

Transformer Basics

Danna Gurari

University of Colorado Boulder

Spring 2025



Review

- Last week:
 - Motivation: machine neural translation for long sentences
 - Encoder
 - Decoder: attention
 - Performance evaluation
 - Final project: ways to find a partner
- Assignments (Canvas):
 - Lab assignment 2 due Tuesday
- Questions?

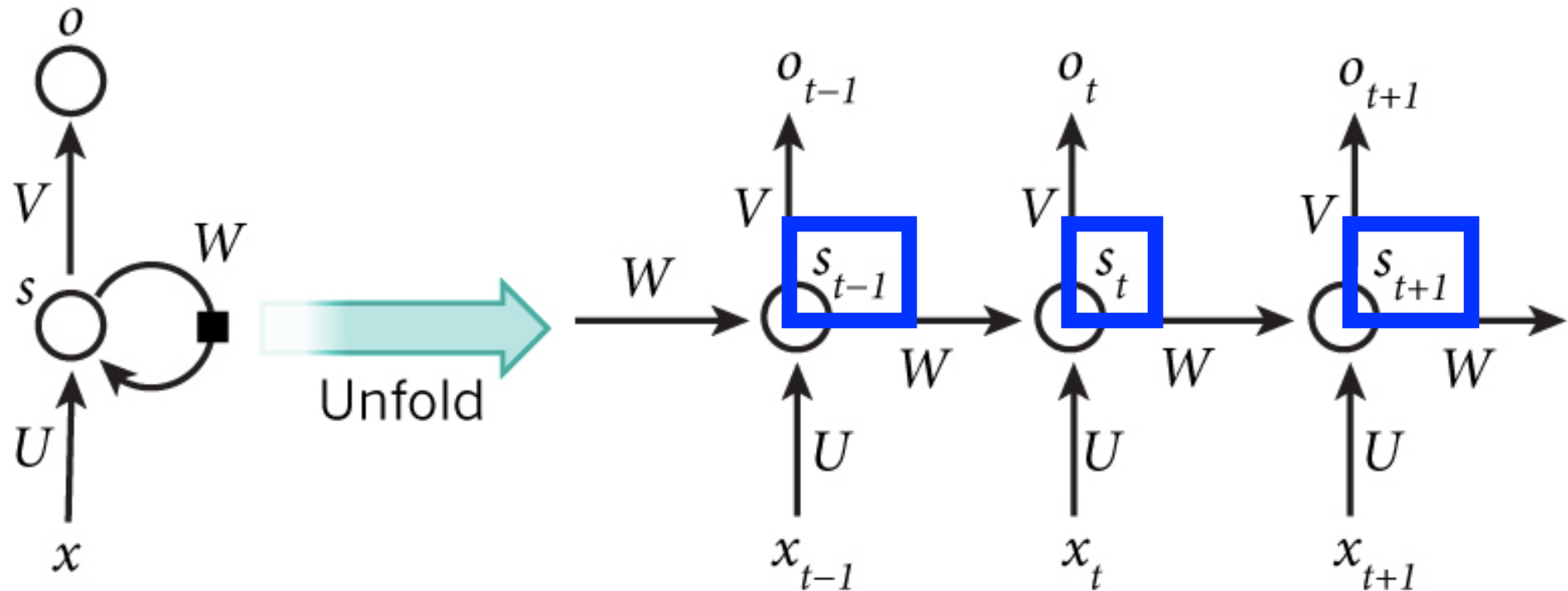
Today's Topics

- Transformer overview
- Self-attention
- Common transformer ingredients
- Pioneering transformer: machine translation
- Programming tutorial

Today's Topics

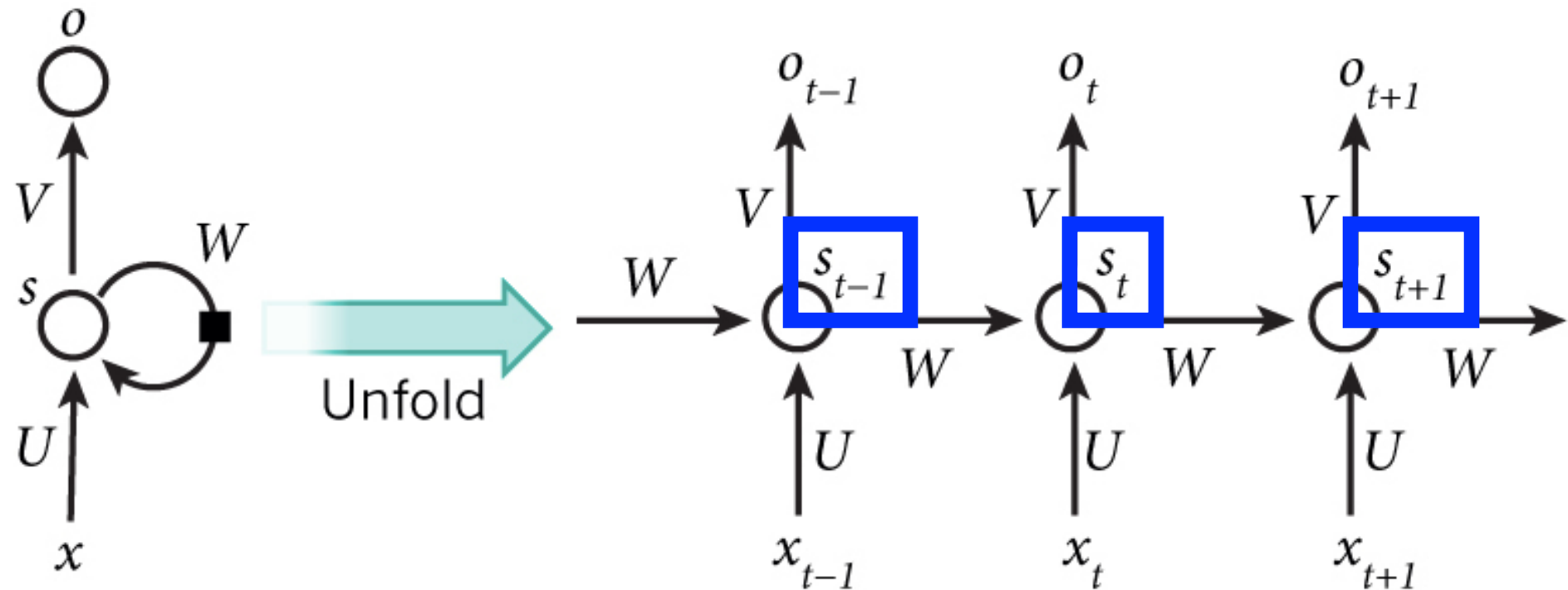
- Transformer overview
- Self-attention
- Common transformer ingredients
- Pioneering transformer: machine translation
- Programming tutorial

Goal: Model Sequential Data (Recall RNN)



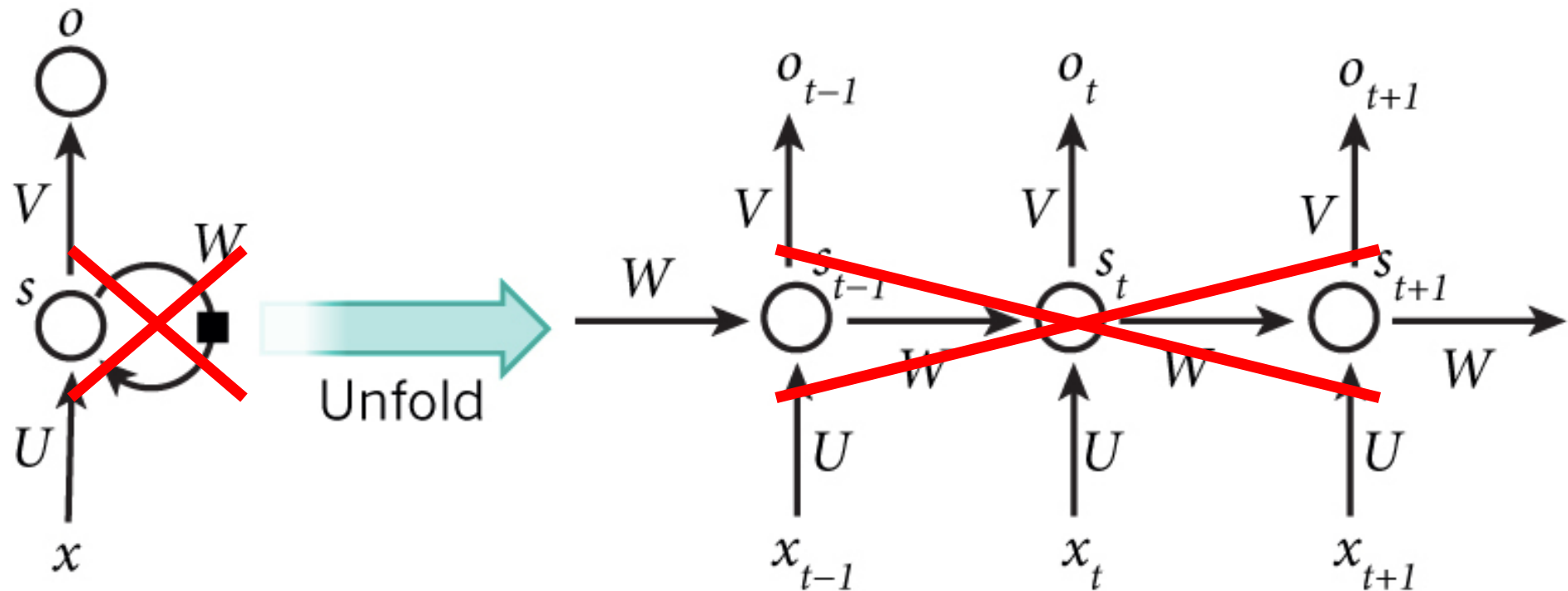
Each hidden state is a function of the previous hidden state

Problem: RNNs Use Sequential Computation



RNNs struggle to carry information through hidden states across many time steps and train/testing is slow

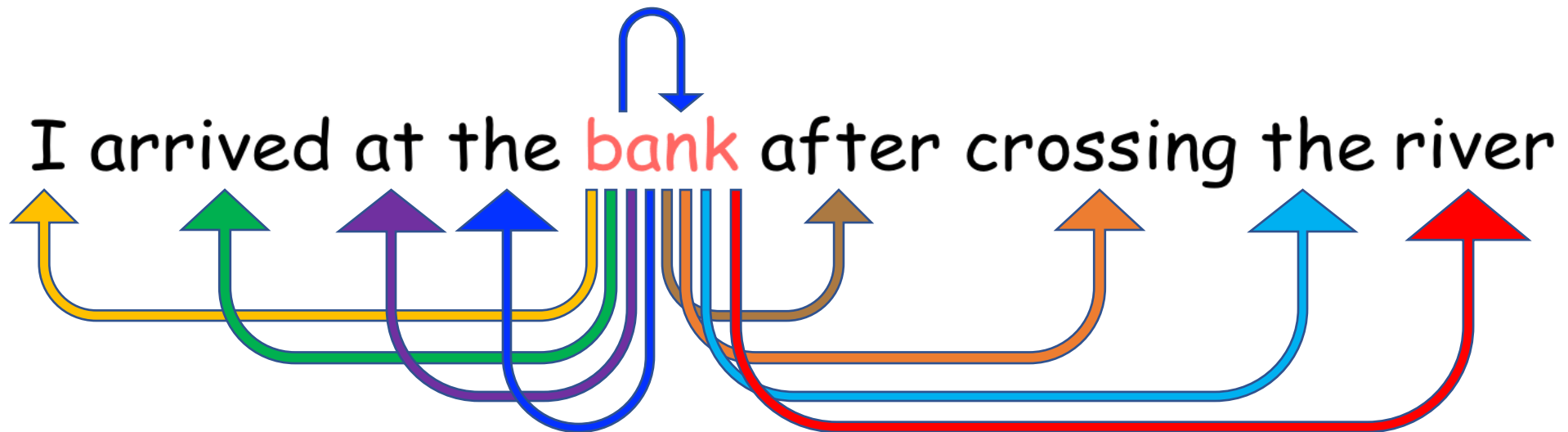
Idea: Model Sequential Data Without Recurrence



Replace sequential hidden states for capturing knowledge of other inputs with a new representation of each input that shows its relationship to all other inputs (i.e., self-attention)

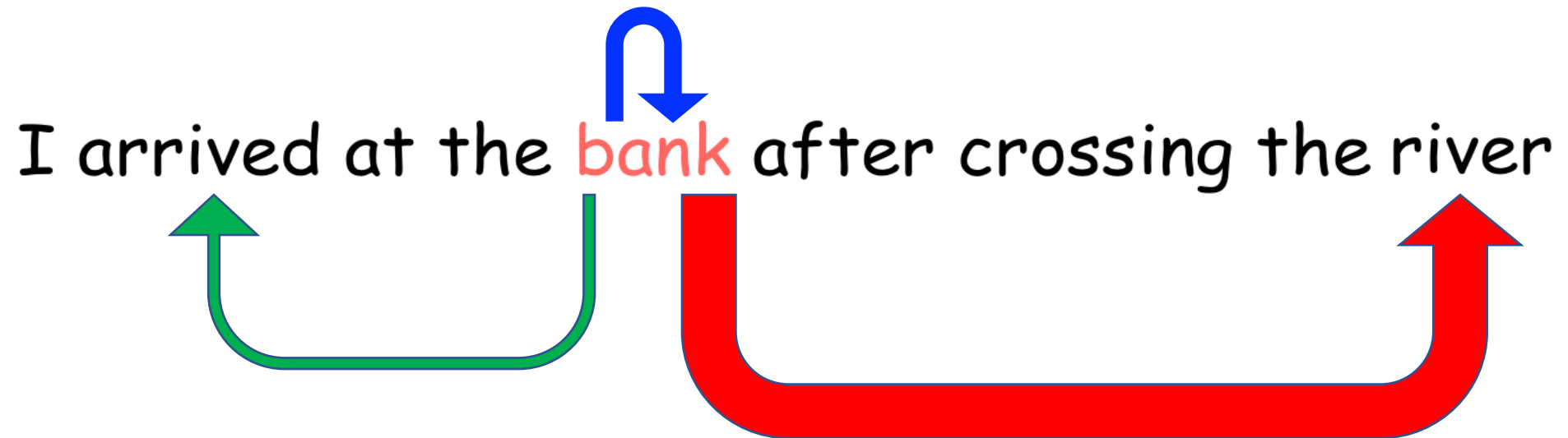
Transformer Key Idea: Self-Attention

New representation of each **token** in a sequence showing its relationship to all tokens; e.g.,



Transformer Key Idea: Self-Attention

New representation of each **token** in a sequence showing its relationship to all tokens; e.g.,



Arrow thickness is indicative of attention weight

Transformer Key Idea: Self-Attention

New representation of each **token** in a sequence showing its relationship to all tokens; e.g.,

I arrived at the **bank** after crossing the river



A large attention score means the other word will strongly inform the new representation of the word

Transformer Intuition

What does **bank** mean in this sentence?

I arrived at the **bank** after crossing the ...

Transformer Intuition

What does **bank** mean in this sentence?

- new word representation disambiguates meaning by identifying other relevant words (e.g., high attention score with “river”)

I arrived at the **bank** after crossing the river

vs

I arrived at the **bank** after crossing the street



Transformer vs RNN (Intuition)

I arrived at the **bank** after crossing the ...

...street? ...river?

What does **bank** mean in this sentence? *Meaning depends on other input words*

Transformer vs RNN (Intuition)

I arrived at the **bank** after crossing the ...

...street? ...river?

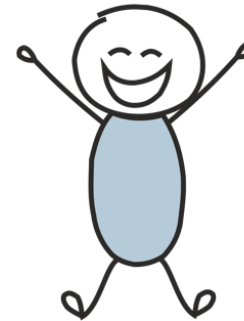
What does **bank** mean in this sentence? Meaning depends on other input words



I've no idea: let's wait until I read the end

RNNs

$O(N)$ steps to process a sentence with length N



I don't need to wait - I see all words at once!

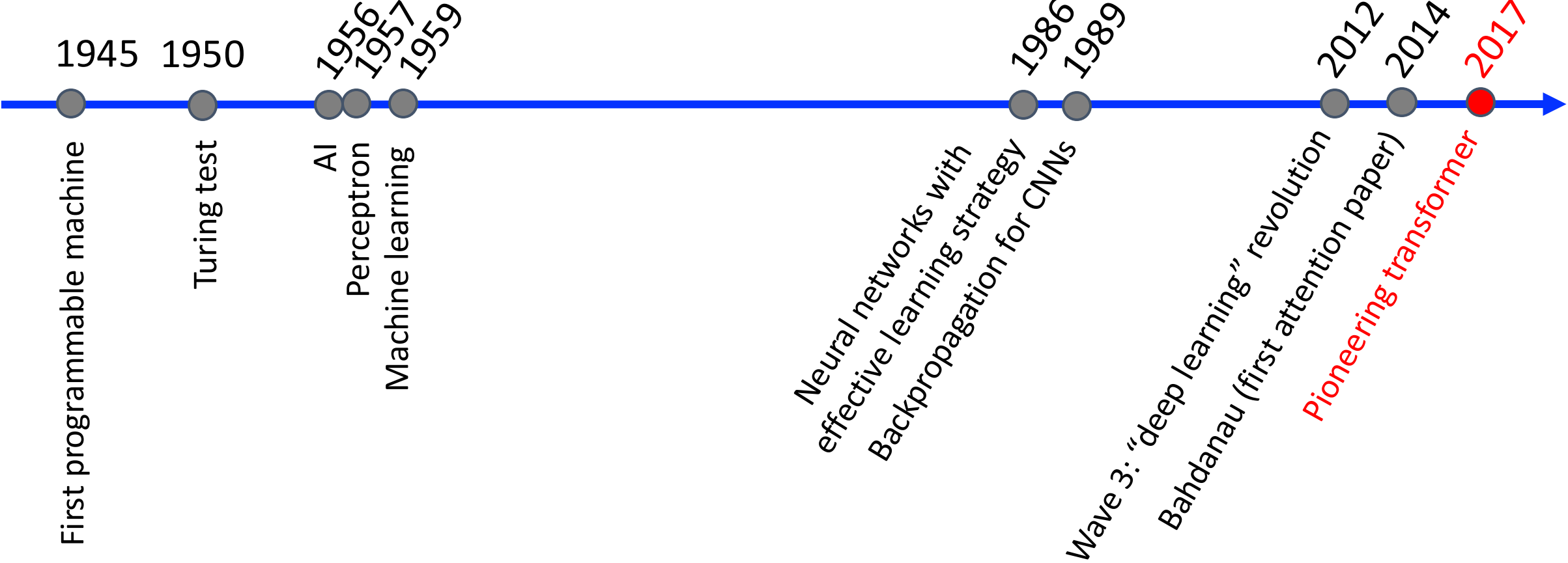
Transformer

Constant number of steps to process any sentence

Transformer: A Suggested Definition

“Any architecture designed to process a connected set of units—such as the tokens in a sequence or the pixels in an image—where the only interaction between units is through self-attention.”

Historical Context: Pioneering Transformer

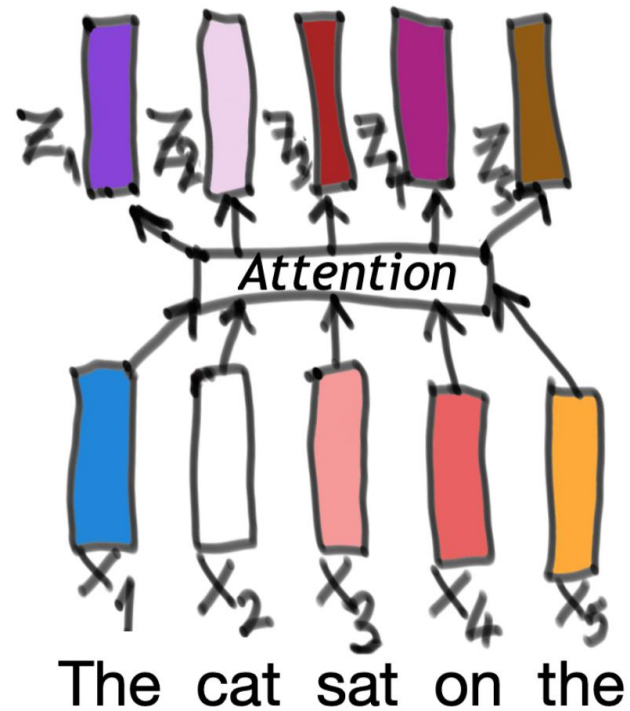


Today's Topics

- Transformer overview
- **Self-attention**
- Common transformer ingredients
- Pioneering transformer: machine translation
- Programming tutorial

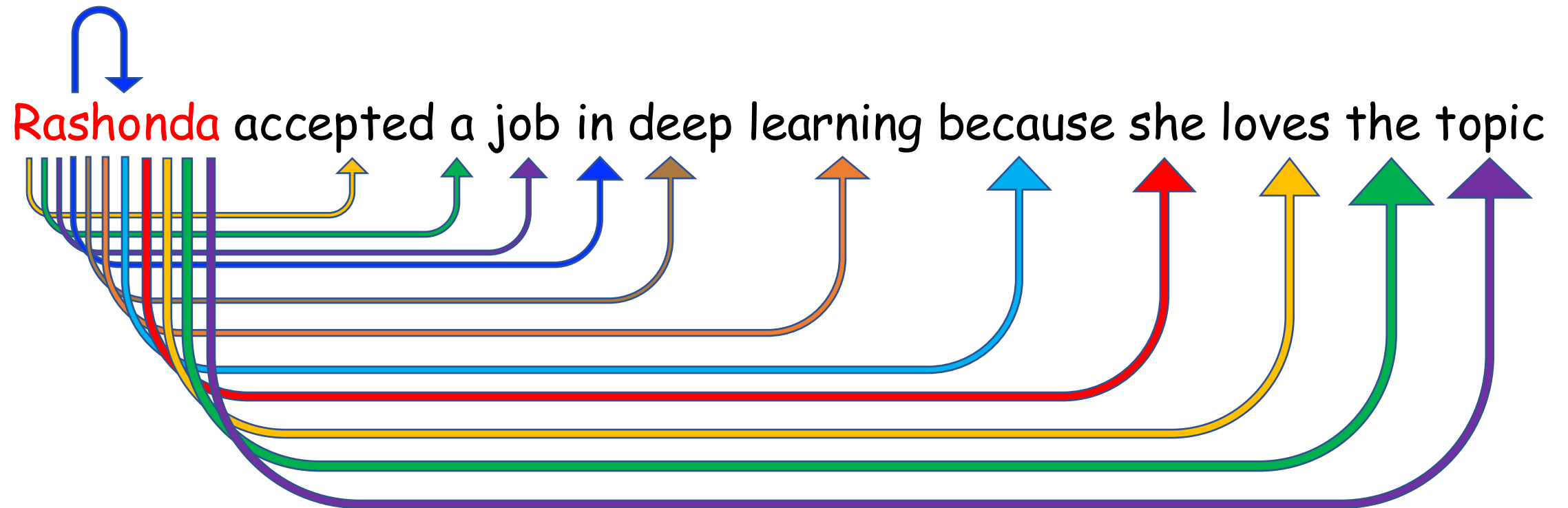
Self-Attention: Outcome

New representation of each **token** in a sequence showing its relationship to all tokens



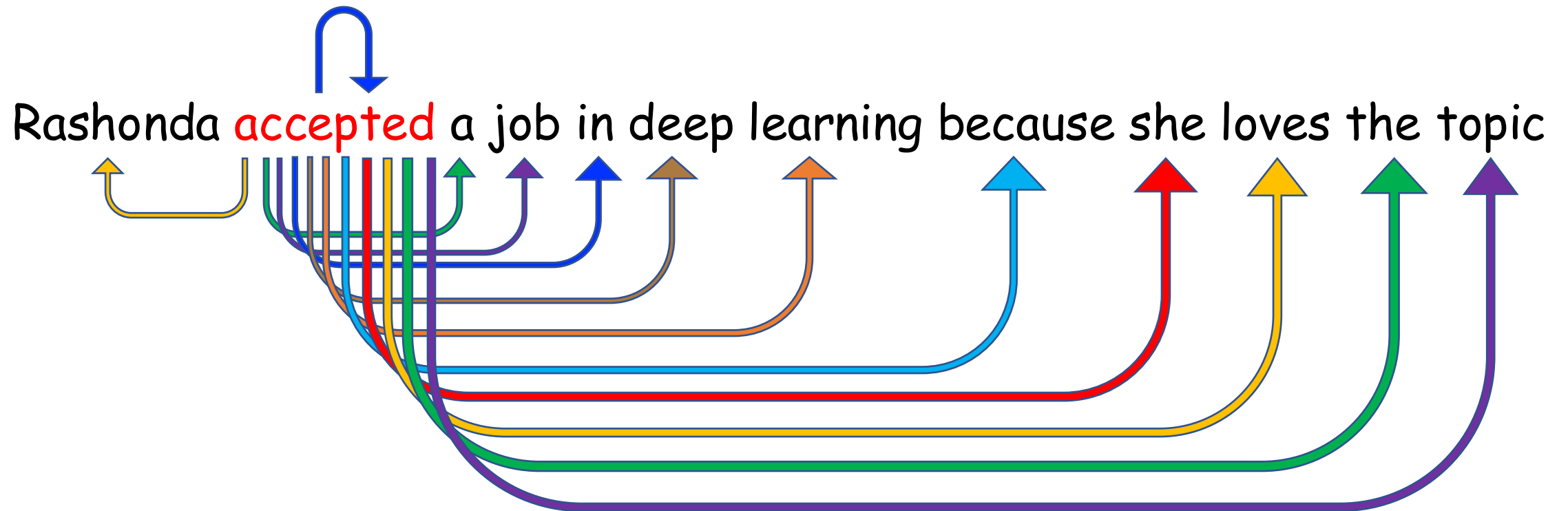
Self-Attention: Outcome

New representation of each **token** in a sequence showing its relationship to all tokens; e.g.,



Self-Attention: Outcome

New representation of each **token** in a sequence showing its relationship to all tokens; e.g.,



Self-Attention: Outcome

New representation of each **token** in a sequence showing its relationship to all tokens; e.g.,

Rashonda accepted **a** job in deep learning because she loves the topic



And so on for remaining words...

Self-Attention: Disambiguates Word Meanings

New representation of each **token** in a sequence showing its relationship to all tokens; e.g.,

Rashonda accepted a job in deep learning because **she** loves the topic



A better representation of "she" would
encode information about "Rashonda"

Self-Attention: Disambiguates Word Meanings

New representation of each **token** in a sequence showing its relationship to all tokens; e.g.,

I arrived at the bank across the river

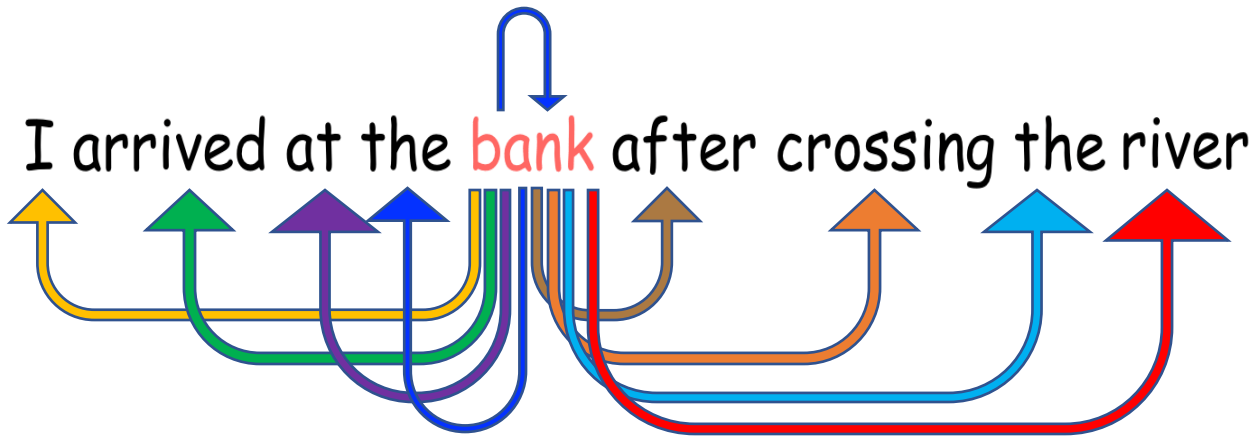


Recall: a better representation of “bank”
would encode information about “river”

Self-Attention vs General Attention

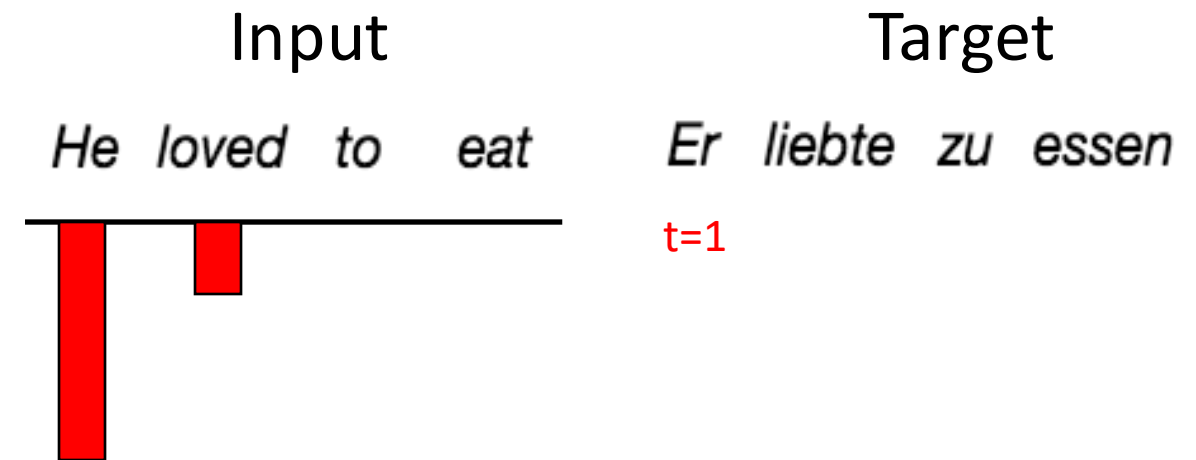
Self-attention

Relates tokens from the same source

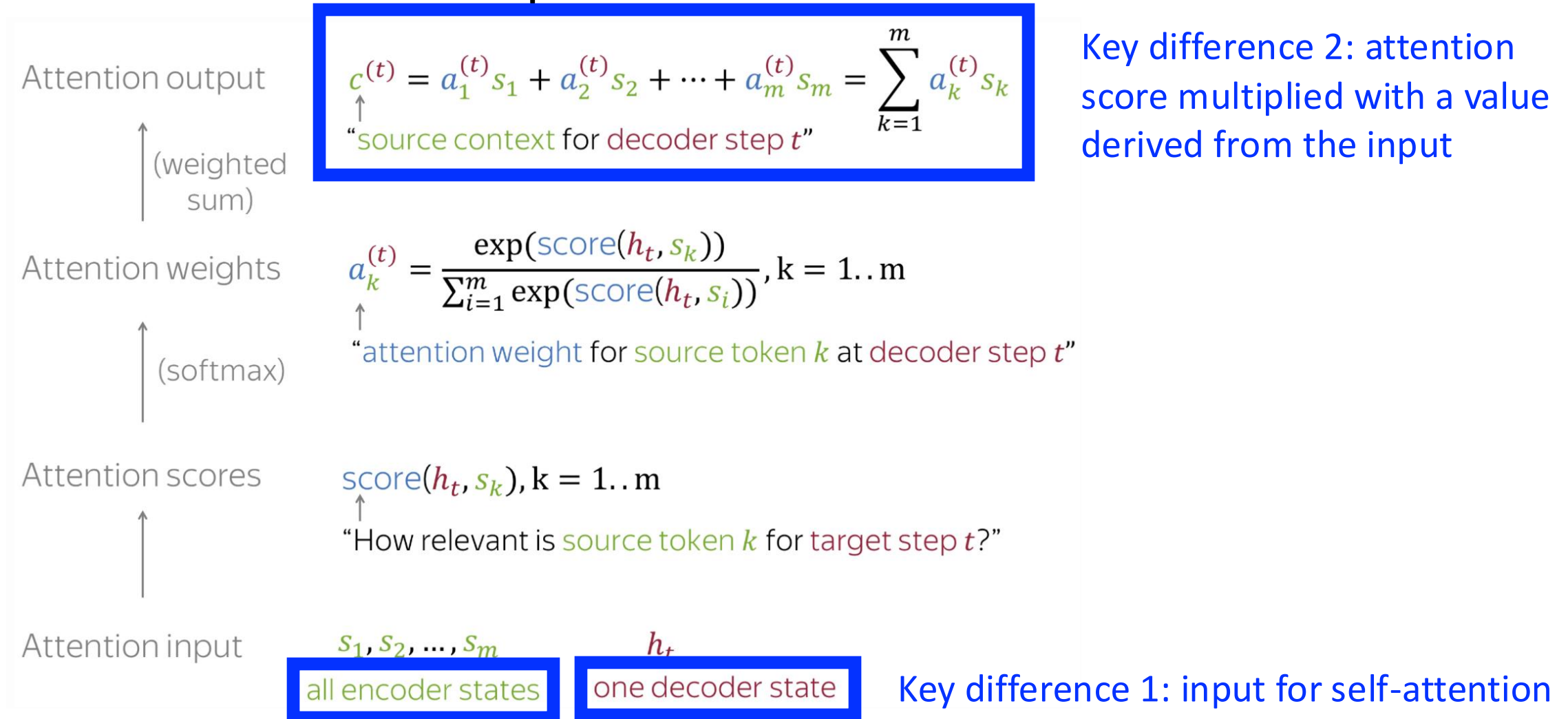


General attention

Relates tokens from different sources

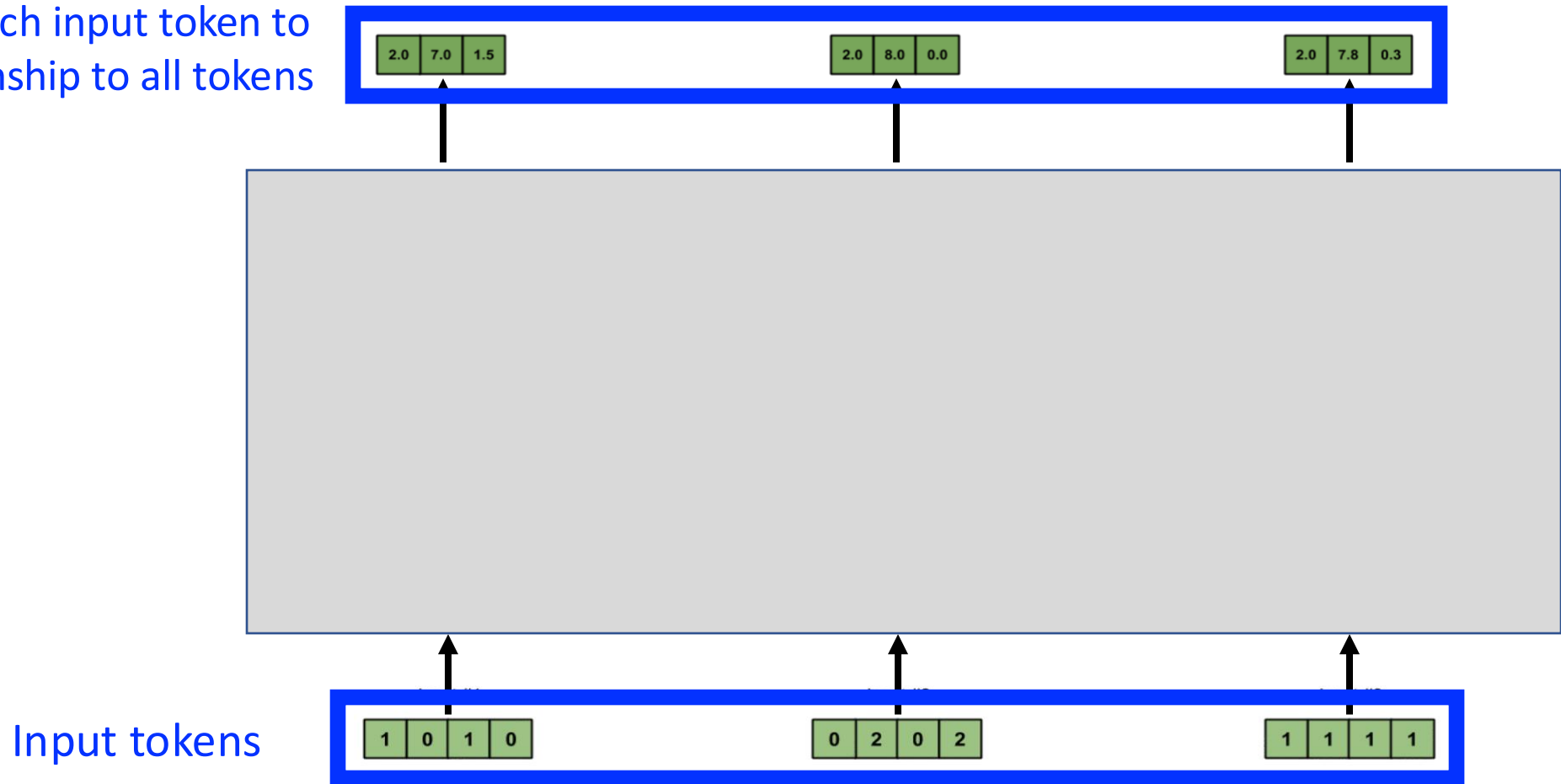


Computing Self-Attention: Similar Approach to How We Compute General Attention

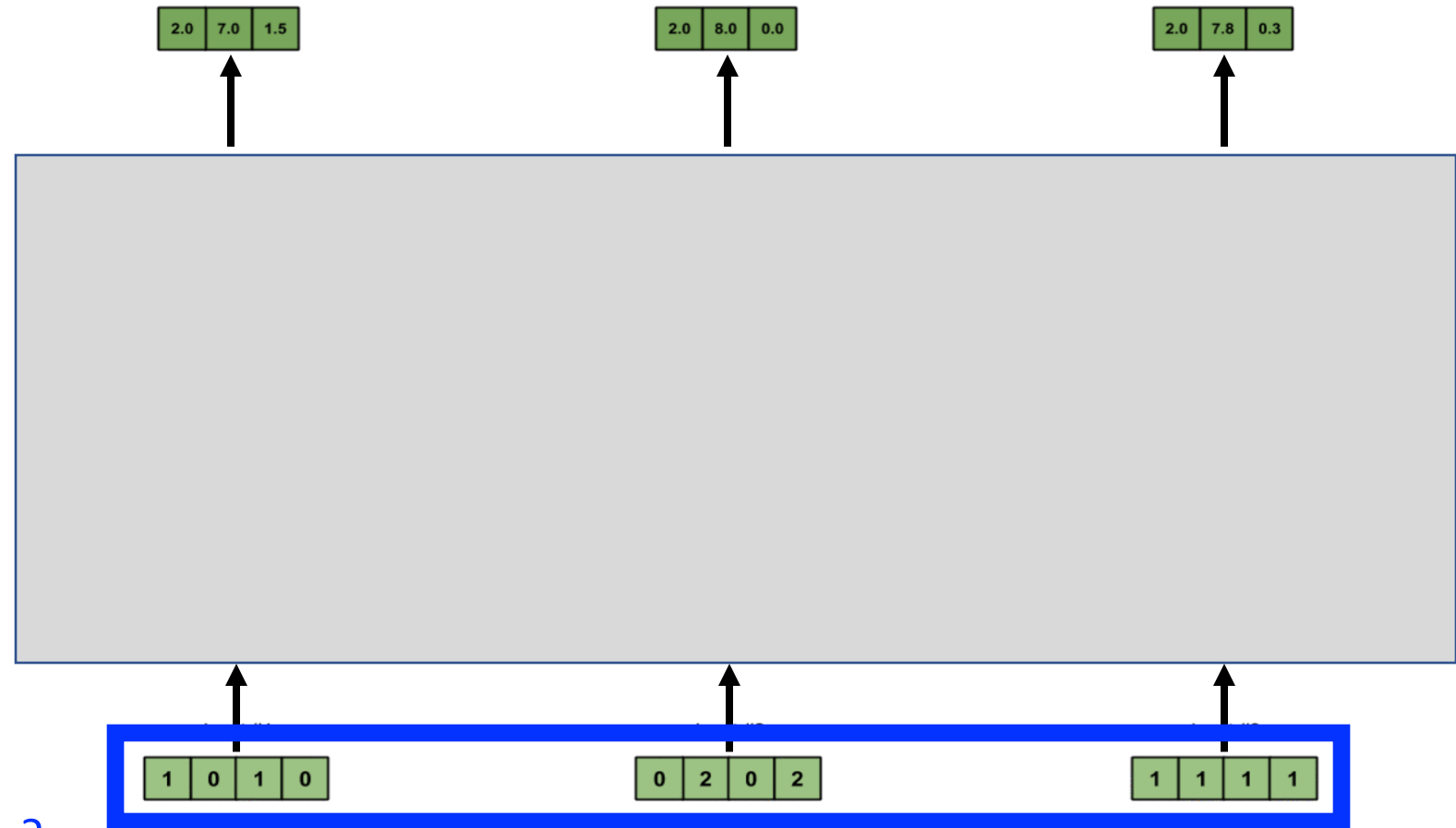


Computing Self-Attention: Example

New representation of each input token to reflect each one's relationship to all tokens



Computing Self-Attention: Example

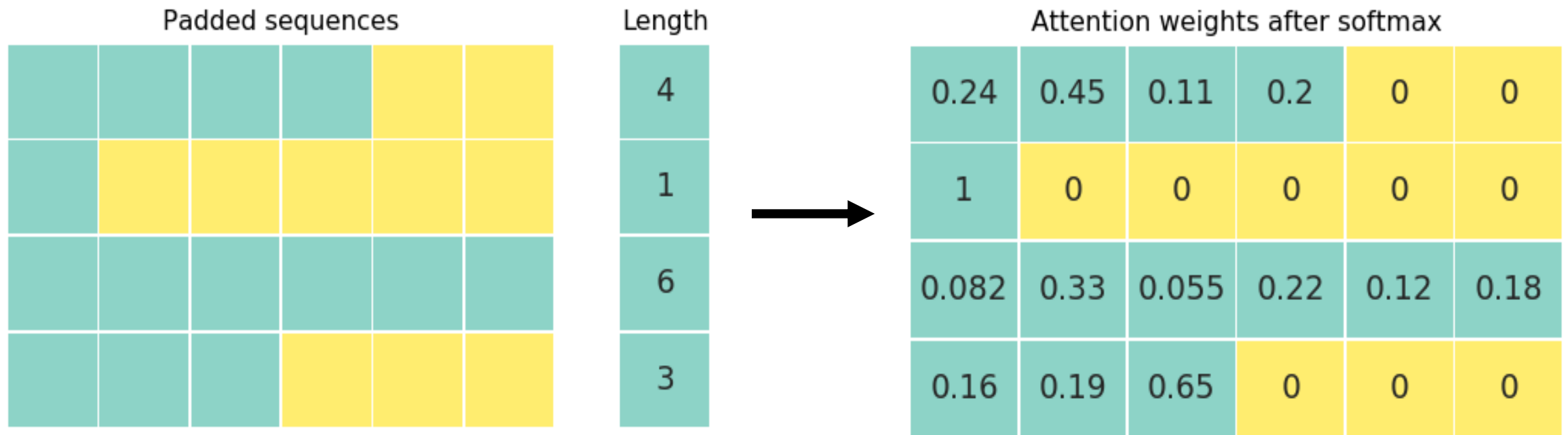


- How many input tokens are there?
- What is each token's dimensionality?
- How to support arbitrary length inputs?

* Input length is a hyperparameter: pad shorter sequences with zeros and truncate longer sequences

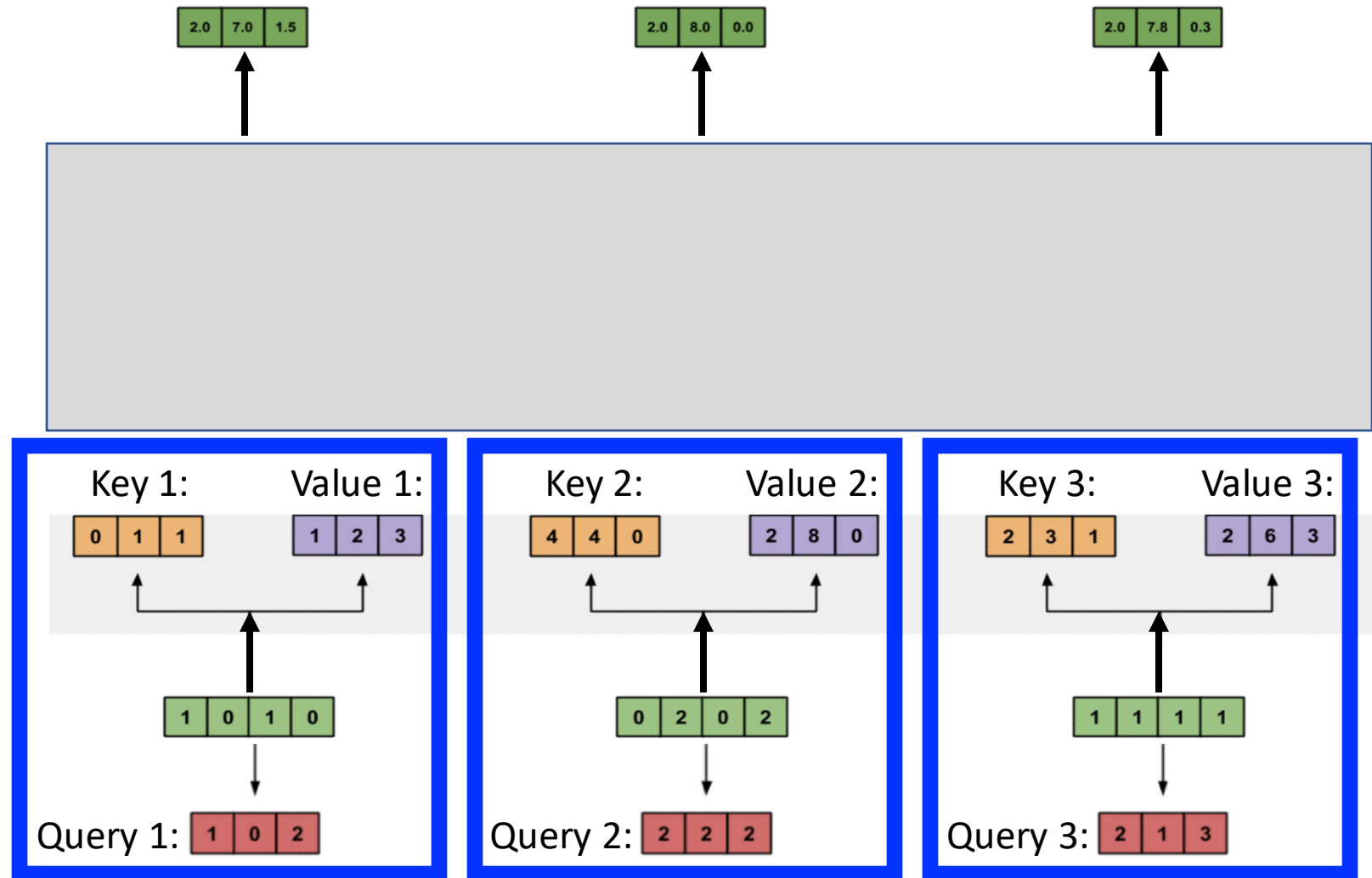
Input Length: Implementation

- **[PAD]** tokens with attention set to 0 enables variable input length



Computing Self-Attention: Example

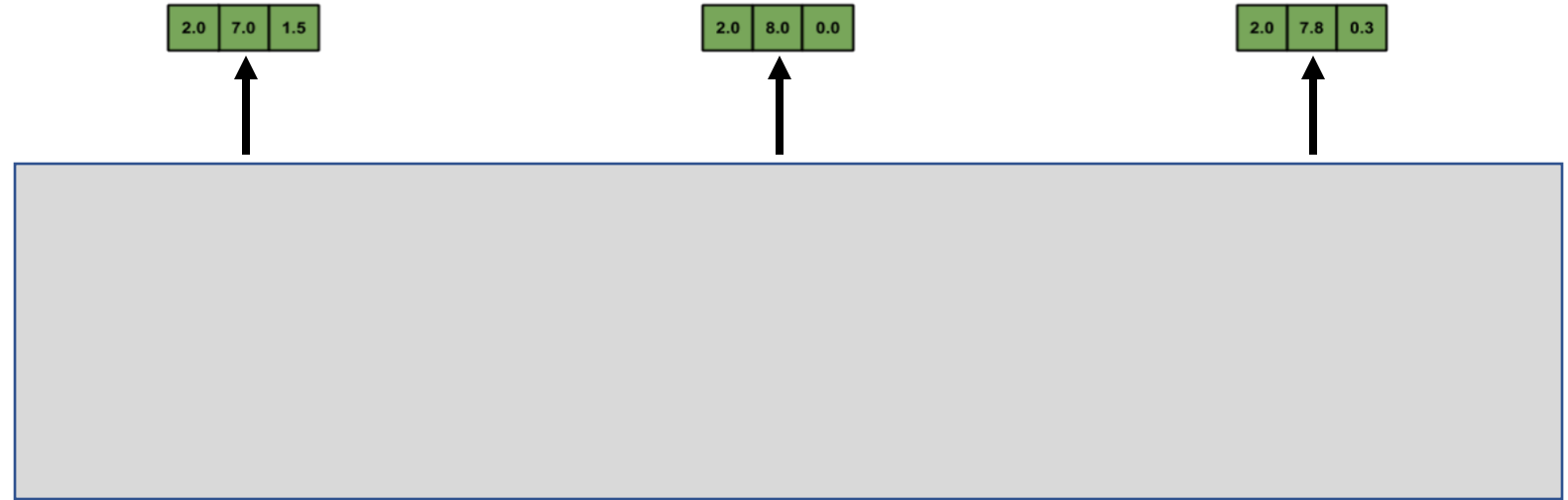
Three vectors are derived for each **input** by multiplying with three weight matrices (learned during training): **query**, **key**, and **value**



Computing Self-Attention: Example

e.g., **key** weights

```
[0, 0, 1]  
[1, 1, 0]  
[0, 1, 0]  
[1, 1, 0]
```



$$\begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$$

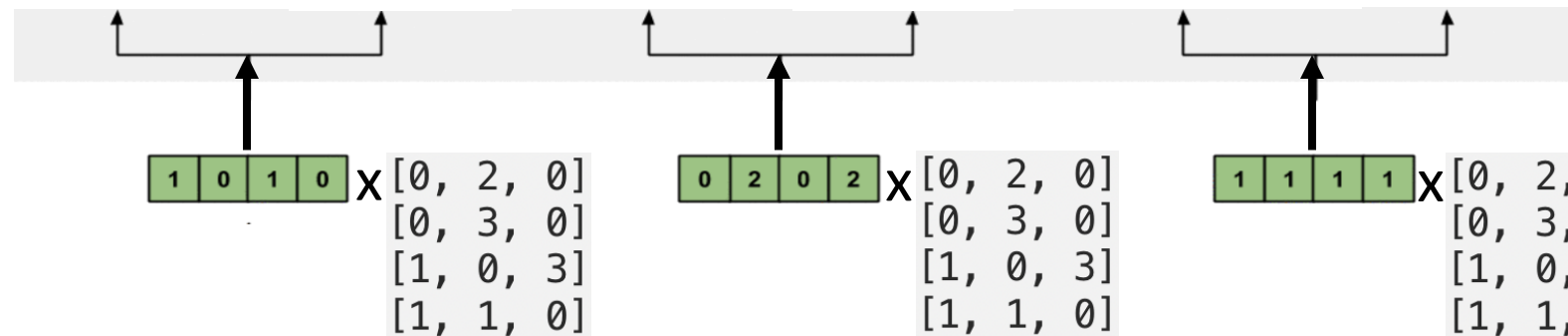
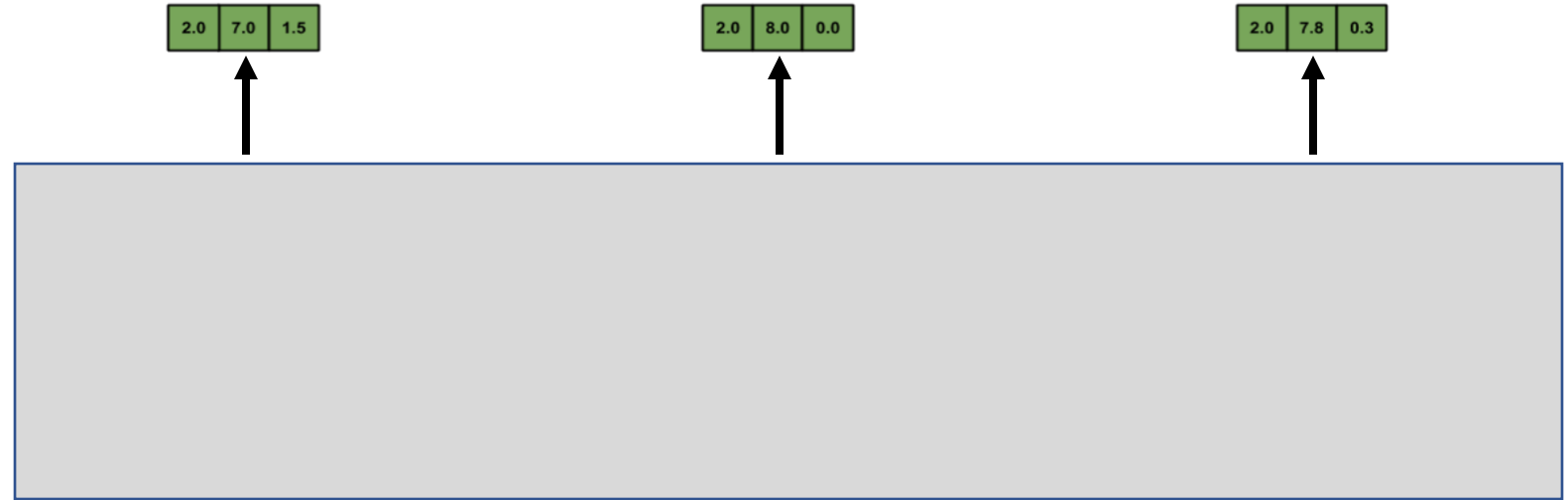
$$\begin{bmatrix} 0 & 2 & 0 & 2 \end{bmatrix} \times \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix}$$

Computing Self-Attention: Example

e.g., **value** weights

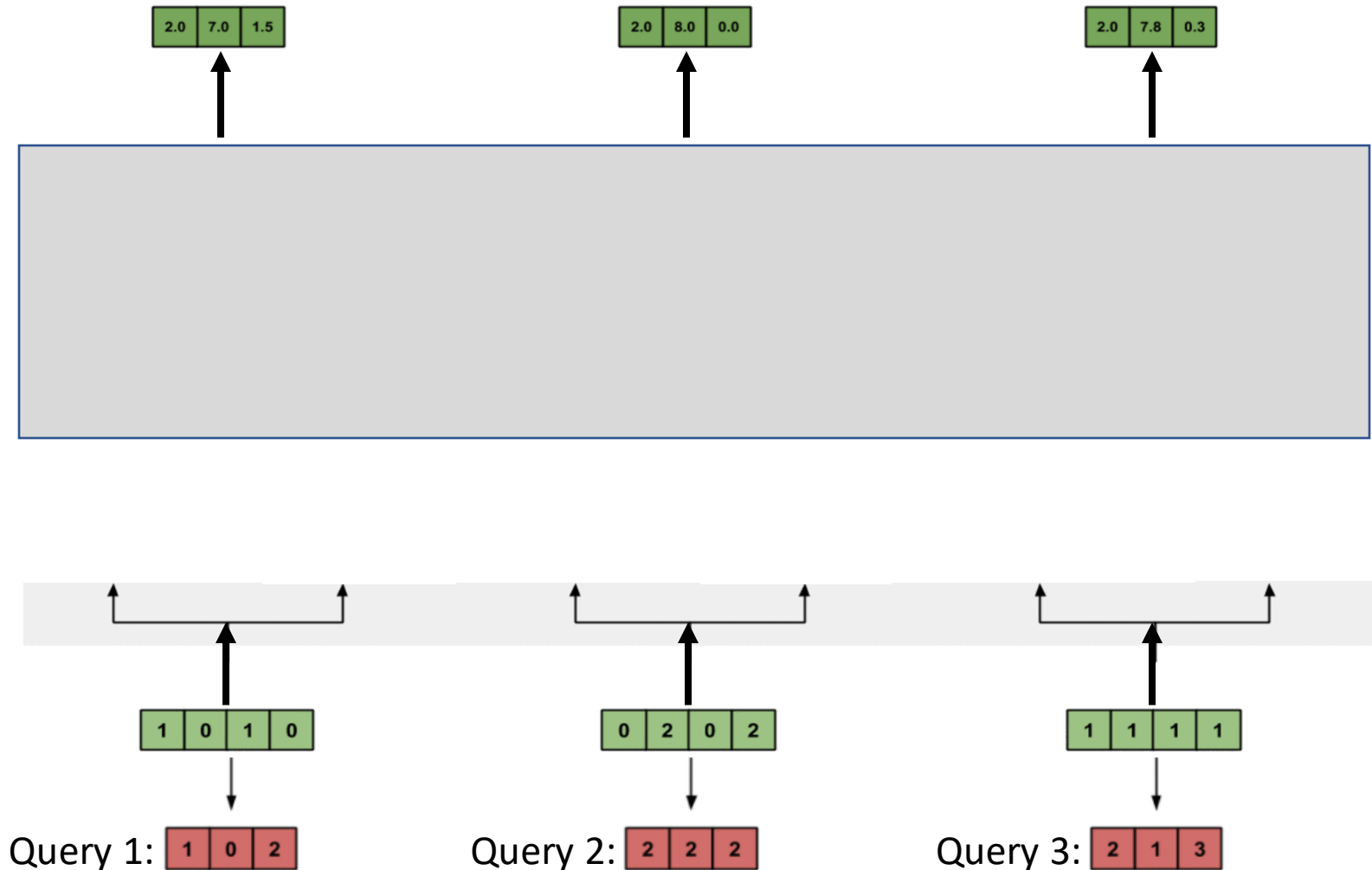
```
[0, 2, 0]  
[0, 3, 0]  
[1, 0, 3]  
[1, 1, 0]
```



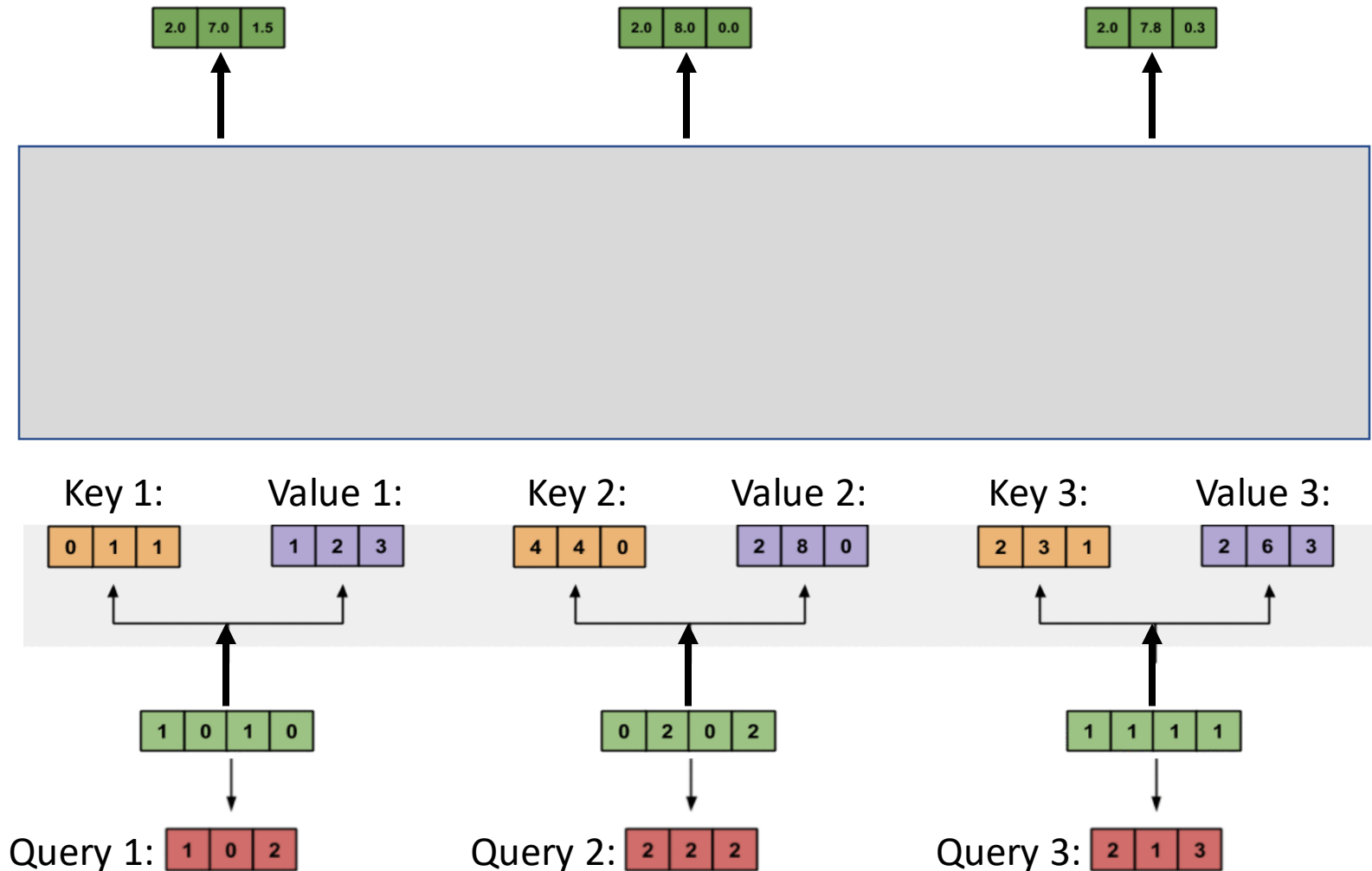
Computing Self-Attention: Example

e.g., **query** weights

```
[1, 0, 1]  
[1, 0, 0]  
[0, 0, 1]  
[0, 1, 1]
```



Computing Self-Attention: Example

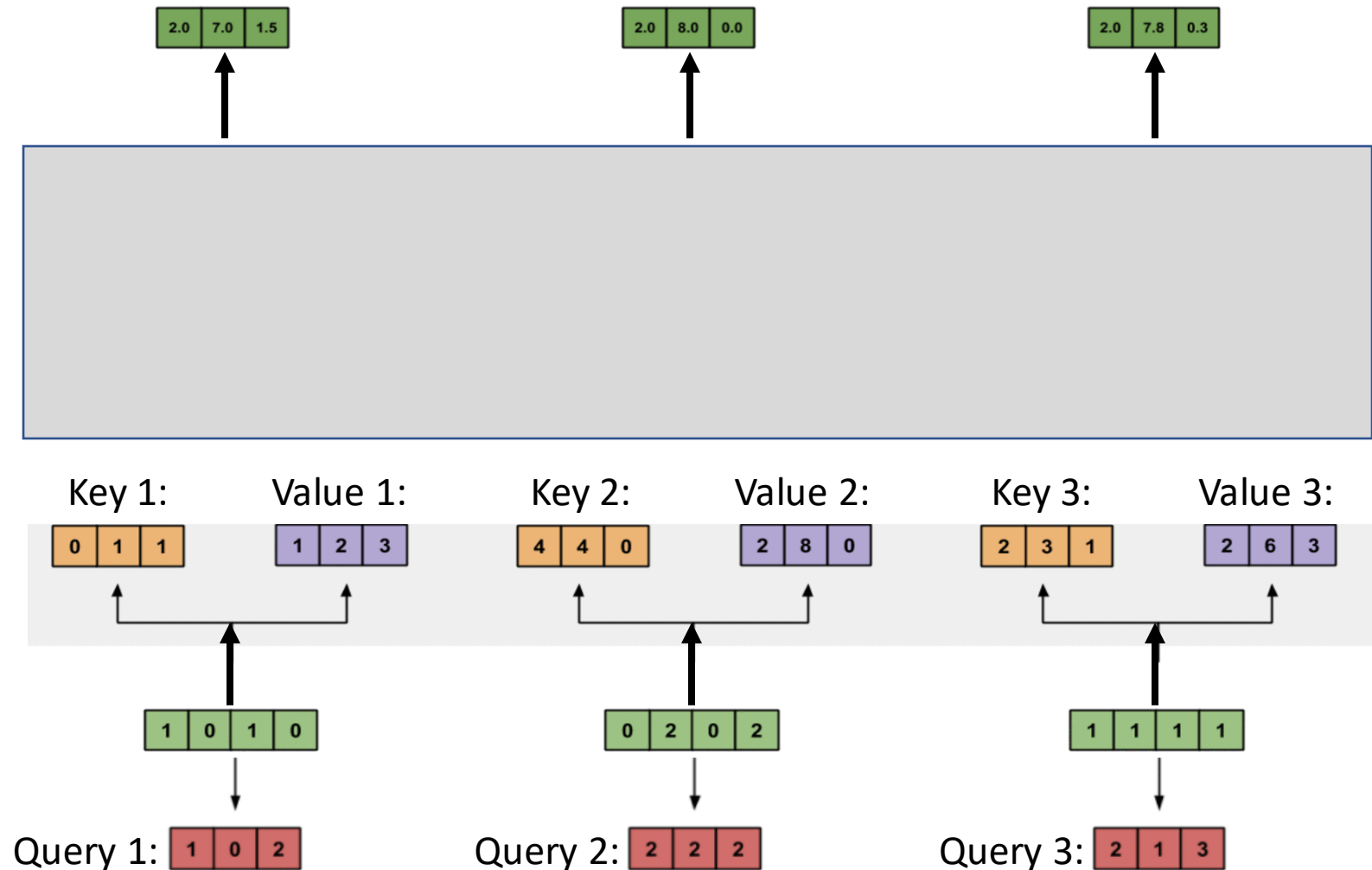


How many weight matrices are learned in this example?

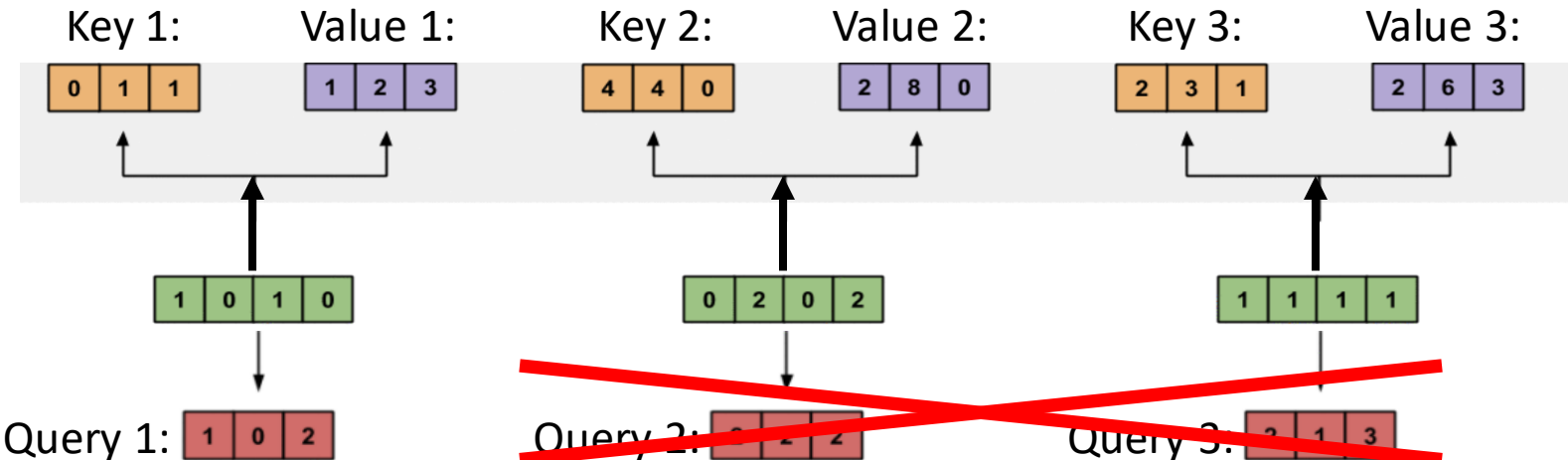
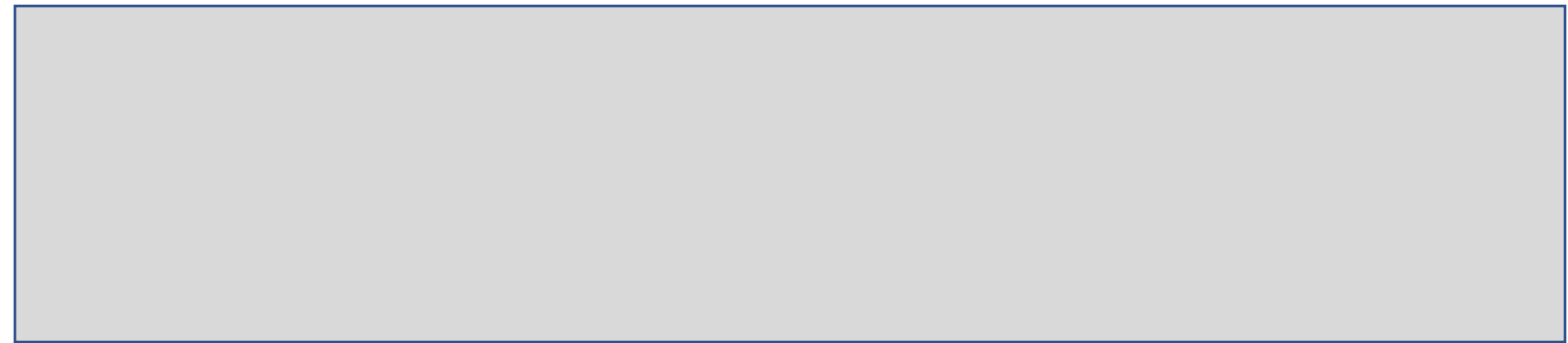
Computing Self-Attention: Example

What is the purpose of the three weight matrices?

For each **input**, 2 of the derived vectors are used to compute **attention weights** (**query** and **key**) and the 3rd is **information** passed on for the new representation (**value**)



Computing Self-Attention: Example

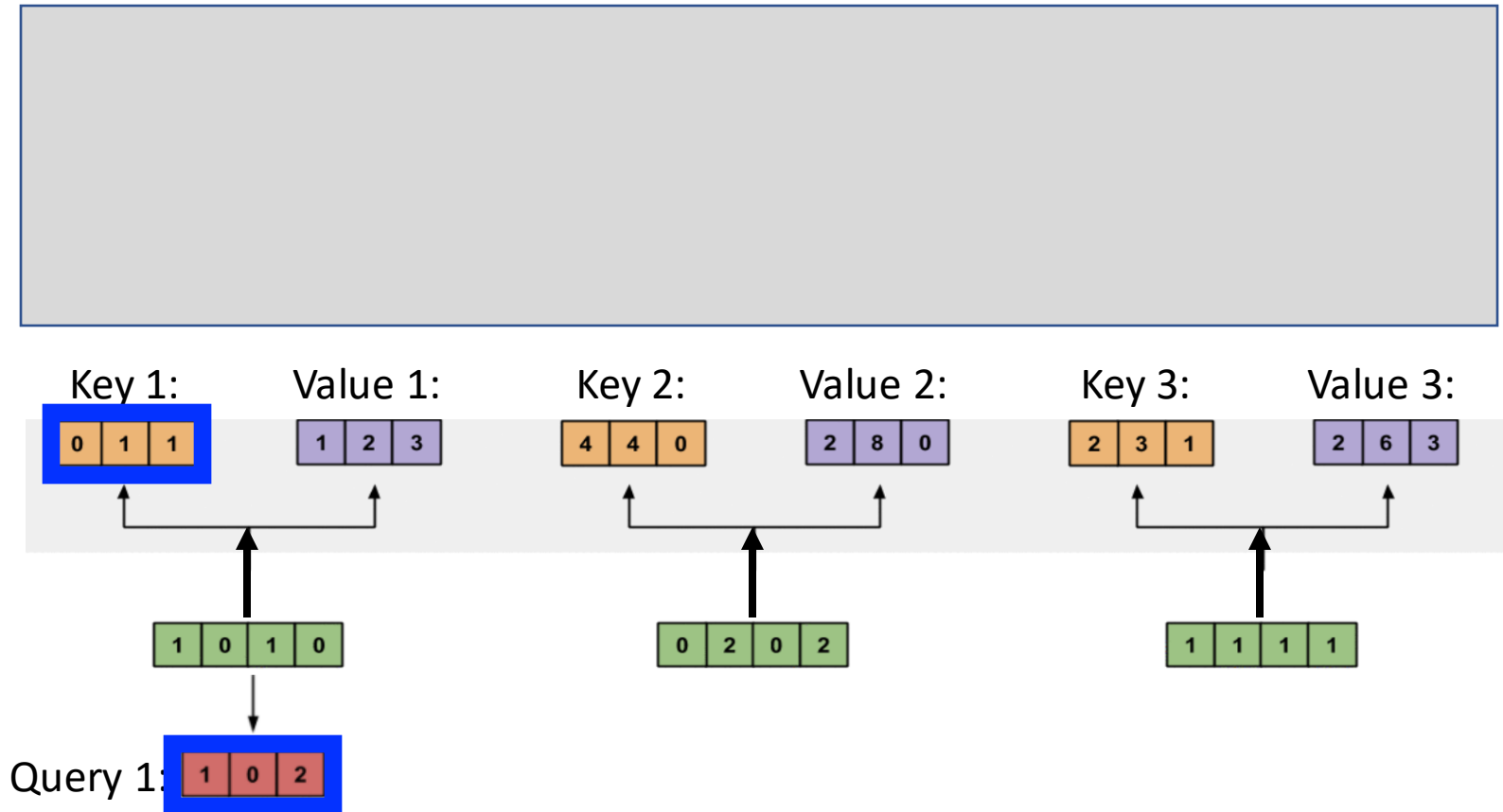


Let's compute the new representation for the inputs...

Computing Self-Attention: Example

Attention score: dot product of **query** (“what am I looking for”) with all **keys** (“what I have”) to identify relevant tokens (higher scores are better matches); e.g.,

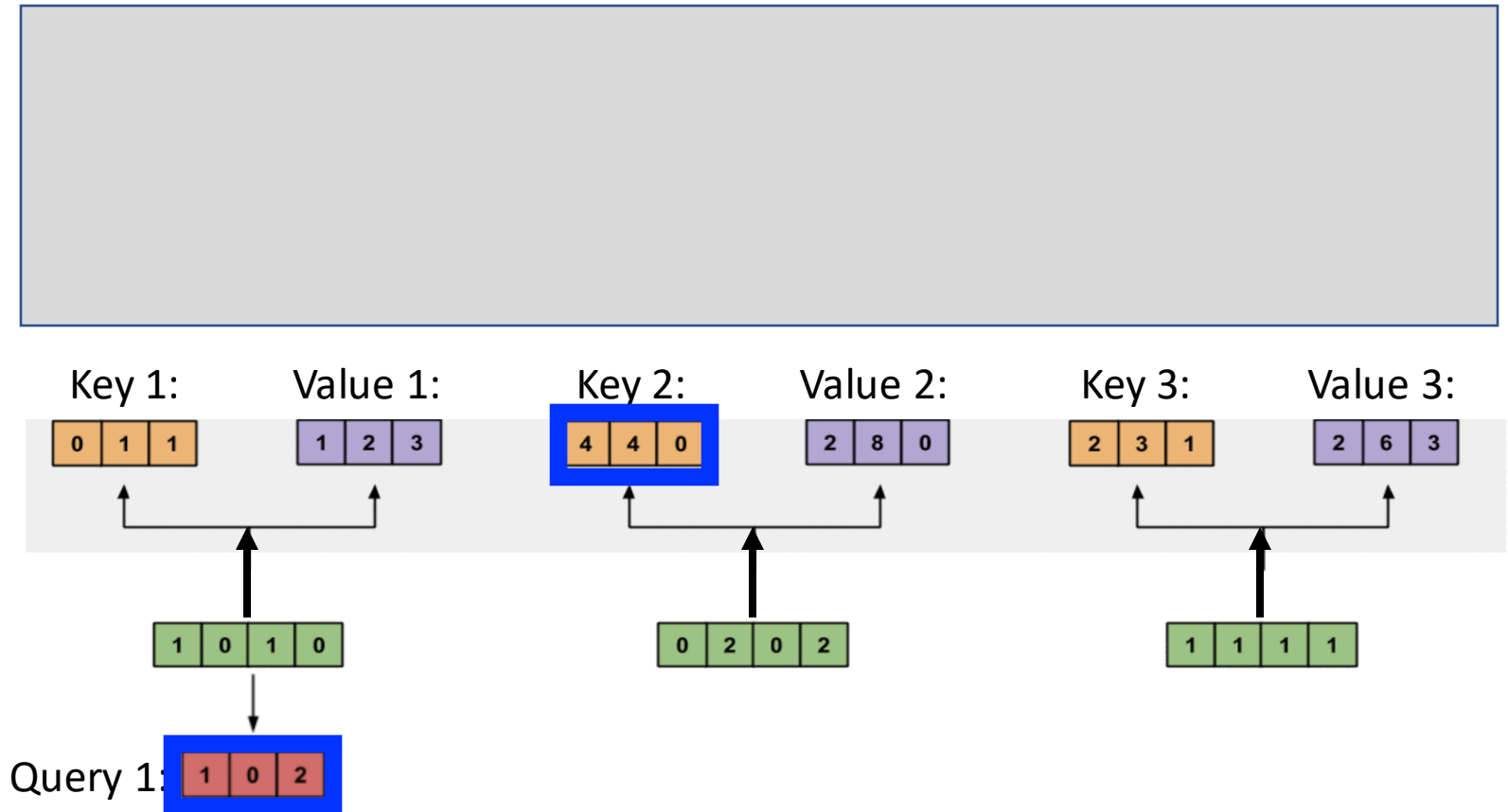
$$\begin{bmatrix} 1 & 0 & 2 \end{bmatrix} \times \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} = ?$$



Computing Self-Attention: Example

Attention score: dot product of **query** (“what am I looking for”) with all **keys** (“what I have”) to identify relevant tokens (higher scores are better matches); e.g.,

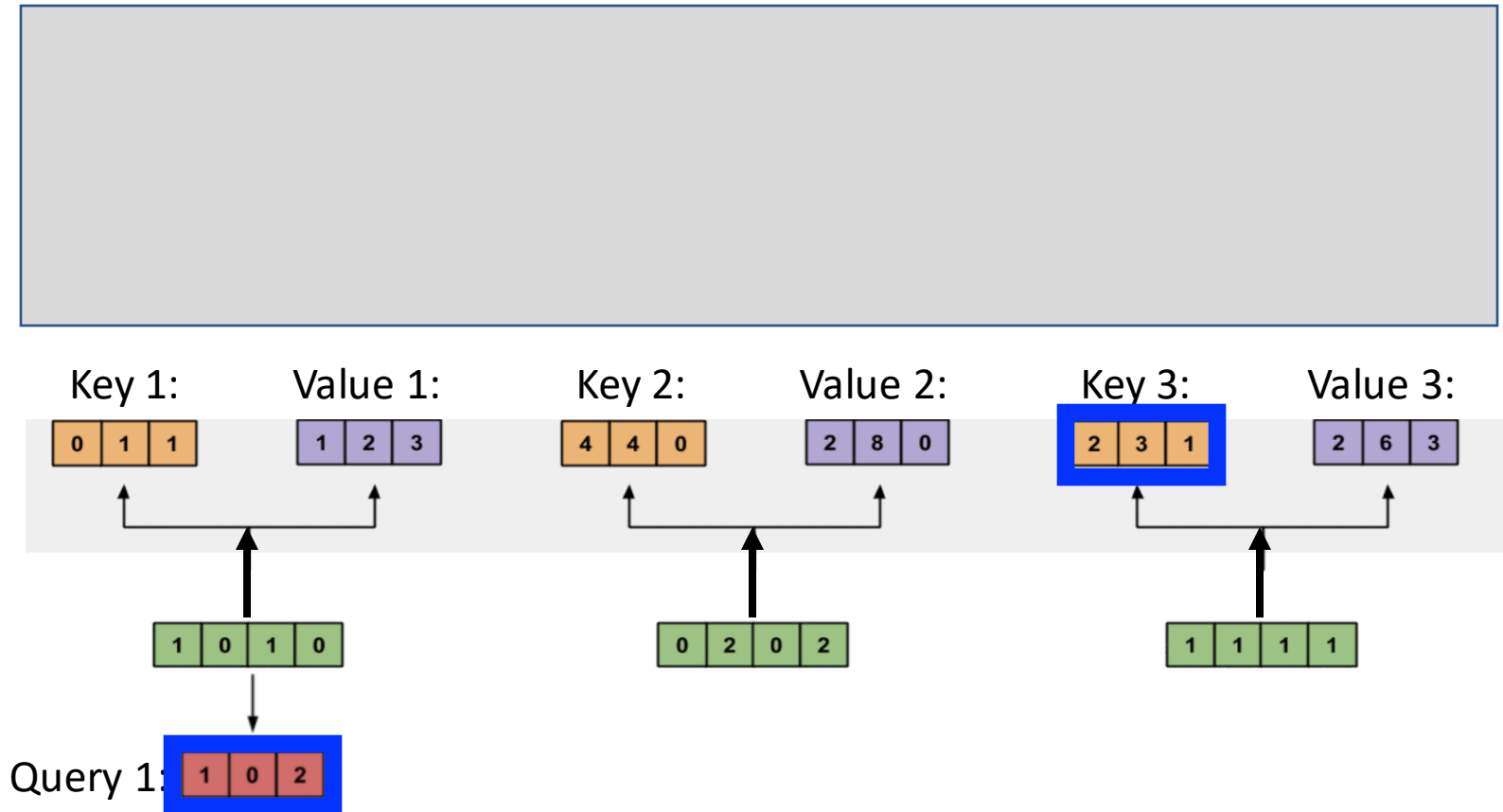
$$\begin{bmatrix} 1 & 0 & 2 \end{bmatrix} \times \begin{bmatrix} 4 \\ 4 \\ 0 \end{bmatrix} = ?$$



Computing Self-Attention: Example

Attention score: dot product of **query** (“what am I looking for”) with all **keys** (“what I have”) to identify relevant tokens (higher scores are better matches); e.g.,

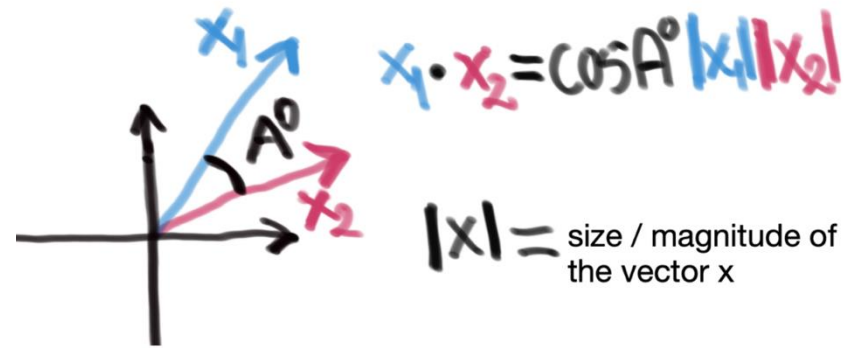
$$\begin{bmatrix} 1 & 0 & 2 \end{bmatrix} \times \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} = ?$$



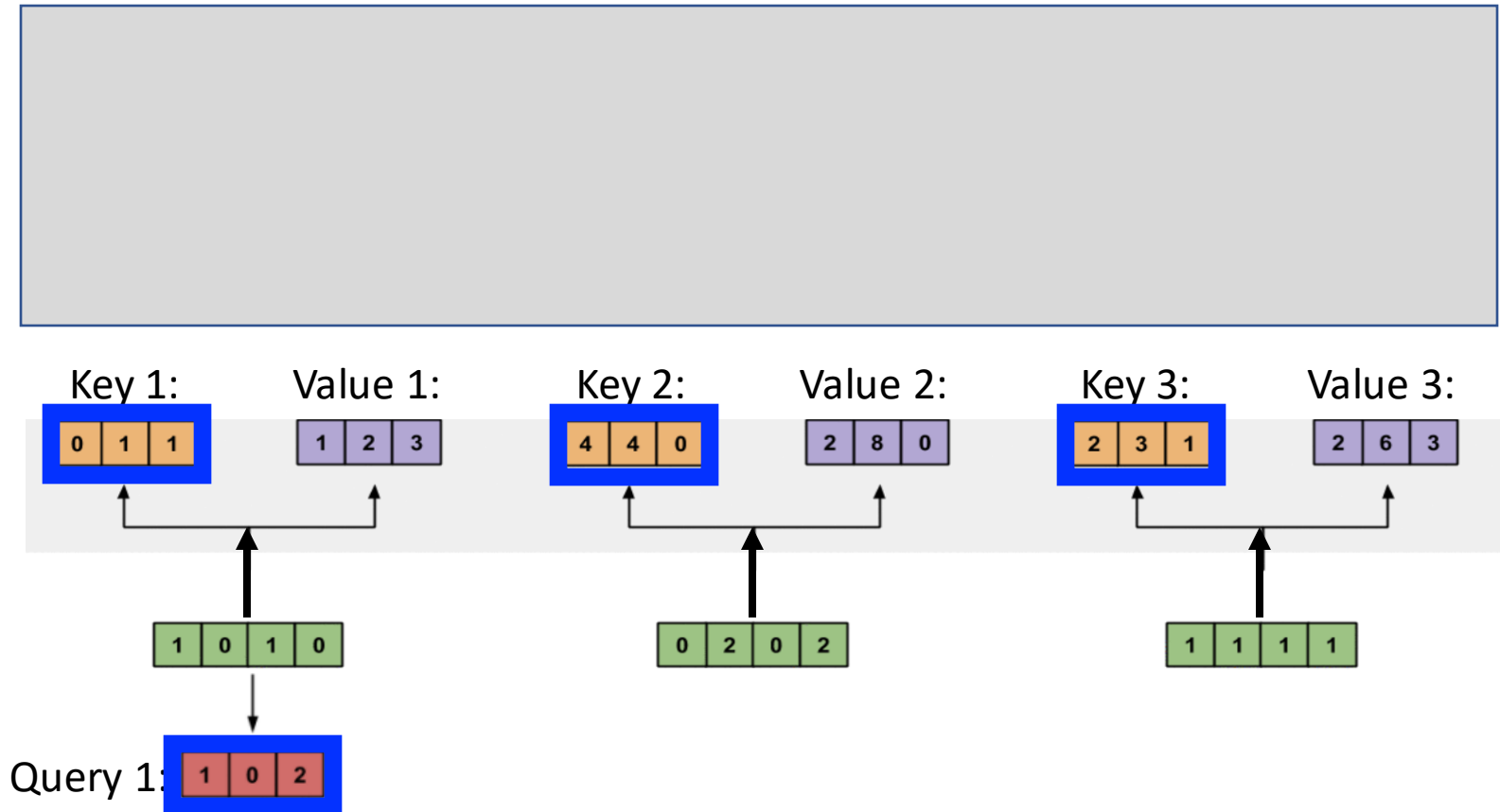
Computing Self-Attention: Example

Why dot product? Indicates similarity of two vectors

- Match = 1 (i.e., $\cos(0)$)
- Opposites = -1 (i.e., $\cos(180)$)



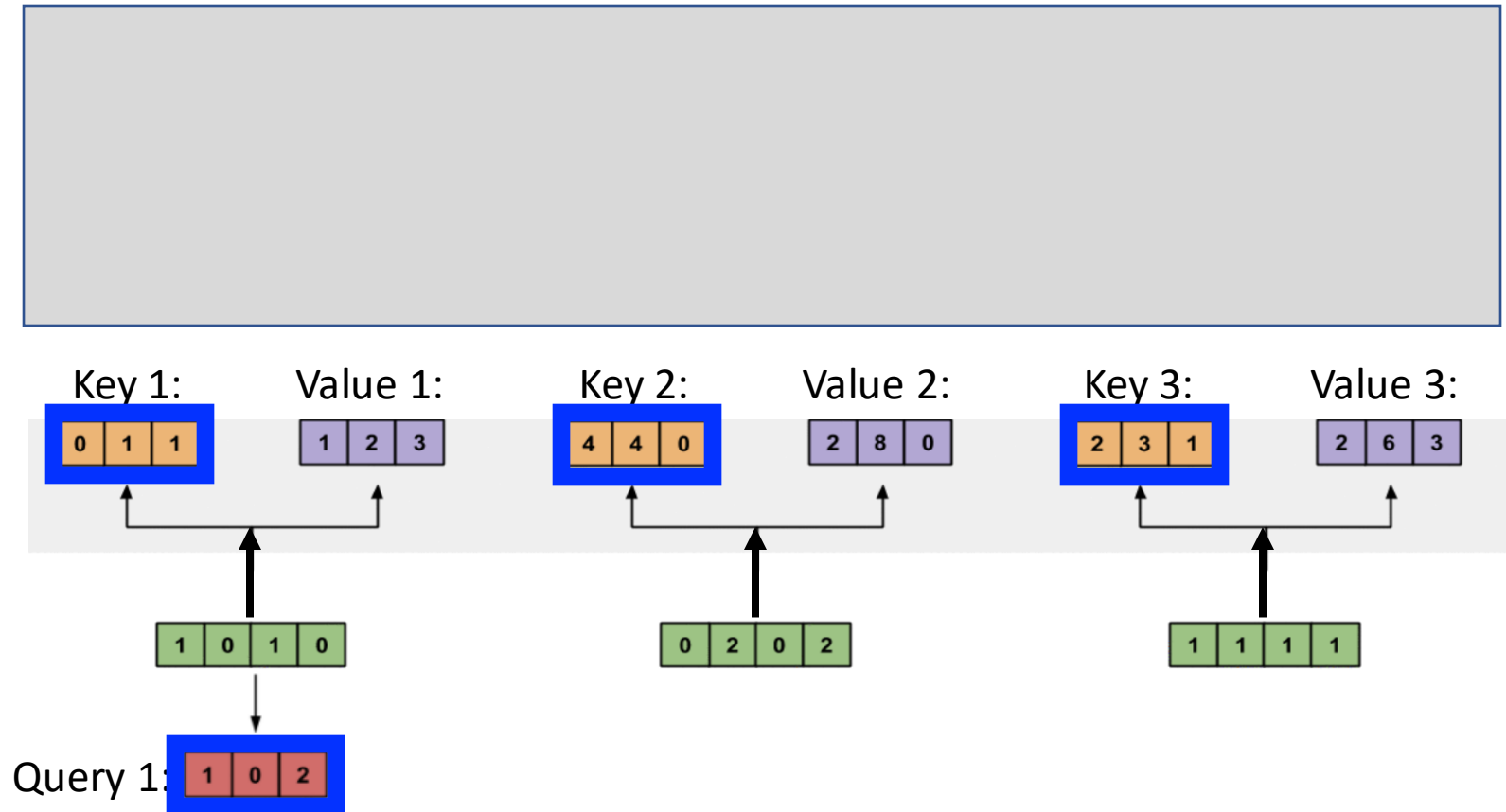
<https://towardsdatascience.com/self-attention-5b95ea164f61>



<https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a>

Computing Self-Attention: Example

Recall: there are many compatibility measures



Computing Self-Attention: Example

Attention weights: softmax scores for all inputs to quantify each token's relevance; e.g.,

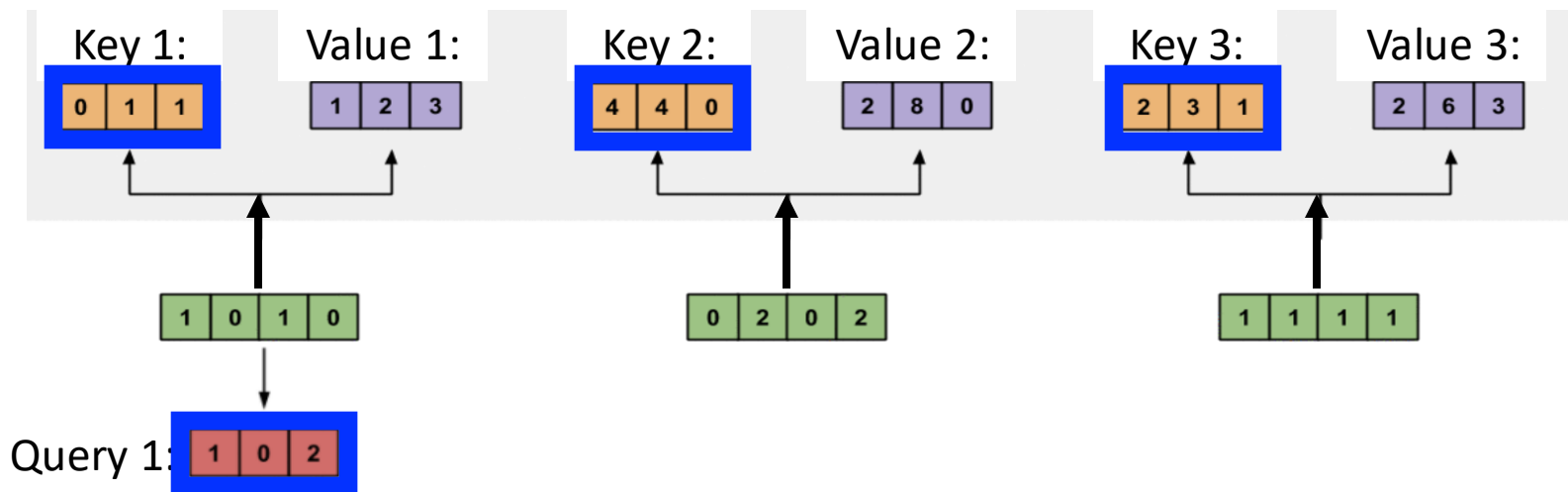
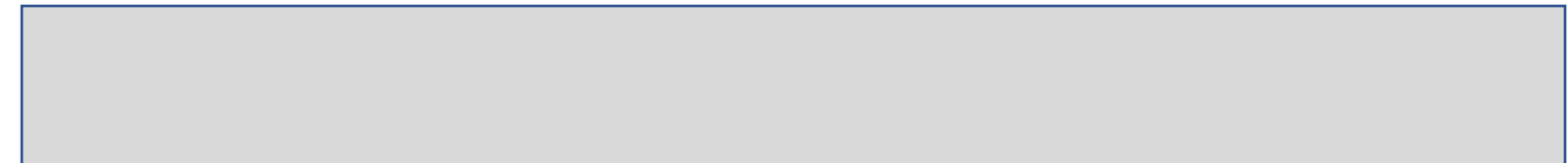
$$= \text{softmax}([2, 4, 4])$$

$$= [0.0, 0.5, 0.5]$$

Note: 0 from softmax can arise from rounding

To which input(s) is input 1 **least** related?

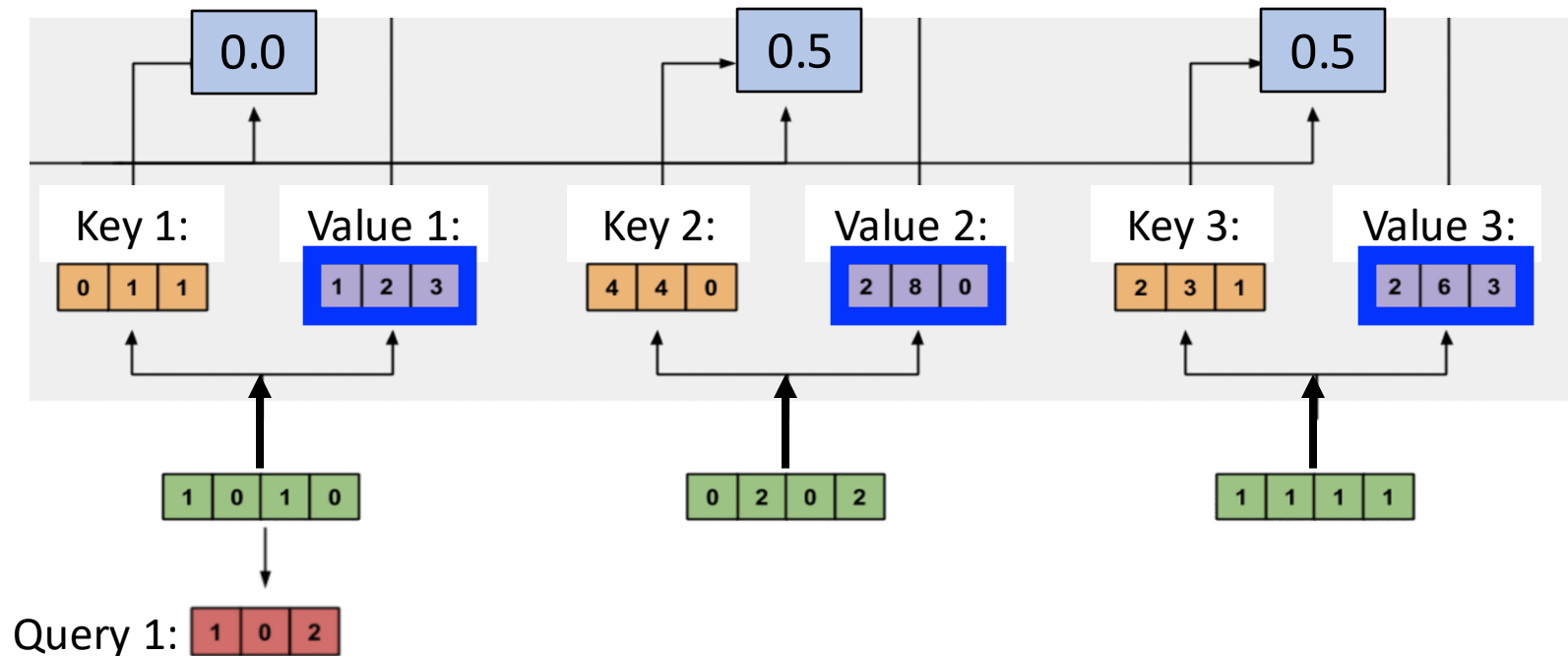
To which input(s) is input 1 **most** related?



Computing Self-Attention: Example

Compute **new representation** of **input token** that reflects entire input:

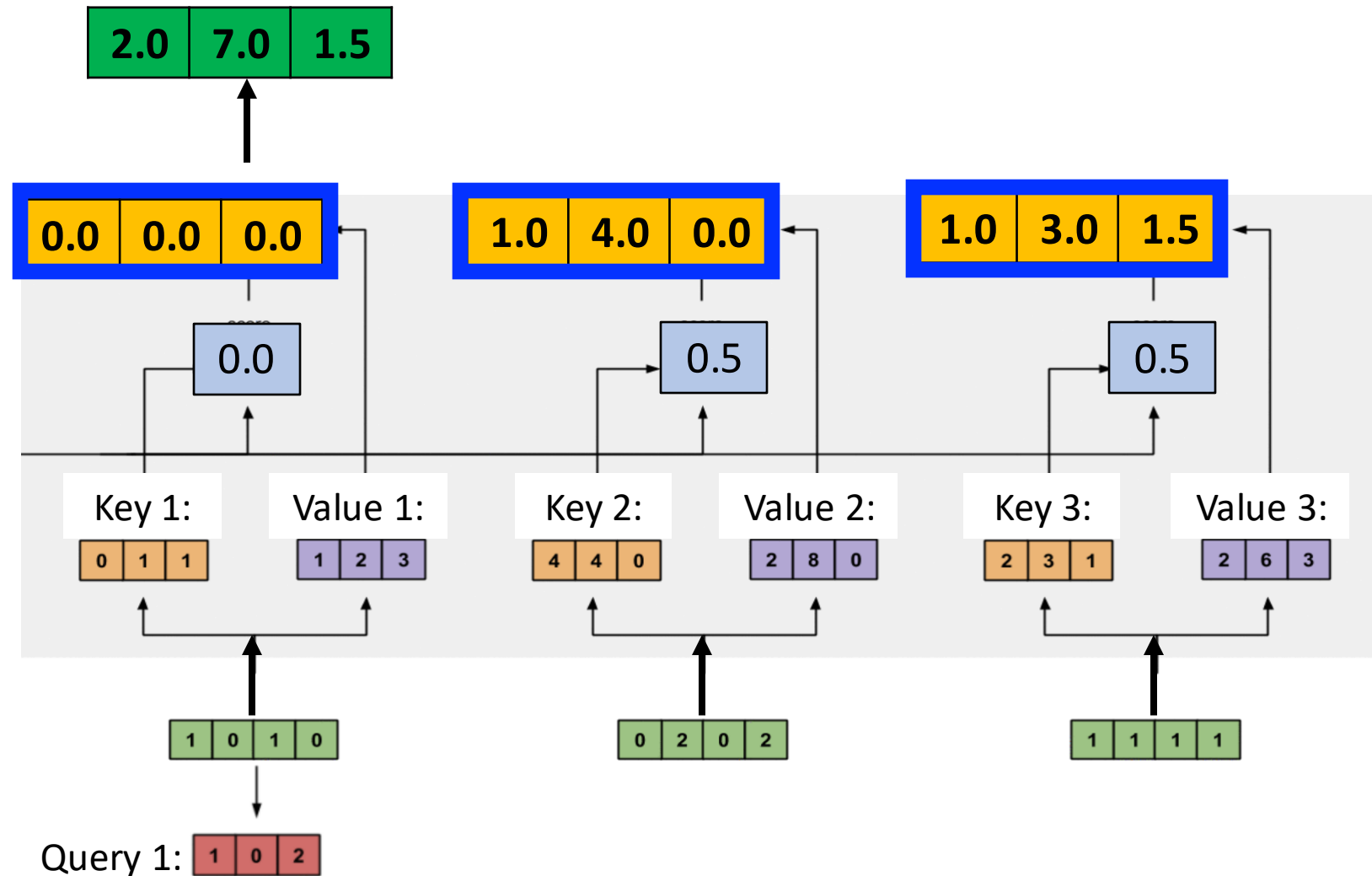
1. Attention weights x Values



Computing Self-Attention: Example

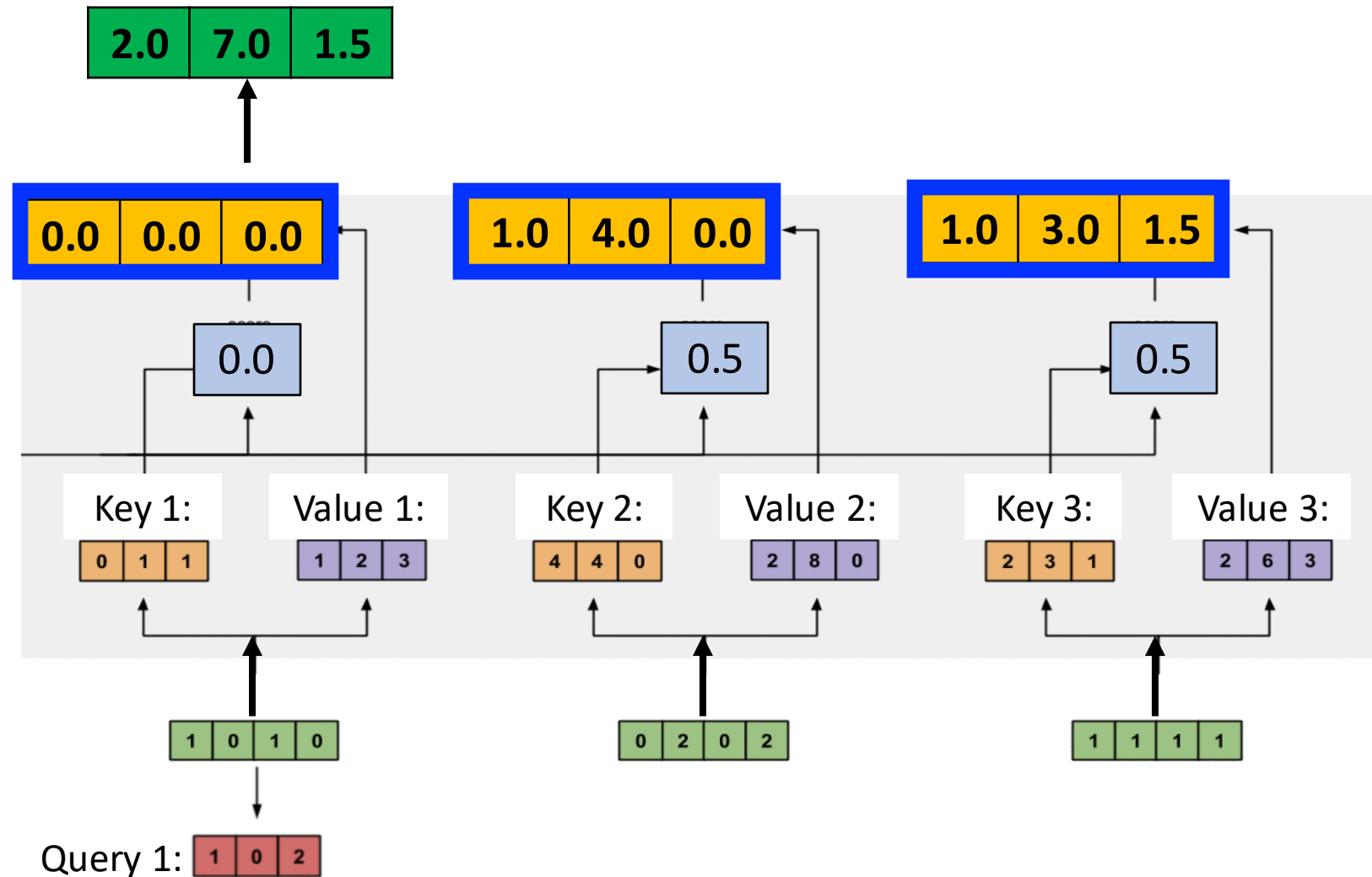
Compute **new representation** of **input token** that reflects entire input:

1. **Attention weights** x **Values**
2. Sum all **weighted vectors**



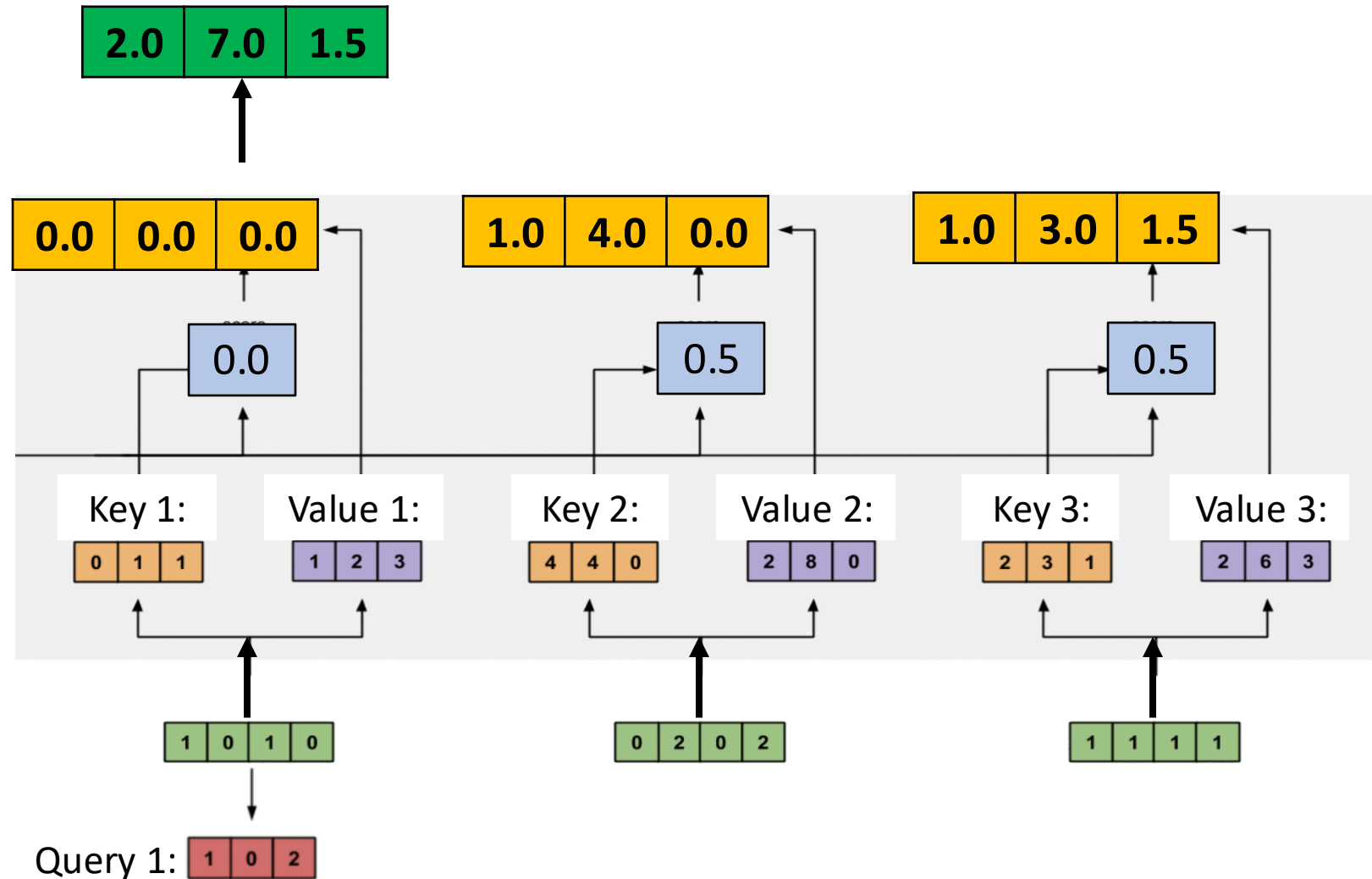
Computing Self-Attention: Example

Attention weights amplify input representations (values) that we want to pay attention to and repress the rest



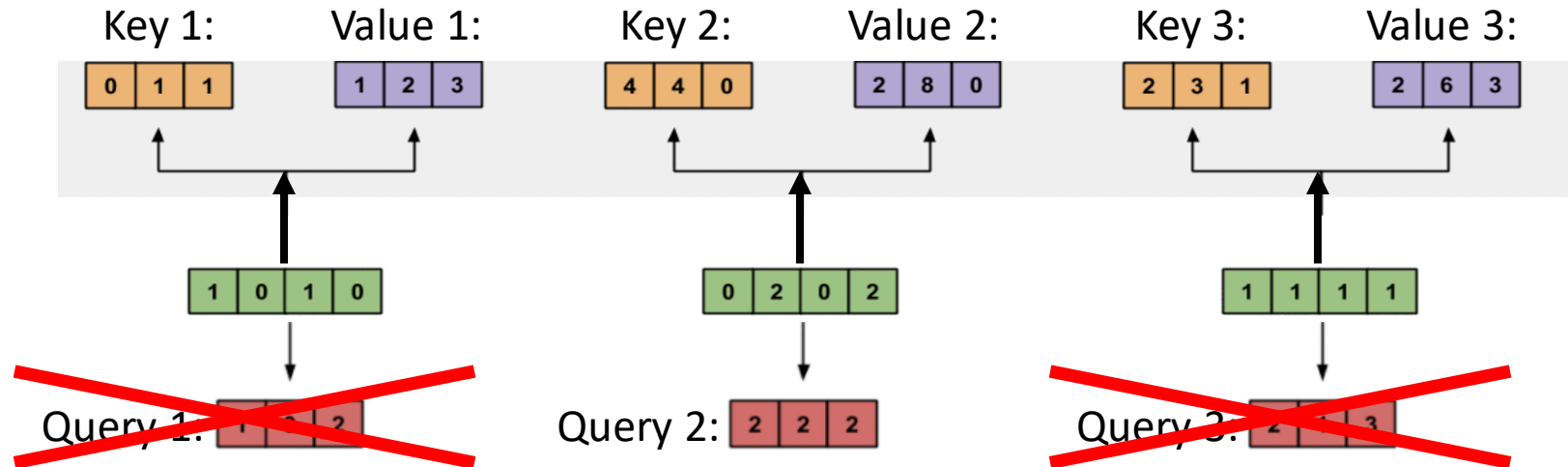
Computing Self-Attention: Example

Attention weights amplify input representations (values) that we want to pay attention to and repress the rest



Computing Self-Attention: Example

Repeat the same process for each remaining input token

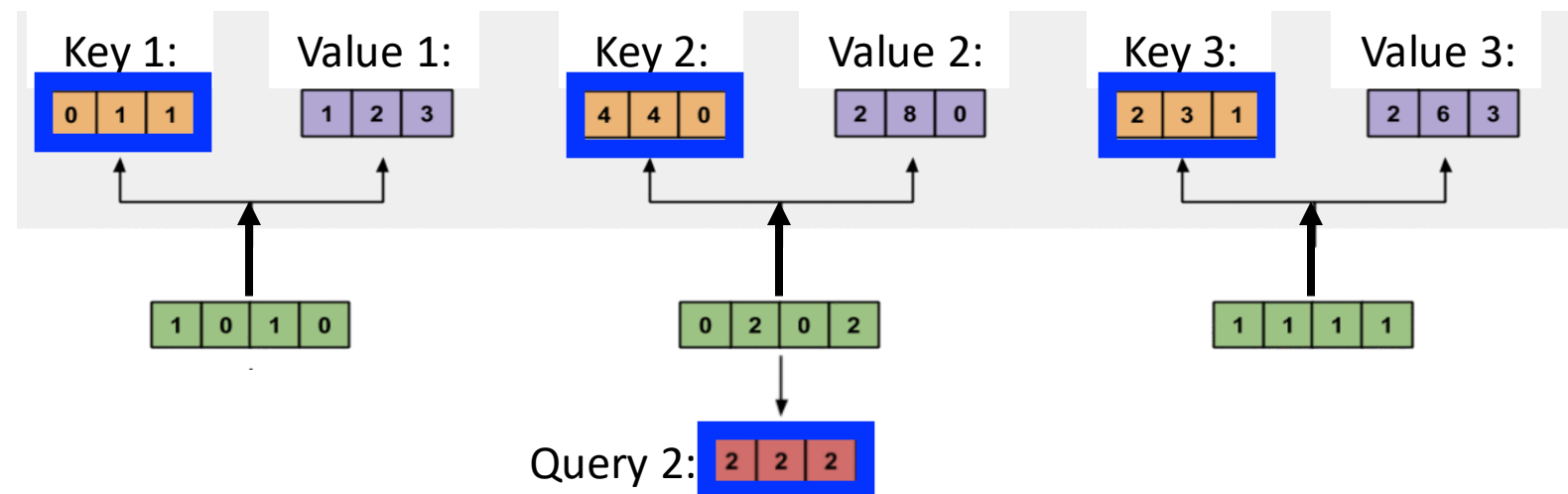


Computing Self-Attention: Example

1. Compute attention weights

- Softmax resulting 3 scores from **query** x **keys**

To which input(s) is input 2 most related?

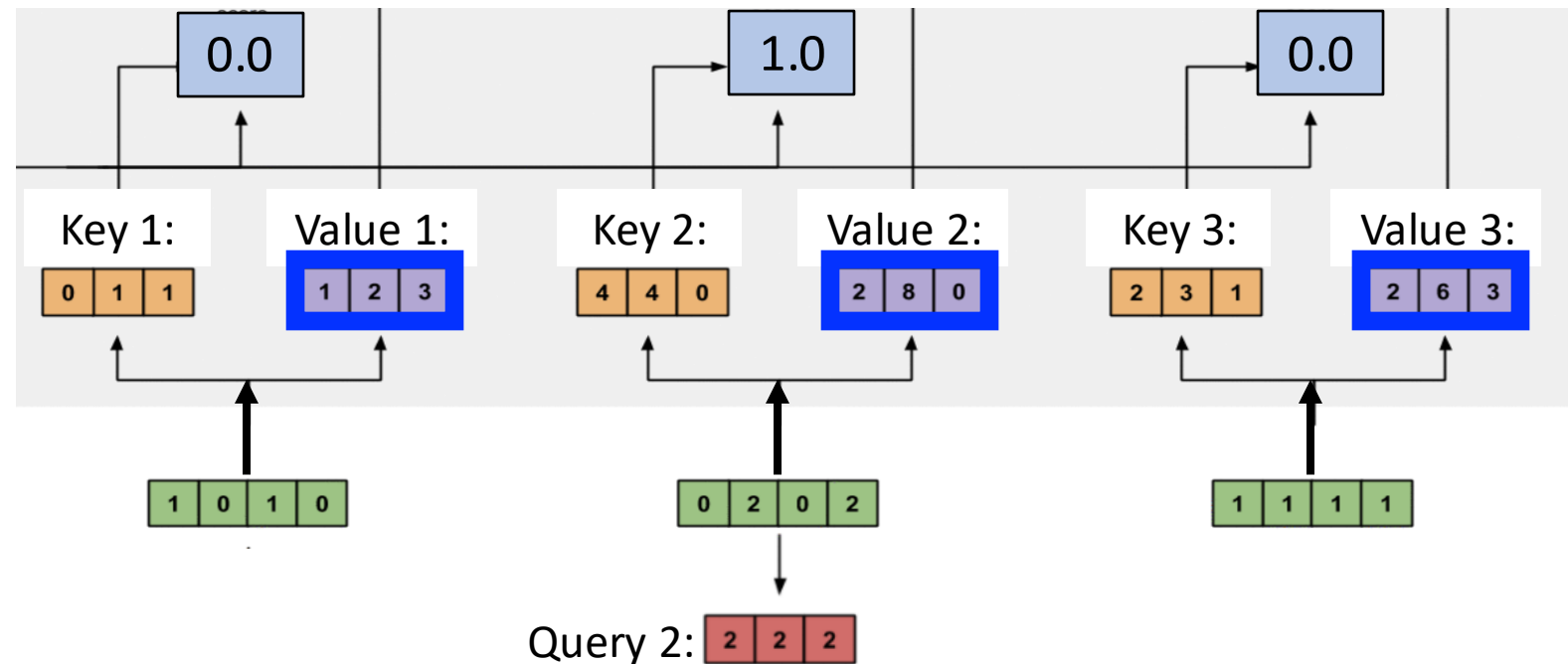


Computing Self-Attention: Example

1. Compute attention weights

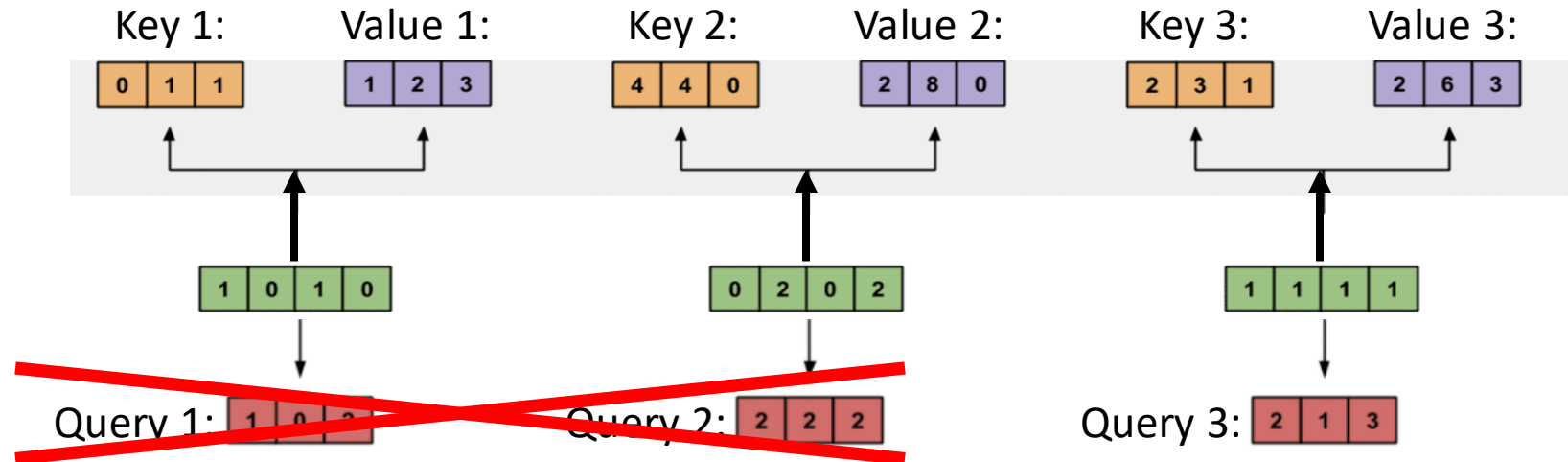
- Softmax resulting 3 scores from query x keys

2. Compute weighted sum of values using attention scores



Computing Self-Attention: Example

Repeat the same process for each remaining input token

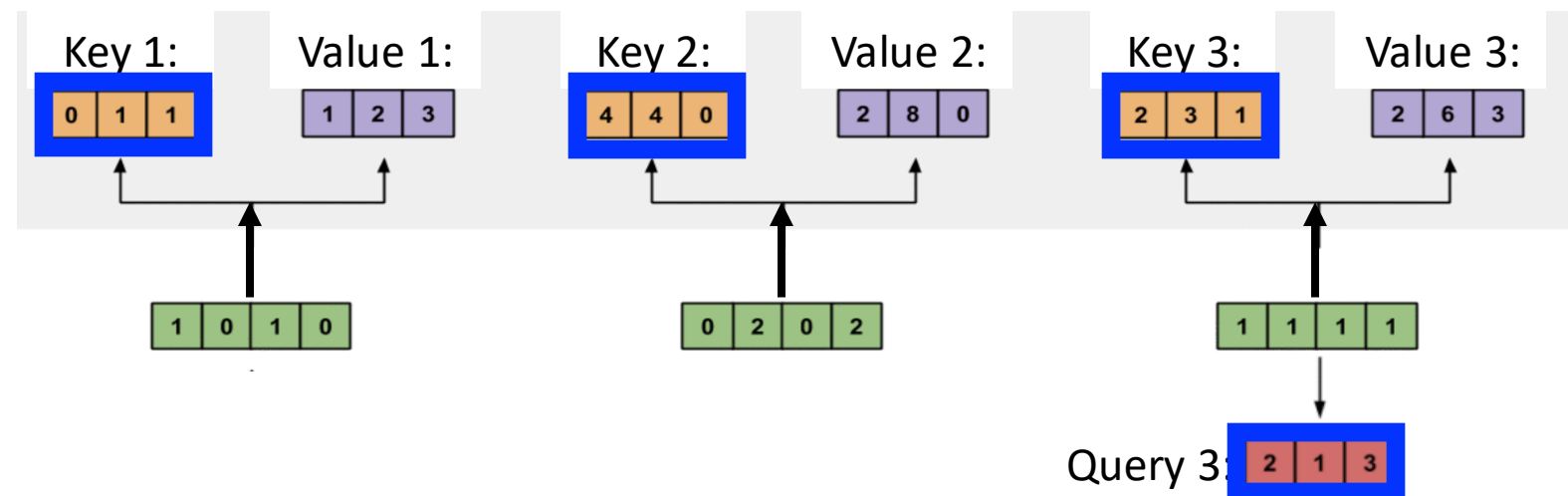


Computing Self-Attention: Example

1. Compute attention weights

- Softmax resulting 3 scores from **query** x **keys**

To which input(s) is input 3 most related?

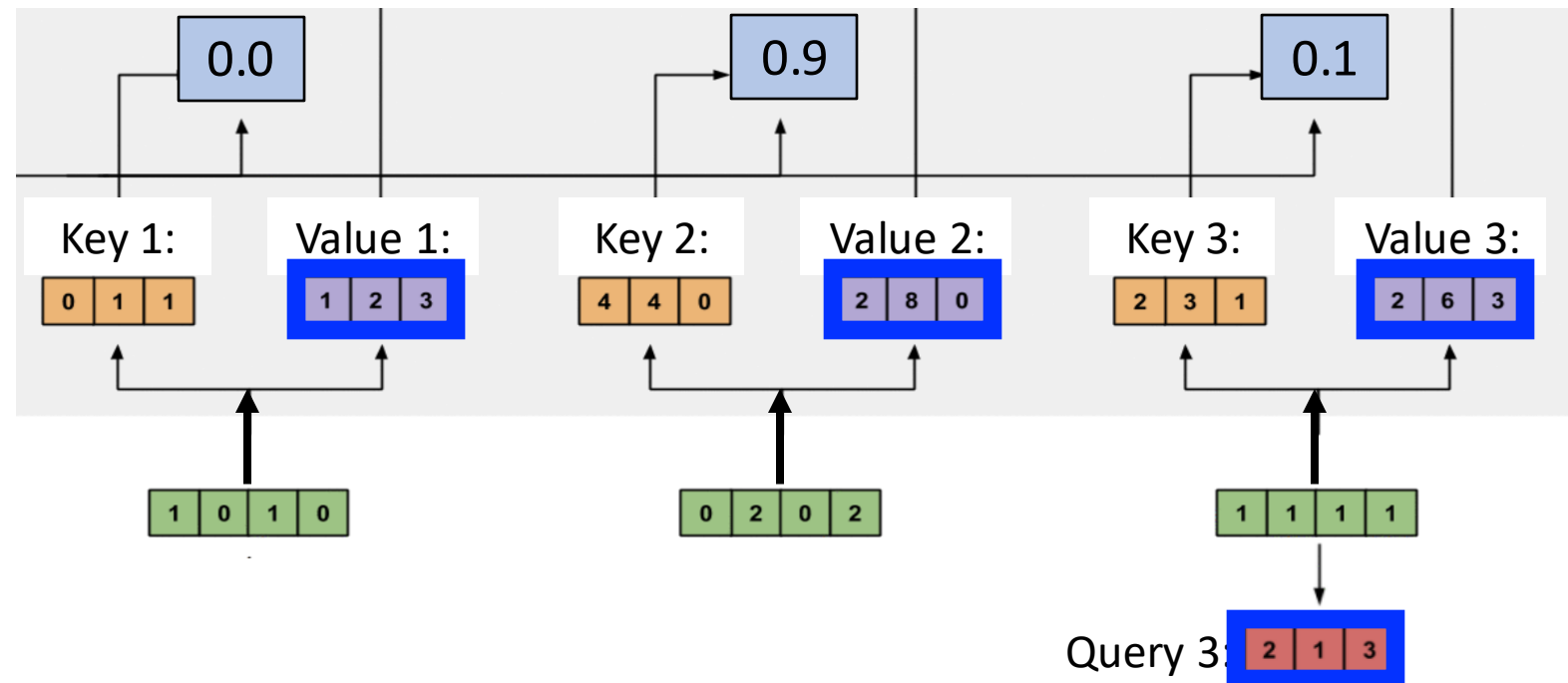


Computing Self-Attention: Example

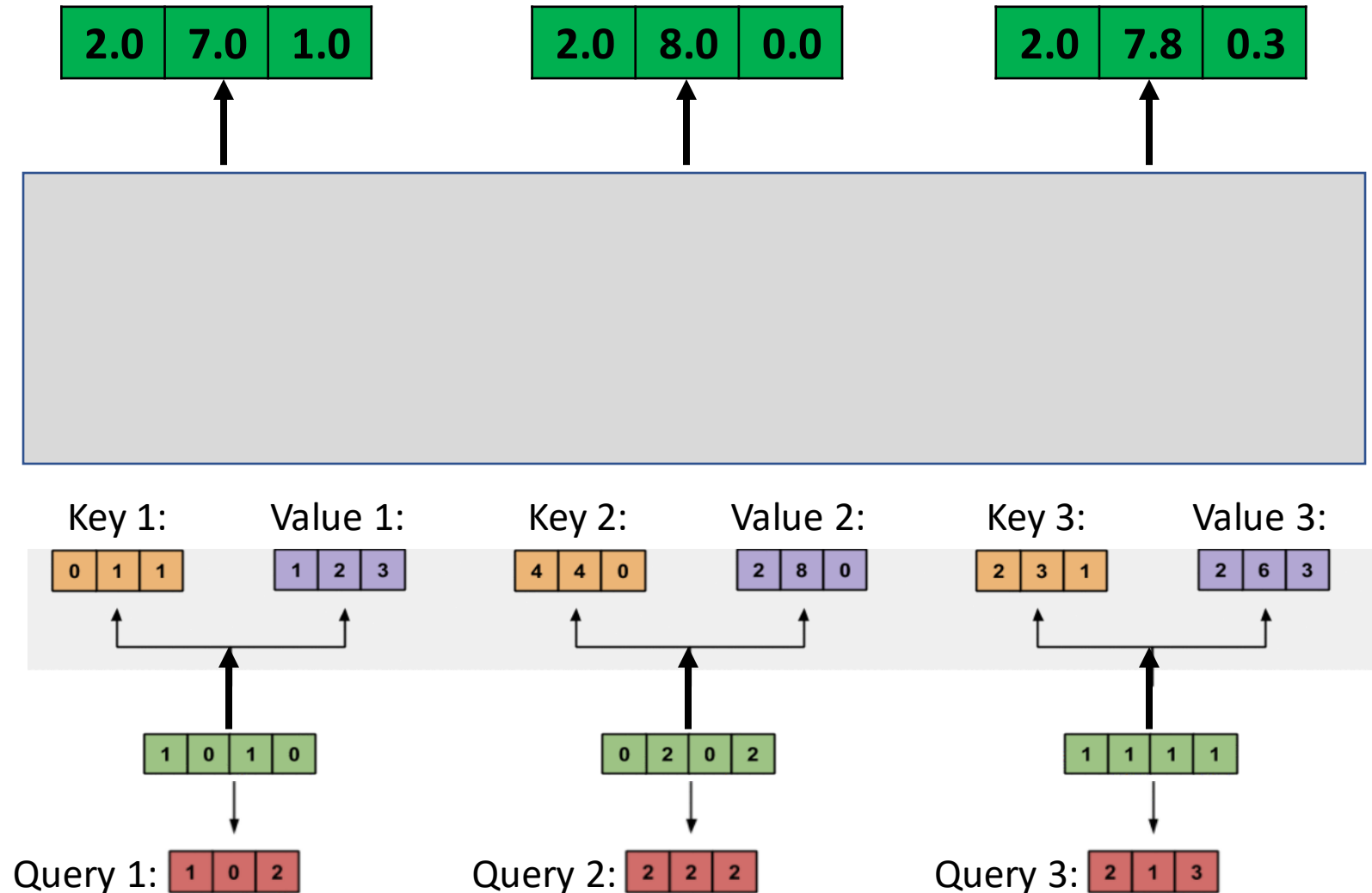
1. Compute attention weights

- Softmax resulting 3 scores from **query** x **keys**

2. Compute weighted sum of values using attention scores



Computing Self-Attention: Example

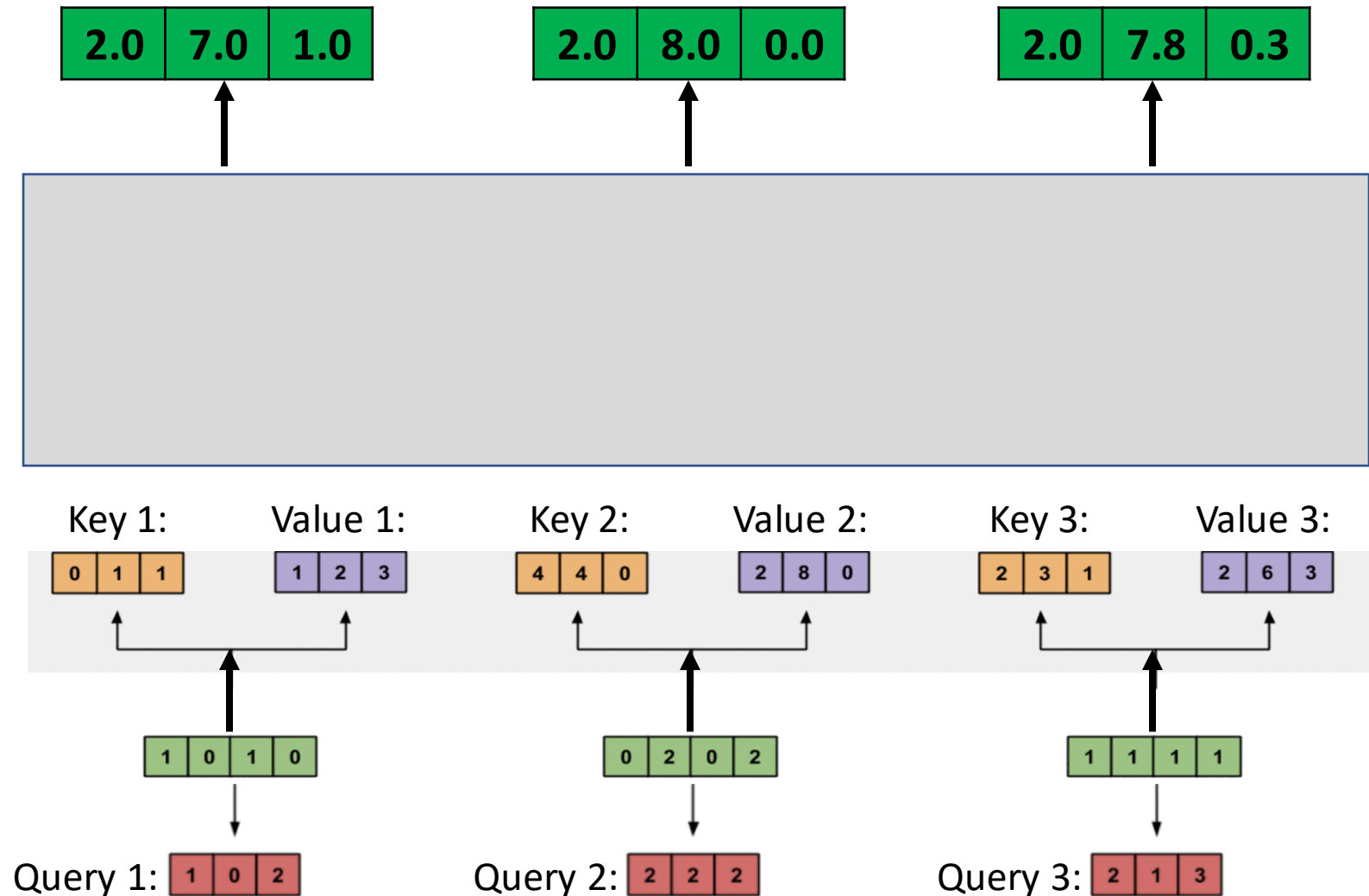


Hyperparameters

A developer chooses **input token length** and **number of matrices' columns** (and thus vector sizes)

Dimension of **query** and **key** must match to assess similarity (e.g., dot product).

Dimension of **value** can differ from that of **query** and **key** and is output dimension.



Efficient Computation for Self-Attention

Step 1

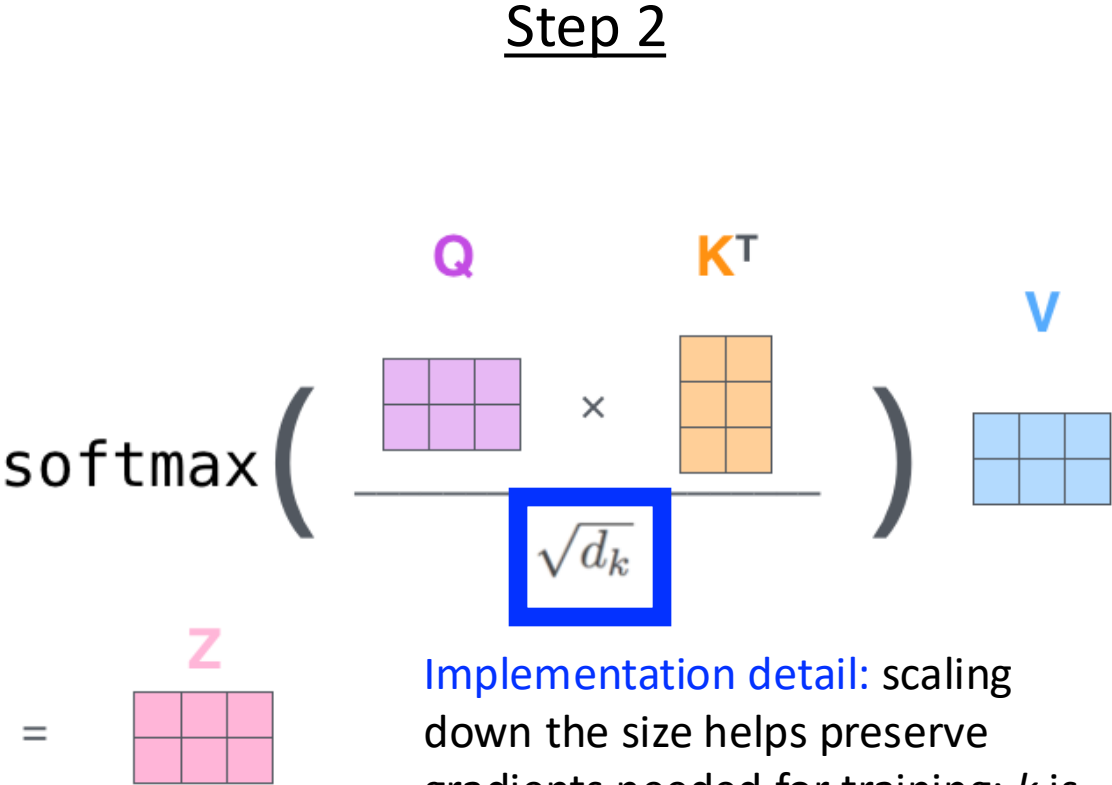
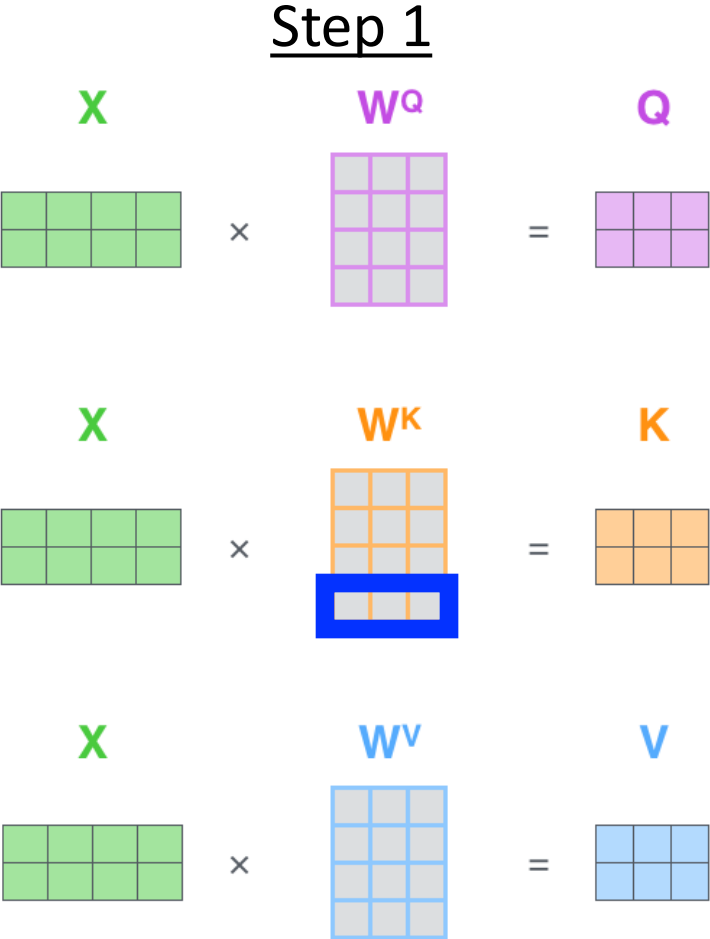
Each row is an
input token:

$$\begin{matrix} \mathbf{X} \\ \text{[Green 2x4 grid]} \end{matrix} \times \begin{matrix} \mathbf{W}^Q \\ \text{[Purple 4x4 grid]} \end{matrix} = \begin{matrix} \mathbf{Q} \\ \text{[Purple 2x4 grid]} \end{matrix} \quad \text{Each row is a query}$$

$$\begin{matrix} \mathbf{X} \\ \text{[Green 2x4 grid]} \end{matrix} \times \begin{matrix} \mathbf{W}^K \\ \text{[Orange 4x4 grid]} \end{matrix} = \begin{matrix} \mathbf{K} \\ \text{[Orange 2x4 grid]} \end{matrix} \quad \text{Each row is a key}$$

$$\begin{matrix} \mathbf{X} \\ \text{[Green 2x4 grid]} \end{matrix} \times \begin{matrix} \mathbf{W}^V \\ \text{[Blue 4x4 grid]} \end{matrix} = \begin{matrix} \mathbf{V} \\ \text{[Blue 2x4 grid]} \end{matrix} \quad \text{Each row is a value}$$

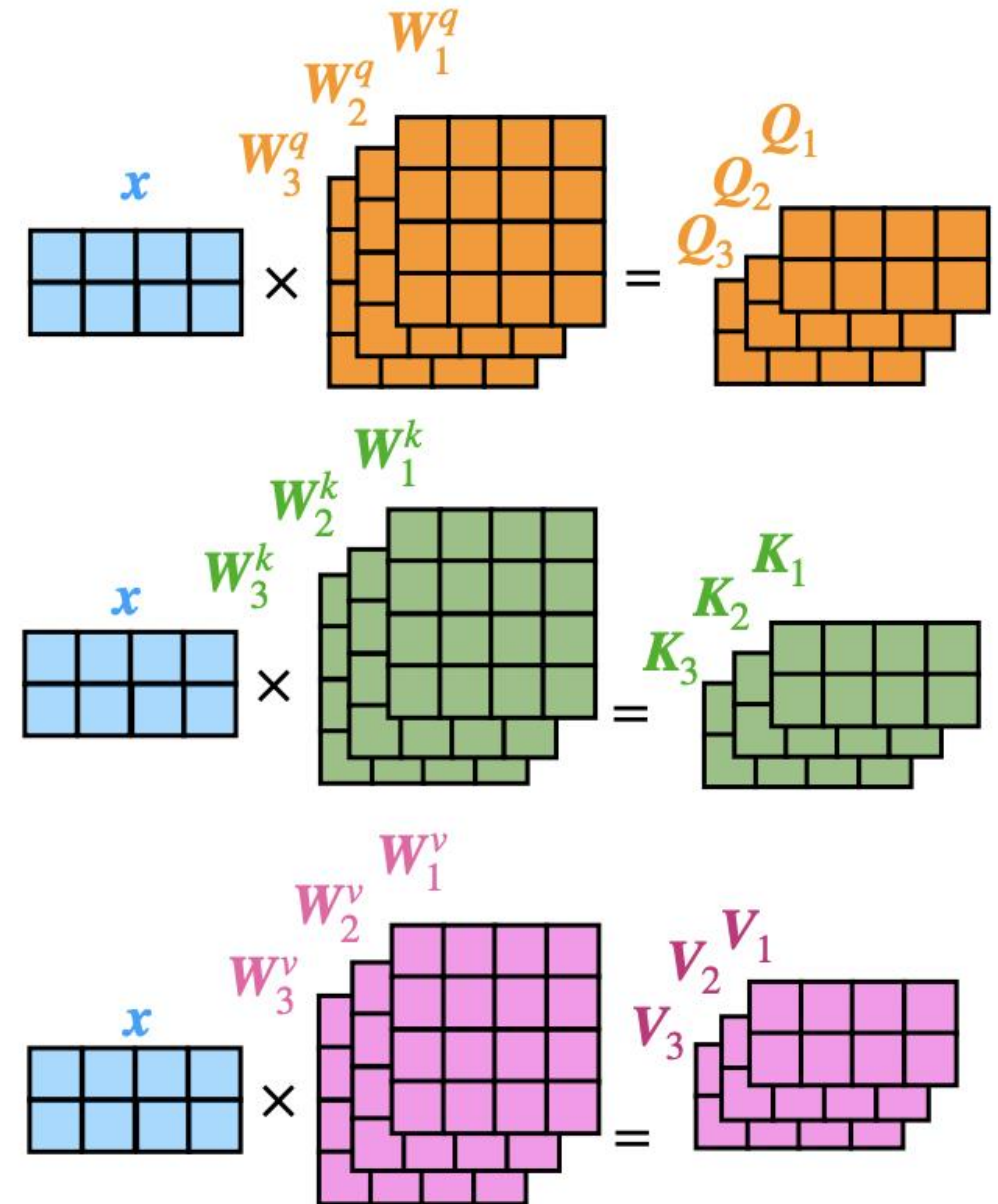
Efficient Computation for Self-Attention



Implementation detail: scaling down the size helps preserve gradients needed for training; k is dimensionality of the key vector

Multi-head Attention

- **Goal:** enable each token to relate to other tokens in multiple ways
- **Key idea:** multiple self-attention mechanisms, each with their own key, value and query matrices

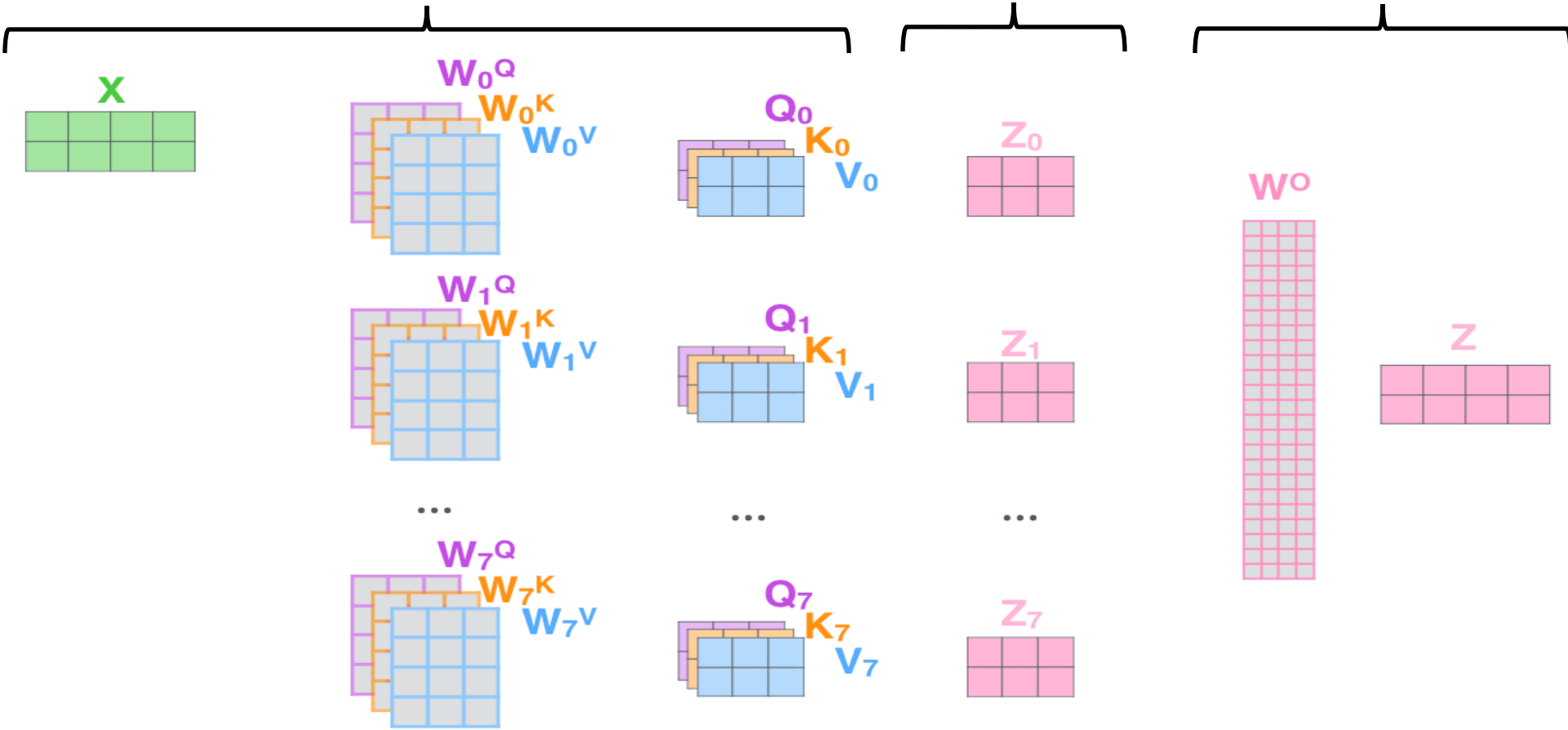


Multi-head Attention

1) Create query, key, and value vectors for all attentions heads

2) Compute new input representations

3) Condense all representations into a single representation by concatenating z-s and multiplying by a weight matrix



Trained Multi-head Attention Examples

Figure shows two columns of attention weights for the first two attention heads

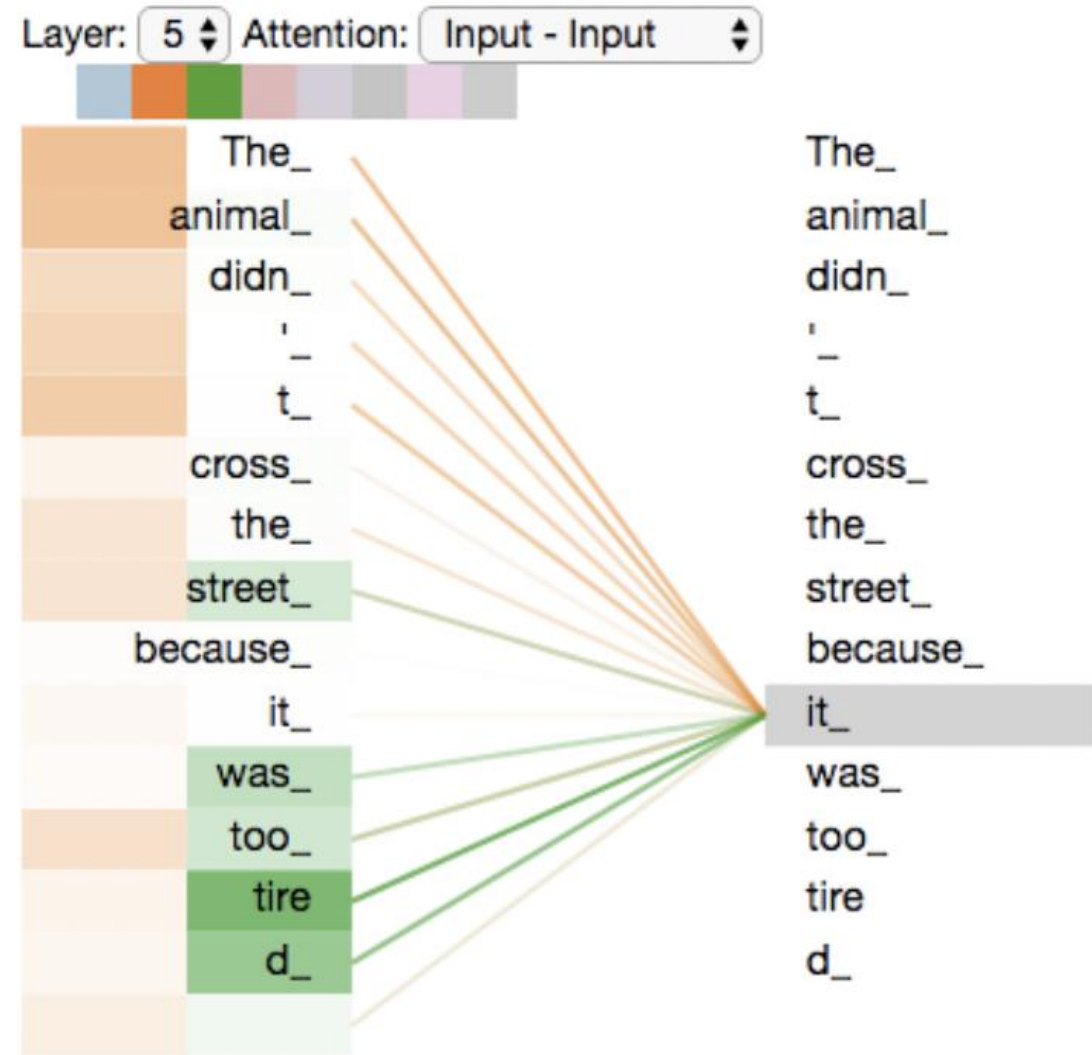
- Darker values signify larger attention scores

What does “it” focus on most in the first attention head?

- The animal (e.g., represents what is “it”)

What does “it” focus on most in the second attention head?

- tired (e.g., represents how “it” feels); note, a tokenizer was used that separates “tire” and “d”

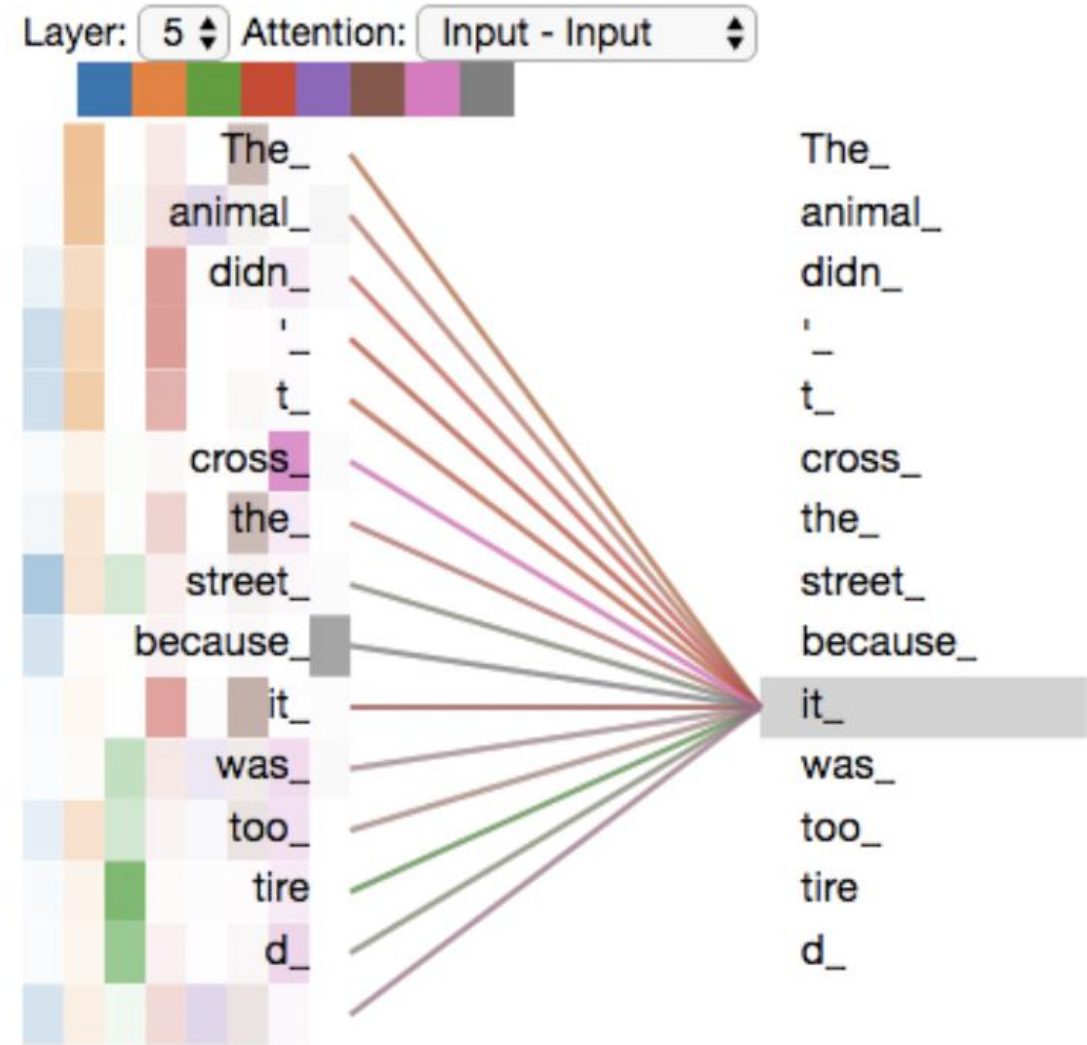


Trained Multi-head Attention Examples

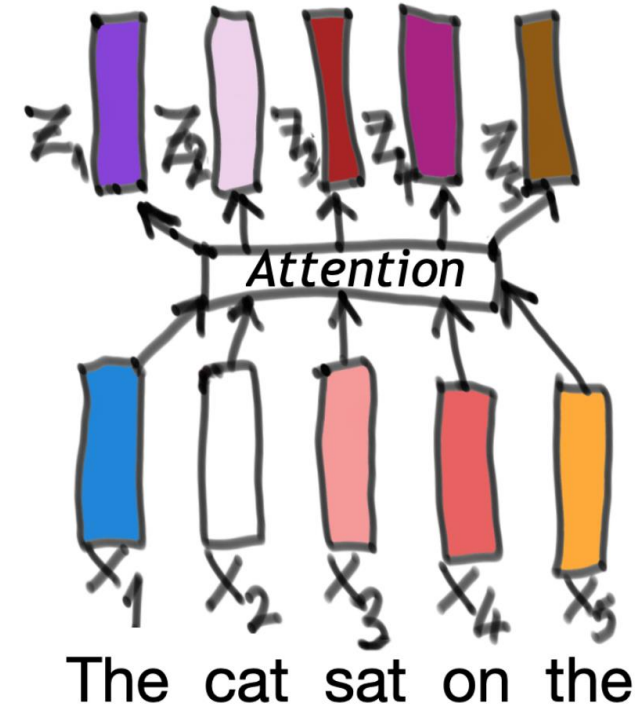
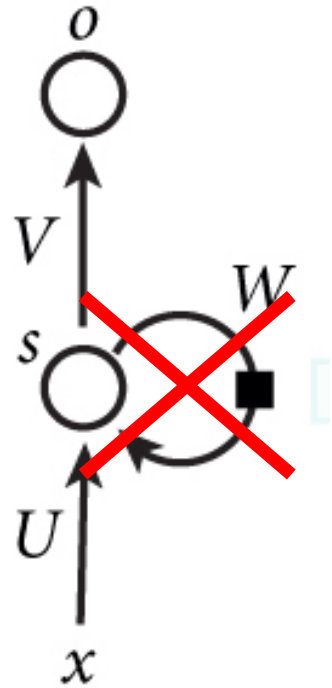
Figure shows five columns of attention weights for five attention heads

- Darker values signify larger attention scores

Attention weights may be hard to interpret



Self-Attention vs RNN: Propagates Information About Other Inputs **Without** Recurrent Units



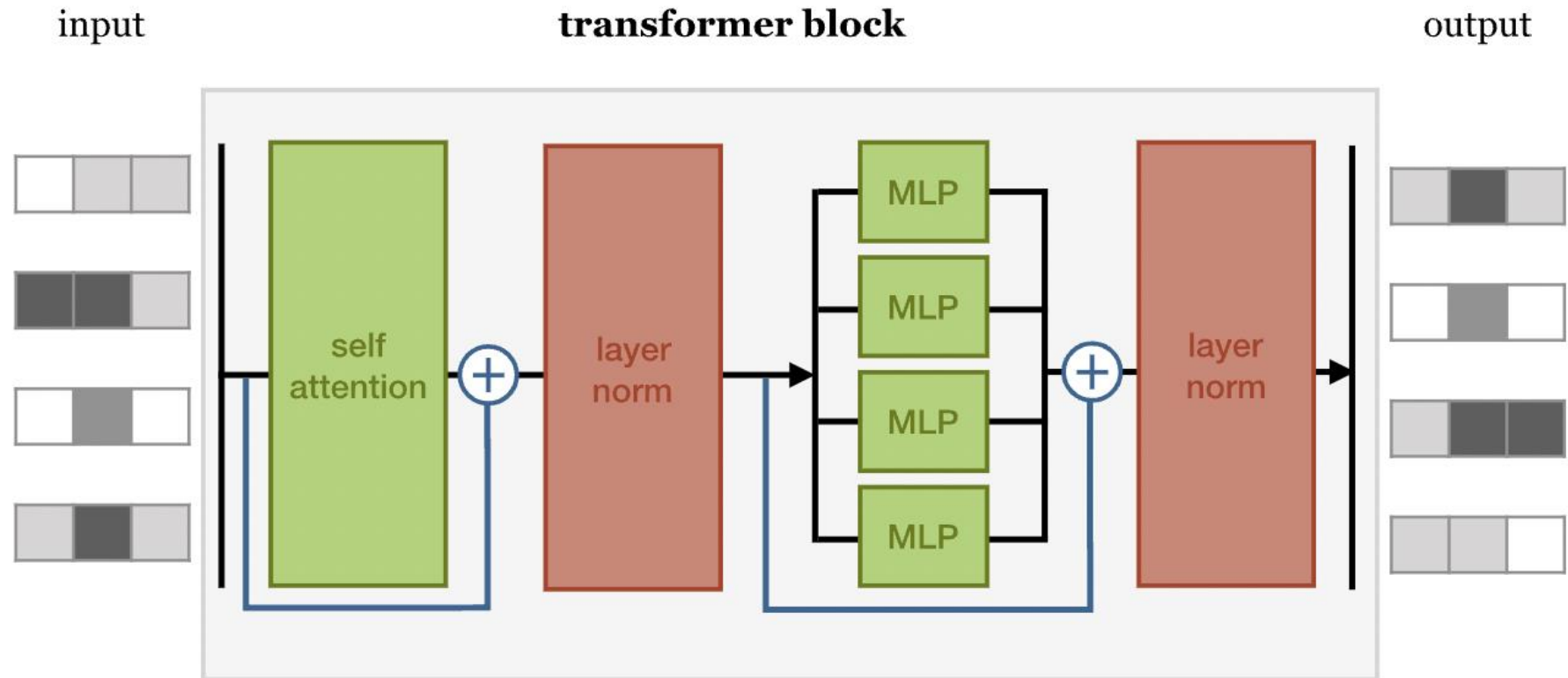
<http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>

<https://towardsdatascience.com/self-attention-5b95ea164f61>

Today's Topics

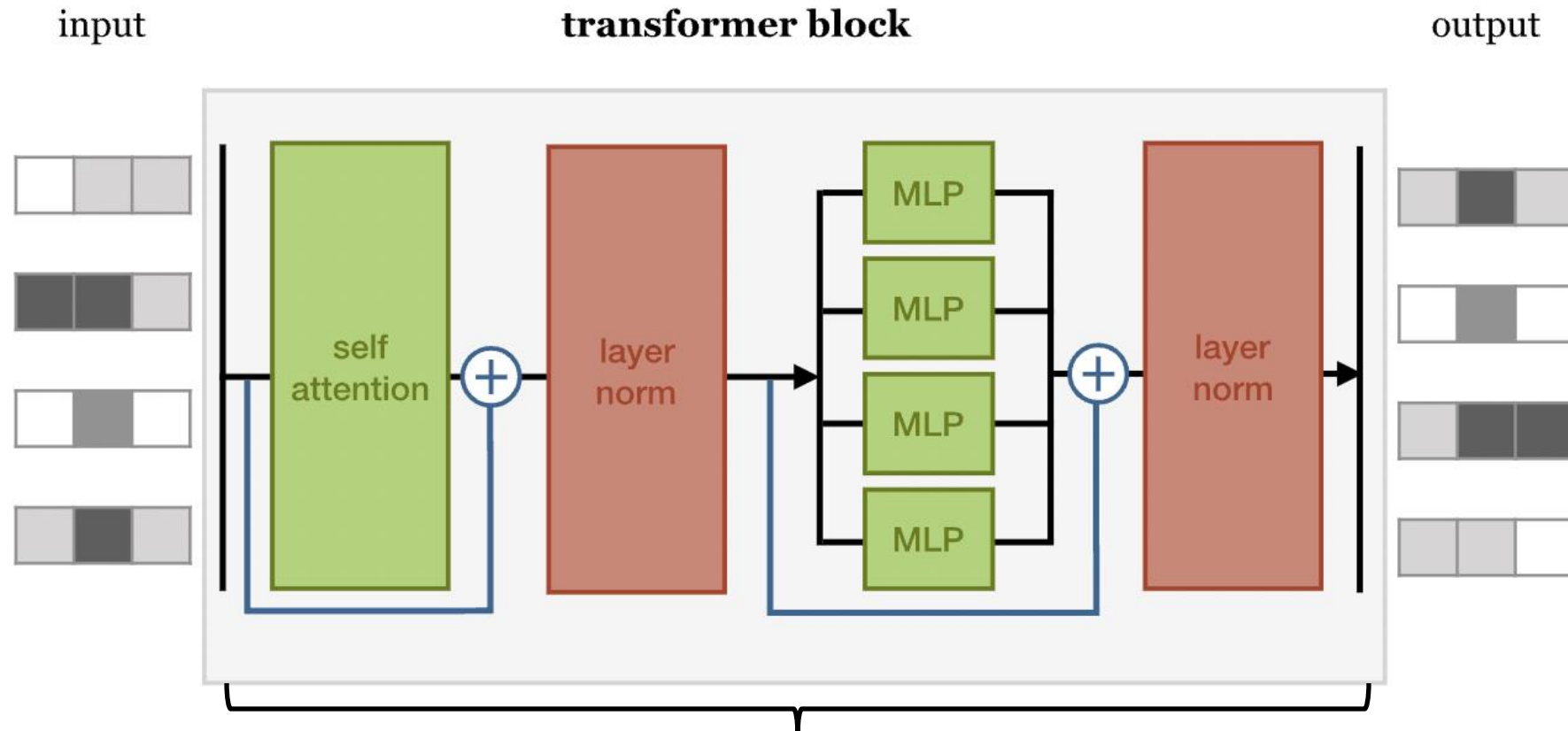
- Transformer overview
- Self-attention
- **Common transformer ingredients**
- Pioneering transformer: machine translation
- Programming tutorial

Typical Transformer Block



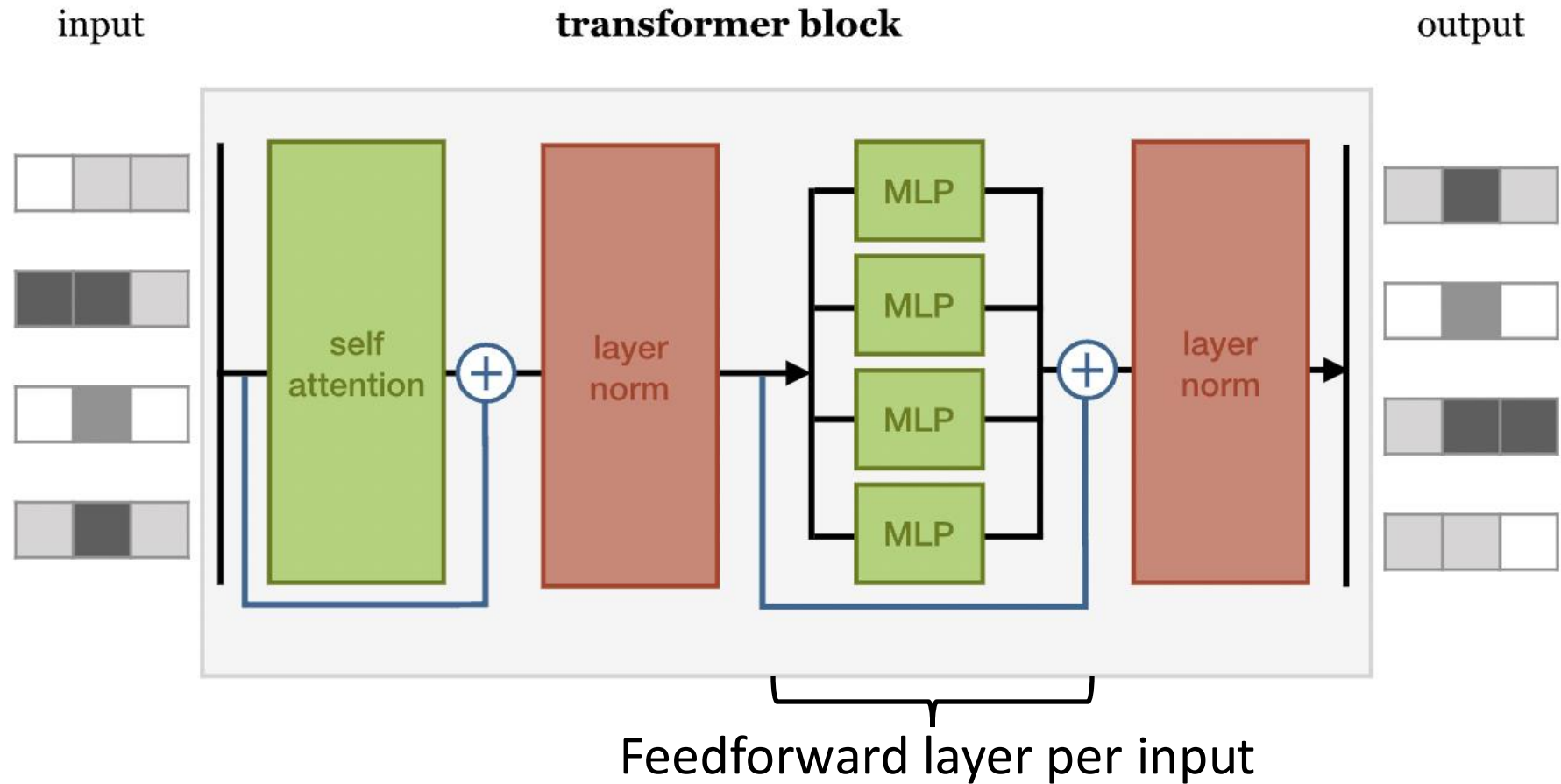
Architectures often chain together multiple transformer blocks, like that shown here

Typical Transformer Block

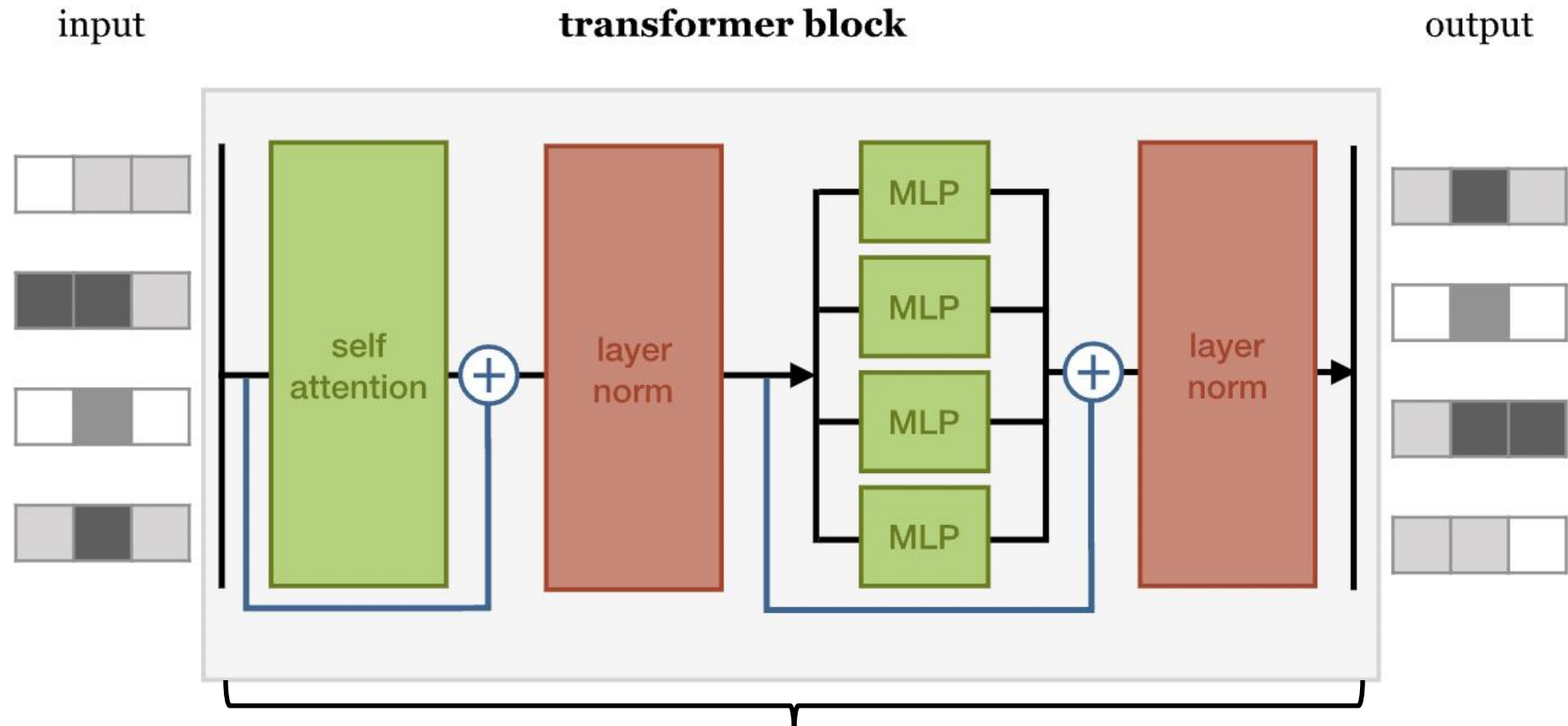


Layer normalization and residual connections improve training (i.e., faster and better results)

Typical Transformer Block



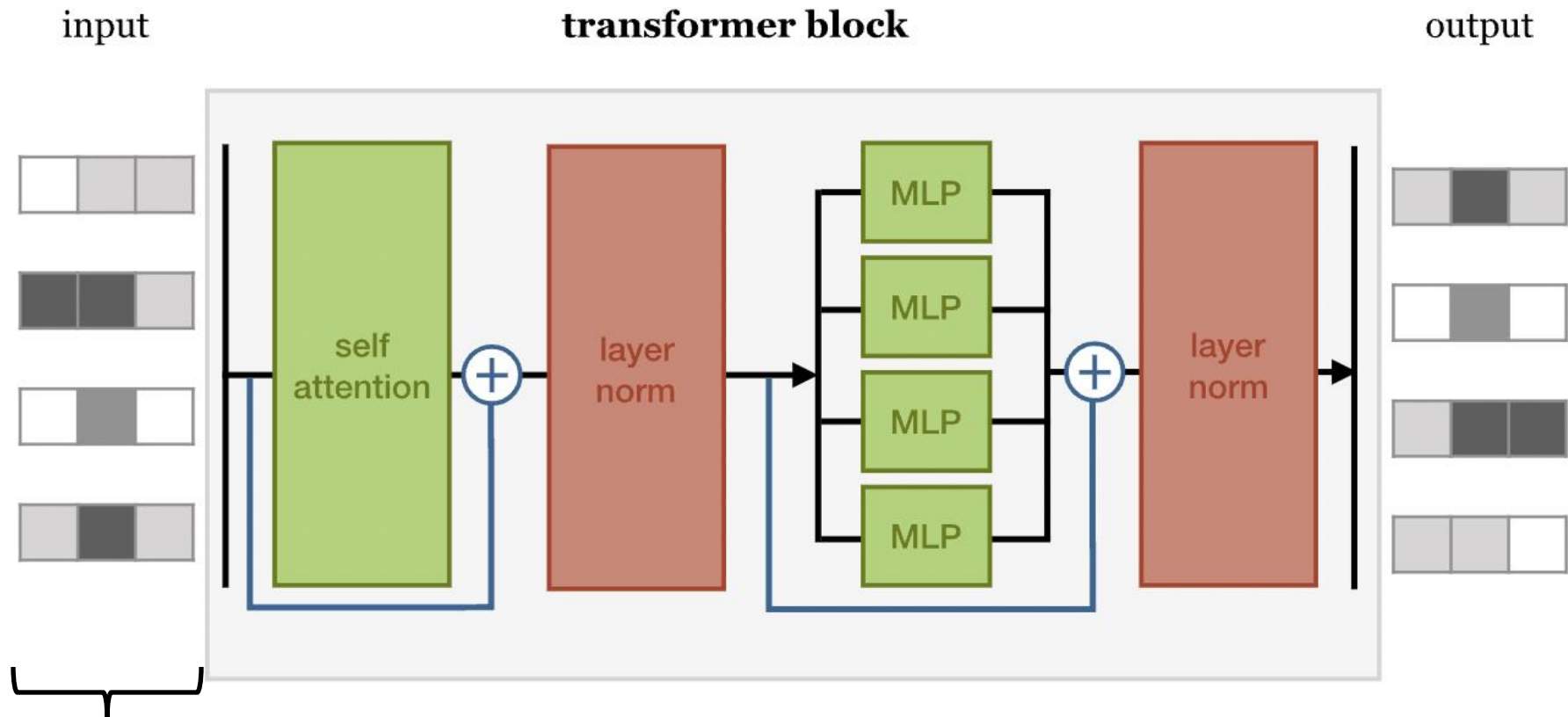
Typical Transformer Block



Where are non-linearities introduced in this block?

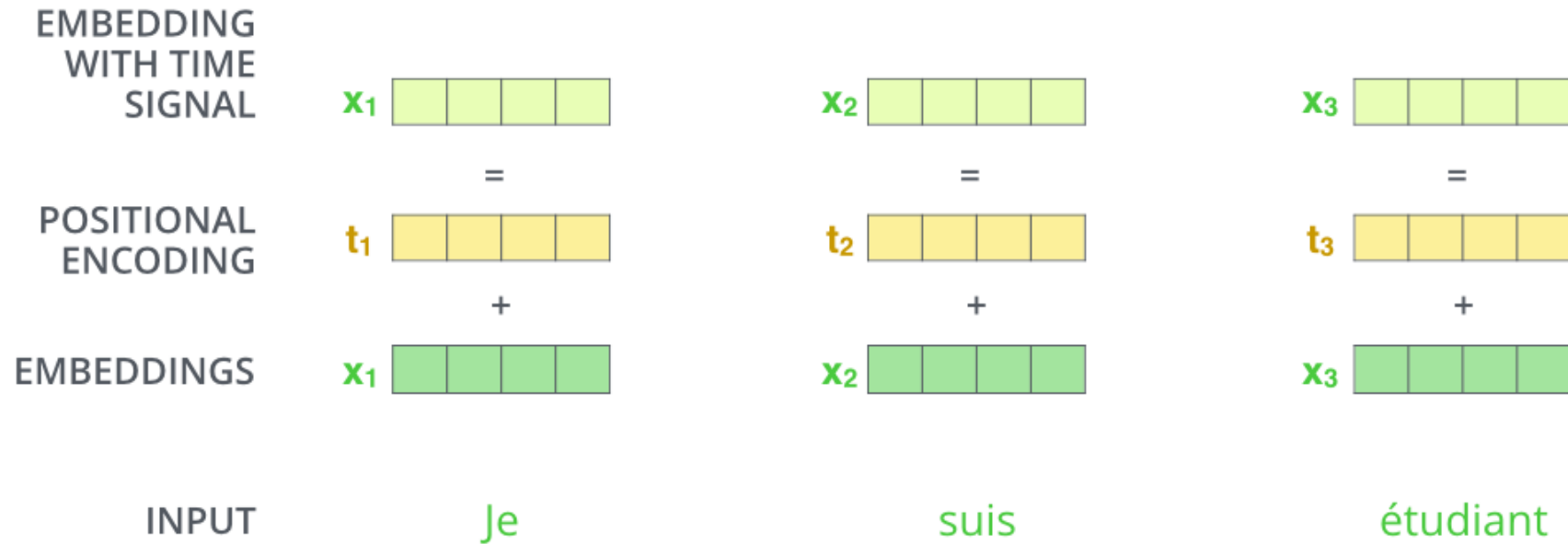
- self-attention's softmax, MLP's activation functions, layer norms

Challenge: Transformers Lack Sensitivity to the Order of the Input Tokens



Input: a *set* and so shuffling order of input tokens results yields same outputs except in the same shuffled order (i.e. self-attention is *permutation equivariant*)

Solution: Add Position as Input to Transformer



- Options:

- **Position embeddings:** created by training with sequences of every length during training
- **Position encodings:** a function mapping positions to vectors that the network learns to interpret (enables generalization to lengths not observed during training)

Today's Topics

- Transformer overview
- Self-attention
- Common transformer ingredients
- **Pioneering transformer: machine translation**
- Programming tutorial

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

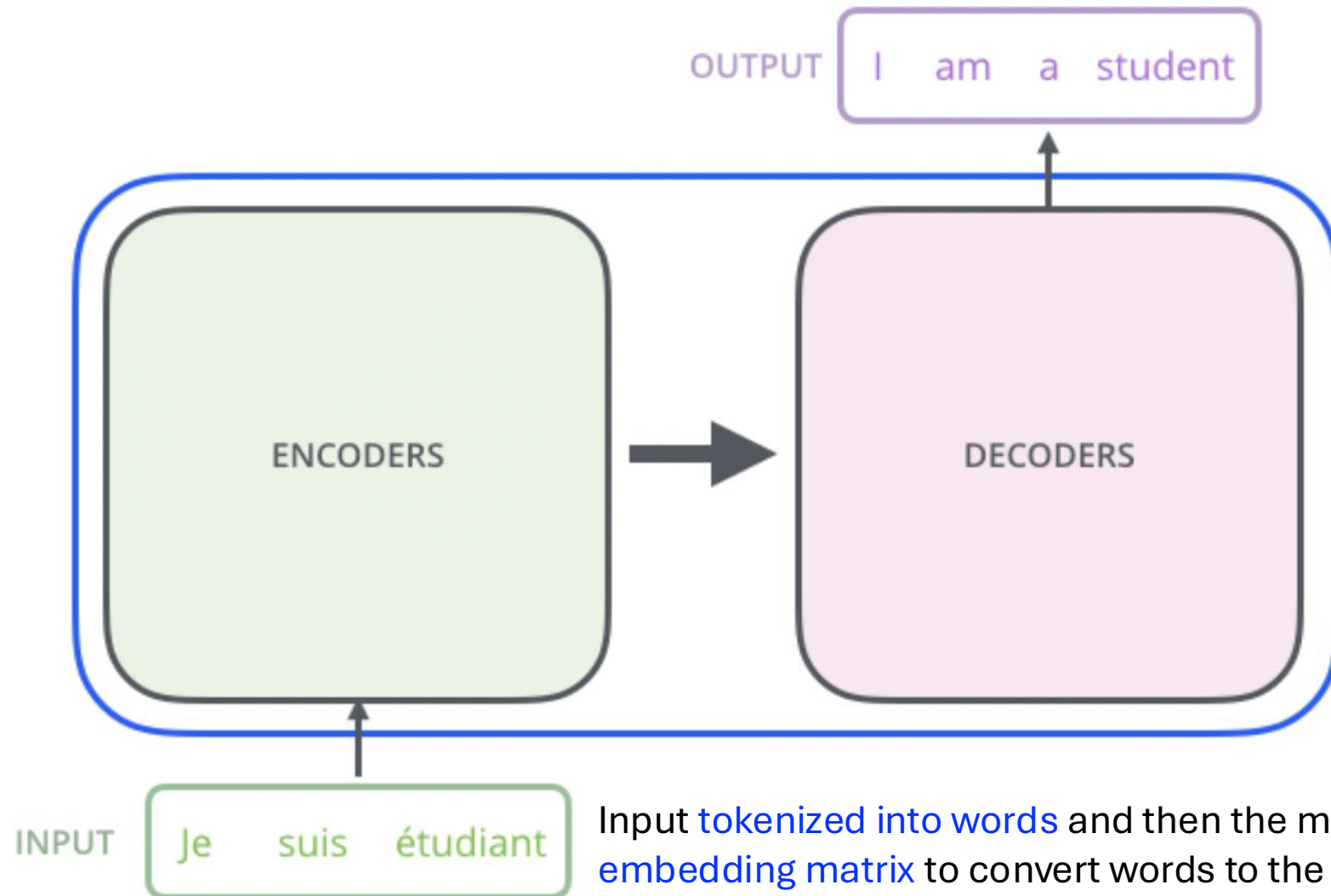
Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

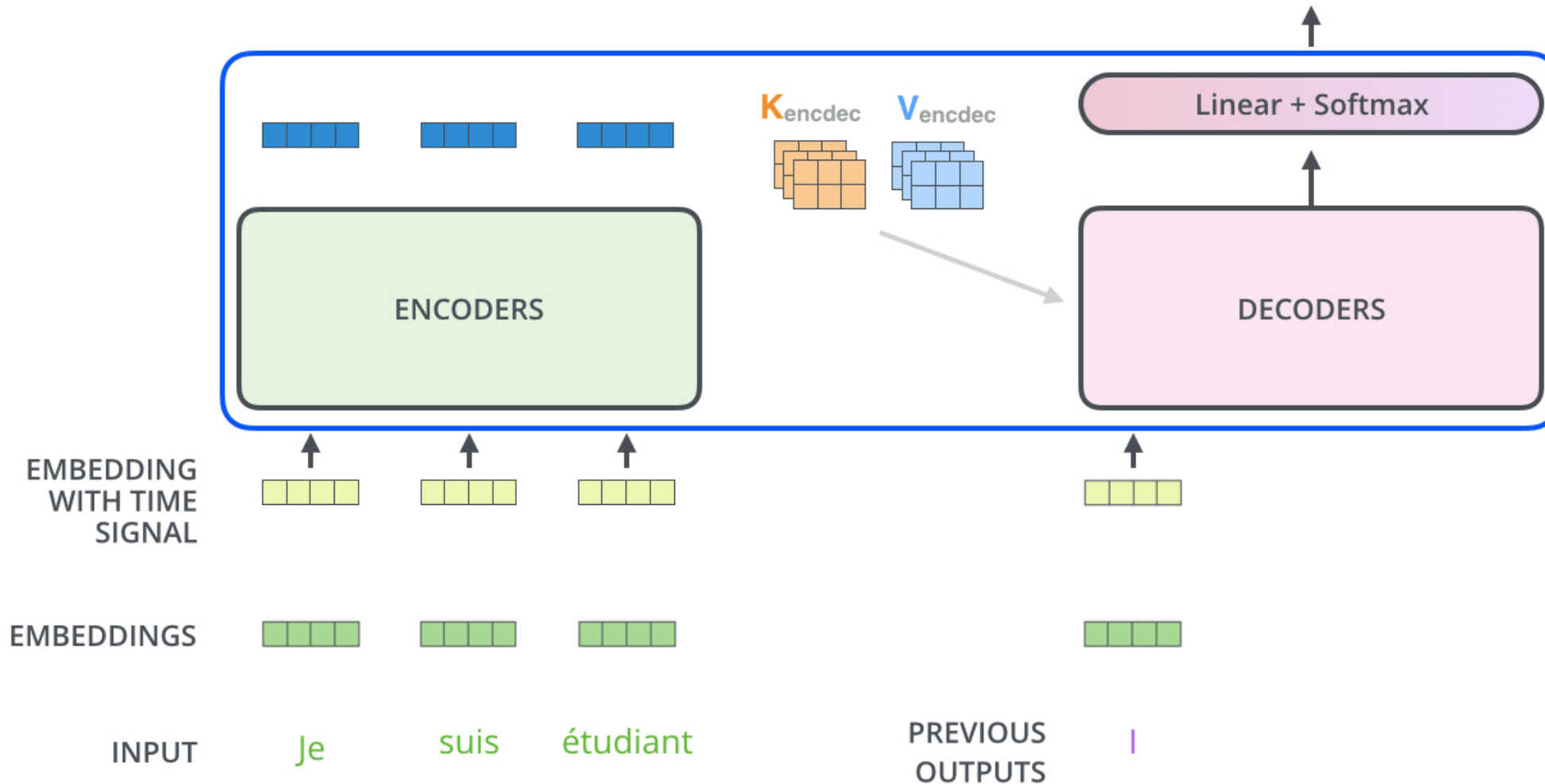
Target Application: Machine Translation



Example: Autoregressive Model

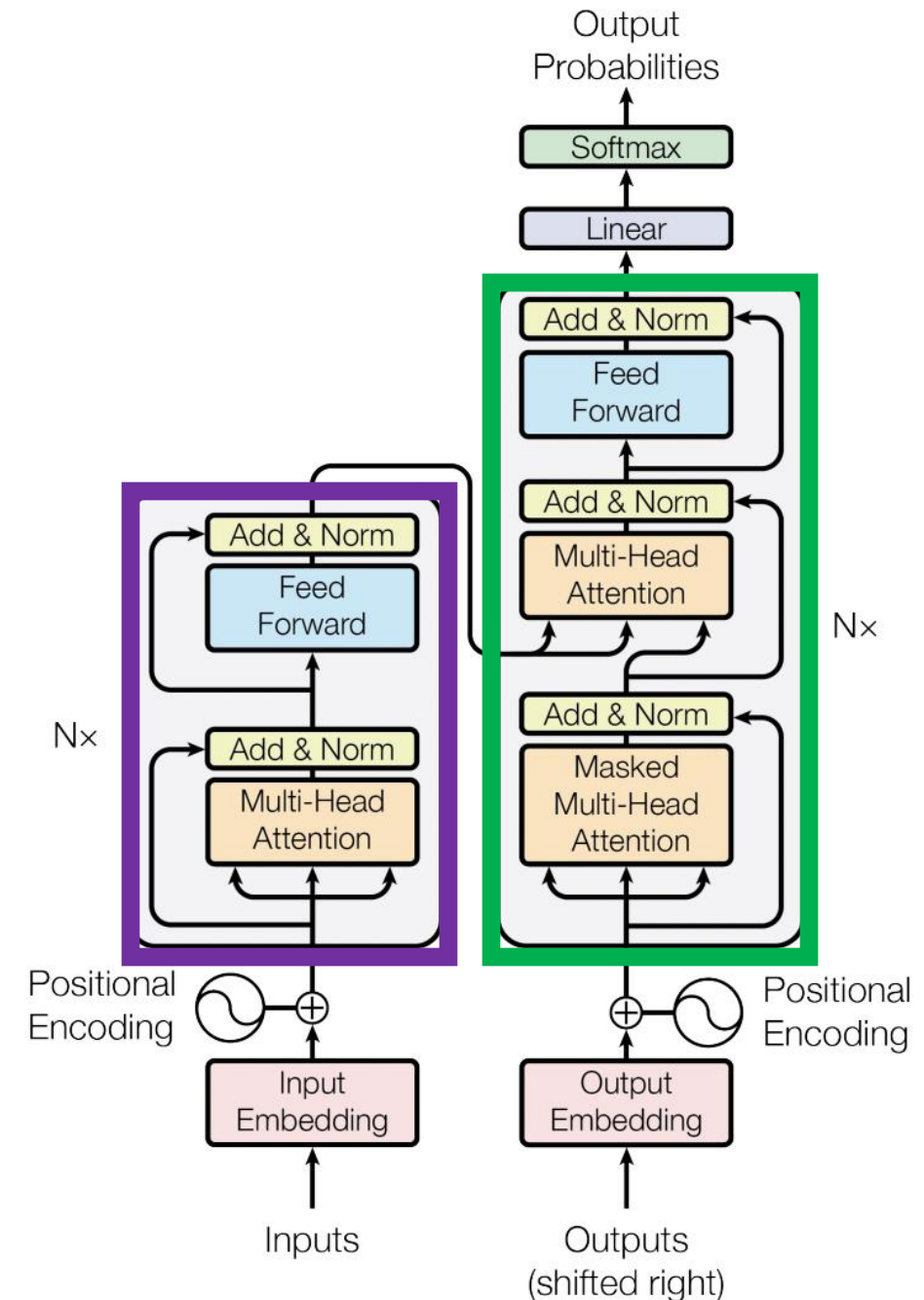
Decoding time step: 1 2 3 4 5 6

OUTPUT |



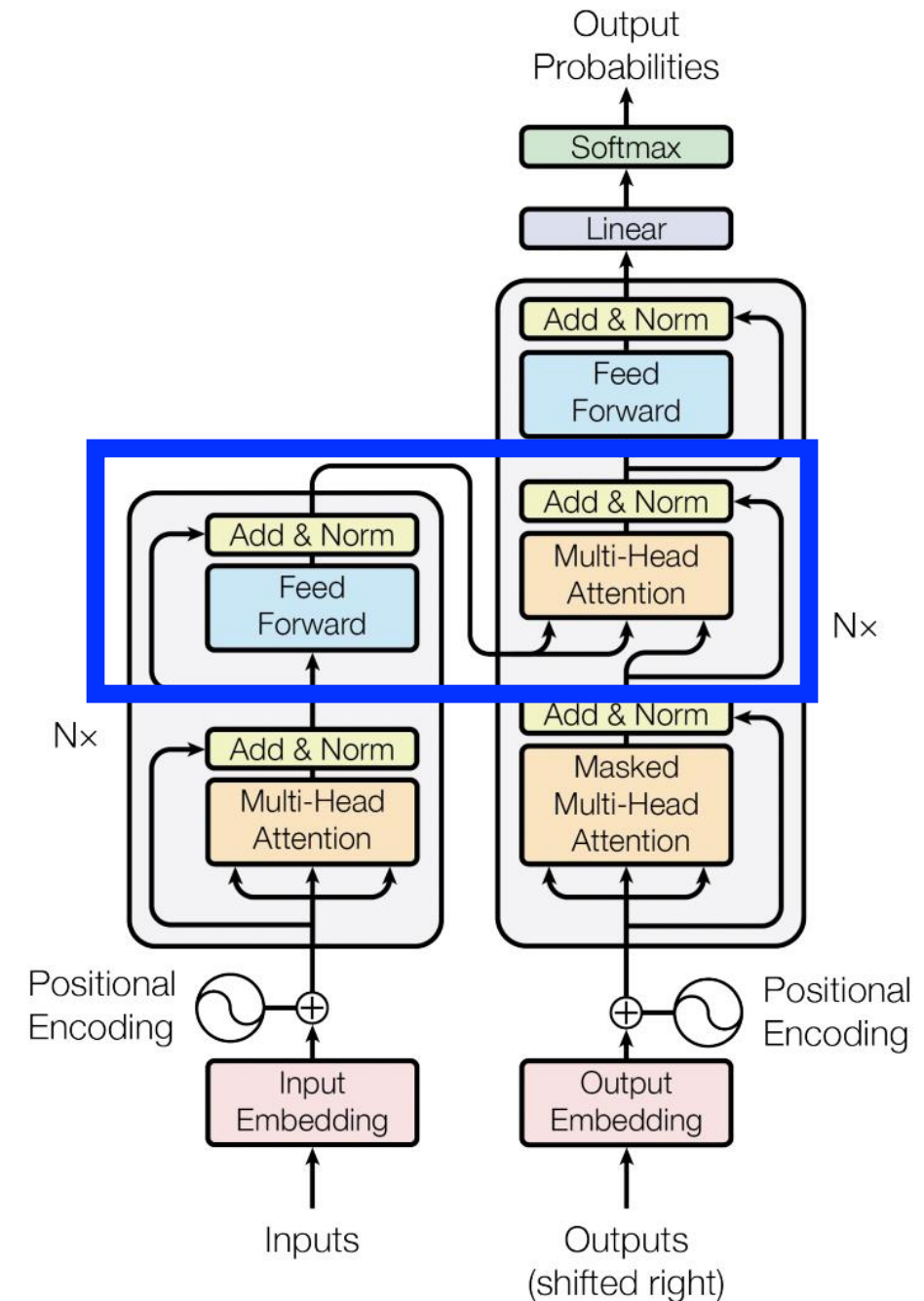
Architecture

- Key Ingredient: **Self-Attention**
 - Used in both the **encoder** (provides context for translation) and **decoder** (translates)
- Other ingredients
 - Positional encoding
 - Layer normalization
 - Residual connections
 - Feed forward layers
- $N_x = 6$ chained blocks (**encoder & decoder**)



Architecture

Decoder can attend to all inputs via **cross-attention** (i.e., **keys** and **values** come from the **encoder** and **query** comes from the **decoder**)

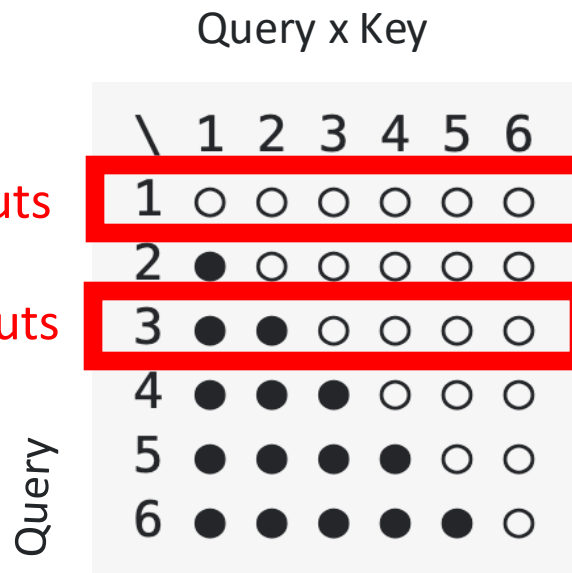


Architecture

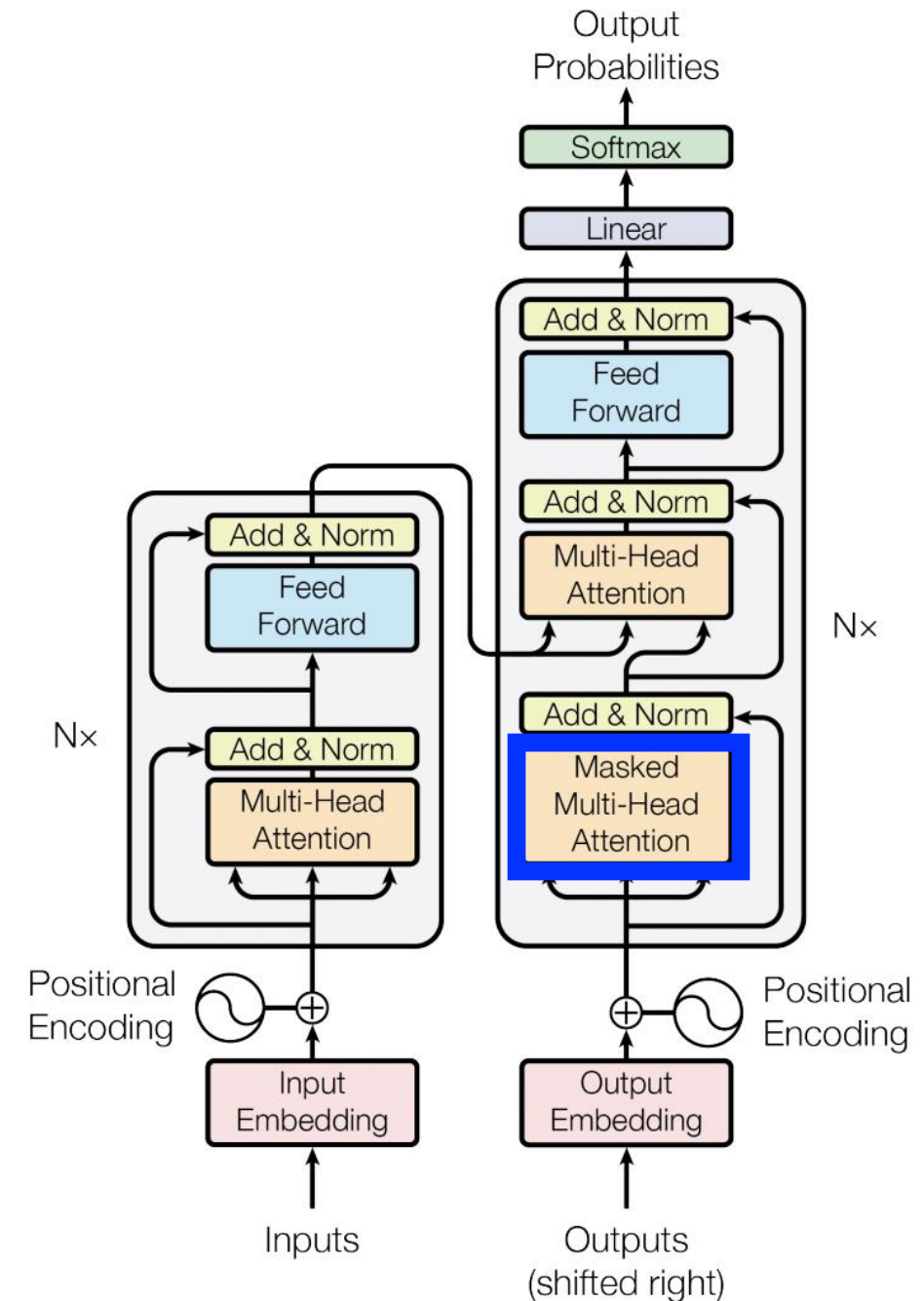
Masking so decoder ONLY sees earlier predictions
(i.e., masked values set to $-\infty$ so softmax output is 0)

e.g., at start, no previous inputs

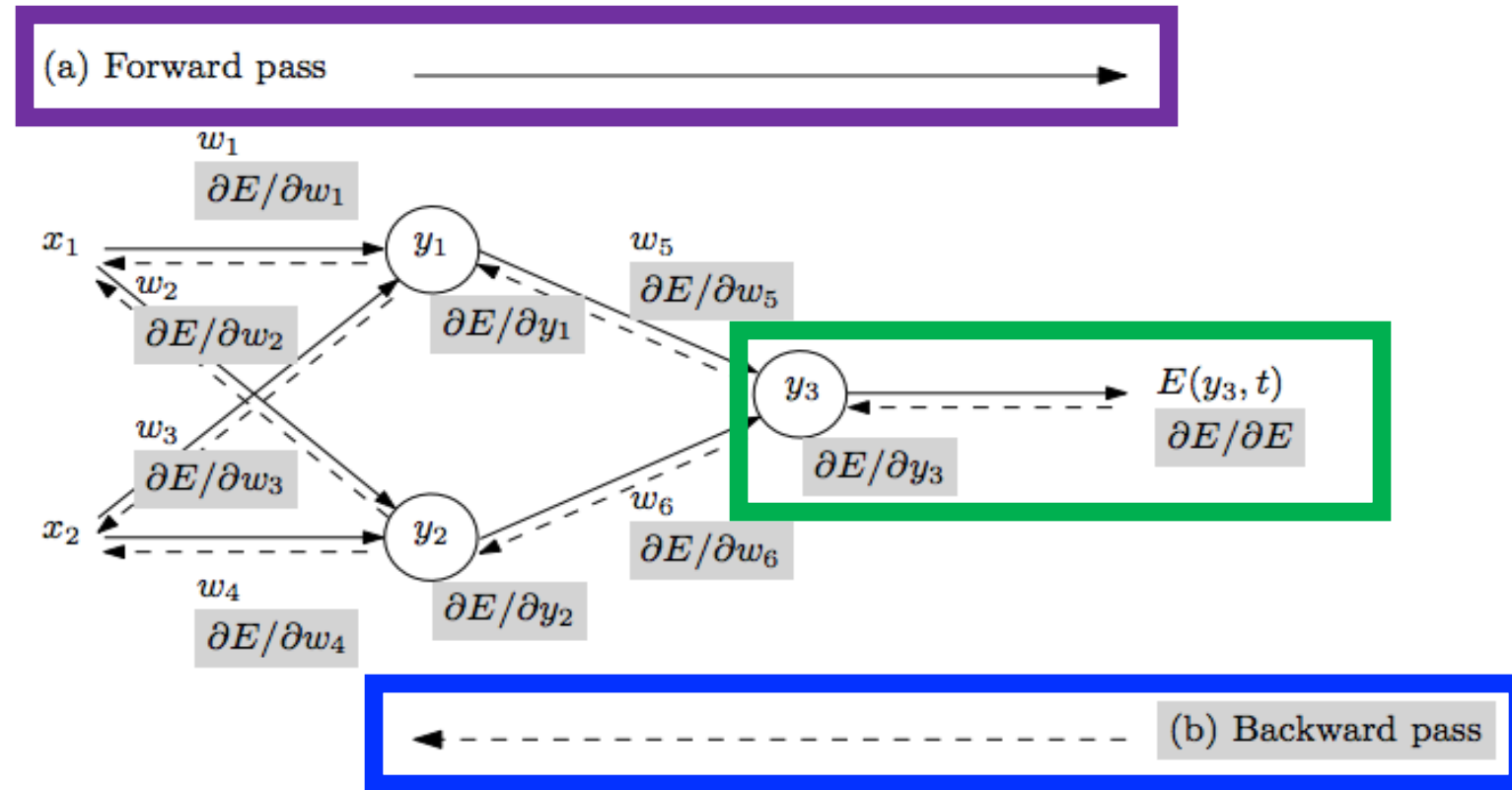
e.g., at 3rd step, two previous inputs



<https://stackoverflow.com/question/s/64799622/how-is-the-gpts-masked-self-attention-is-utilized-on-fine-tuning-inference>



Training Procedure: 3.5 days on 8 NVIDIA P100s



- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make predictions
 2. **Error quantification:** measure error of the model's predictions on training data using a loss function
 3. **Backward pass:** calculate gradients to determine how each model parameter contributed to model error
 4. Update each parameter using calculated gradients

For **two tested datasets**, achieved **state-of-the-art performance**:

- (1) English-German, with ~ 4.5 million sentence pairs (byte-pair encoded)
- (2) English-French, with 36M sentences

Today's Topics

- Transformer overview
- Self-attention
- Common transformer ingredients
- Pioneering transformer: machine translation
- **Programming tutorial**

Today's Topics

- Transformer overview
- Self-attention
- Common transformer ingredients
- Pioneering transformer: machine translation
- Programming tutorial



The End