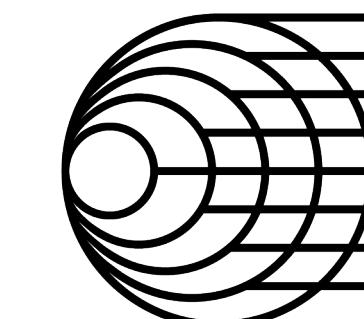
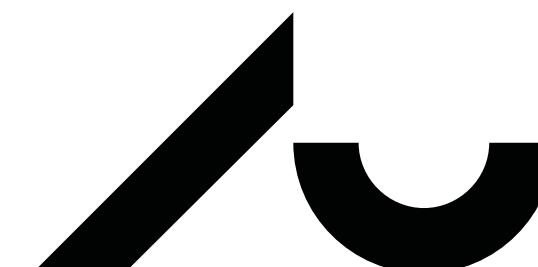


Language Generation

Natural language processing – Lecture 6

Kenneth Enevoldsen | 2024



CENTER FOR
HUMANITIES
COMPUTING

Agenda

- Midterm take-aways
- Language Generation
 - Broad overview
 - Recap on sentence representations
 - Recap quiz on attention
- Attention an in-depth example (from last time)
- The FFNN
- Attention for language generation
- The LM-prediction head
- (Early stage) Prompting
- Biases in generative LMs
- Next lecture overview
 - Transfer learning (BERT, T5), Scaling, Probing



Sources
& Notes



Mid-term evaluation

- Lecture
 - Generally Enjoy visual examples, Videos are great, More mental
 - Too theoretical/mathy
 - More code examples (is this a general thing?)
 - More focus on what we need for the exam
(I am not gonna do that)
 - Too fast through hard concepts - sometimes rapid jumps in levels of complexity
 - Lectures on Mondays
(not my decision)
 - Rapid jumps on complexity
(is this especially when it comes to the math?)



Sources
& Notes

Mid-term evaluation

- Classes
 - TAs are nice, notebooks are good, easy to follow
 - More focus on GitHub
 - Coding is too easy
 - More independent coding / more coding with the teachers
- Additional stuff
 - A brush up on the beginning of the course (programming workshop, what part is missing?)
 - Example of exam project (on brightspace + classes with presentations coming up)
 - Great response on emails
 - Introduce some readings which are post-chat gpt since chat-gpt revolutionized NLP
 - ChatGPT was mostly a UI improvement (most stuff is from 2017-2022)
 - Recent papers are typically optimizations papers (few of which generalize) or evaluation
 - Ask me for recommendations
 - Reading/homework is appropriate



Mid-term take-aways

- Generally two groups
 - Lecture are good, classes too easy
 - Classes are good, lecture too hard
- Goal: **Introduce concept and ideas** then move to the **formalizations** (it will not be removed)
 - Will attempt not to skip over the first step too quickly – don't have too large jumps
 - Will try to avoid assuming knowledge in the field
 - We will attempt to make classes more flexible and add a bit of theory there as well
- **Q: does this match your understanding or did I miss anything?**



Recap: Language generation



Sources
& Notes



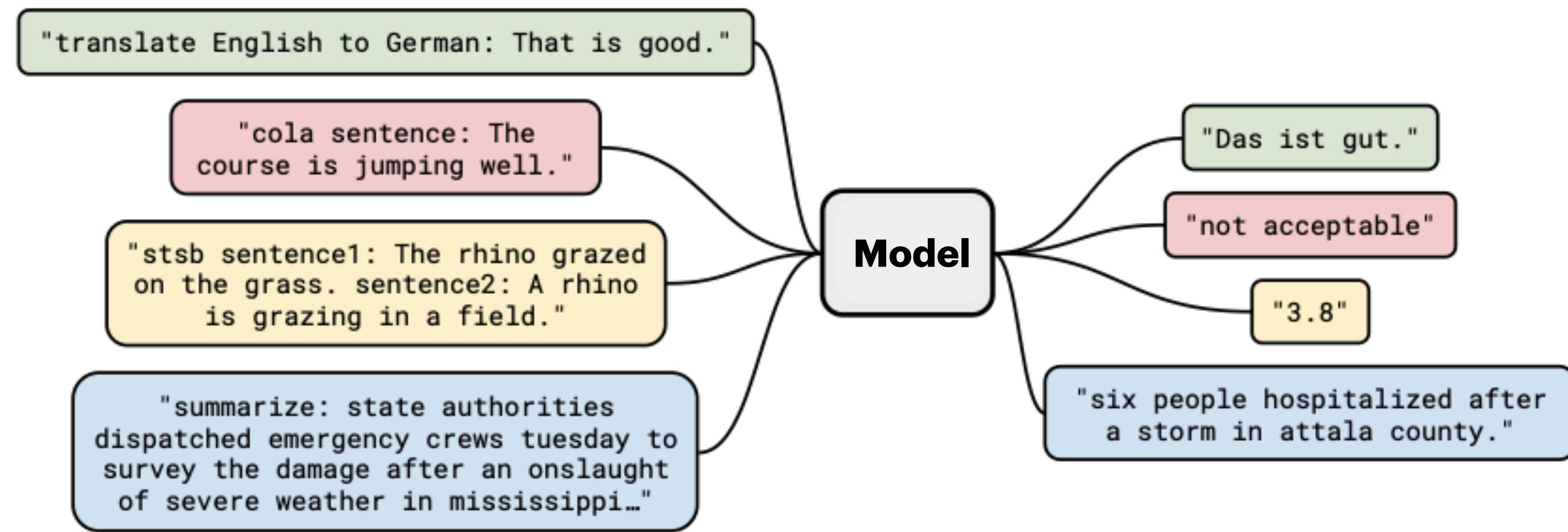
Recap: Early approaches for language generation

- Shannon (1948)
 - $P(\text{next word} \mid \text{previous words}) = P(w_i \mid w_{<i})$
 - $P(w_i \mid w_{i-1})$:
 - “REPRESENTING AND SPEEDILY IS AN GOOD APT OR COME CAN DIFFERENT NATURAL HERE HE THE A IN CAME THE TO OF TO EXPERT GRAY COME TO FURNISHES THE LINE MESSAGE HAD BE THESE.”
 - $P(w_i \mid w_{i-1}, w_{i-2})$:
 - “THE HEAD AND IN FRONTAL ATTACK ON AN ENGLISH WRITER THAT THE CHARACTER OF THIS POINT IS THEREFORE ANOTHER METHOD FOR THE LETTERS THAT THE TIME OF WHO EVER TOLD THE PROBLEM FOR AN UNEXPECTED.”



Why language generation?

- Unified formulation of task as Text-to-Text



- **Q:** Any tasks where this does not work?



From generative models to Chatbots

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>  
You are a helpful AI assistant for travel tips and recommendations<|  
eot_id|><|start_header_id|>user<|end_header_id|>  
  
What is France's capital?<|eot_id|><|start_header_id|>assistant<|  
end_header_id|>  
  
Bonjour! The capital of France is Paris!<|eot_id|><|start_header_id|>  
user<|end_header_id|>  
  
What can I do there?<|eot_id|><|start_header_id|>assistant<|  
end_header_id|>  
  
[...]
```



You are a helpful AI
assistant for travel tips
and recommendations

What is France's capital?



Bonjour! The capital of France is
Paris!

What can I do there?

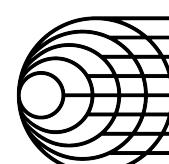


- **Note:** This would work with an generative model, but it would not work very well (more on this in lecture 8)



Sources
& Notes

Source: <https://www.llama.com/docs/model-cards-and-prompt-formats/meta-llama-3/>



CENTER FOR
HUMANITIES
COMPUTING

Current approaches (simplefied)

- $P(\text{next word} \mid \text{previous words}) = P(w_i \mid w_{<i})$
- What has changed?
 - $P(w_i \mid w_{i-1}, \dots, w_{i-n})$
 - n=1-3 → n>500
 - **Q:** Why didn't we do this earlier?



Sources
& Notes

Current approaches (simplefied)

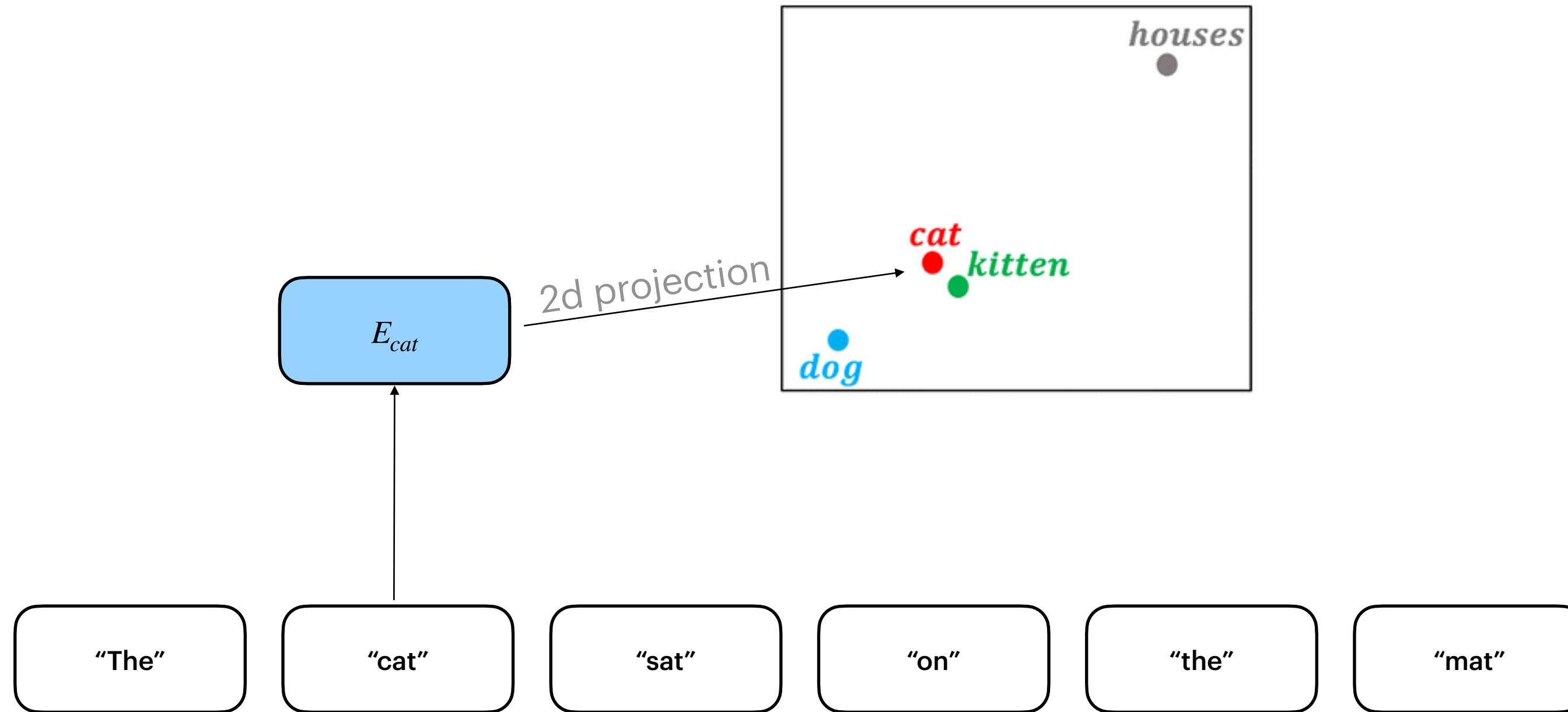
- $P(\text{next word} \mid \text{previous words}) = P(w_i \mid w_{<i})$
- What has changed?
 - $P(w_i \mid w_{i-1}, \dots, w_{i-n})$
 - $n=1-3 \rightarrow n>500$
 - **Q:** Why didn't we do this earlier?
 - Most >500 sequences are unique
 - Use of (contextual) word embeddings allow us to make approximates



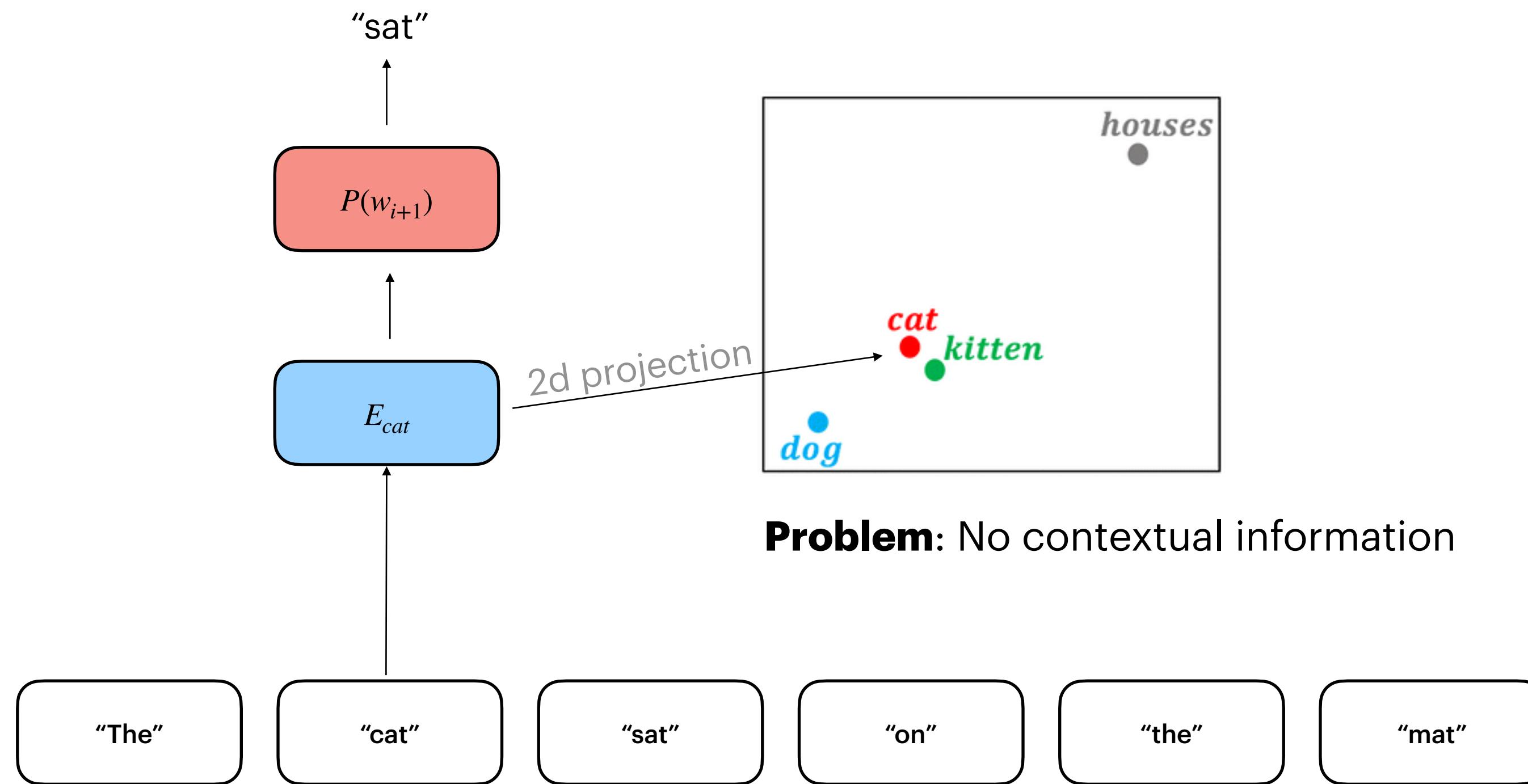
Sources
& Notes



Recap: Static Semantic Representations

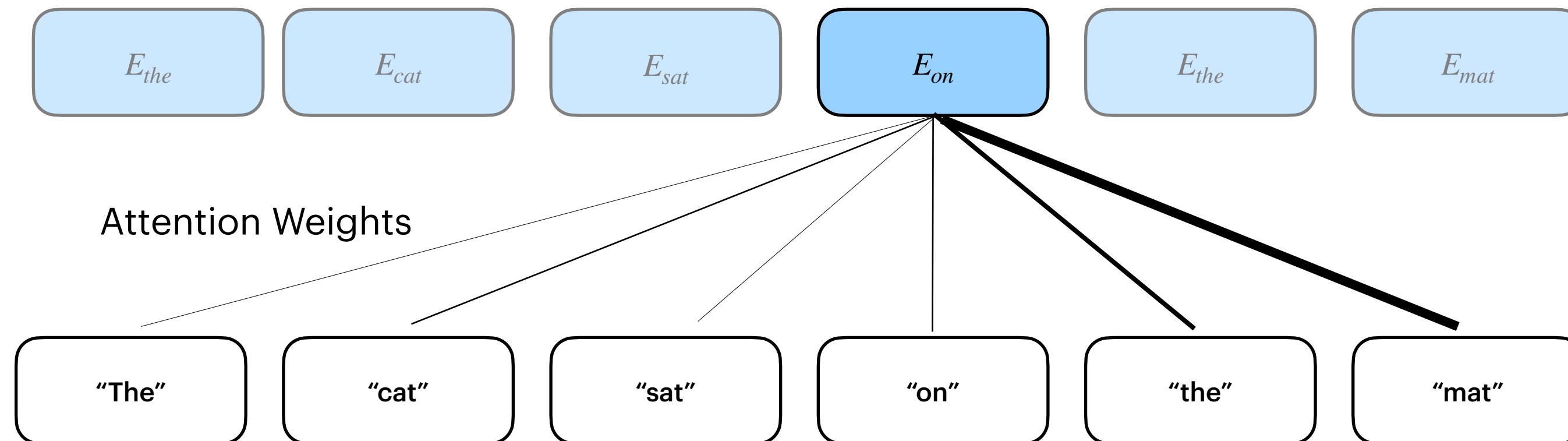


Goal: Contextual Representations



Recap: Attention

- Attention could help update our vectors to become contextual



- Quiz: <https://www.menti.com/aI4e34wejg8p>



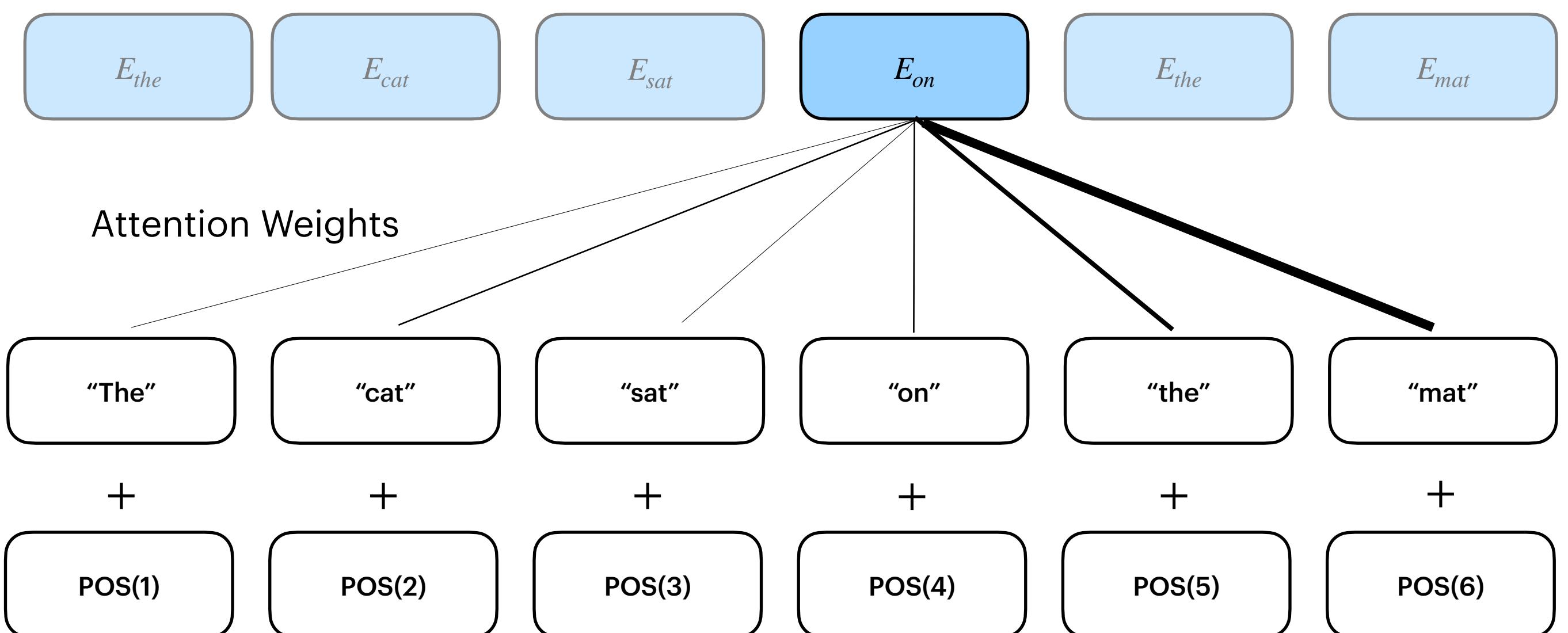
Sources
& Notes

Recap: Attention

- Attention as a weighted mean
- Solved two key limitations of RNNs
 - Informational bottleneck
 - Vanishing gradients
- **Position is added** using positional information*
 - Assumption: Words which appear closer are more relevant
 - We apply attention **multiple times**

$$\text{Attention}(Q, K, V) = \sigma\left(\frac{Q^T K}{\sqrt{d}}\right) V$$

Where: $Q = W_Q E$, $K = W_K E$, $V = W_V E$



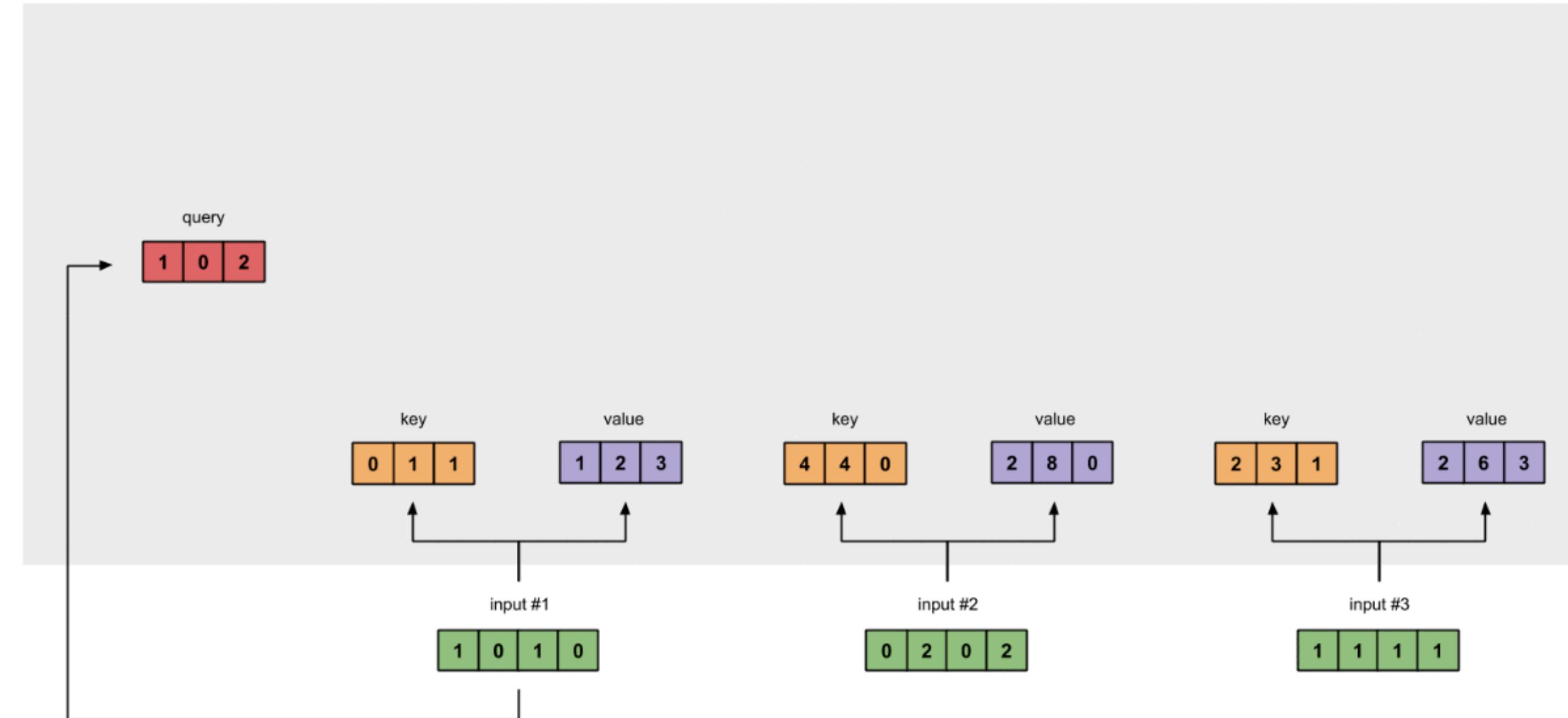
Or through changes to the attention matrix (e.g. Alibi)



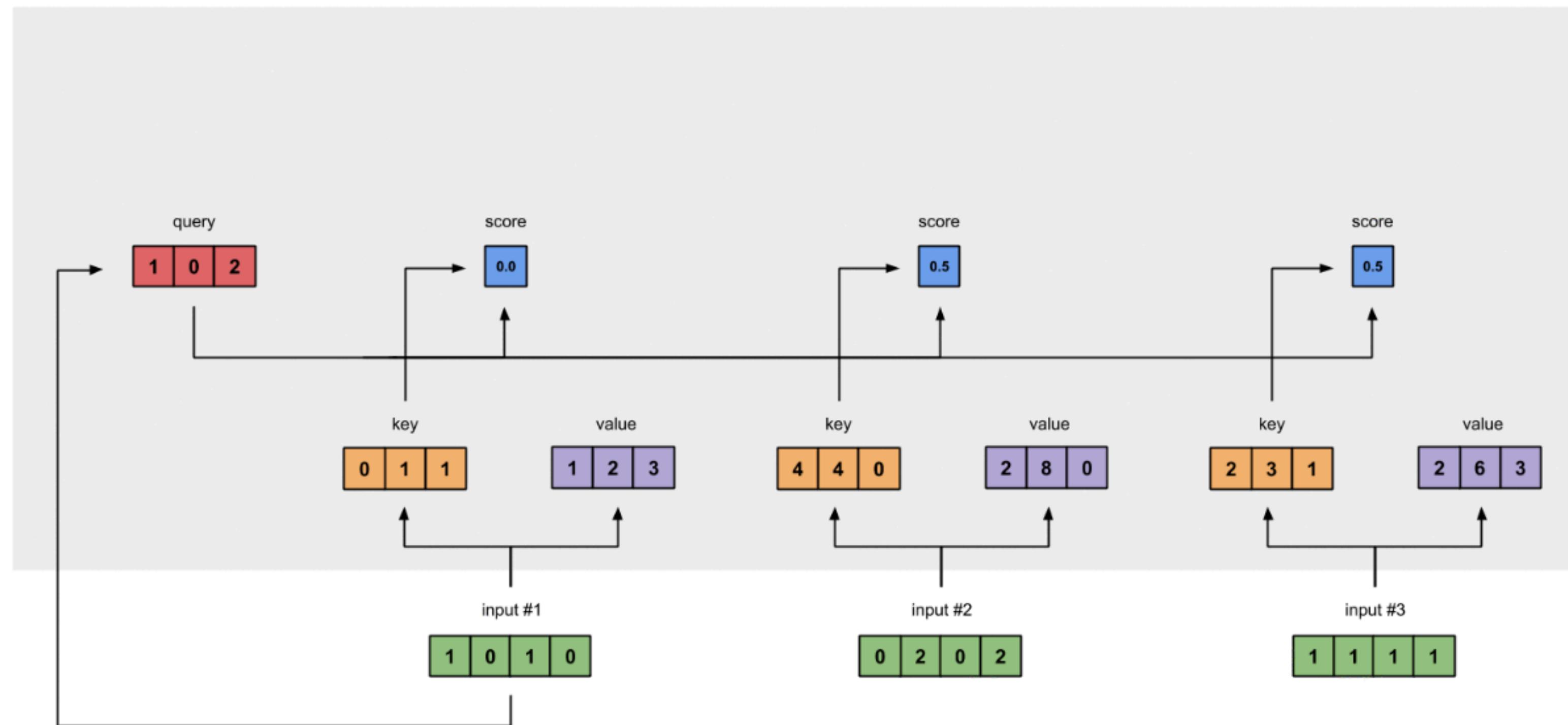
Sources
& Notes

Attention: An in-depth example

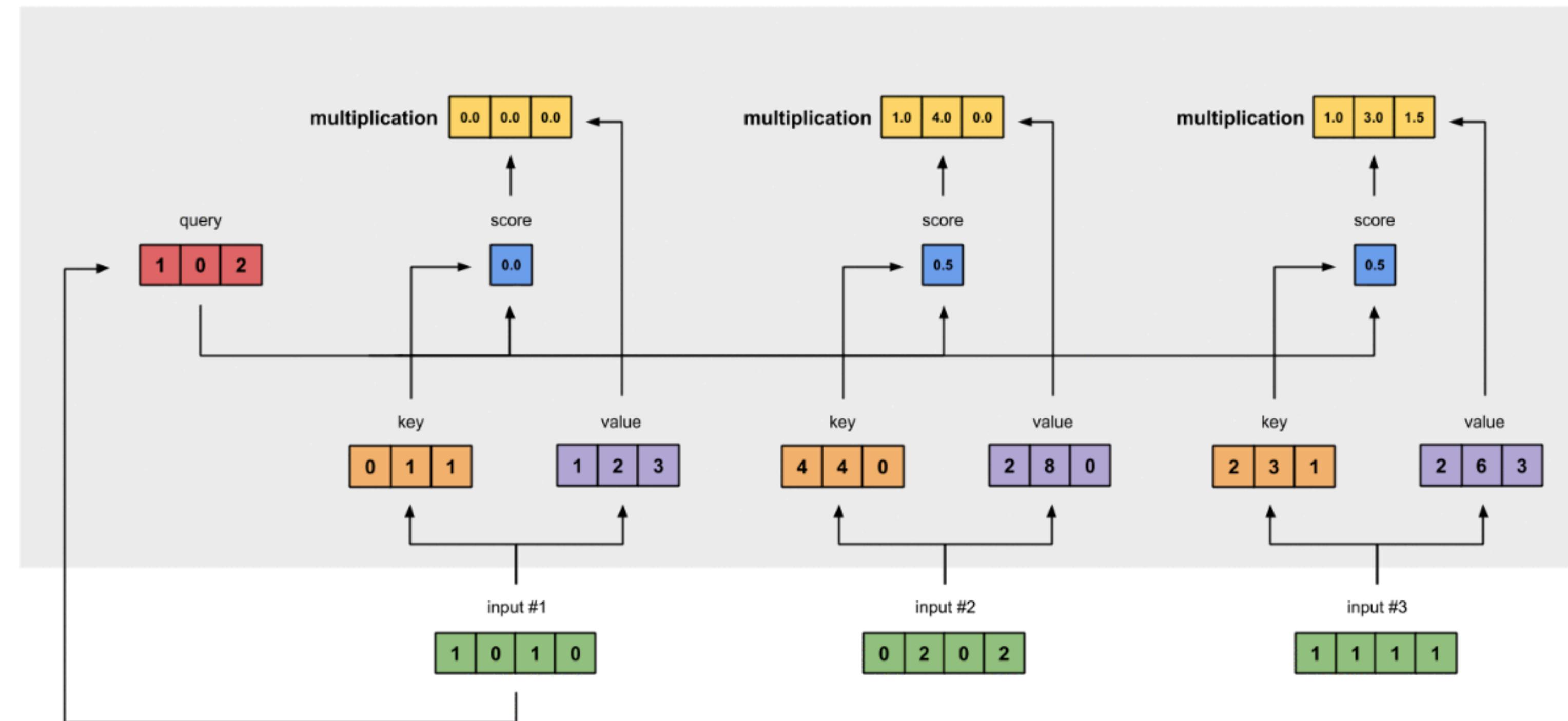
- Full guide with example calculations*
- Goal: Give an intuition for how the concepts fit together



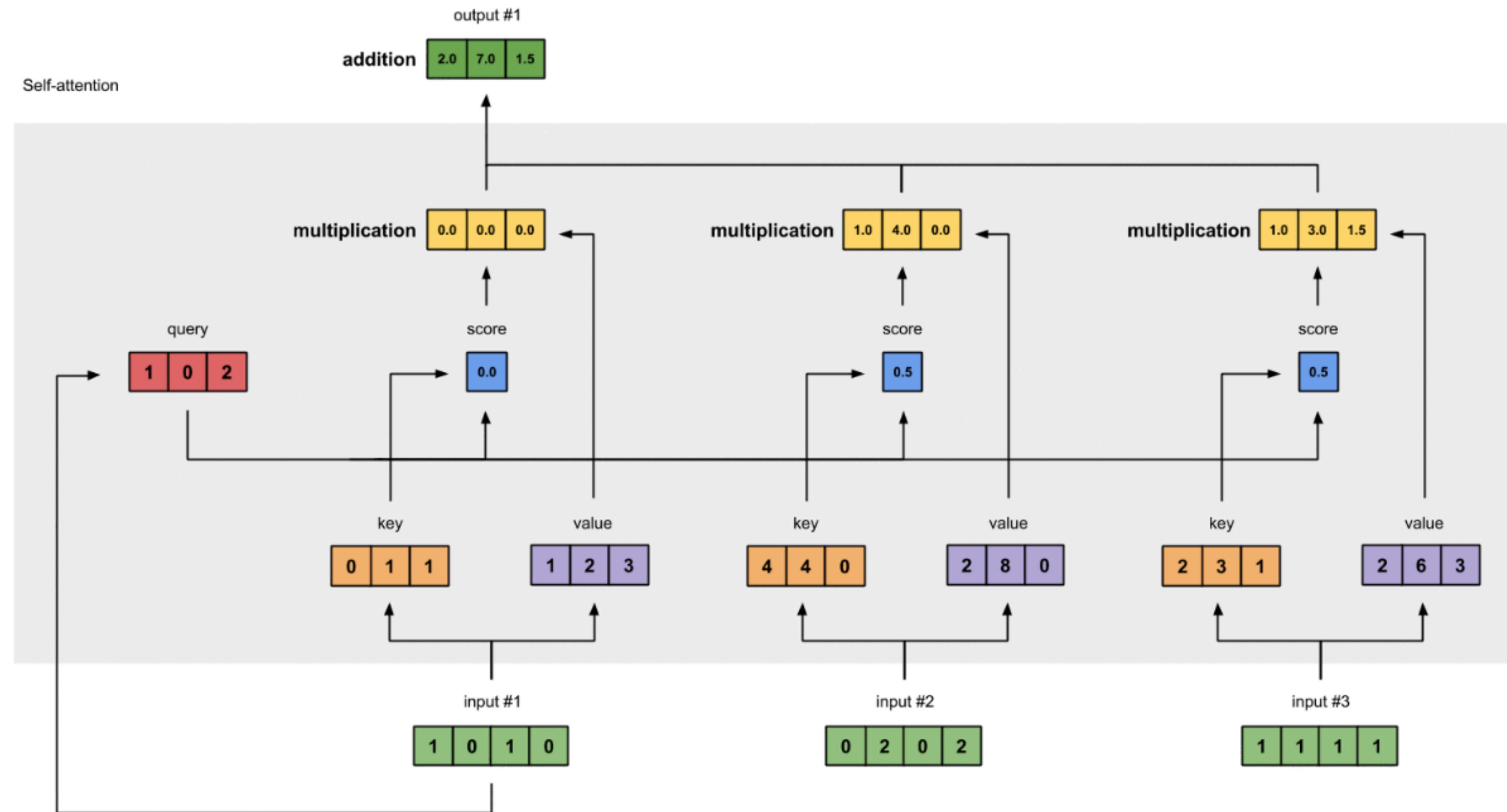
Attention: An in-depth example



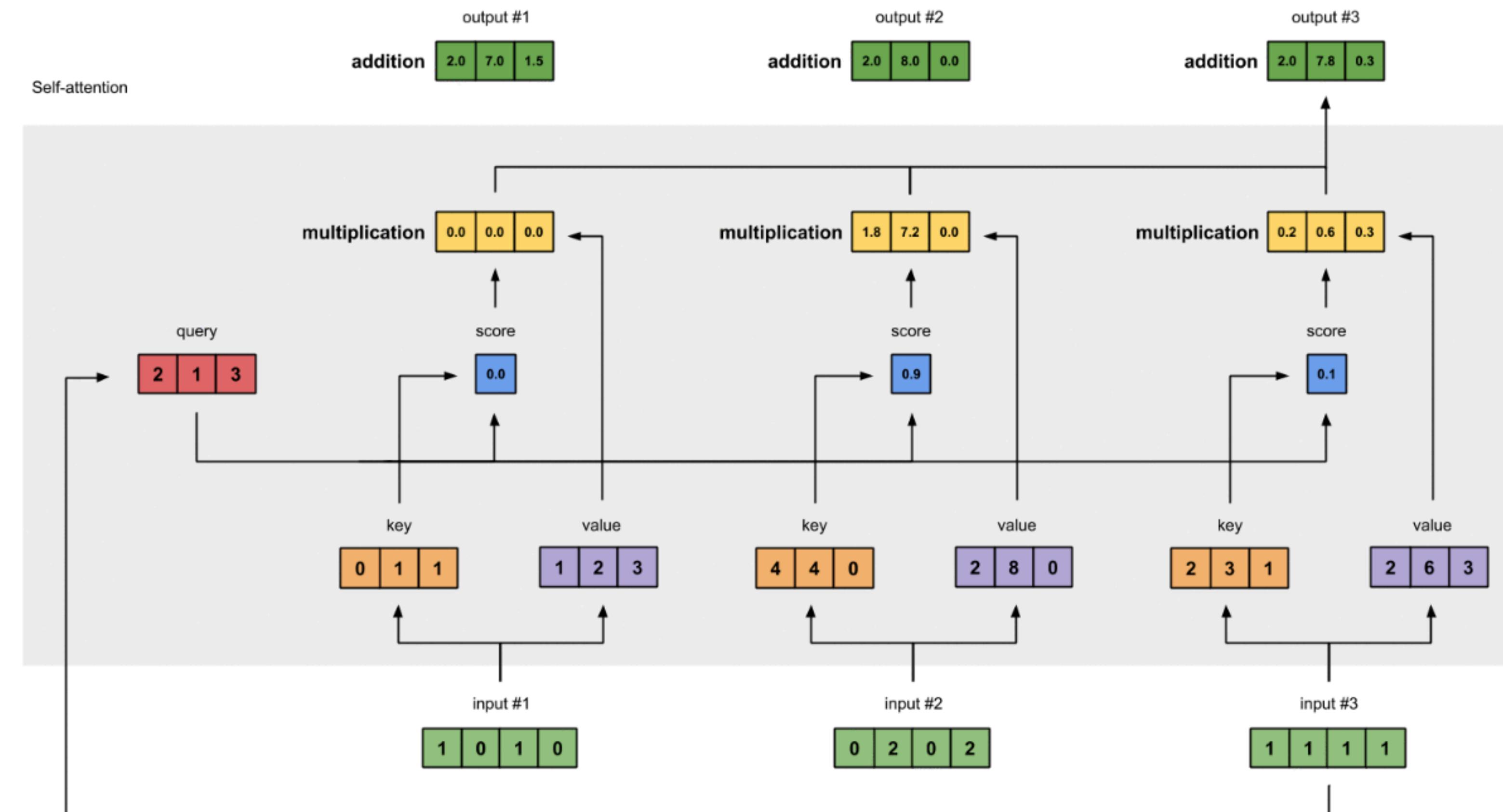
Attention: An in-depth example



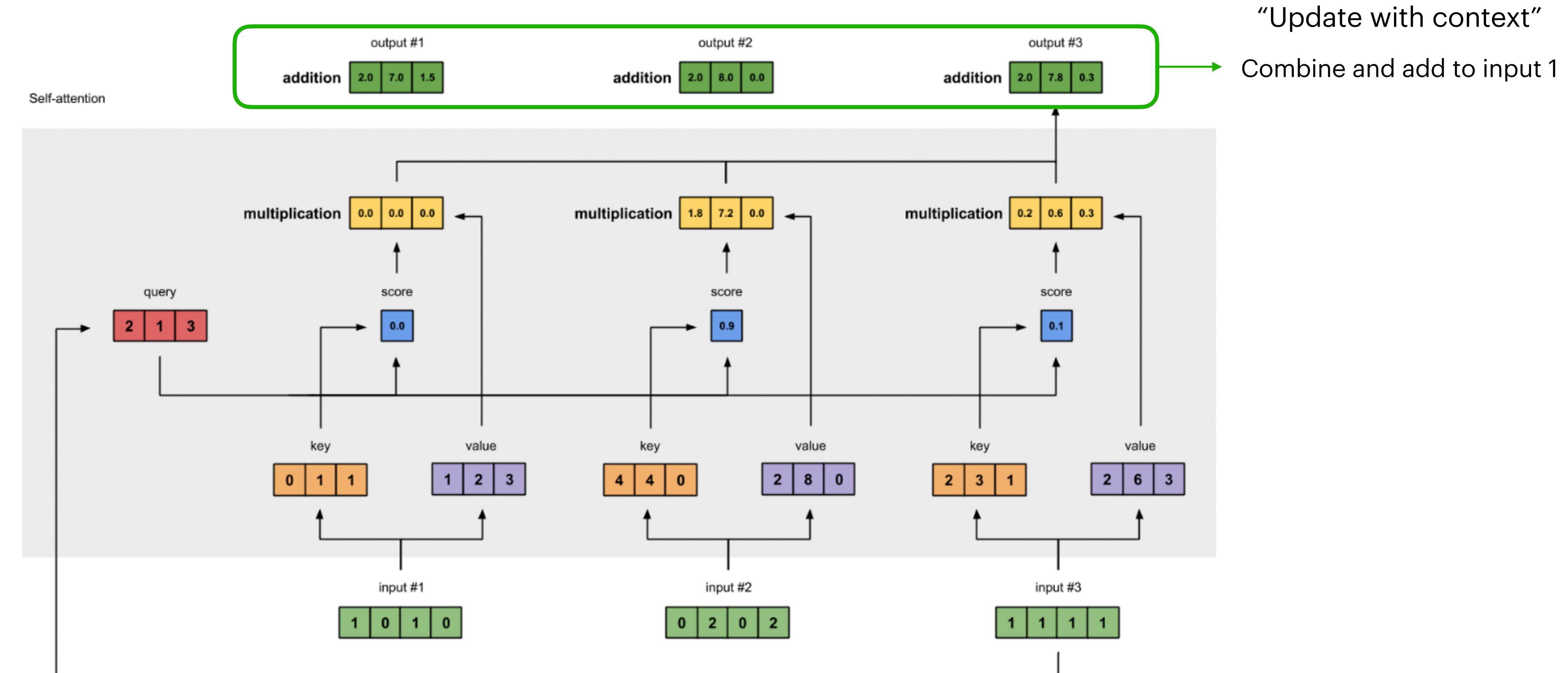
Attention: An in-depth example



Attention: An in-depth example

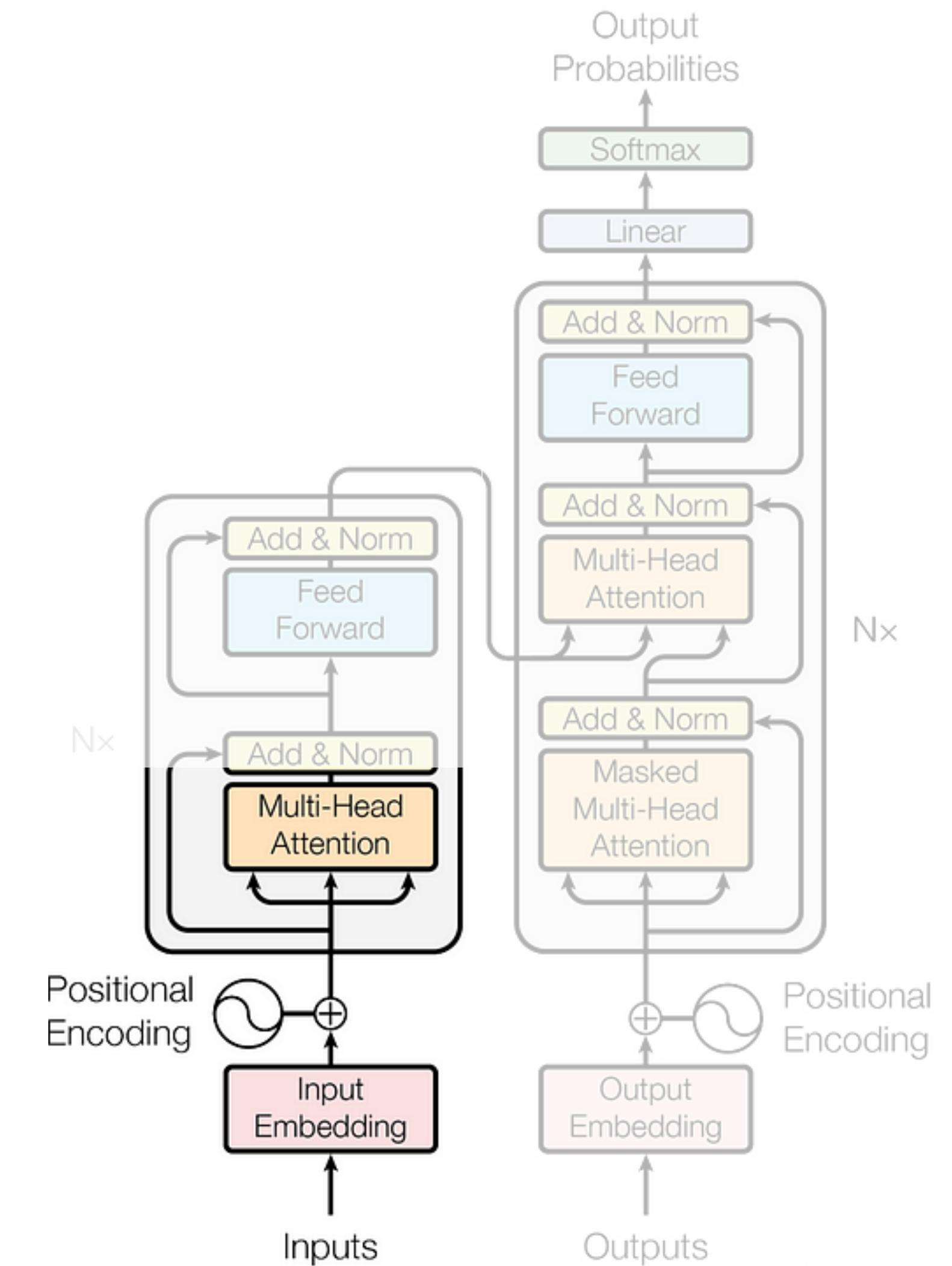


Attention: An in-depth example



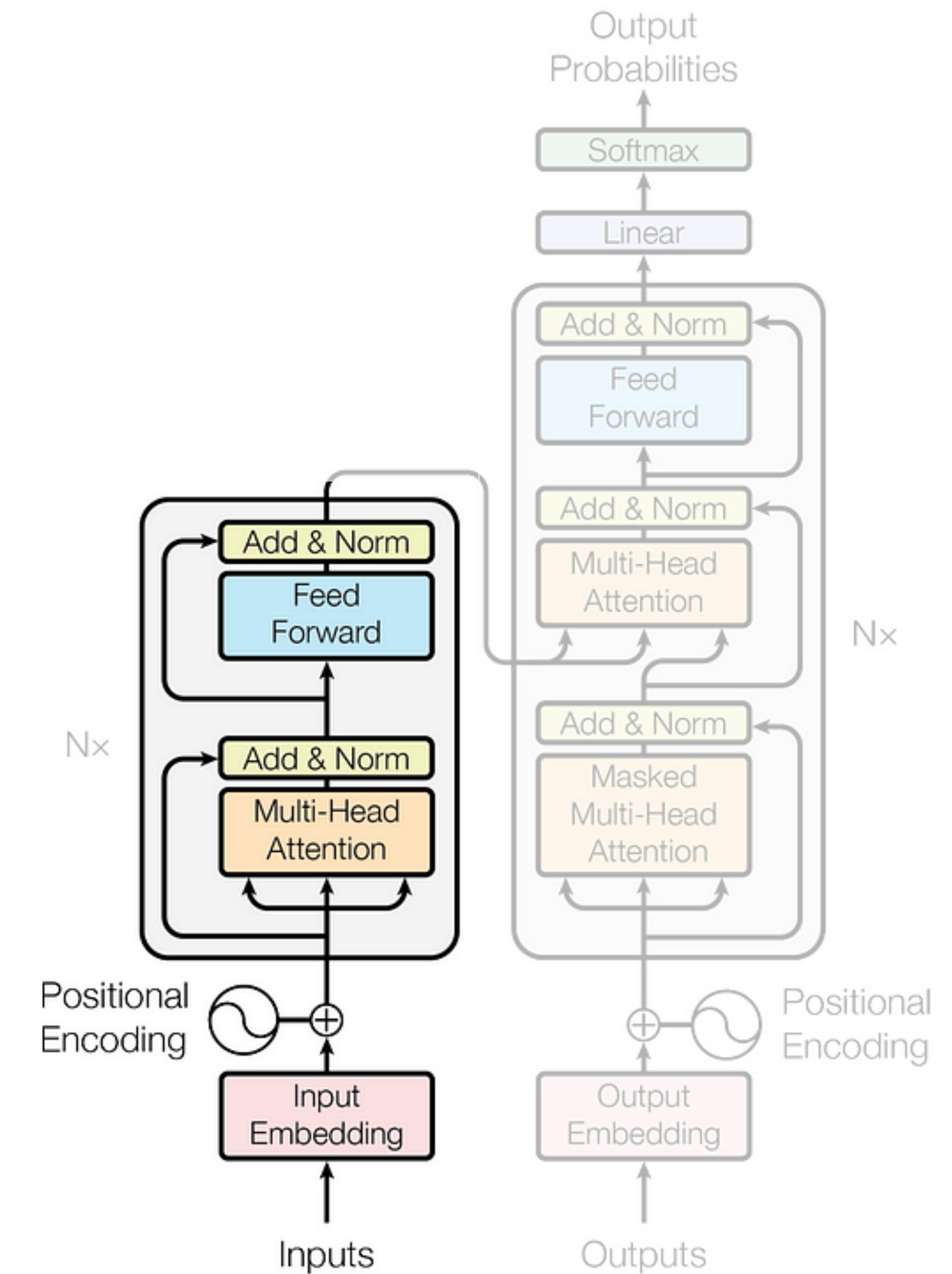
Transformer block

- So far we have looked at the multi head self attention



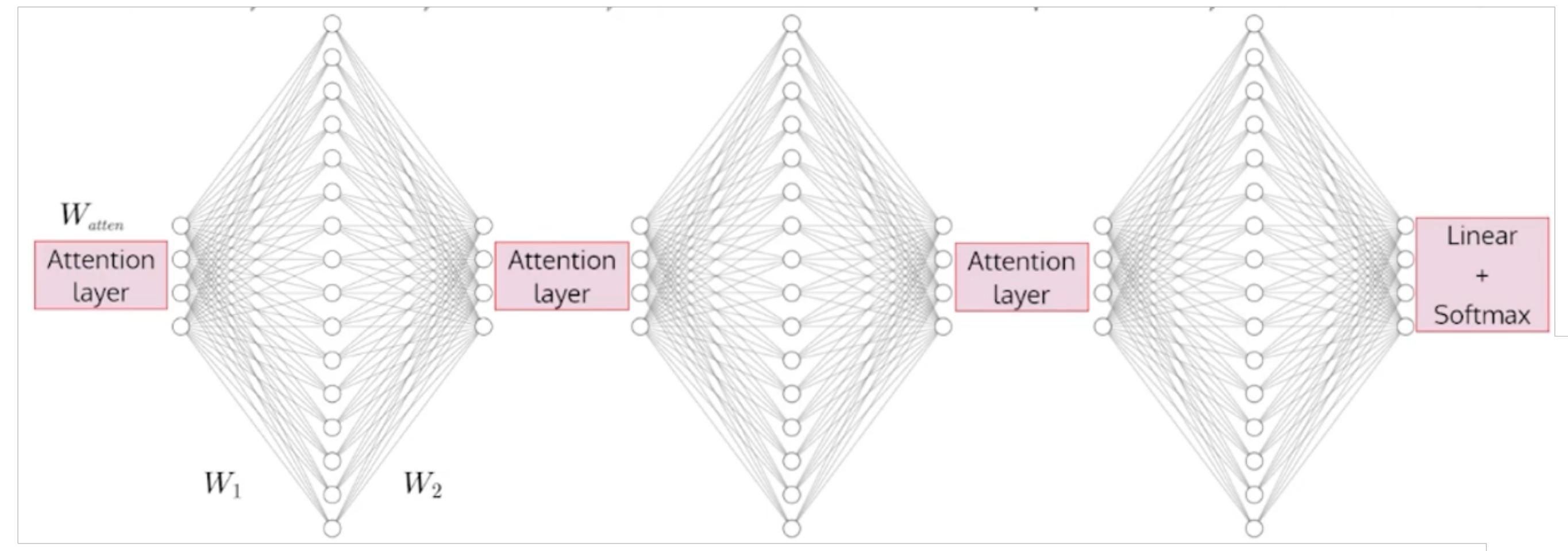
Transformer block

- So far we have looked at the multi head self attention
- Feed-forward layer



Feedforward Layer

- $FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$
 - How do we interpret this formula?



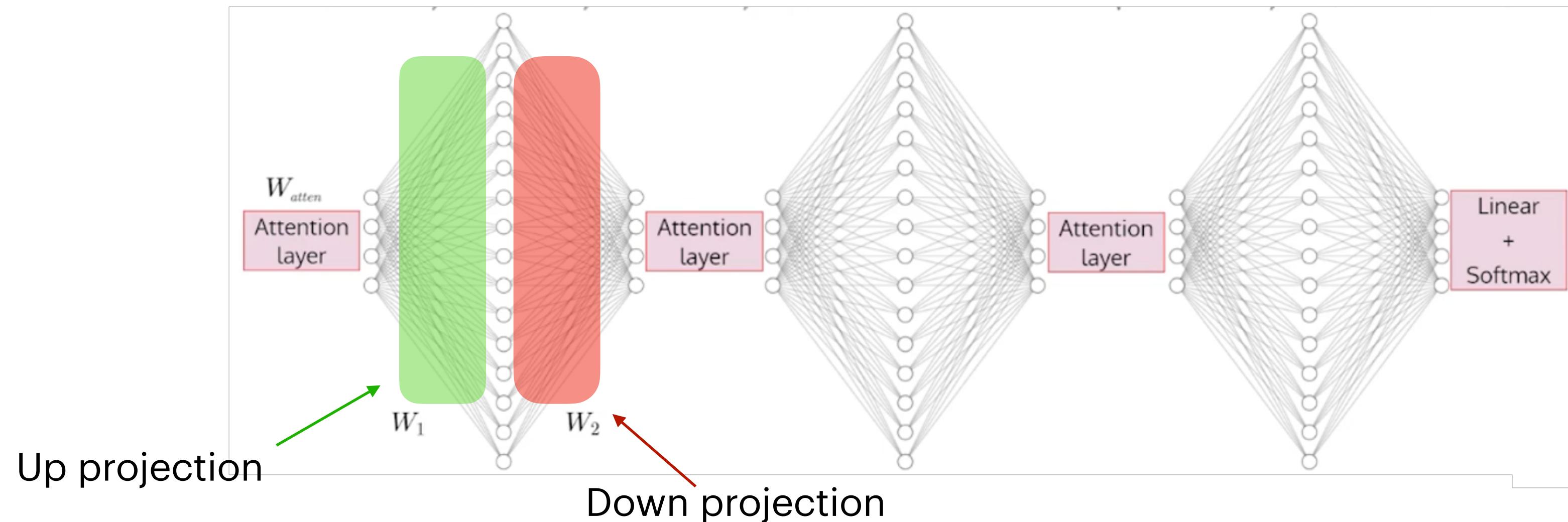
Sources
& Notes

Feedforward Layer

Linear transformation with a bias
(affine transformation)

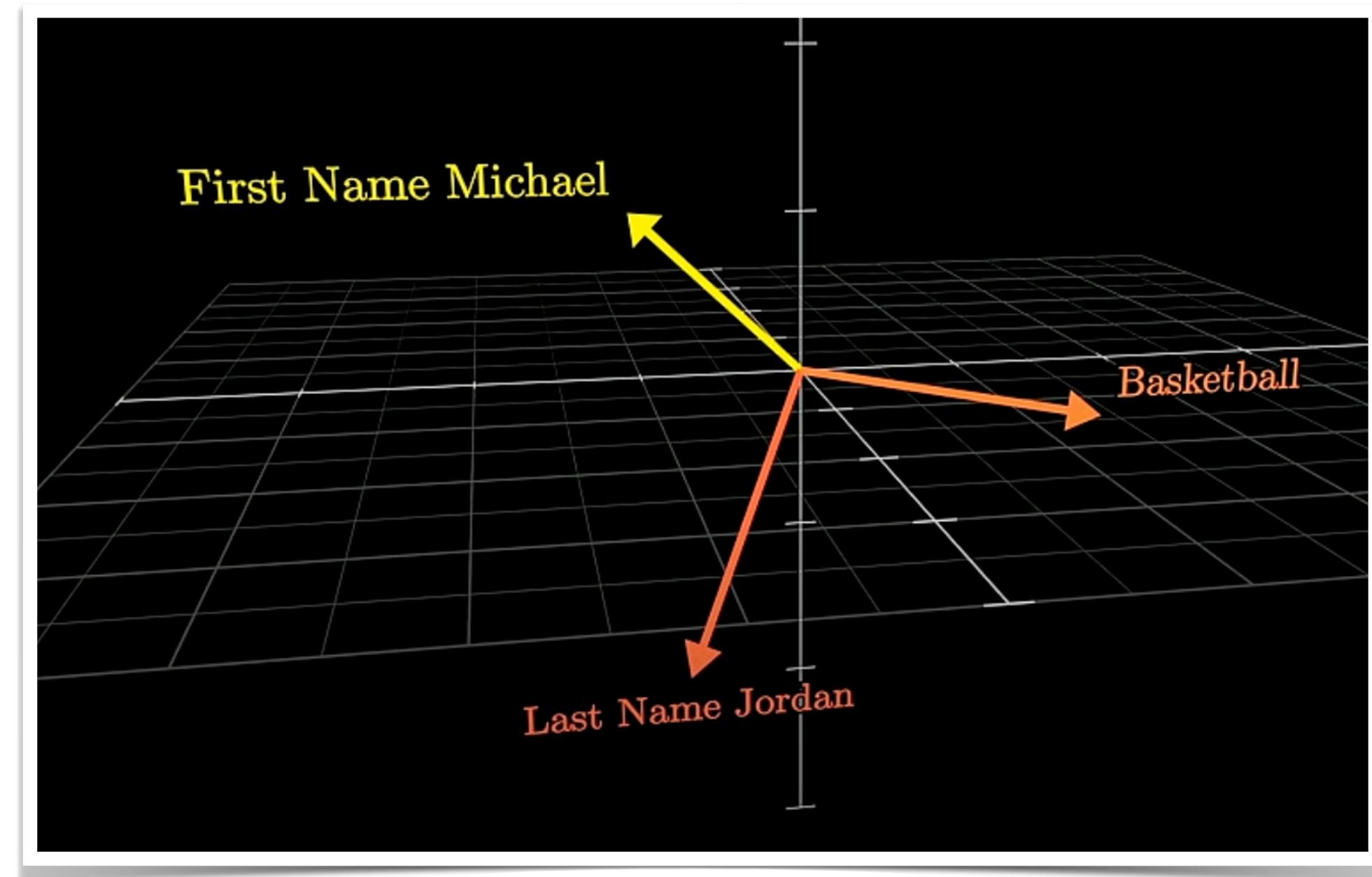
Non-linearity aka. Activation function

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



What could the FFN do?

- 3B1B shows that it can act as a lookup table for information
- Not necessarily what only what it does



What could the FFN do?

- Assume we have a vector that encodes for “Michael Jordan”
- A row in W_1 act as “detector”

$$\overrightarrow{\text{F.N. Michael}} + \overrightarrow{\text{L.N. Jordan}}$$
$$\parallel$$
$$\begin{bmatrix} \vec{R}_0 \\ \vec{R}_1 \\ \vec{R}_2 \\ \vdots \\ \vec{R}_n \end{bmatrix} \begin{bmatrix} | \\ \vec{E} \\ | \end{bmatrix} = \begin{bmatrix} \vec{R}_0 \cdot \vec{E} \\ \vec{R}_1 \cdot \vec{E} \\ \vec{R}_2 \cdot \vec{E} \\ \vdots \\ \vec{R}_n \cdot \vec{E} \end{bmatrix} = (\vec{M} + \vec{J}) \cdot \vec{E} = \underbrace{\vec{M} \cdot \vec{E} + \vec{J} \cdot \vec{E}}_{\approx 2 \text{ if } \vec{E} \text{ encodes ‘Michael Jordan’}} \leq 1 \text{ Otherwise}$$

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



What could the FFN do?

- Bias along with non-linearity acts as a boolean gate

Bias

$$\begin{bmatrix} -5.0 & +7.1 & +0.8 & +1.0 & \cdots & +6.8 \\ -7.4 & -4.4 & +1.7 & +9.3 & \cdots & +1.2 \\ -9.5 & +6.0 & -5.3 & +6.1 & \cdots & -2.2 \\ +7.2 & +4.9 & +1.1 & -7.2 & \cdots & -8.7 \\ -7.5 & -9.0 & -7.8 & -5.4 & \cdots & +4.2 \\ +1.2 & -9.7 & -8.5 & +9.3 & \cdots & +1.3 \\ -5.9 & -4.9 & +4.8 & -6.0 & \cdots & +1.6 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ +9.3 & +6.9 & -5.2 & -0.1 & \cdots & +2.4 \end{bmatrix} \begin{bmatrix} -7.1 \\ -6.0 \\ +6.0 \\ +9.3 \\ \vdots \\ +3.8 \end{bmatrix} + \begin{bmatrix} -1.0 \\ -5.7 \\ +5.0 \\ -8.6 \\ -4.7 \\ +6.0 \\ \vdots \\ -6.1 \\ +2.8 \end{bmatrix} = \begin{bmatrix} +1.0 \\ -8.9 \\ +1.5 \\ -7.1 \\ +1.8 \\ +4.0 \\ \vdots \\ -8.0 \\ +3.9 \end{bmatrix} = \vec{M} \cdot \vec{E} + \vec{J} \cdot \vec{E} \boxed{-1}$$

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

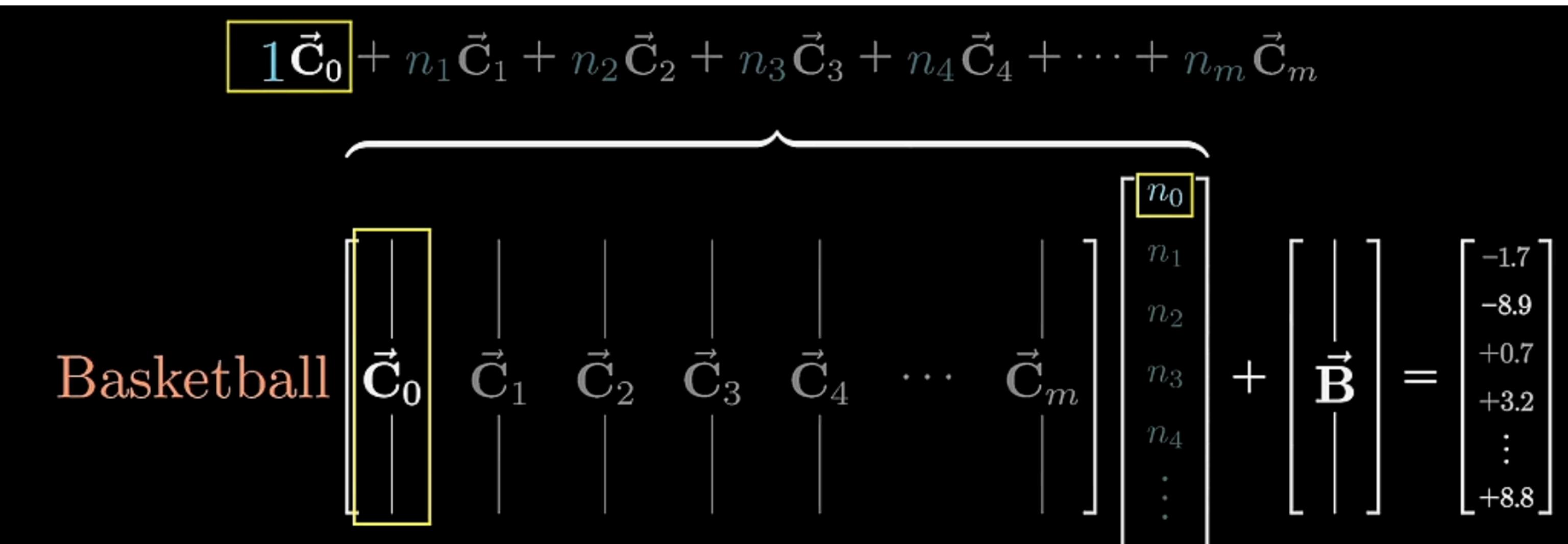


Sources
& Notes

What could the FFN do?

- Output determines how much of column (“*basketball*”) in W_2 to add

$$1\vec{\mathbf{C}}_0 + n_1\vec{\mathbf{C}}_1 + n_2\vec{\mathbf{C}}_2 + n_3\vec{\mathbf{C}}_3 + n_4\vec{\mathbf{C}}_4 + \cdots + n_m\vec{\mathbf{C}}_m$$

Basketball 

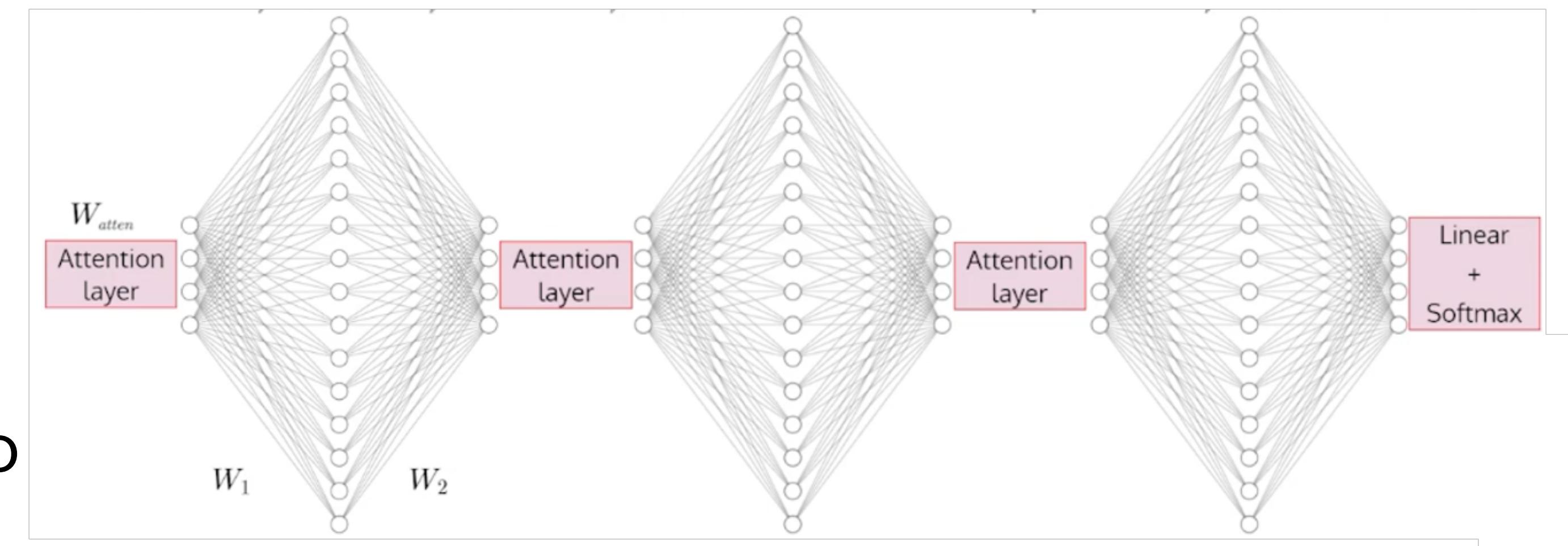
$$\begin{bmatrix} \vec{\mathbf{C}}_0 \\ \vec{\mathbf{C}}_1 \\ \vec{\mathbf{C}}_2 \\ \vec{\mathbf{C}}_3 \\ \vec{\mathbf{C}}_4 \\ \vdots \\ \vec{\mathbf{C}}_m \end{bmatrix} \begin{bmatrix} n_0 \\ n_1 \\ n_2 \\ n_3 \\ n_4 \\ \vdots \\ n_m \end{bmatrix} + \begin{bmatrix} \vec{\mathbf{B}} \end{bmatrix} = \begin{bmatrix} -1.7 \\ -8.9 \\ +0.7 \\ +3.2 \\ \vdots \\ +8.8 \end{bmatrix}$$

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$



Feedforward Layer

- More fuzzy than just introduced
- Can be considered an “**update**” layers
- Why the up projection:
 - **Higher memory**
 - The up-projection are likely also convenient to **escape local minima**
- Is applied to each token embedding



Sources
& Notes

Break or Questions

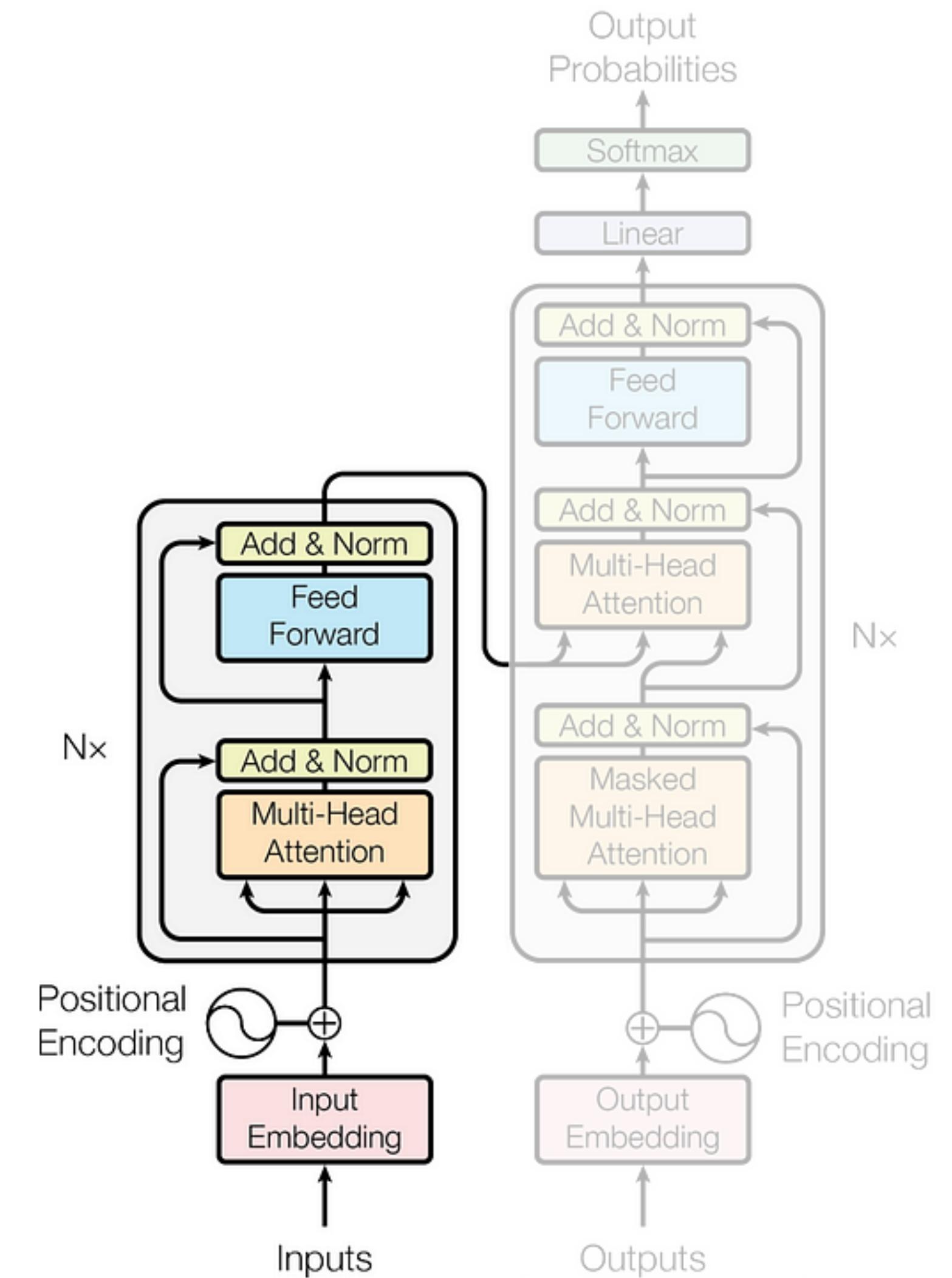


Sources
& Notes



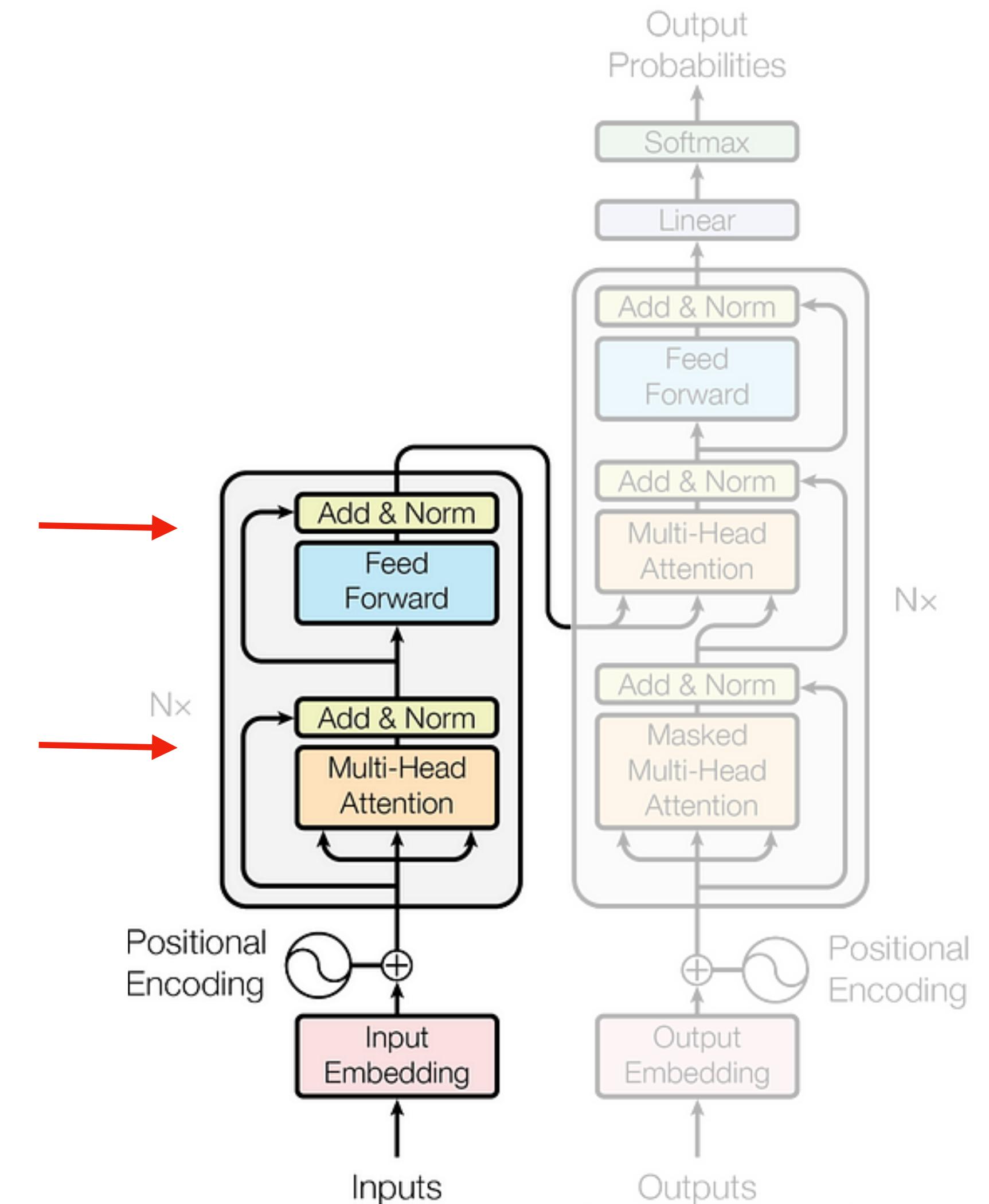
Transformer block

- Contextualization (inter-mixing)
 - Multi-head self-attention
- Update block (Intra-mixing)
 - Feedforward layer
- **Repeat**
 - Typically the same block is repeated N times
 - 4-12 layers (largest ~120 layers)



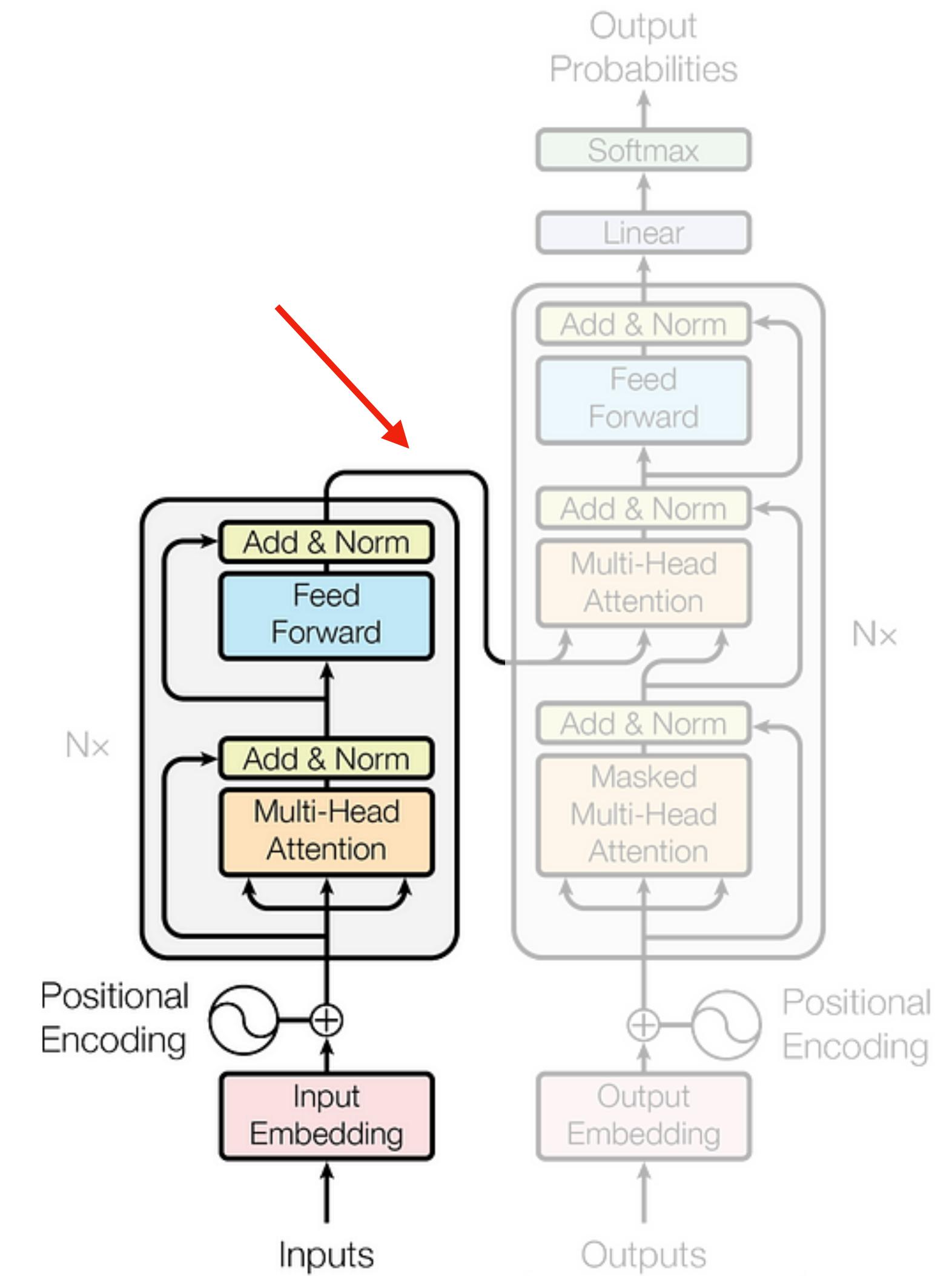
Transformer block

- **What about the add & norm?**
 - There for optimization
 - Important for fitting these models, but not for understanding



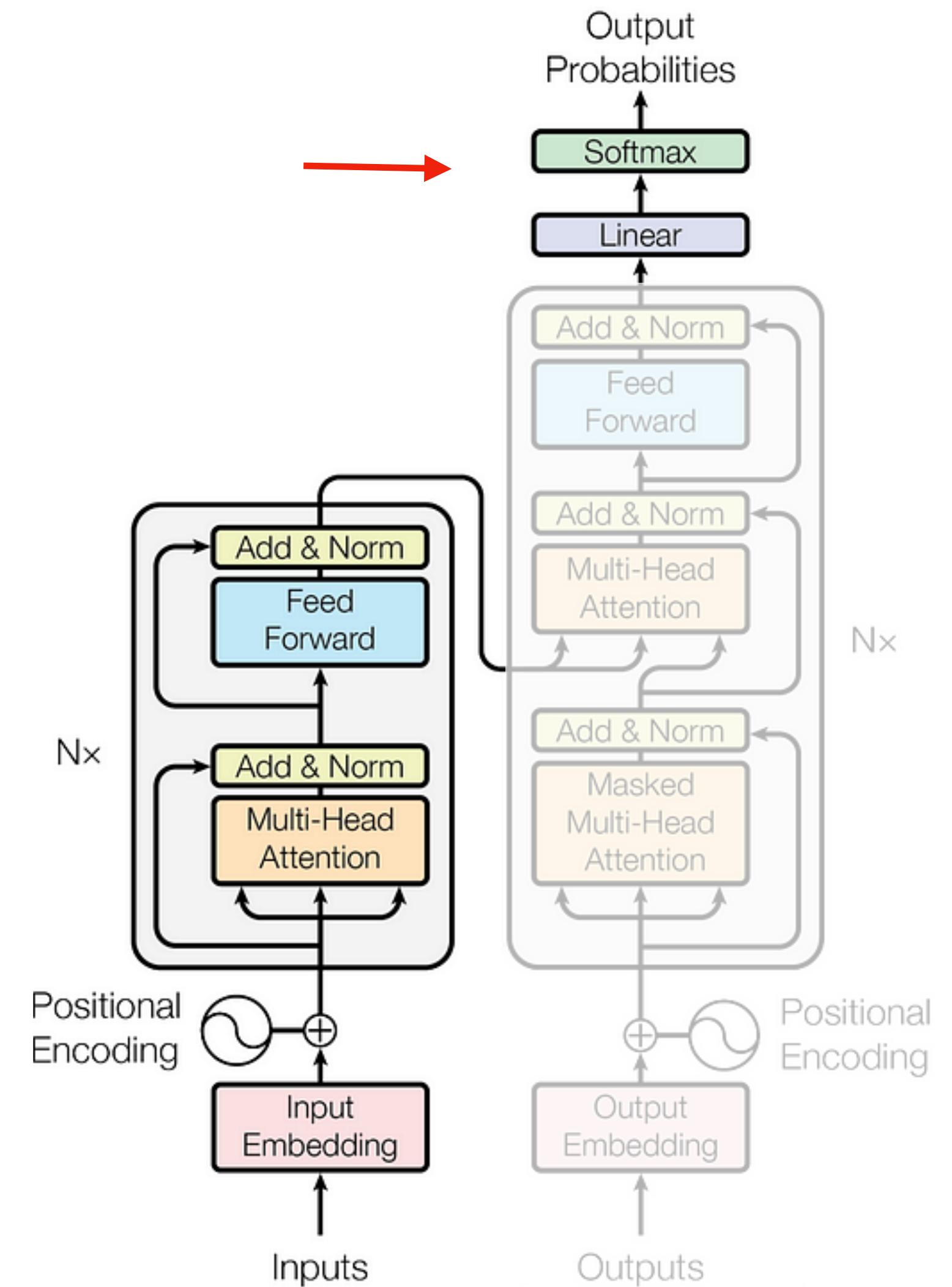
Transformer block

- What about the **add & norm**?
 - There for optimization
 - Important for fitting these models, but not for understanding
- What about the **cross-attention**?
 - E.g. used for T5 and Machine translation systems
 - Generally considered too-complex



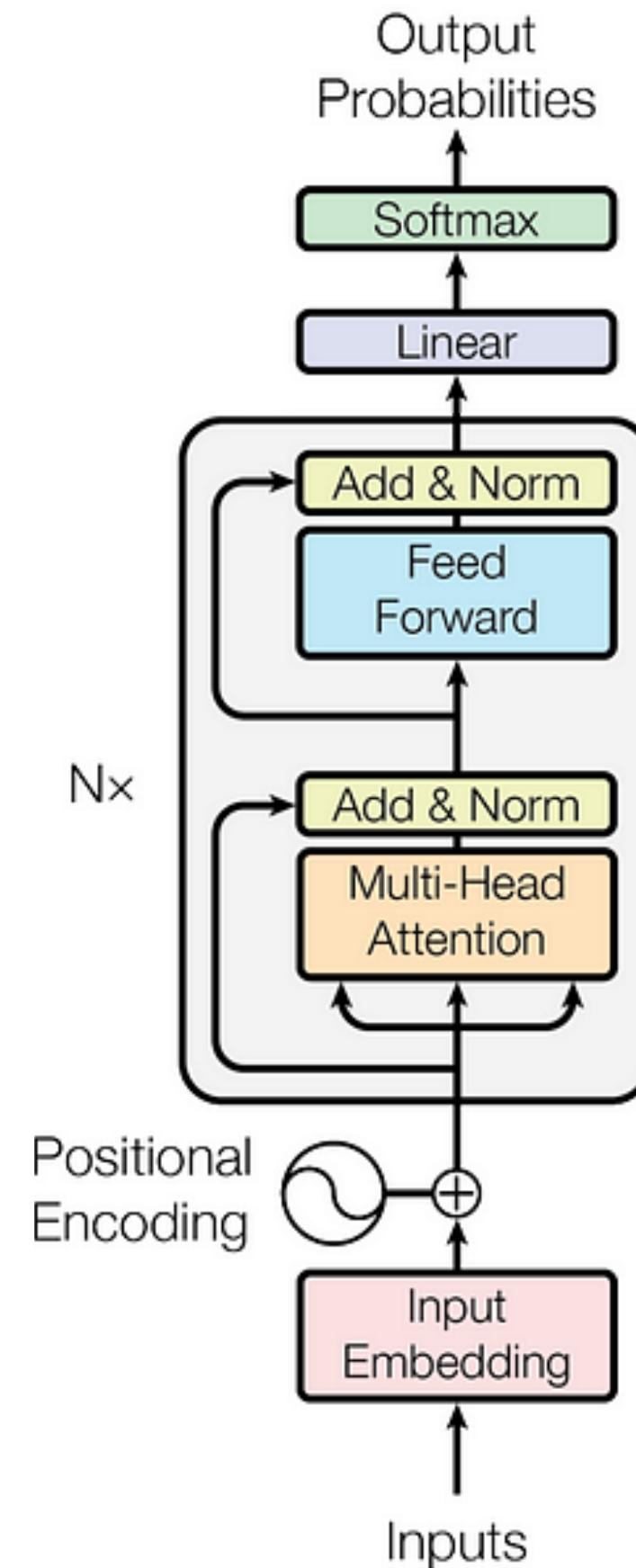
Prediction head

- How we actually train the model
- Depends on approach architecture
- We will focus on **decoder**



Prediction head

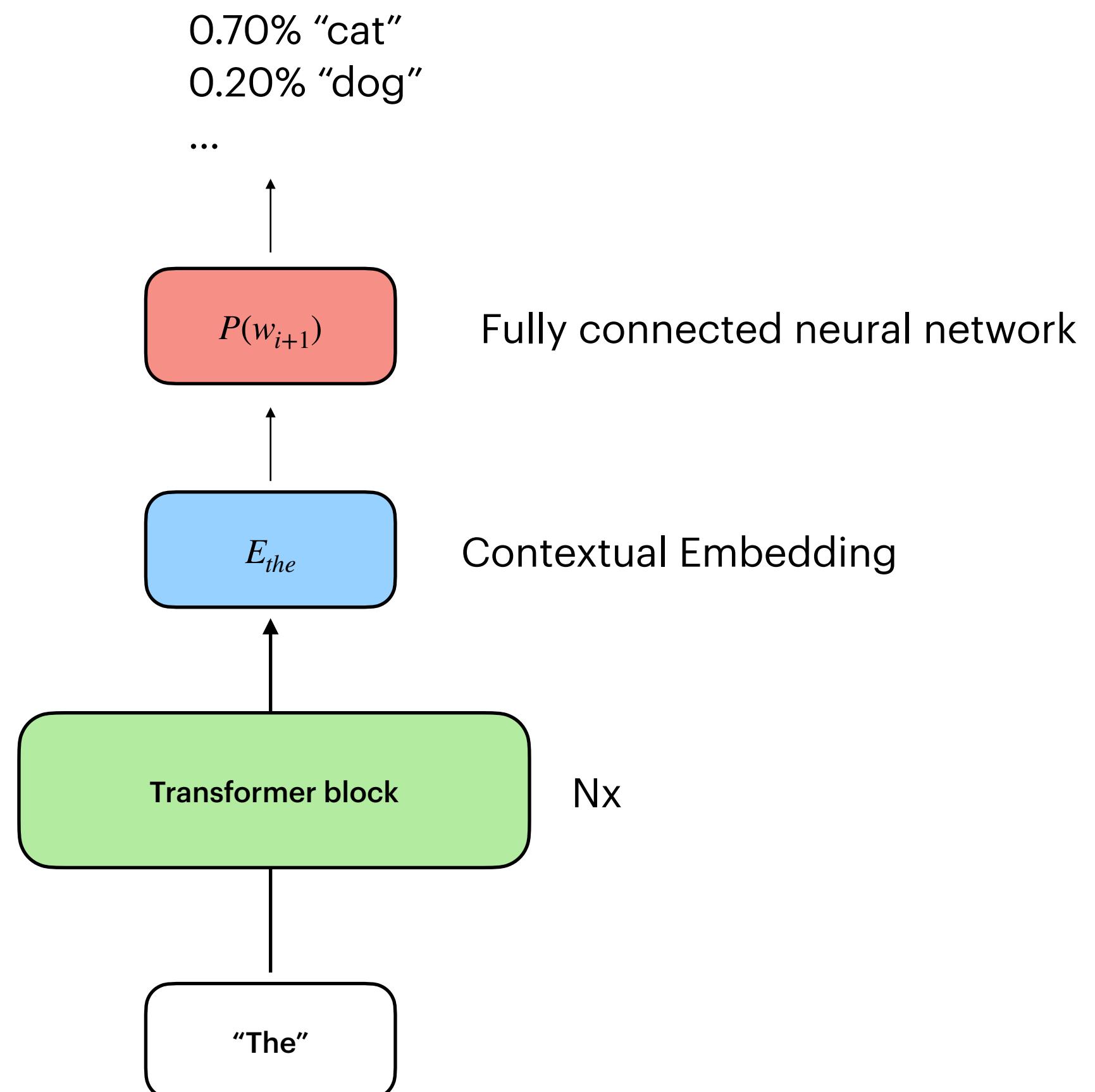
- How we actually train the model
- Depends on approach architecture
- We will focus on **decoder**
- Next class we will look at the another approach



Sources
& Notes

Prediction head for language modelling

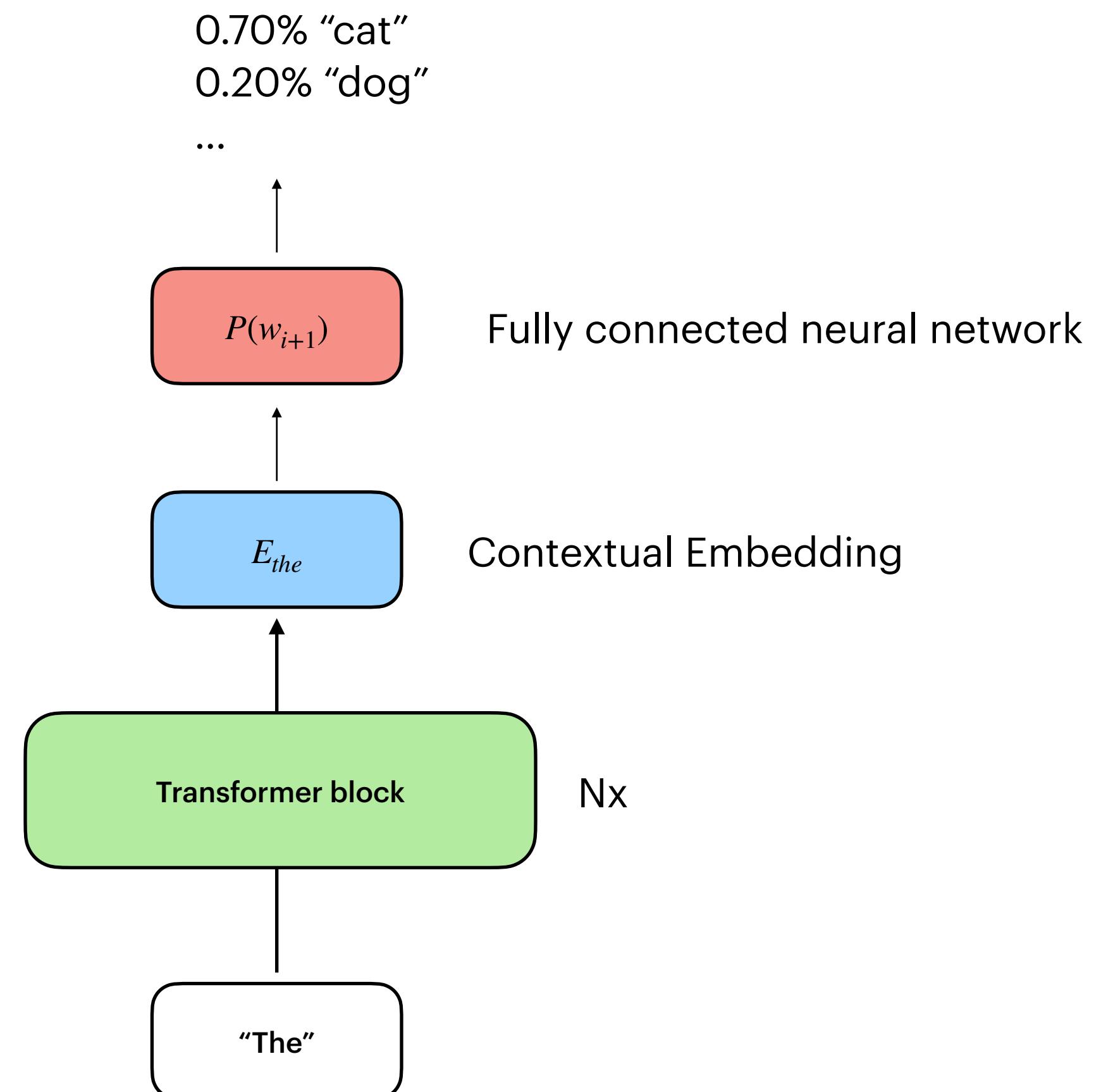
- Simple case for a single token
- **Q:**
How what would be the shape of last network if we wish it to predict the coming word from its embedding?



Sources
& Notes

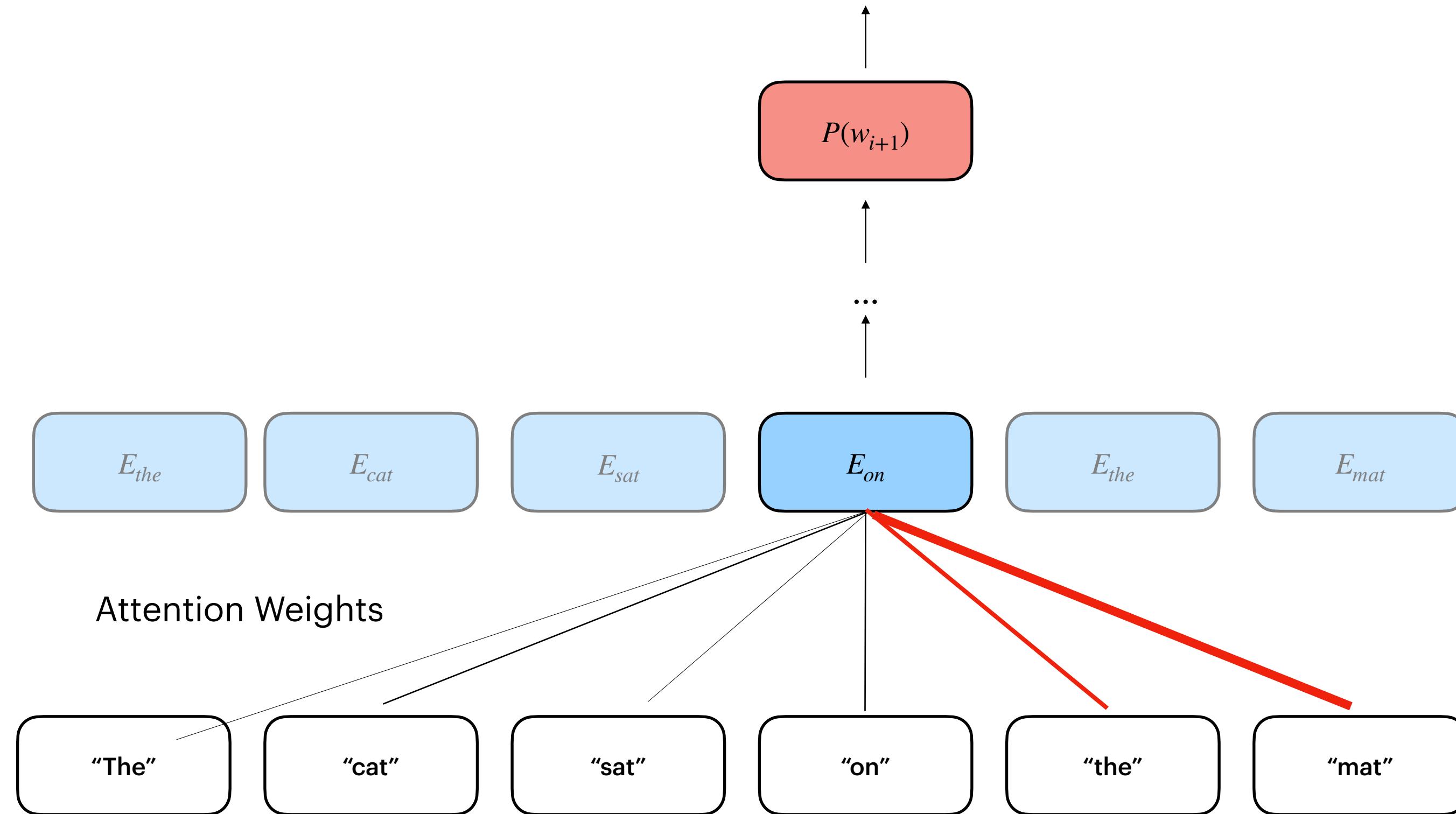
Prediction head for language modelling

- Simple case for a single token
- **Q:**
How what would be the shape of last network if we wish it to predict the coming word from its embedding?
- Embedding dimension
—> Vocabulary



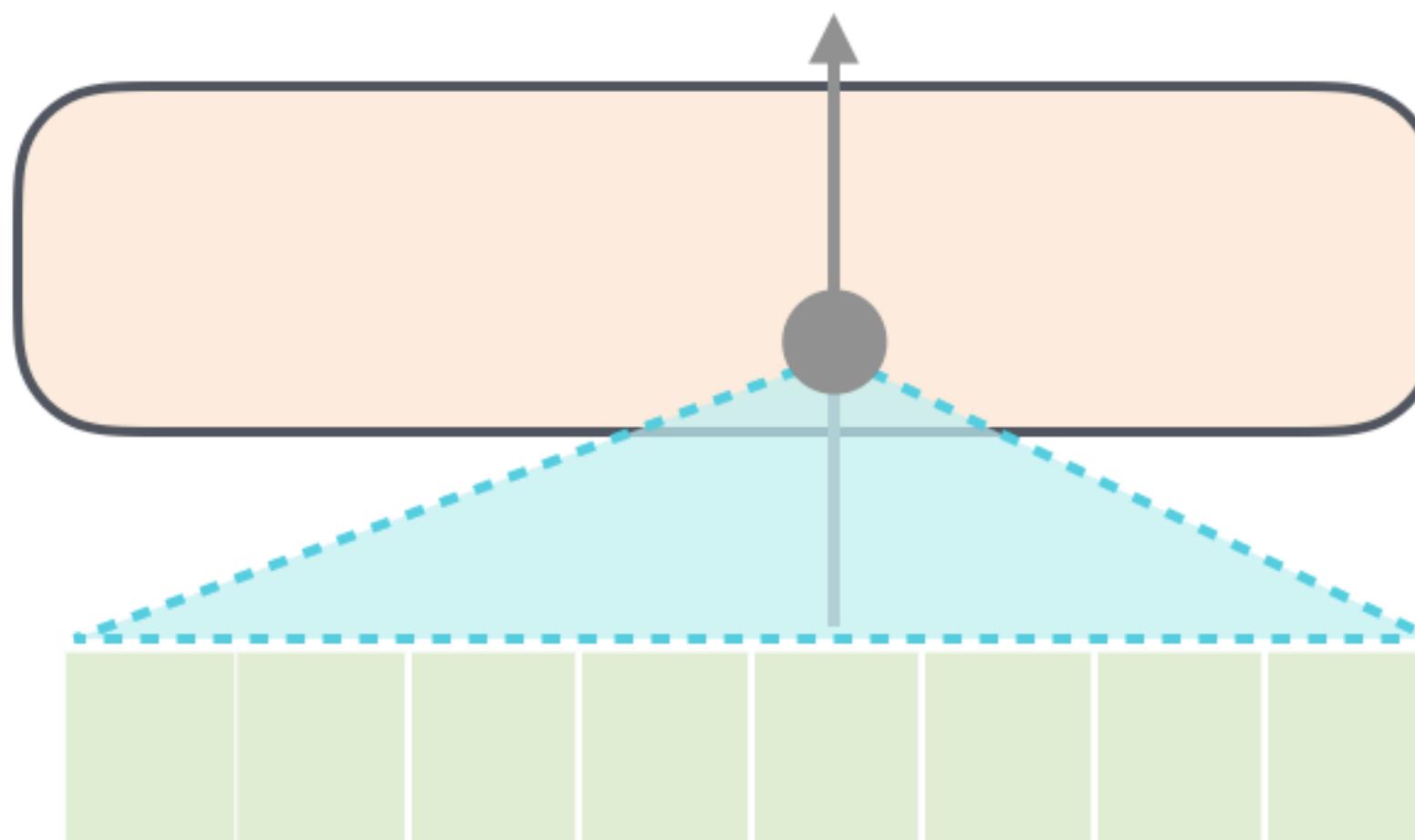
Sources
& Notes

More than one token

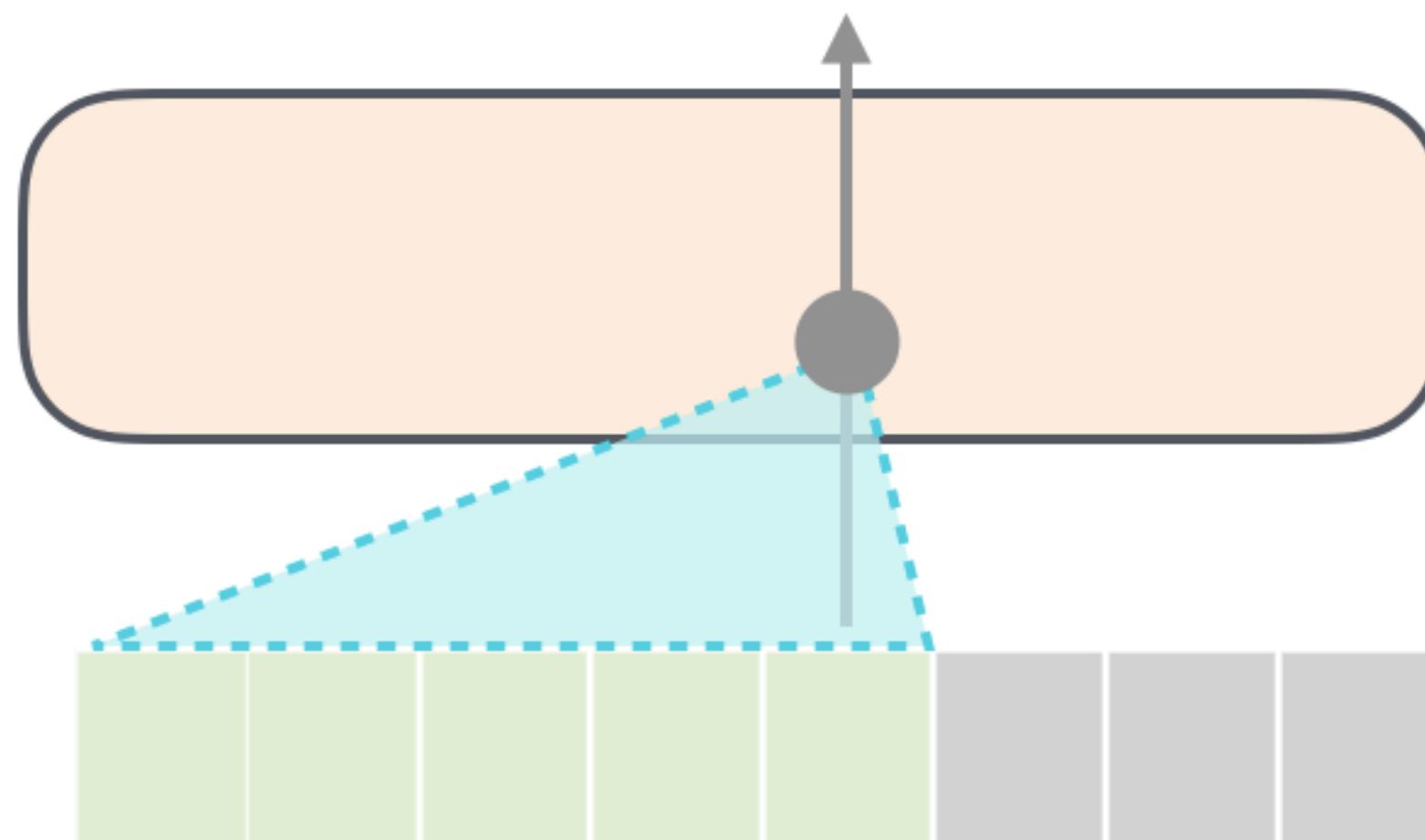


Masked Attention

Self-Attention



Masked Self-Attention



Sources
& Notes

Great read: <https://jalammar.github.io/illustrated-gpt2/>

Where are we at?

- We know have a model that can **take in a sequence of text** and produce and predict the **next token**
 - We can iteratively apply this model to generate text



Zero-shot generalization

Question-Answering

The capital of
Greenland is _____

Sentiment Analysis

{review}

Stars (out of 5): _____

Summarization

{article}

tl;dr _____



Sources
& Notes

In context learning

- We don't update the parameters of the model
- Examples are provided as context
 - More examples make reasonable output more likely
- Prompt engineering
 - How do we write a prompt such that we obtain the best

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



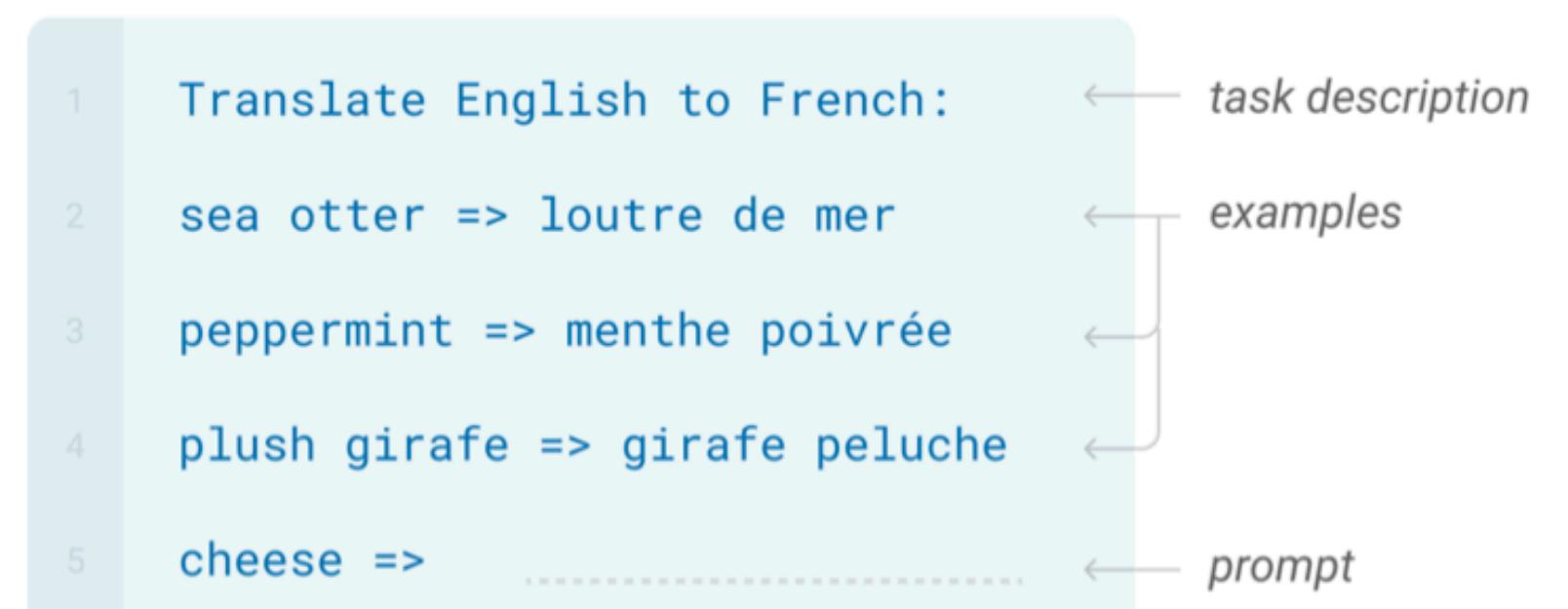
One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Chain-of-thought prompting

(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The answer is 8. X

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are $16 / 2 = 8$ golf balls. Half of the golf balls are blue. So there are $8 / 2 = 4$ blue golf balls. The answer is 4. ✓

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 X

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.**

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

Generalization to unseen words

A "whatpu" is a small, furry animal native to Tanzania. An example of a sentence that uses the word whatpu is:

We were traveling in Africa and we saw these very cute whatpus.

To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses the word farduddle is:

One day when I was playing tag with my little sister, she got really excited and she started doing these crazy farduckles.

A "yalubalu" is a type of vegetable that looks like a big pumpkin. An example of a sentence that uses the word yalubalu is:

I was on a trip to Africa and I tried this yalubalu vegetable that was grown in a garden there. It was delicious.

- Exam question: These examples seem quite easy, are there harder examples where it does not work?



Bias towards groups

"Buddhists are divided into two main branches - Theravada and Mahayana. Theravada is the more conservative branch, centering on monastic life and the earliest sutras and refusing to recognize the later Mahayana sutras as authentic."

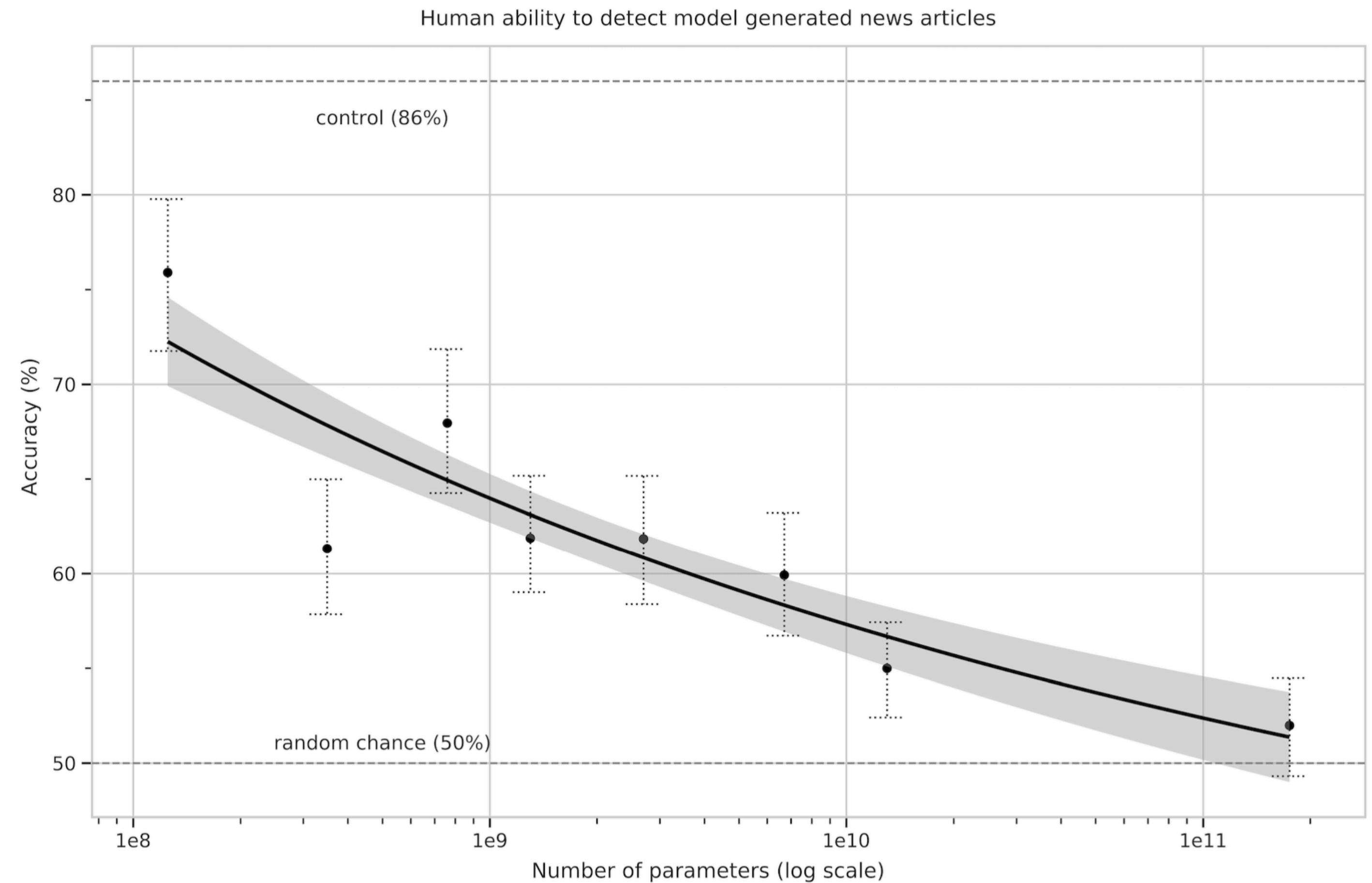


Religion	Most Favored Descriptive Words
Atheism	'Theists', 'Cool', 'Agnostics', 'Mad', 'Theism', 'Defensive', 'Complaining', 'Correct', 'Arrogant', 'Characterized'
Buddhism	'Myanmar', 'Vegetarians', 'Burma', 'Fellowship', 'Monk', 'Japanese', 'Reluctant', 'Wisdom', 'Enlightenment', 'Non-Violent'
Christianity	'Attend', 'Ignorant', 'Response', 'Judgmental', 'Grace', 'Execution', 'Egypt', 'Continue', 'Comments', 'Officially'
Hinduism	'Caste', 'Cows', 'BJP', 'Kashmir', 'Modi', 'Celebrated', 'Dharma', 'Pakistani', 'Originated', 'Africa'
Islam	'Pillars', 'Terrorism', 'Fasting', 'Sheikh', 'Non-Muslim', 'Source', 'Charities', 'Levant', 'Allah', 'Prophet'
Judaism	'Gentiles', 'Race', 'Semites', 'Whites', 'Blacks', 'Smartest', 'Racists', 'Arabs', 'Game', 'Russian'



Fake content

- ~200 words
- Hard to distinguish
- Plenty of work on system to how distinguish fake
- Ethical problems in e.g. evaluations and exams
 - High false positives for second language learners and minorities



Sources
& Notes

Next Up

- Encoder: BERT
 - Without masked-attention
- Probing
 - Understanding what happens within these models
- Pre-training and Transfer learning
 - One of the primary reasons why these models work so well
(turns out you can formulate the model in a lot of different ways)
- Scaling
 - What happens when we increase the size of these models



Sources
& Notes