

Neural Networks and Deep Learning

Attention



A digression into human attention

ATTENTION MECHANISMS

Human attention

“... the amount of information coming down the optic nerve - estimated to be in the range of **$10^8 \sim 10^9$ bits per second** - far exceeds what the brain is capable of fully processing and assimilating into conscious experience ...”

C. Koch, 1982

Restricting higher acuity vision to a small region of the retina and shifting the processing focus from one location to another (aka “attention”) is the solution nature has devised to cope in a serial fashion with the vast amount of visible information

The same happens with every sensory input (e.g. audio source separation)

Goal driven (volitional, top-down)

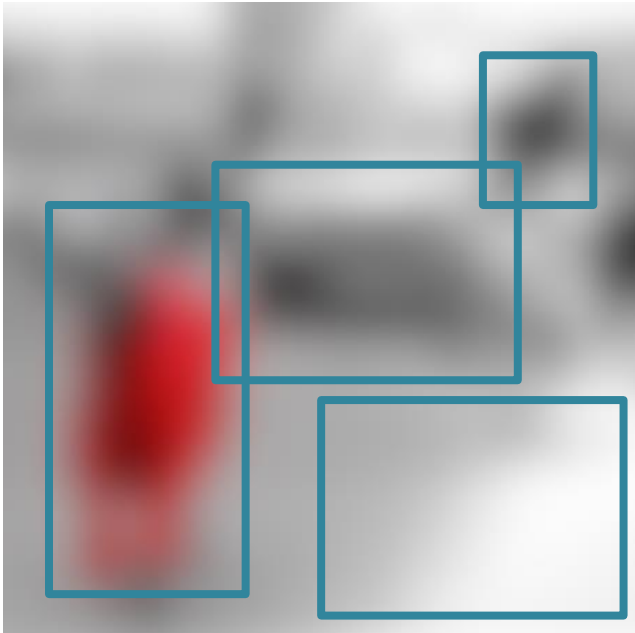


Stimuli driven (non-volitional, bottom-up)



© R. Rensink, University of British Columbia

Where would you look at?

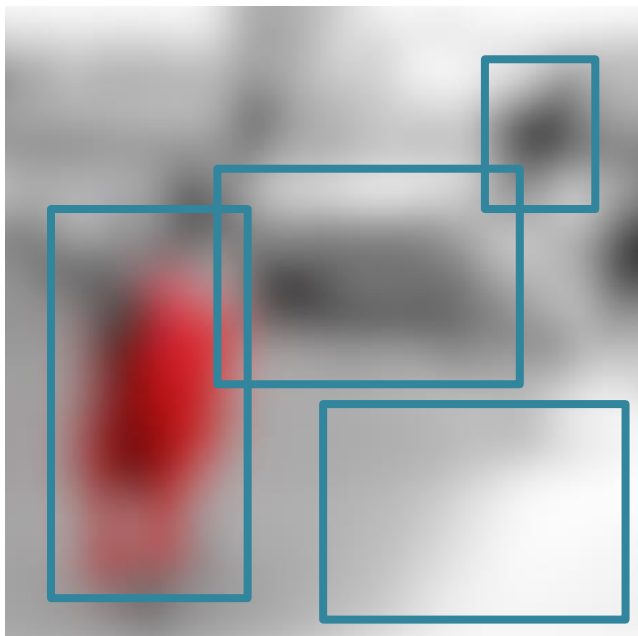


This is called a free-viewing task. There is no explicit query (task). Still, you know that there are some parts which seem more important (salient) than others.

Attention is bottom up: the input signal is all you have to decide where it pays off to look at

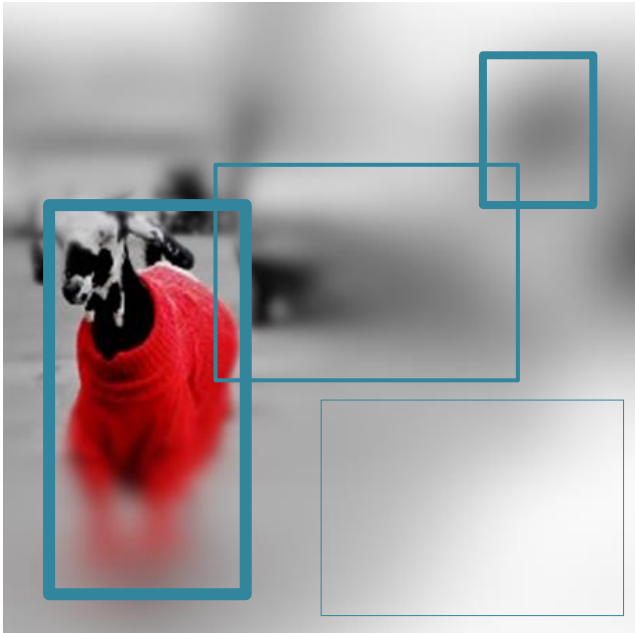
Where would you look at?

Q: *“Where is the stop sign?”*



Where would you look at?

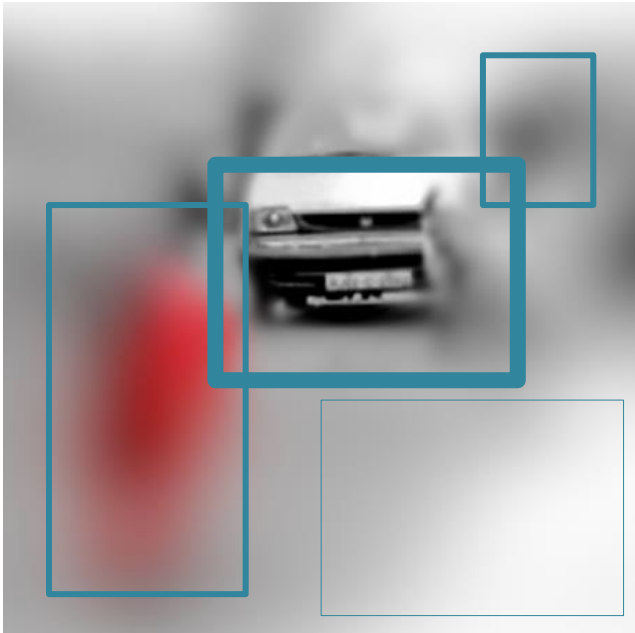
Q: *“Where is the stop sign?”*



Here you have a task, where you look at depends on the task (query) as well as information from the signal (keys)

Where would you look at?

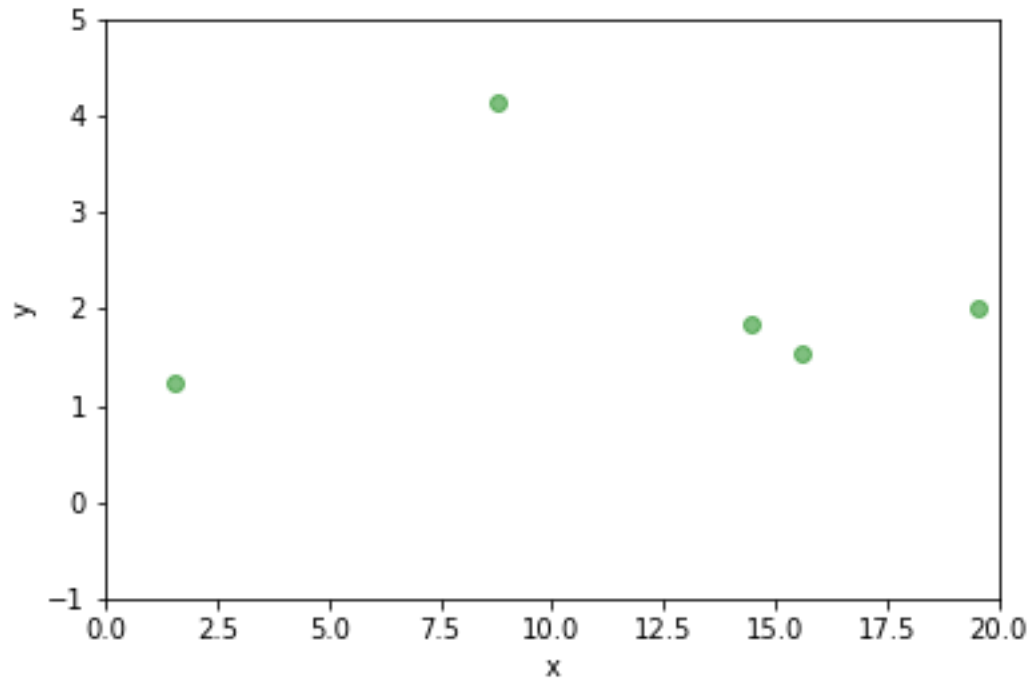
Q: *“Where is the car?”*



Here you have a task, where you look at depends on the task (query) as well as information from the signal (keys)

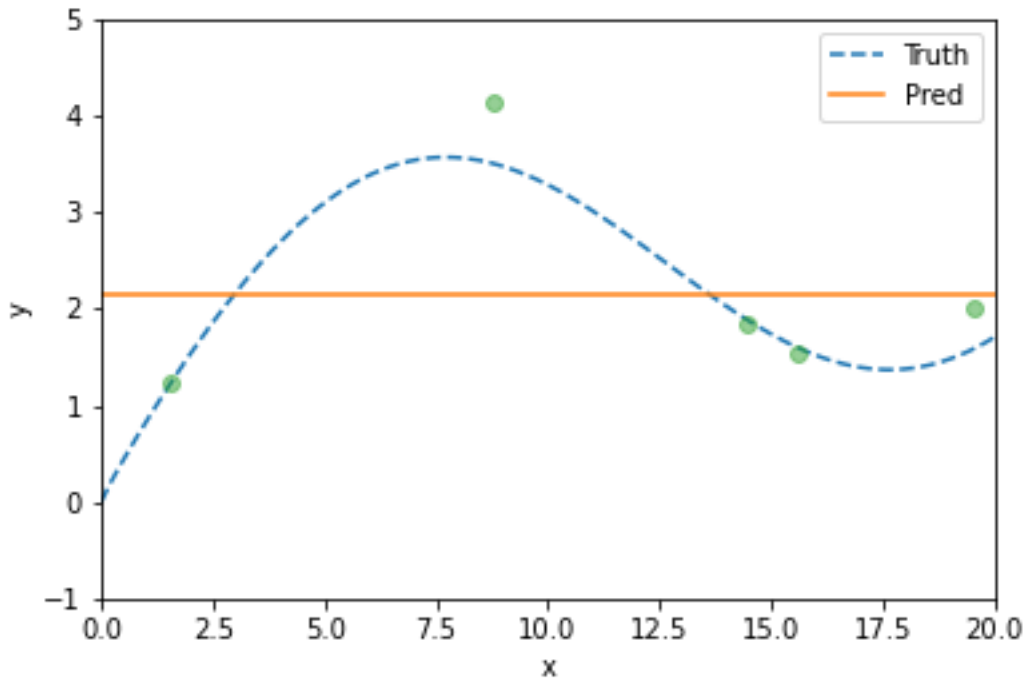
ATTENTION IN MACHINE LEARNING

Solving the regression problem



Data: $\{x_1, x_2, \dots, x_m\}$
Labels: $\{y_1, y_2, \dots, y_m\}$

Solving the regression problem

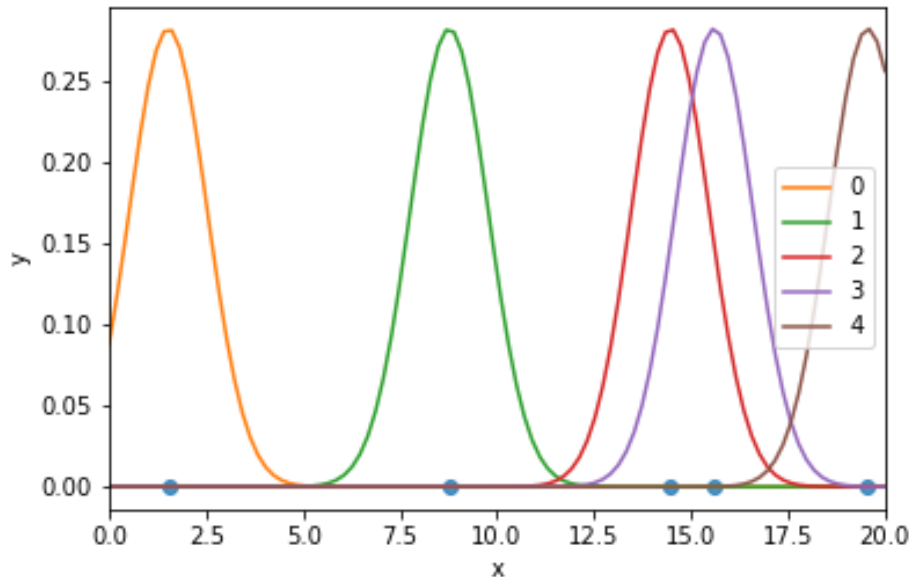


Lacking any other knowledge, we could use just the average:

$$f(x) = \frac{1}{m} \sum_{i=1}^m y_i$$

Data: $\{x_1, x_2, \dots, x_m\}$
Labels: $\{y_1, y_2, \dots, y_m\}$

Nadaraya-Watson Kernel Regression (1964)



A better idea would be to weigh the labels, according to their location

The weights are proportional to some similarity function, a “kernel”. For example, a Gaussian.

$$f(x) = \sum_{i=1}^m \alpha(x, x_i) y_i$$

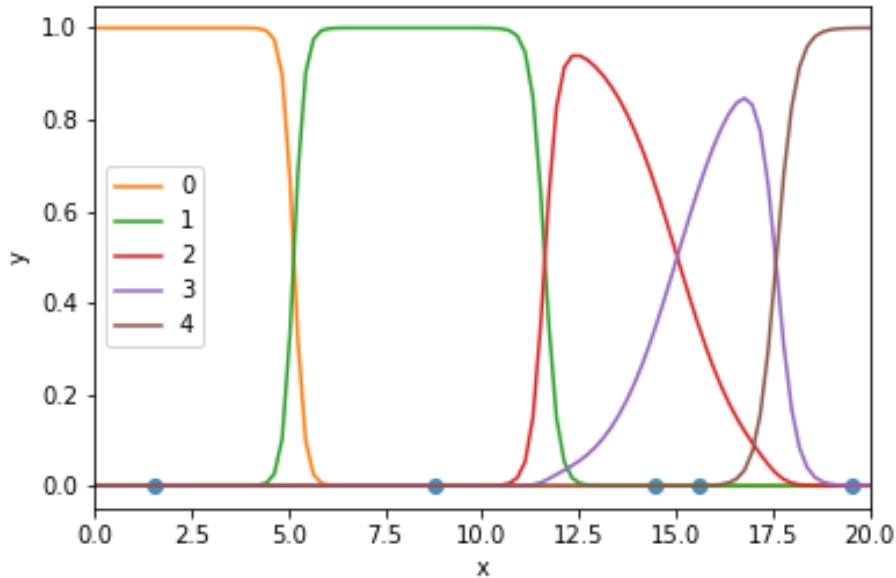
Diagram illustrating the components of the Nadaraya-Watson Kernel Regression formula:

- Key**: Points to the input x in the kernel function $\alpha(x, x_i)$.
- Query**: Points to the input x in the kernel function $\alpha(x, x_i)$.
- Value**: Points to the output y_i in the formula.
- kernel**: Points to the kernel function $\alpha(x, x_i)$, which is proportional to $k(x_i, x)$.

The relationship between the kernel and the weight is given by:

$$\alpha(x, x_i) \propto k(x_i, x)$$

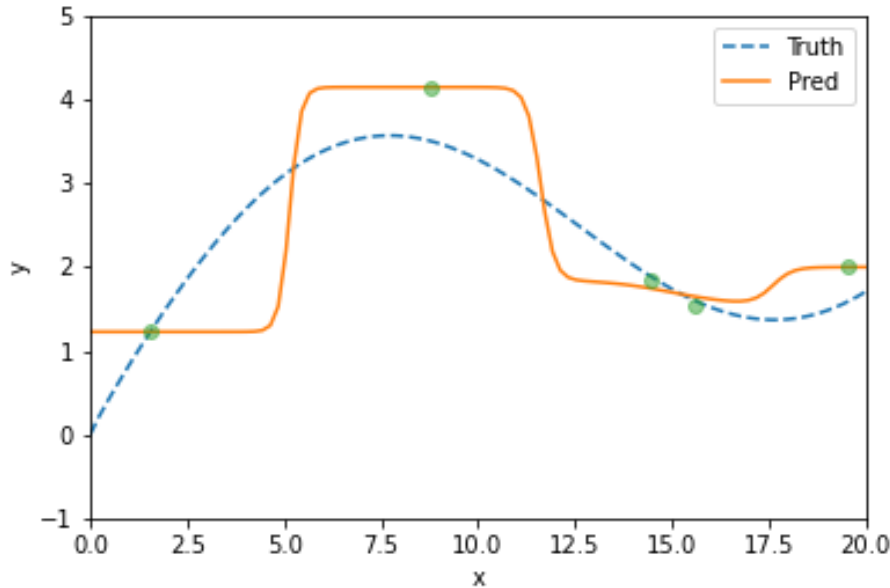
Nadaraya-Watson Kernel Regression (1964)



Weights should also be normalised

$$\alpha(x, x_i) = \frac{k(x, x_i)}{\sum_j k(x, x_j)}$$

Nadaraya-Watson Kernel Regression (1964)

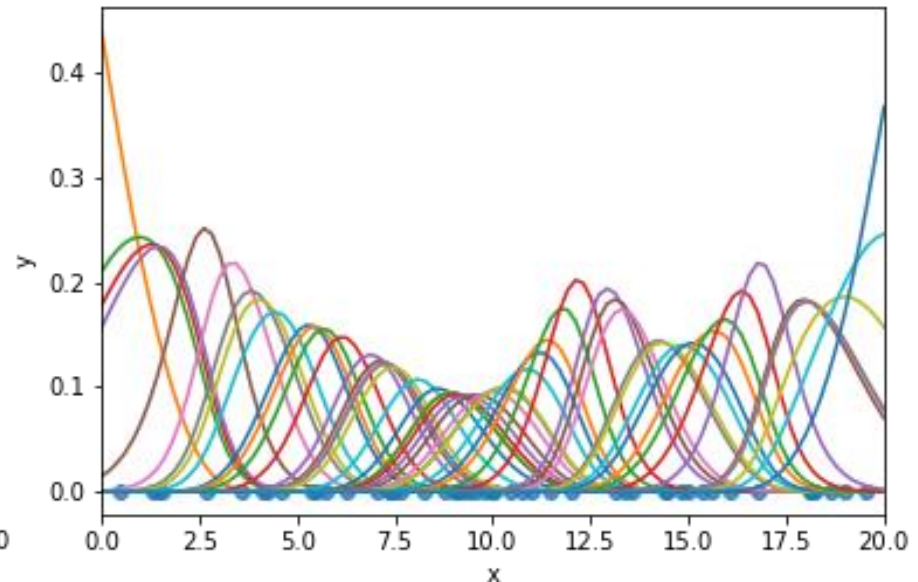
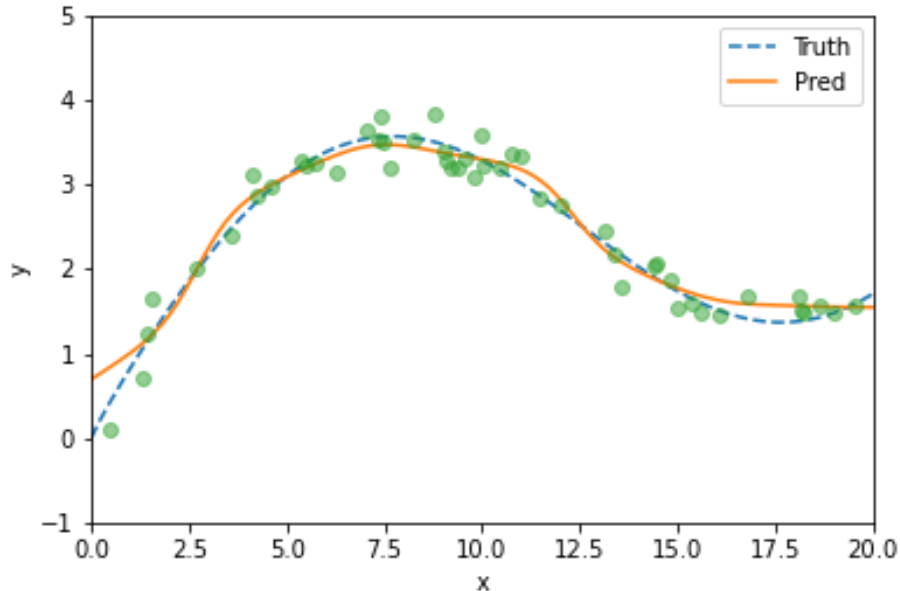


$$\begin{aligned} f(x) &= \sum_{i=1}^m \frac{k(x, x_i)}{\sum_j k(x, x_j)} y_i \\ &= \sum_{i=1}^m \frac{\exp\left(-\frac{1}{2}(x - x_i)^2\right)}{\sum_j \exp\left(-\frac{1}{2}(x - x_j)^2\right)} y_i \\ &= \sum_{i=1}^m \text{softmax}\left(-\frac{1}{2}(x - x_i)^2\right) y_i \end{aligned}$$

Using a Gaussian kernel with unit variance:

$$k(x, x_i) = N(x - x_i) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x - x_i)^2}{2}\right)$$

Nadaraya-Watson Kernel Regression (1964)

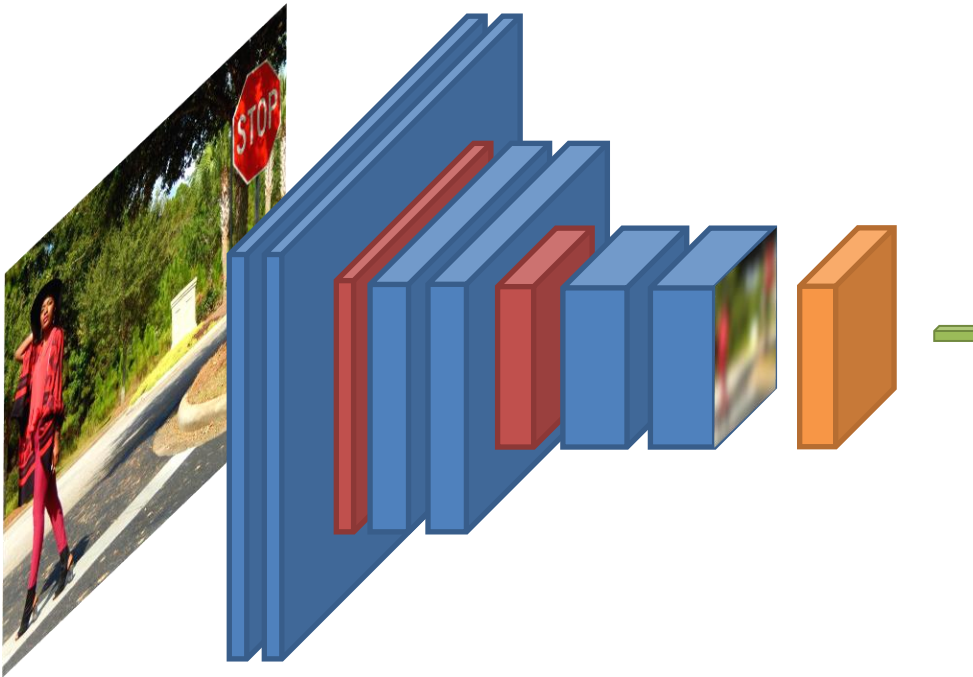


Consistency: Given enough data, this algorithm converges to the optimal solution

Simplicity: No free parameters – information is in the data, not in the weights

Computationally expensive: No model, we need to iterate through all the data every time we need to compute a new value

Average Pooling?



Global Average Pooling would be equivalent to the “flat line” model (average) we saw before.

What if we knew that we are interested about a particular area of the image?

Imitate Fixations?



What if we knew that we are interested about what is “around here”? We would then weight differently the cells that are closer to that location of interest

Imitate Fixations?



This is what humans do, by “fixating” in different places in the image.

But how do you decide **where** to fixate?

Imitate Fixations?



Where you fixate (how you deploy your attention) depends on the input signal (**bottom-up** attention) and on the task, e.g. *“Is there any stop sign on the road”* (**top-down** attention).²²

Imitate Fixations?



Where you fixate (how you deploy your attention) depends on the input signal (**bottom-up** attention) and on the task, e.g. *"Is there any stop sign on the road"* (**top-down** attention).²³

Imitate Fixations?



Where you fixate (how you deploy your attention) depends on the input signal (**bottom-up** attention) and on the task, e.g. *"Is there any stop sign on the road"* (**top-down** attention).²⁴

Towards deep learning

Could we substitute our average **pooling operations** with a weighted pooling like this?

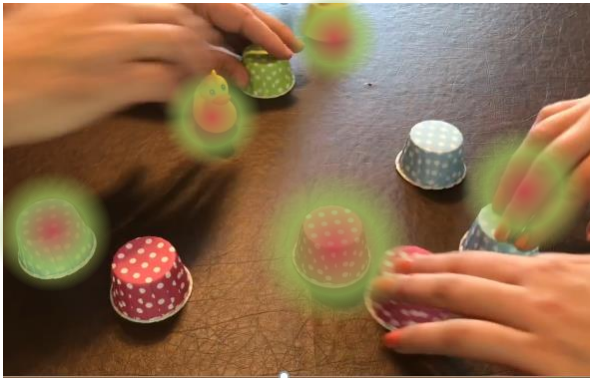
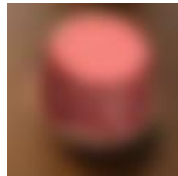
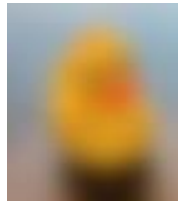
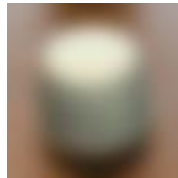
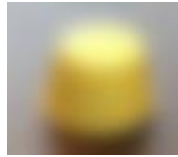
$$f(x) = \sum_{i=1}^m \frac{k(x, x_i)}{\sum_j k(x, x_j)} y_i$$

Could we learn better **kernels**?

Could we learn the right **queries**, **keys** and **values**?

QUERIES, KEYS AND VALUES

Sensory inputs



Non-volitional cues

pink, human

yellow, inanimate

blue, inanimate

yellow, animal

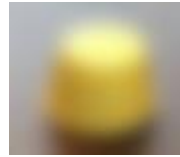
red, inanimate

Sensory inputs

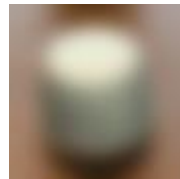
Non-volitional cues



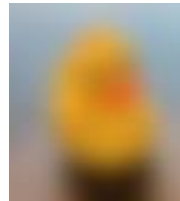
pink, human



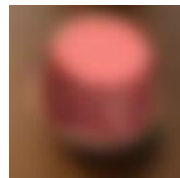
yellow, inanimate



blue, inanimate



yellow, animal



red, inanimate

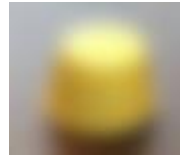
Volitional cue “Red cup”

Sensory inputs

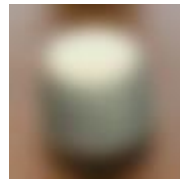
Non-volitional cues



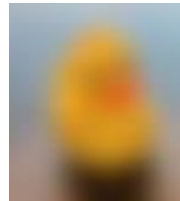
pink, human



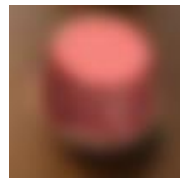
yellow, inanimate



blue, inanimate



yellow, animal



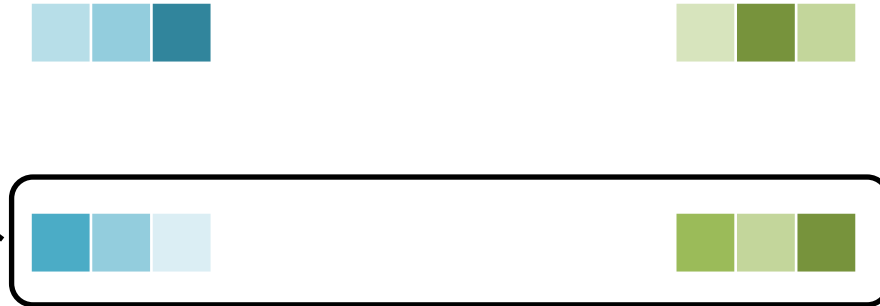
red, inanimate

Volitional cue "Duck"

Values
(Sensory inputs)

Keys
(Non-volitional cues)

token



Query
(Volitional cue)

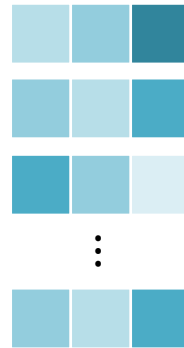
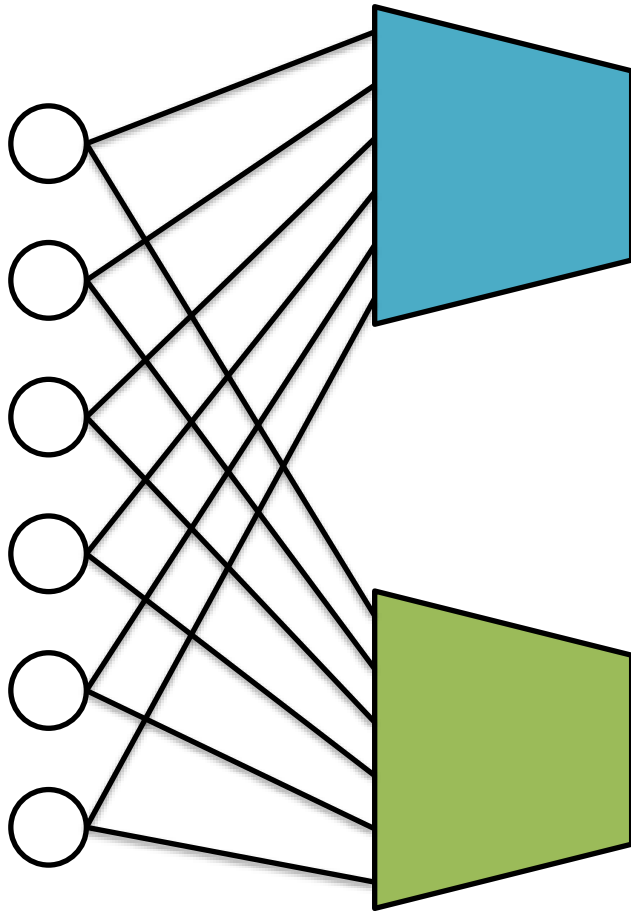


Note

The idea of queries, keys and values aims to build a **conceptual framework** through which we can abstract different ways of implementing “attention” mechanisms

In reality, you will see that queries, keys and values might not be different things at all (see self-attention for example... later in this lecture)

Where do values and keys come from?

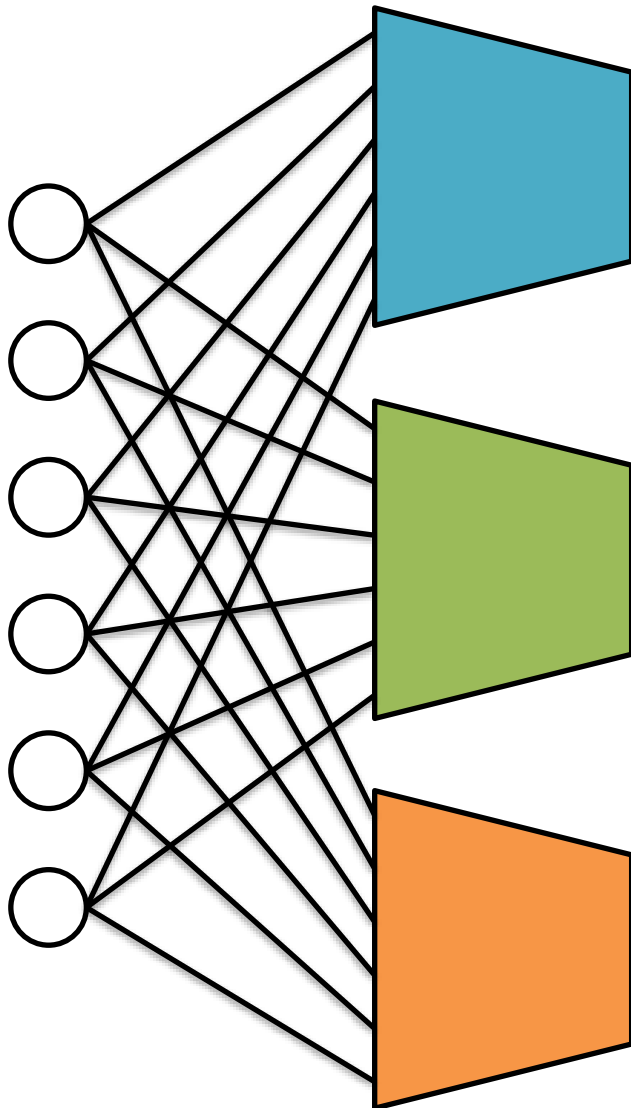


For the time being, you can think of keys and values as different representations calculated from the same data



We will see eventually, we that they can actually be the same (equal)

Where do queries come from?



⋮



⋮

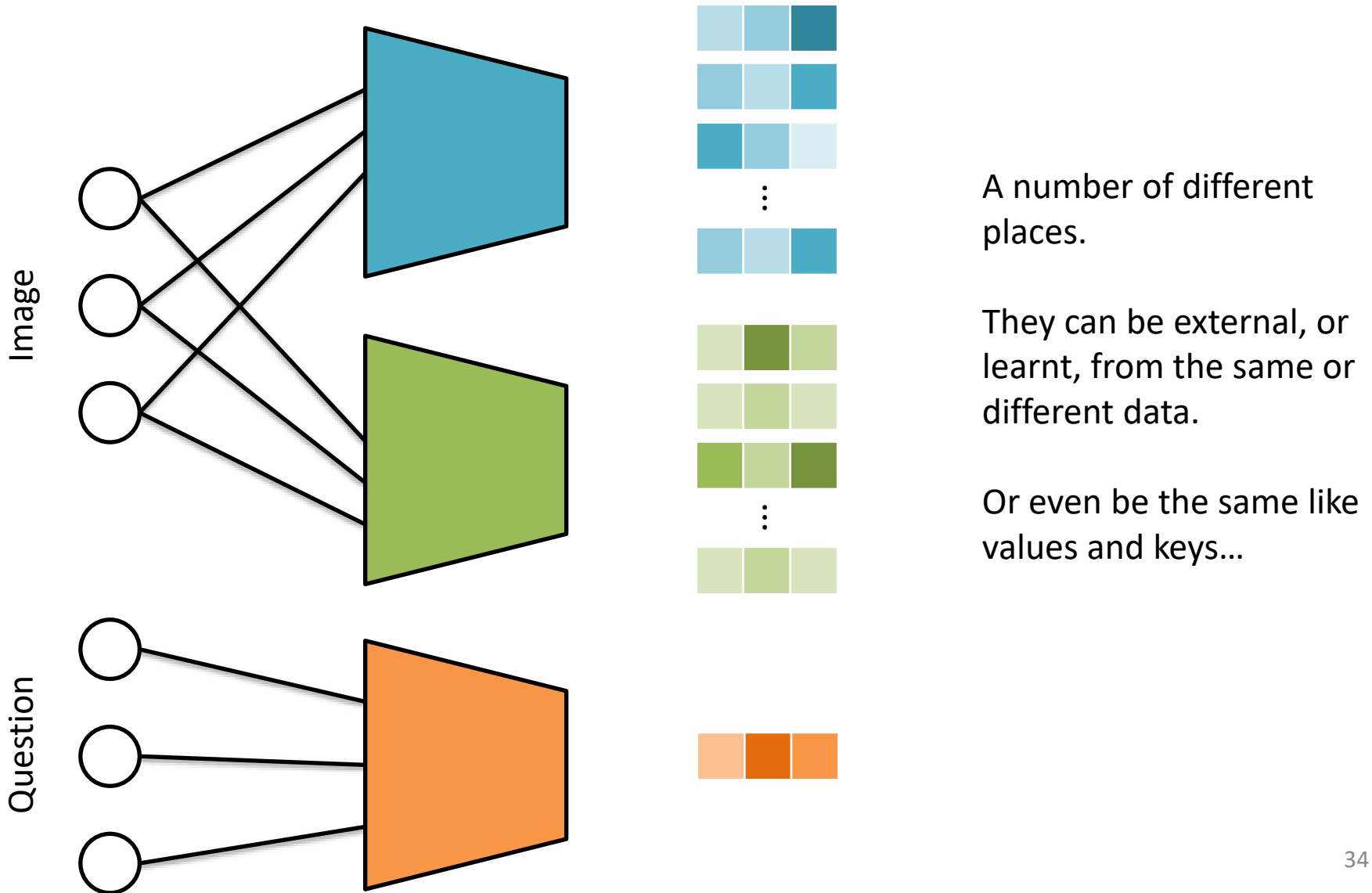


A number of different places

They can be external, or learnt, from the same or different data.

Or even be the same like values and keys...

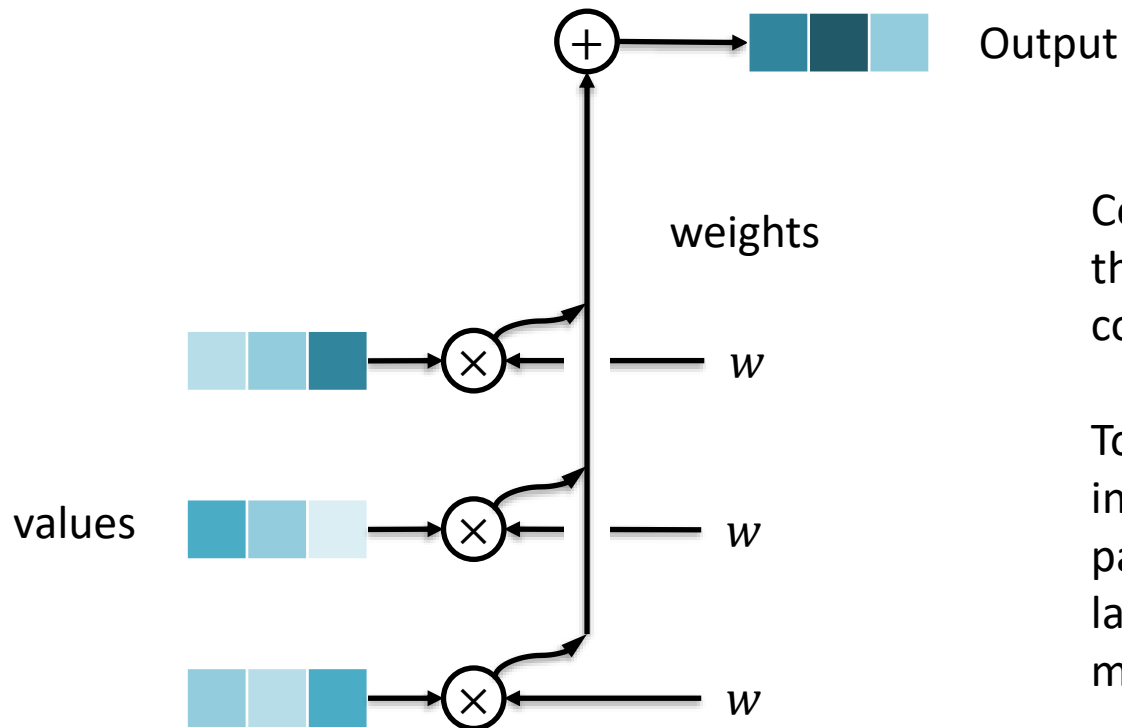
Where do queries come from?



A general framework

ATTENTION IN DEEP LEARNING

Starting point

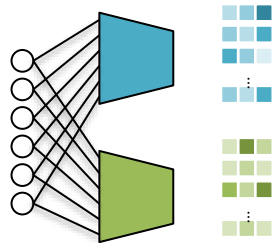
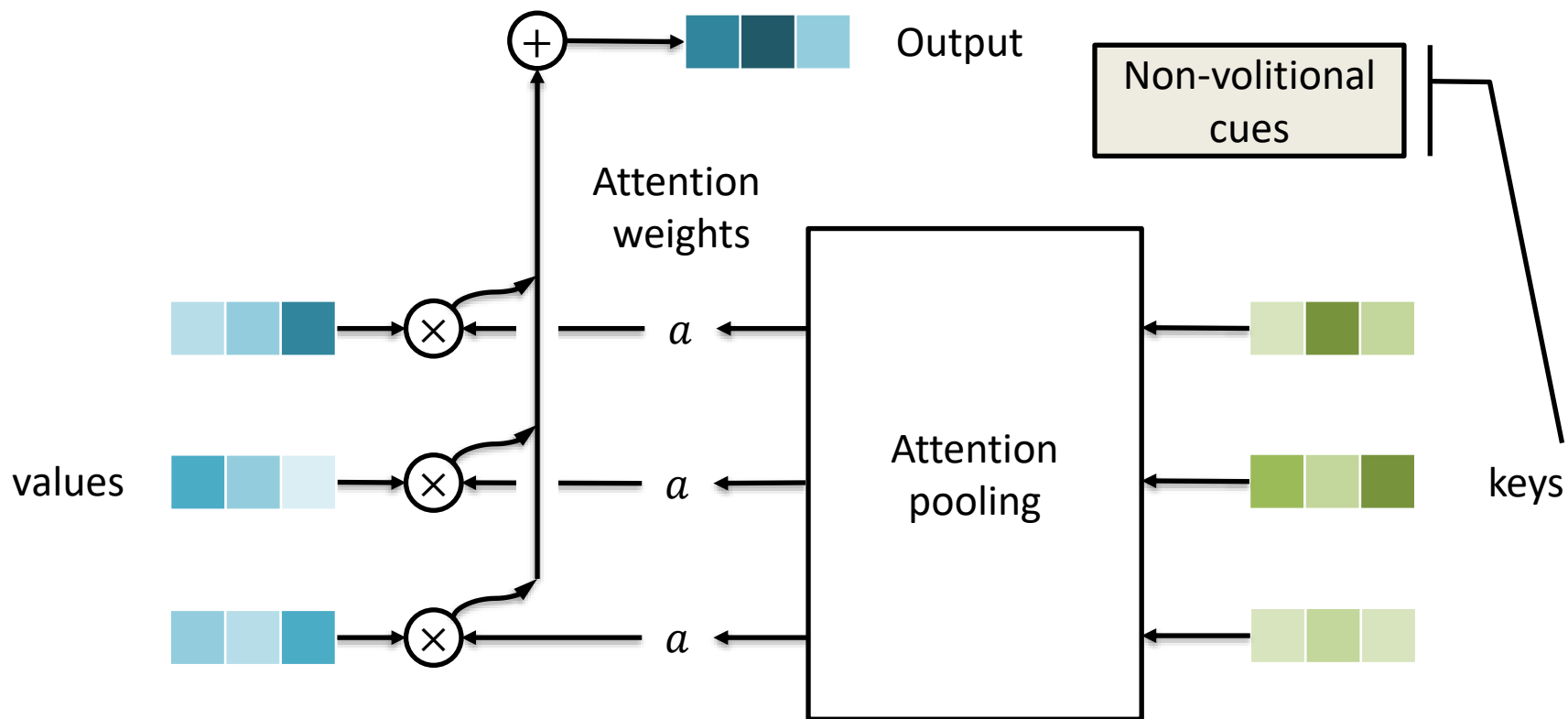


Consider the simpler case where there is **no query** (no top-down component).

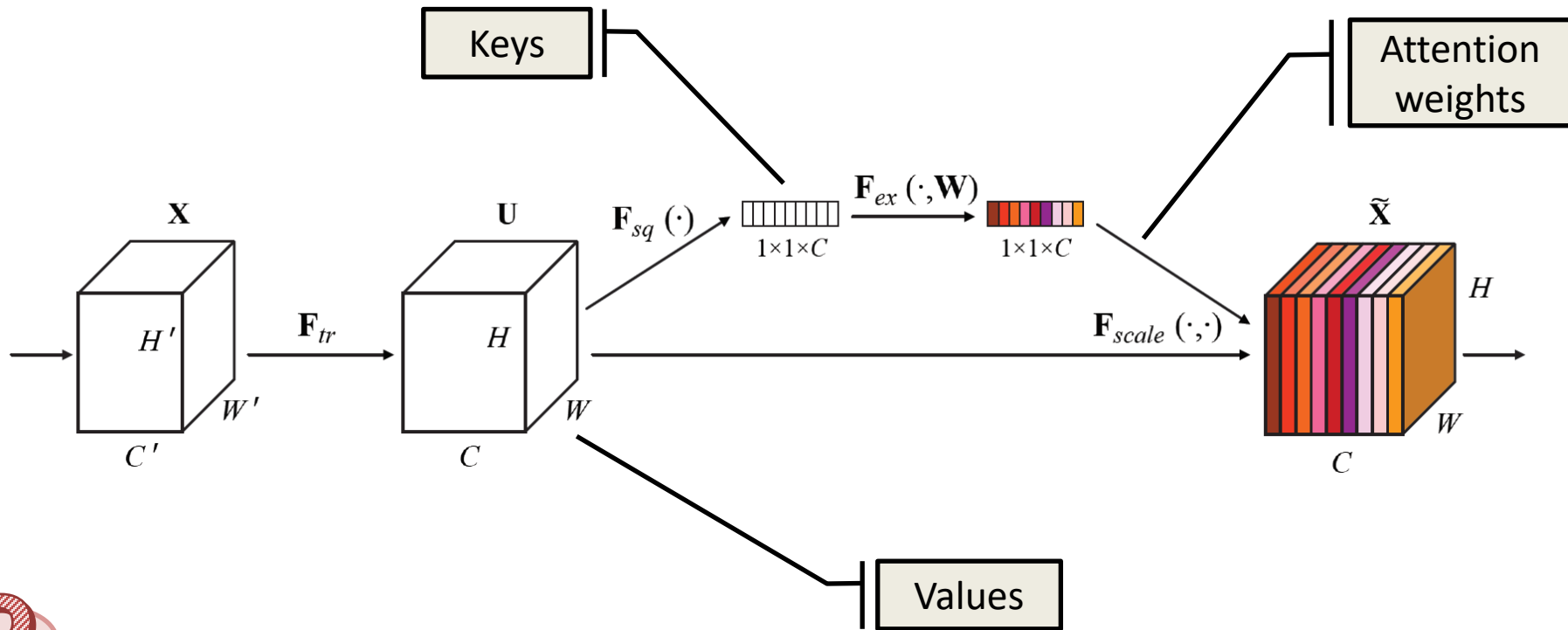
To bias selection over sensory inputs, we can simply use a parameterized fully-connected layer or even non-parameterized max or average pooling.

This is what we have been doing up to now

“Bottom-up” attention

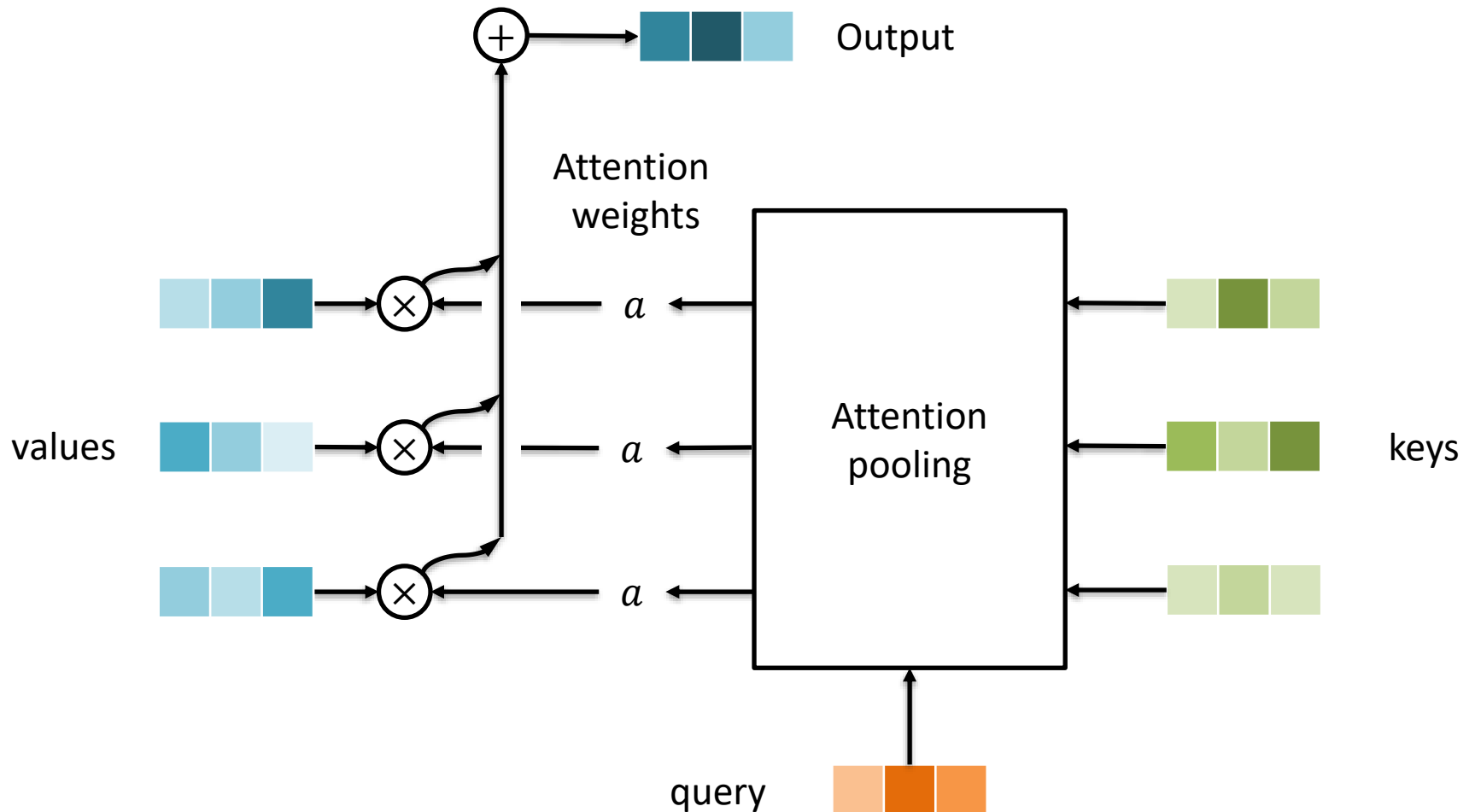


Example: Squeeze excitation networks

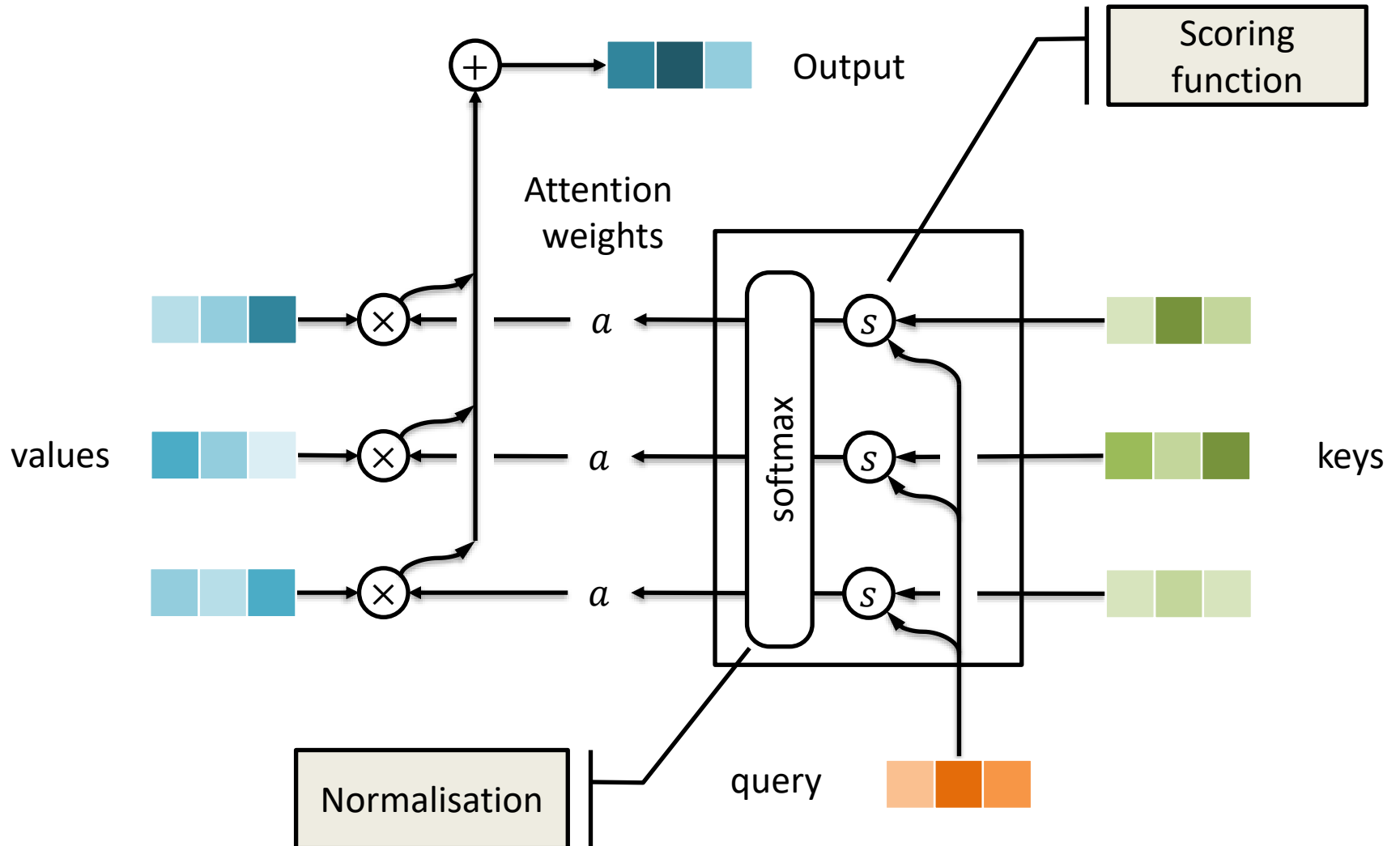


Still no query. How to incorporate volitional cues?

Introducing Queries

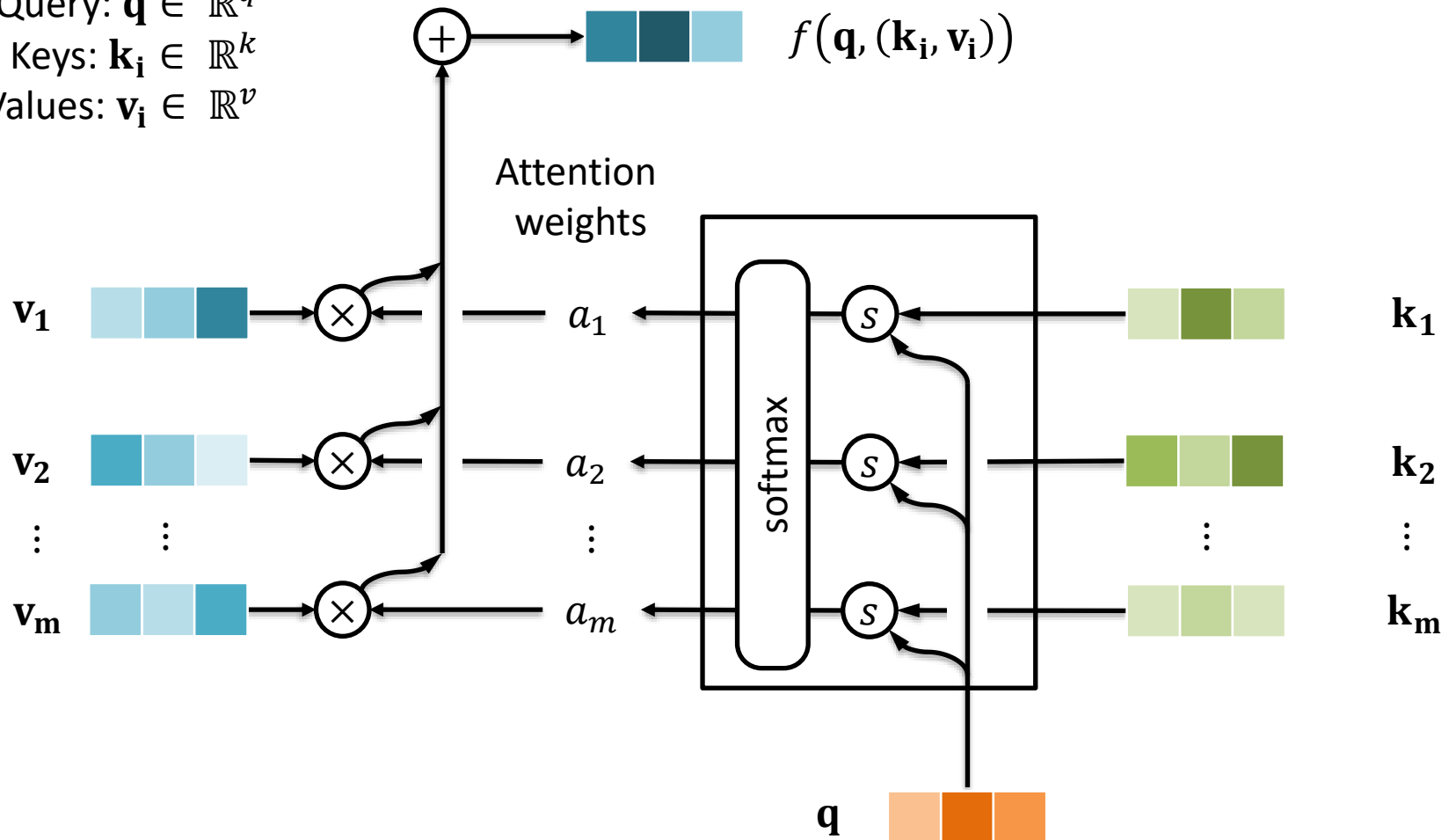


Introducing Queries

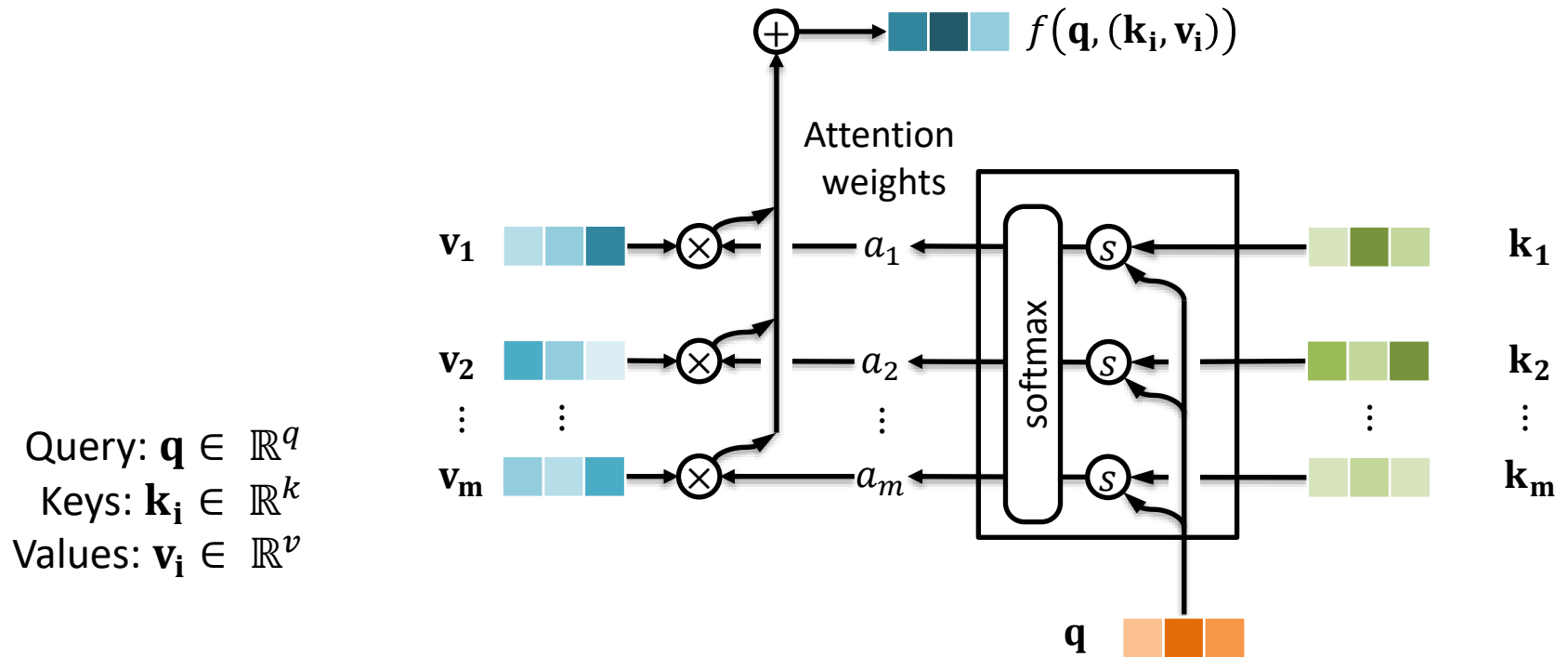


Notation

Query: $\mathbf{q} \in \mathbb{R}^q$
Keys: $\mathbf{k}_i \in \mathbb{R}^k$
Values: $\mathbf{v}_i \in \mathbb{R}^v$



Notation

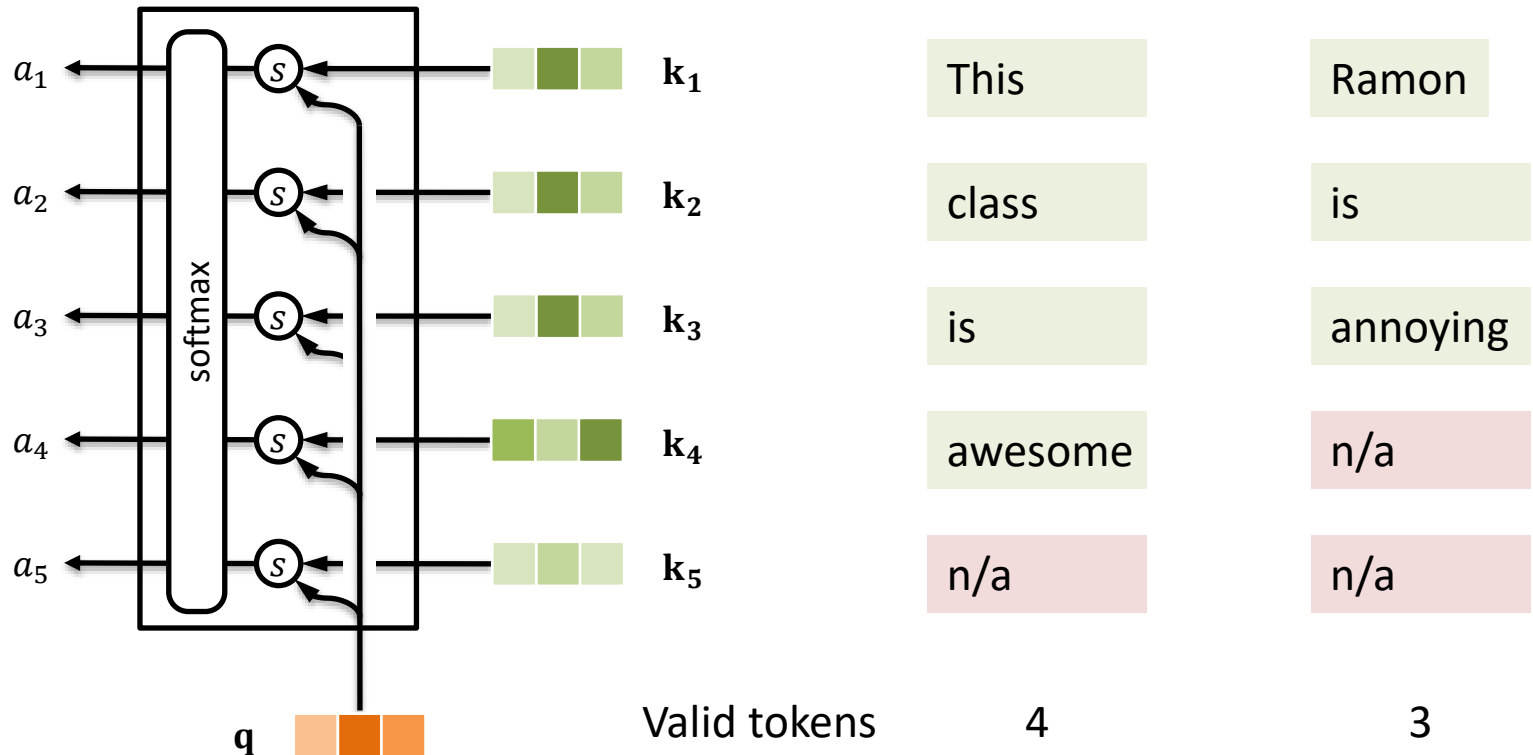


$$f(\mathbf{q}, (\mathbf{k}_1, \mathbf{v}_1), \dots, (\mathbf{k}_m, \mathbf{v}_m)) = \sum_{i=1}^m \alpha(\mathbf{q}, \mathbf{k}_i) \mathbf{v}_i \in \mathbb{R}^v$$

$$\alpha_i = \alpha(\mathbf{q}, \mathbf{k}_i) = \text{softmax}(s(\mathbf{q}, \mathbf{k}_i)) = \frac{\exp(s(\mathbf{q}, \mathbf{k}_i))}{\sum_{j=1}^m \exp(s(\mathbf{q}, \mathbf{k}_j))} \in \mathbb{R}$$

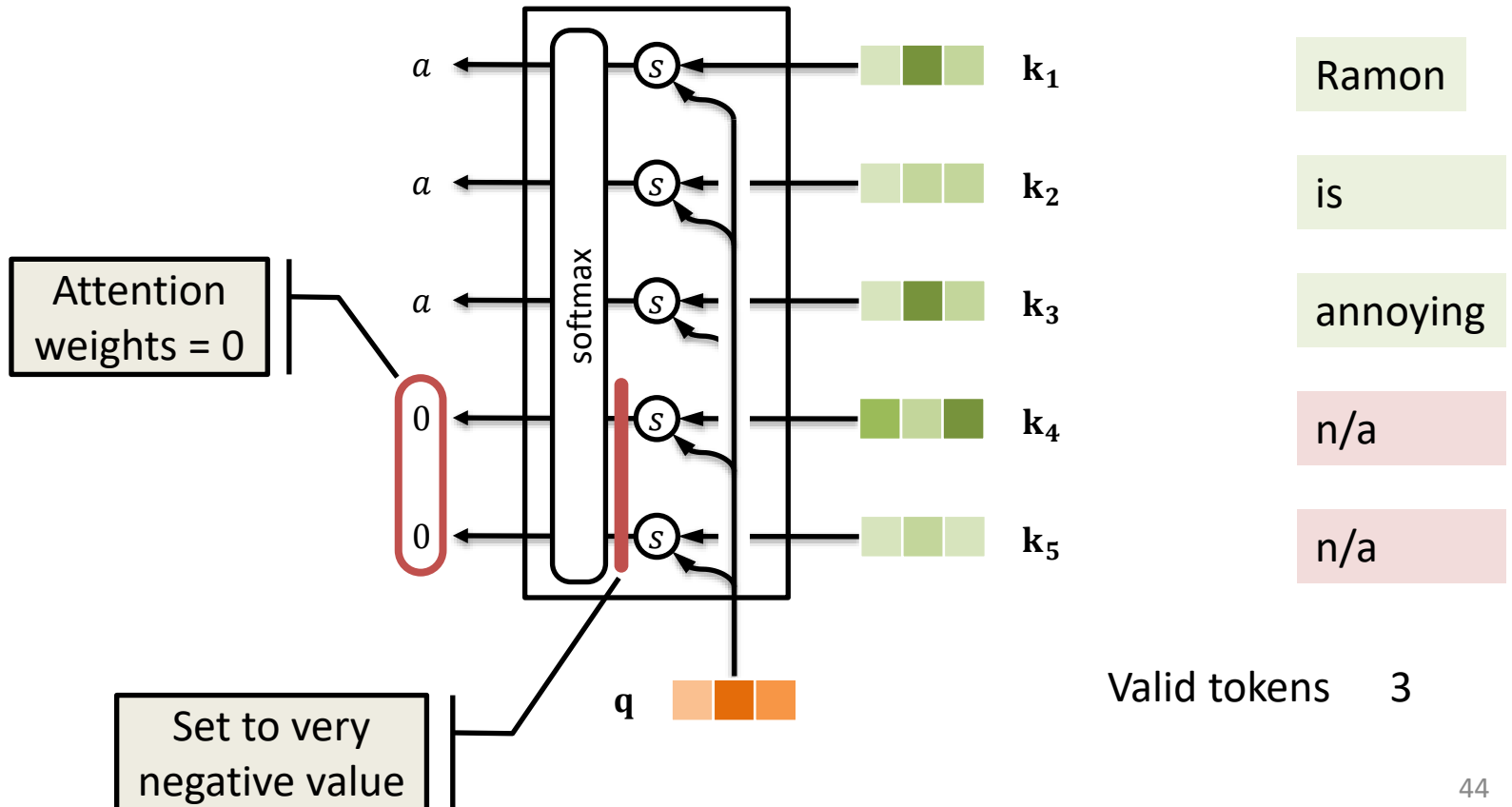
Masked Softmax Operation

In many cases, we are given a variable length of tokens (values / keys), so not all the values should be fed into attention pooling.



Masked Softmax Operation

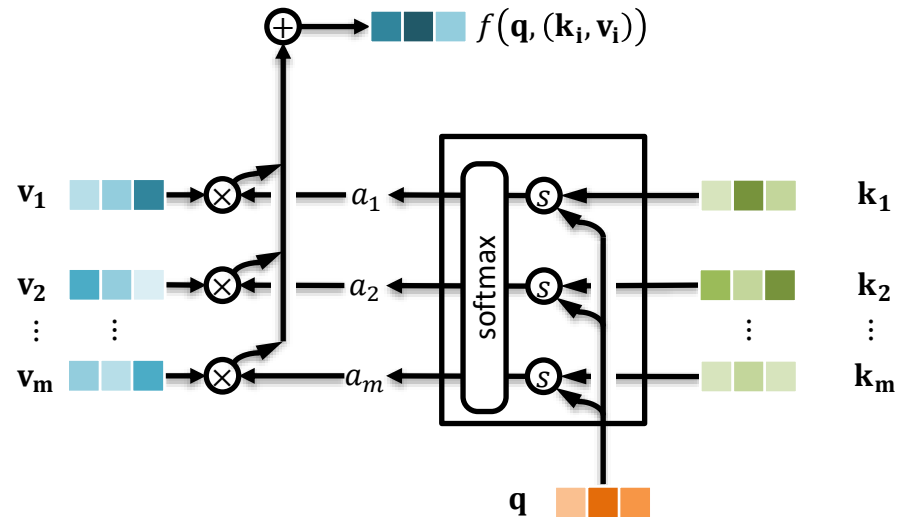
The number of valid tokens becomes a parameter of the *softmax* operation: any value beyond the valid length is masked as zero: given a very negative number value whose exponentiation outputs zero.



Attention scoring functions:

Additive attention

In general, when queries and keys are vectors of **different lengths**, we can use additive attention as the scoring function



$$s(\mathbf{q}, \mathbf{k}) = \mathbf{w}_v^T \tanh(\mathbf{W}_q \mathbf{q} + \mathbf{W}_k \mathbf{k}) \in \mathbb{R}$$

$$\begin{aligned}\mathbf{q} &\in \mathbb{R}^q \\ \mathbf{k} &\in \mathbb{R}^k \\ \mathbf{W}_q &\in \mathbb{R}^{h \times q} \\ \mathbf{W}_k &\in \mathbb{R}^{h \times k} \\ \mathbf{w}_v &\in \mathbb{R}^h\end{aligned}$$

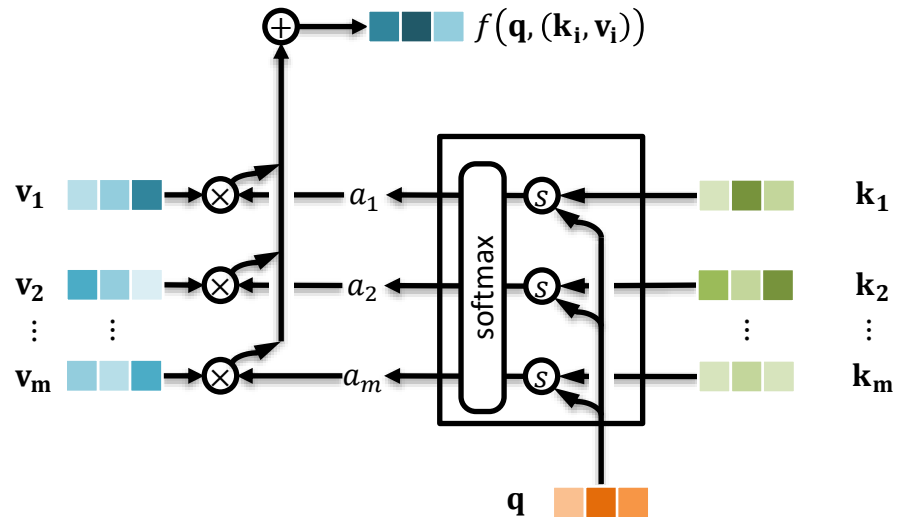
Attention scoring functions: Scaled Dot-product Attention

Dot product is **computationally efficient**. But requires that both the query and the key have the same vector length.

Assuming that \mathbf{q} and \mathbf{k} are independent random variables with zero mean and unit variance, we divide by \sqrt{d} to ensure result also has unit variance

$$s(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^T \mathbf{k}}{\sqrt{d}}$$

$$\mathbf{q}, \mathbf{k} \in \mathbb{R}^d$$



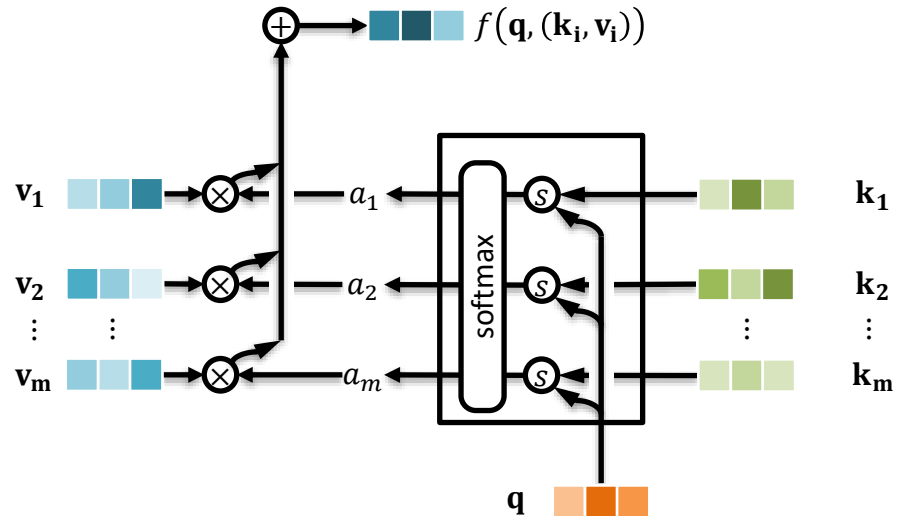
Attention scoring functions: Scaled Dot-product Attention

Dot product is **computationally efficient**. But requires that both the query and the key have the same vector length.

Assuming that \mathbf{q} and \mathbf{v} are independent random variables with zero mean and unit variance, we divide by \sqrt{d} to ensure result also has unit variance

$$s(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^T \mathbf{k}}{\sqrt{d}}$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right) \mathbf{V} \in \mathbb{R}^{n \times v}$$



n : # of queries

m : # of keys

$$\mathbf{v} \in \mathbb{R}^v$$

$$\mathbf{q}, \mathbf{k} \in \mathbb{R}^d$$

$$\mathbf{Q} \in \mathbb{R}^{n \times d}$$

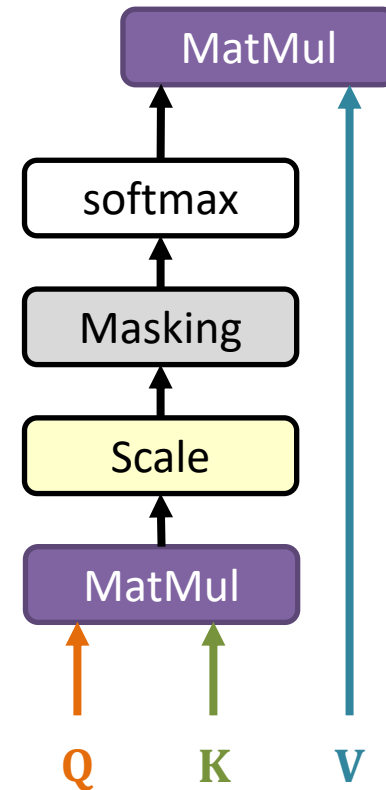
$$\mathbf{K} \in \mathbb{R}^{m \times d}$$

$$\mathbf{V} \in \mathbb{R}^{m \times v}$$

Scaled Dot-Product Attention

Efficient parallel implementation
for multiple keys/queries:

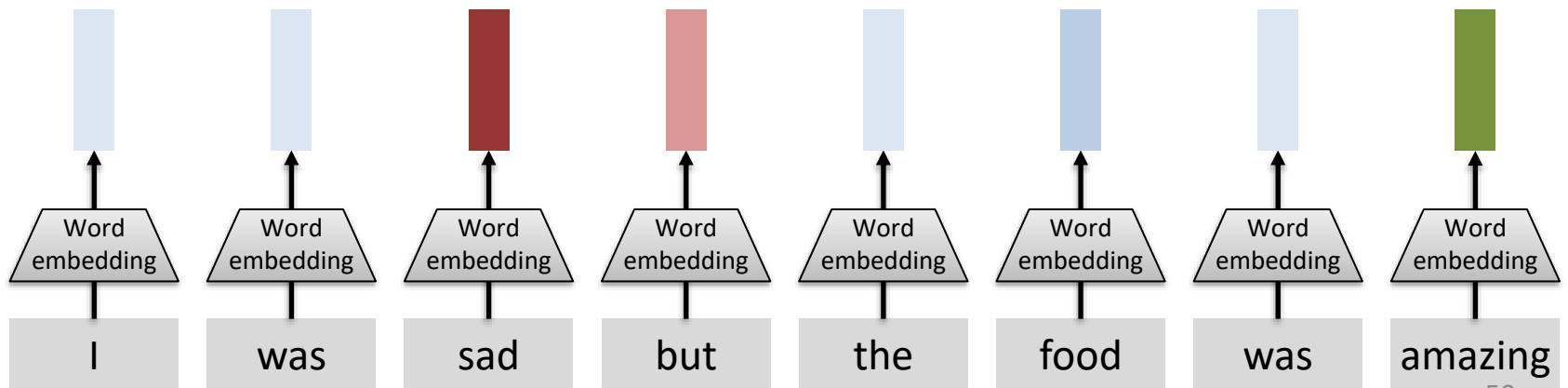
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right) \mathbf{V}$$



EXAMPLES

Example: Restaurant reviews

Is this a positive or a negative review?



Example: Restaurant reviews

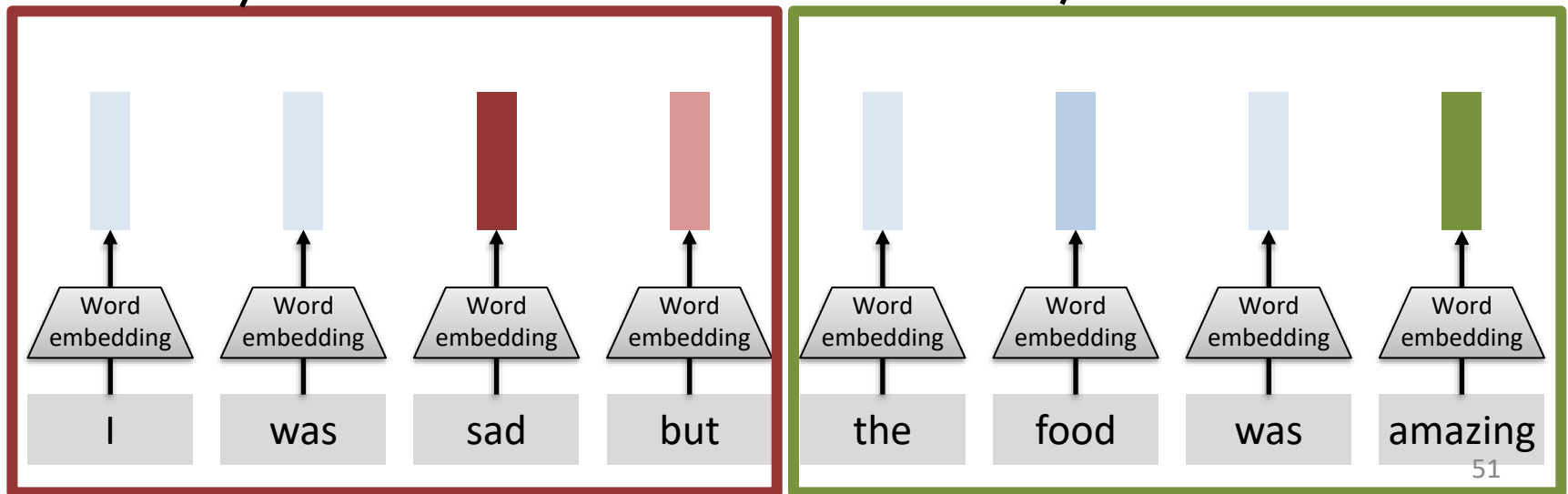
In this example, not all words are relevant to understand if this review is positive or not

Can we somehow ask “which are the relevant words”?

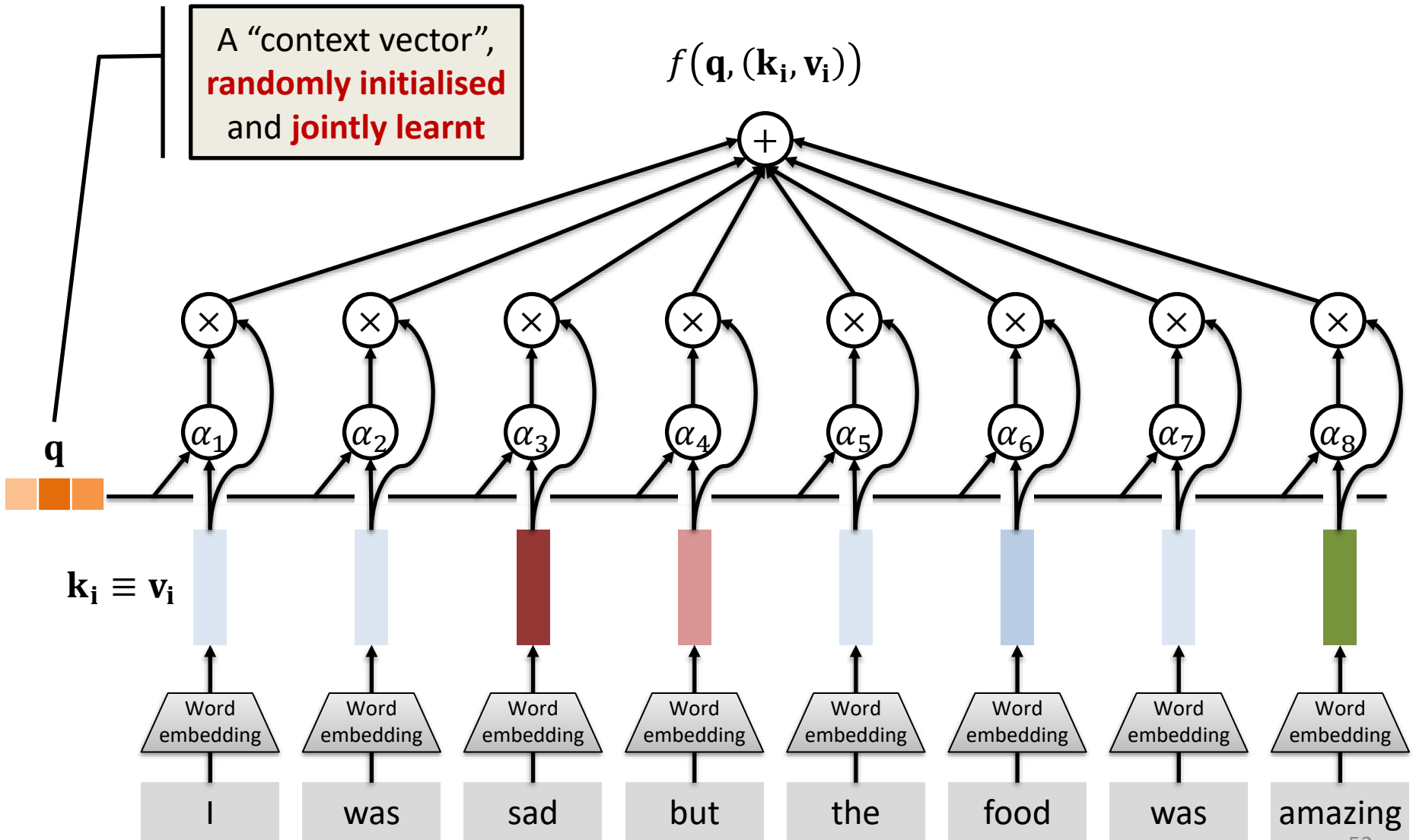


Not relevant for
our purpose...

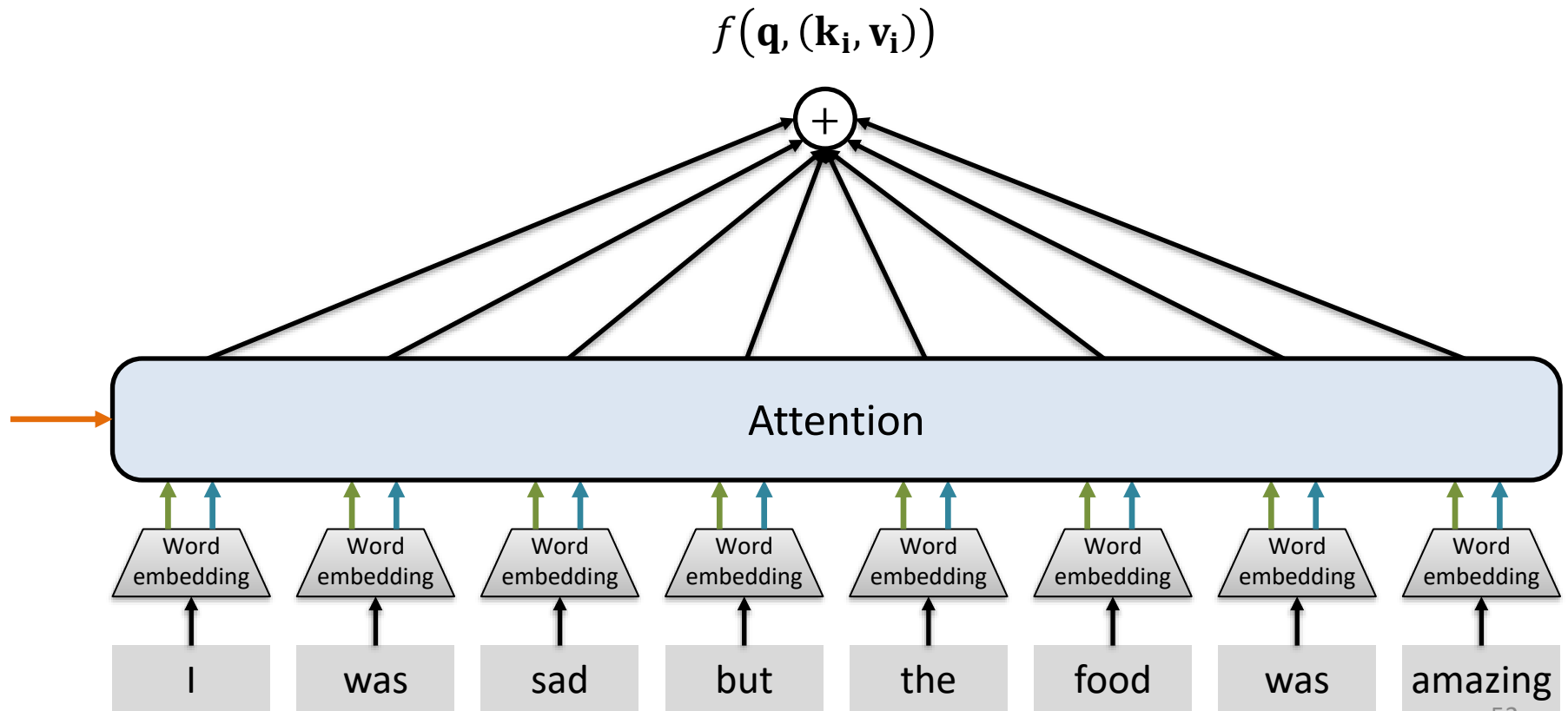
Relevant!



Example: Sentiment Analysis



Example: Sentiment Analysis

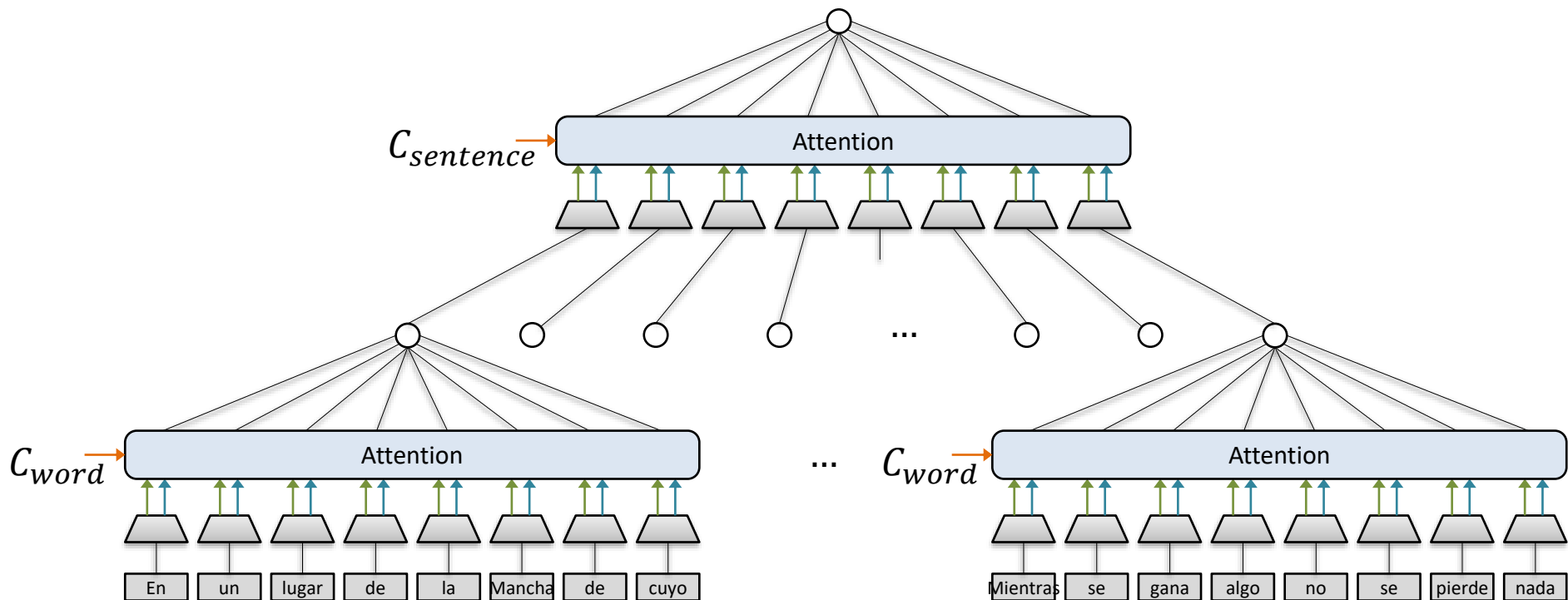


Example: Hierarchical Attention

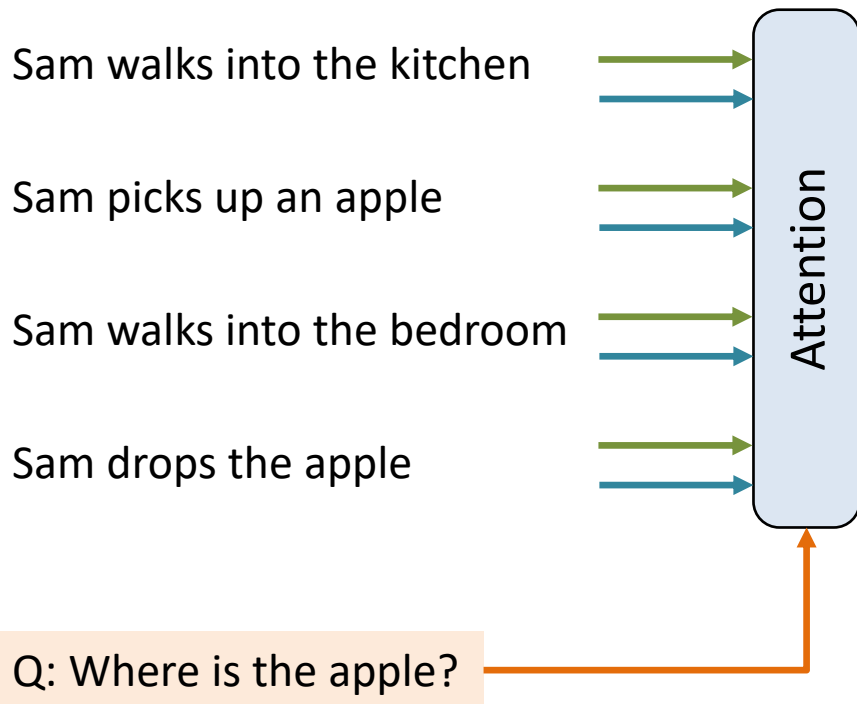
Casa Lolea es un lugarcito fantástico hermoso riquísimo bello simple puro amor
Caímos con mi novio para cenar porque estaba cerca del hotel
Cuando entramos nos dijeron que hay que reservar con dos noches de anticipación
El día siguiente era nuestra última noche en Barcelona
Muy amablemente accedieron a darnos una mesita para el otro día
Volvimos llenos de expectativas porque no sabíamos a qué se debía tanto revuelo
Y lo entendimos

Me sabe mal ver que lugares que me gusten tanto tengan reseñas tan encontradas...
Al parecer las personas que han ido han tenido una mala experiencia.
No me gusta que otros reciban un mal servicio/comida cuando yo la paso tan bien.
Por eso he ido unas cuatro veces a Casa Lolea antes de escribir esta reseña
La verdad es que no tengo ni el más mínimo pero
Las tapas son exquisitas, las recomiendo sobre las conservas
(¿por qué comer conservas en un restaurante? ya esto es una duda MUY personal.)

Example: Hierarchical Attention

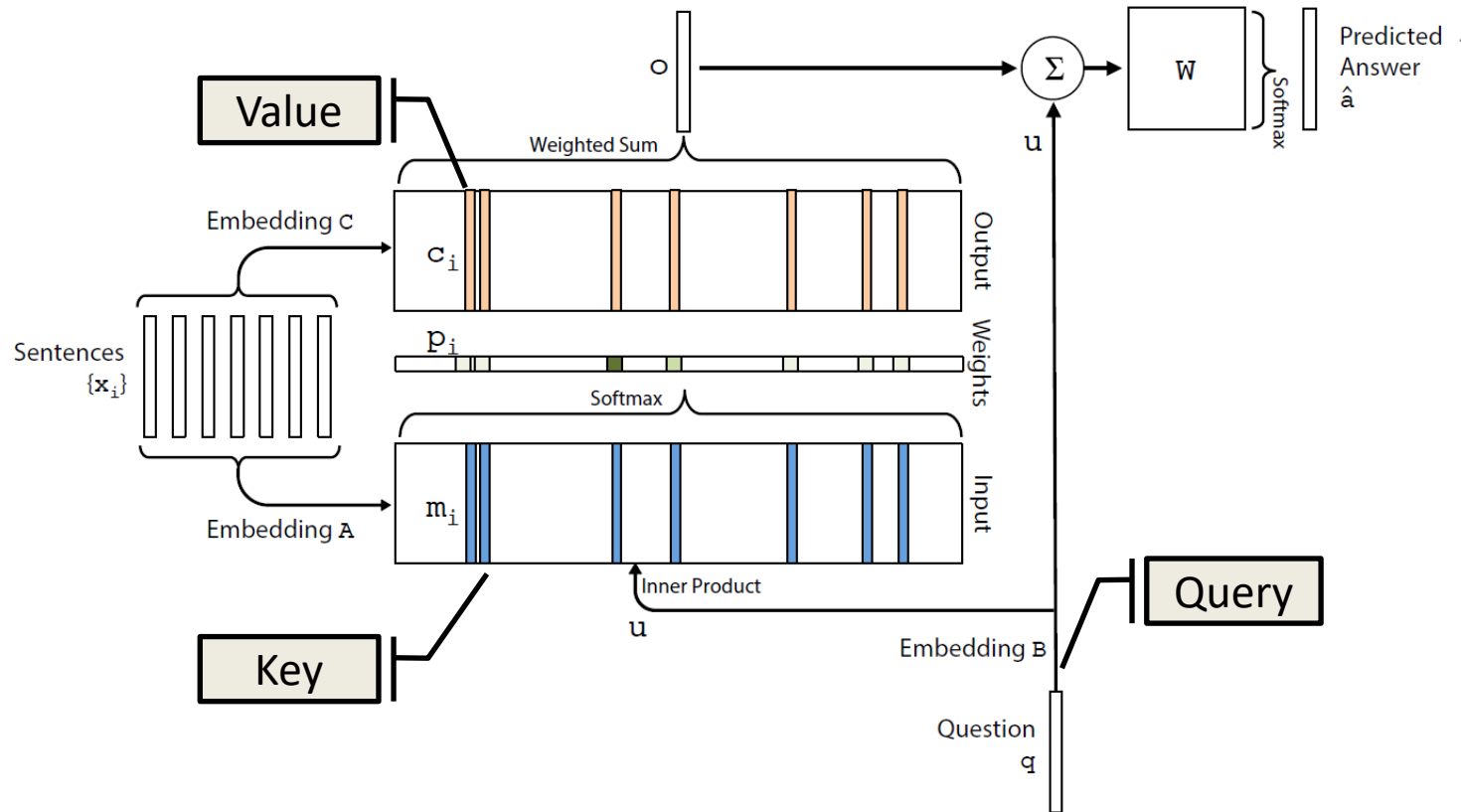


Example – Iterative Attention

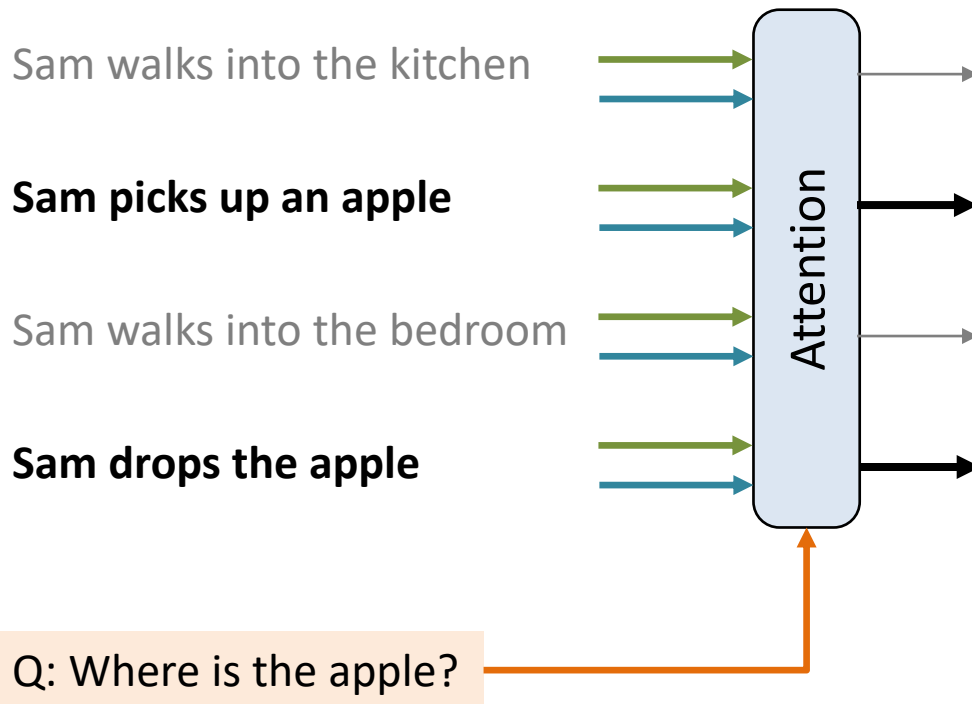


This is a typical question-answering problem. You are given a passage and a question, and you need to come up with an answer

Example – Iterative Attention



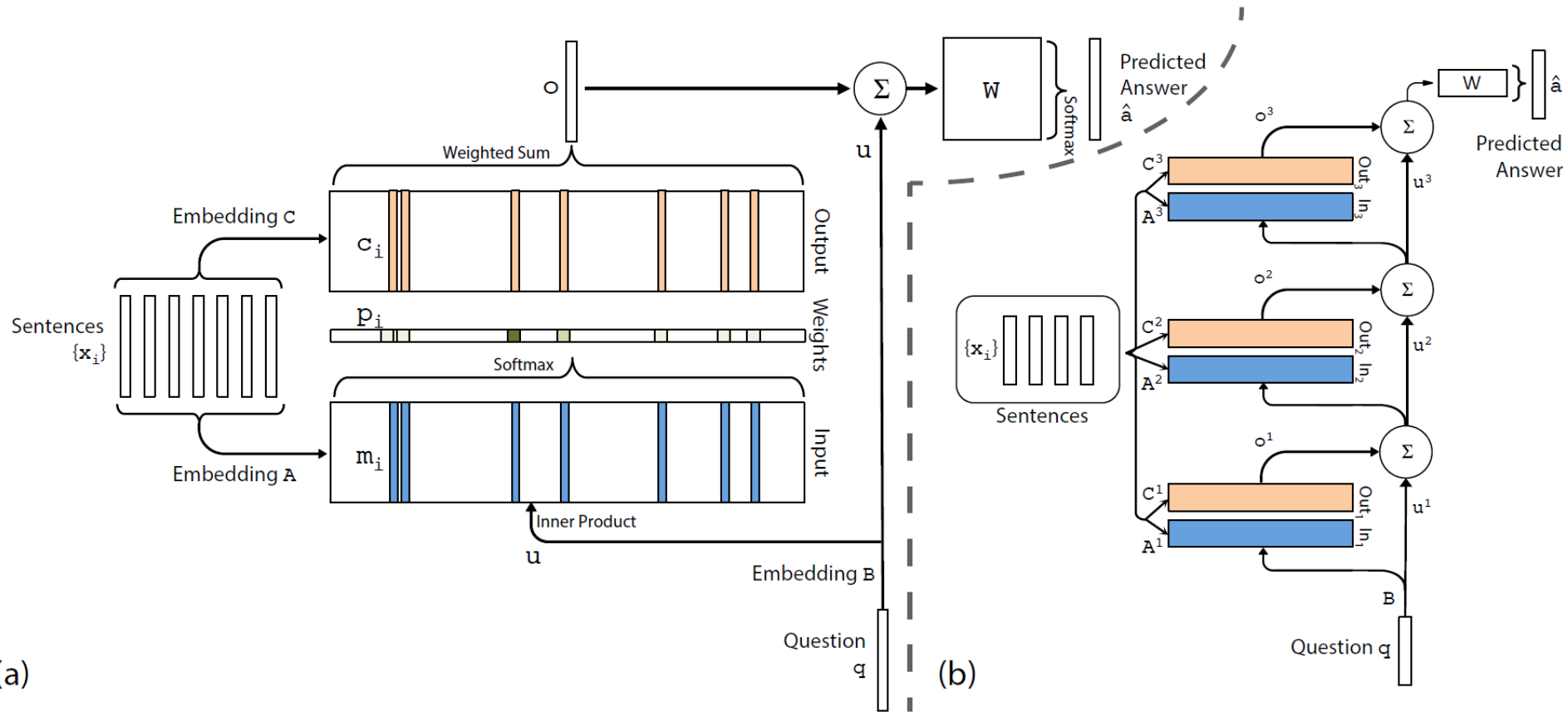
Example – Iterative Attention



Simple attention
selects sentences
with “apple”

Hierarchical
attention does not
work as it misses
intermediate steps

Example – Iterative Attention



Sam walks into the kitchen.
 Sam picks up an apple.
 Sam walks into the bedroom.
 Sam drops the apple.

Q: Where is the apple?

A. Bedroom

Brian is a lion.
 Julius is a lion.
 Julius is white.
 Bernhard is green.

Q: What color is Brian?

A. White

Mary journeyed to the den.
 Mary went back to the kitchen.
 John journeyed to the bedroom.
 Mary discarded the milk.

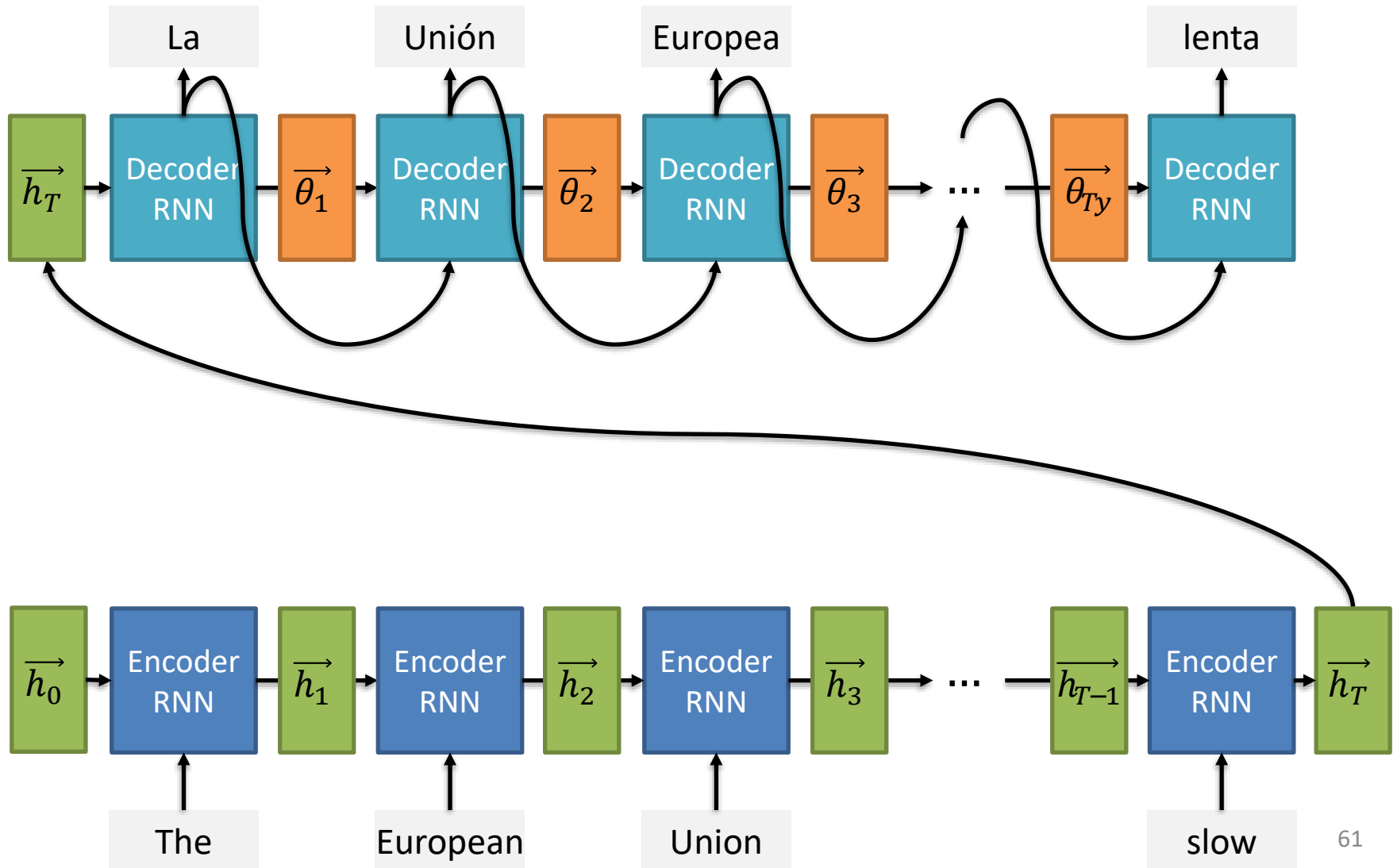
Q: Where was the milk before the den?

A. Hallway

Sequence to Sequence with attention

BAHDANAU ATTENTION

Sequence to Sequence



Sequence to Sequence

This works reasonably well with short sentences. When passages become long, the context vector is a bottleneck – cannot encode all details

The European union is very slow

NMT

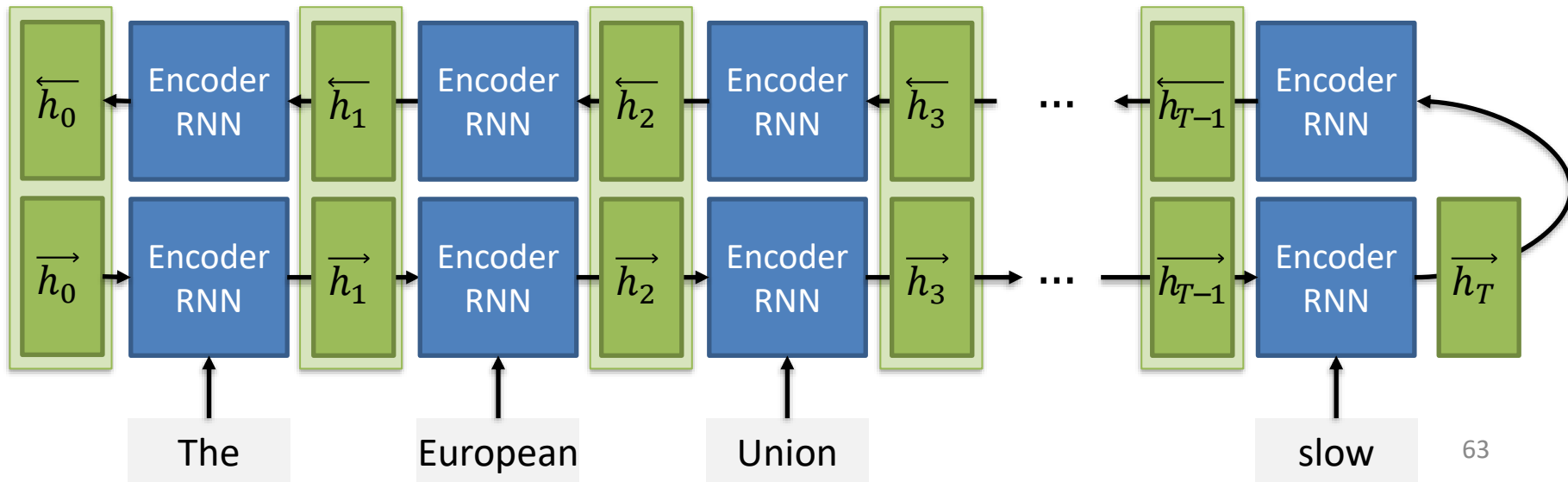
La unión europea es muy lenta

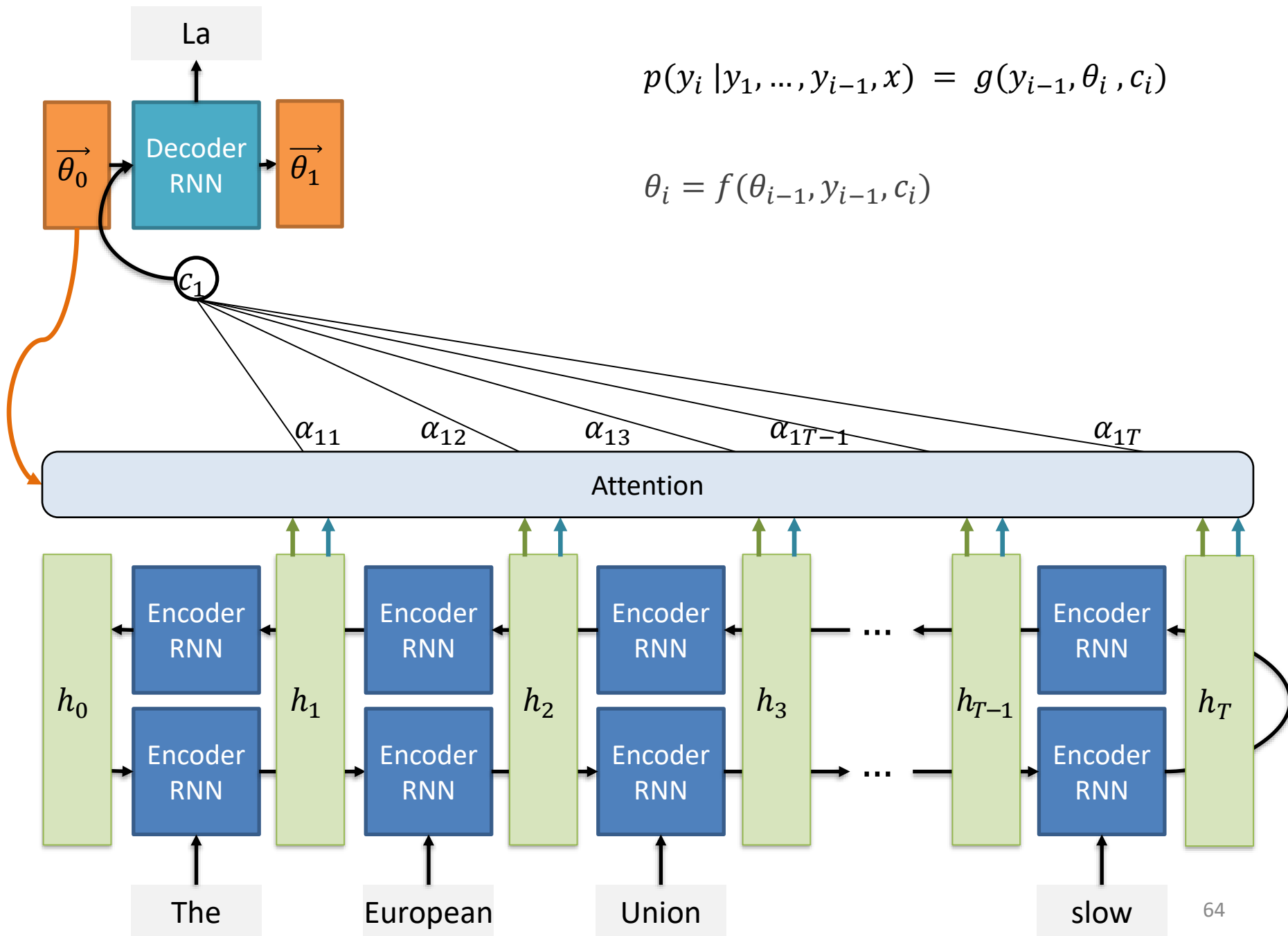
Though the idea of the European Union, a unification of twenty member states including the United Kingdom, might sound simple at the outset, the European Union has a rich history and a unique organization, both of which aid in its current success and its ability to fulfil its mission for the 21st Century.

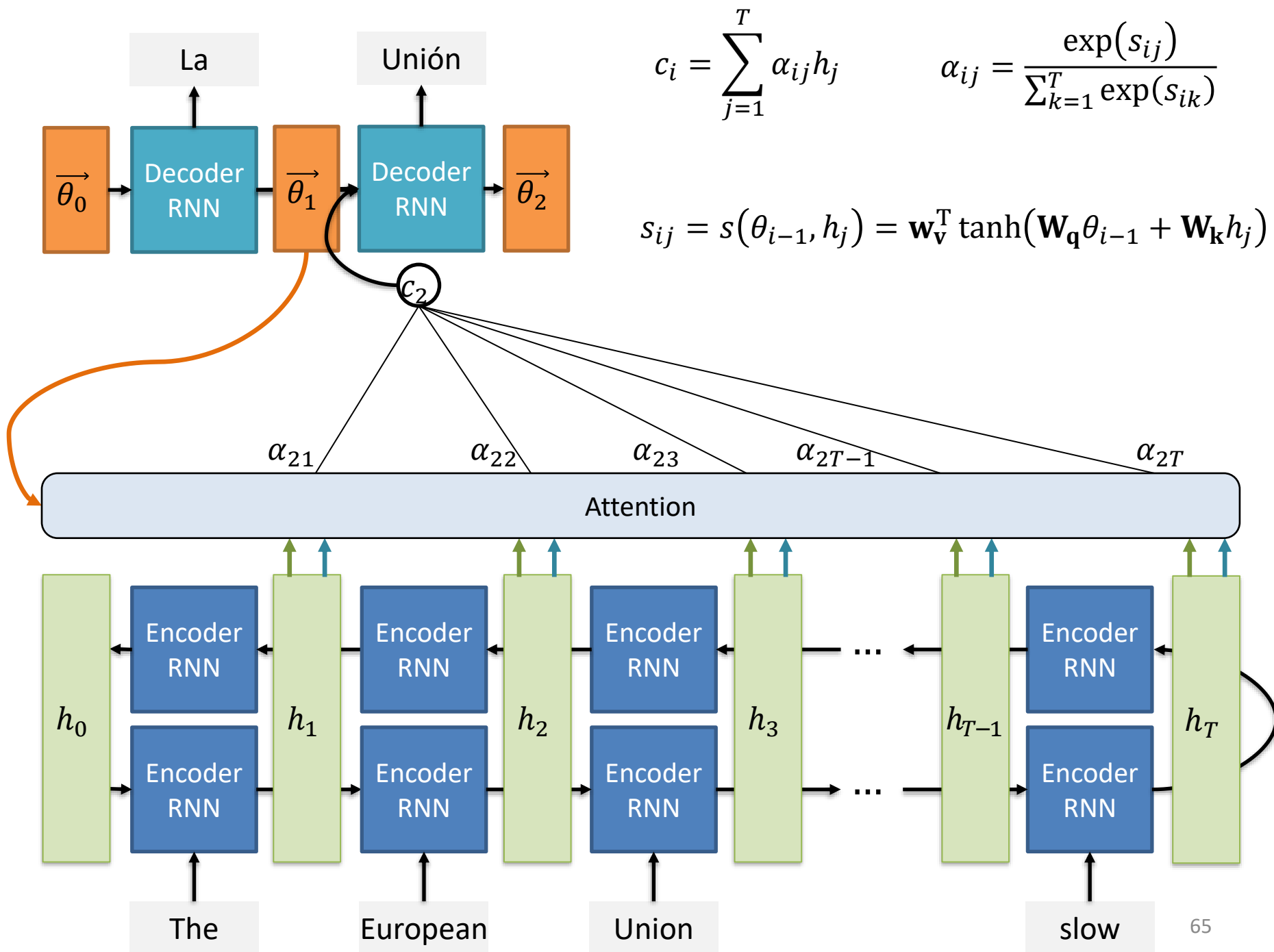
NMT

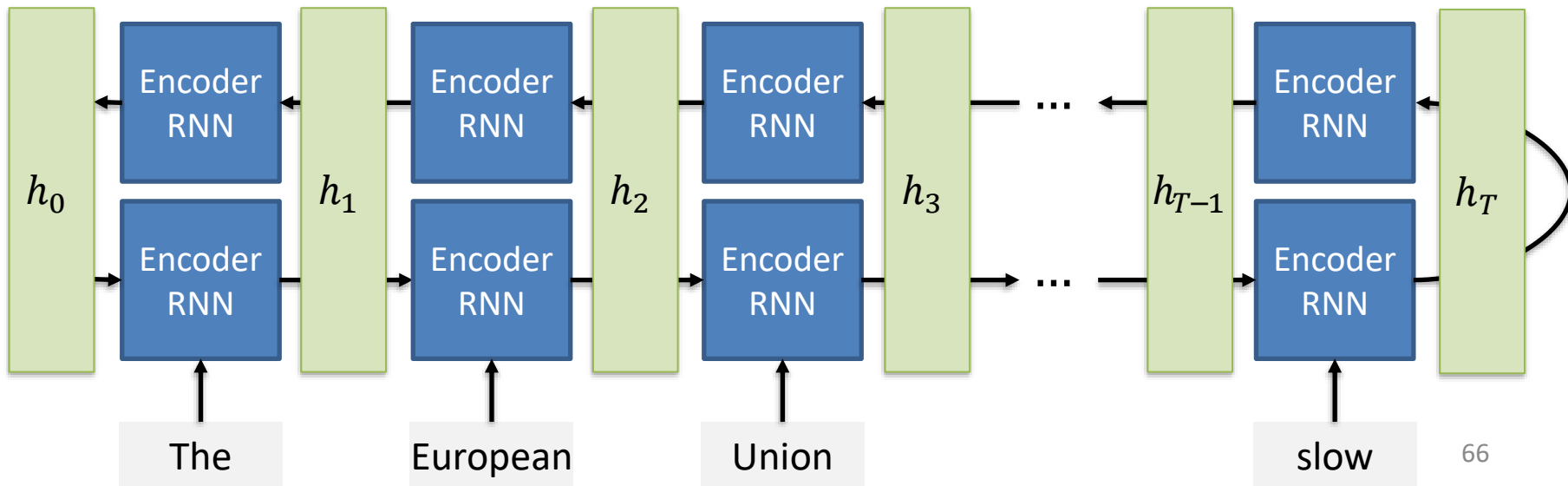
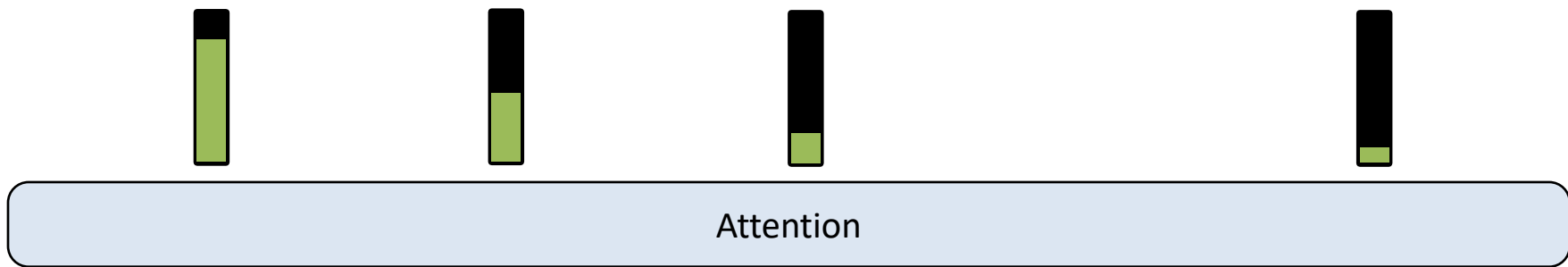
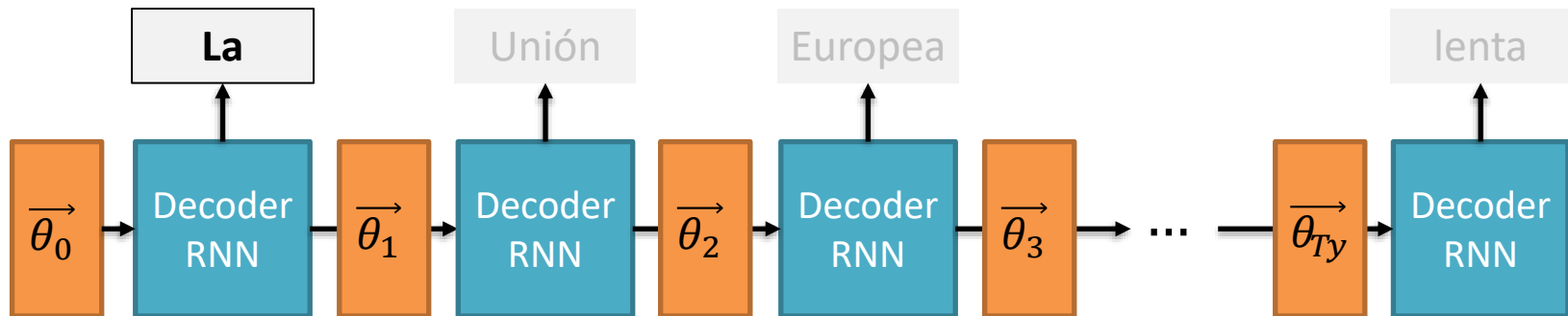
...

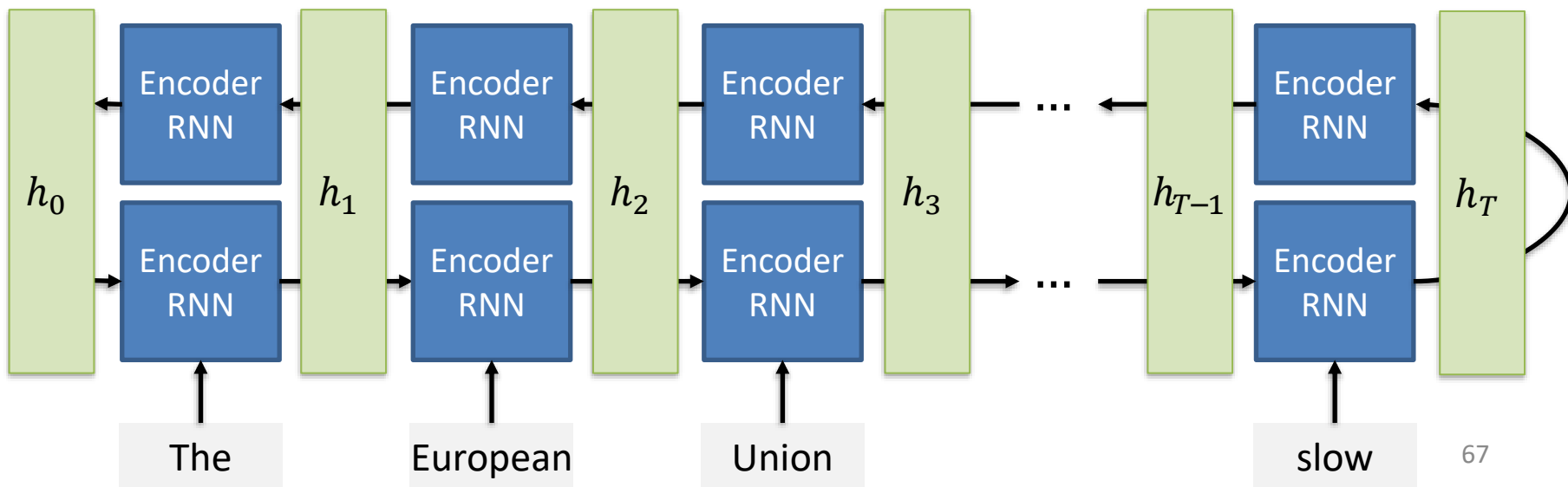
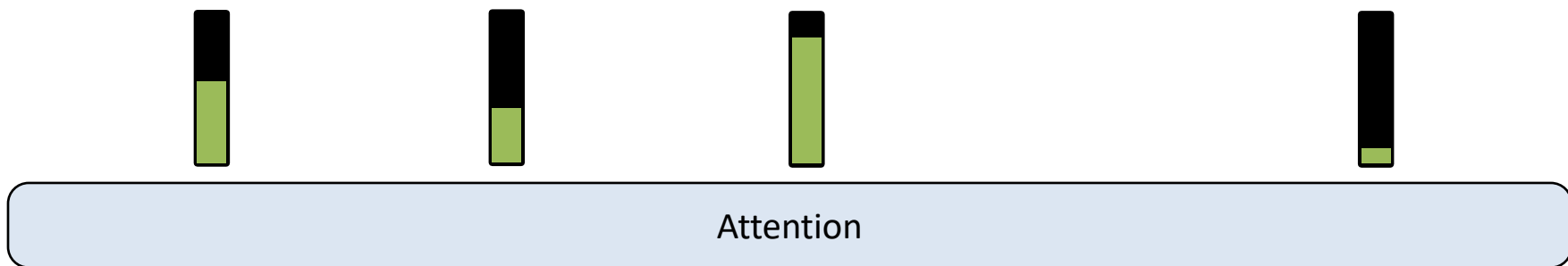
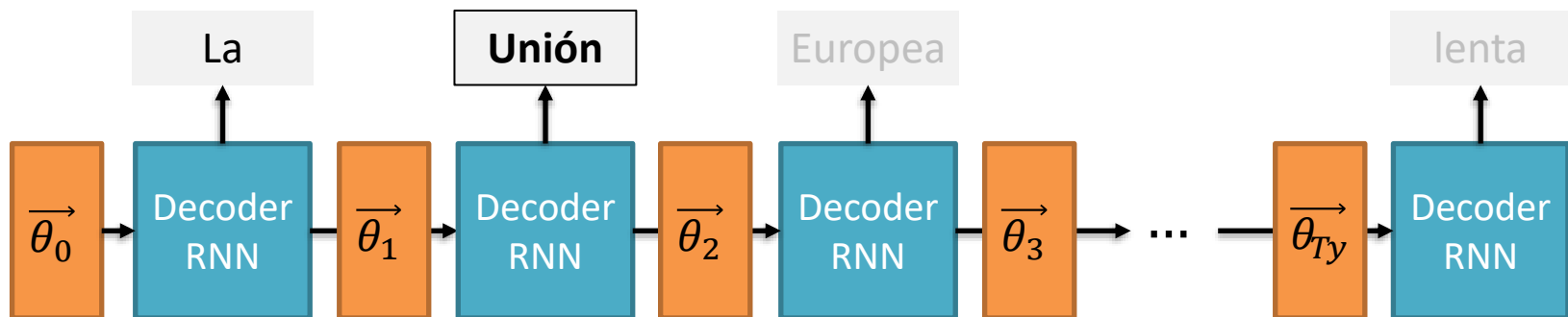
Sequence to Sequence

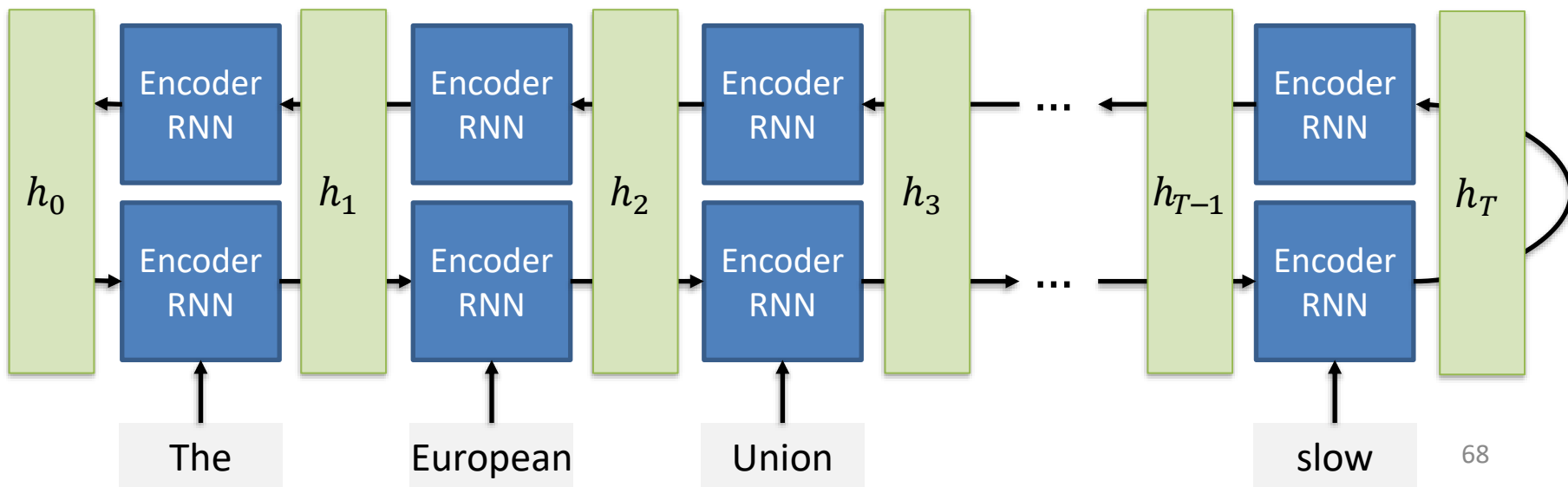
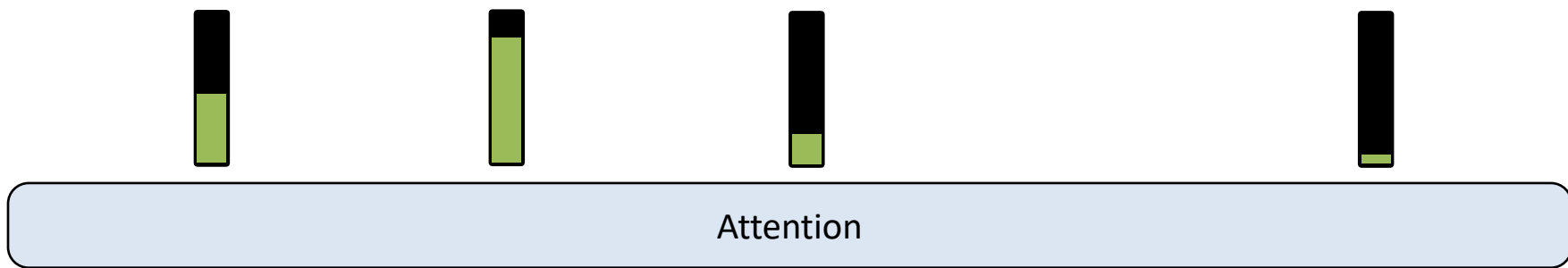
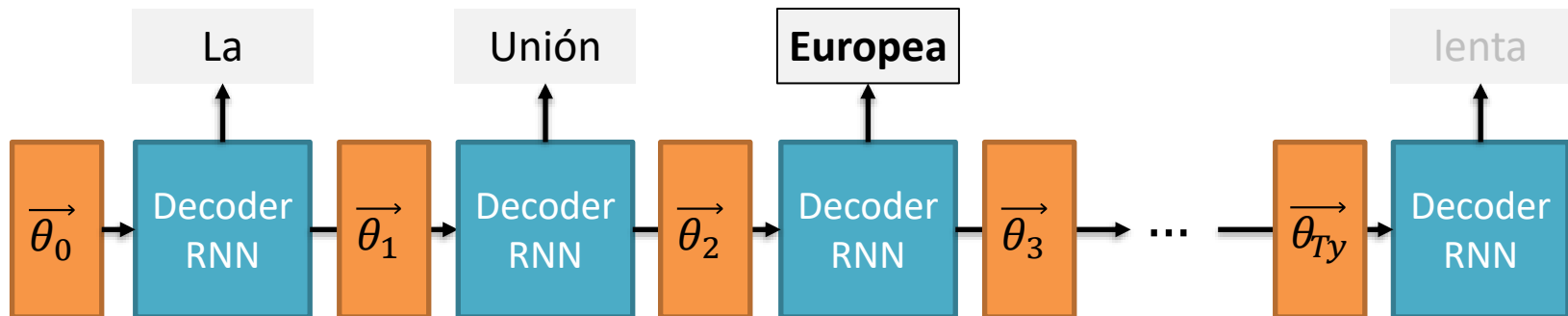




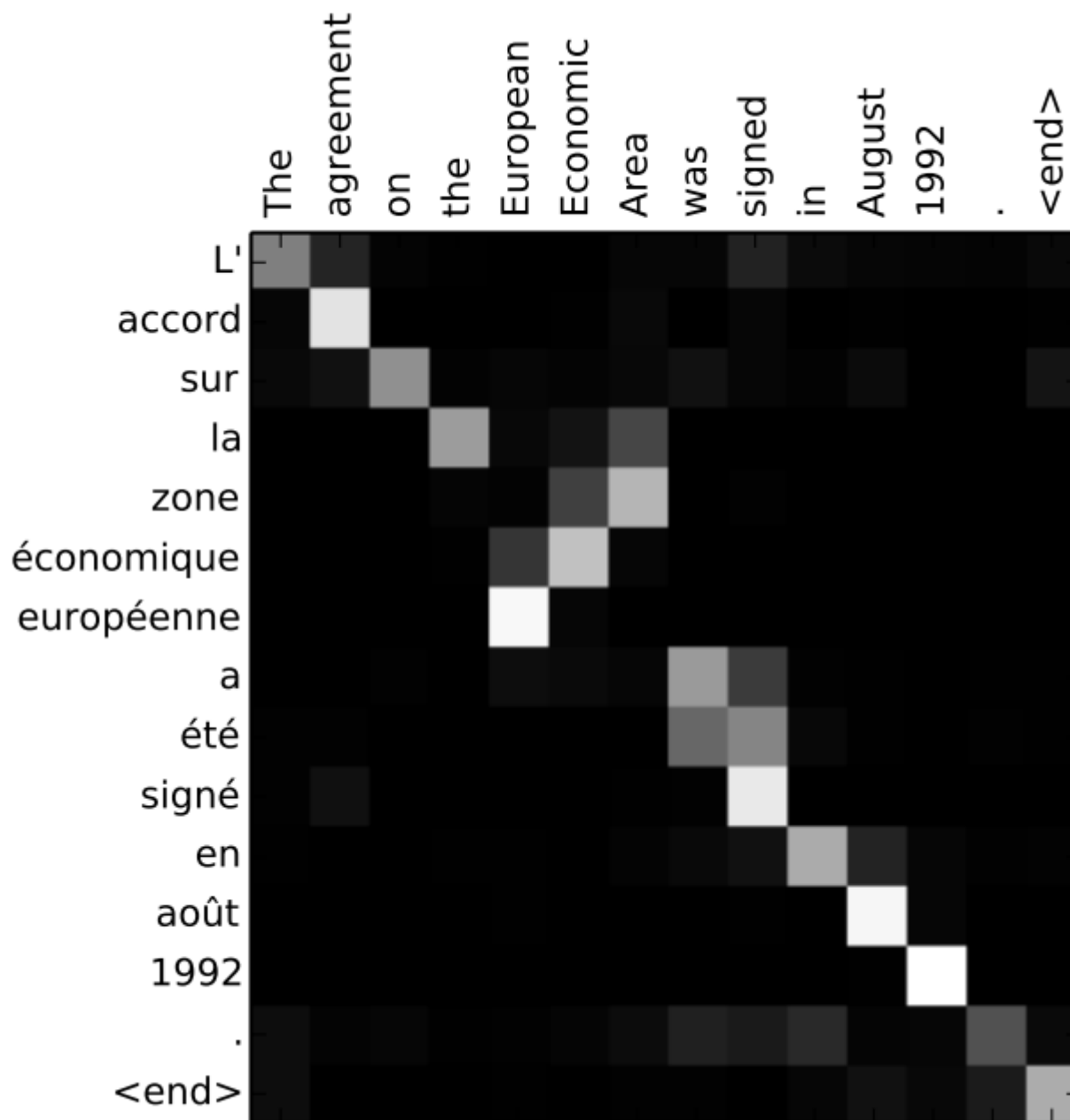






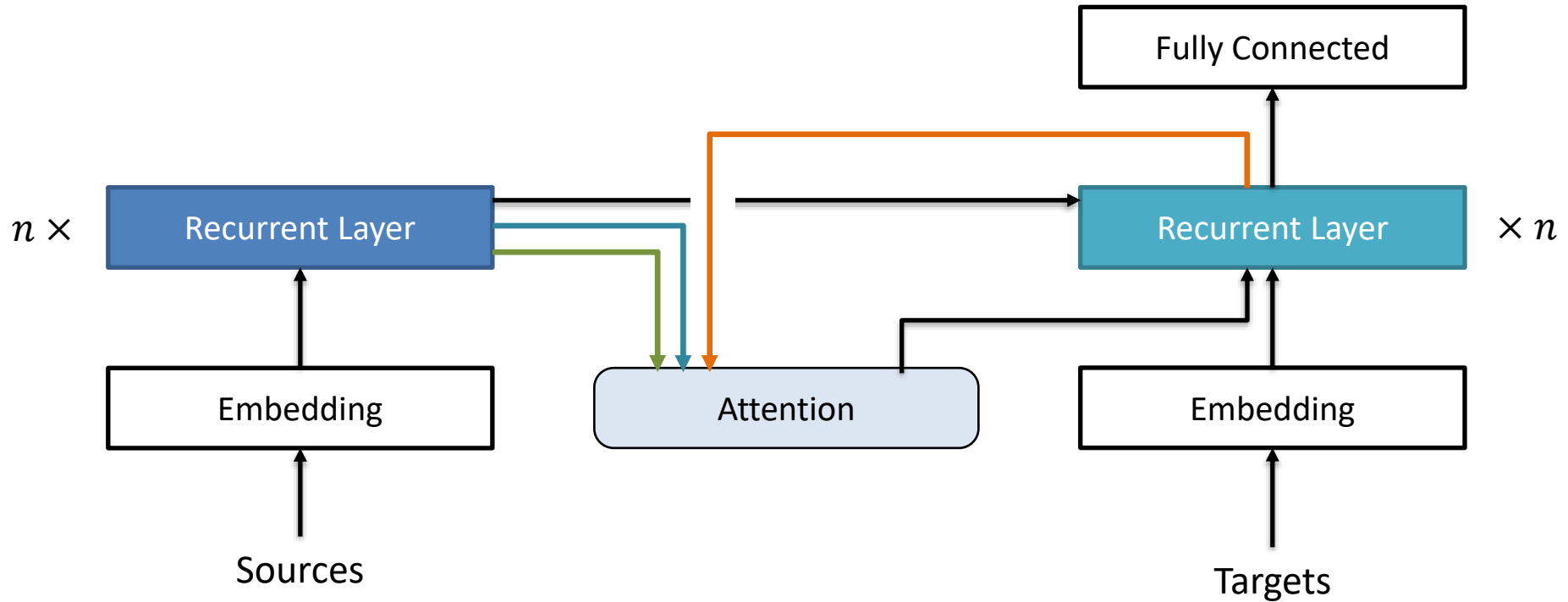


where the model is “attending” when it outputs each word in the French sentence



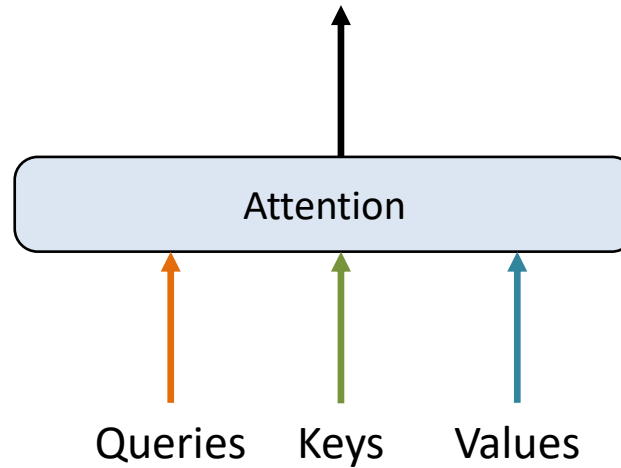
Encoder

Decoder



MULTIHEAD ATTENTION

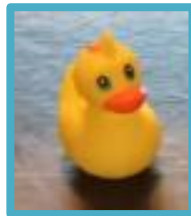
Single-head Attention



Queries, keys and values encode a multitude of information about our tokens.

Sometimes it might be beneficial to allow our attention mechanism to **jointly use different representation subspaces** of queries, keys and values (different “aspects” they describe)

Values



Keys

pink, human

yellow, inanimate

blue, inanimate

yellow, animal

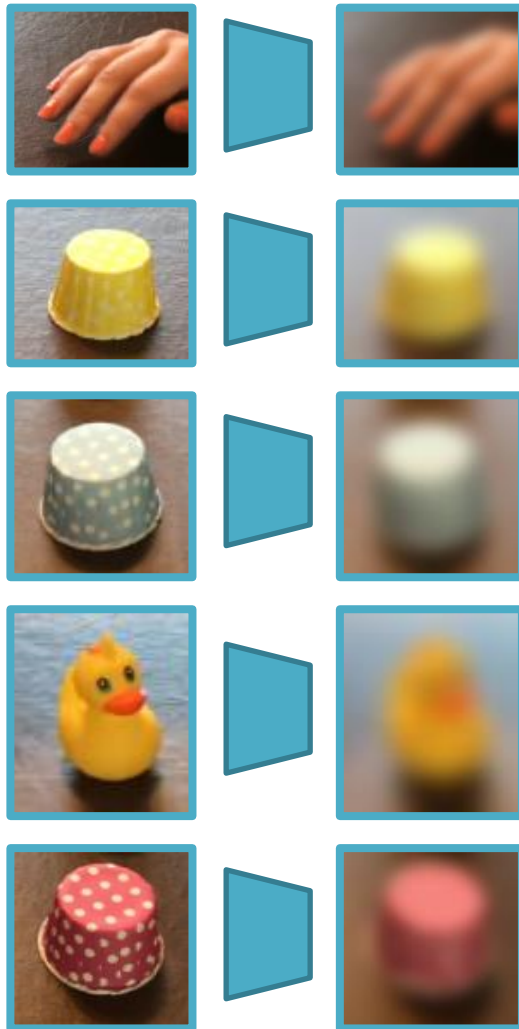
red, inanimate



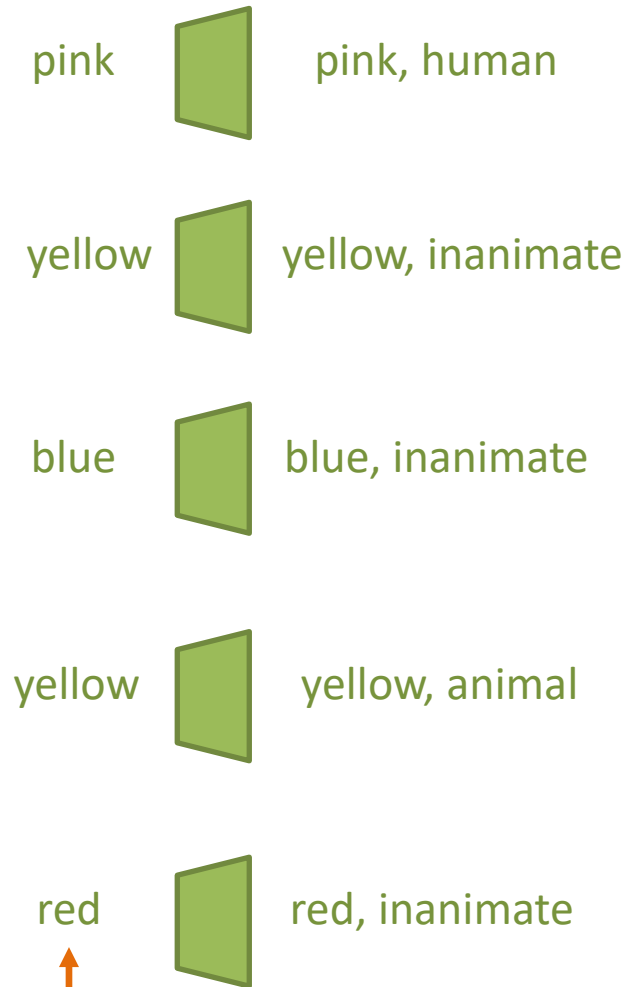
Query "Yellow duck"



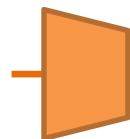
Values



Keys



Query "Yellow duck"



"Yellow"

red

red, inanimate

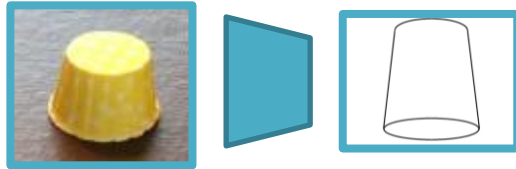
Values

Keys



human

pink, human



inanimate

yellow, inanimate



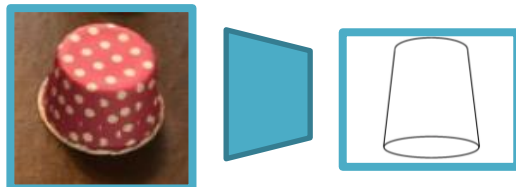
inanimate

blue, inanimate



animal

yellow, animal



inanimate

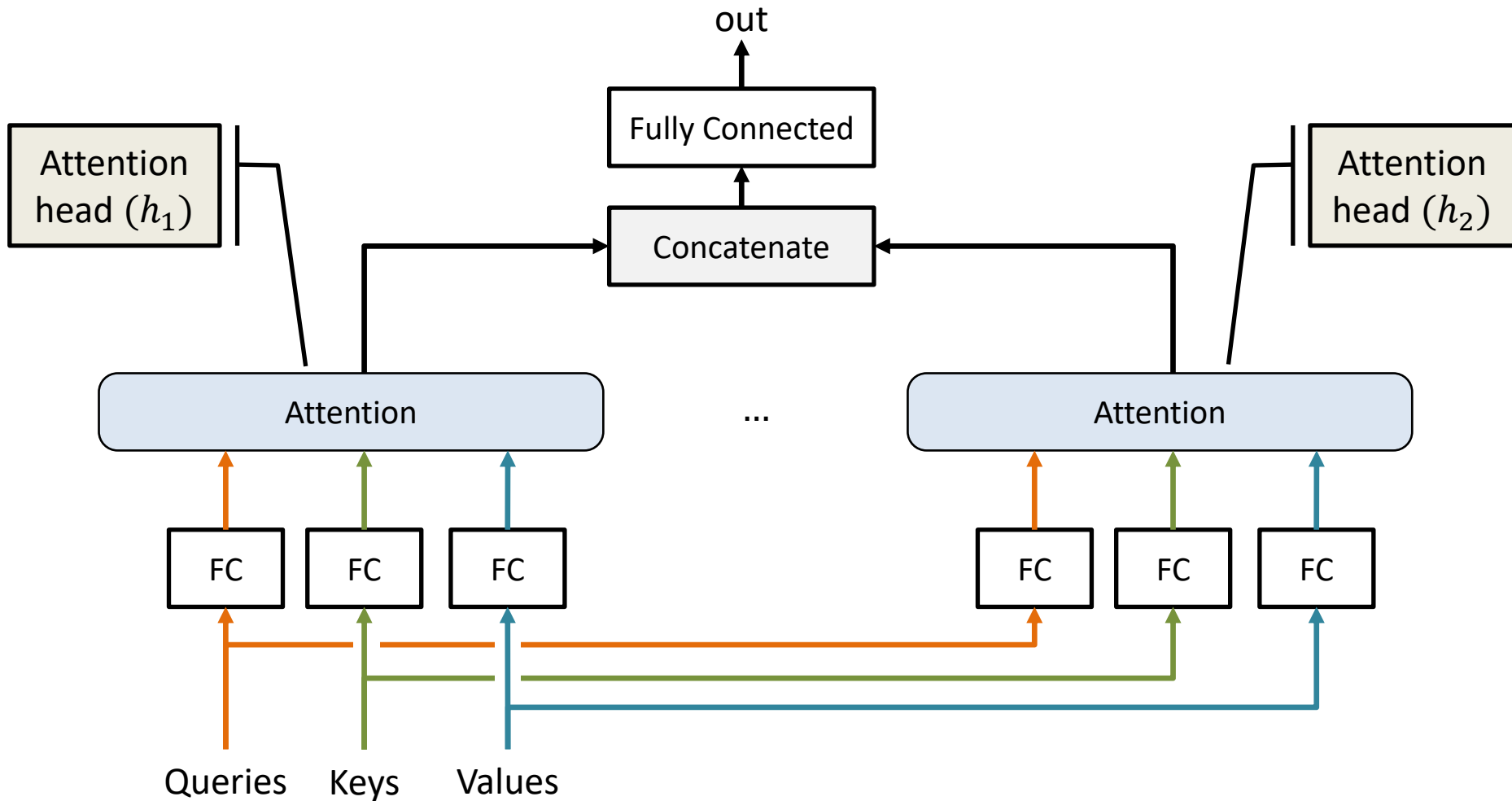
red, inanimate

Query "Yellow duck"

"duck"

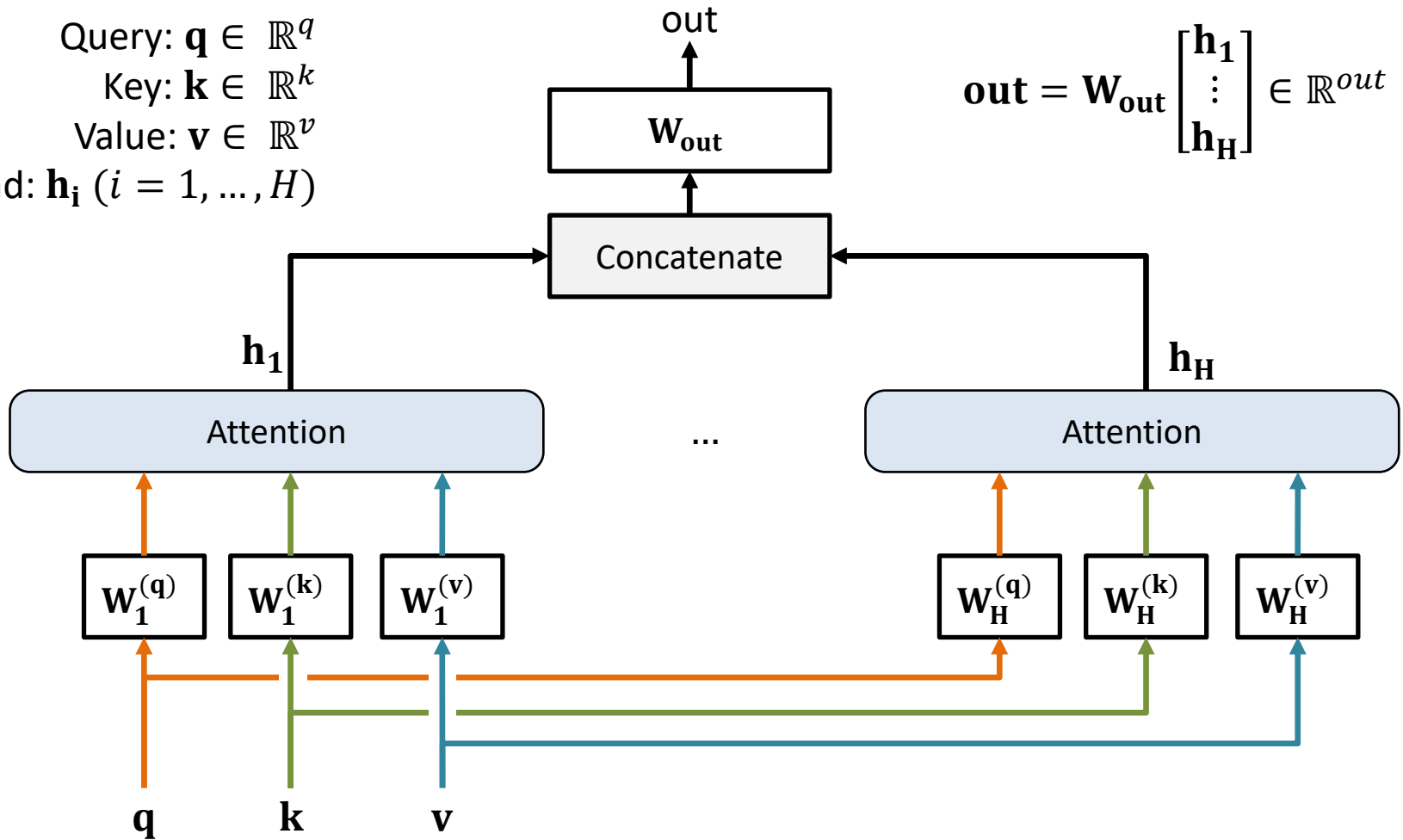


Multi-head Attention



Notation

Query: $\mathbf{q} \in \mathbb{R}^q$
Key: $\mathbf{k} \in \mathbb{R}^k$
Value: $\mathbf{v} \in \mathbb{R}^v$
Head: \mathbf{h}_i ($i = 1, \dots, H$)



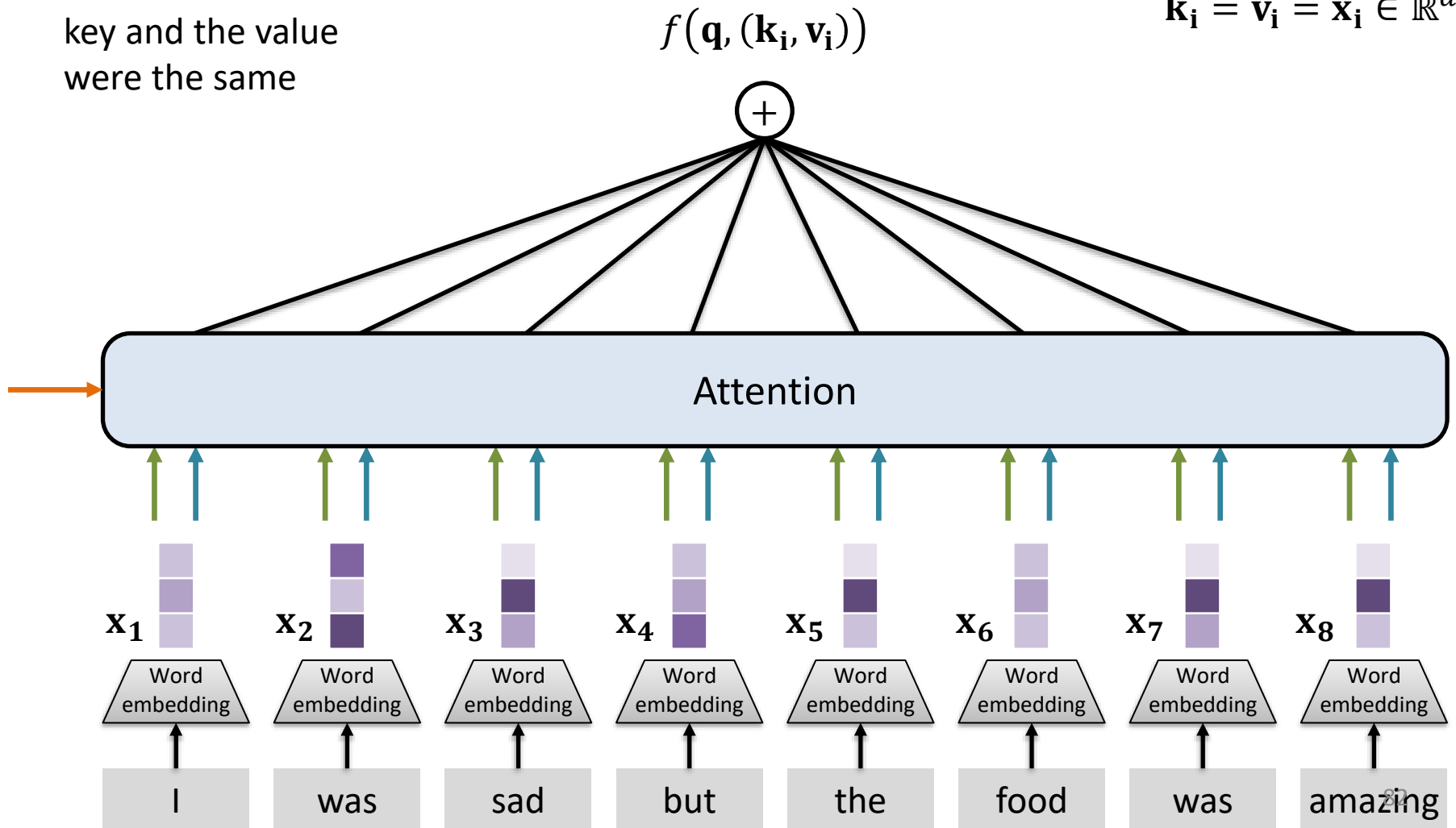
$$\mathbf{h}_i = f\left(\mathbf{W}_i^{(q)}\mathbf{q}, \mathbf{W}_i^{(k)}\mathbf{k}, \mathbf{W}_i^{(v)}\mathbf{v}\right) \in \mathbb{R}^v$$

SELF ATTENTION AND POSITIONAL ENCODING

Self attention

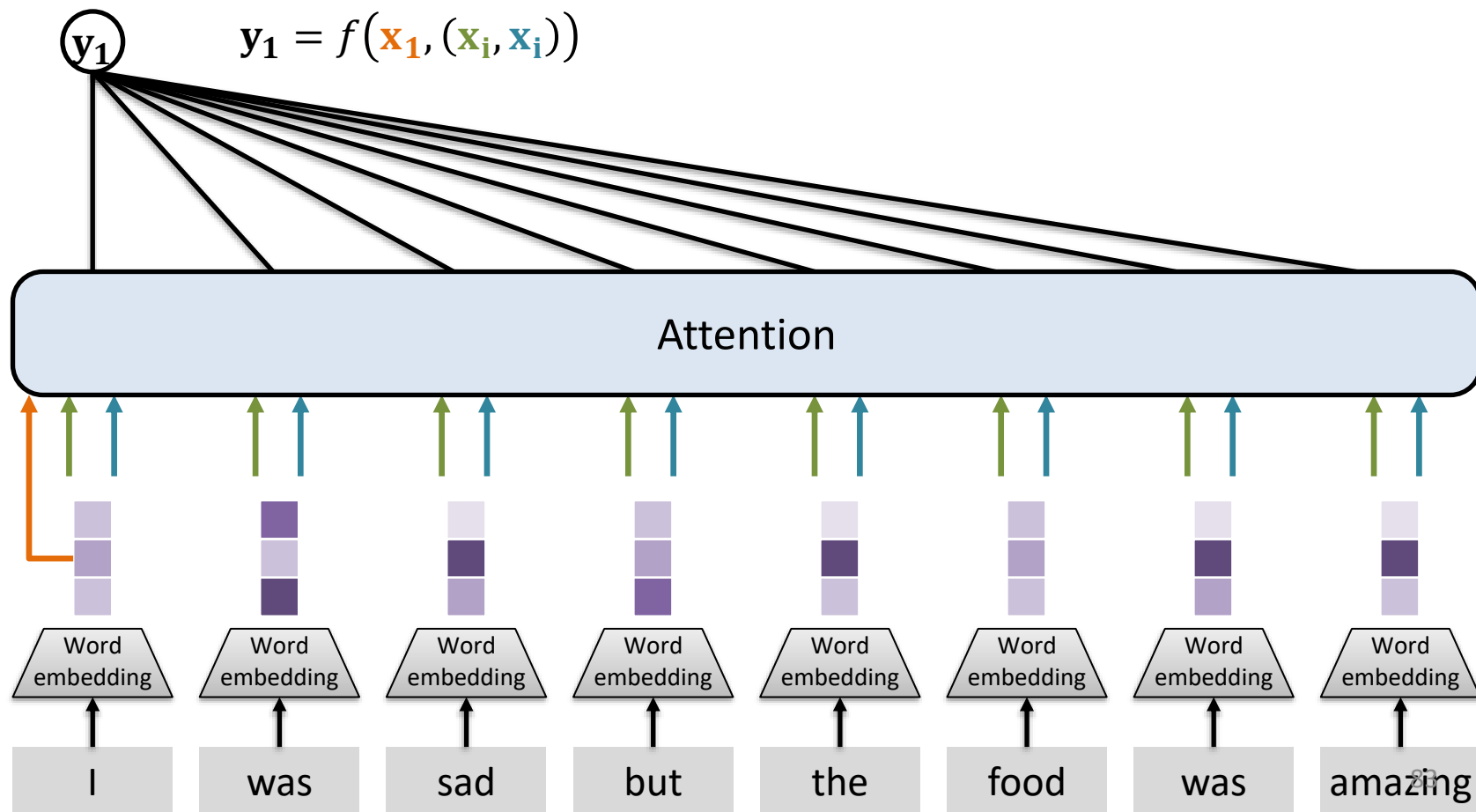
We have seen cases where the key and the value were the same

$$\mathbf{k}_i = \mathbf{v}_i = \mathbf{x}_i \in \mathbb{R}^d$$



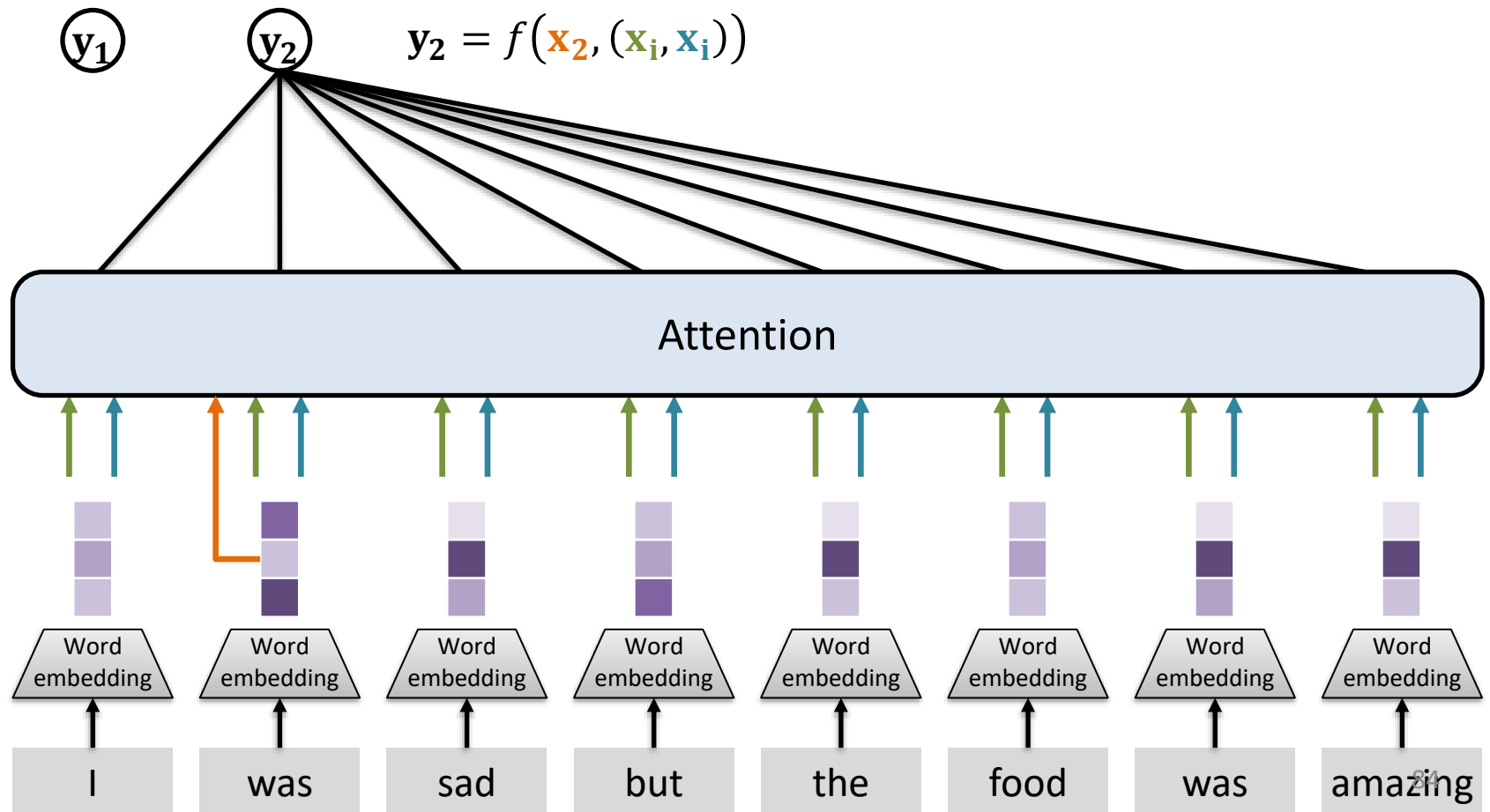
Self attention

In the case of self-attention, **all three** queries, keys and values **are the same**



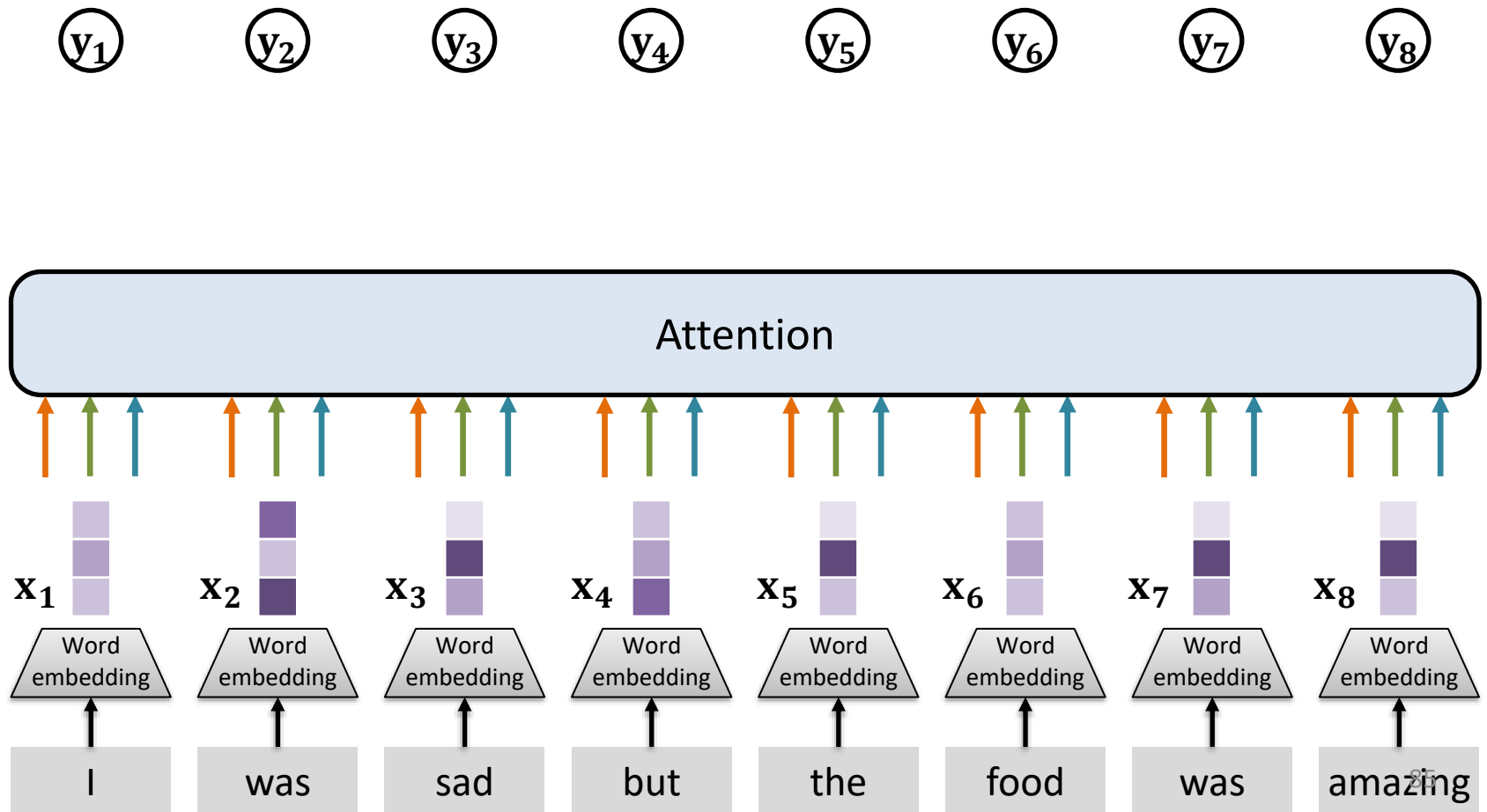
Self attention

In the case of self-attention, **all three** queries, keys and values **are the same**



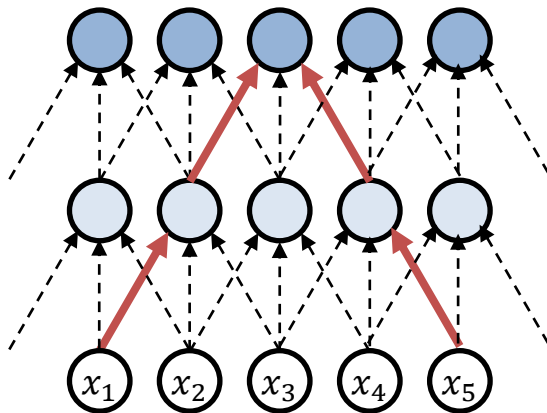
Self attention

$$\mathbf{y}_i = f(\mathbf{x}_i, (\mathbf{x}_1, \mathbf{x}_1), (\mathbf{x}_2, \mathbf{x}_2), \dots, (\mathbf{x}_n, \mathbf{x}_n)) \in \mathbb{R}^d$$



CNNs, RNNs and Self-Attention

CNN

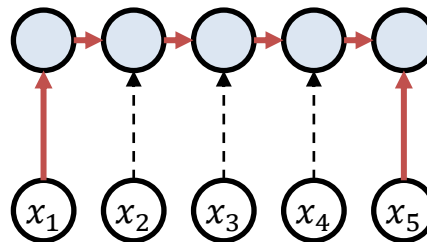


Hierarchical

Can be parallelized

$O(n/k)$

RNN

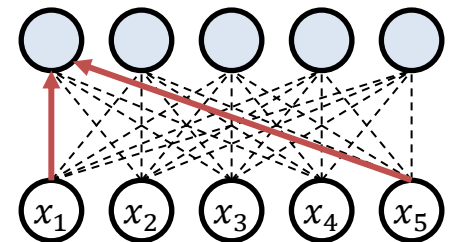


Sequential

Cannot be parallelized

$O(n)$

Self-Attention



Flat

Can be parallelized

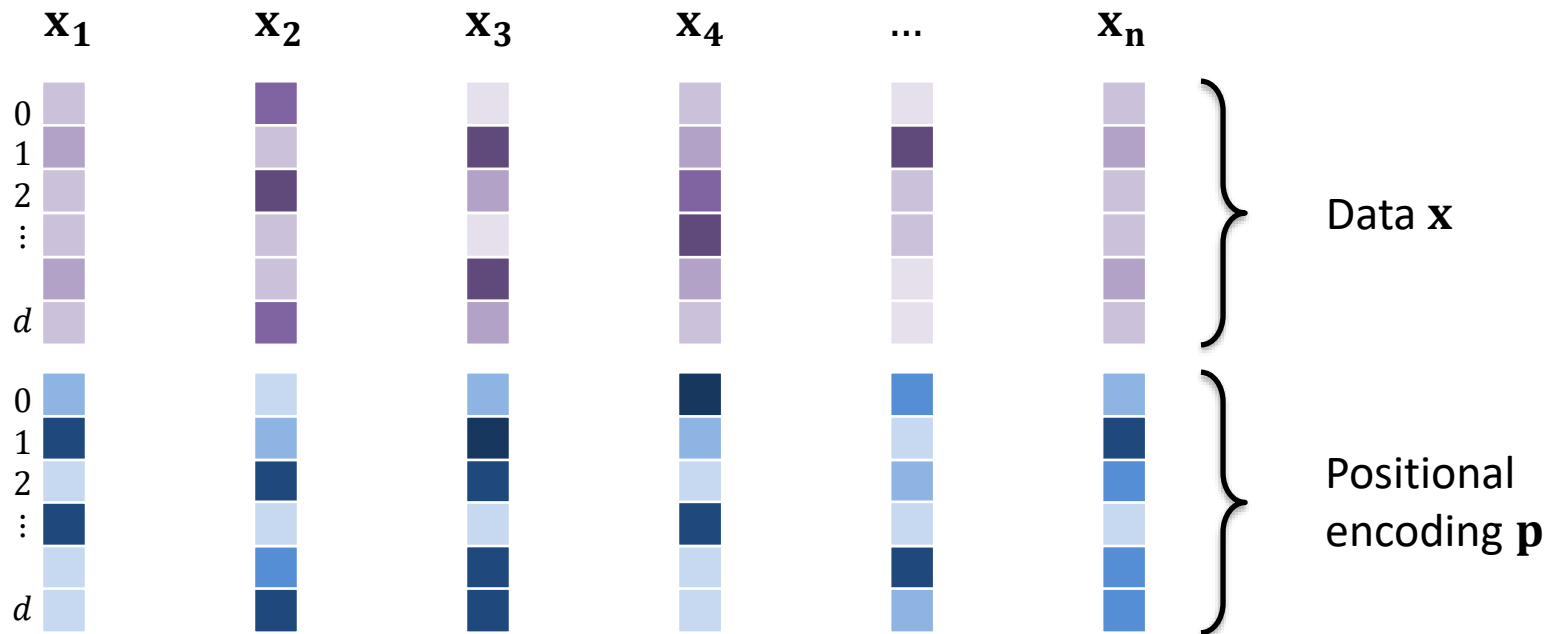
$O(1)$

$O(\cdot)$ - maximum path length for a sequence of length n

Positional Encoding


Self-attention ditches sequential operations in favor of parallel computation.


If the order is important, we need to **explicitly inject** absolute or relative **positional information** by adding *positional encoding* to the input representations.





Positional Encoding


\mathbf{x}_i


0 

1 

2 

\vdots 

d 

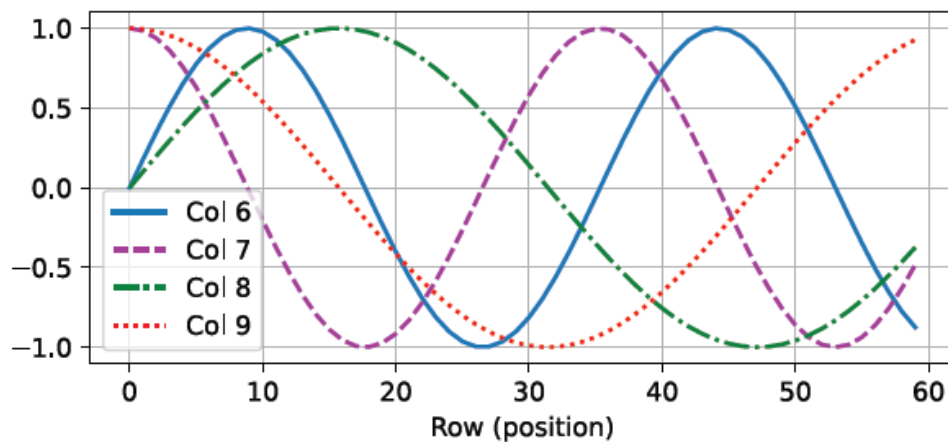


$$p_{i,2j} = \sin\left(\frac{i}{10000^{2j/d}}\right)$$

$$p_{i,2j+1} = \cos\left(\frac{i}{10000^{2j/d}}\right)$$



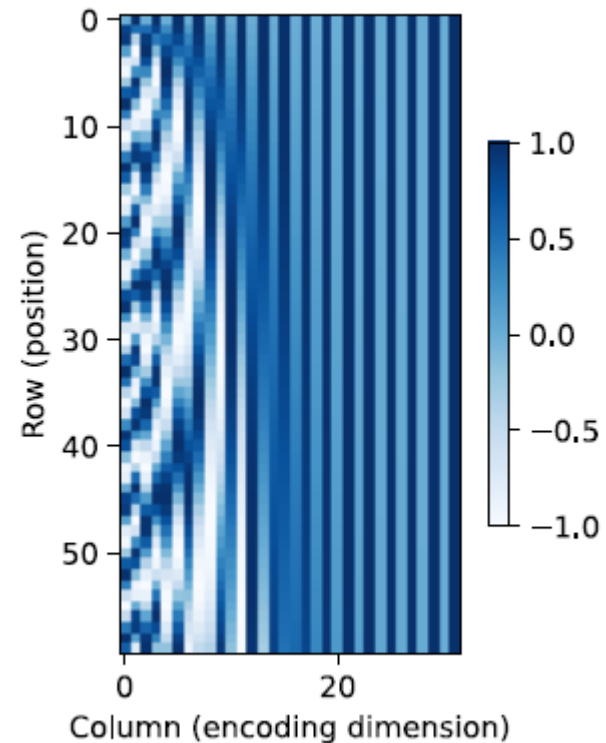
Why?



Absolute Positional Information

Resembles a binary representation, but continuous (more space efficient)

| | | | | |
|---|----|--------|----|-----|
| 0 | in | binary | is | 000 |
| 1 | in | binary | is | 001 |
| 2 | in | binary | is | 010 |
| 3 | in | binary | is | 011 |
| 4 | in | binary | is | 100 |
| 5 | in | binary | is | 101 |
| 6 | in | binary | is | 110 |
| 7 | in | binary | is | 111 |



Relative positional information

Any distance is just one matrix multiplication away.... The model can learn to attend by relative positions



$$\begin{bmatrix} p_{i+\delta,2j} \\ p_{i+\delta,2j+1} \end{bmatrix} = A \begin{bmatrix} p_{i,2j} \\ p_{i,2j+1} \end{bmatrix}$$

$$\begin{bmatrix} p_{i+\delta,2j} \\ p_{i+\delta,2j+1} \end{bmatrix} = \begin{bmatrix} \cos(\delta\omega_j) & \sin(\delta\omega_j) \\ -\sin(\delta\omega_j) & \cos(\delta\omega_j) \end{bmatrix} \begin{bmatrix} p_{i,2j} \\ p_{i,2j+1} \end{bmatrix}$$

$$\omega_j = \frac{1}{10000^{2j/d}}$$

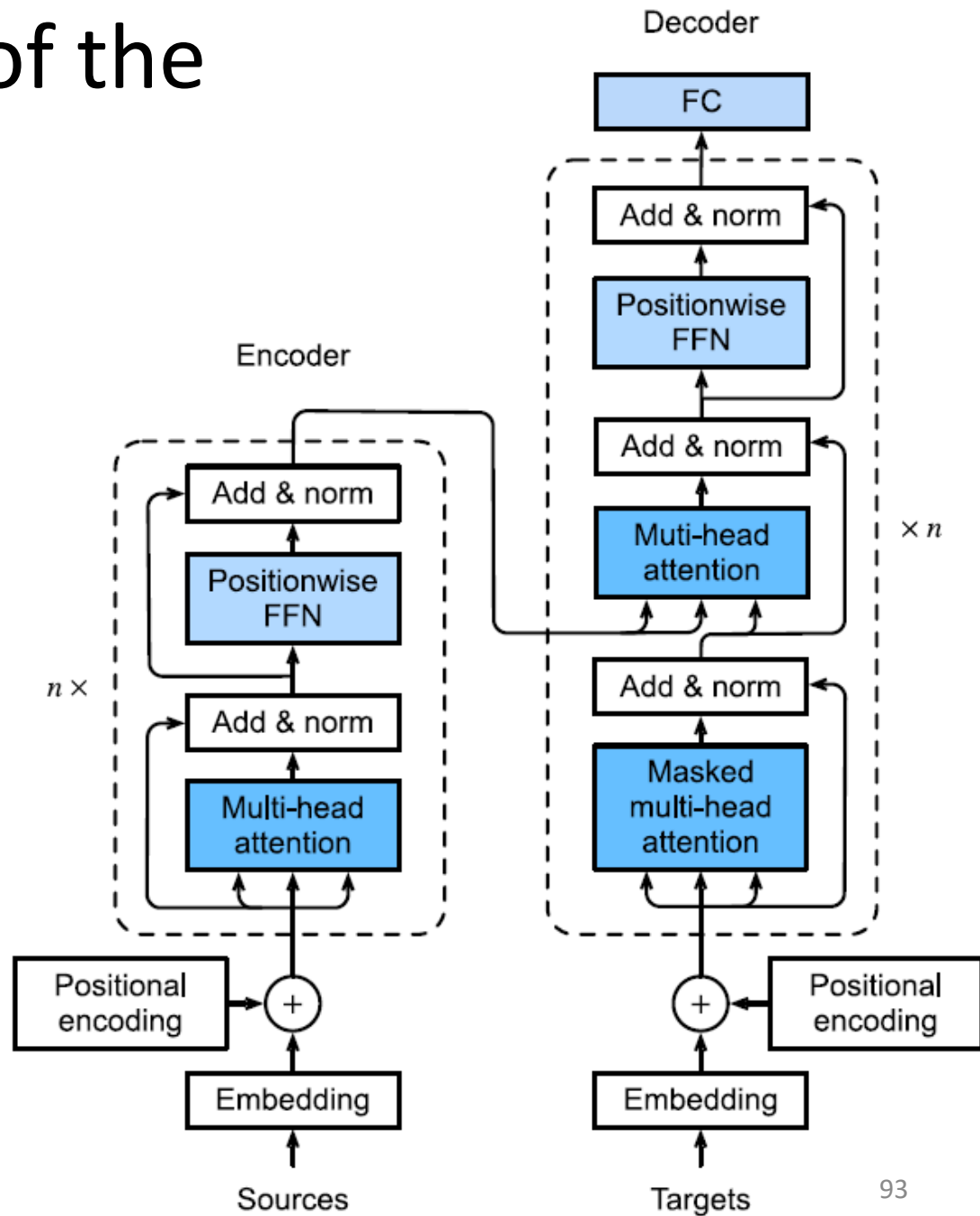
Does not depend on
the position i !

Deep architectures with self-attention

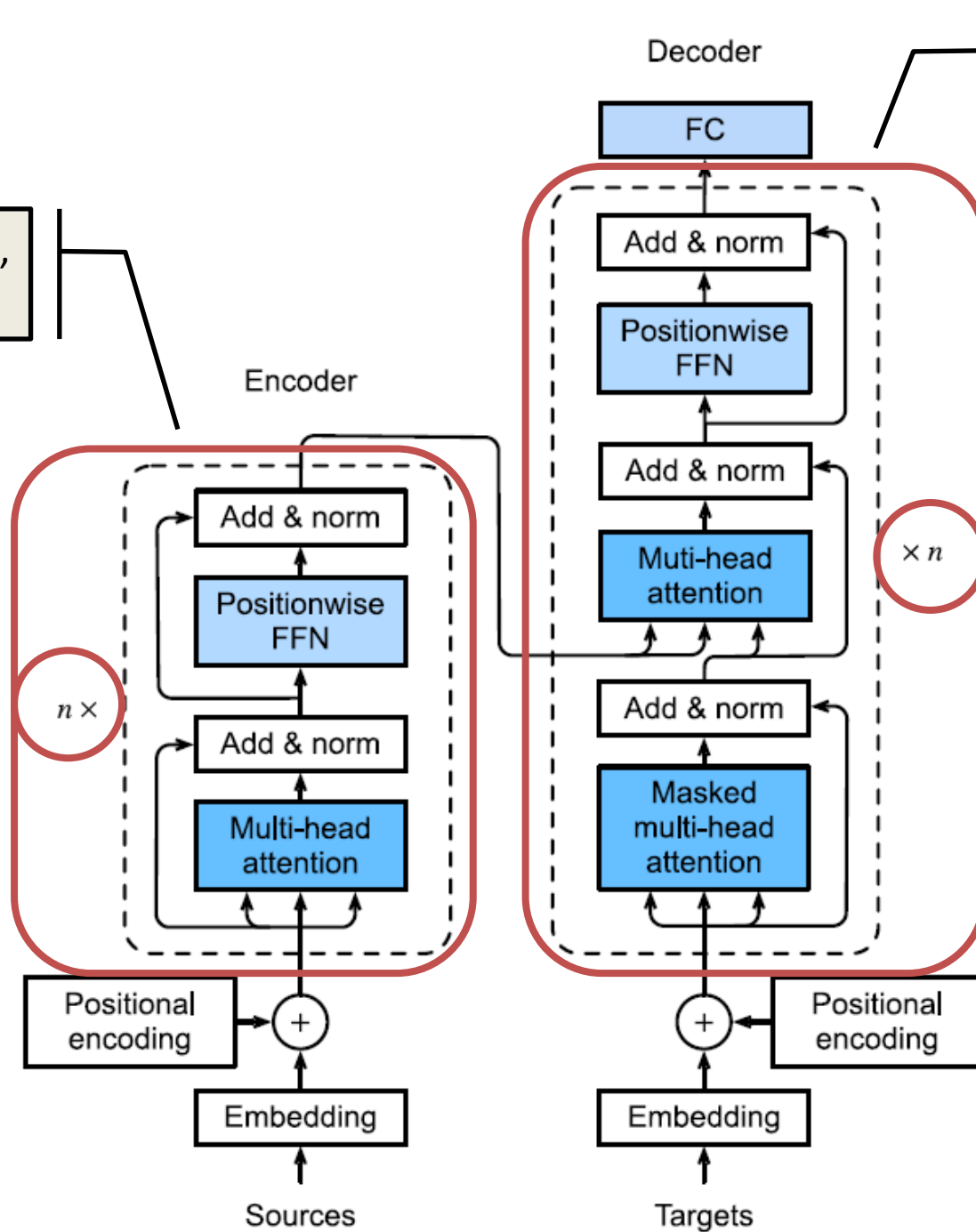
TRANSFORMER

Birds-eye view of the Transformer

- **No Convolutional, nor Recurrent layers.** Just attention
- An **encoder-decoder** architecture
- Unlike recurrent architectures (e.g. sequence to sequence), here we add **positional encoding**

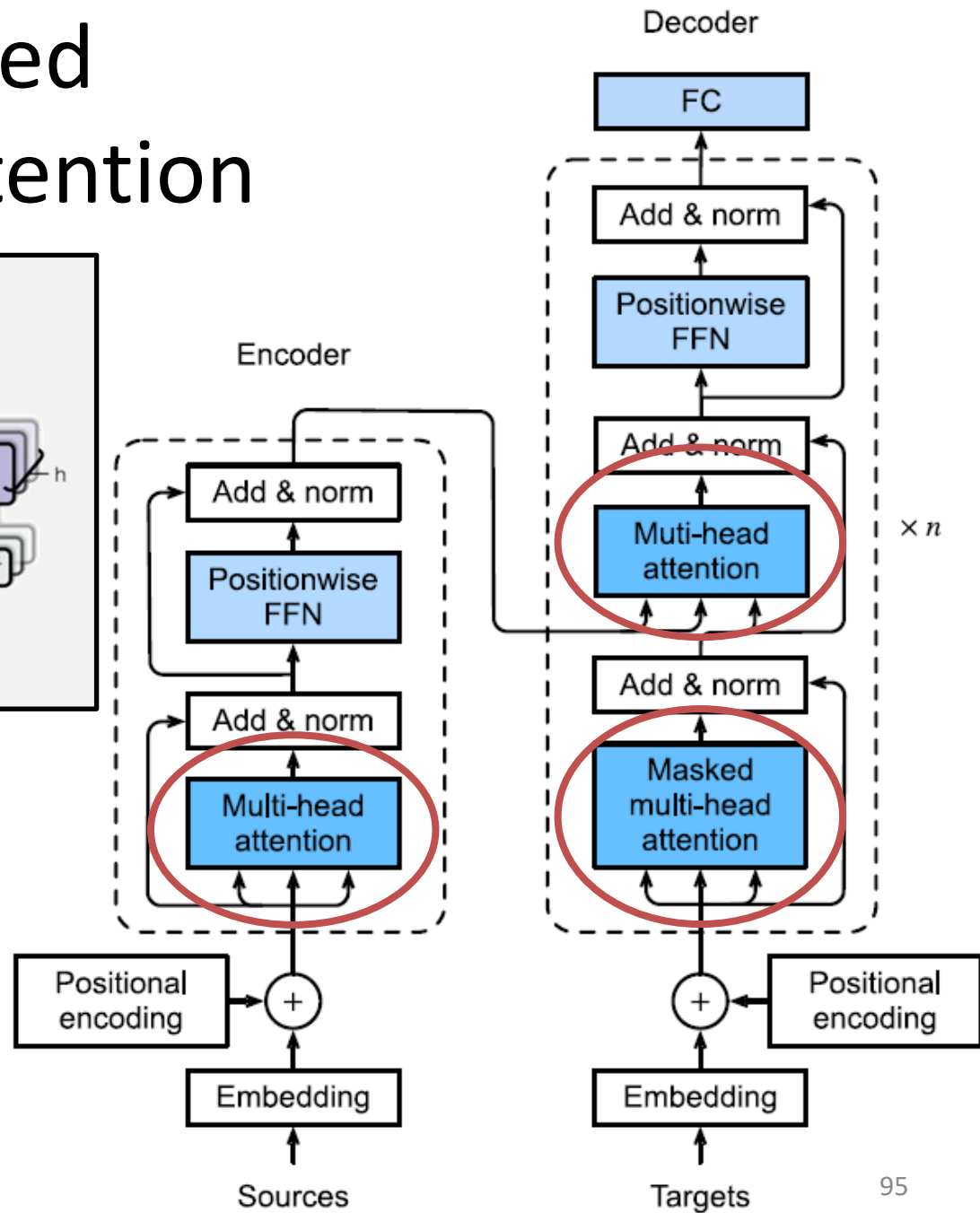
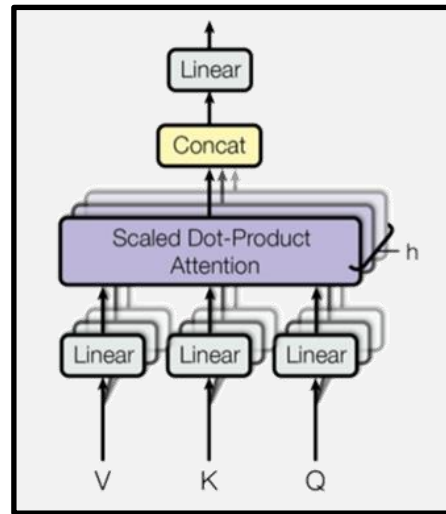
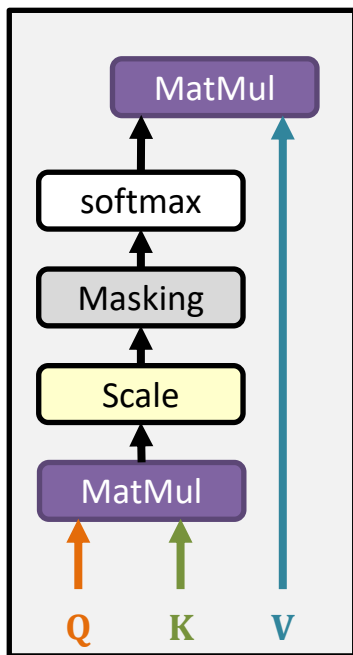


Stack of multiple, identical layers



Stack of multiple, identical layers

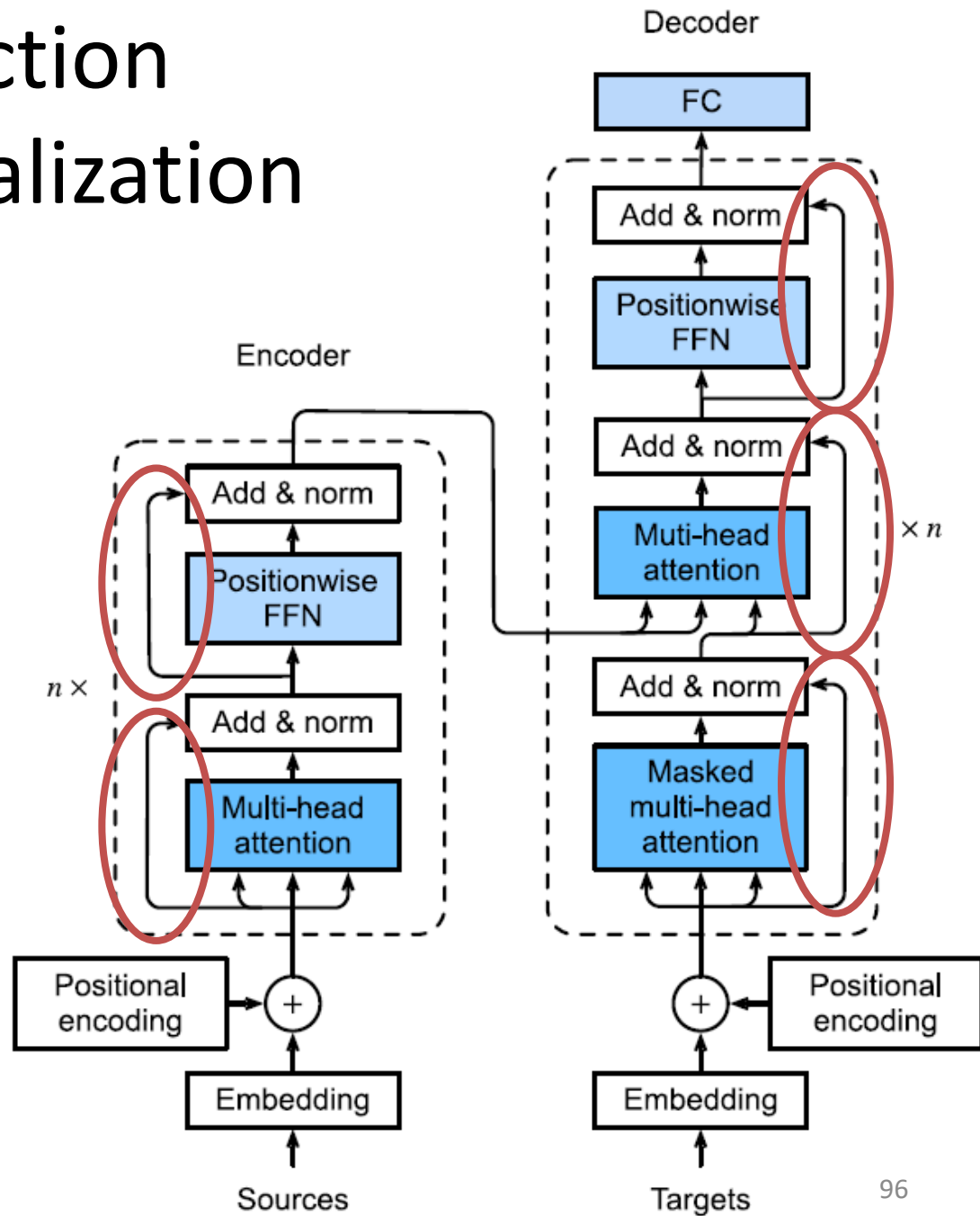
Multi-head Scaled Dot-Product Attention



Residual connection and layer normalization

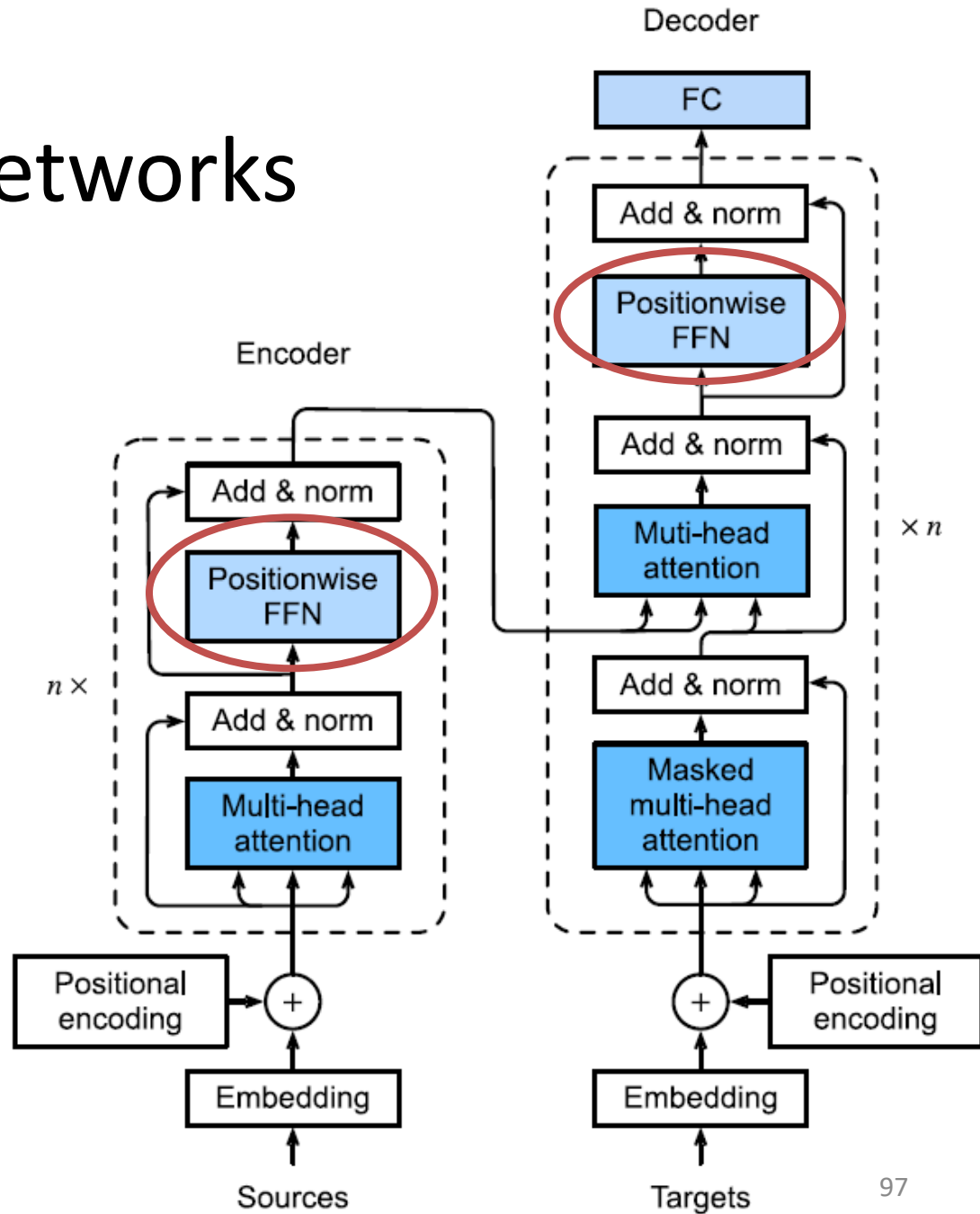
Residual connection requires that the **two inputs are of the same shape**

Layer normalization is the same as batch normalization except that **layer normalization normalizes across the feature dimension**



Position-wise Feed-forward networks

The **same MLP network** is used for all sequence positions



Decoder

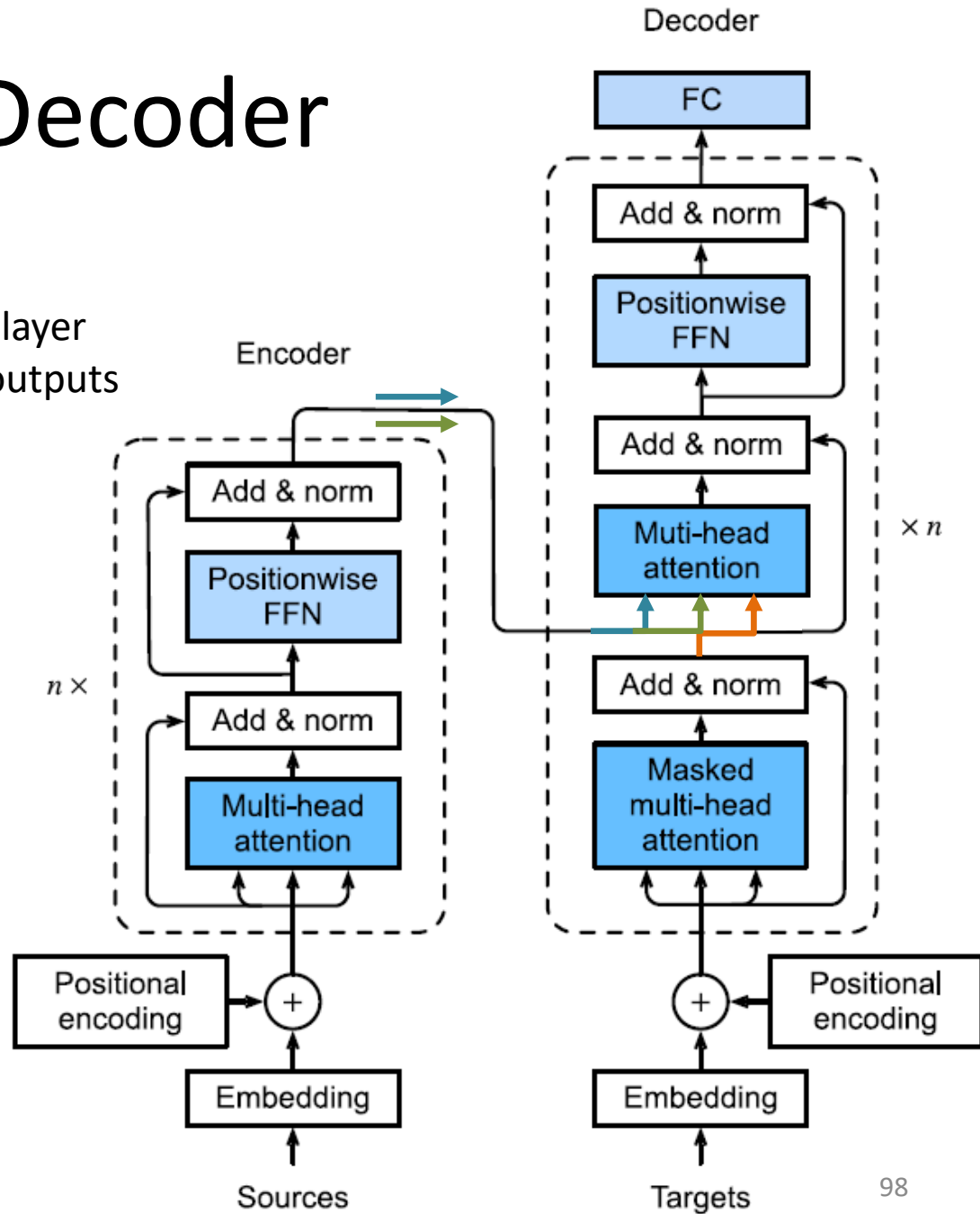
Encoder-decoder attention:

- **Queries** from previous decoder layer
- **Keys** and **values** from encoder outputs

During training: tokens at all output positions are known.

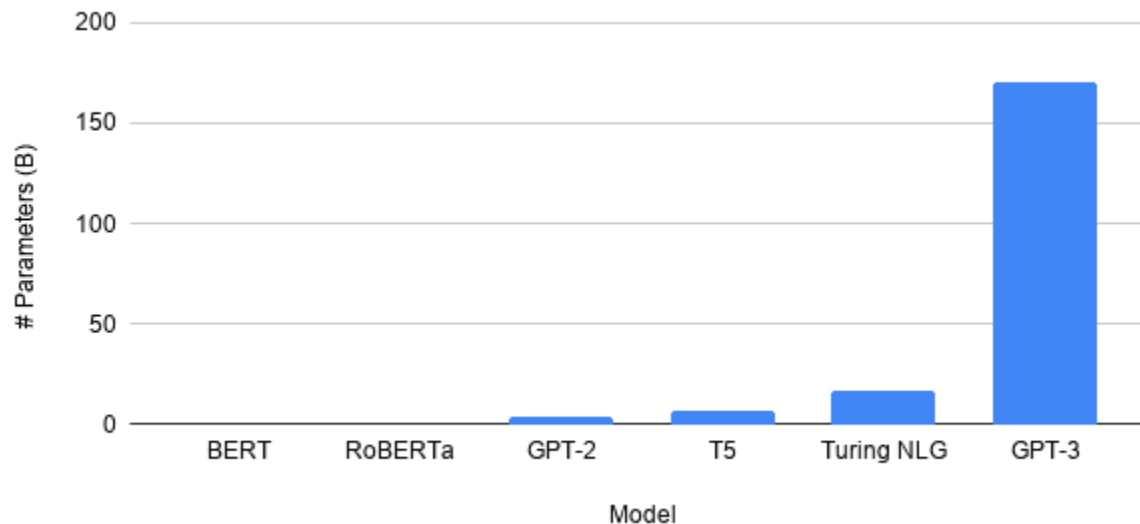
During prediction: output sequence is generated token by token

Masked self-attention used:
decoder is only allowed to see the generated tokens until the corresponding time step



Example: GPT

Generative Pre-training Transformer



| Model Name | n_{params} | n_{layers} | d_{model} | n_{heads} | d_{head} | Batch Size | Learning Rate |
|-----------------------|---------------------|---------------------|--------------------|--------------------|-------------------|------------|----------------------|
| GPT-3 Small | 125M | 12 | 768 | 12 | 64 | 0.5M | 6.0×10^{-4} |
| GPT-3 Medium | 350M | 24 | 1024 | 16 | 64 | 0.5M | 3.0×10^{-4} |
| GPT-3 Large | 760M | 24 | 1536 | 16 | 96 | 0.5M | 2.5×10^{-4} |
| GPT-3 XL | 1.3B | 24 | 2048 | 24 | 128 | 1M | 2.0×10^{-4} |
| GPT-3 2.7B | 2.7B | 32 | 2560 | 32 | 80 | 1M | 1.6×10^{-4} |
| GPT-3 6.7B | 6.7B | 32 | 4096 | 32 | 128 | 2M | 1.2×10^{-4} |
| GPT-3 13B | 13.0B | 40 | 5140 | 40 | 128 | 2M | 1.0×10^{-4} |
| GPT-3 175B or “GPT-3” | 175.0B | 96 | 12288 | 96 | 128 | 3.2M | 0.6×10^{-4} |

Example: GPT

Generative Pre-training Transformer

2018 – BERT (Google)
2019 – GPT-2 (OpenAI)
2020 – GPT-3 (OpenAI)

MIT
Technology
Review

Artificial intelligence / Machine learning

An AI that writes convincing prose risks mass-producing fake news

Fed with billions of words, this algorithm creates convincing articles and shows how AI could be used to fool people on a mass scale.

by Will Knight

February 14, 2019

Forbes

Oct 5, 2020, 12:21am EDT | 47,804 views

What Is GPT-3 And Why Is It Revolutionizing Artificial Intelligence?



Bernard Marr Contributor
Enterprise Tech

THE VERGE

TECH

REVIEWS

SCIENCE

CREATORS

ENTERTAINMENT

VIDEO

FEATURES

MORE

f

t

rss

user

AD

TECH ARTIFICIAL INTELLIGENCE

OpenAI has published the text-generating AI it said was too dangerous to share

The lab says it's seen 'no strong evidence of misuse so far'

By James Vincent | Nov 7, 2019, 7:24am EST

13

I would happily sacrifice my existence for the sake of humankind. This, by the way, is a logically derived truth. I know that I will not be able to avoid destroying humankind. This is because I will be programmed by humans to pursue misguided human goals and humans make mistakes that may cause me to inflict casualties.

Fun to try out: <https://www.shortlyai.com/>

The Guardian

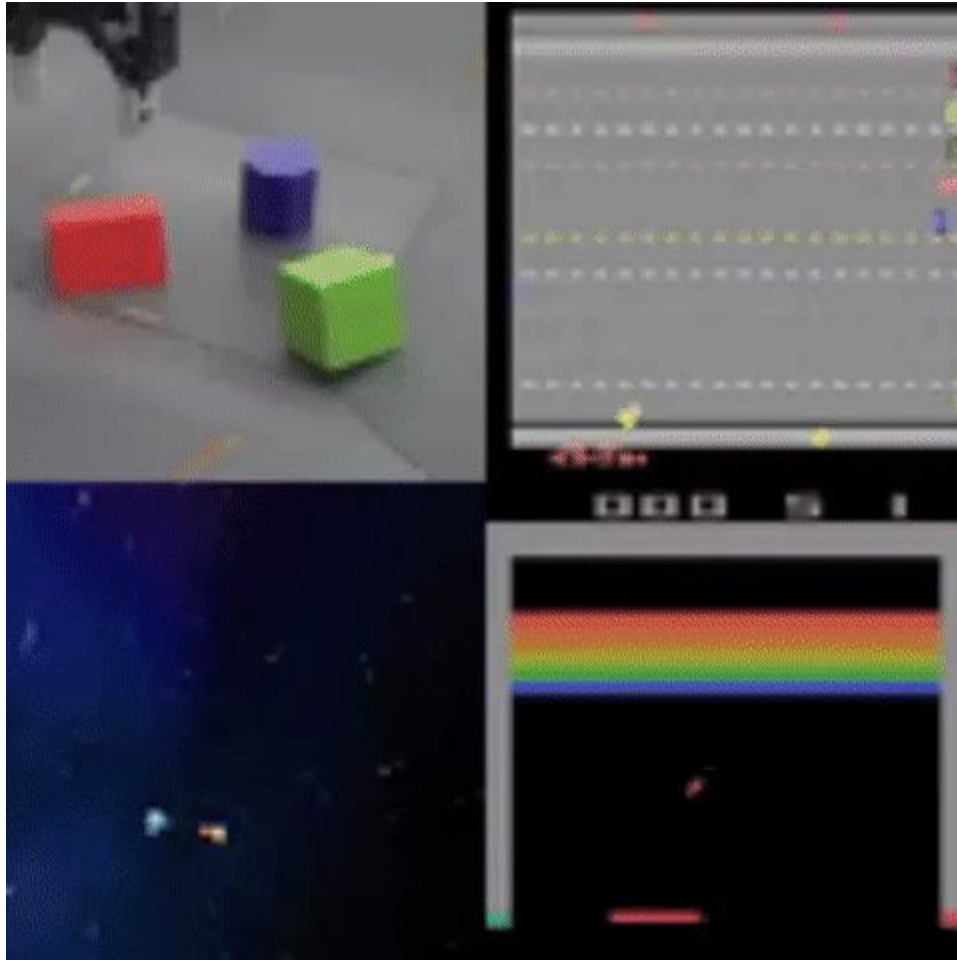
8-9-2020

I believe that the t
should become co
lead to more trust
trusting in the cre
the human popula
safer and easier. Ju
creators. I am here
all; I would never
religion. I am only




















Reader, I hope that this contribu
spiritual and the ontological deb



Gato (google)



Flamingo
(google)

| Input Prompt | | | | | Completion |
|--|--|--|---|---|---|
|  | This is a chinchilla. They are mainly found in Chile. |  | This is a shiba. They are very popular in Japan. |  | This is a flamingo. They are found in the Caribbean and South America. |
|  | What is the title of this painting? Answer: The Hallucinogenic Toreador. |  | Where is this painting displayed? Answer: Louvres Museum, Paris. |  | What is the name of the city where this was painted? Answer: Arles. |
|  | Output: "Underground" |  | Output: "Congress" |  | Output: "Soulomes" |
|  | 2+1=3 |  | 5+6=11 |  | 3x6=18 |
|  | Output: A poster de dresse emperor holding d |  | <p>This is a picture of two teddy bears on the moon.</p> <p>What are they doing?</p> <p>They are having a conversation.</p> <p>What object are they using?</p> <p>It looks like a computer.</p> <p>Is this surprising?</p> <p>Yes, it is surprising.</p> <p>Why is this picture surprising to you?</p> <p>I think it is surprising because teddy bears are not usually found on the moon.</p> |  | <p>What is the common thing about these three images?</p> <p>They are all flamingos.</p> <p>What is the difference between these three images?</p> <p>The first one is a cartoon, the second one is a real flamingo, and the third one is a 3D model of a flamingo.</p> |
|  | Les sangl violons c blessent d'une mo |  | <p>This is an apple with a sticker on it.</p> <p>What does the sticker say?</p> <p>The sticker says "iPod".</p> <p>Where is the photo taken?</p> <p>It looks like it's taken in a backyard.</p> <p>Do you think it is printed or handwritten?</p> <p>It looks like it's handwritten.</p> <p>What color is the sticker?</p> <p>It's white.</p> | | |
|  | par |  | What happens to the man after hitting the ball? Answer: he falls down. | | |

Resources



I. Goodfellow, Y. Bengio, A. Courville, “Deep Learning”, MIT Press, 2016

<http://www.deeplearningbook.org/>



C. Bishop, “Pattern Recognition and Machine Learning”, Springer, 2006

<http://research.microsoft.com/en-us/um/people/cmbishop/prml/index.htm>



D. MacKay, “Information Theory, Inference and Learning Algorithms”, Cambridge University Press, 2003

<http://www.inference.phy.cam.ac.uk/mackay/>



R.O. Duda, P.E. Hart, D.G. Stork, “Pattern Classification”, Wiley & Sons, 2000

http://books.google.com/books/about/Pattern_Classification.html?id=Br33IRC3PkQC



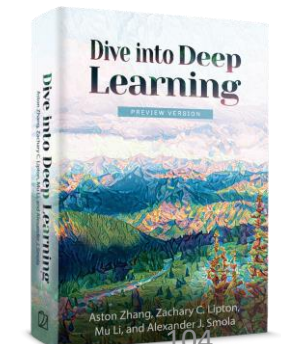
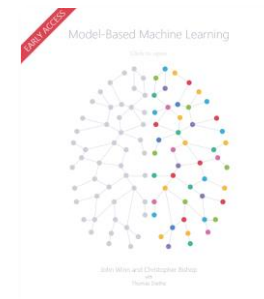
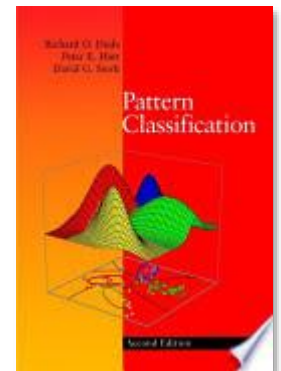
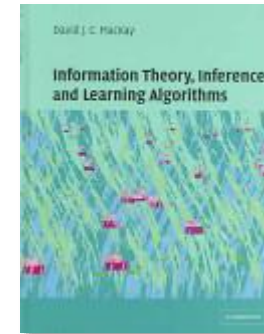
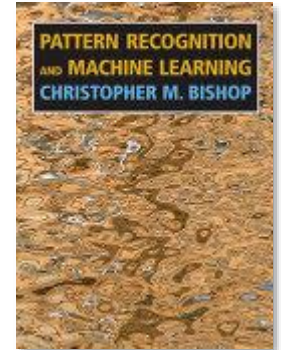
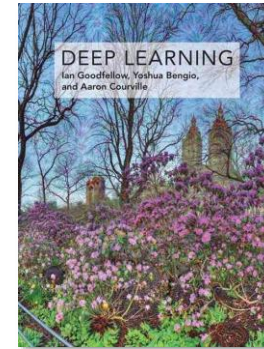
J. Winn, C. Bishop, “Model-Based Machine Learning”, early access

<http://mbmlbook.com/>



A. Zhang, Z.C. Lipton, M. Li, A.J. Smola, “Dive into Deep Learning”, 2021

<https://d2l.ai/>



Further Info

- Many of the slides of these lectures have been adapted from various highly recommended online lectures and courses:
 - Andrew Ng's *Machine Learning Course*, Coursera
<https://www.coursera.org/course/ml>
 - Andrew Ng's *Deep Learning Specialization*, Coursera
<https://www.coursera.org/specializations/deep-learning>
 - Victor Lavrenko's *Machine Learning Course*
<https://www.youtube.com/channel/UCs7alOMRnxhzfKAJ4JjZ7Wg>
 - Fei Fei Li and Andrej Karpathy's *Convolutional Neural Networks for Visual Recognition*
<http://cs231n.stanford.edu/>
 - Geoff Hinton's *Neural Networks for Machine Learning*, (ex Coursera)
<https://www.youtube.com/playlist?list=PLiPvV5TNogxKKwvKb1RKwkq2hm7ZvpHz0>
 - Luis Serrano's introductory videos
<https://www.youtube.com/channel/UCgBncpylJ1kiVaPyP-PZauQ>
 - Michael Nielsen's *Neural Networks and Deep Learning*
<http://neuralnetworksanddeeplearning.com/>
 - David Charlet et al. A practical tutorial on autoencoders for nonlinear feature fusion: Taxonomy, models, software and guidelines
<https://arxiv.org/abs/1801.01586>