

Transformers

Jose González-Abad
Instituto de Física de Cantabria (CSIC-UC)

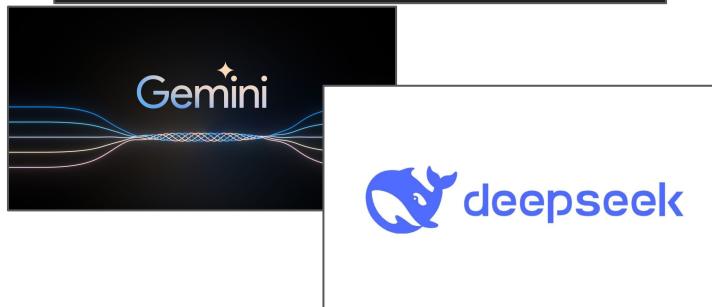
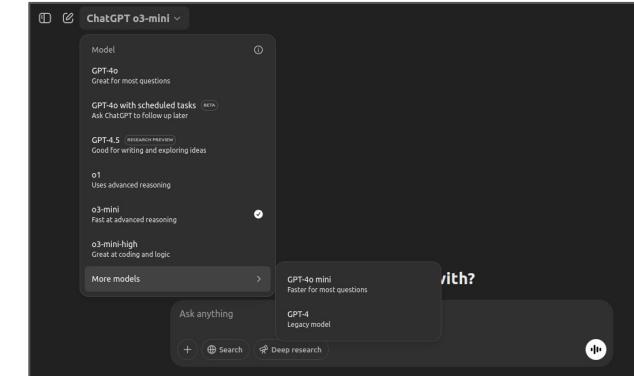
Master Data Science (UIMP-UC)
March 2025

 gonzabad@ifca.unican.es
 [jgonzalezab](https://github.com/jgonzalezab)

Origin of the Attention Mechanism

Most of the revolutionary advances in deep learning in recent years rely on **one simple yet powerful mechanism**

Large language models



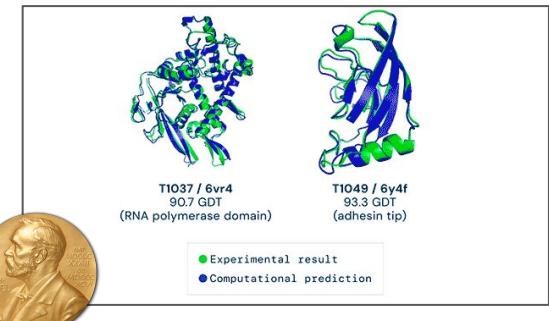
Video generation



Image segmentation



Protein folding



Origin of the Attention Mechanism



It all started in **2017** at the **Neural Information Processing System (NIPS) conference** in Long Beach, California



Where a group of researchers from **Google** and **UoT** presented
a **10-pages paper**

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer* **Niki Parmar***
Google Brain Google Research
noam@google.com 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

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

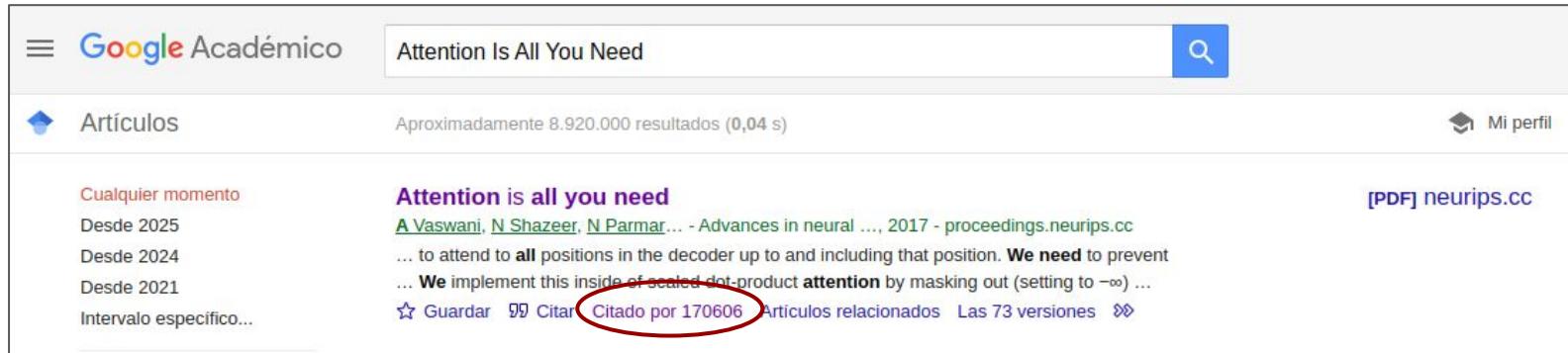
Illia Polosukhin* [†]
illia.polosukhin@gmail.com

Origin of the Attention Mechanism

Typically, research papers receive between 10 and 100 citations*

Origin of the Attention Mechanism

Typically, **research papers receive between 10 and 100 citations***; however, this particular paper has accumulated over **150,000 citations** to date



The screenshot shows a Google Scholar search results page for the paper "Attention Is All You Need". The search query is "Attention Is All You Need". The results count is approximately 8,920,000. The first result is the paper by Vaswani et al., published in 2017. A red oval highlights the citation count "Citado por 170606". Other visible elements include filters for "Artículos", search bar, user profile icon, and PDF download link.

Google Académico

Attention Is All You Need

Artículos

Aproximadamente 8.920.000 resultados (0,04 s)

Cualquier momento

Desde 2025

Desde 2024

Desde 2021

Intervalo específico...

Attention is all you need

A Vaswani, N Shazeer, N Parmar... - Advances in neural ..., 2017 - proceedings.neurips.cc

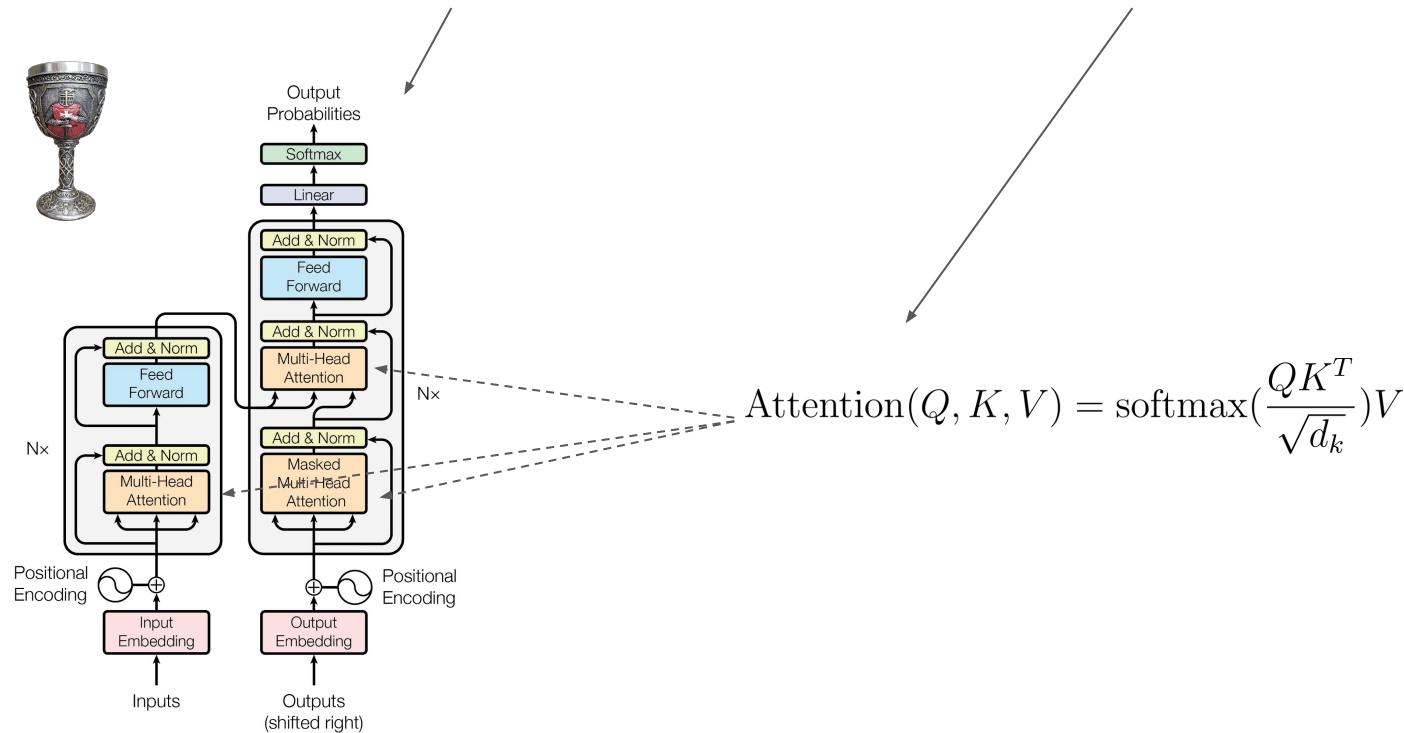
... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent
... **We** implement this inside of scaled dot-product **attention** by masking out (setting to $-\infty$) ...

☆ Guardar 99 Citar Citado por 170606 Artículos relacionados Las 73 versiones

[PDF] neurips.cc

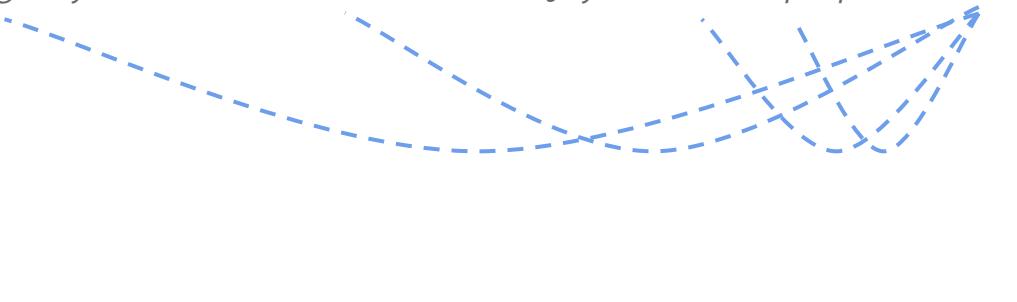
Origin of the Attention Mechanism

Its main contribution are **transformers**, which are based on the **attention mechanism**



What is attention?

After a long day at work, she sat down to enjoy a warm cup of -----



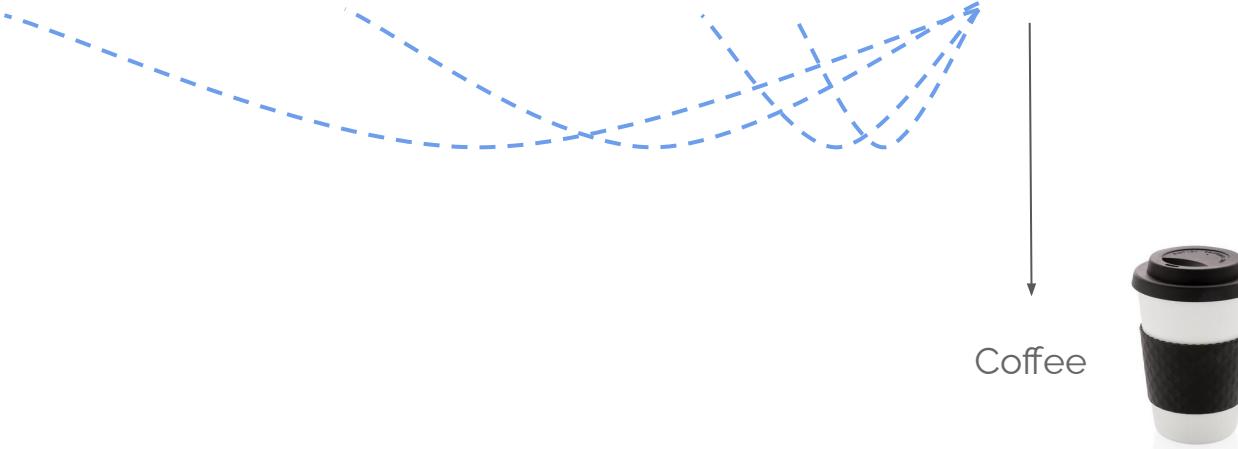
Coffee



The **attention mechanism** enables the deep learning model to determine **which parts of the context to focus on**

What is attention?

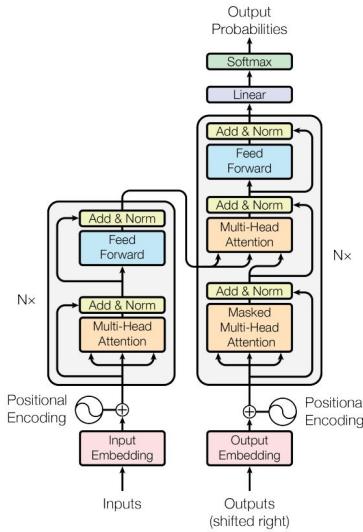
After a long day at work, she sat down to enjoy a warm cup of -----



In the following, we will explain how the transformer is applied to learn this task



Tokenization

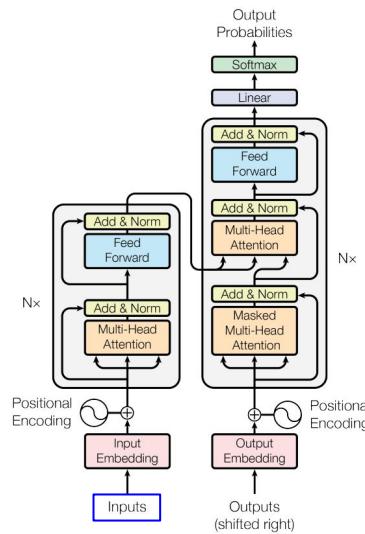


To understand the transformer architecture, we will examine each of its elements in detail as we complete the **following task**

After a long day at work, she sat down to enjoy a warm cup of -----

Figure 1: The Transformer - model architecture.

Tokenization



After a long day at work, she sat down to enjoy a warm cup of -----

22 0 432 150 15 1700 898 900 99 1200 501 0 1895 203 332

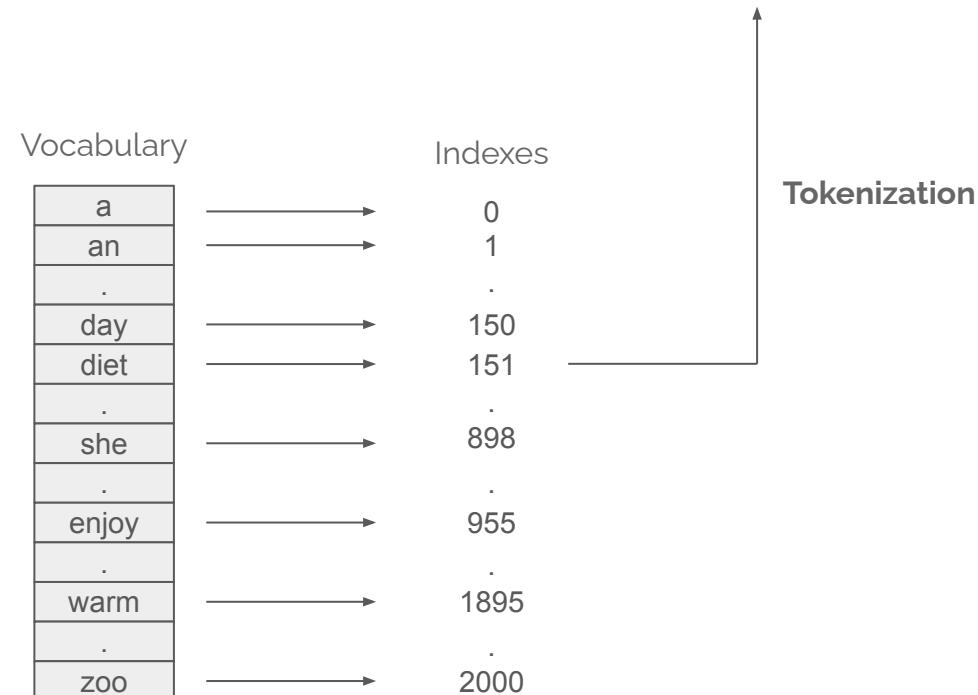
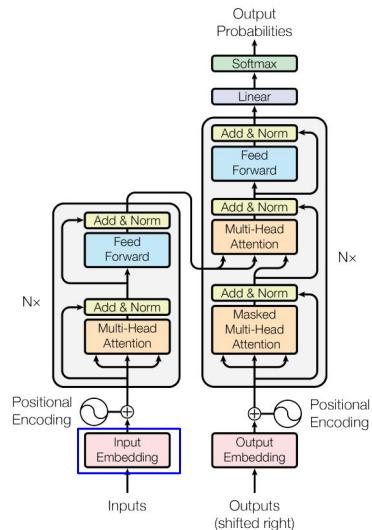


Figure 1: The Transformer - model architecture.

Input embedding



After a long day at work, she sat down to enjoy a warm cup of _____

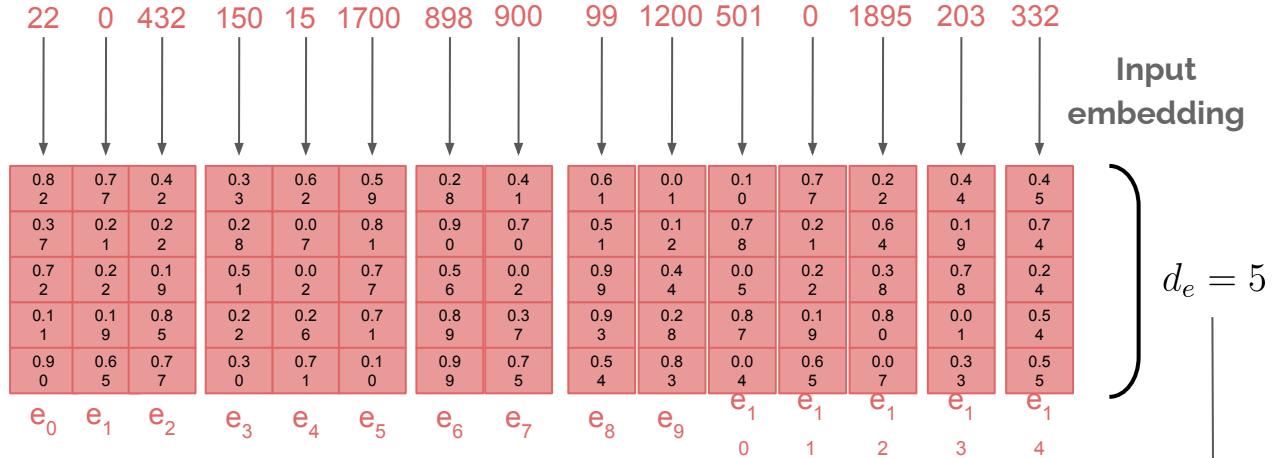
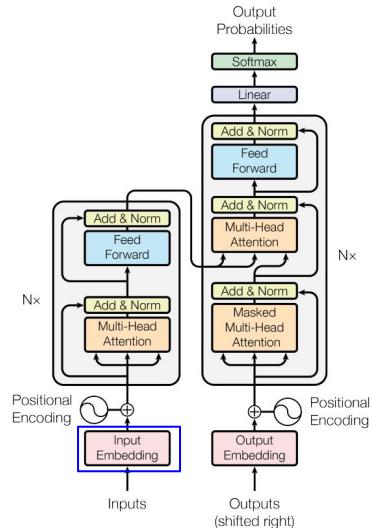


Figure 1: The Transformer - model architecture.

This embeddings are generally **learned** through a matrix \mathbf{W}_e

The embedding dimension tends to be higher

Input embedding



What do these embeddings **represent**?

$$E(\text{aunt}) - E(\text{uncle}) \approx E(\text{woman}) - E(\text{man})$$

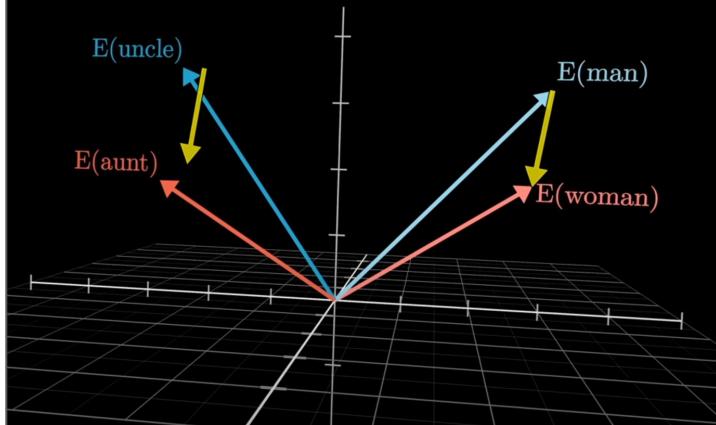
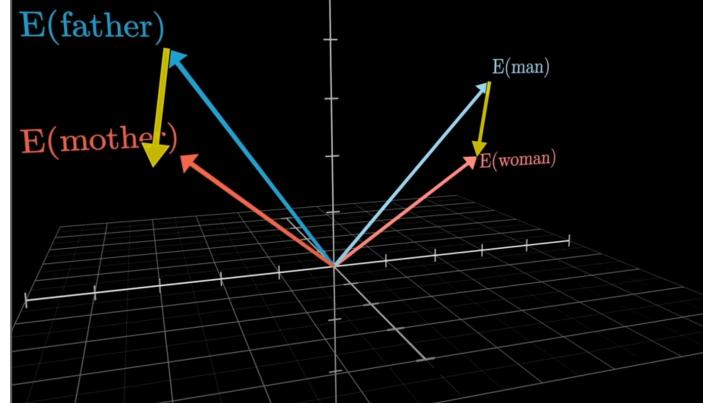
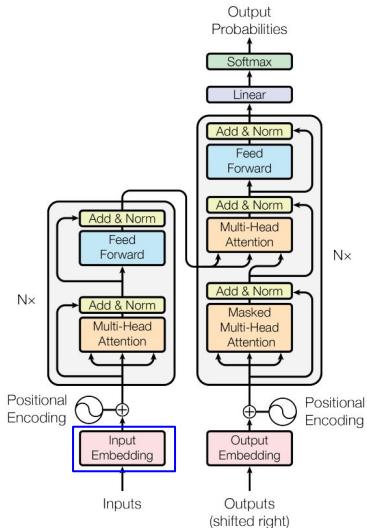


Figure 1: The Transformer - model architecture.

$$E(\text{mother}) - E(\text{father}) \approx E(\text{woman}) - E(\text{man})$$



Input embedding



What do these embeddings **represent**?

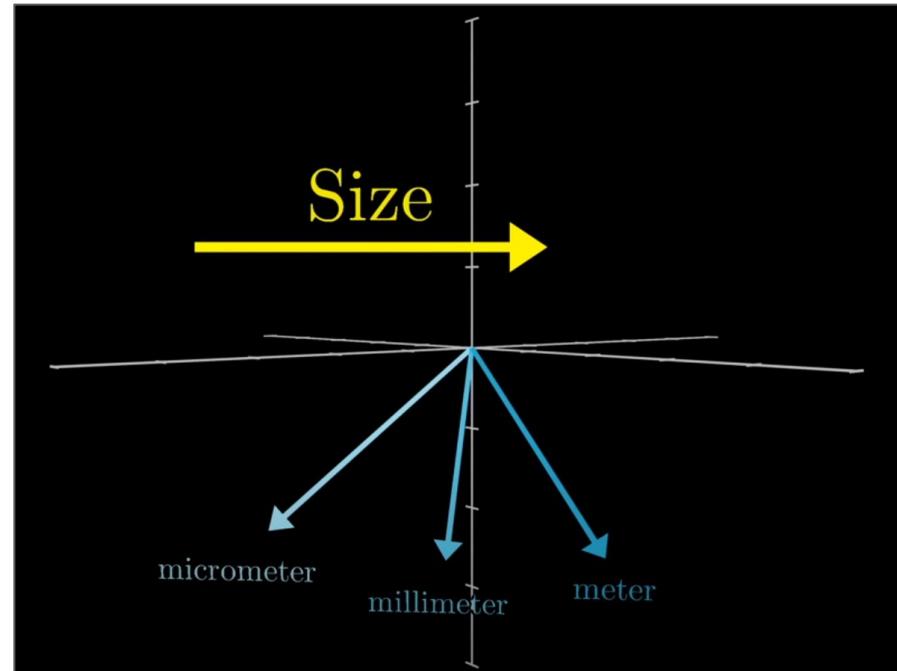
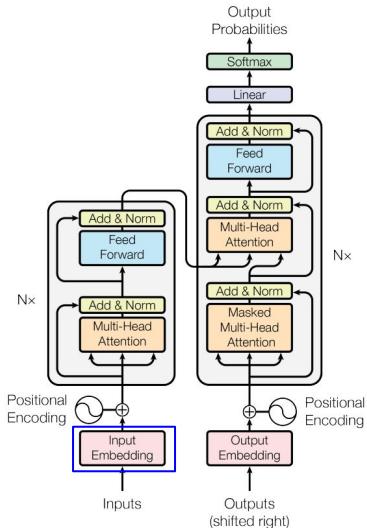


Figure 1: The Transformer - model architecture.

Input embedding



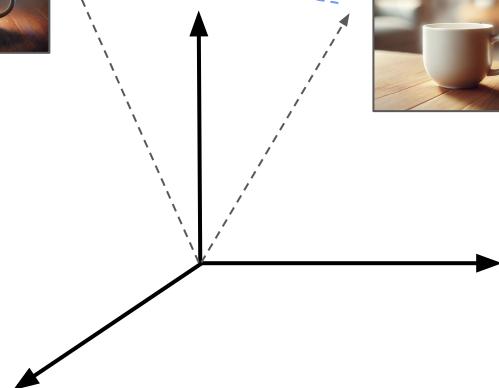
After applying the **attention mechanism**, we expect the embeddings to capture a more **complex representation of the context**

*After a long day at work, she sat down to enjoy a **warm cup** of -----*

Warm cup



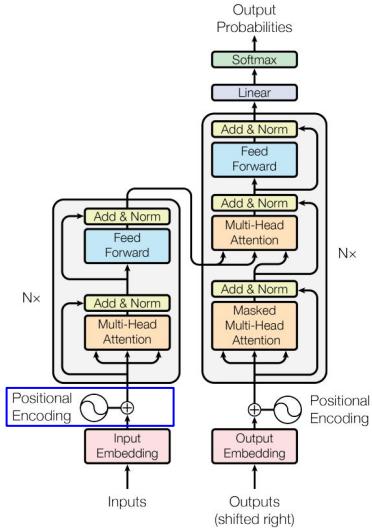
Cup



We **transform** the meaning of the **embedding** based on the **context**

Figure 1: The Transformer - model architecture.

Positional encoding



Other models that process sequences can take the order of inputs into account

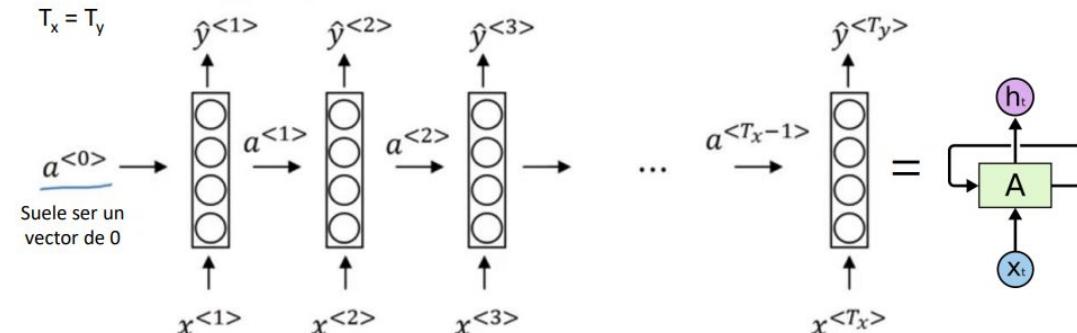
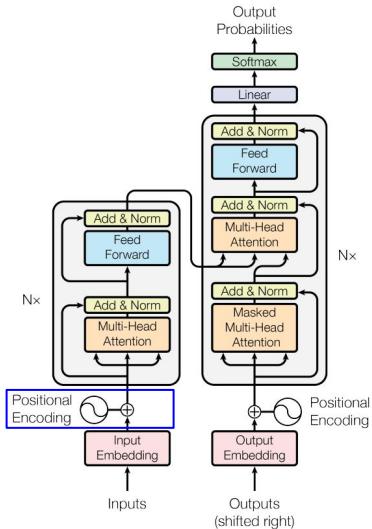


Figure 1: The Transformer - model architecture.

However, transformers consider **all inputs at once**

How can we **enforce an order on the inputs** for transformers?

Positional encoding



After a long day at work, she sat down to enjoy a warm cup of _____

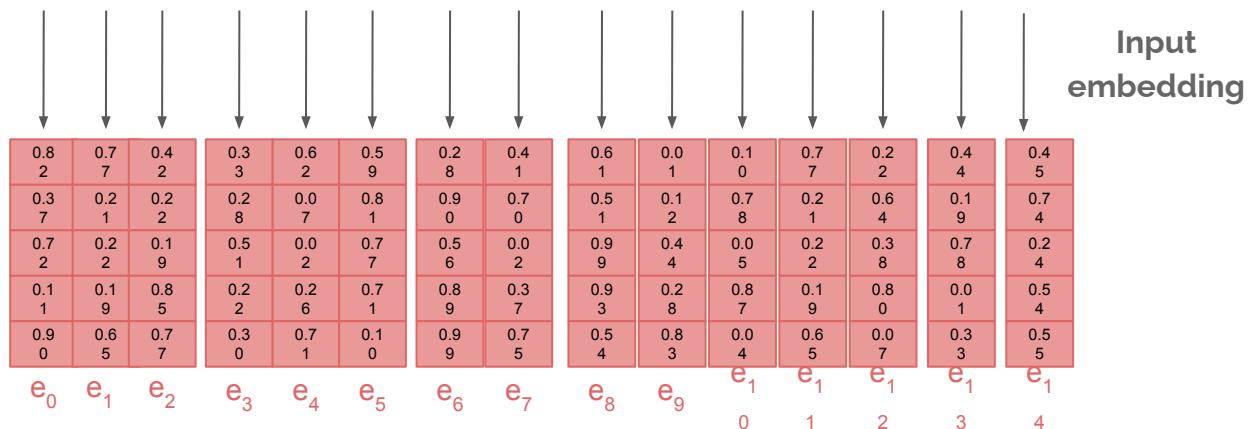
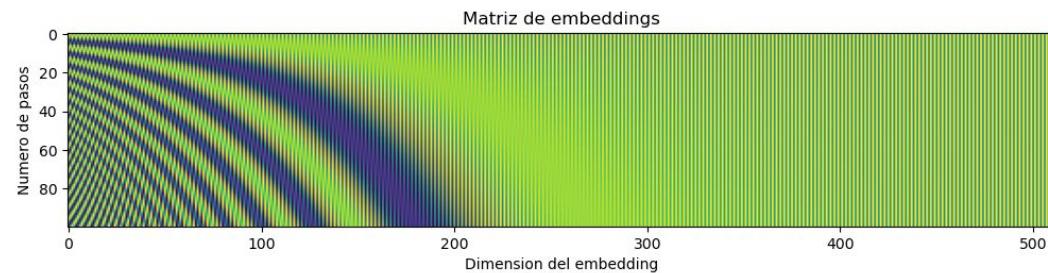
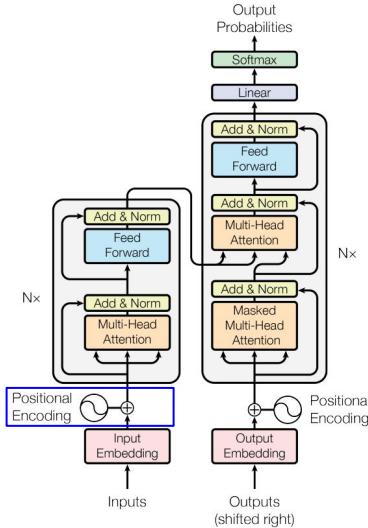


Figure 1: The Transformer - model architecture.



Positional encoding



After a long day at work, she sat down to enjoy a warm cup of -----

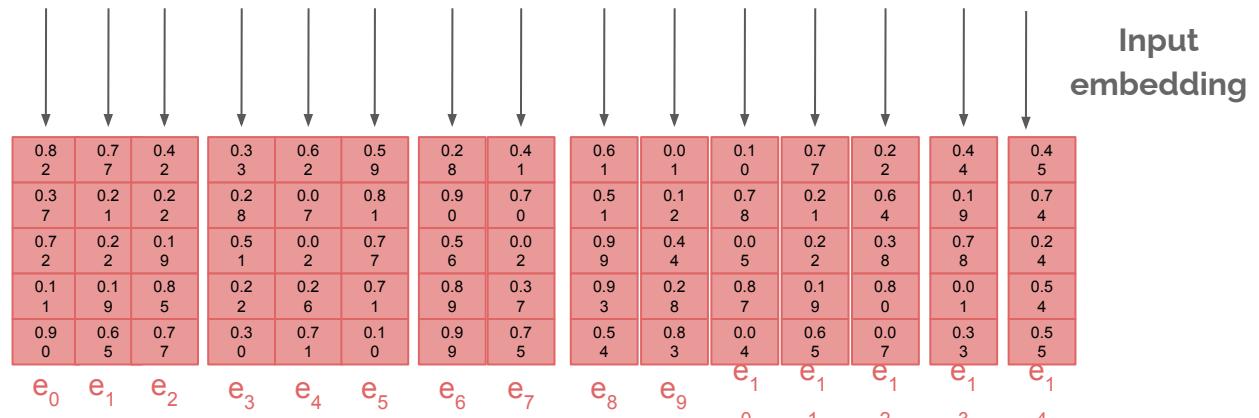
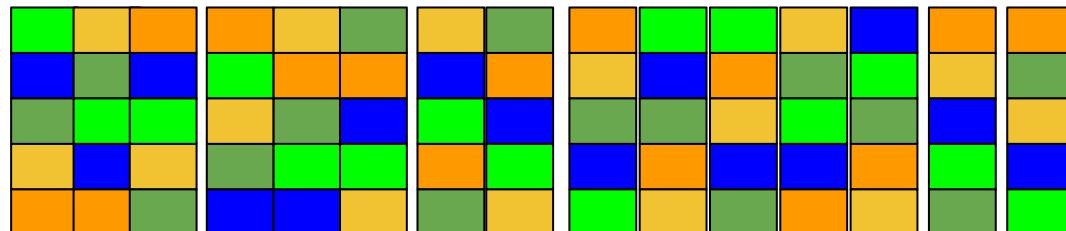
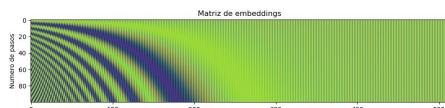


Figure 1: The Transformer - model architecture.



Word embeddings **aware of its position** within the sequence

Attention mechanism

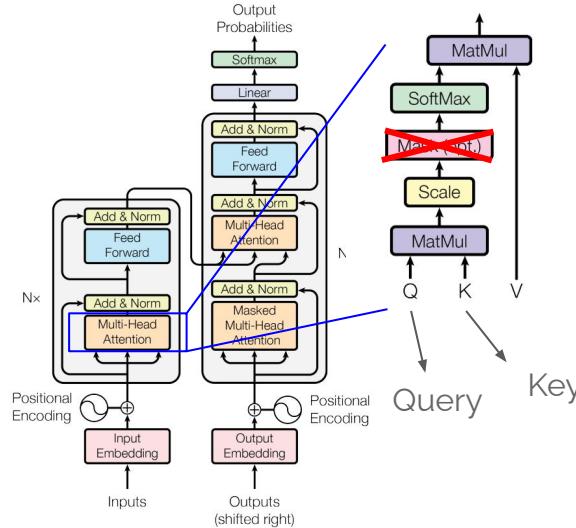


Figure 1: The Transformer - model architecture.



Attention mechanism

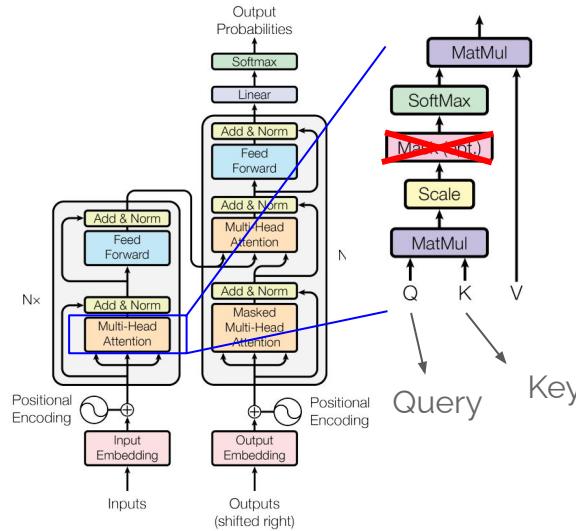
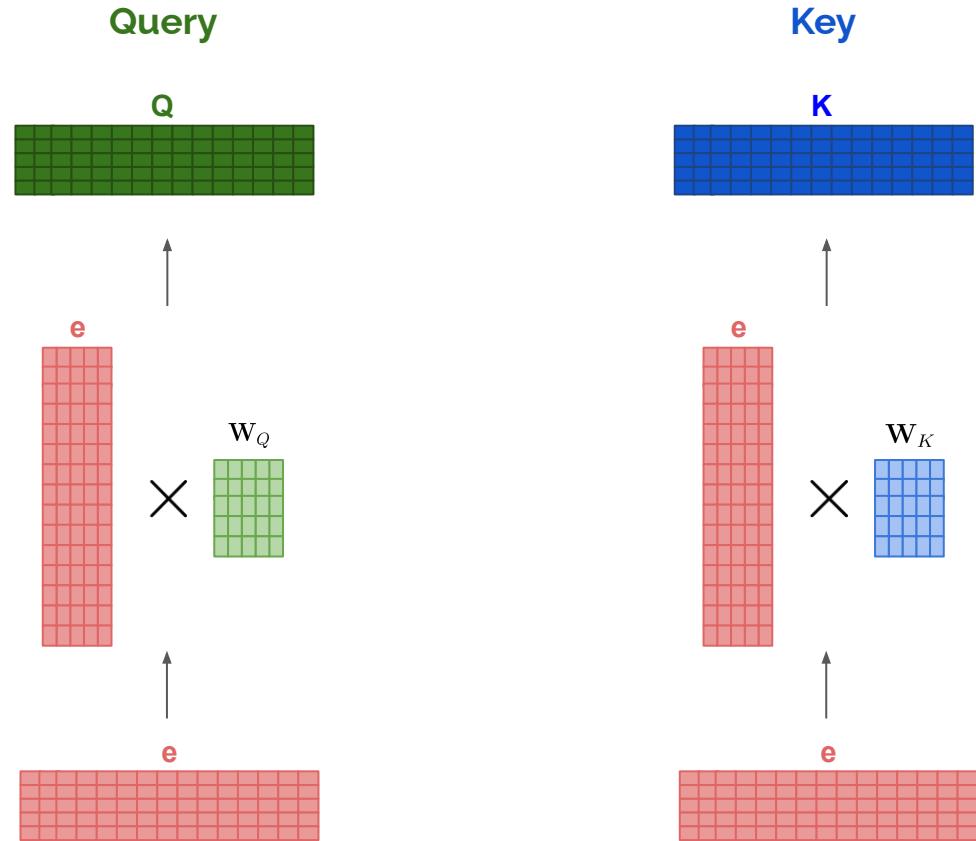
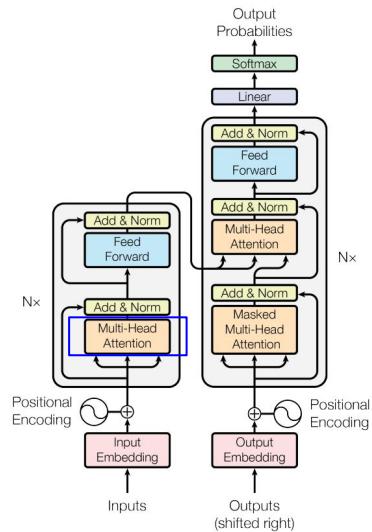


Figure 1: The Transformer - model architecture.



Attention mechanism

What do **queries** and **keys** represent?



After a long day at work, she sat down to enjoy a warm cup of -----

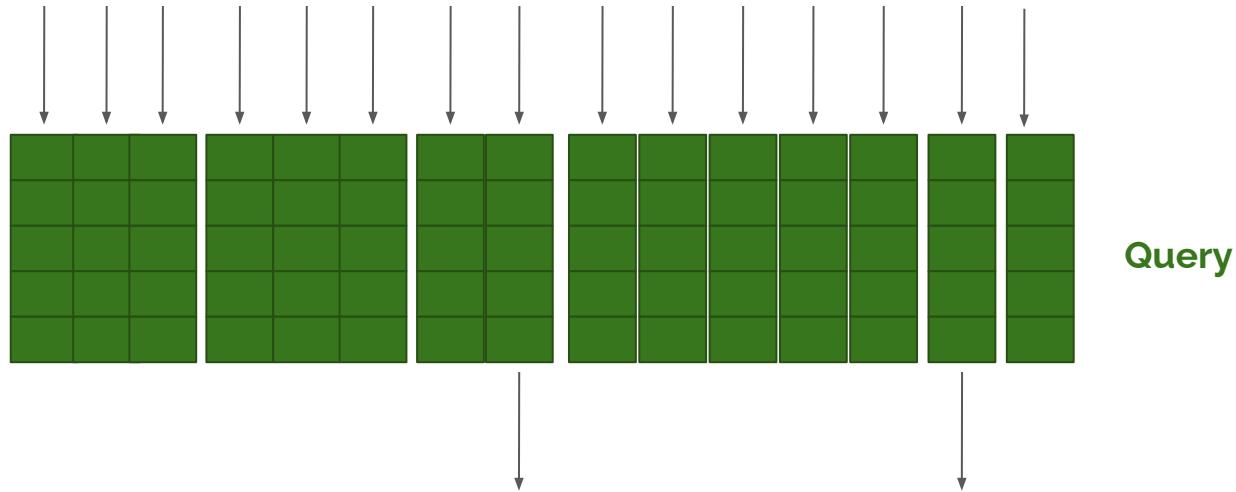


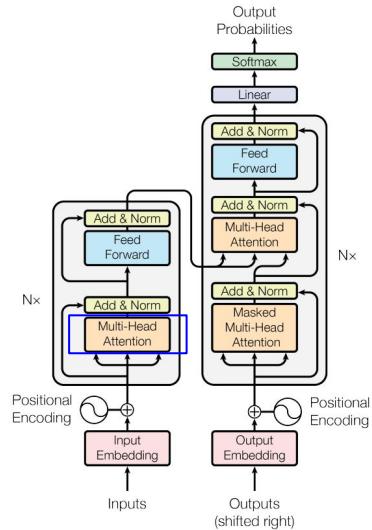
Figure 1: The Transformer - model architecture.

Any information about
who performed the
 action?

Any **adjectives** in front
 of me?

Attention mechanism

What do **queries** and **keys** represent?



After a long day at work, she sat down to enjoy a warm cup of -----

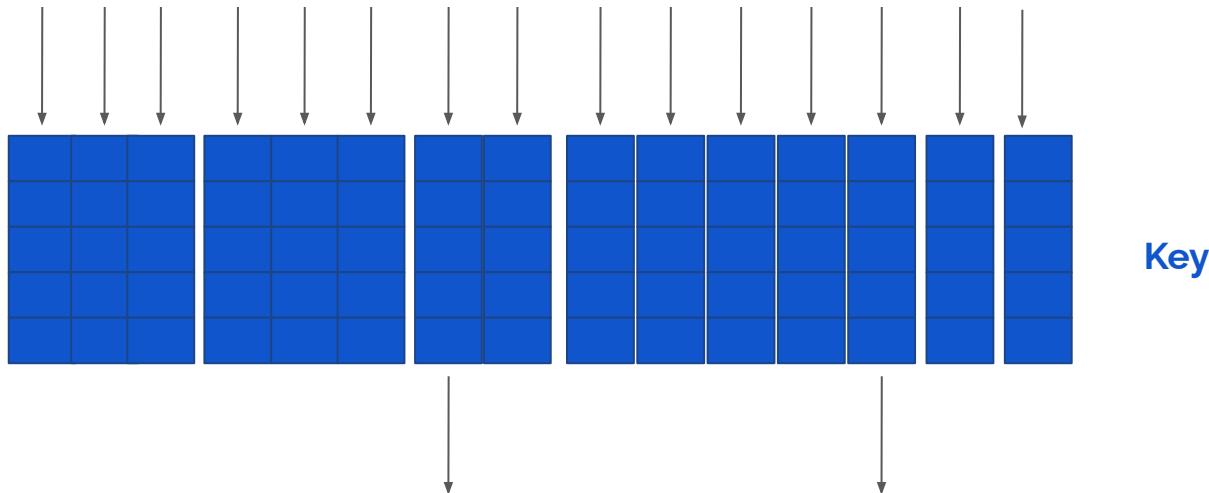
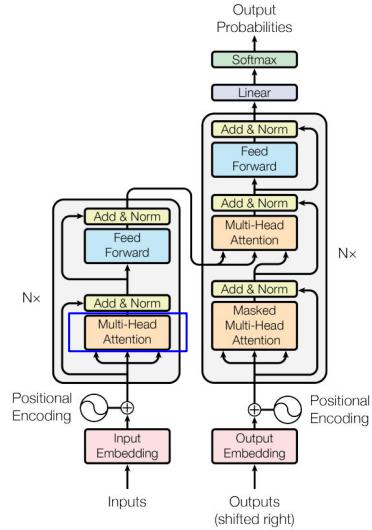


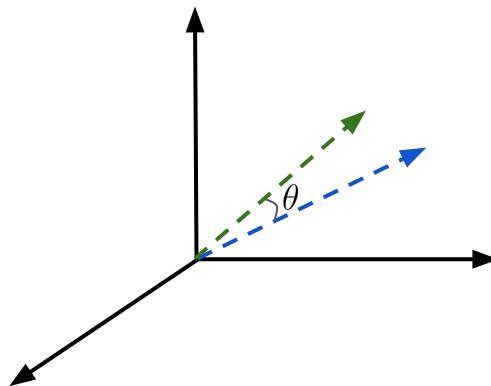
Figure 1: The Transformer - model architecture.

Attention mechanism

How can we compute similarity between **query** and **key** vectors?



High similarity



Low similarity

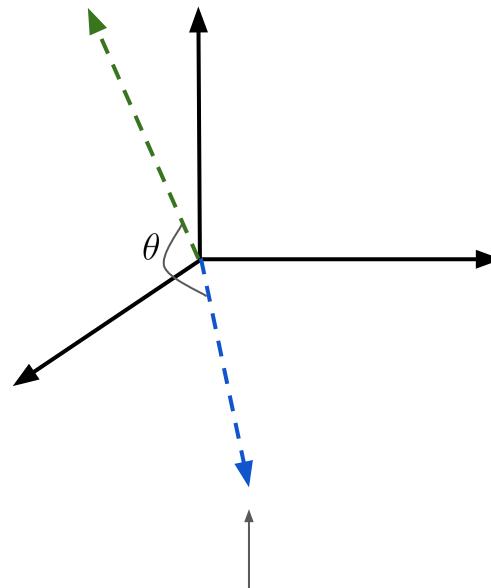
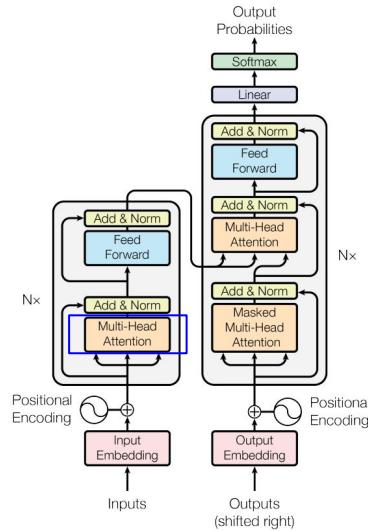


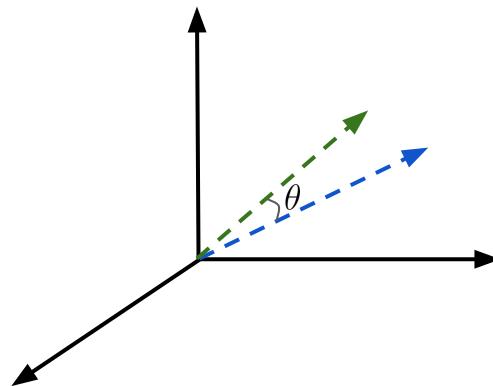
Figure 1: The Transformer - model architecture.

Attention mechanism

How can we compute similarity between **query** and **key** vectors?



High similarity



Low similarity

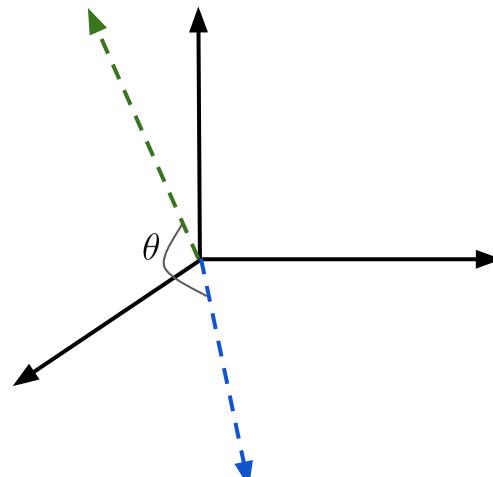
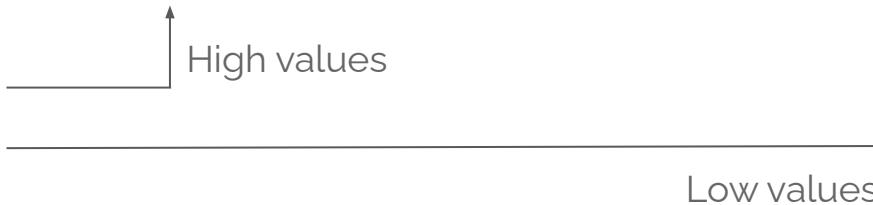


Figure 1: The Transformer - model architecture.

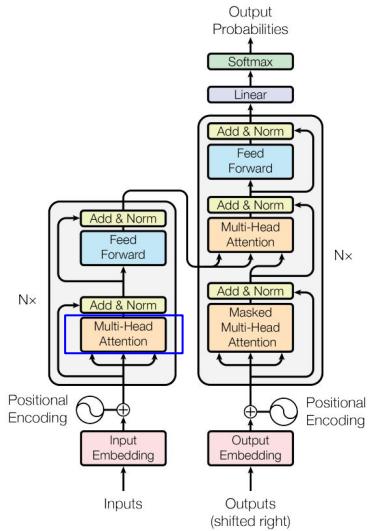
Cosine similarity

$$\text{Cos}(A, B) = \frac{A \cdot B^T}{|A| \cdot |B|}$$



Attention mechanism

How can we compute similarity between **query** and **key** vectors?



$$\text{Cos}(A, B) = \frac{A \cdot B^T}{|A| \cdot |B|} \longrightarrow \text{Similarity}(Q, K) = \frac{QK^T}{\sqrt{d_K}}$$

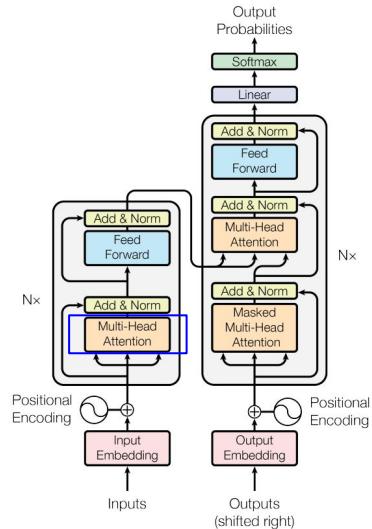


Dimension of the key

Figure 1: The Transformer - model architecture.

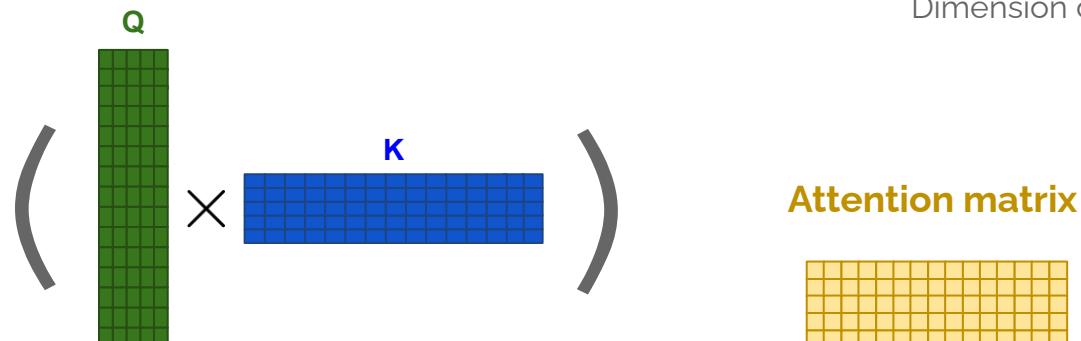
Attention mechanism

How can we compute similarity between **query** and **key** vectors?

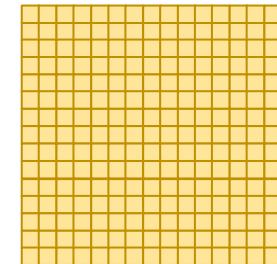


$$\text{Cos}(A, B) = \frac{A \cdot B^T}{|A| \cdot |B|} \longrightarrow \text{Similarity}(Q, K) = \frac{QK^T}{\sqrt{d_K}}$$

Dimension of the key



$$\text{Similarity}(Q, K) = \frac{\dots}{\sqrt{d_K}} =$$



Attention mechanism

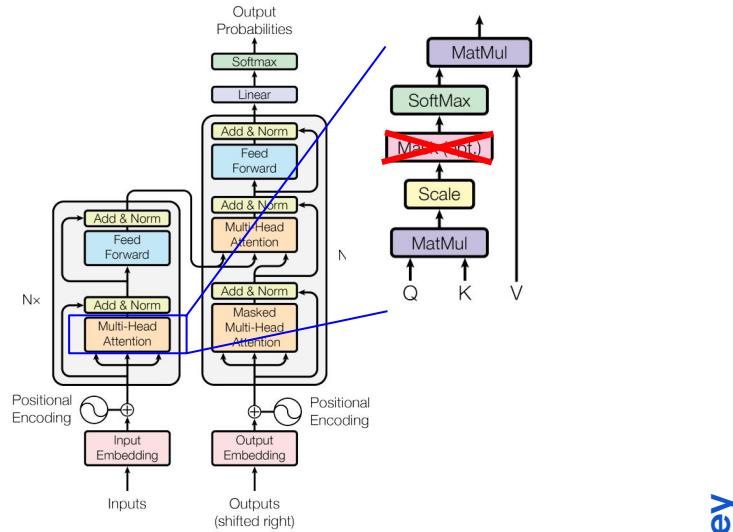
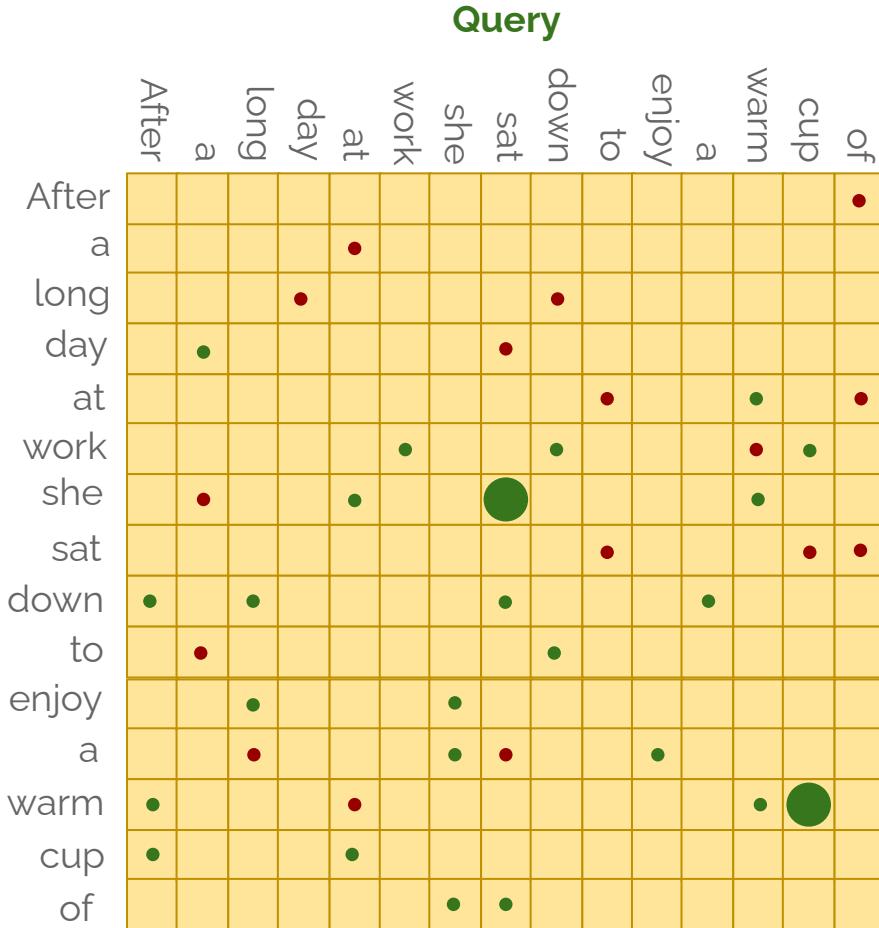


Figure 1: The Transformer - model architecture.

$$\frac{QK^T}{\sqrt{d_K}}$$



Attention mechanism

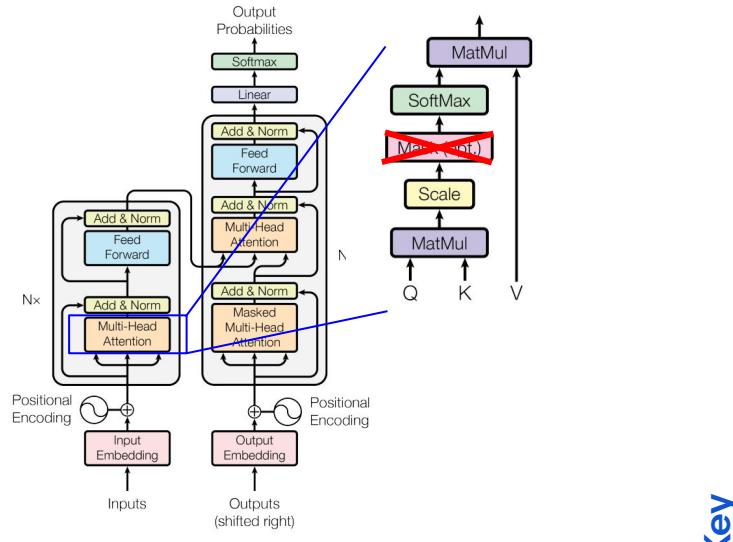
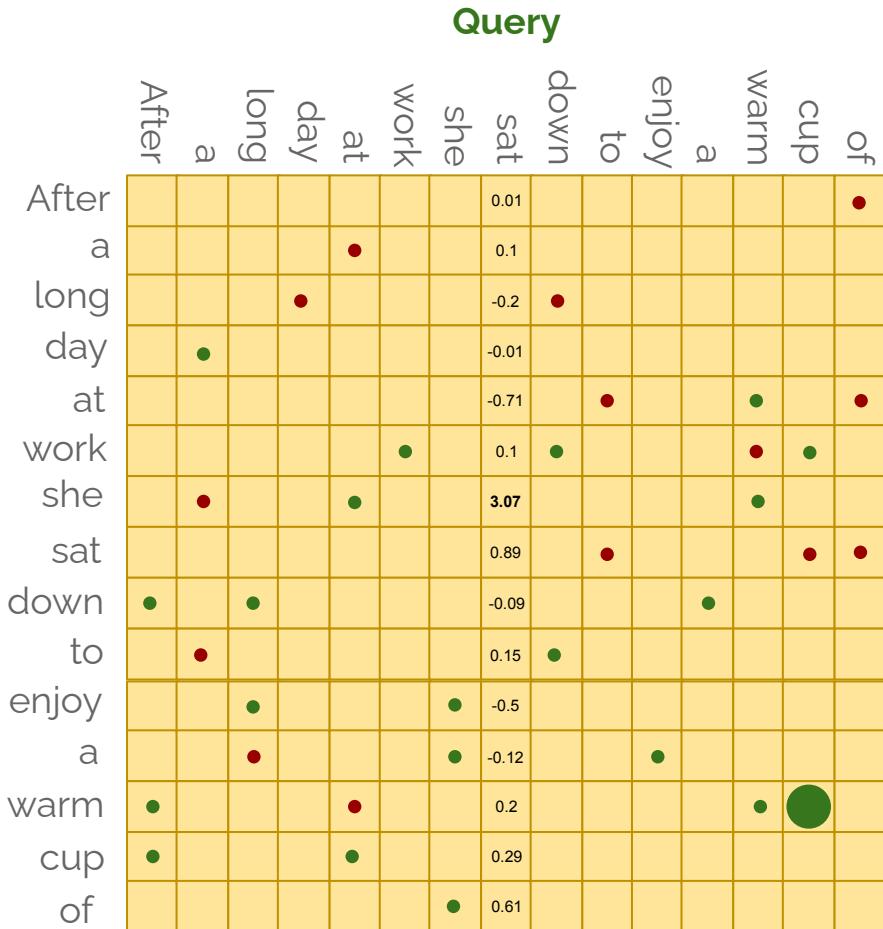


Figure 1: The Transformer - model architecture.

$$\frac{QK^T}{\sqrt{d_K}}$$



Attention mechanism

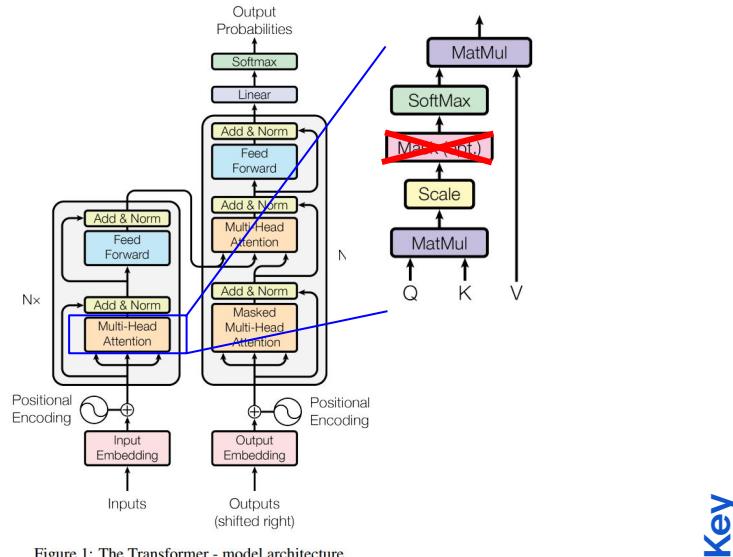
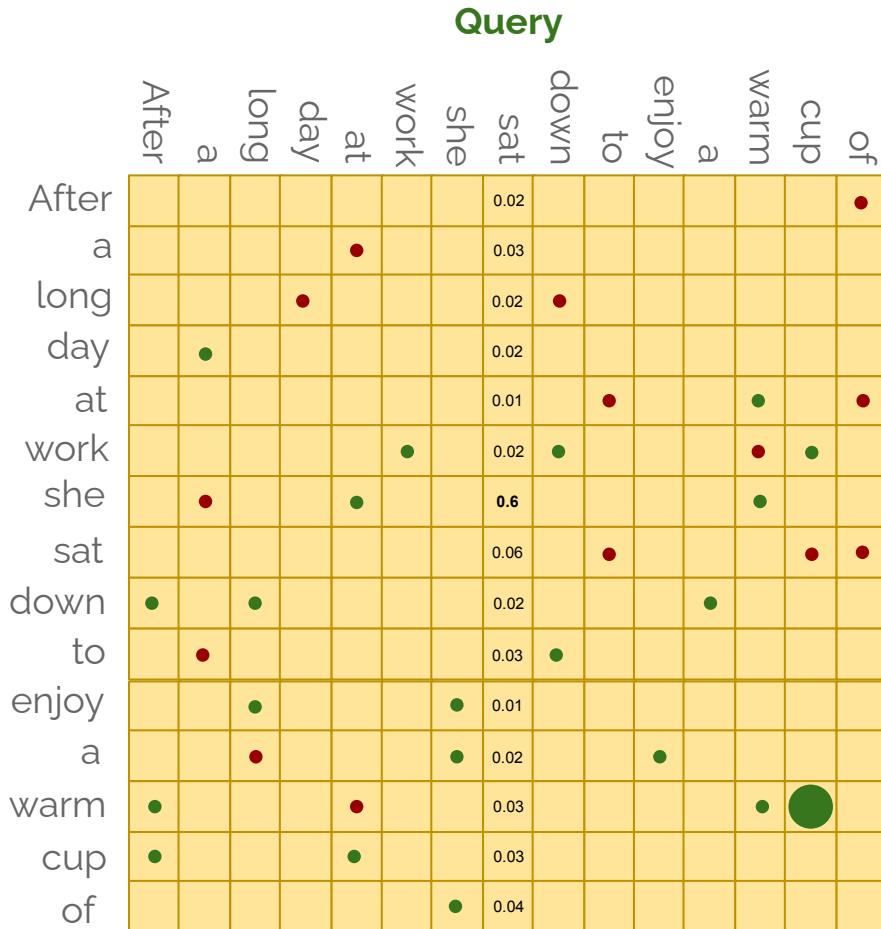


Figure 1: The Transformer - model architecture.

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)$$

Key



Attention mechanism

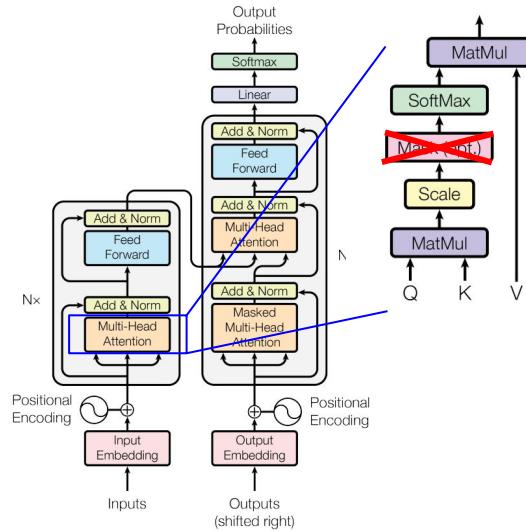
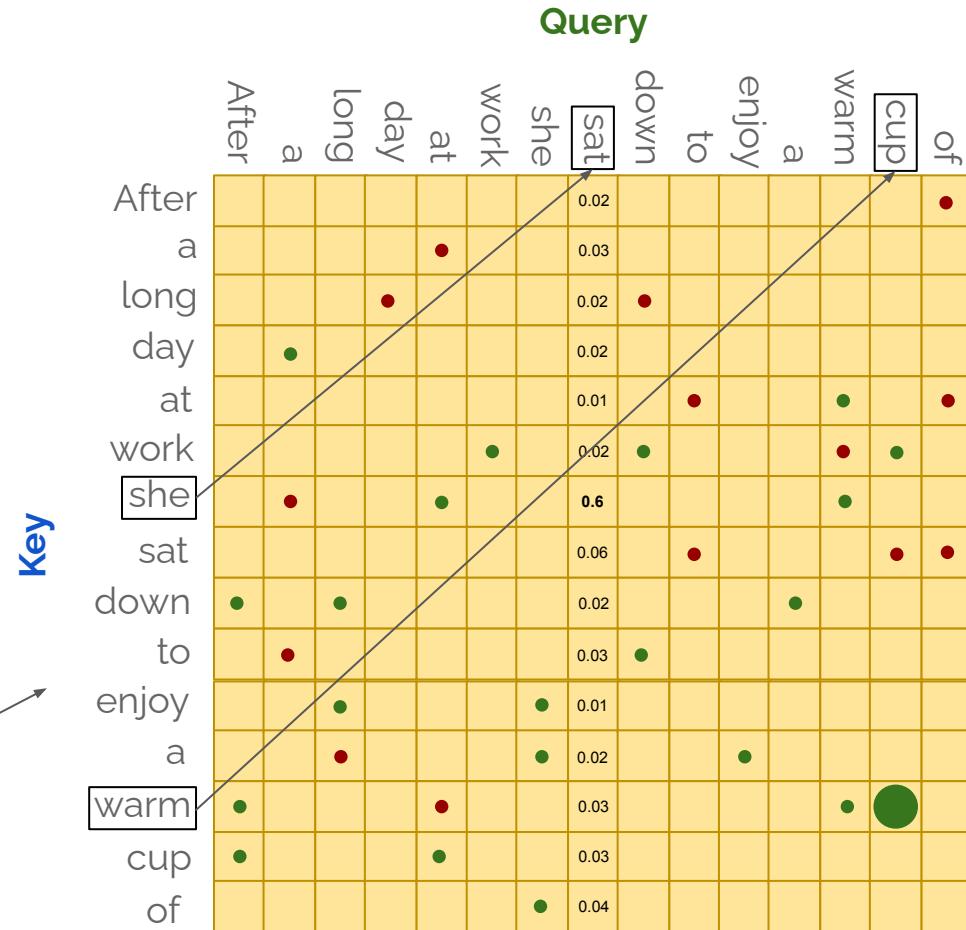


Figure 1: The Transformer - model architecture.

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)$$



Attention mechanism

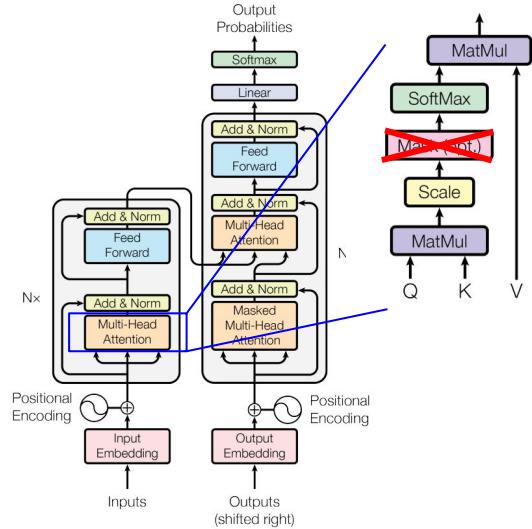


Figure 1: The Transformer - model architecture.

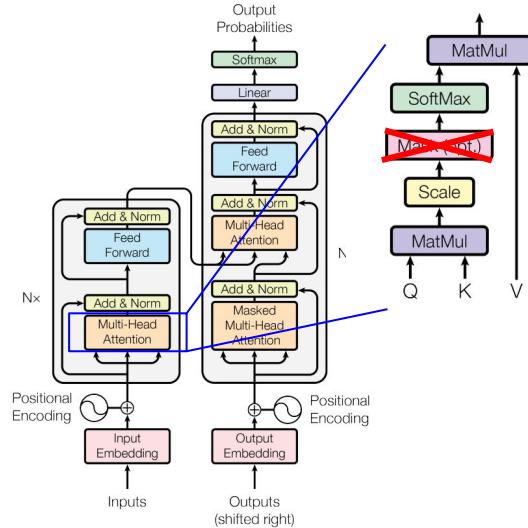
If we look more closely, we can see that we have almost **derived the attention formula**

Attention matrix \longrightarrow
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)$$

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

But, what about the **V**?

Attention mechanism

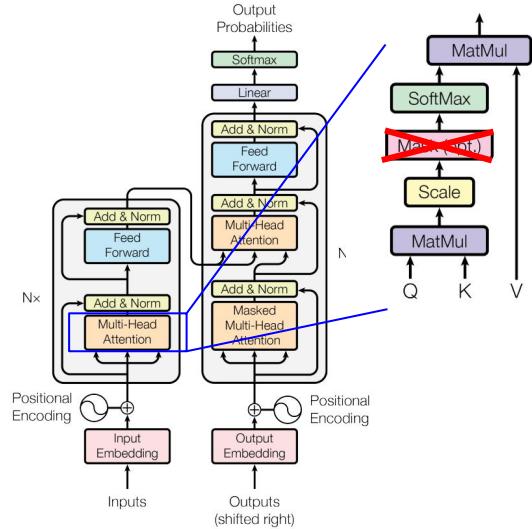


Similarly to how we generate the **query** and **key** vectors, we also generate the **value** vectors

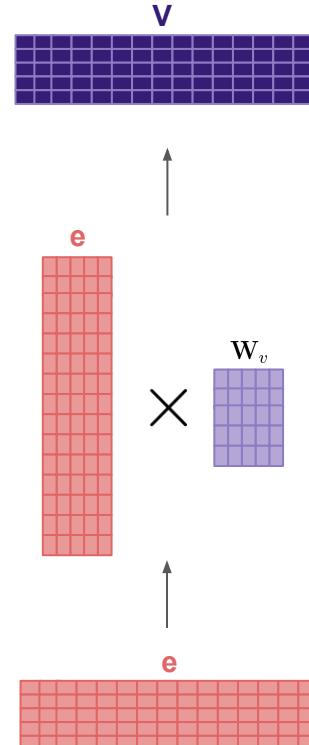


Figure 1: The Transformer - model architecture.

Attention mechanism



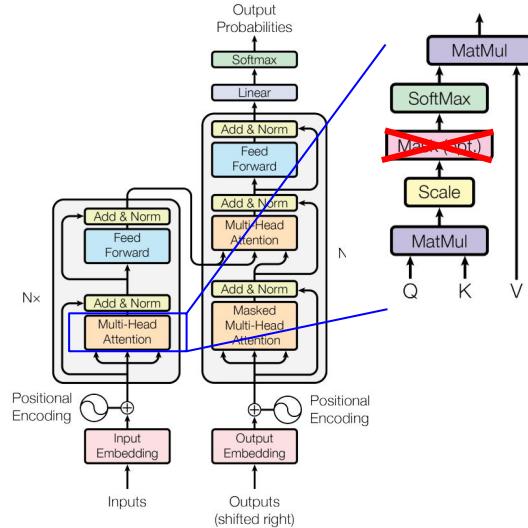
Similarly to how we generate the **query** and **key** vectors, we also generate the **value** vectors



A more refined and **complex** version of e

Figure 1: The Transformer - model architecture.

Attention mechanism



The **value** represents the step that the embeddings need to take to update their meaning based on the context

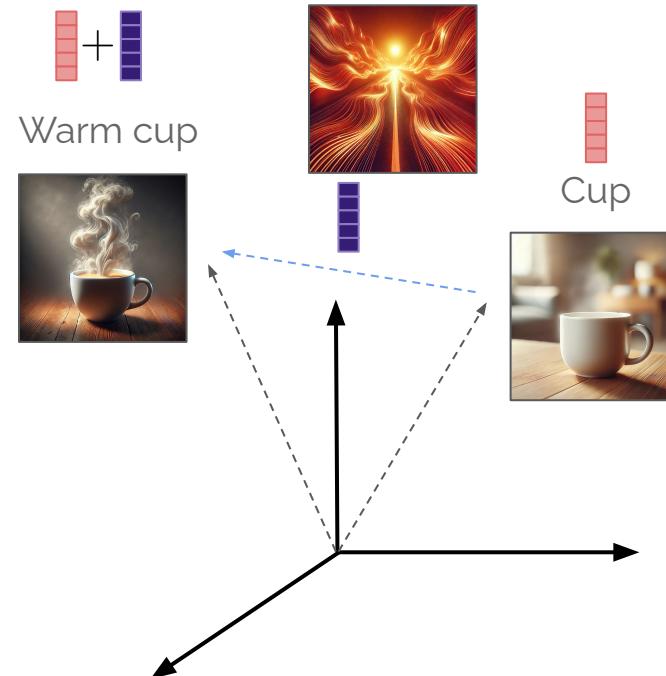


Figure 1: The Transformer - model architecture.

After a long day at work, she sat down to enjoy a warm cup of -----

Attention mechanism

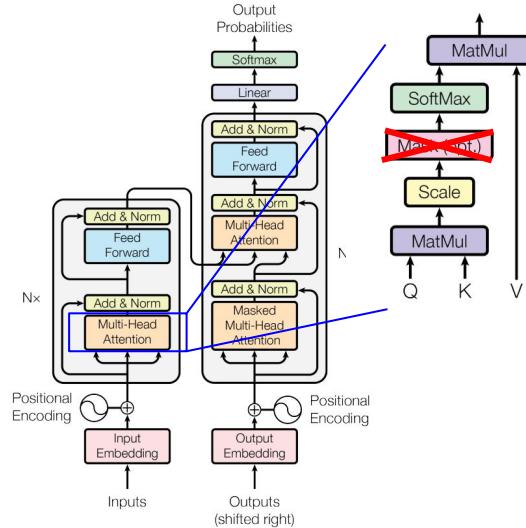
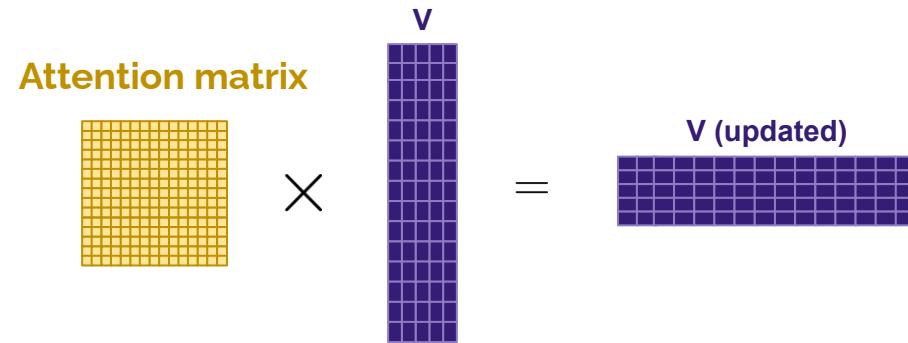


Figure 1: The Transformer - model architecture.

Before adding the **values** to the embedding we **multiply** them by the **attention matrix**



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

Attention mechanism

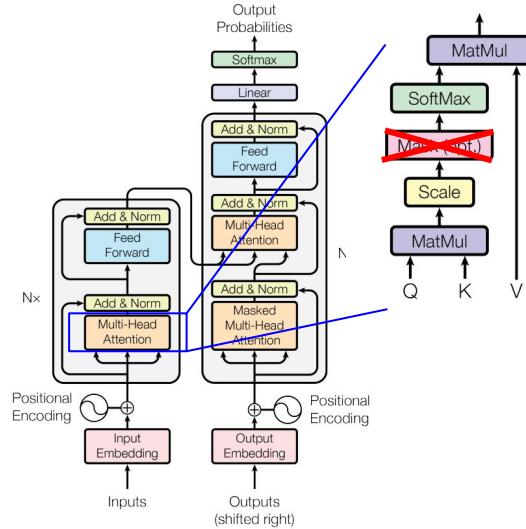
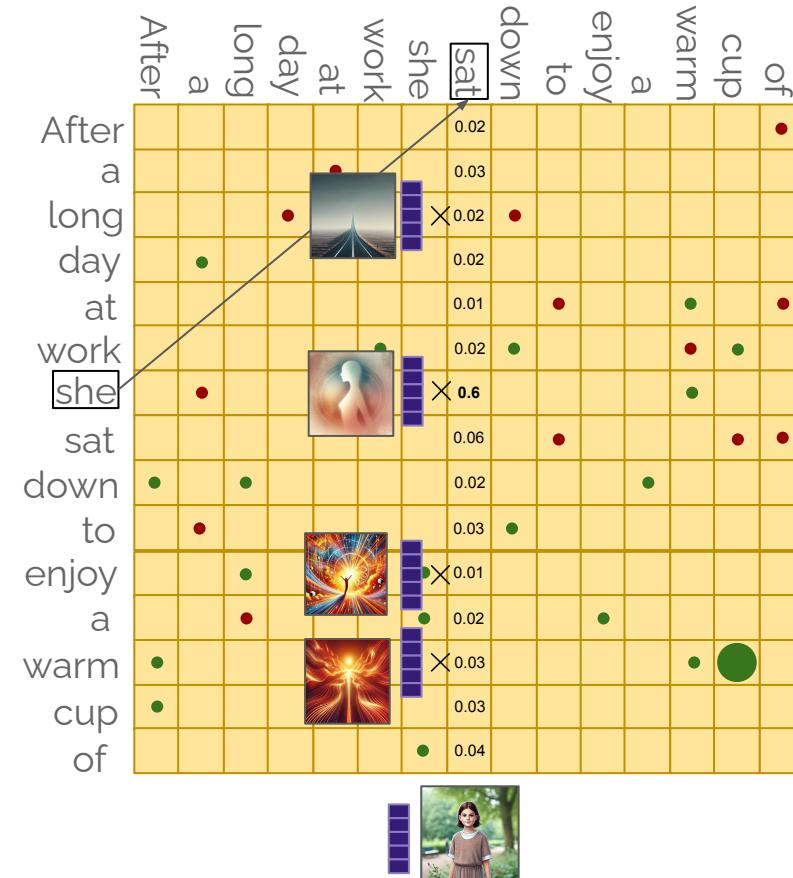


Figure 1: The Transformer - model architecture.



Attention mechanism

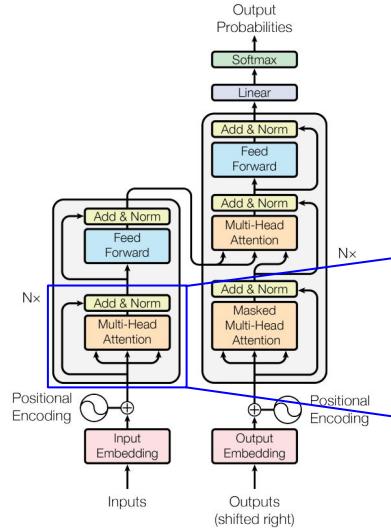
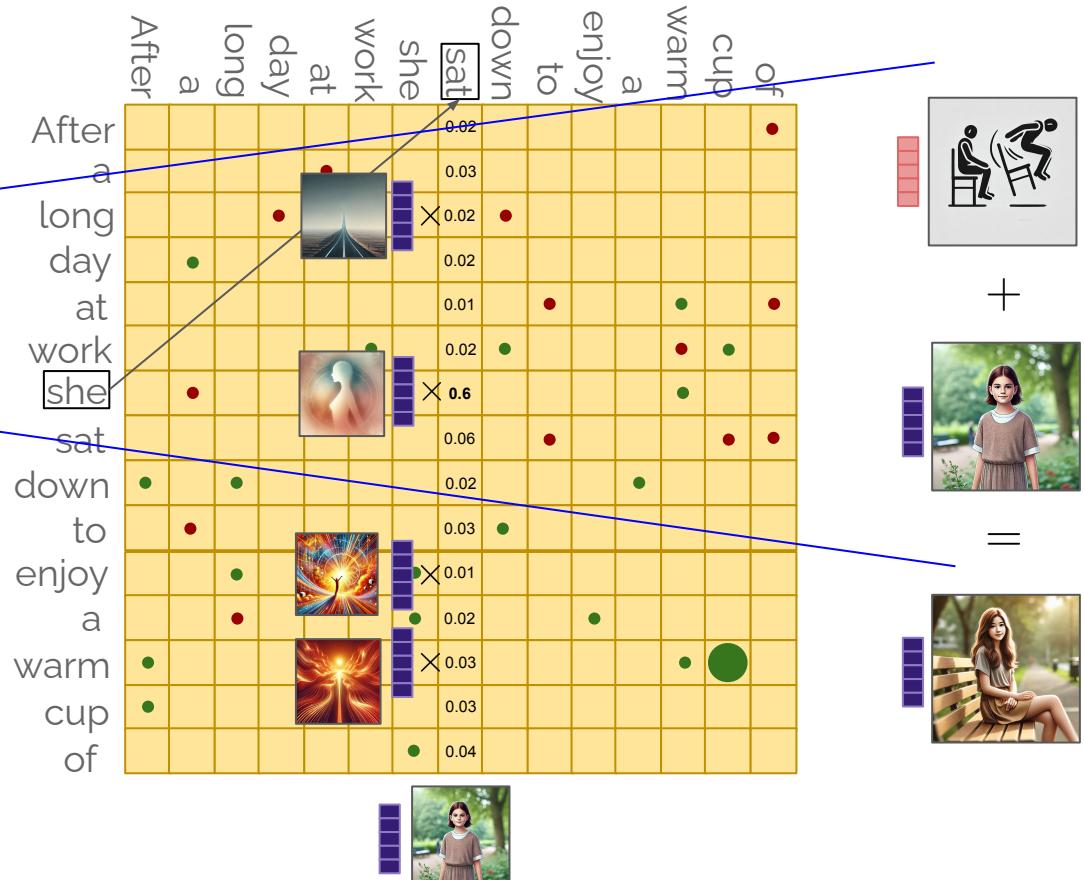
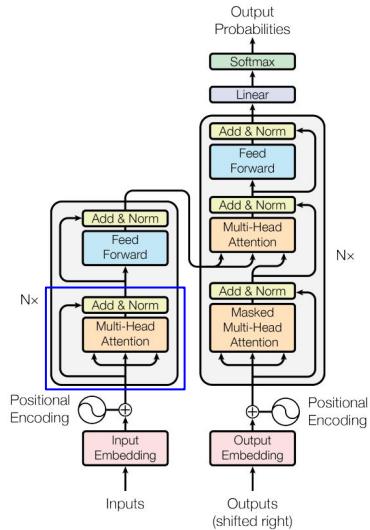


Figure 1: The Transformer - model architecture.



Attention mechanism



With the attention mechanism we **transform** the meaning of the **embedding** based on the **context**

*After a long day at work, she sat down to enjoy a **warm cup** of -----*

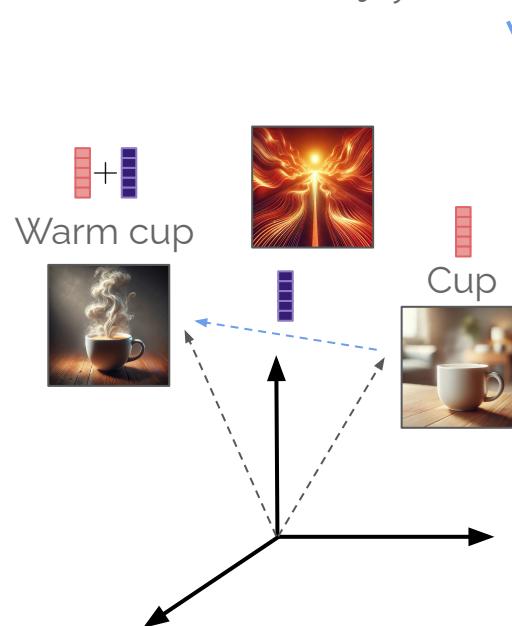
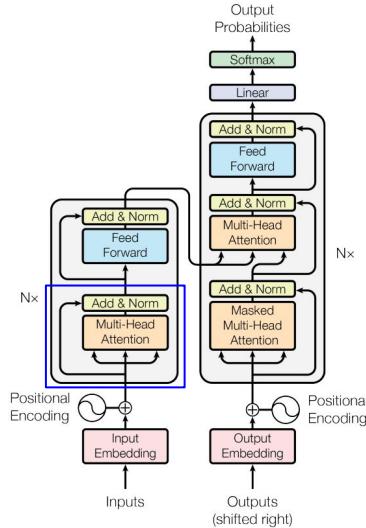


Figure 1: The Transformer - model architecture.

Attention mechanism



We can apply the attention mechanism to the **output of previous layers**, thereby building **more complex representations** of relationships between concepts

*After a long day at work, **she** sat down to enjoy a **warm cup** of -----*

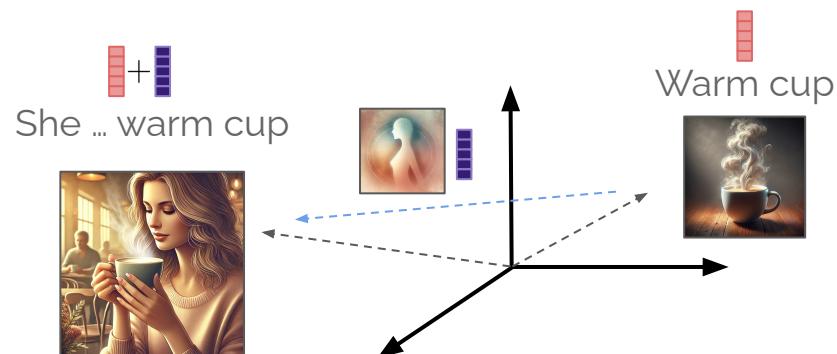
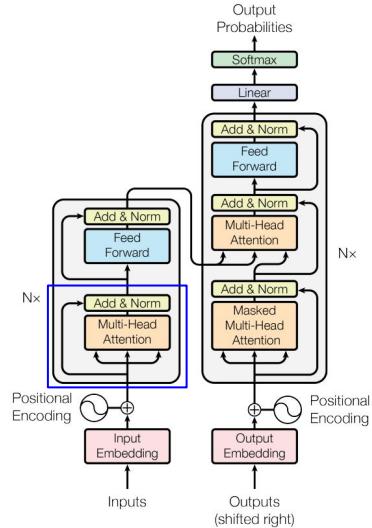


Figure 1: The Transformer - model architecture.

Multi-Head Attention



In transformers, multiple **attention mechanisms** are computed in **parallel**, a process known as **multi-head attention**

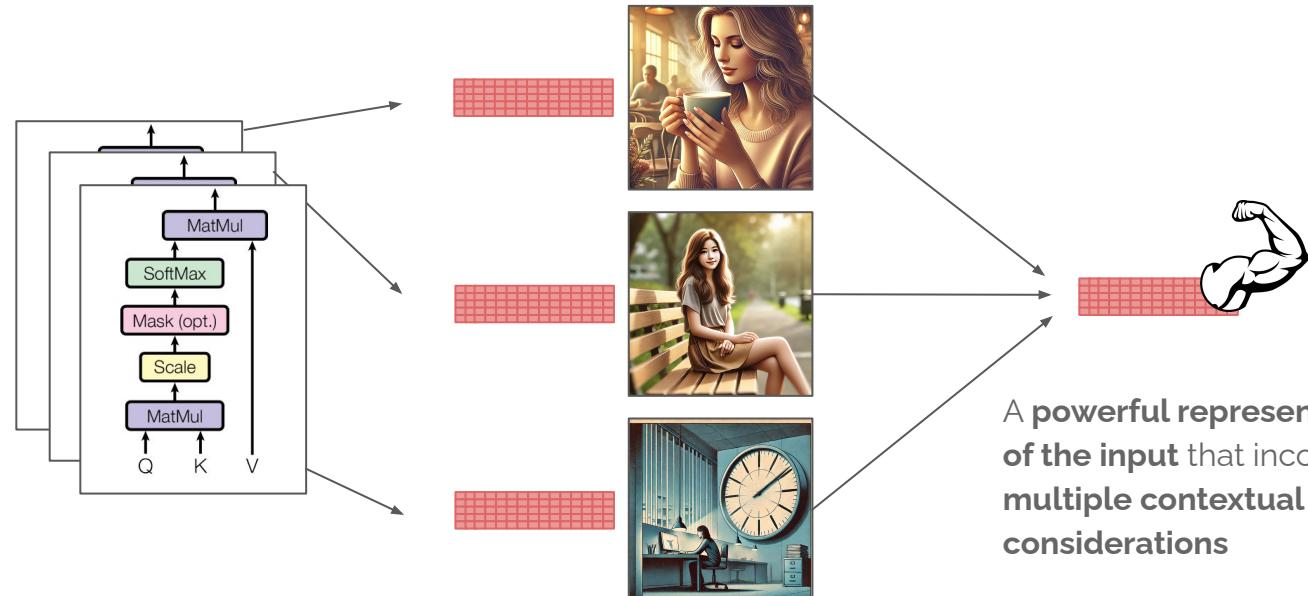
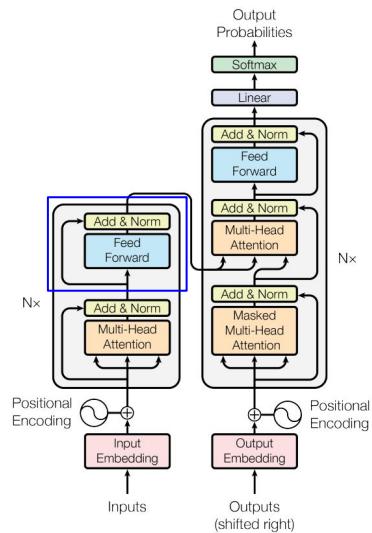


Figure 1: The Transformer - model architecture.

Encoder



Finally, this representation is passed through a **fully-connected layer** through a **residual connection**

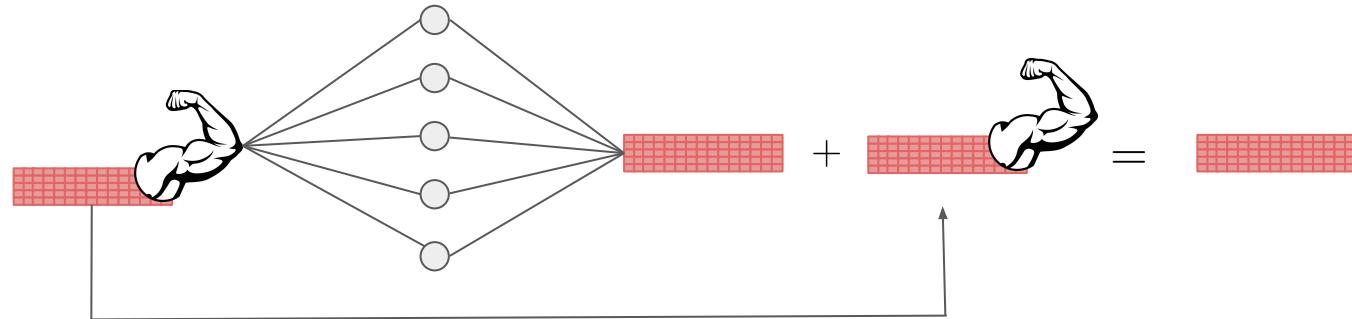


Figure 1: The Transformer - model architecture.

Encoder

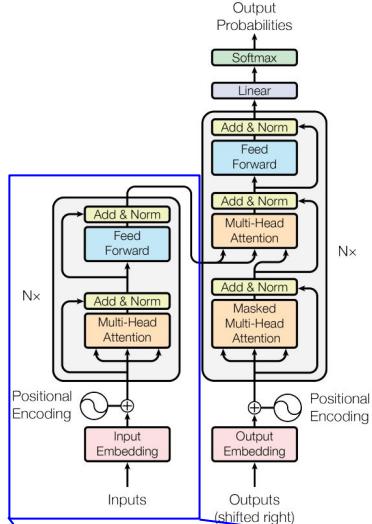
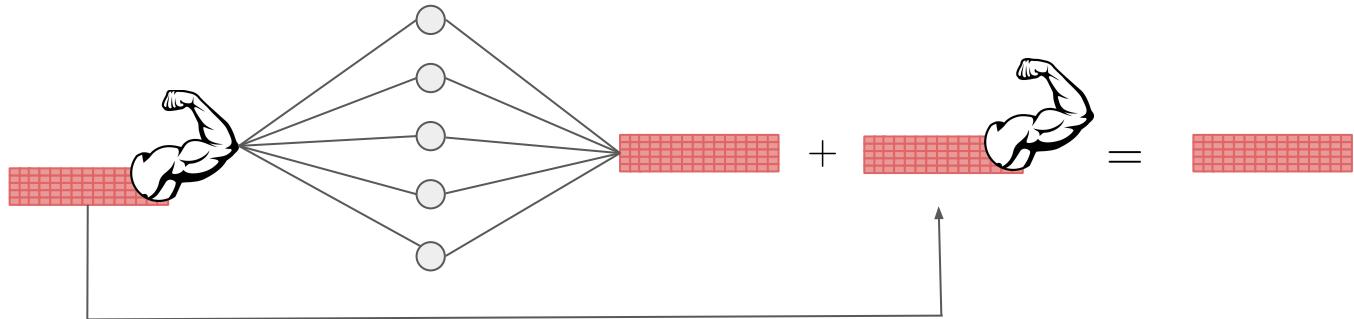


Figure 1: The Transformer - model architecture.

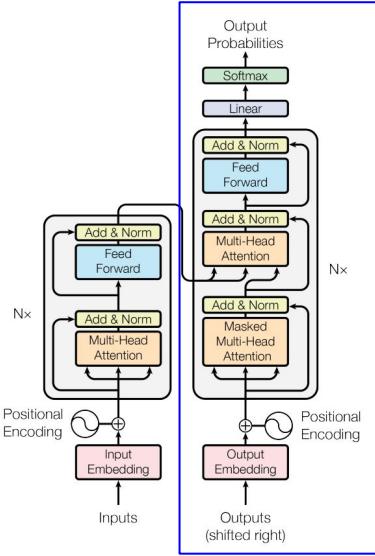
Finally, this representation is passed through a **fully-connected layer** through a **residual connection**



Encoder

This part of the transformer **encodes** the input while considering the context in multiple ways

Decoder



The other part of the transformer is generally referred to as the **decoder**. Its purpose is to **generate the next word in the sequence**

After a long day at work, she sat down to enjoy a warm cup of -----

To achieve this, it relies on the **embeddings processed and updated by the encoder**

Figure 1: The Transformer - model architecture.

Decoder

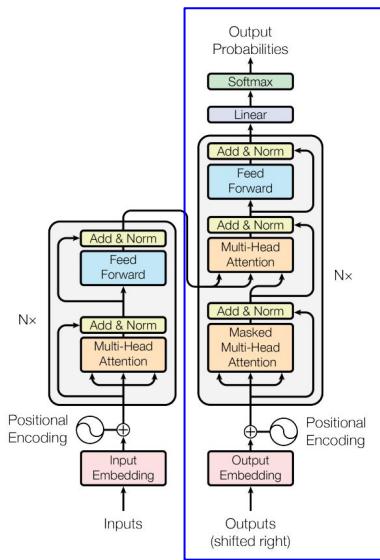


Figure 1: The Transformer - model architecture.

The other part of the transformer is generally referred to as the **decoder**. Its purpose is to **generate the next word in the sequence**

After a long day at work, she sat down to enjoy a warm cup of -----

To achieve this, it relies on the **embeddings processed and updated by the encoder**

Most of the components of the decoder have already been explored in previous discussions

Output embedding

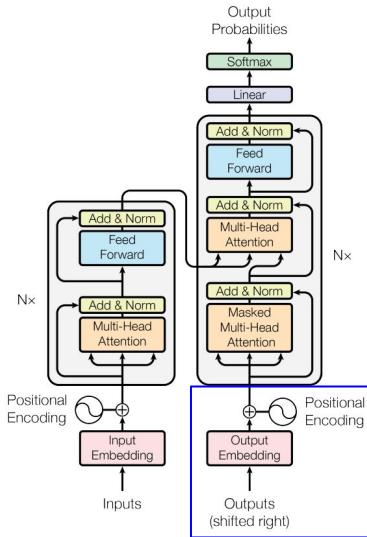
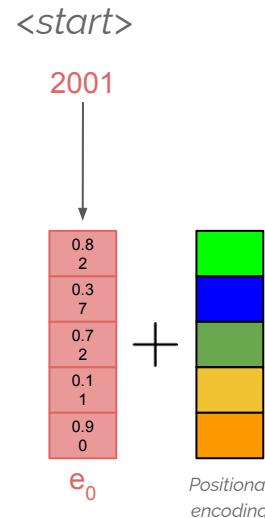


Figure 1: The Transformer - model architecture.

When generating text, we begin by passing the **<start> token** to the decoder, indicating the start of the generation



Multi-Head Attention

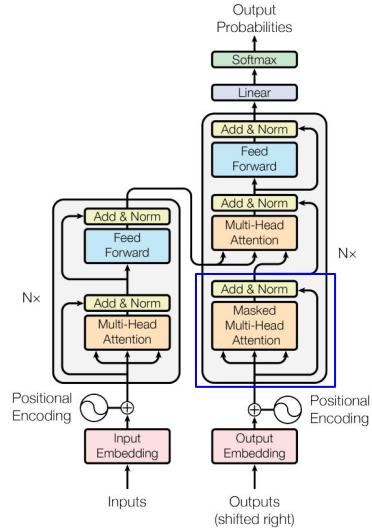
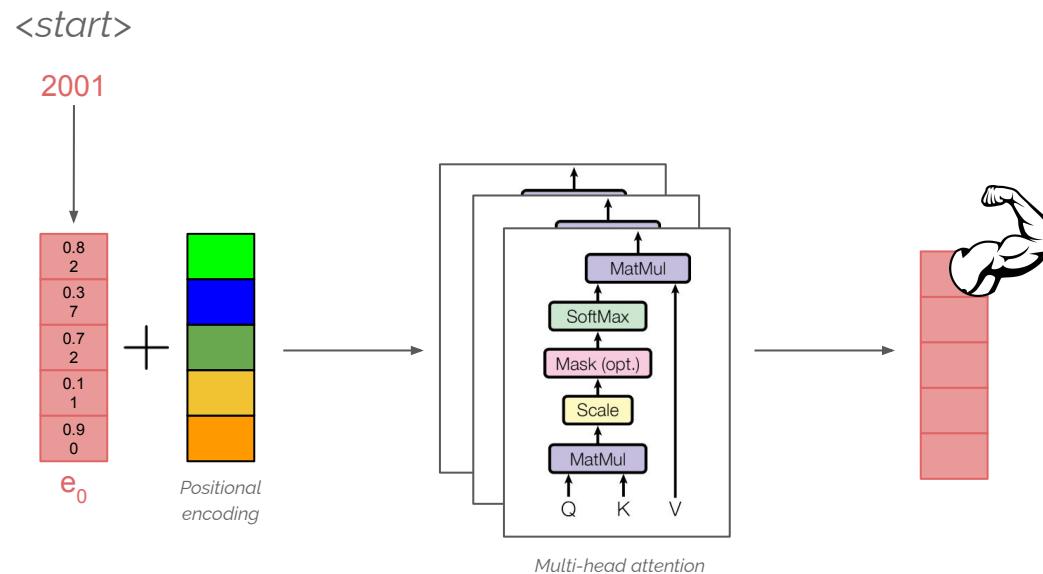
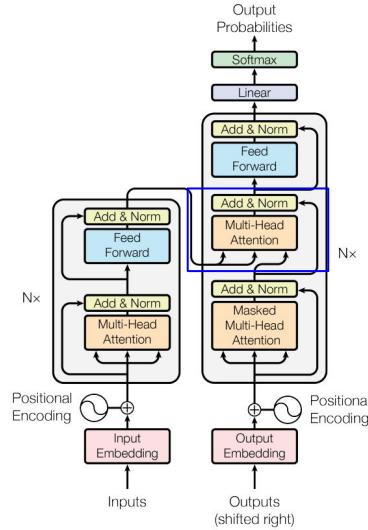


Figure 1: The Transformer - model architecture

Then, as with the encoder, we apply **multi-head attention** to generate more **complex embeddings** by learning from the context



Integration of the Encoder



Then, we combine information from the encoder and the decoder to take both into account for generating the next token

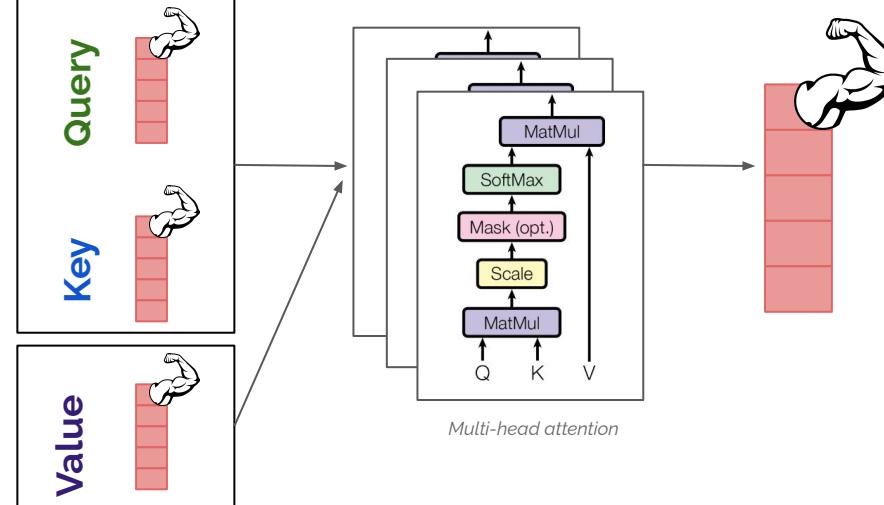
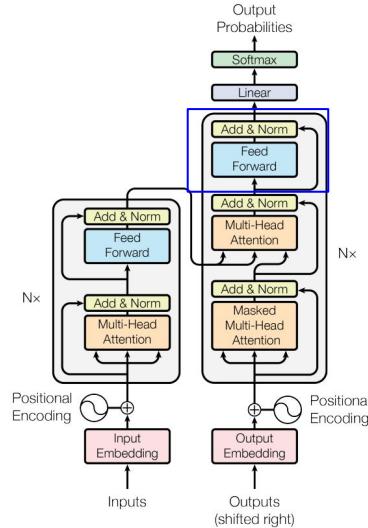


Figure 1: The Transformer - model architecture.

Decoder



Then, we combine information from the encoder and the decoder to take both into account for generating the next token

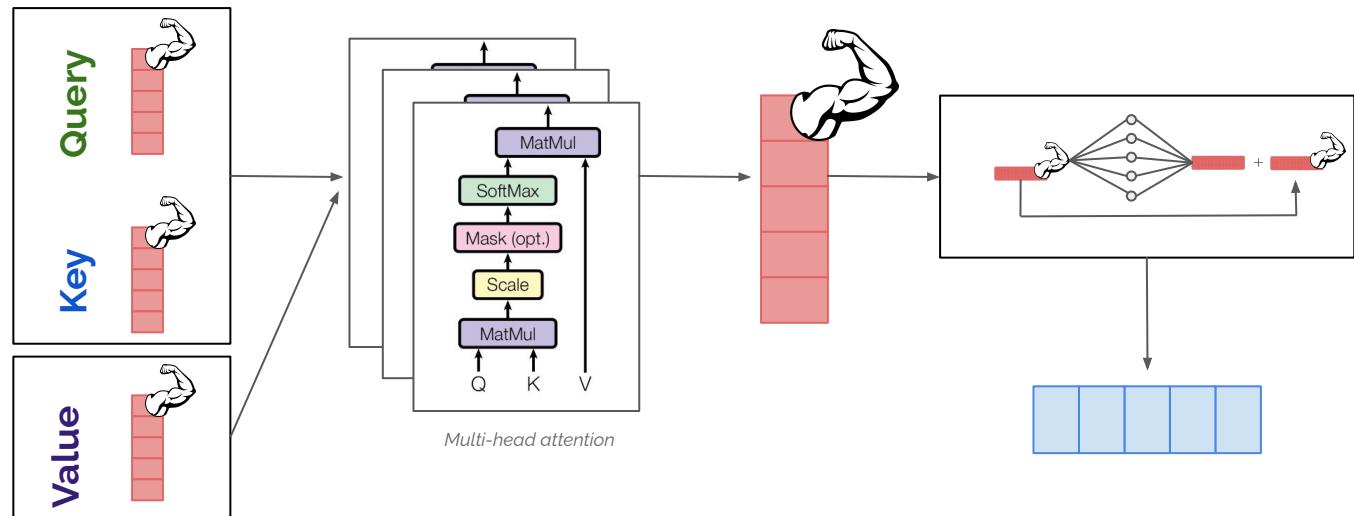
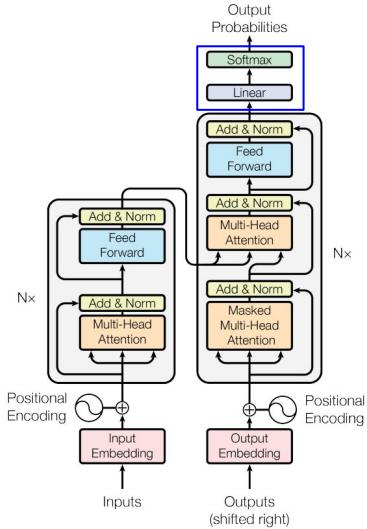


Figure 1: The Transformer - model architecture.

Final layer



Finally, we apply a **fully-connected layer** to match the shape of the **vocabulary** and use a **softmax** function to generate a **probability distribution** over it

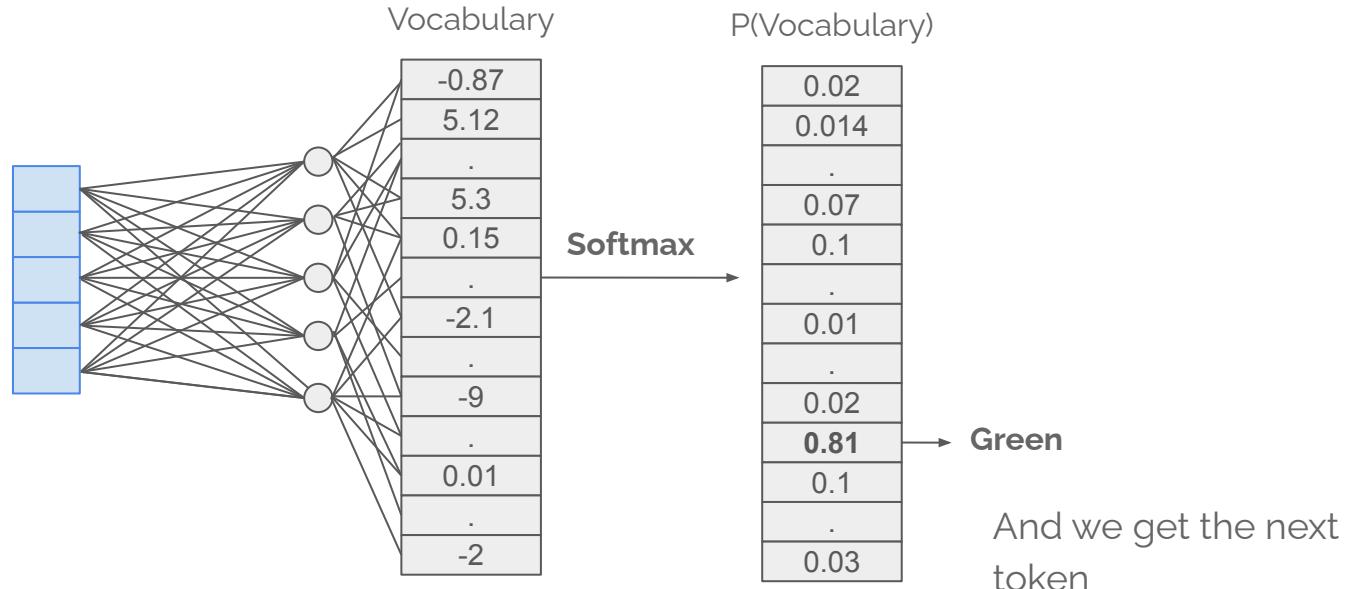


Figure 1: The Transformer - model architecture.

Token generator

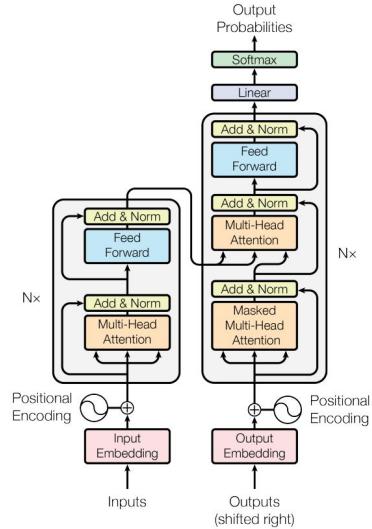
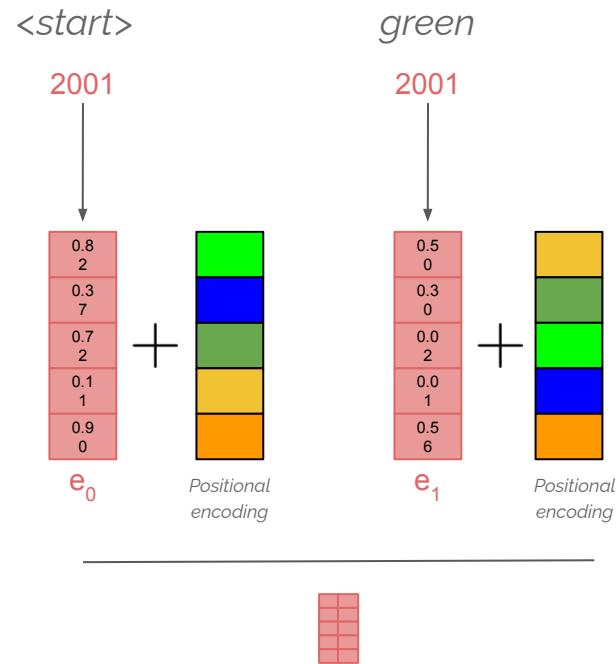
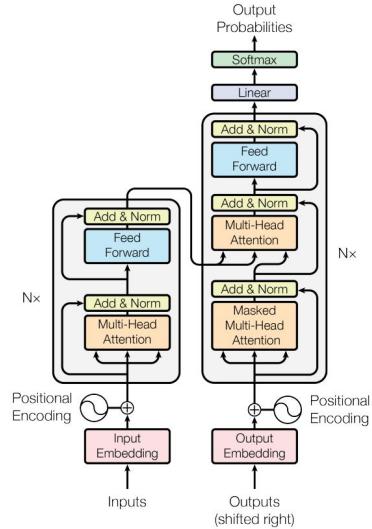


Figure 1: The Transformer - model architecture.

To continue generating tokens, we **add the most recently generated token** and **repeat the process**



Token generator



To continue generating tokens, we **add the most recently generated token** and **repeat the process**

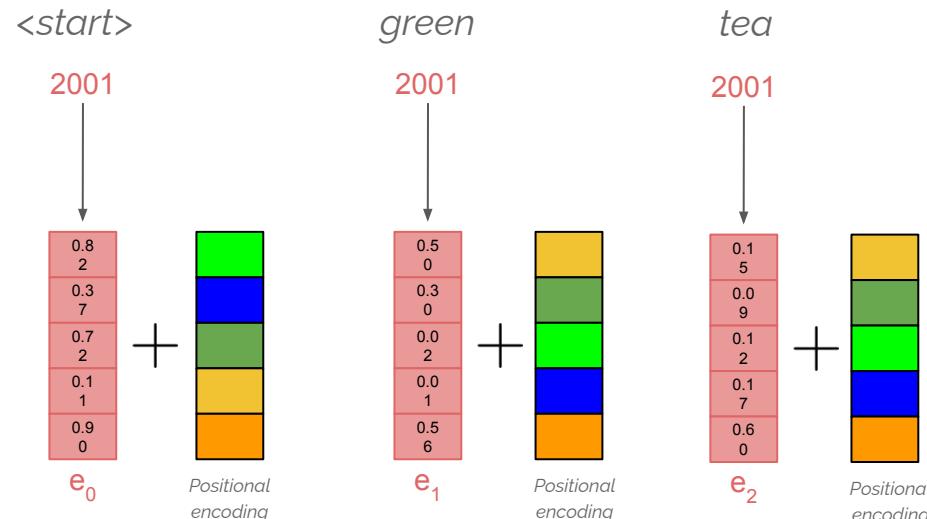
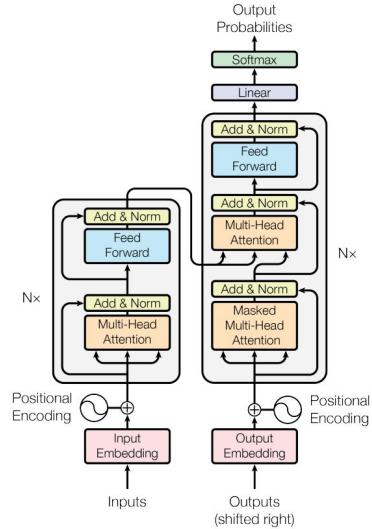


Figure 1: The Transformer - model architecture.



Token generator



The generation process continues until the transformer produces the **<stop>** token

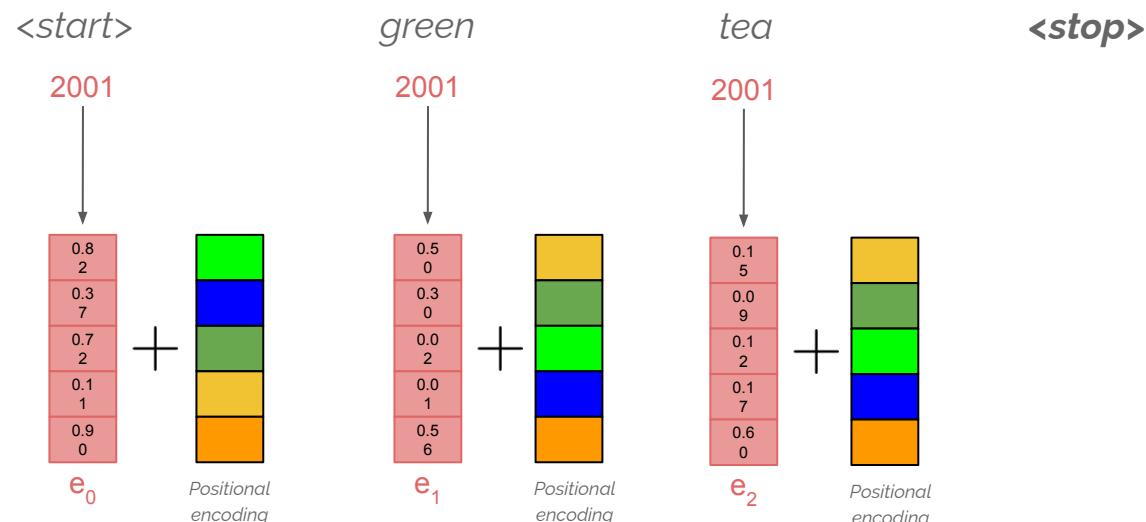
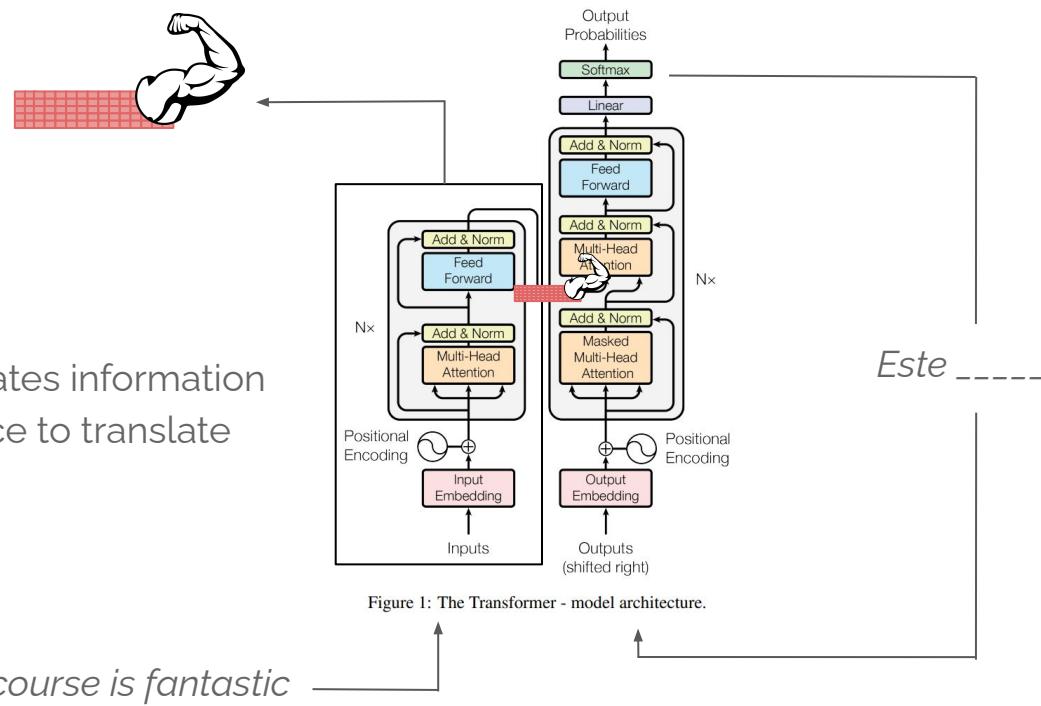


Figure 1: The Transformer - model architecture.



Encoder and decoder

Encoder and decoder can be used, for instance, for **translation**



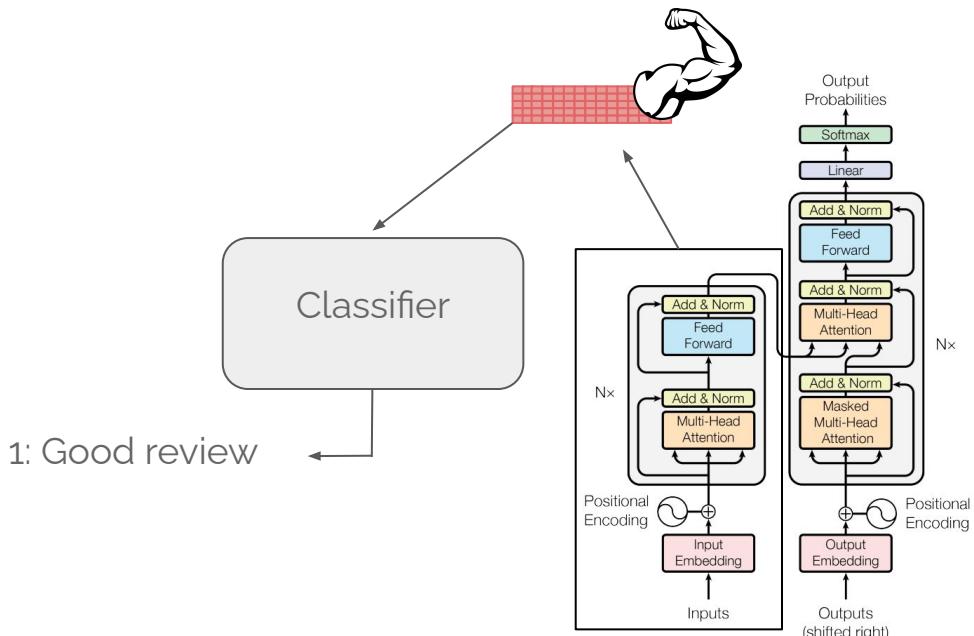
The encoder generates information about the sentence to translate

Figure 1: The Transformer - model architecture.

This course is fantastic

Encoder and decoder

But, depending on the task **you might not need an encoder and a decoder**



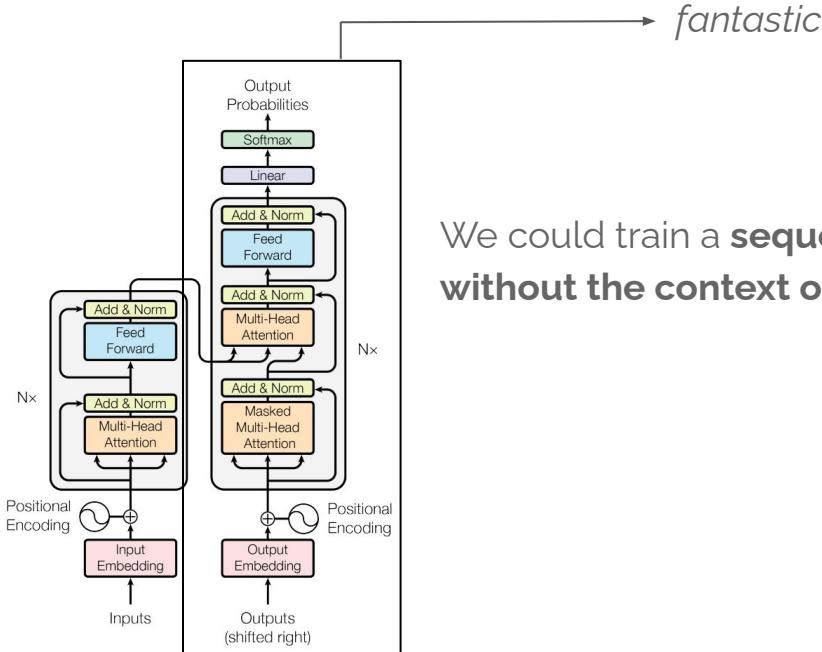
We could train a **classifier** over the embeddings of the **encoding**



Figure 1: The Transformer - model architecture.

Encoder and decoder

But, depending on the task **you might not need an encoder and a decoder**



We could train a **sequence generator**
without the context of the encoder

Figure 1: The Transformer - model architecture.

This course is _____

Additional observations

There are some **observations** worth mentioning regarding transformers

- Masking
- Scalability
- Temperature

Masking

When training a transformer, we can generate **multiple training samples from a single text input**

After a long day at work, she sat down to enjoy a warm cup of coffee.

After _____

After a _____

After a long _____

After a long day _____

After a long day at _____

After a long day at work _____

After a long day at work, she _____

After a long day at work, she sat _____

After a long day at work, she sat down _____

After a long day at work, she sat down to _____

After a long day at work, she sat down to enjoy _____

After a long day at work, she sat down to enjoy a warm _____

After a long day at work, she sat down to enjoy a warm cup of _____

Masking

When training a transformer, we can generate **multiple training samples from a single text input**

After a long day at work, she sat down to enjoy a warm cup of coffee.

After _____

After a _____

After a long _____

After a long day _____

After a long day at _____

After a long day at work _____

After a long day at work, she _____

After a long day at work, she sat _____

After a long day at work, she sat down _____

After a long day at work, she sat down to _____

After a long day at work, she sat down to enjoy _____

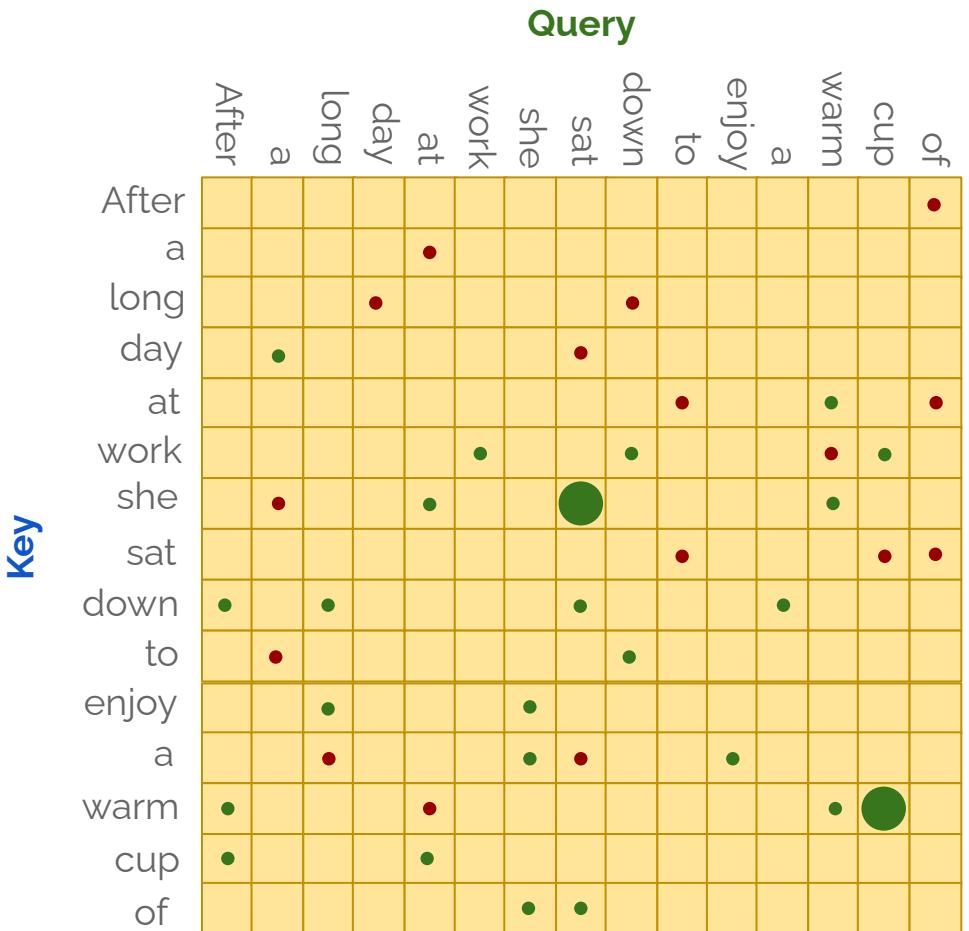
After a long day at work, she sat down to enjoy a warm _____

After a long day at work, she sat down to enjoy a warm cup of _____

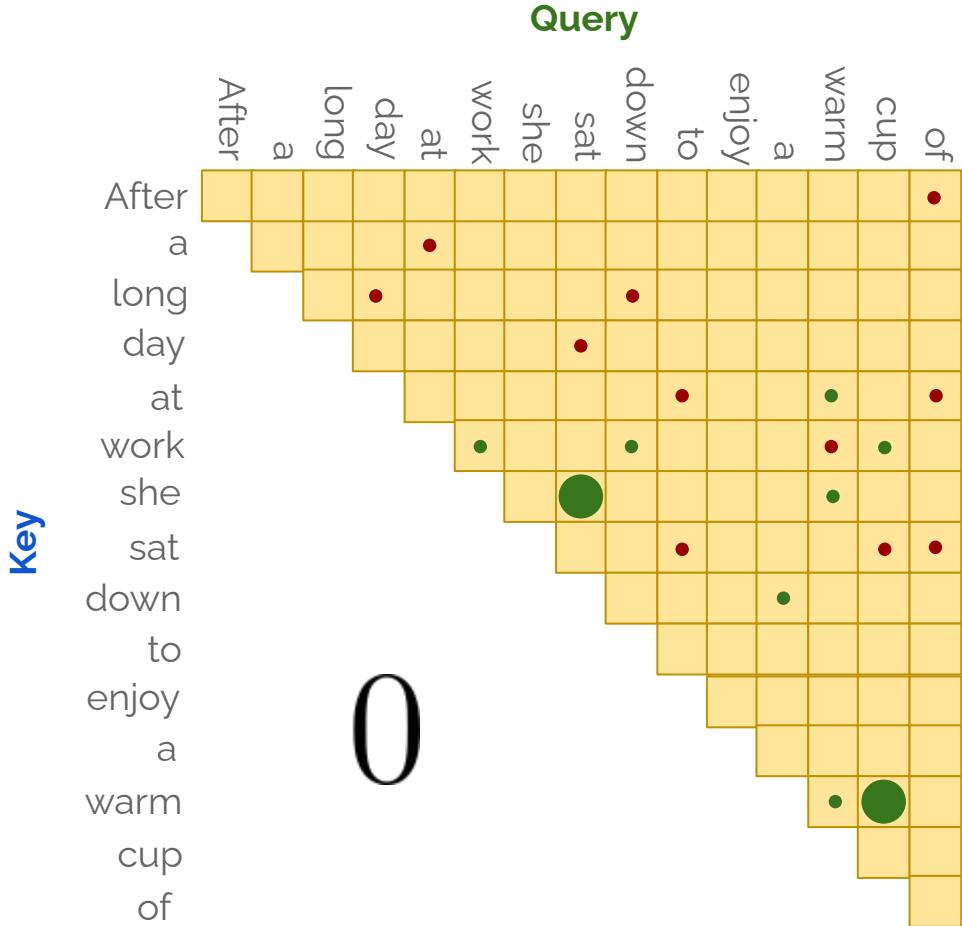
How can we **optimize** this process?

Masking

How can we force the attention mechanism to **only attend to previous words?**



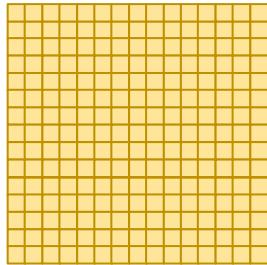
Masking



In this way, the transformer can only update the embeddings **based on the context of previous words**, thus allowing for a **efficient training**

Scalability

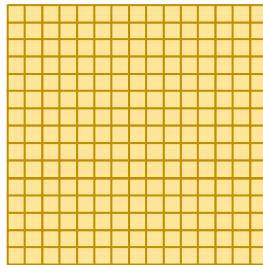
One of the main concerns of transformers concerns the **attention matrix**



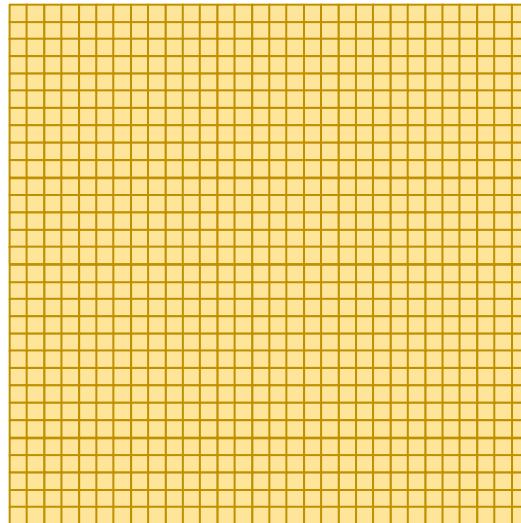
15 tokens

Scalability

One of the main concerns of transformers concerns the **attention matrix**



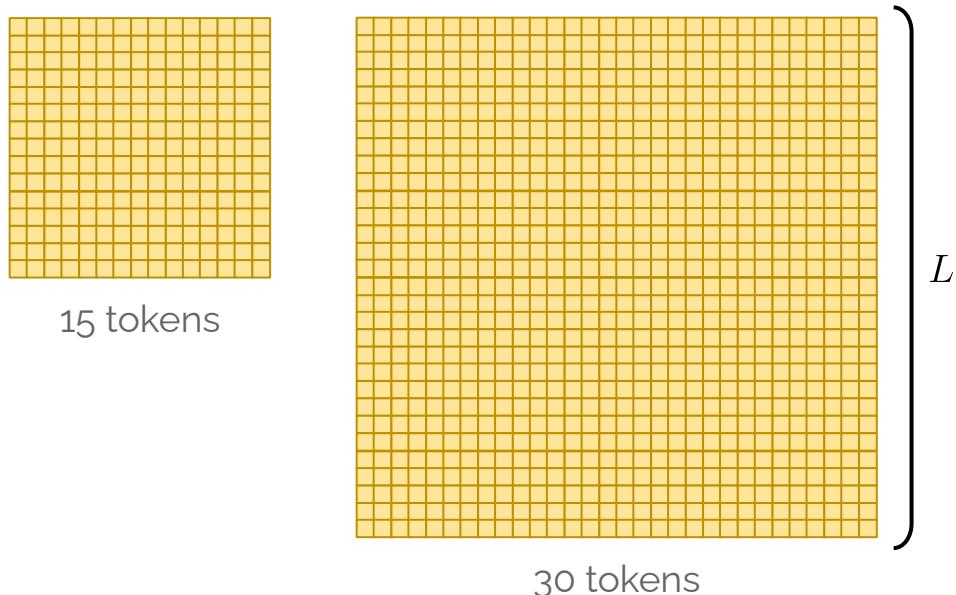
15 tokens



30 tokens

Scalability

One of the main concerns of transformers concerns the **attention matrix**



Number of computations: $O(L^2)$

Memory requirement: $O(L^2)$

Transformers **do not scale well** with respect to the **size of their context**

Scalability

One of the main concerns of transformers concerns the **attention matrix**



Number of computations: $O(L^2)$

Memory requirement: $O(L^2)$

Transformers **do not scale well** with respect to the **size of their context**

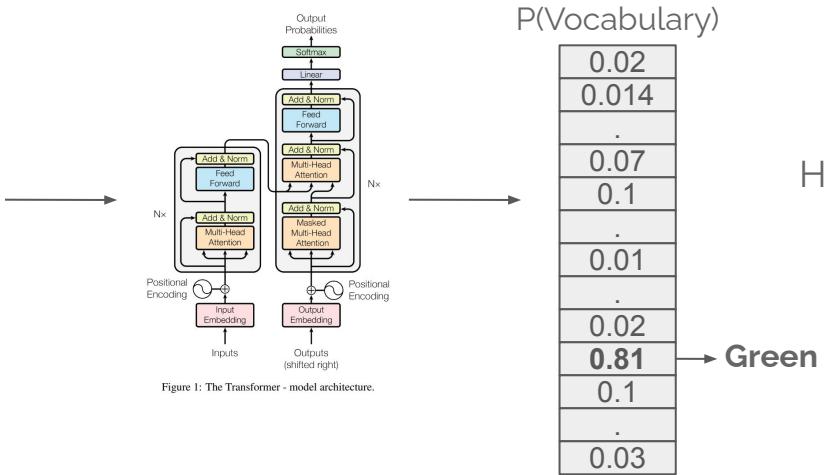
However, there exist solutions



$L \sim 10^6$

Temperature

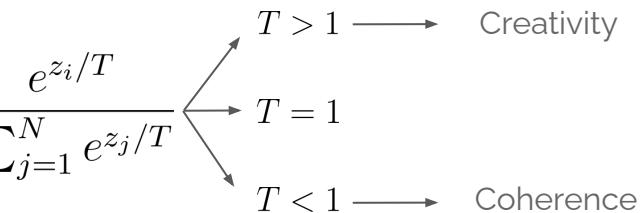
After a long day at work,
she sat down to enjoy a
warm cup of -----



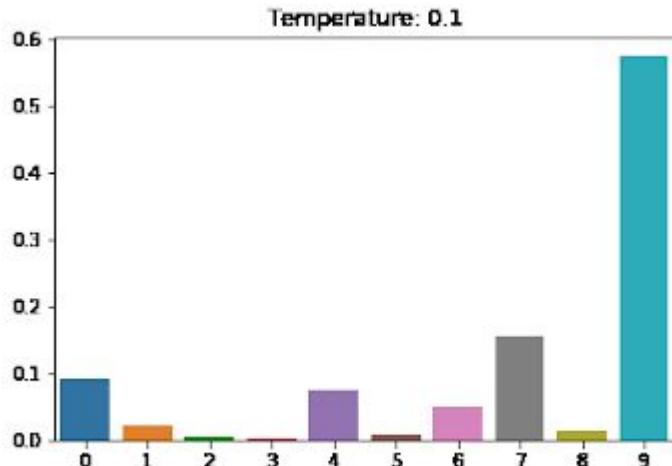
$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}$$

Temperature

$$\text{softmax}(z)_i = \frac{e^{z_i/T}}{\sum_{j=1}^N e^{z_j/T}}$$



Temperature



$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}}$$

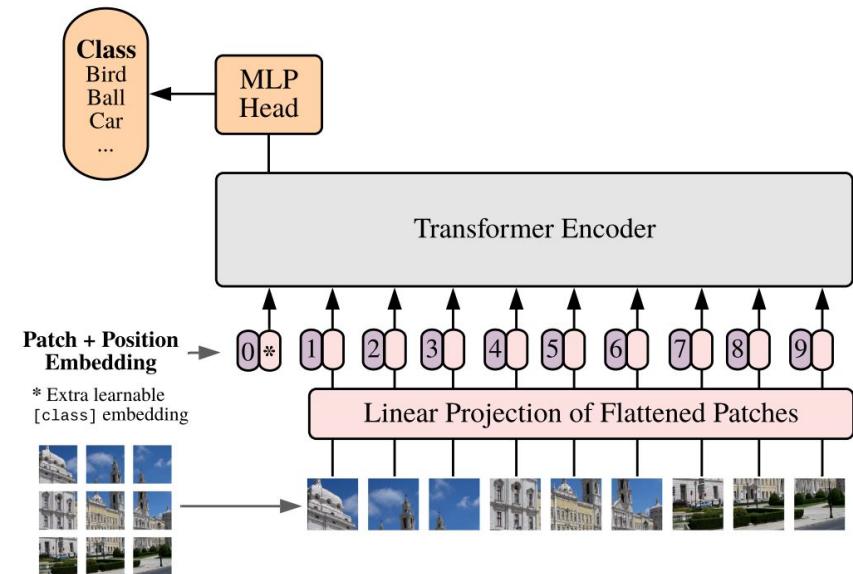
Temperature

$$\longrightarrow \text{softmax}(z)_i = \frac{e^{z_i/T}}{\sum_{j=1}^N e^{z_j/T}}$$

$T > 1 \longrightarrow$ Creativity
 $T = 1$
 $T < 1 \longrightarrow$ Coherence

Beyond NLP - Vision Transformers

- Vision Transformers (ViT) is the application of Transformers for image classification
- The ViT team knew they couldn't exactly mimic the language approach since self-attention on every pixel would be prohibitively expensive in computing time.
- Instead, they **divided the larger image into square units, or tokens**. The size is arbitrary, as the tokens could be made larger or smaller depending on the resolution of the original image (the default is 16 pixels on a side).



CNN vs ViT

- In CNNs, you start off being very local and slowly get a global perspective. But in transformers, with self-attention, even the very first layer of information processing makes connections between distant image locations.
- CNNs assume that nearby pixels are related (locality) and that different parts of an image are processed similarly (weight sharing) (*inductive bias*). These assumptions help CNNs learn effectively with limited training data. But this means, that given enough data, ViTs will learn better global relationships (because they don't assume locality).
- CNNs are smaller, so they are preferred for applications that need high latency

Video resources

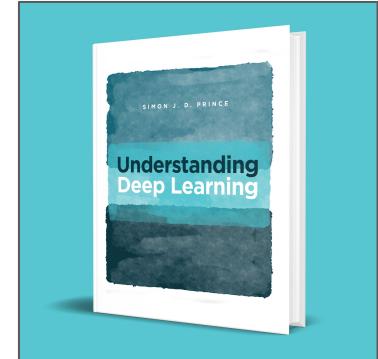
- [3Blue1Brown series](#)
- [Karpathy - Deep Dive into LLMs like ChatGPT](#)
- [Ari Seff - What are Transformer Neural Networks?](#)

Text resources

- [Understanding Deep Learning - Simon Prince](#)

Code resources

- [Minimal implementation of GPT \(PyTorch\)](#)



Questions

Jose González-Abad
Instituto de Física de Cantabria (CSIC-UC)



gonzabad@ifca.unican.es



[@jgonzalezab](https://twitter.com/jgonzalezab)



Image sources