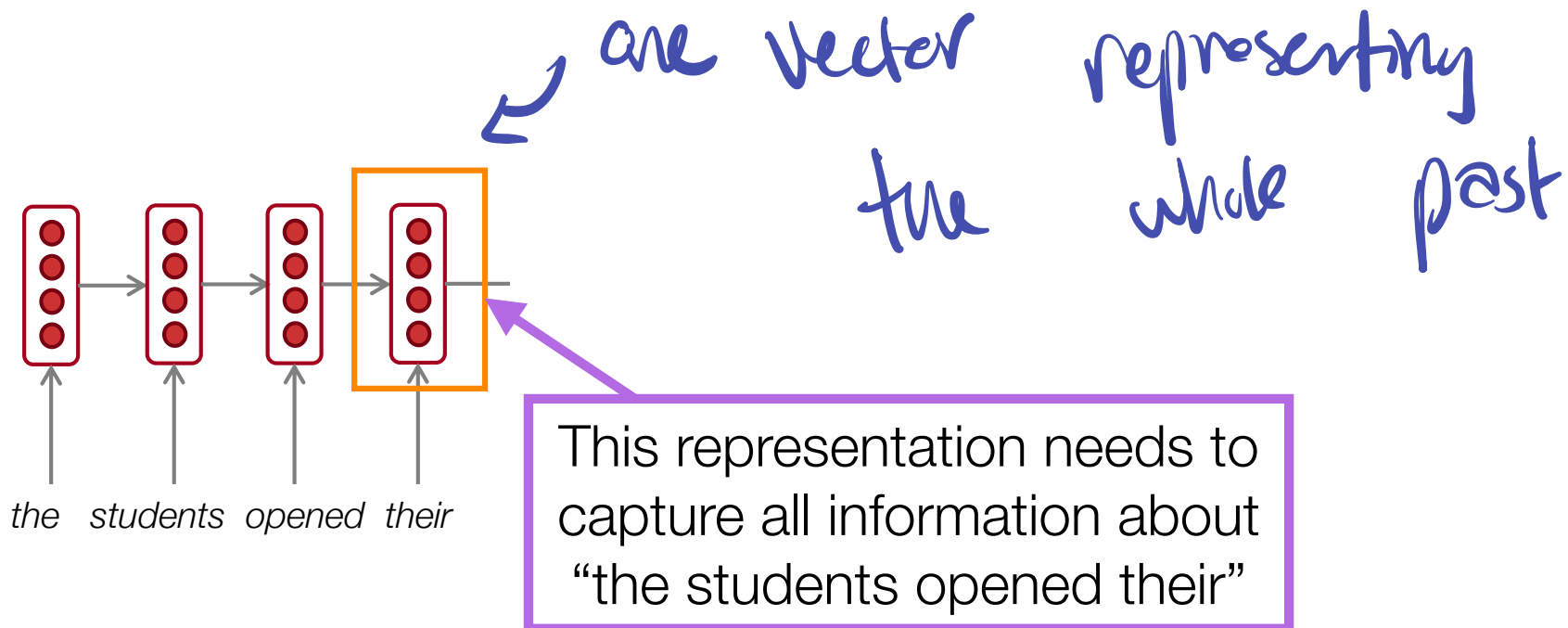
“you can’t cram the meaning  
of a whole %&@#&ing  
sentence into a single  
\$\*(&@ing vector!”

— Ray Mooney (NLP professor at UT Austin)

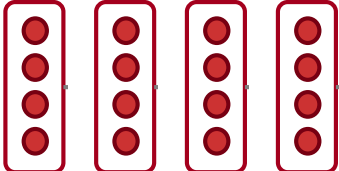
idea: what if we use multiple vectors?



idea: what if we use multiple vectors?



Instead of this, let's try:

the students opened their =  (all 4 hidden states!)

multiple vectors representing the past

# 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

# French

Le chat est joyeuse

English	The			
	cat			
	is			
	happy			

**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 Is Attention?

Goal: learn a task-specific vector  $v$

Intuition: think of  $v$  as an "important word" vector

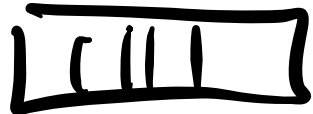
Step 1: Measure the importance of each input vector  $x$  by computing its similarity to  $v$ . Dot Product

$$r_1 = v \cdot x_1$$



$x_1$ : The

$$r_2 = v \cdot x_2$$



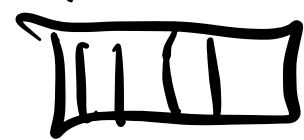
$x_2$ : students

$$r_3 = v \cdot x_3$$



$x_3$ : opened

$$r_4 = v \cdot x_4$$



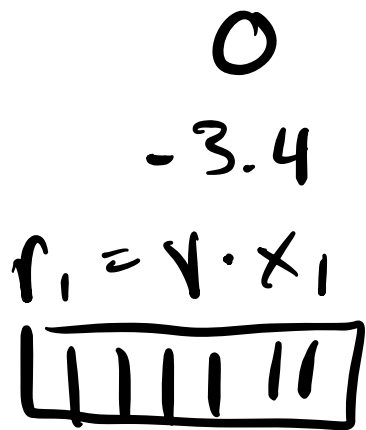
$x_4$ : their



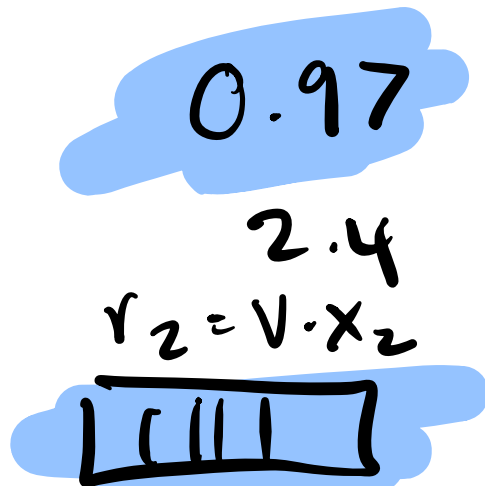
# What Is Attention?

Step 2: Take  $r$  & normalize it  
using softmax

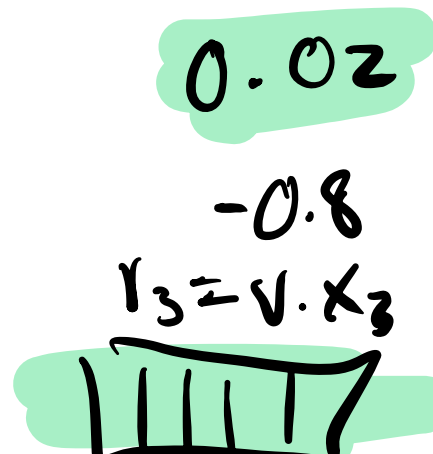
$$a = \text{softmax}(r)$$



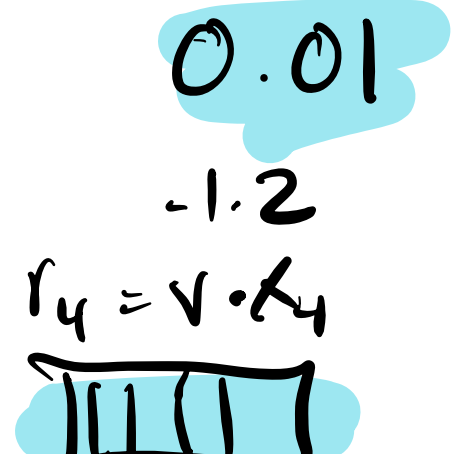
$x_1$ : The



$x_2$ : students



$x_3$ : opened

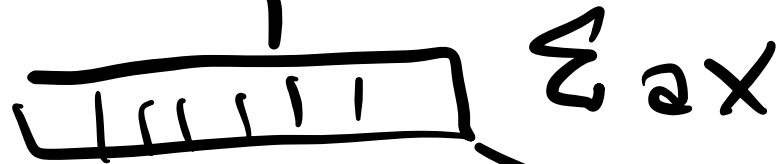


$x_4$ : their

# What Is Attention?

Step 3: Compute a weighted average of  $X$ .

$$\hat{y} = \sum a_i x_i$$



$$a_1 = 0$$

$$-3.4$$

$$r_1 = v \cdot x_1$$

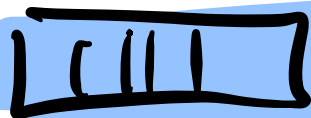


$x_1$ : The

$$a_2 = 0.97$$

$$2.4$$

$$r_2 = v \cdot x_2$$



$x_2$ : students

$$a_3 = 0.02$$

$$-0.8$$

$$r_3 = v \cdot x_3$$



$x_3$ : opened

$$a_4 = 0.01$$

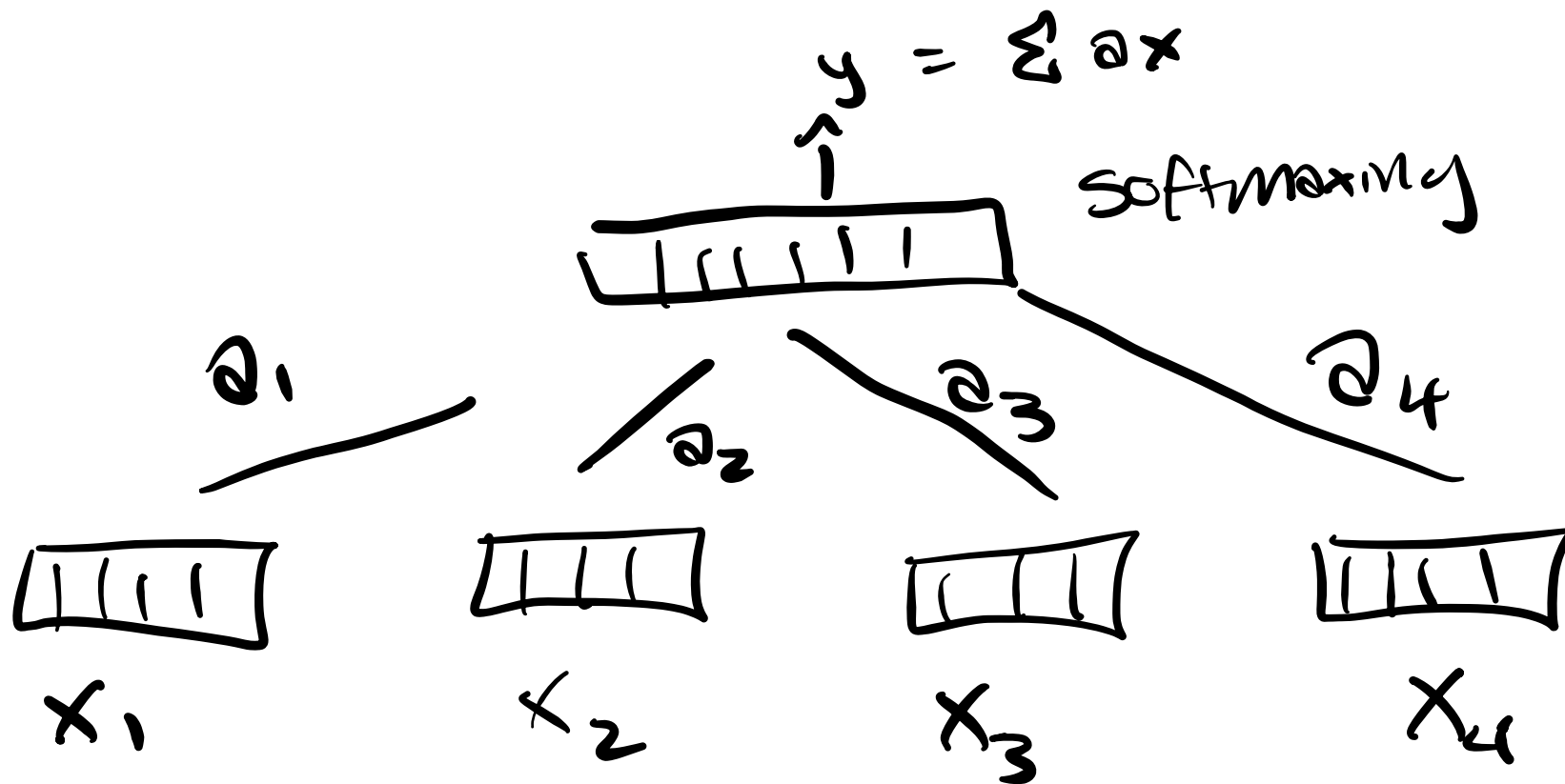
$$-1.2$$

$$r_4 = v \cdot x_4$$



$x_4$ : their

# What Is Attention?



$$y = \text{softmax}(ax)$$



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

6:20 PM · Feb 22, 2023 · 88.1K Views

# 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.sonoma.edu/users/f/forham/sonoma/projects/ca/labview/index.htm>

In-class Lab 1: Introduction to LabVIEW

A- Read <http://sonoma.edu/content/19837/latest/>

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.at.com.au/presentation/us/labview/snap/default.htm>

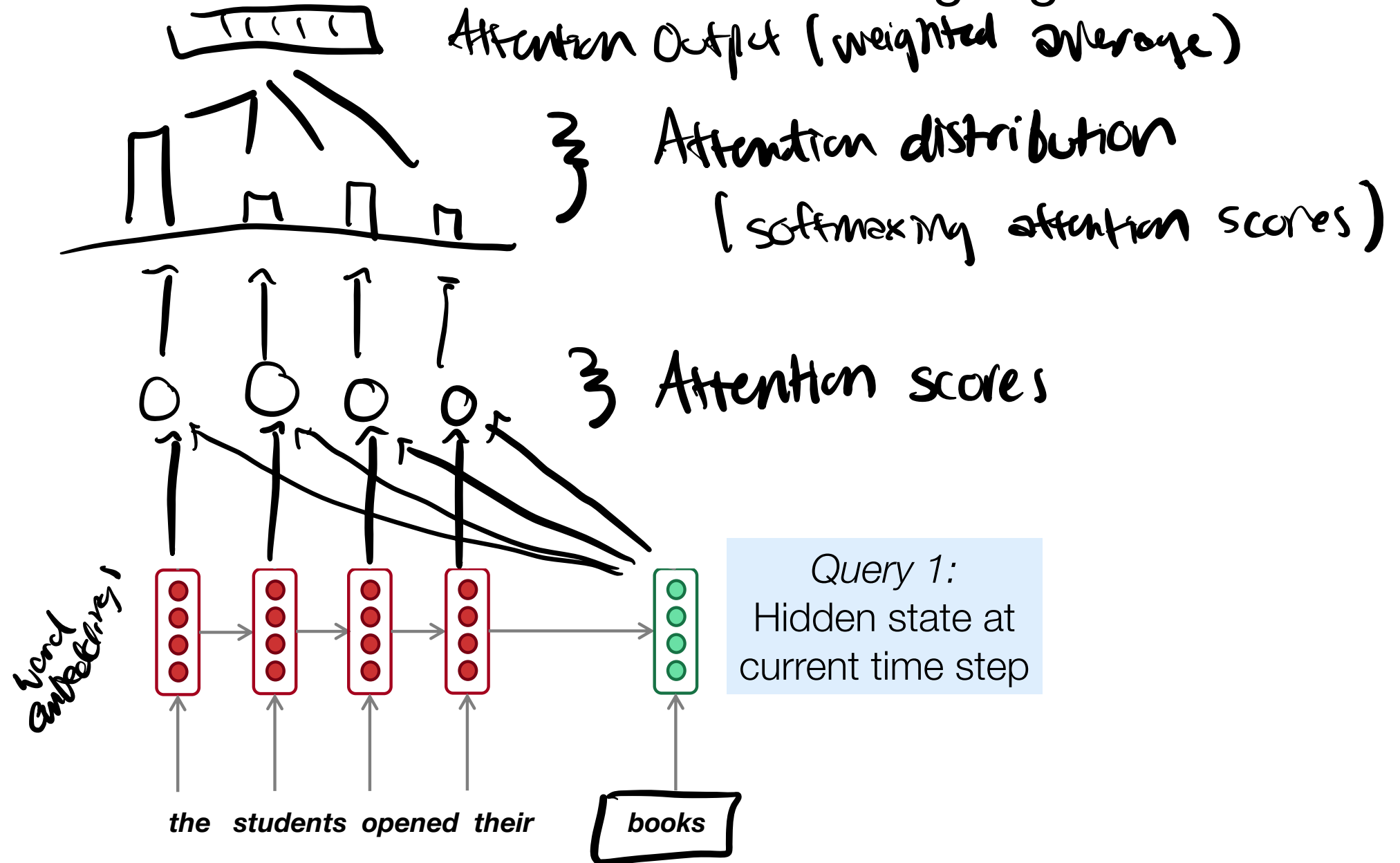
Assignment 1:  
Create a simple program that can convert a temperature from the Celsius scale to the Fahrenheit scale: <http://www.utevan.edu/~scdm/firbybytes/lab1.htm>. Take a snap shot of the Front panel and Diagram. Place the figures in the table below.

Figure 1. Front Panel VI for Temperature Converter.
Snap shot here!
Figure 1. Block Diagram for Temperature Converter.
Snap shot here!

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.

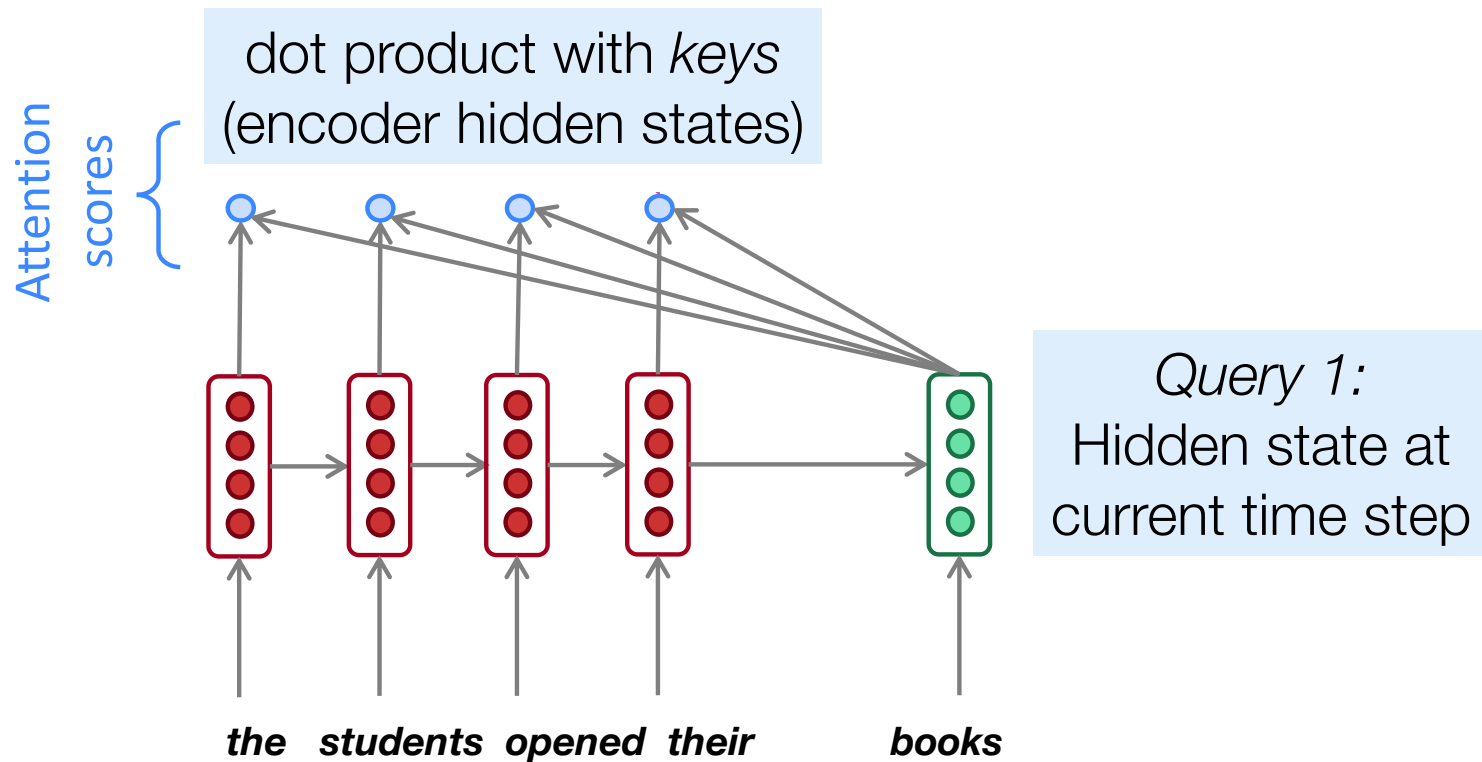
Figure 1. Front Panel VI for Random Number Generator.
Snap shot here!
Snap shot here!

# Attention mechanisms in neural language models

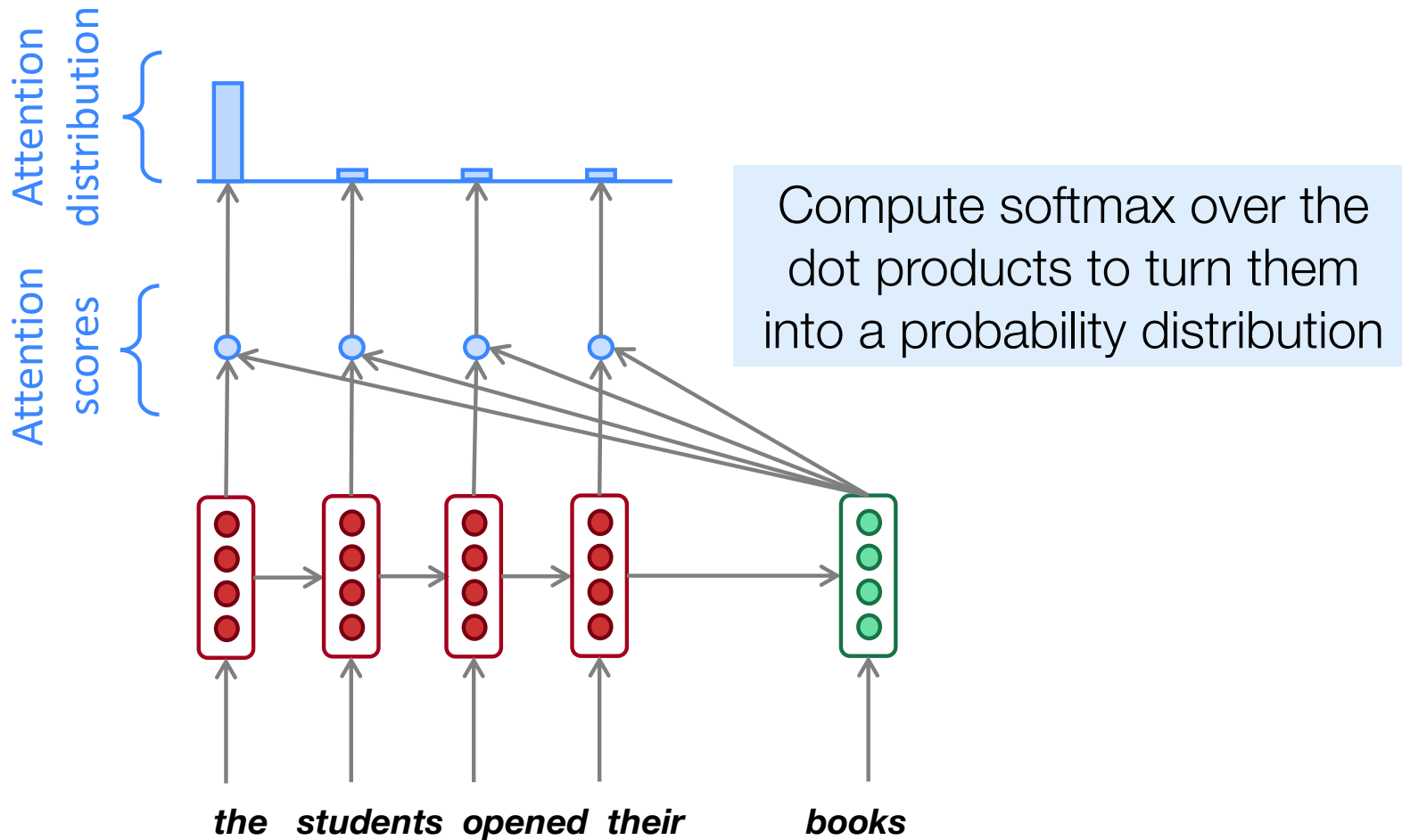




# Attention mechanisms in neural language models



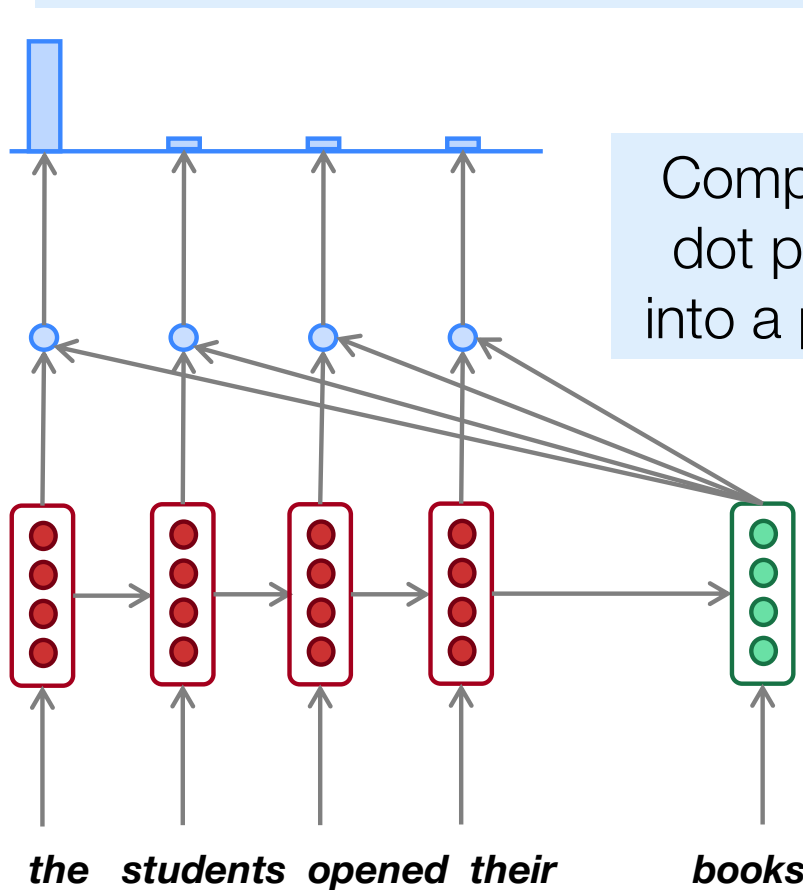
# Attention mechanisms in neural language models



# Attention mechanisms in neural language models

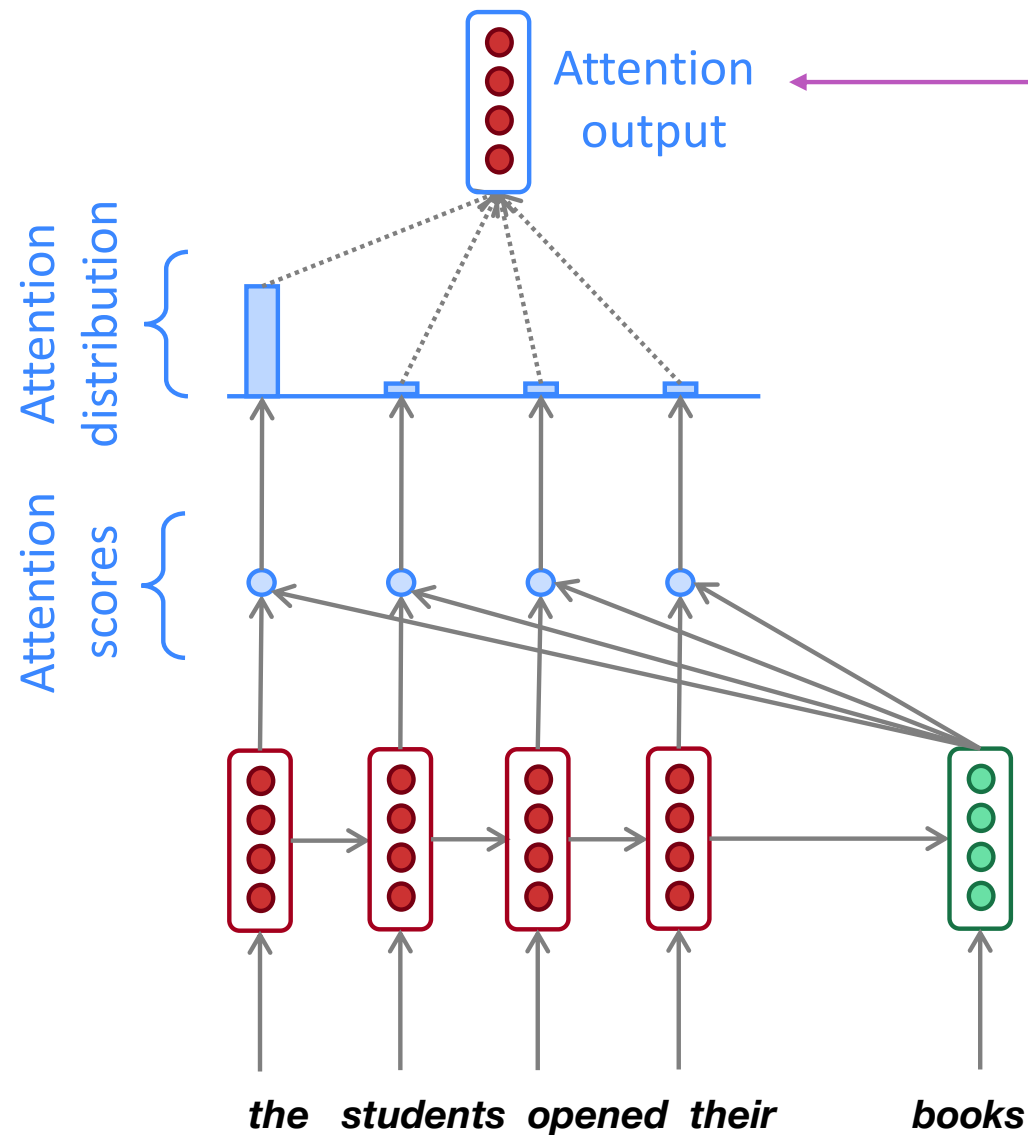
At this time step, the attention distribution is focused on the first word of the sequence (“the”)

Attention scores  
Attention distribution



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

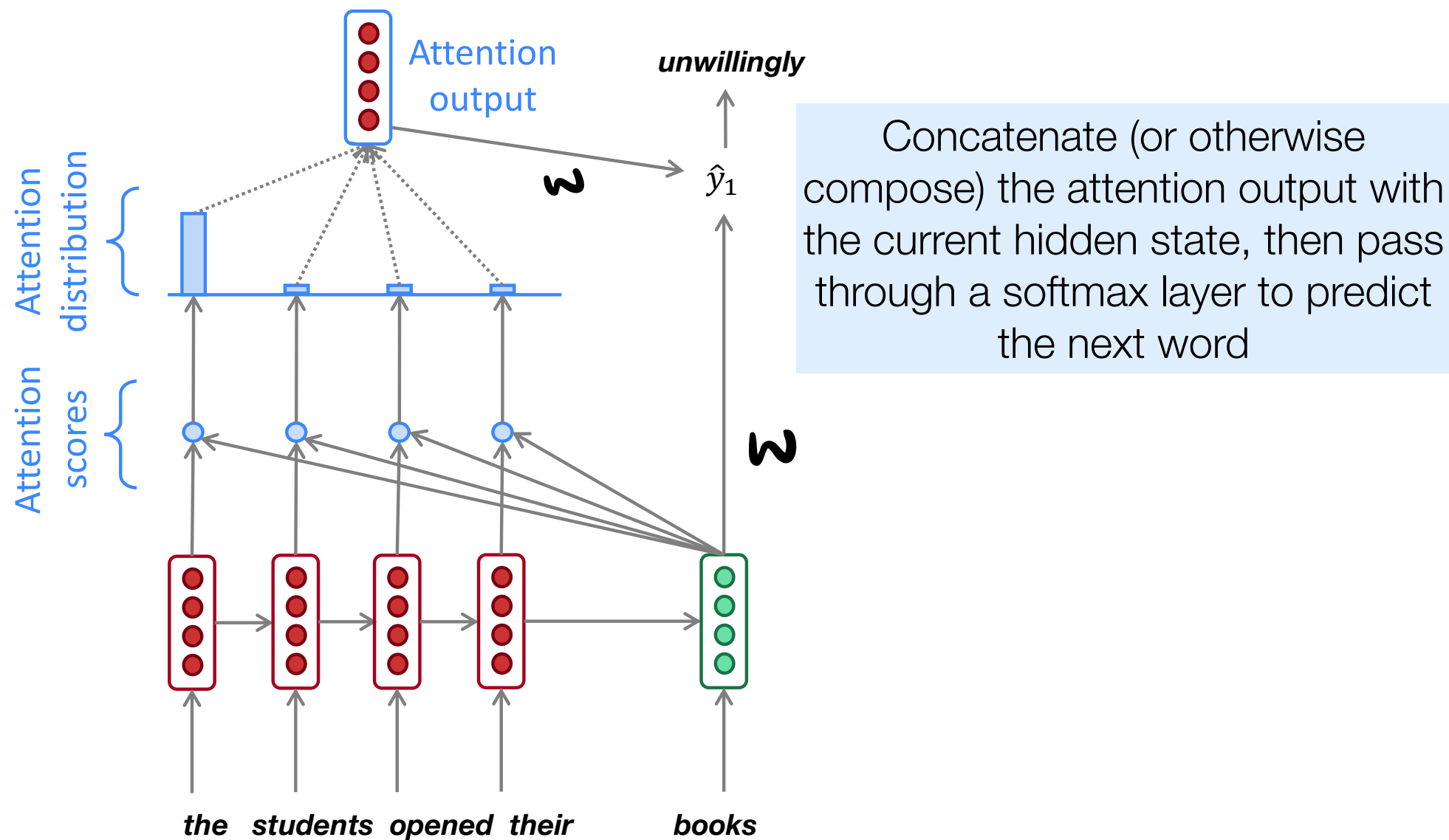
# Attention mechanisms in neural language models



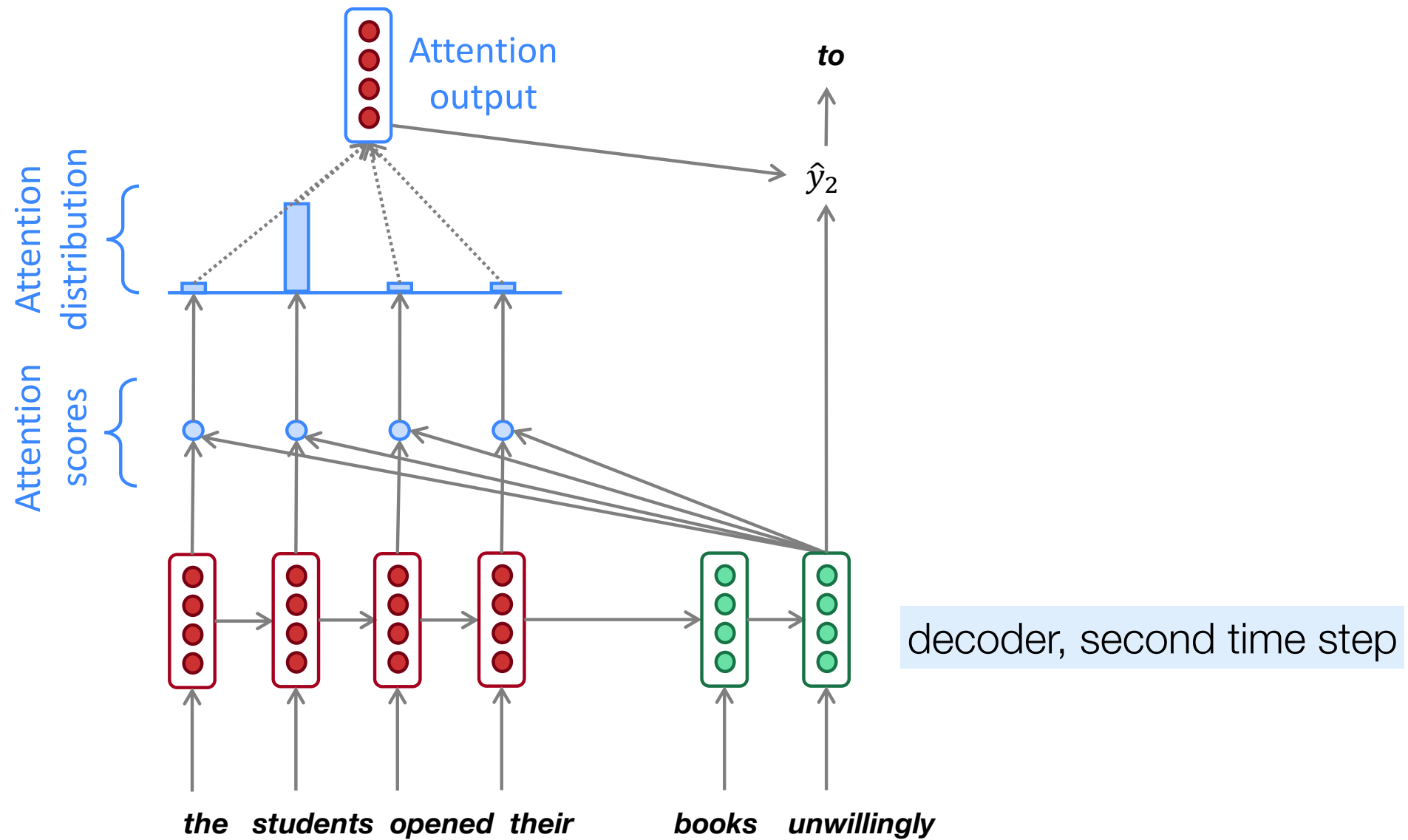
We use the attention distribution to compute a weighted average of the hidden states.

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

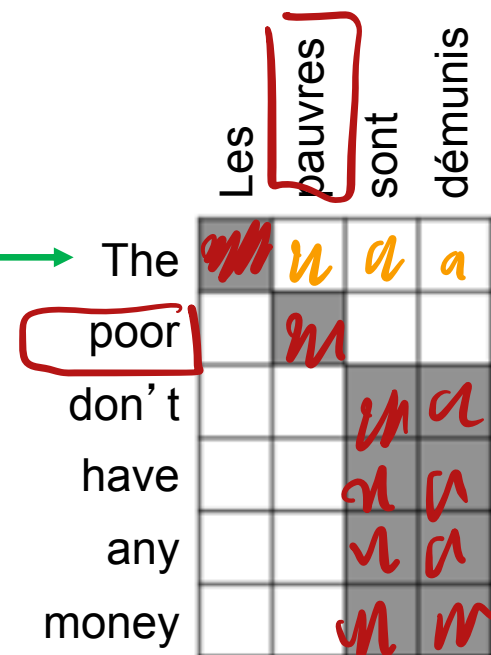
# Sequence-to-sequence with attention



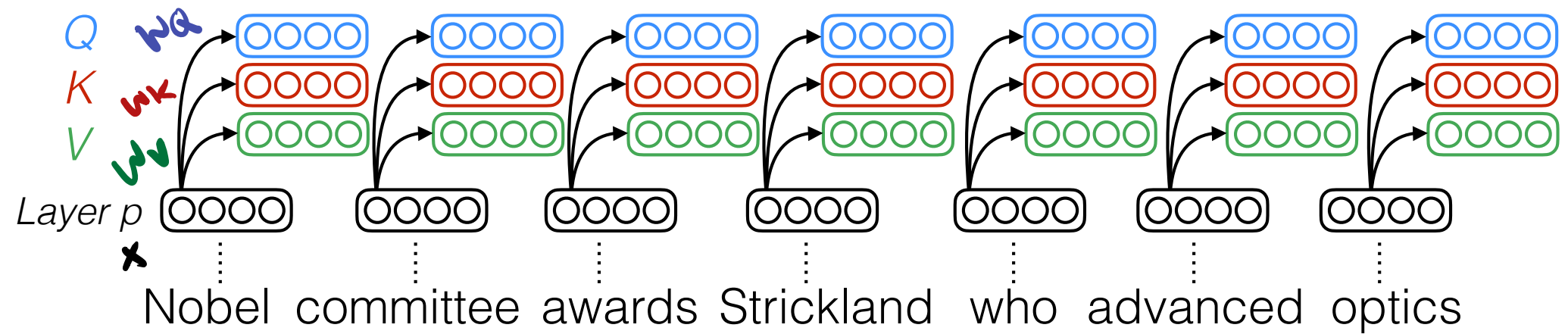
# Sequence-to-sequence with attention



- 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

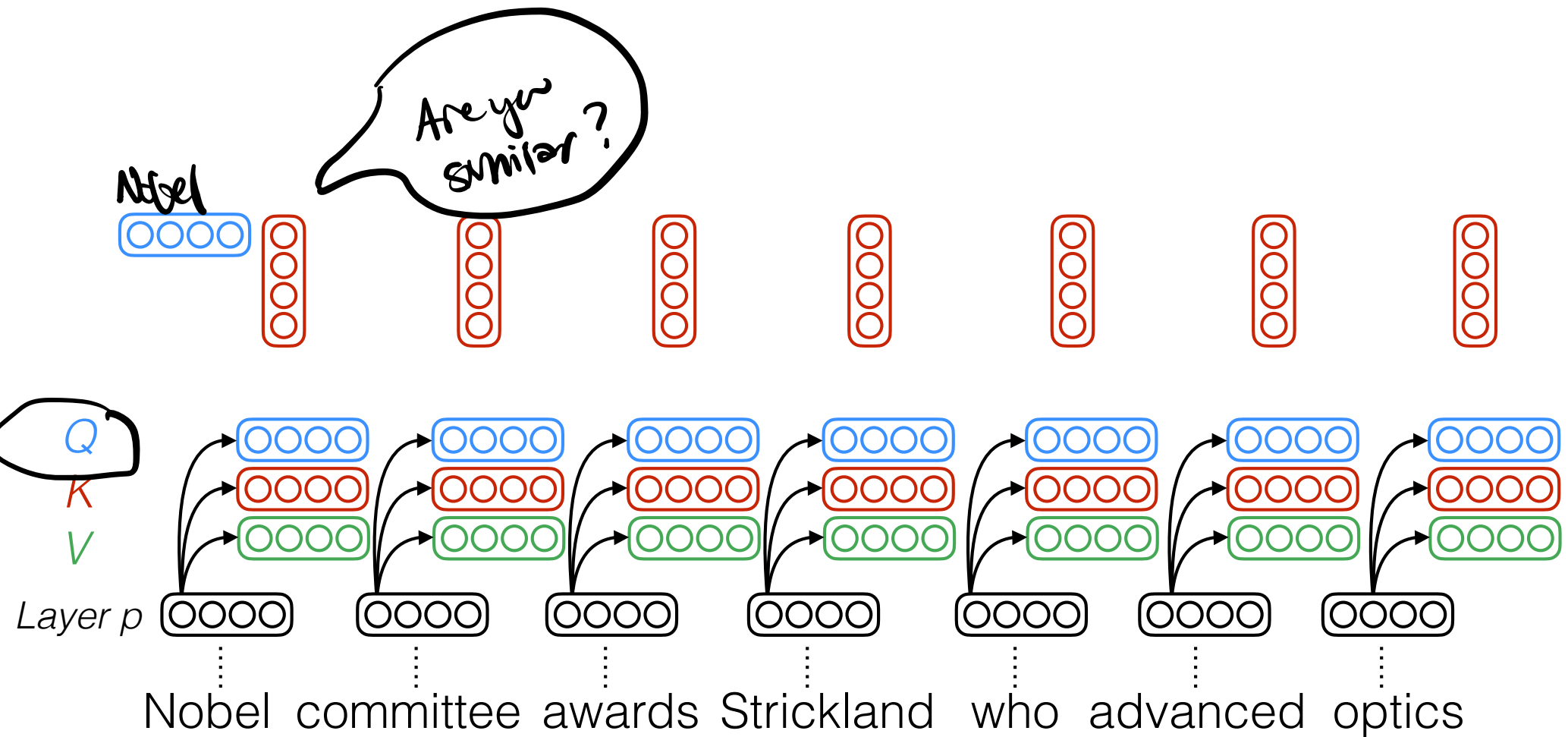


# Self-attention

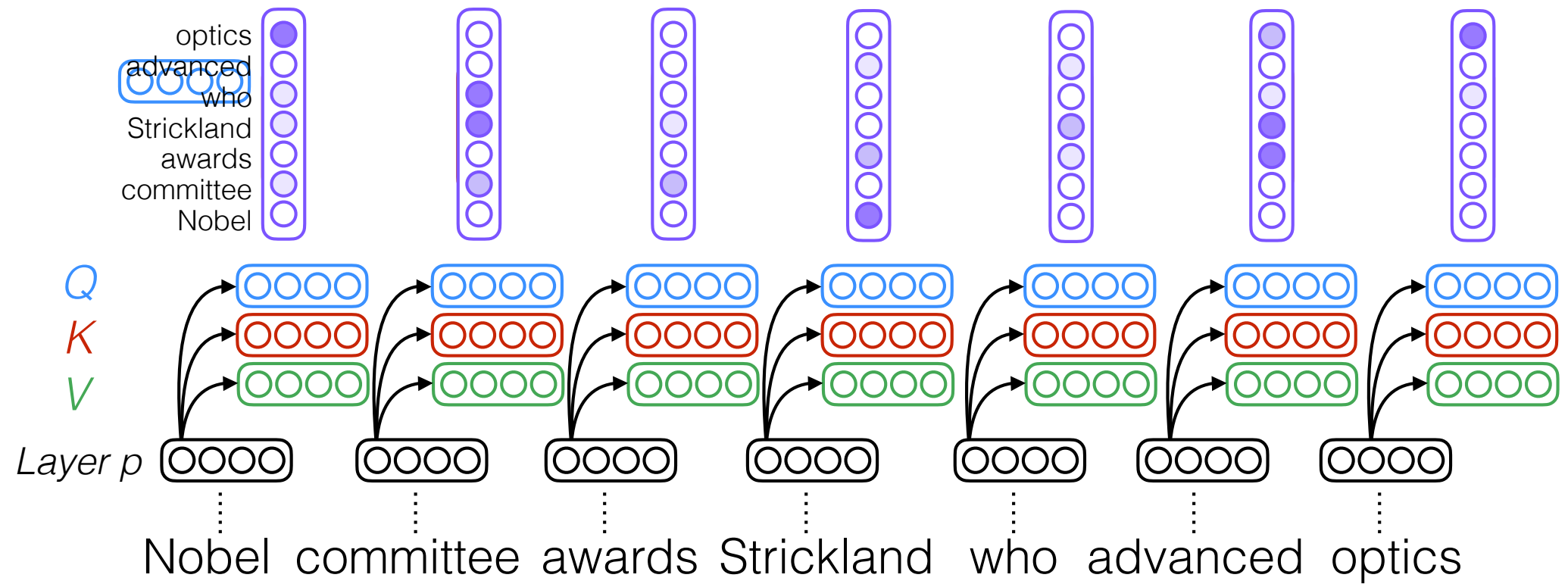




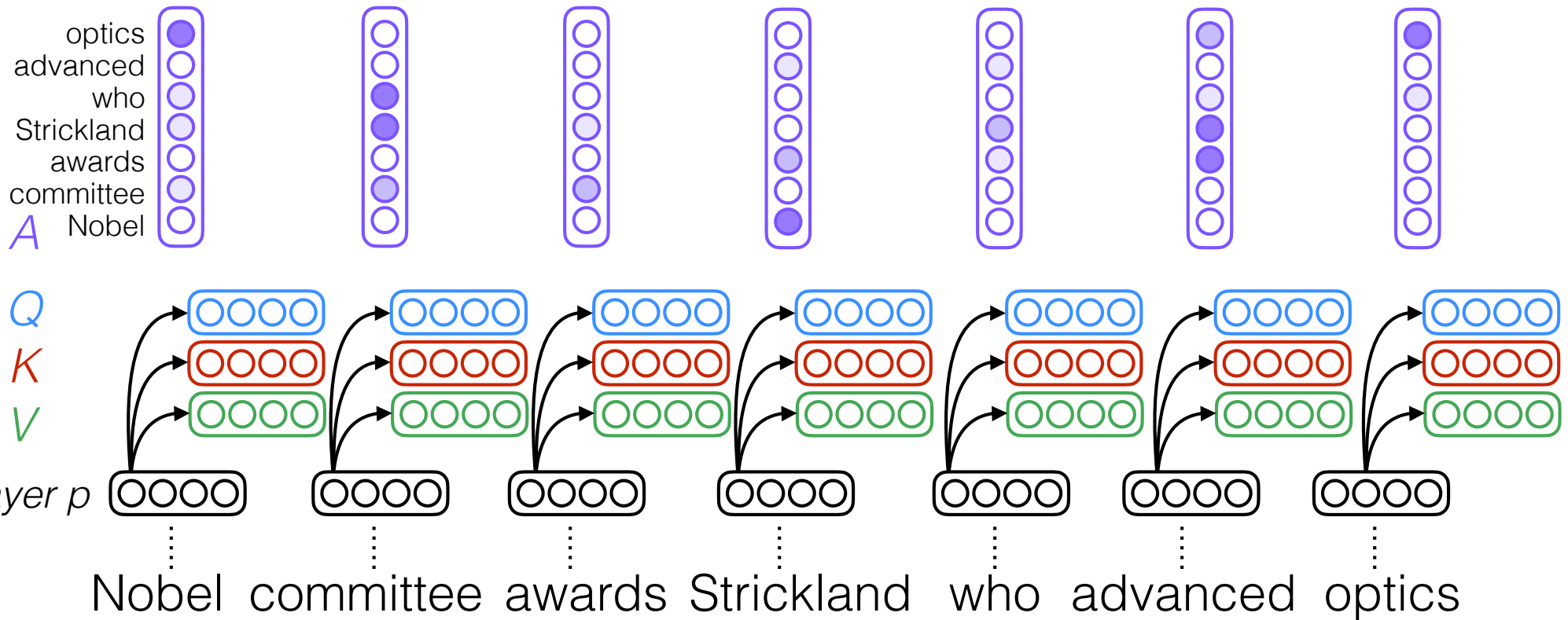
# Self-attention



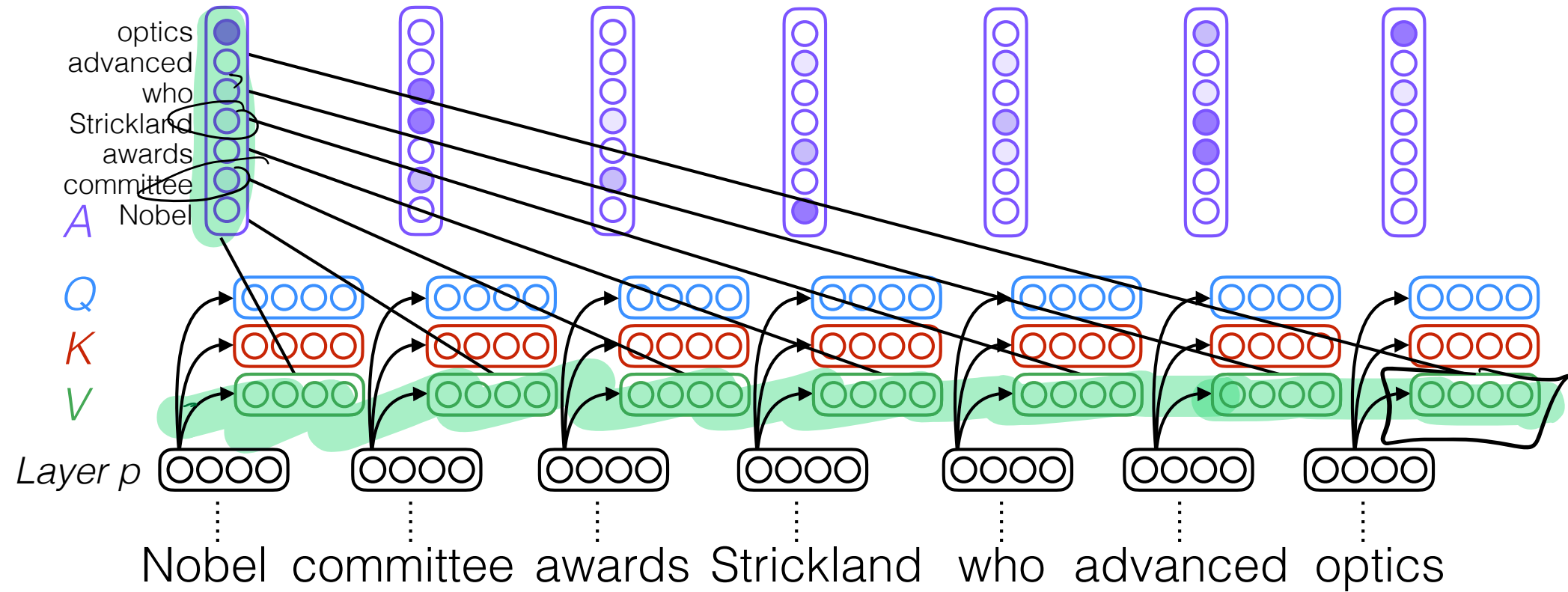
# Self-attention



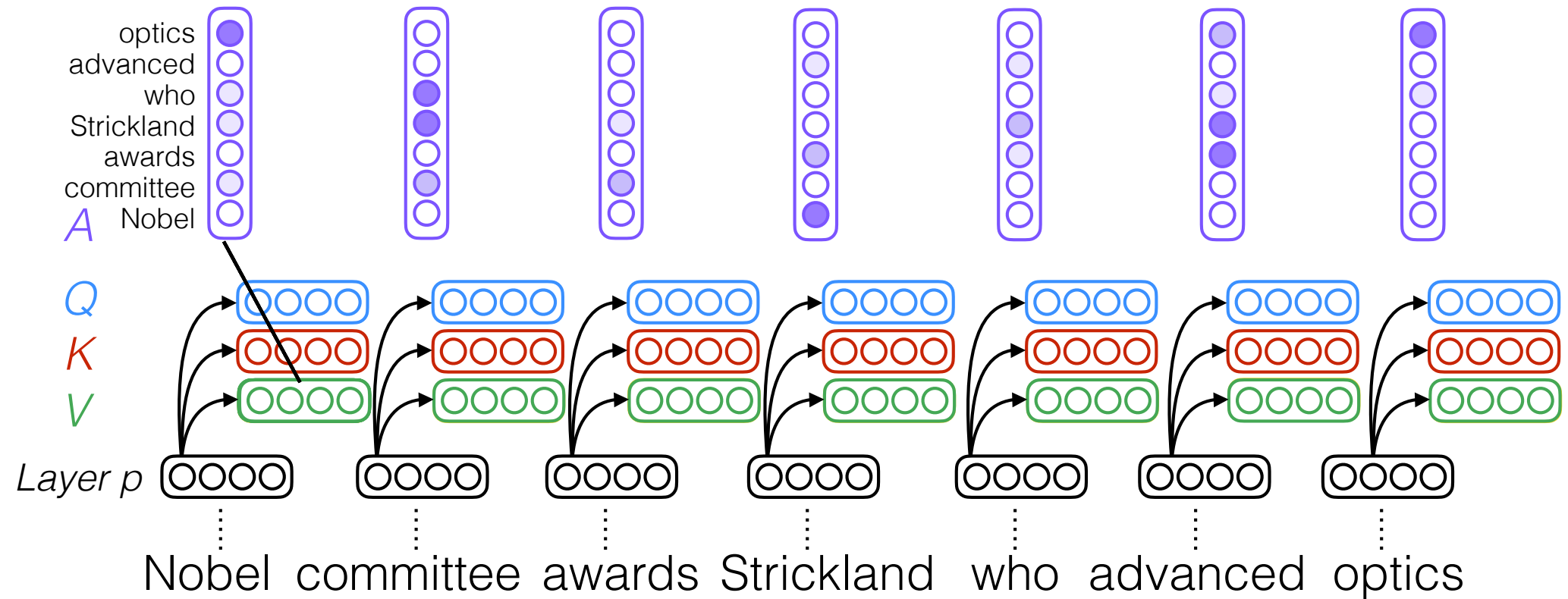
# Self-attention



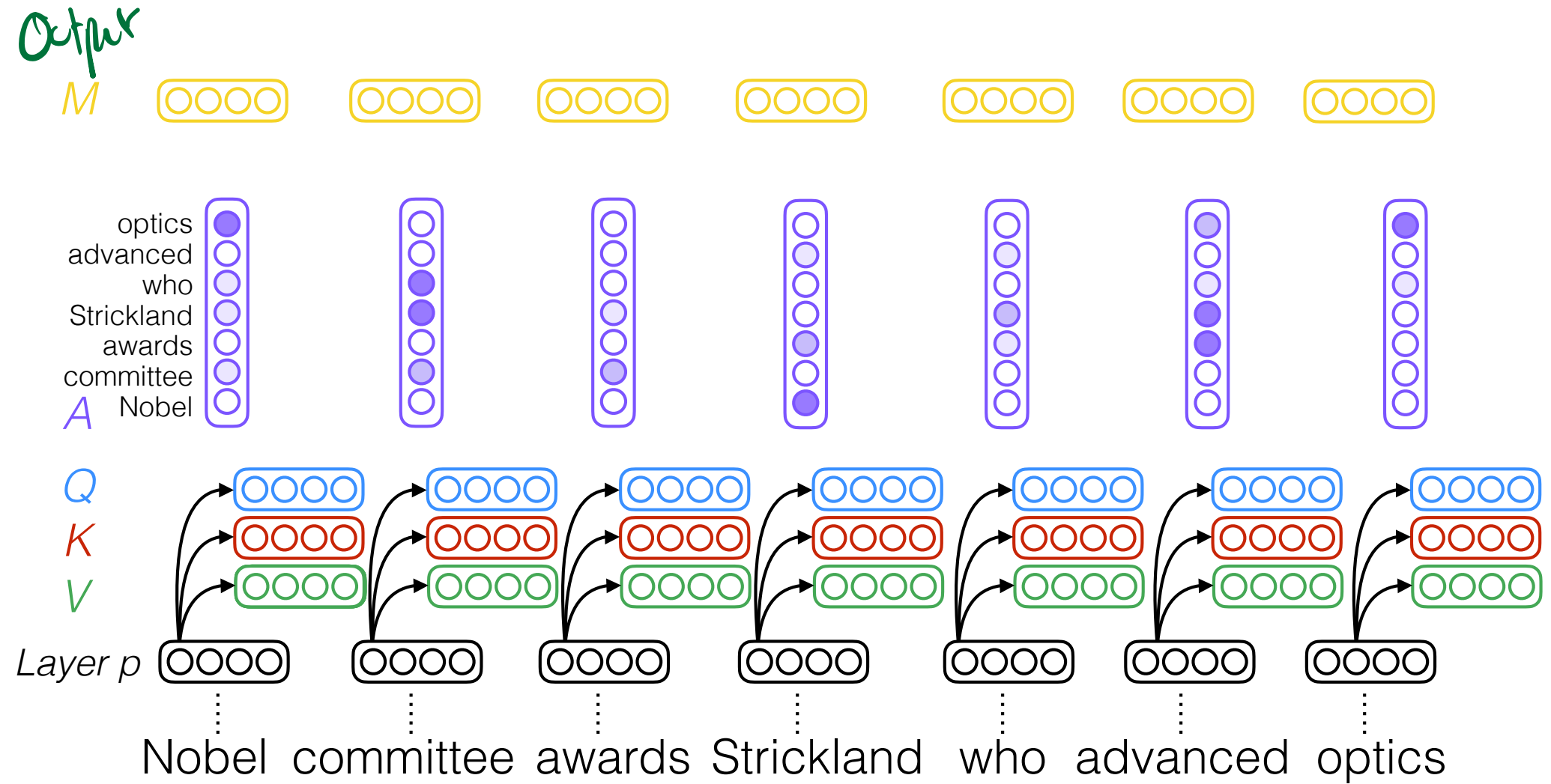
# Self-attention



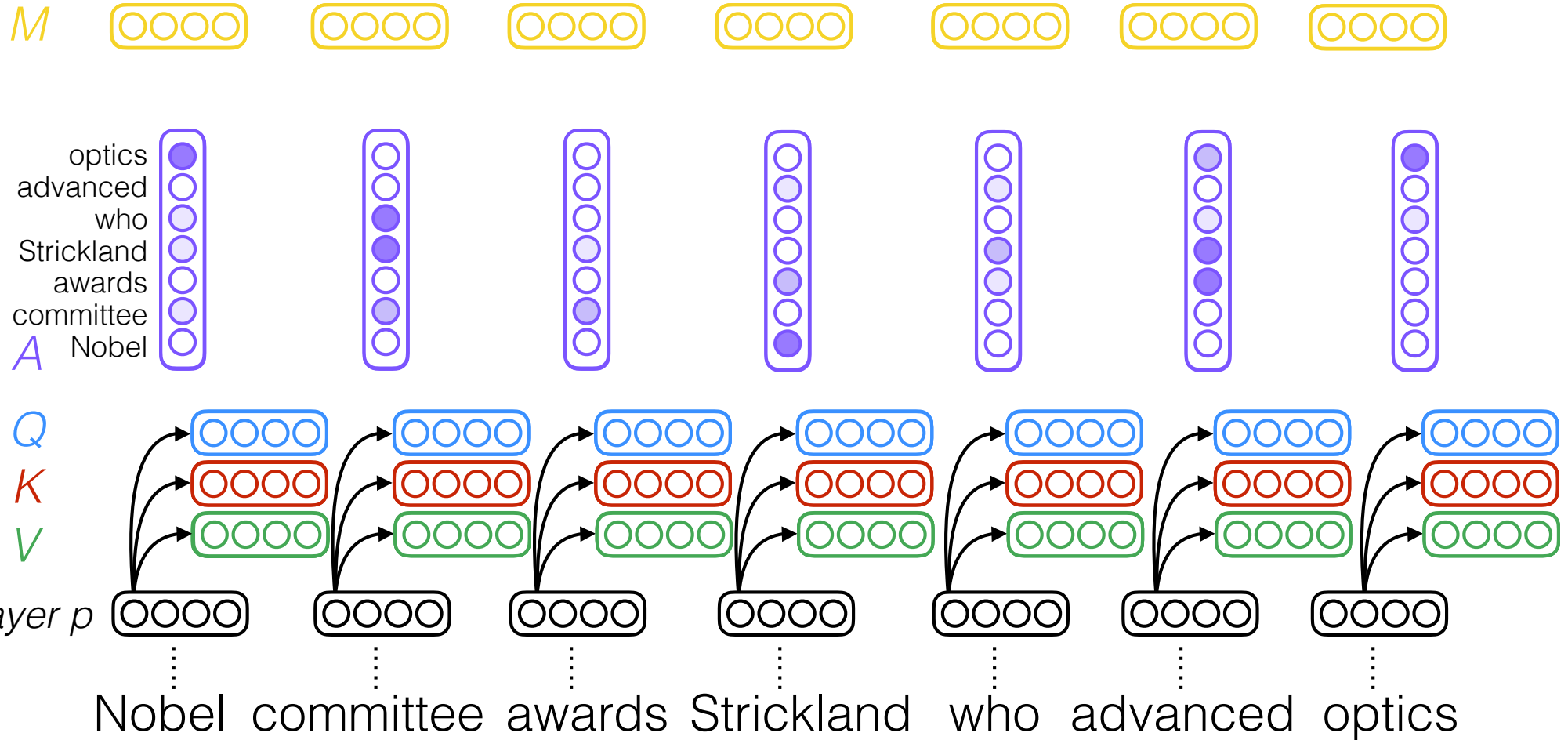
# Self-attention



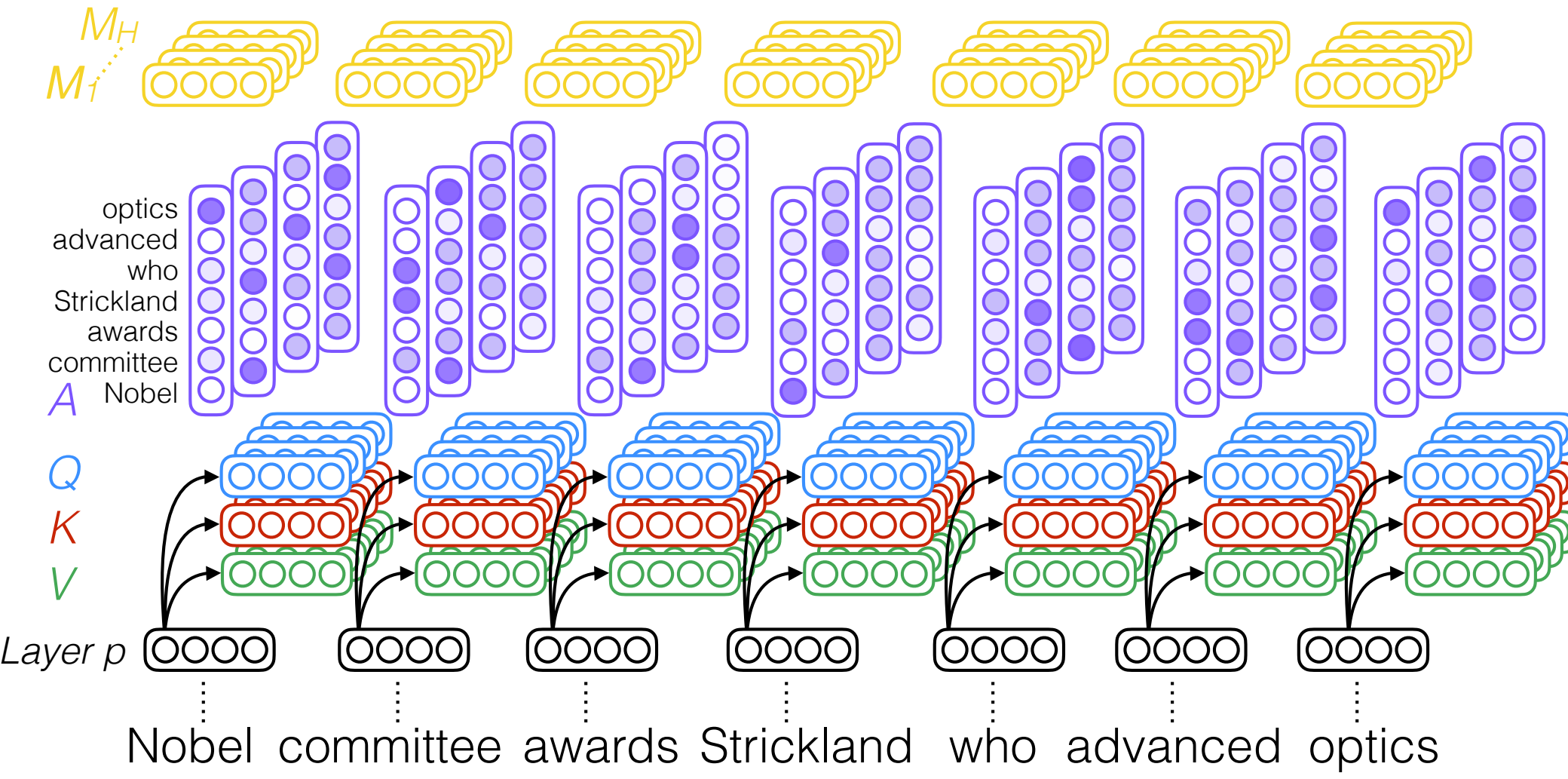
# Self-attention



# Self-attention

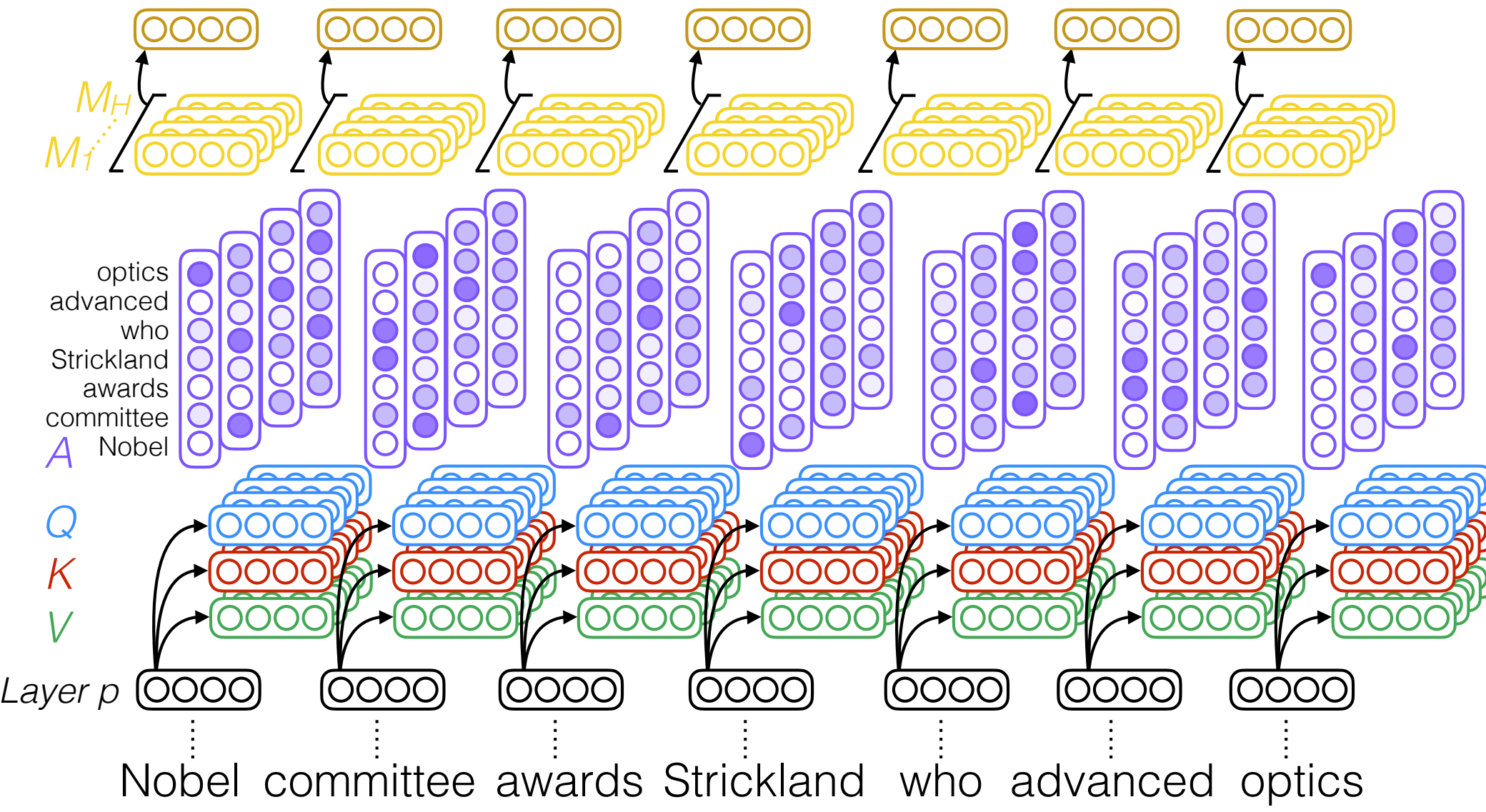


# Multi-head self-attention

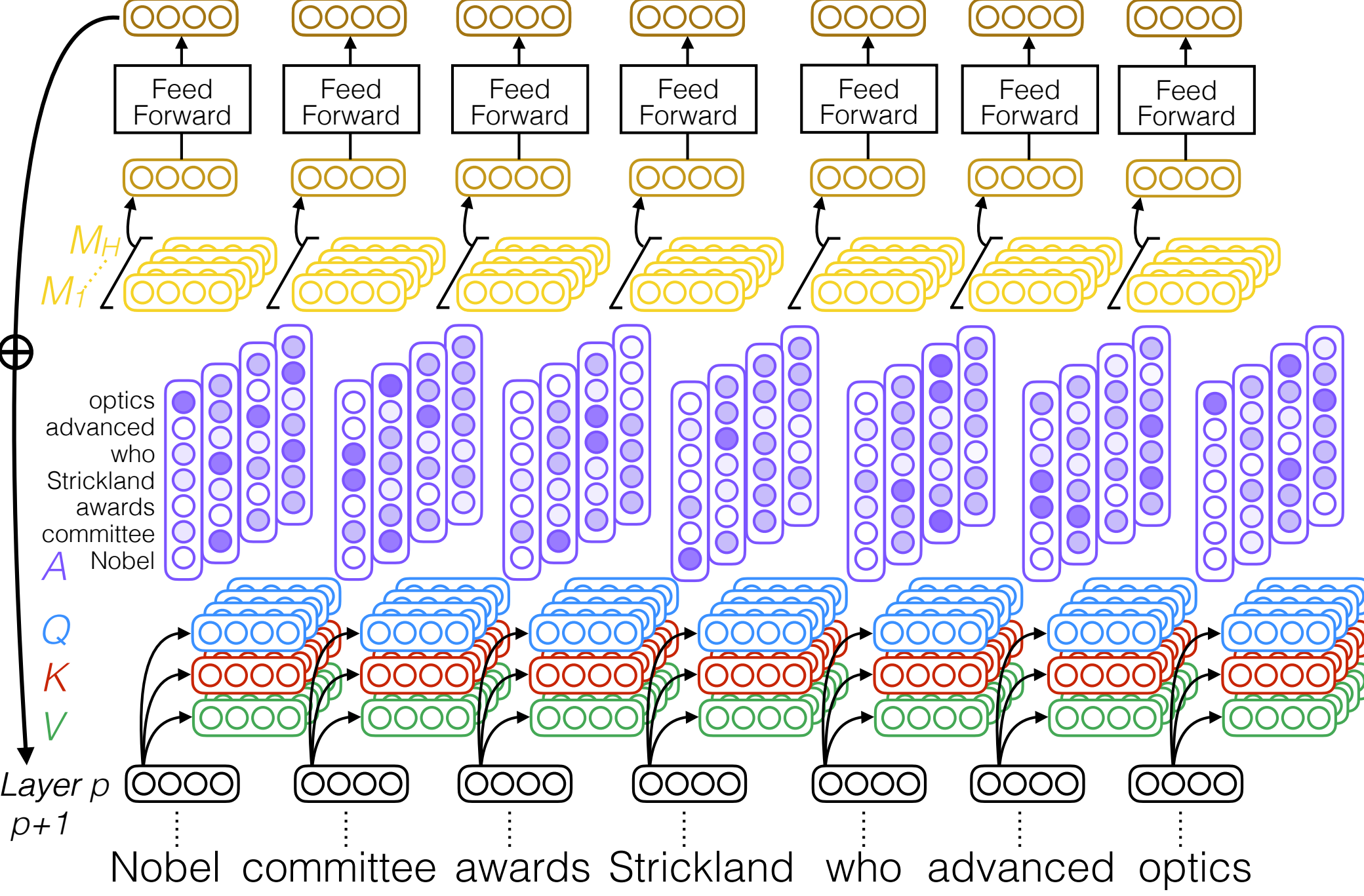




# Multi-head self-attention



# Multi-head self-attention



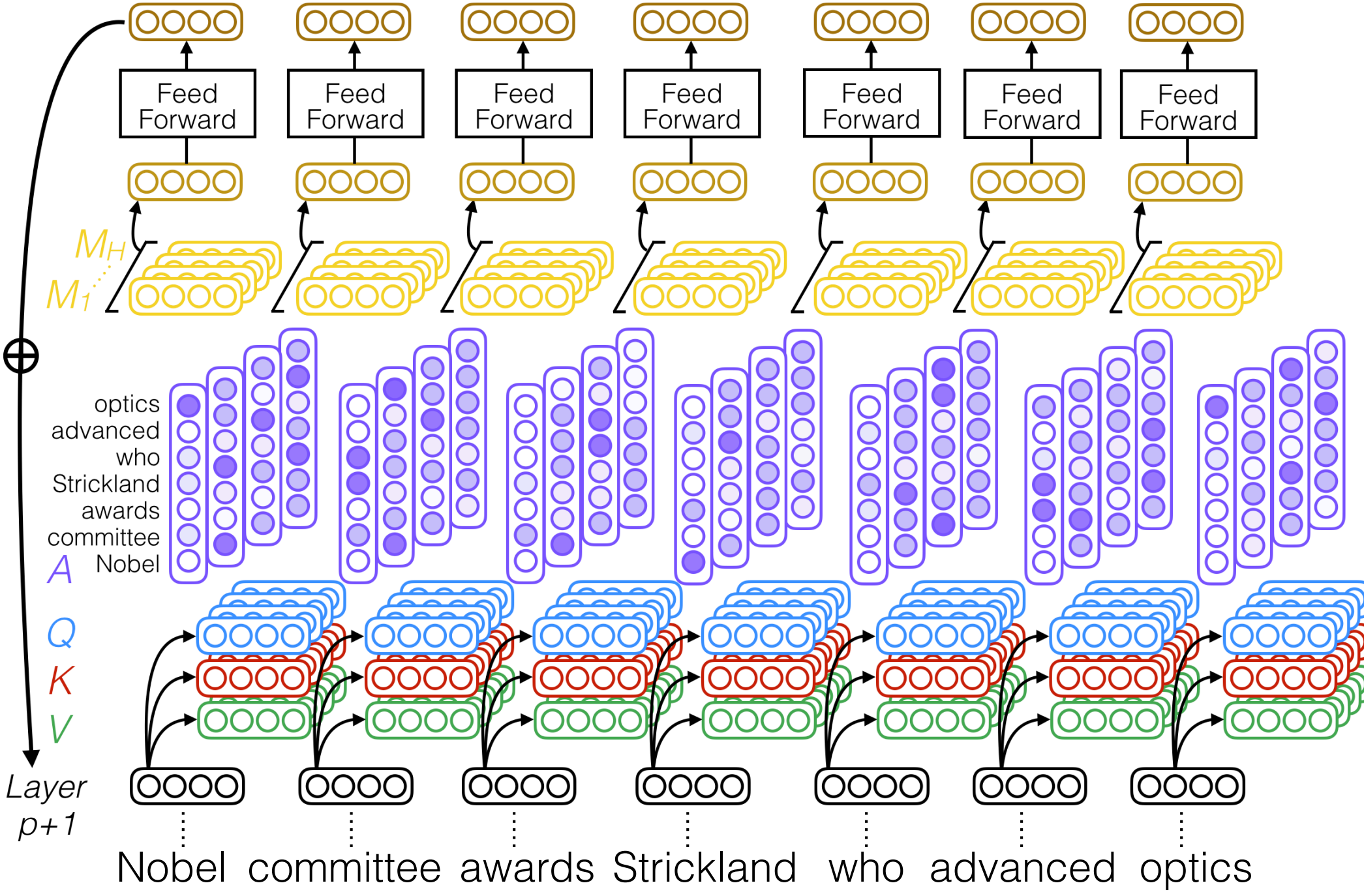
optics  
advanced  
who  
Strickland  
awards  
committee  
Nobel

$Q$   
 $K$   
 $V$

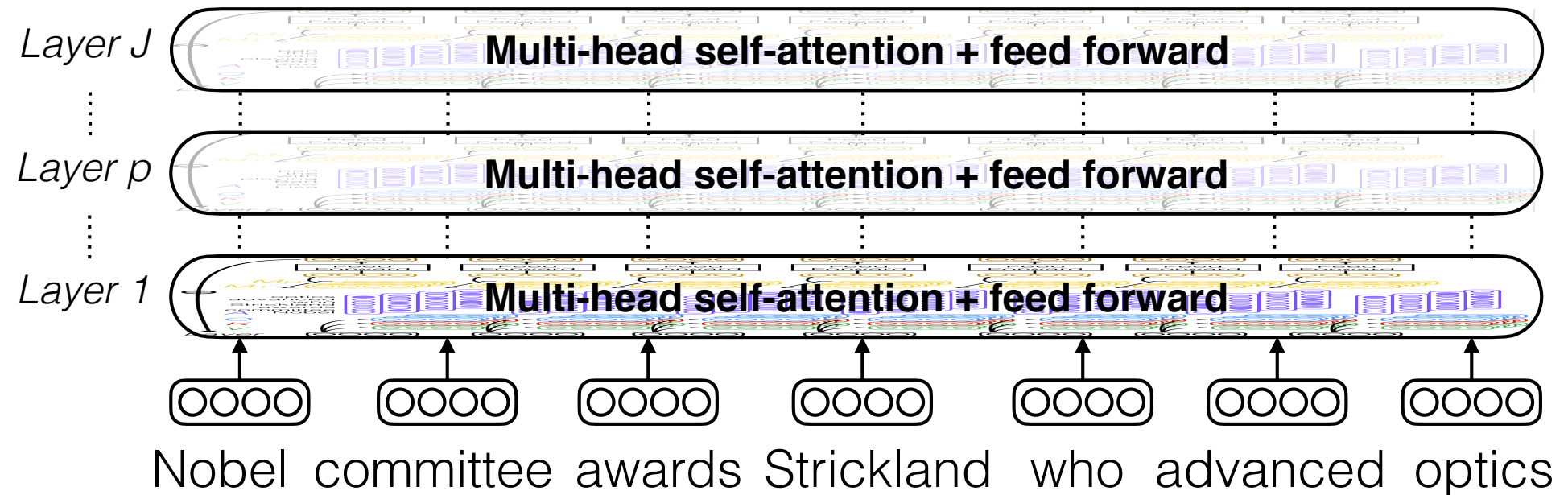
Layer  $p$   
 $p+1$

Nobel committee awards Strickland who advanced optics

# Multi-head self-attention

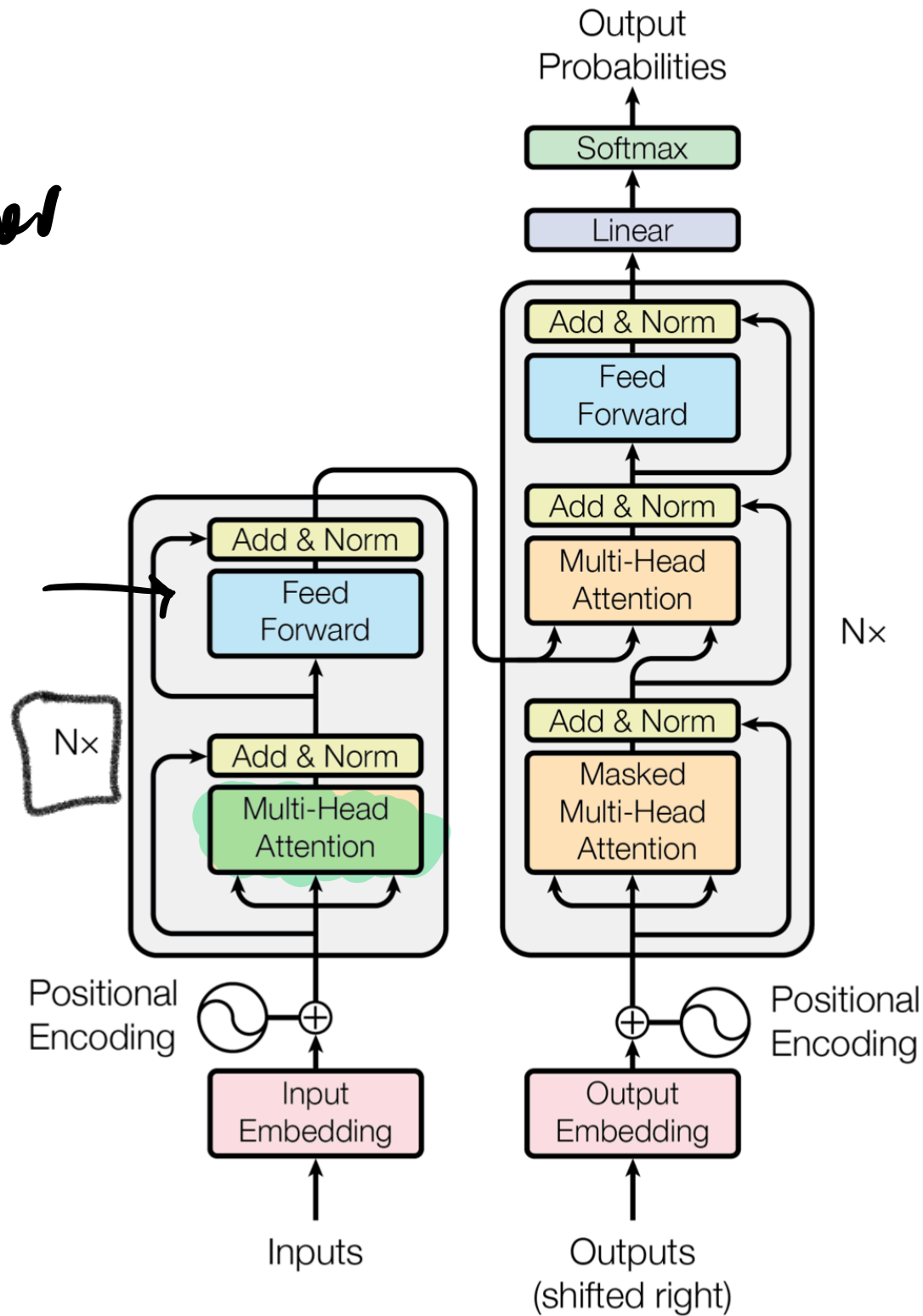


# Multi-head self-attention

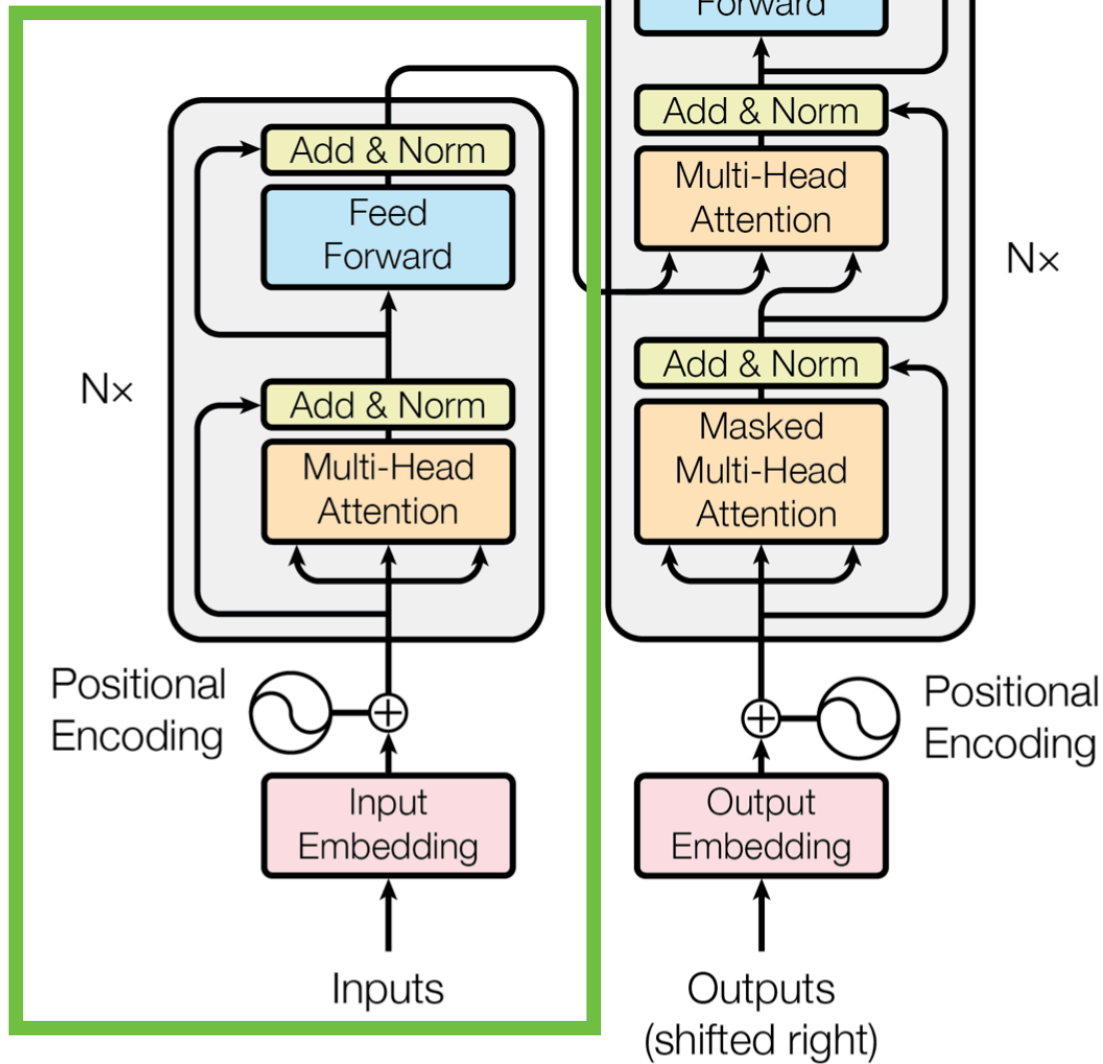


# Transformers

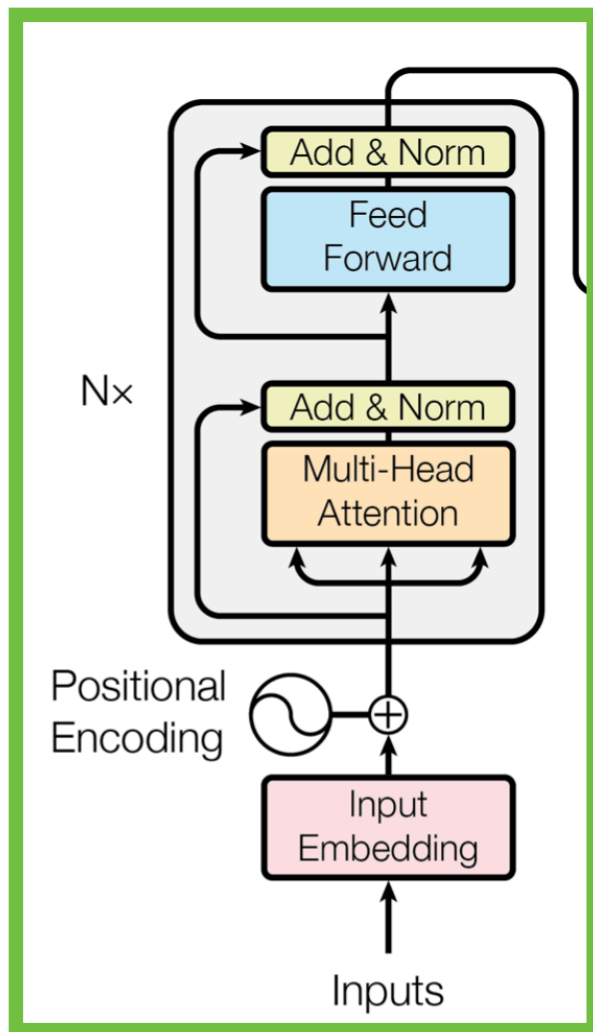
# Transformer



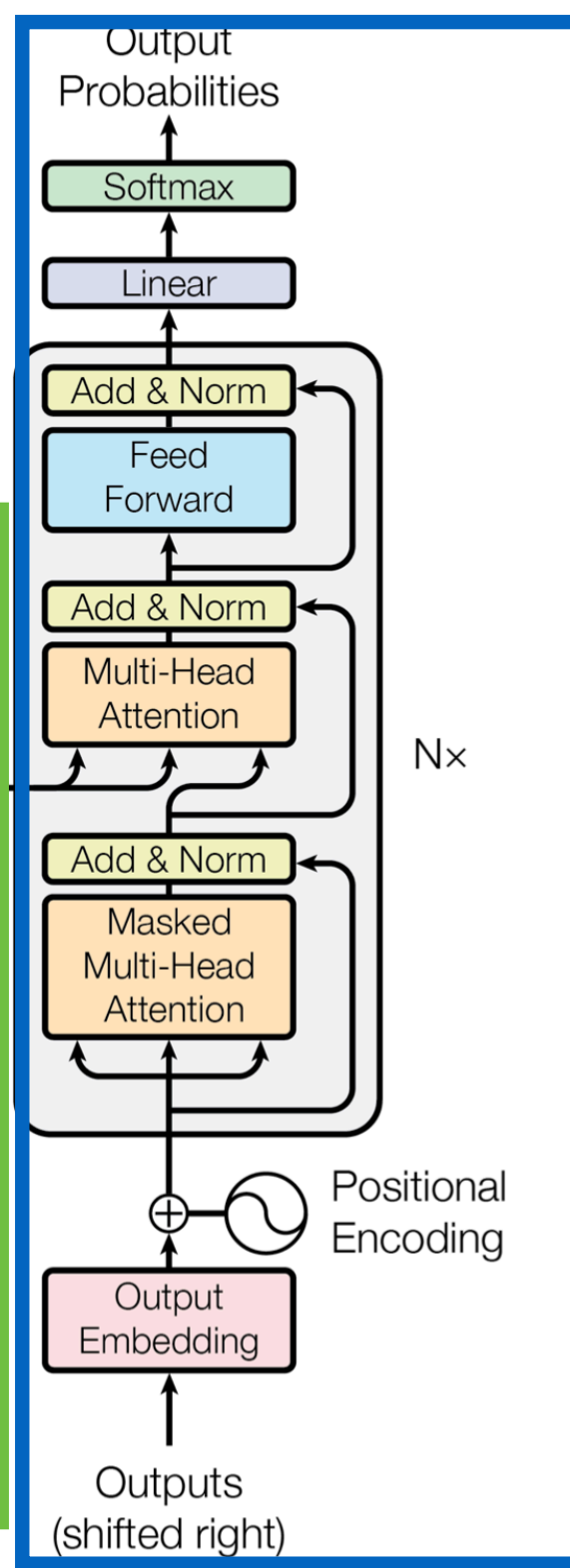
*encoder*



*encoder*



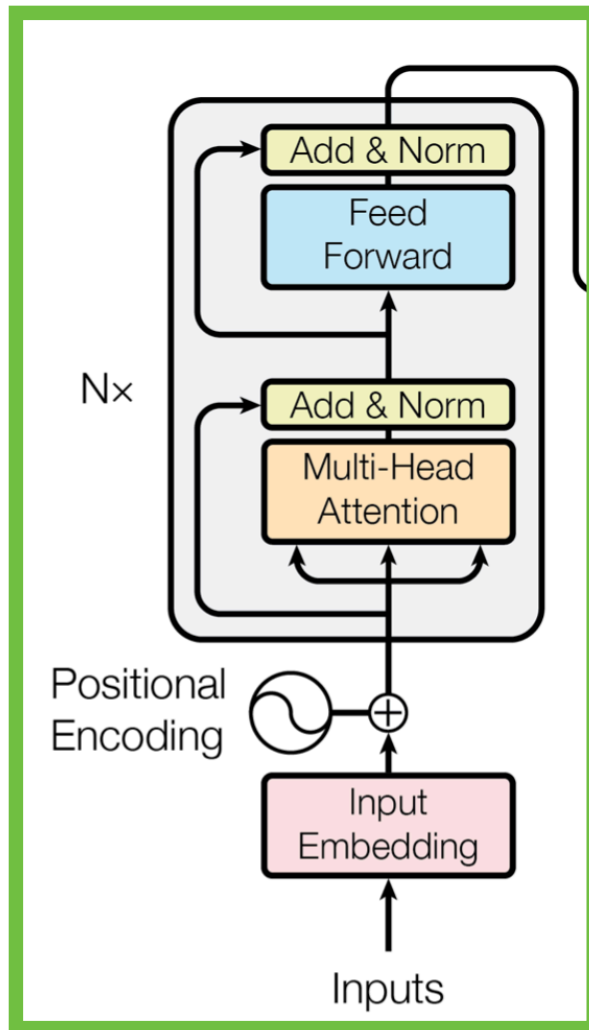
*decoder*



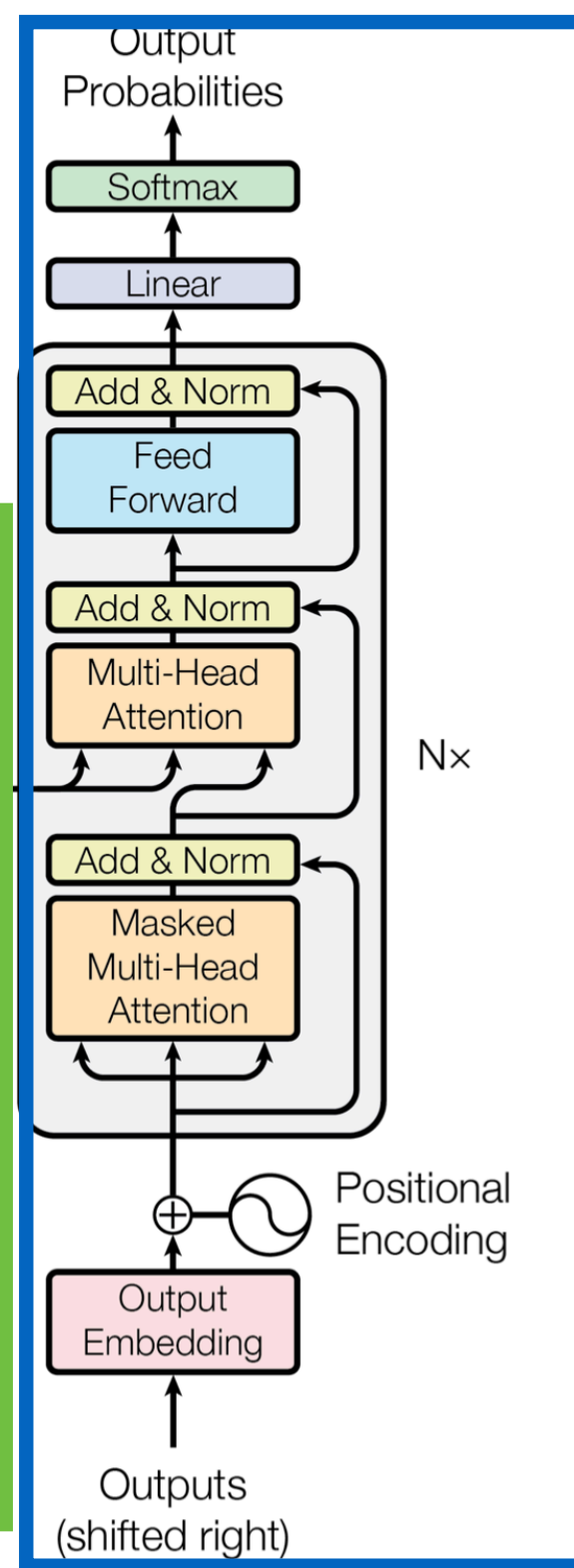


So far we've just talked about self-attention... what is all this other stuff?

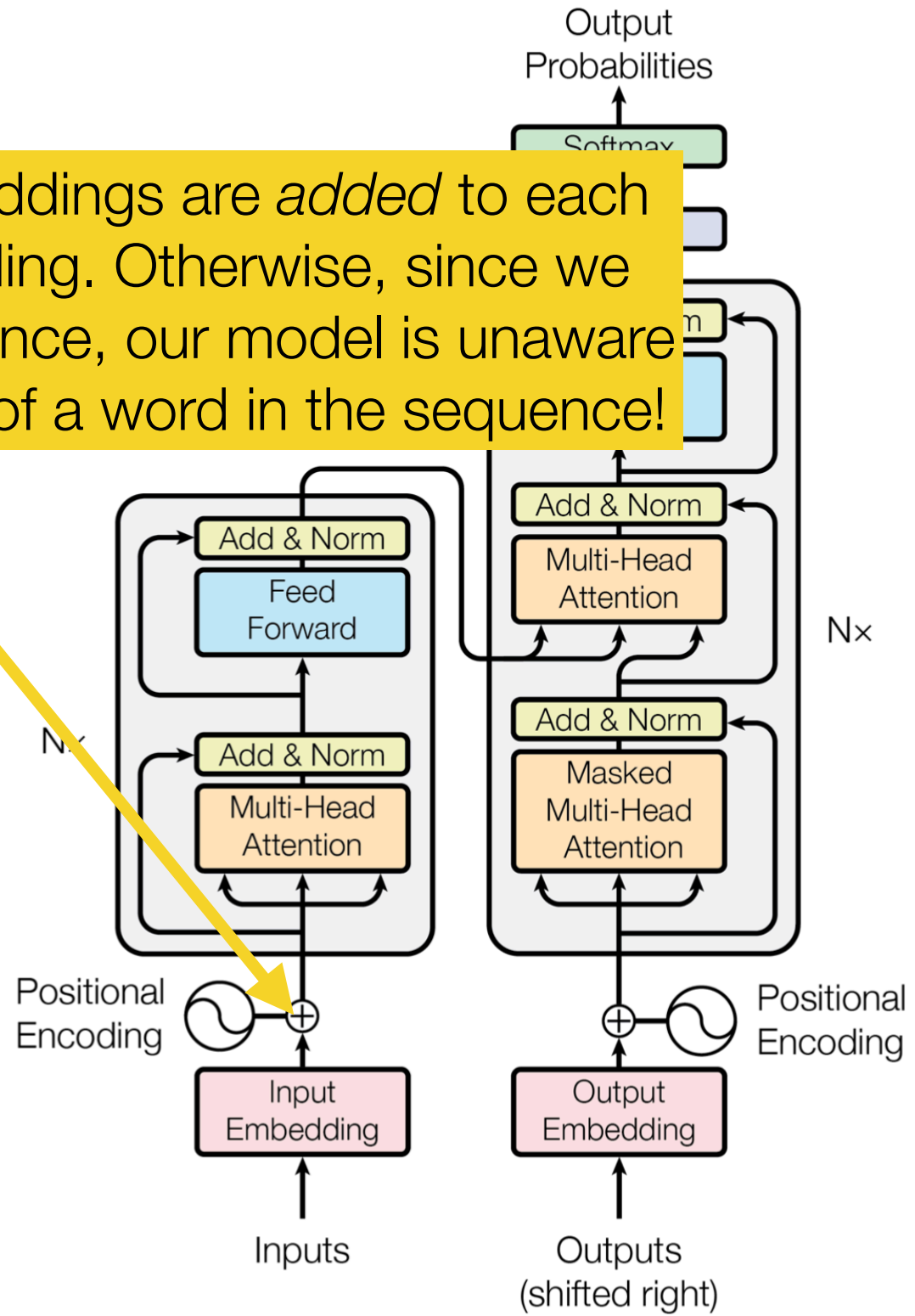
*encoder*



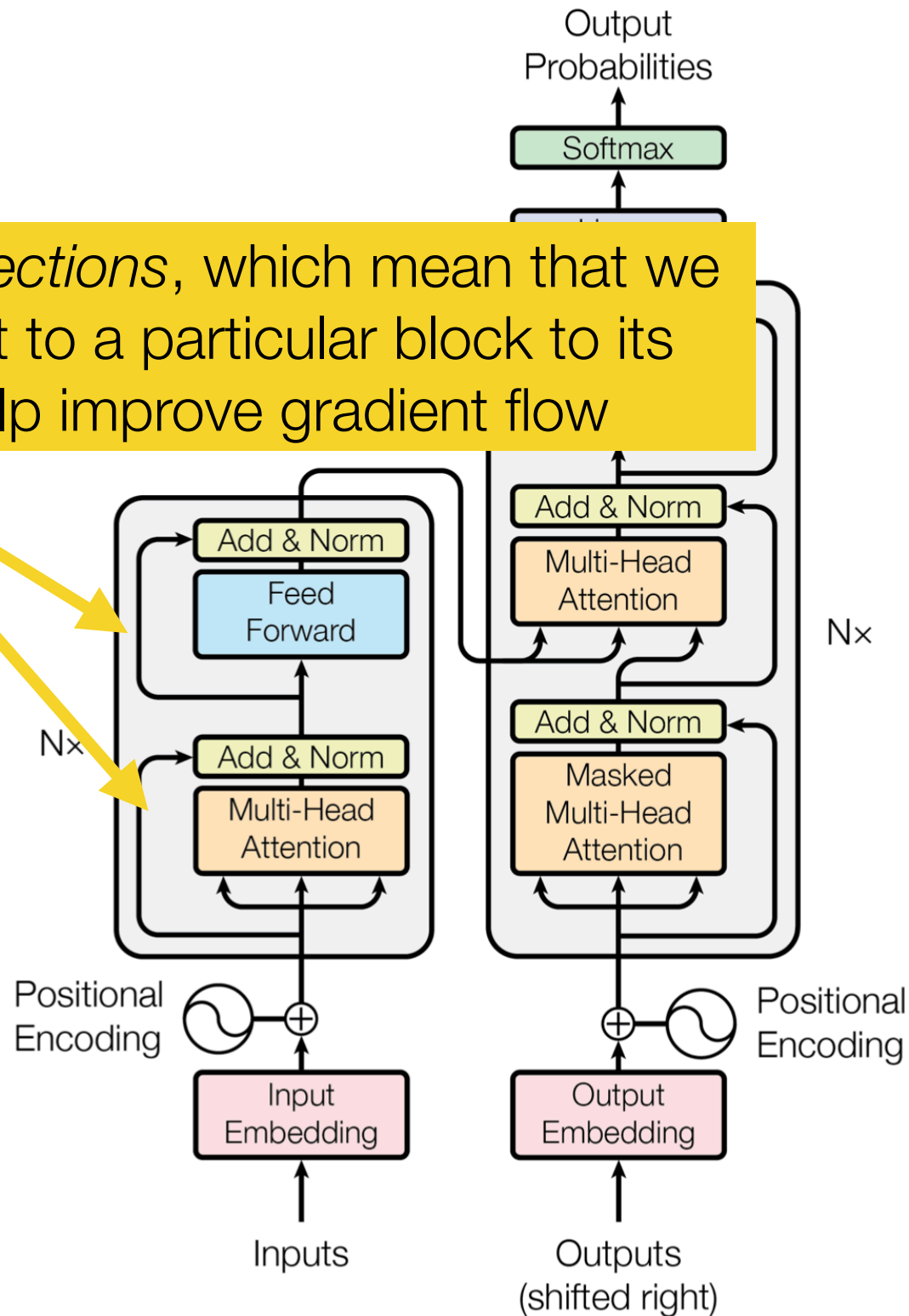
*decoder*



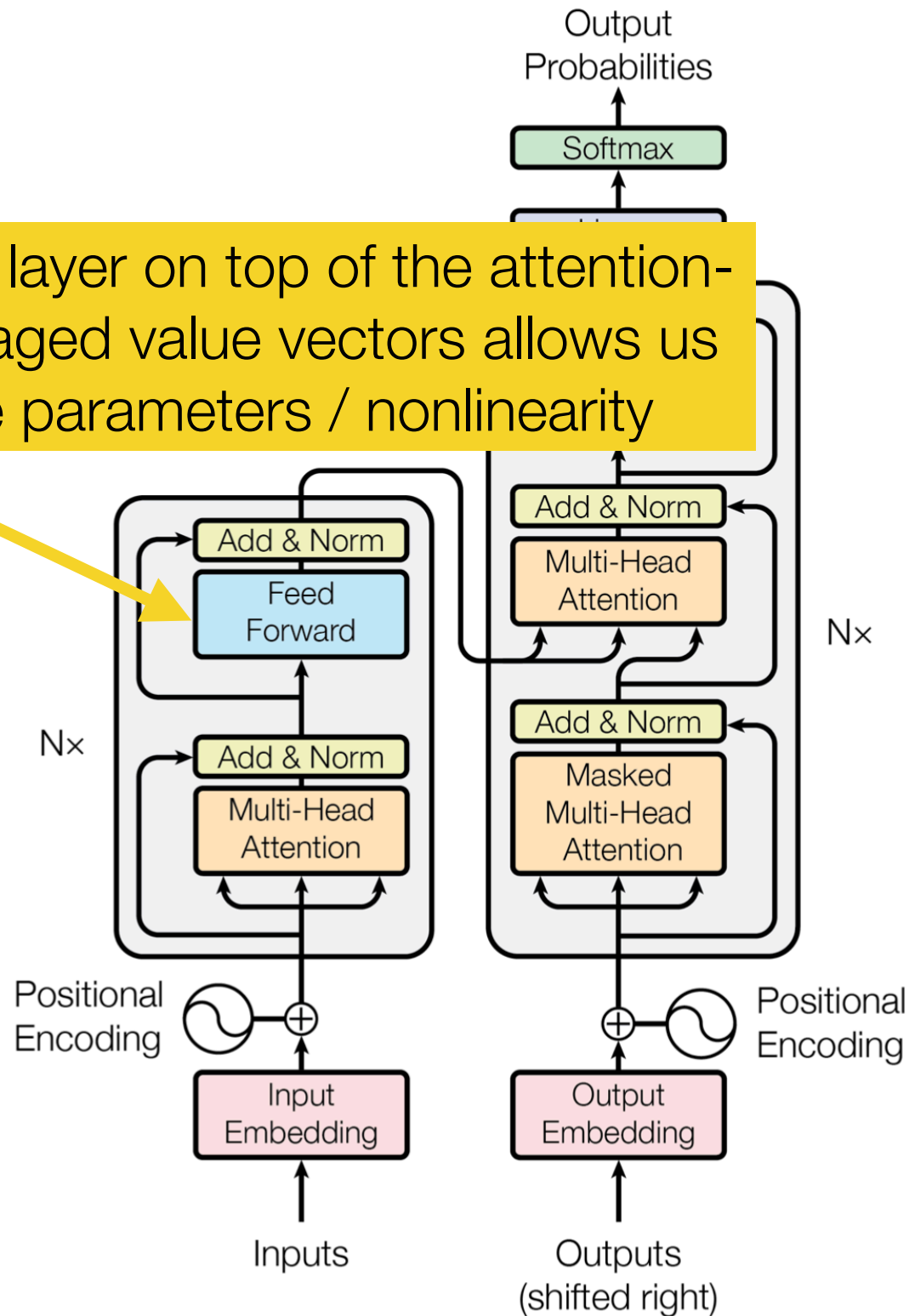
Position embeddings are *added* to each word embedding. Otherwise, since we have no recurrence, our model is unaware of the position of a word in the sequence!



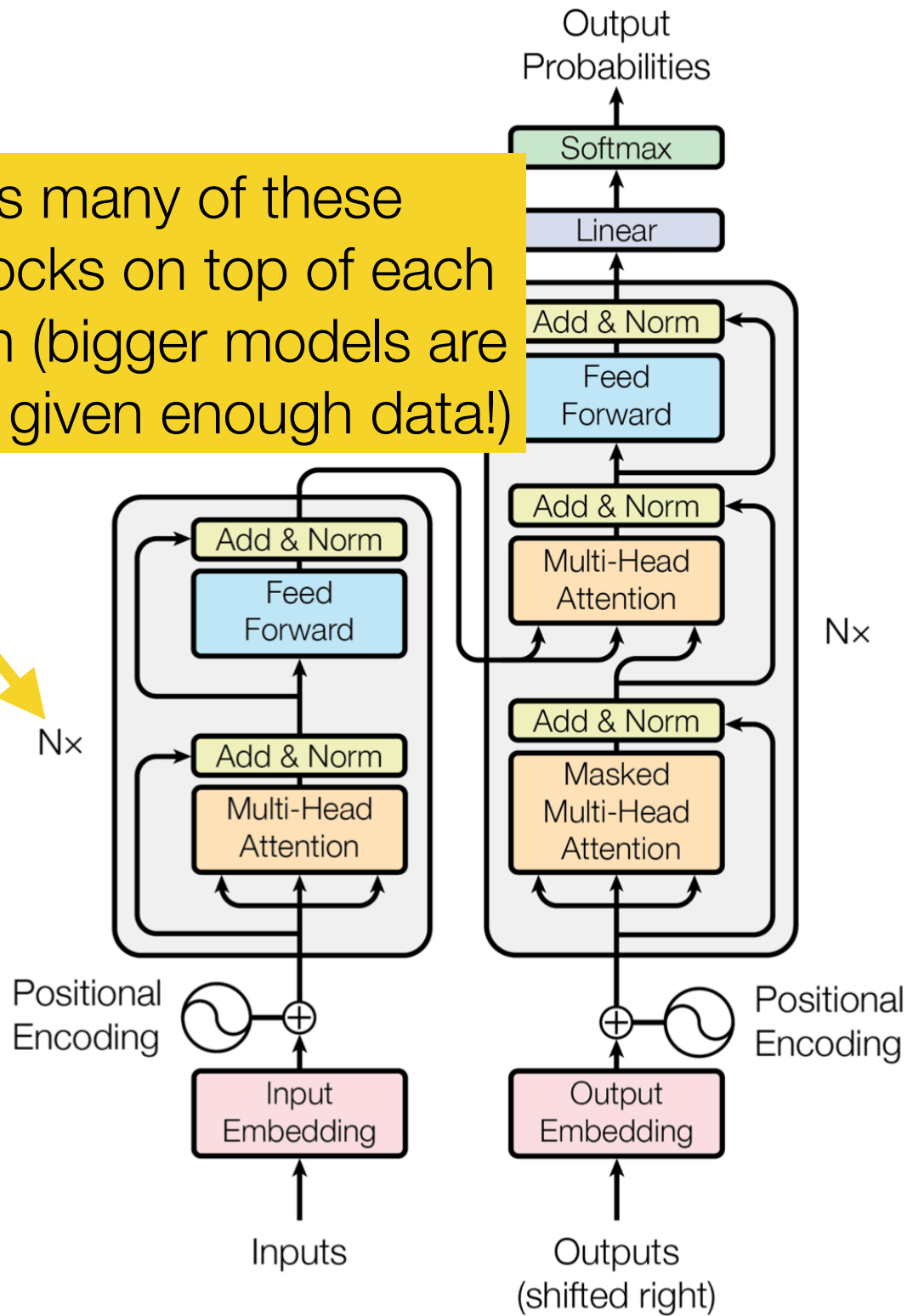
*Residual connections*, which mean that we add the input to a particular block to its output, help improve gradient flow



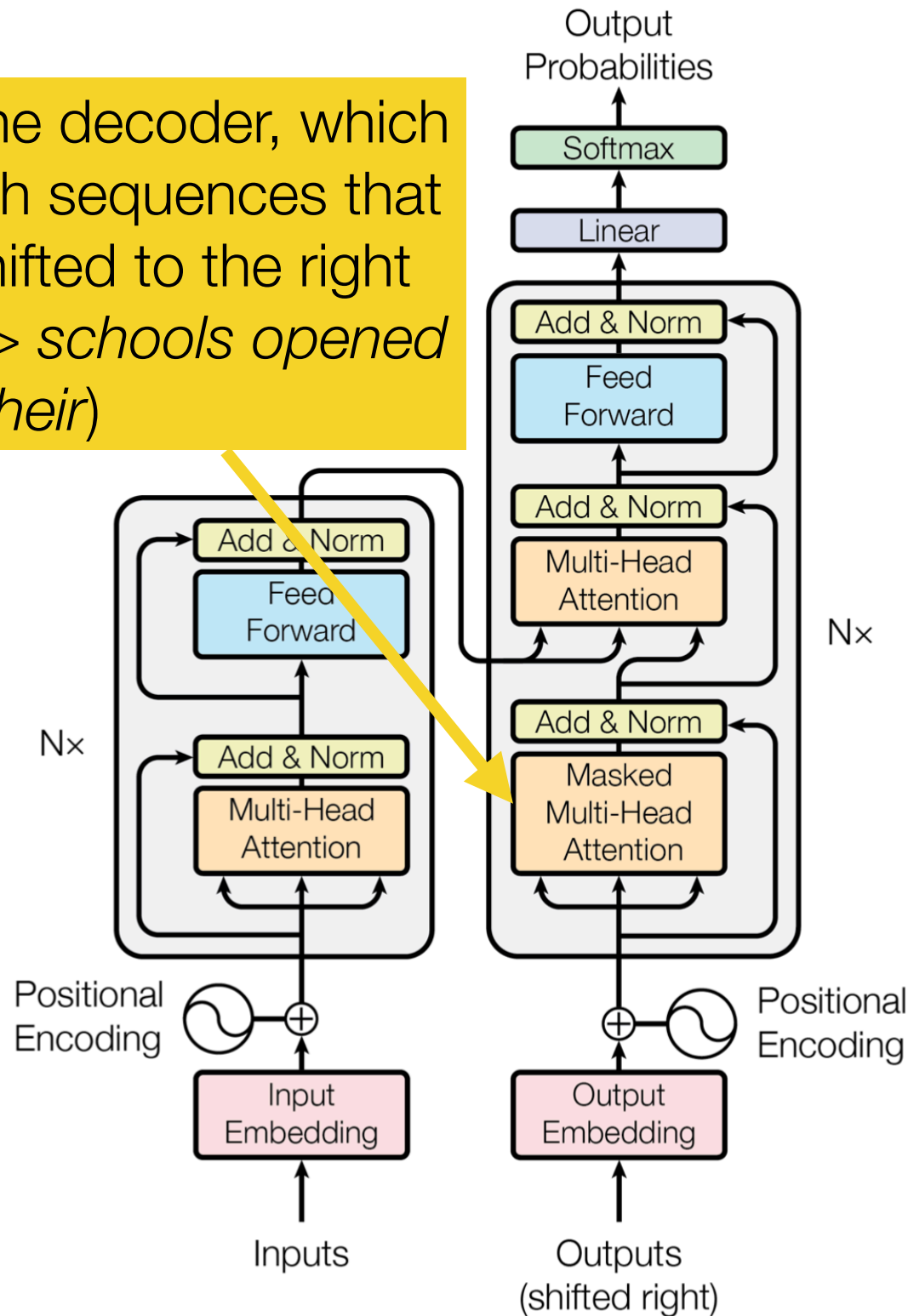
A feed-forward layer on top of the attention-weighted averaged value vectors allows us to add more parameters / nonlinearity



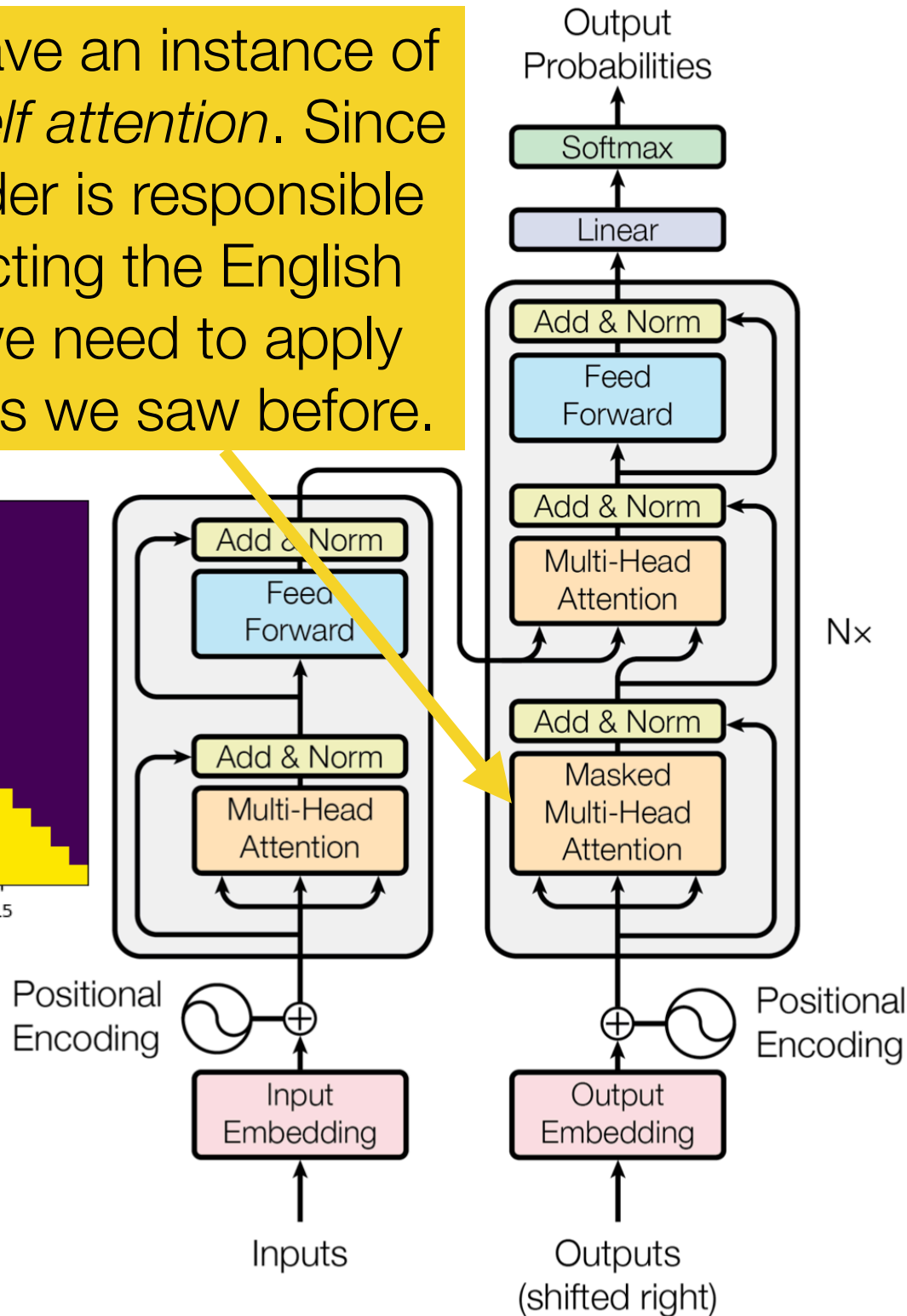
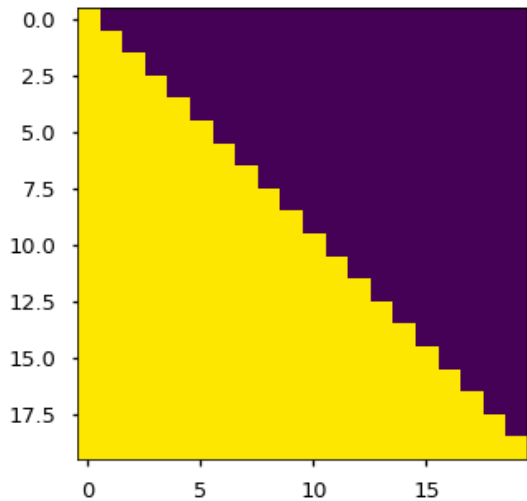
We stack as many of these *Transformer* blocks on top of each other as we can (bigger models are generally better given enough data!)



Moving onto the decoder, which takes in English sequences that have been shifted to the right (e.g., *<START> schools opened their*)

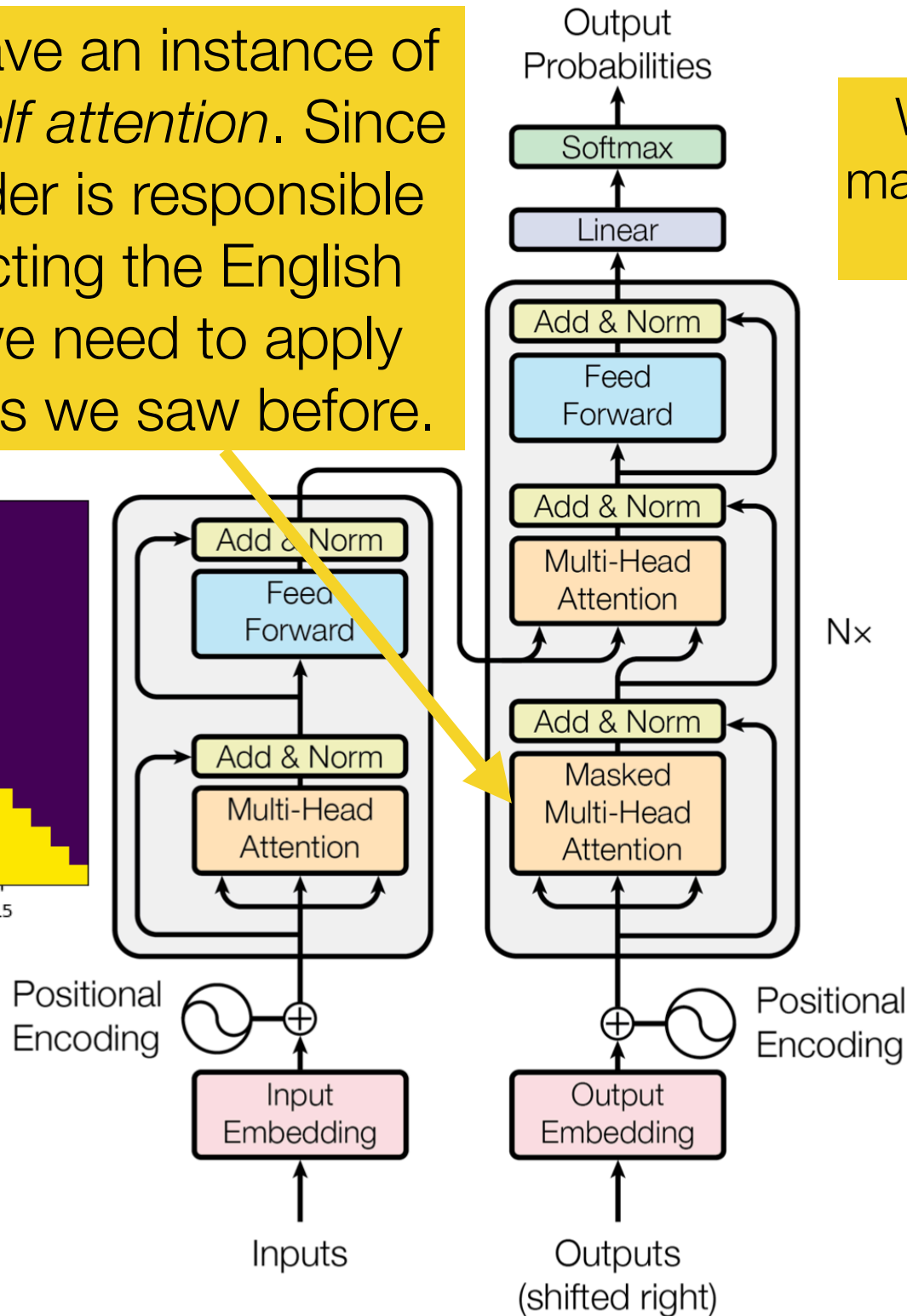
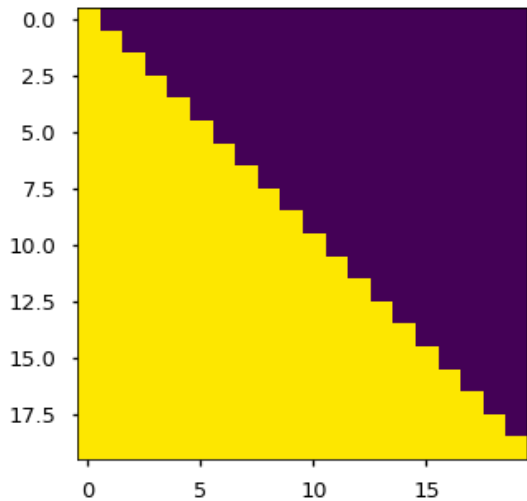


We first have an instance of *masked self attention*. Since the decoder is responsible for predicting the English words, we need to apply masking as we saw before.



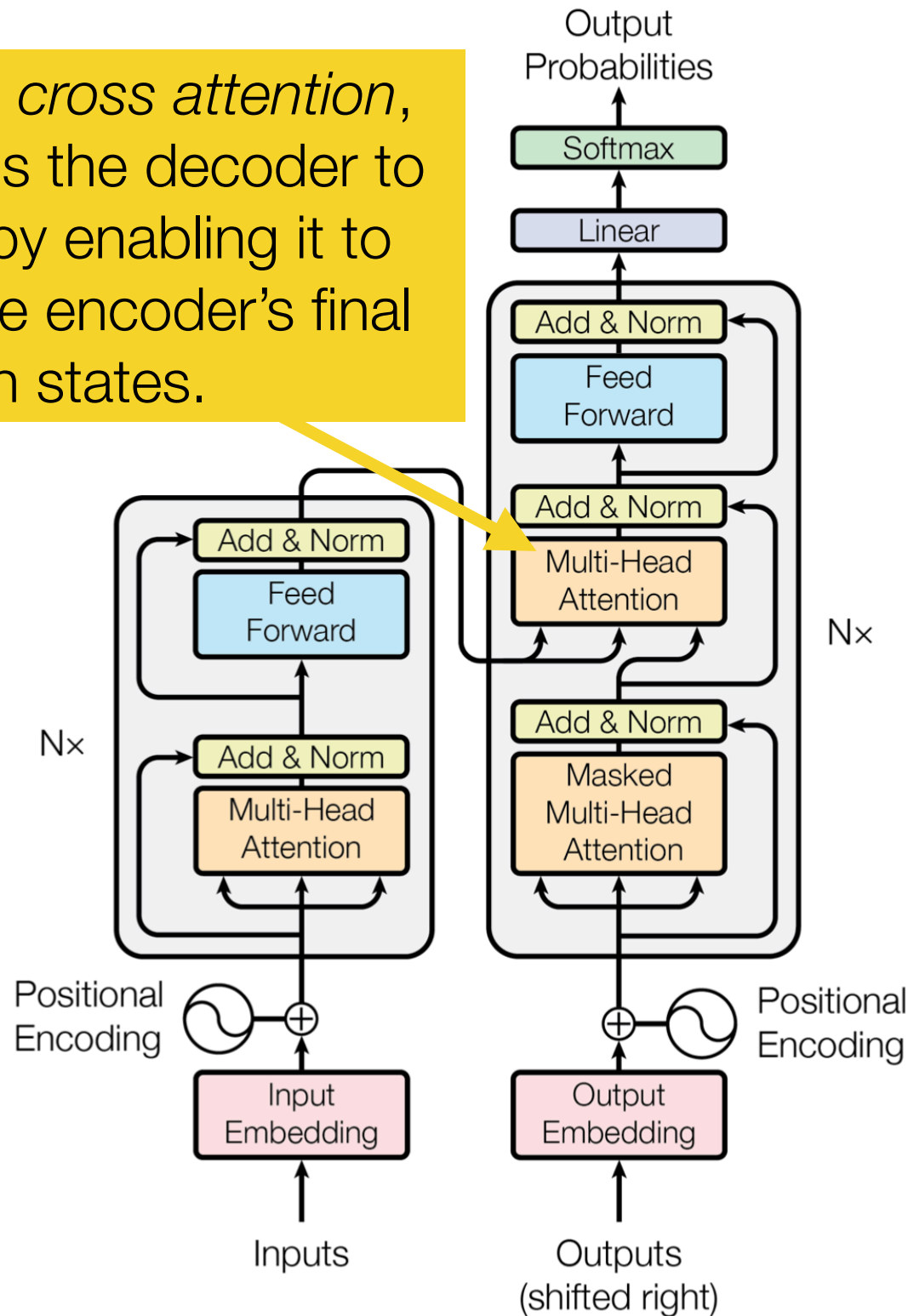
We first have an instance of *masked self attention*. Since the decoder is responsible for predicting the English words, we need to apply masking as we saw before.

Why don't we do masked self-attention in the encoder?





Now, we have *cross attention*, which connects the decoder to the encoder by enabling it to attend over the encoder's final hidden states.



After stacking a bunch of these decoder blocks, we finally have our familiar Softmax layer to predict the next English word

