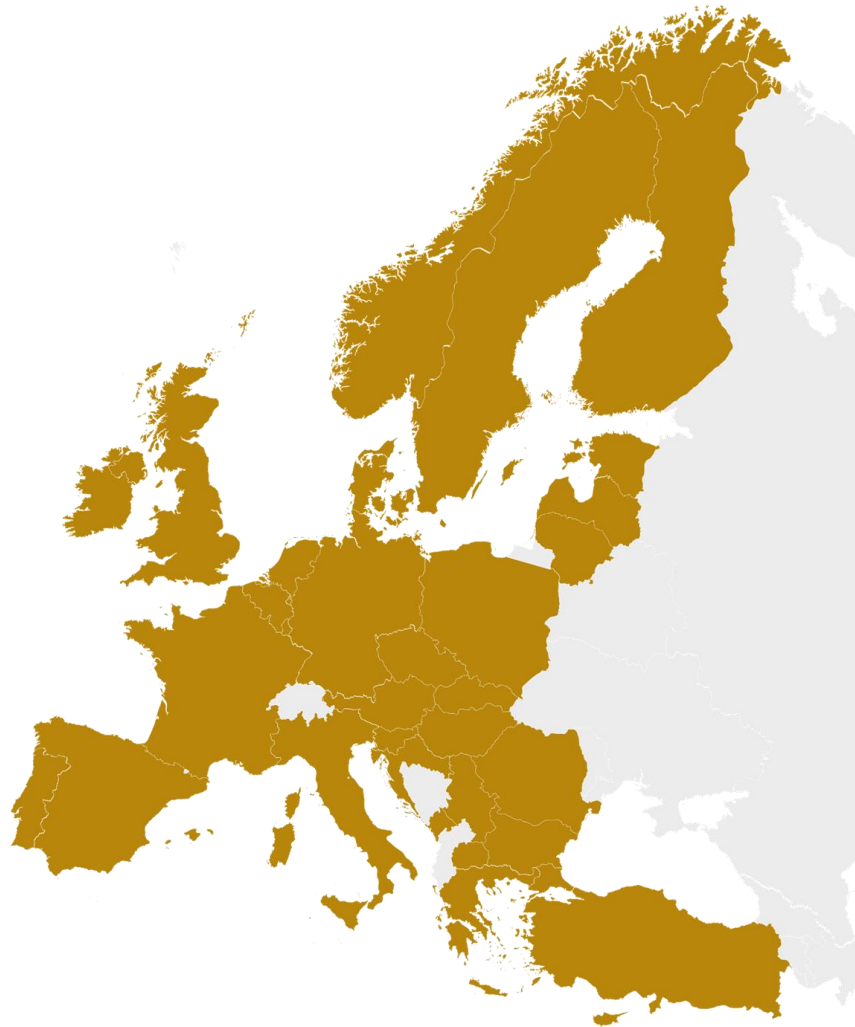# Large Language Models on Supercomputers

A brief overview

Speaker: Simeon Harrison
Trainer at EuroCC Austria

# EuroCC

## Fully funded EU project

- EuroCC is EU-funded international initiative aimed to support the uptake of AI and High-Performance Computing (HPC) in Europe

- Set up of 32 National Competence Centres (NCCs) across Europe

- EuroCC Austria is one of them

- Service Provider for AI, HPC and HPDA

# Need More Compute-Power?

## LUMI

- Third most powerful supercomputer in Europe and the 8th globally (Nov 2024)

- Sustained computing power (HPL) is 380 petaflops

- Over 262 000 AMD EPYC CPU cores

- Equipped with AMD Radeon Instinct MI250X GPUs

  https://www.lumi-supercomputer.eu/

## Leonardo

- 4th most powerful supercomputer in Europe and the 9th globally (Nov 24)

- Sustained computing power (HPL) is 239 petaflops

- Intel new gen Sapphire Rapids 56 cores

- Equipped with custom NVIDIA A100 SXM6 64GB GPUs

  https://leonardo-supercomputer.cineca.eu/

# European HPC Landscape

## EuroHPC JU systems

Different access modes:
Calls for Proposals
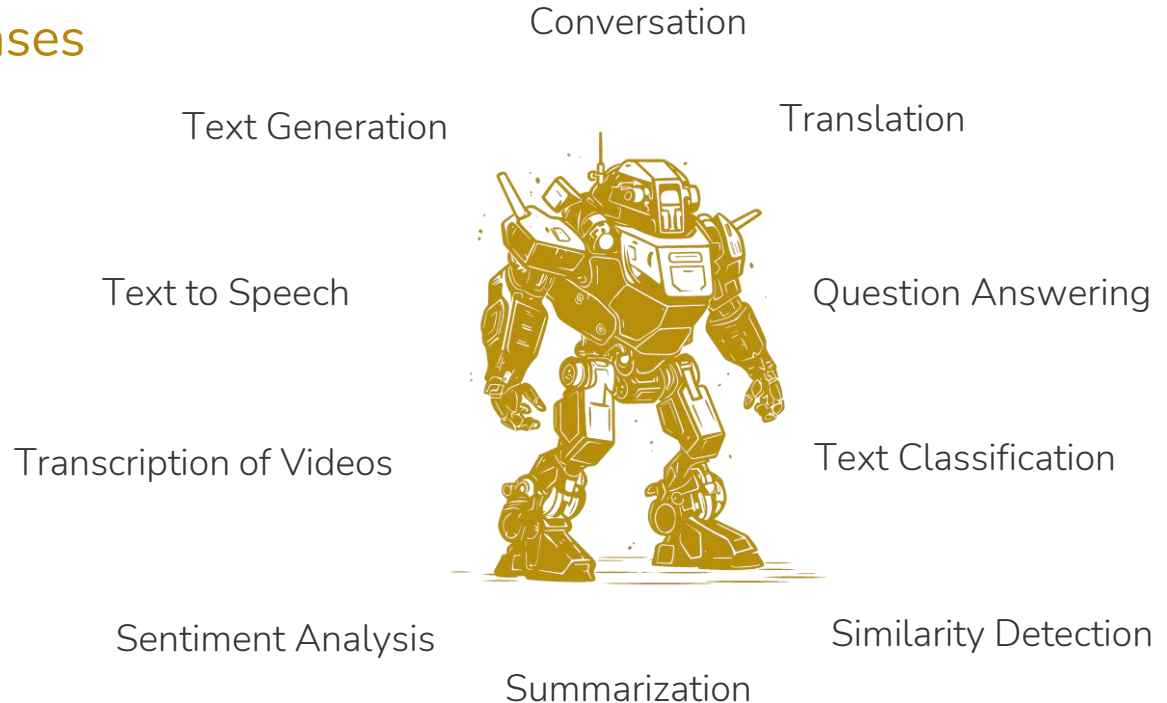
EuroHPC development access:
Opportunity to test the system

Applicants can request a small number of node hours to get acquainted with the supercomputers to further develop their software.
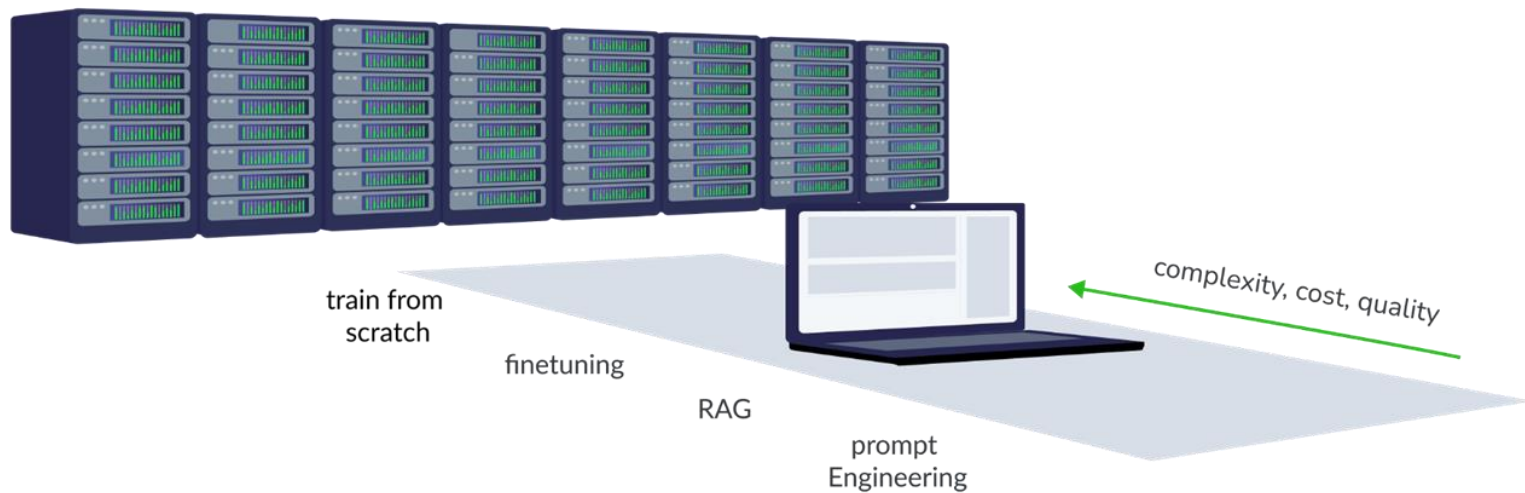
EURO AUSTRIA

LUMI
Finland

MELUXINA
Luxemburg

KAROLINA
Czech Republic

VEGA
Slovenia

DEUCALION
Portugal

DISCOVERER
Bulgaris

MARENOSTRUM 5
Spain

LEONARDO
Italy

# What can LLMs be used for?

## Many different use-cases

- Made possible by the transformer architecture

- Choose your model according to the use-case

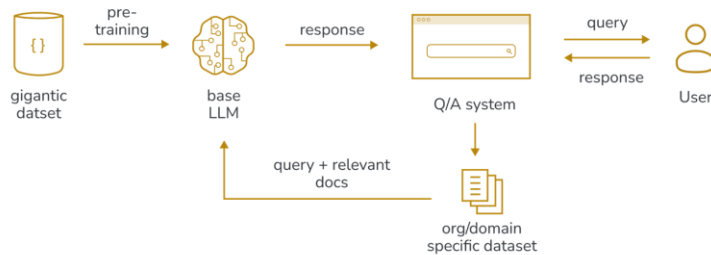- Is there a pre-trianed model for your use-case?

Conversation

Text Generation

Translation

Text to Speech

Question Answering

Transcription of Videos

Text Classification

Sentiment Analysis

Similarity Detection

Summarization

# How can you influence LLMs?



train from scratch

finetuning

RAG

prompt Engineering

complexity, cost, quality

# How can you use LLMs with your data?
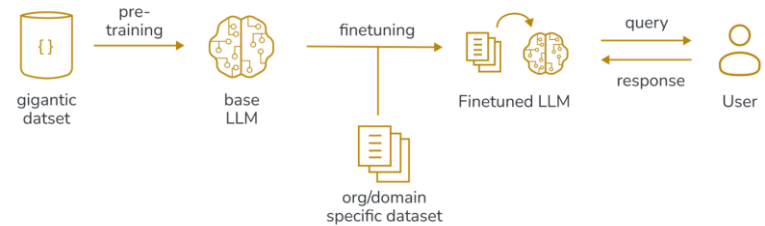
## RAG:
## Retrieval Augmented Generation



## Finetuning



- Ideal for tapping into company's knowledge DBs
- Minimises hallucinations by grounding response on retrieved evidence
- Can quickly adapt to changing data
- Makes it easier to interpret result

- Ideal if plenty of labelled data is available
- Teaches model domain specific vocabulary
- Company's writing/answer style is „baked" into model through fine-tuned parameters

# Prepare your Data

## Garbage in – garbage out

- Most underrated aspect of AI

- Most time consuming aspect of AI. Time spent in data preparation reflects in the quality of the product

- For fine-tuning you need labelled data

- Remember, that you are going to change the models parameters with your data
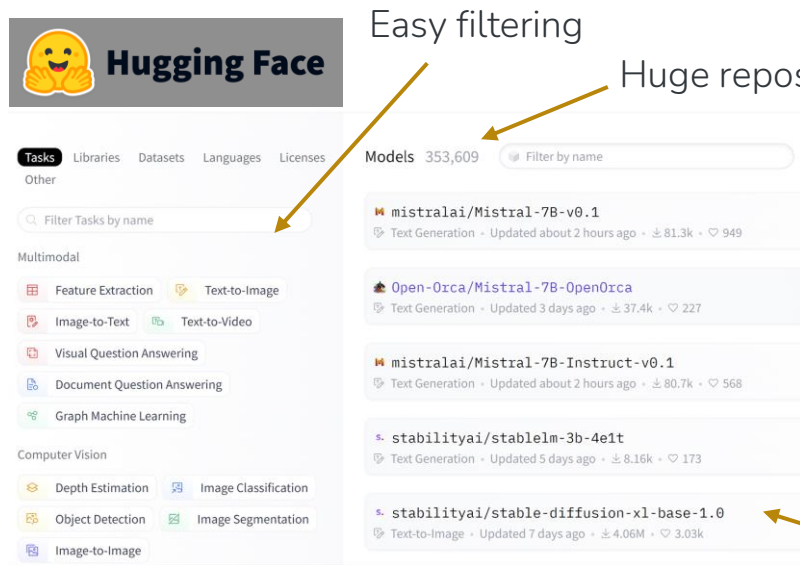
# The Hugging Face Ecosystem

From 🤗Transformers to the 🤗Hub

Speaker: Simeon Harrison
Trainer at EuroCC Austria

EURO
AUSTRIA

# Transformer Models

Spoilt for Choice at https://huggingface.co/



Easy filtering

Huge repository

All the relevant info

# Pick the Right Model

**EURO AUSTRIA**

mistralai / **Mistral-7B-Instruct-v0.2** ♡ like 1.04k

Text Generation  Transformers  PyTorch  Safetensors  mistral  finetuned  conversational

arxiv:2310.06825  License: apache-2.0

📦 Model card  ›≡ Files and versions  👏 Community 61

✎ Edit model card

## Model Card for Mistral-7B-Instruct-v0.2

The Mistral-7B-Instruct-v0.2 Large Language Model (LLM) is an improved instruct fine-tuned version of Mistral-7B-Instruct-v0.1.

For full details of this model please read our paper and release blog post.

# Hub and libraries



🤗

**Hugging Face Hub**

| Models | Datasets | Metrics | Docs |

**Hugging Face libraries**

| Tokenizers | Transformers | Datasets |

| Accelerate |

EURO
AUSTRIA

# Libraries

## Tokenizers

Very fast at tokenizing text, thanks to the Rust bakend.

Takes care of pre- and postprocessing, such as normalization and transformation to required format.

You can load tokenizers just like you load models.

## Transformers

Provides a standardized interface to a wide range of transformer models.

Supports PyTorch, TensorFlow and JAX and makes sitching between them easy.

Provides task-specific heads to fine-tune LLMs on downstream tasks

EURO
AUSTRIA

# Libraries

## Datasets

Simplifies the process of downloading and transforming data into required format.

Interface to thousands of datasets that can be found on the hub.

Interoperable with most common data wrangling tools and frameworks such as Pandas, NumPy etc.

## Accelerate

Adds a layer of abstraction to training loops with the goal of making code portable (local laptop to cluser)

Makes training on multi-GPU, multi-node easier to accomplish, as you can switch between parallelization strategies.

# Transformer Anatomy

Attention is really all you need?

Speaker: Simeon Harrison
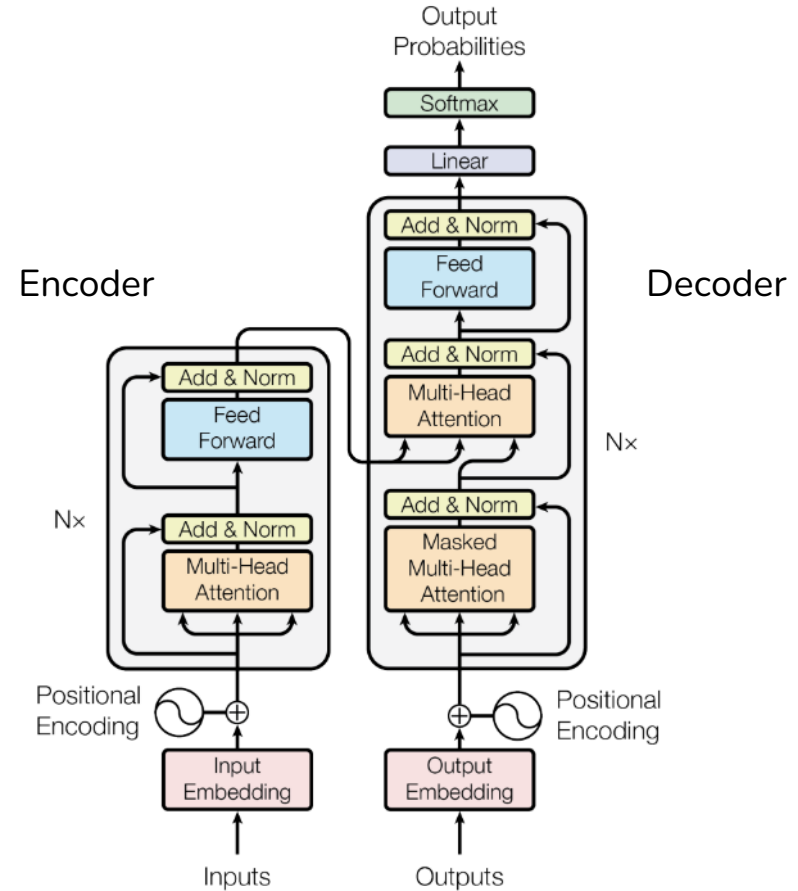Trainer at EuroCC Austria

EURO
AUSTRIA

# Transformer Anatomy

## The original architecture

A transforer consists of a encoder and/or decoder block.

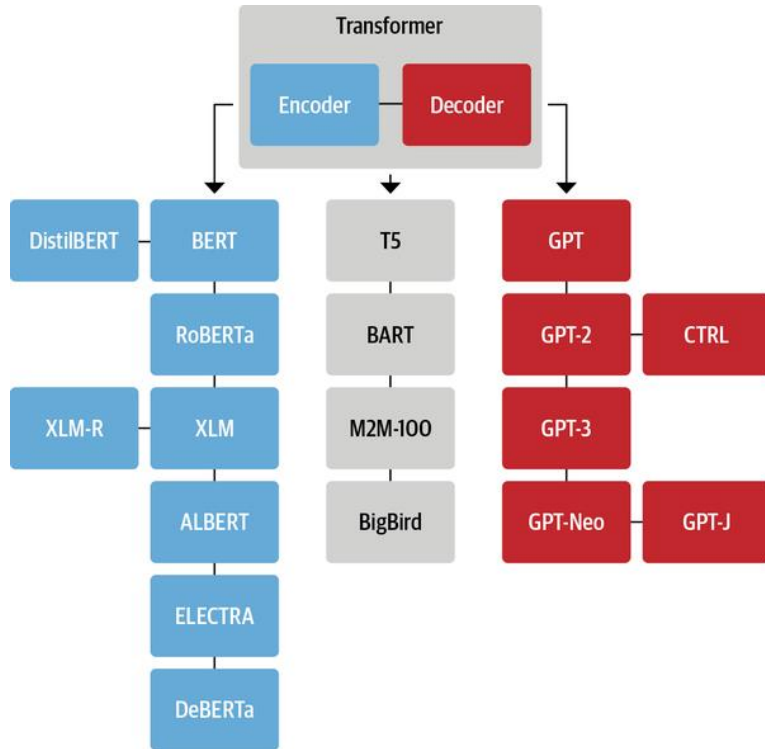Words (tokens) are input as numerical representations (embeddings).

About 1/3 of all parameters are in the multi-head attention blocks

About 2/3 of all parameters are in the feed forward networks (also known as multi layer perceptron)



Source: "Attention Is All You Need", Vaswani et al.

# Transformer Family



**Encoder only:**
These models ecxel at text classification, named entity recognition, and question answering

**Decoder only:**
Very good at predicting the next word in a sequence, therefore mostly used for text generation

**Encoder-Decoder:**
These models are often used for machine translation or summerization tasks.

Source: Natural Language Processing with Transformers, O'Reilly

# Context Is All You Need

## Embeddings

Here, we will refer to "word" instead of "token", as it makes the content easier to explain.

A word embedding comes as a multi dimensional vector (e.g. 12.000 dim).

The initial word embedding in all of the examples of the word „mole" is the same.

The European mole is a mammal

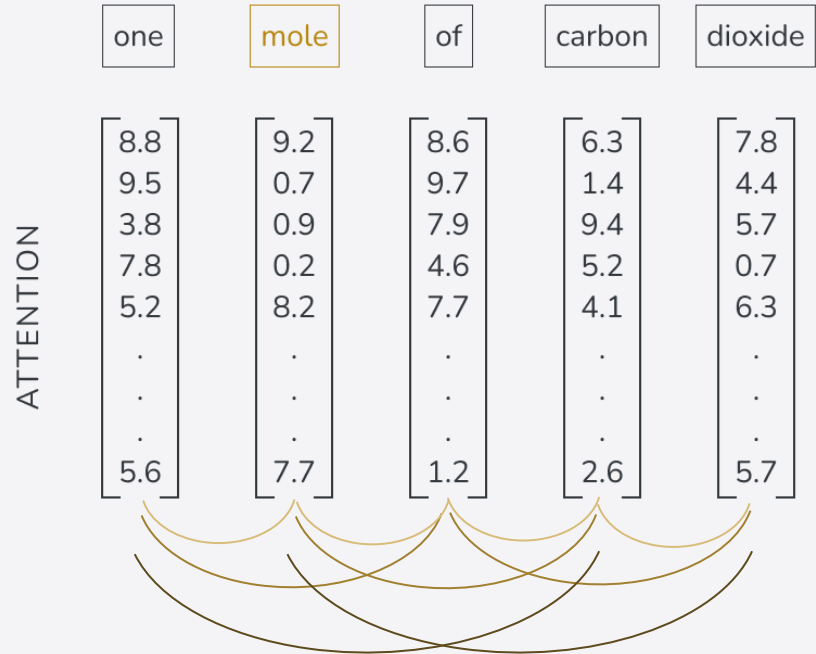$6.02 \times 10^{23}$

One mole of carbon dioxide

Take a biopsy of the mole

EURO
AUSTRIA

# Context Is All You Need

## Attention

The word „mole" should be represented by a **unique vector** in the embedding space, depending on ist context.

An **attention** block should **compute the vectors that you need to add** to the original, generic vector to get it to the correct, meaningfull, rich representation, depending on the context in which the word is used.

| one | mole | of | carbon | dioxide |
|---|---|---|---|---|

ATTENTION

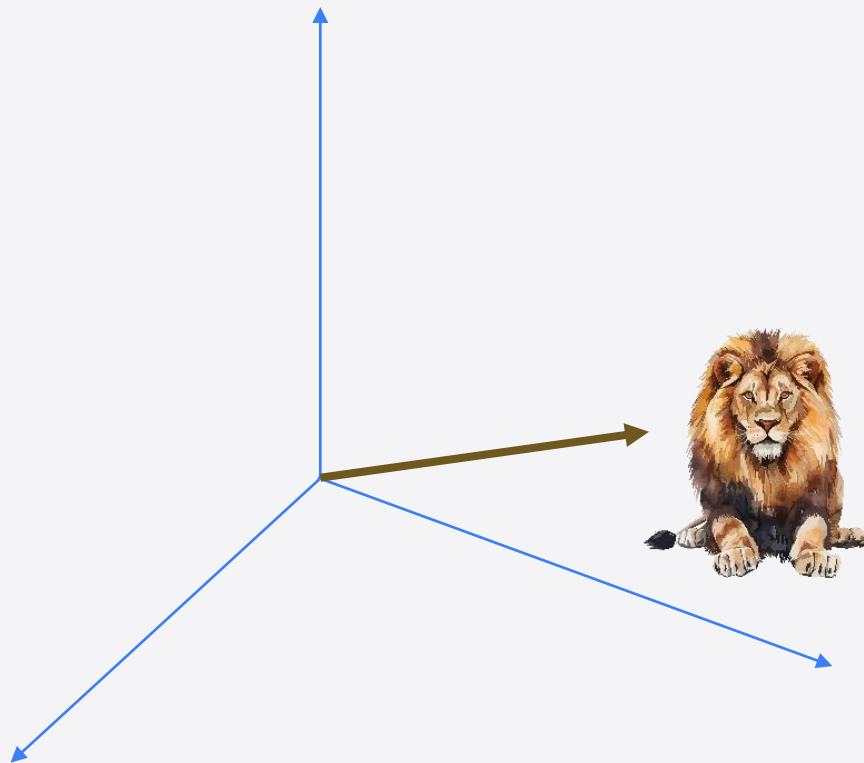| 8.8 | 9.2 | 8.6 | 6.3 | 7.8 |
|---|---|---|---|---|
| 9.5 | 0.7 | 9.7 | 1.4 | 4.4 |
| 3.8 | 0.9 | 7.9 | 9.4 | 5.7 |
| 7.8 | 0.2 | 4.6 | 5.2 | 0.7 |
| 5.2 | 8.2 | 7.7 | 4.1 | 6.3 |
| . | . | . | . | . |
| . | . | . | . | . |
| . | . | . | . | . |
| 5.6 | 7.7 | 1.2 | 2.6 | 5.7 |

# Context Is All You Need

## Lion

We associate the word „lion" with a big cat, living wild on the African continent.
We probably imagine a majestic predator with a big mane.

The embedding of the word „lion" is a vector with a certain length and direction within the embedding space.
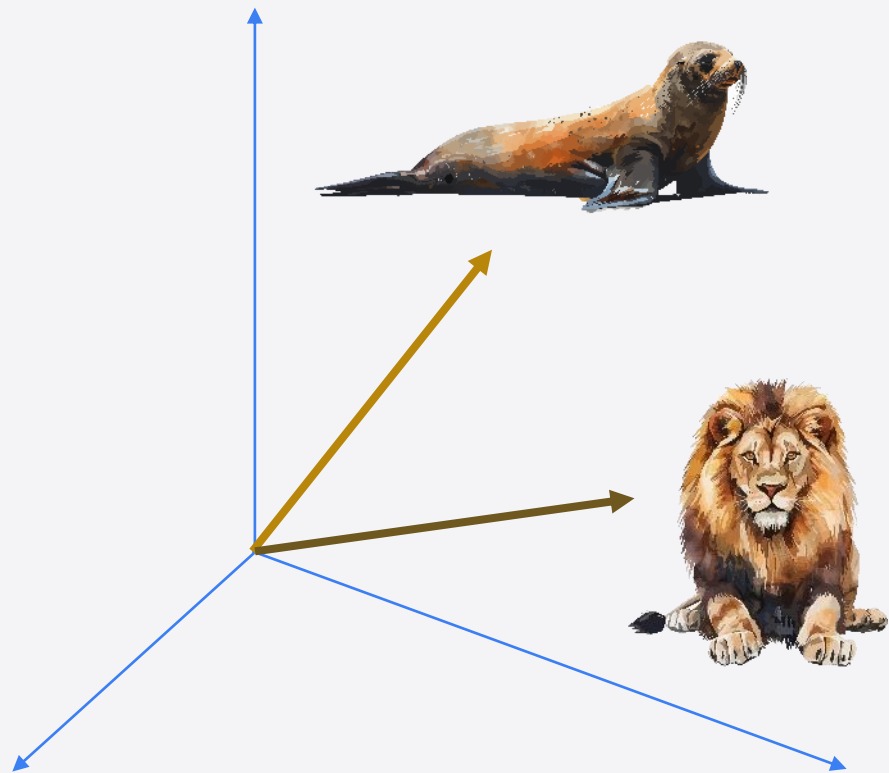
# Context Is All You Need

## Sea Lion

However, as soon we add the word „sea" infront of „lion" we imagine a totally different animal.

The same goes for the embedding. The attiontion mechanism needs to update the direction and length of the vector so that it represents the animal in question correctly.
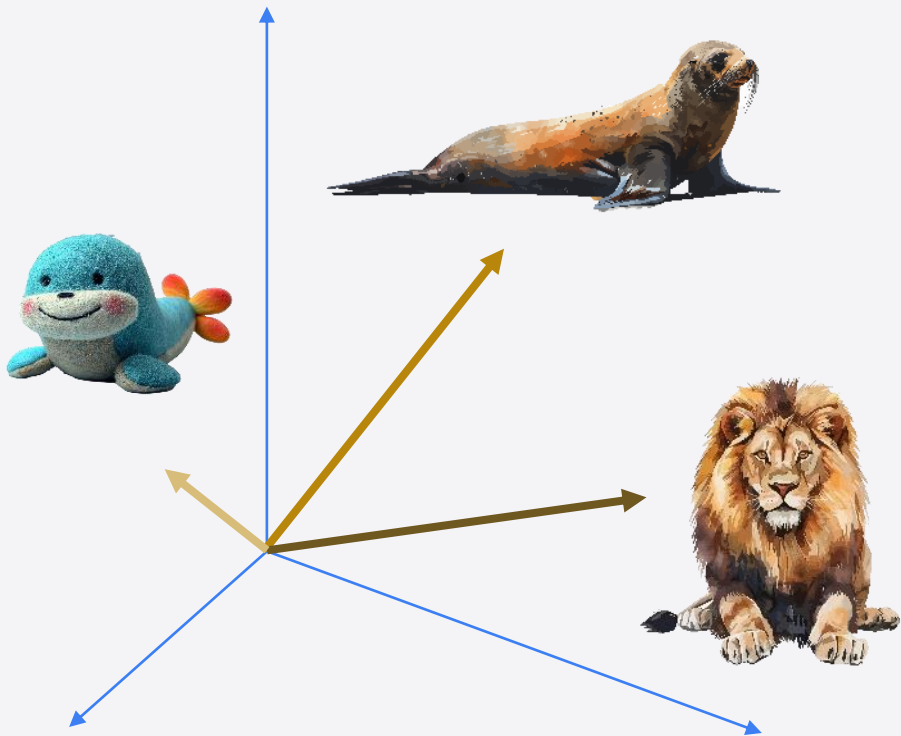
# Context Is All You Need

## Sea Lion Cuddly Toy

The context depends on more than just the immediate words to the left and right.
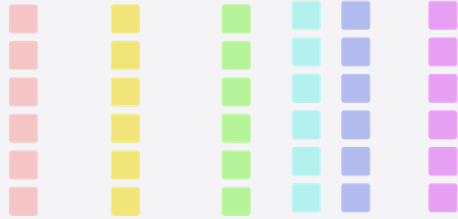
The embedding of „sea lion cuddly toy" will certainly be very different of just „lion".

In order to achieve that the vector for „lion" needs to attend to all the other words in the input (context size).



EURO AUSTRIA

# Self Attention

# Self Attention



the fat ginger cat sat on a mat

$$\vec{E}_1 \quad \vec{E}_2 \quad \vec{E}_3 \quad \vec{E}_4 \quad \vec{E}_5 \quad \vec{E}_6 \quad \vec{E}_7 \quad \vec{E}_8$$

# Attention! Queries, Keys and Values

The Attention Score

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

Q.....query matrix

K.....key matrix

V.....value matrix

d.....dimension of (smaller) query-key space

# Attention! Queries, Keys and Values

## Query

The query vector is obtained by multiplying a matrix full of trainable parameters with the embedding vector of that word (token).
Each word has its own query vector:

$$\vec{Q_i} = W_Q \vec{E_i}$$

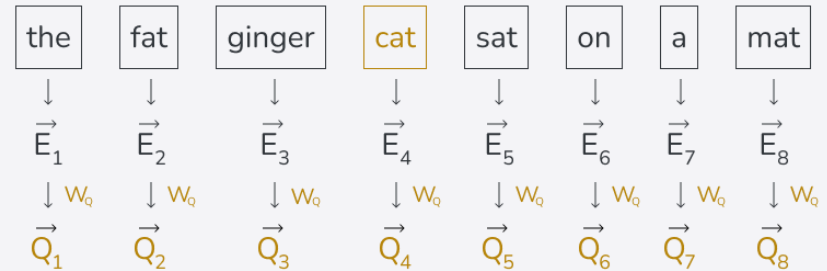It mappes the embedding vector to a much smaller dimensional space (e.g. 128 dimensions)

| the | fat | ginger | cat | sat | on | a | mat |
|---|---|---|---|---|---|---|---|
| ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| $\vec{E_1}$ | $\vec{E_2}$ | $\vec{E_3}$ | $\vec{E_4}$ | $\vec{E_5}$ | $\vec{E_6}$ | $\vec{E_7}$ | $\vec{E_8}$ |
| ↓ $W_Q$ | ↓ $W_Q$ | ↓ $W_Q$ | ↓ $W_Q$ | ↓ $W_Q$ | ↓ $W_Q$ | ↓ $W_Q$ | ↓ $W_Q$ |
| $\vec{Q_1}$ | $\vec{Q_2}$ | $\vec{Q_3}$ | $\vec{Q_4}$ | $\vec{Q_5}$ | $\vec{Q_6}$ | $\vec{Q_7}$ | $\vec{Q_8}$ |

Any adjectives in front of me?

**Imagine it like this:**
One particular attention head is focussing on nouns. The query vector is like looking for labels with „adjective" on them to better understand the noun.
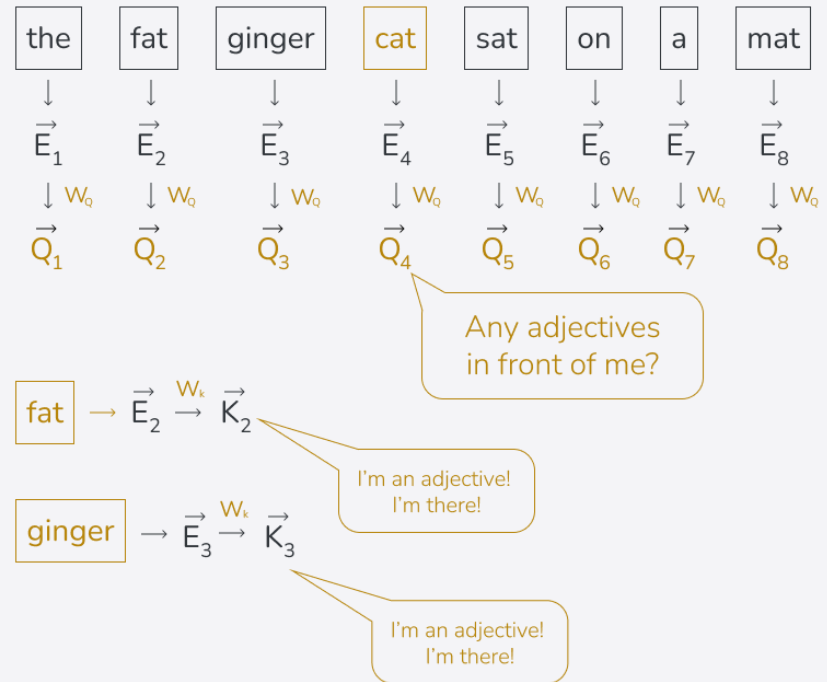In reality, this is much more abstract.

EURO AUSTRIA

# Attention! Queries, Keys and Values

## Key

The key vector is also obtained by multiplying a matrix full of trainable parameters with the embedding vector of that word (token). Each word also has its own key vector: $\vec{K_i} = W_K\vec{E_i}$

**Imagine it like this:**
We are still dealing with the particular attention head that is focussing on nouns. The key is answering the question the query raised.

# Attention! Queries, Keys and Values
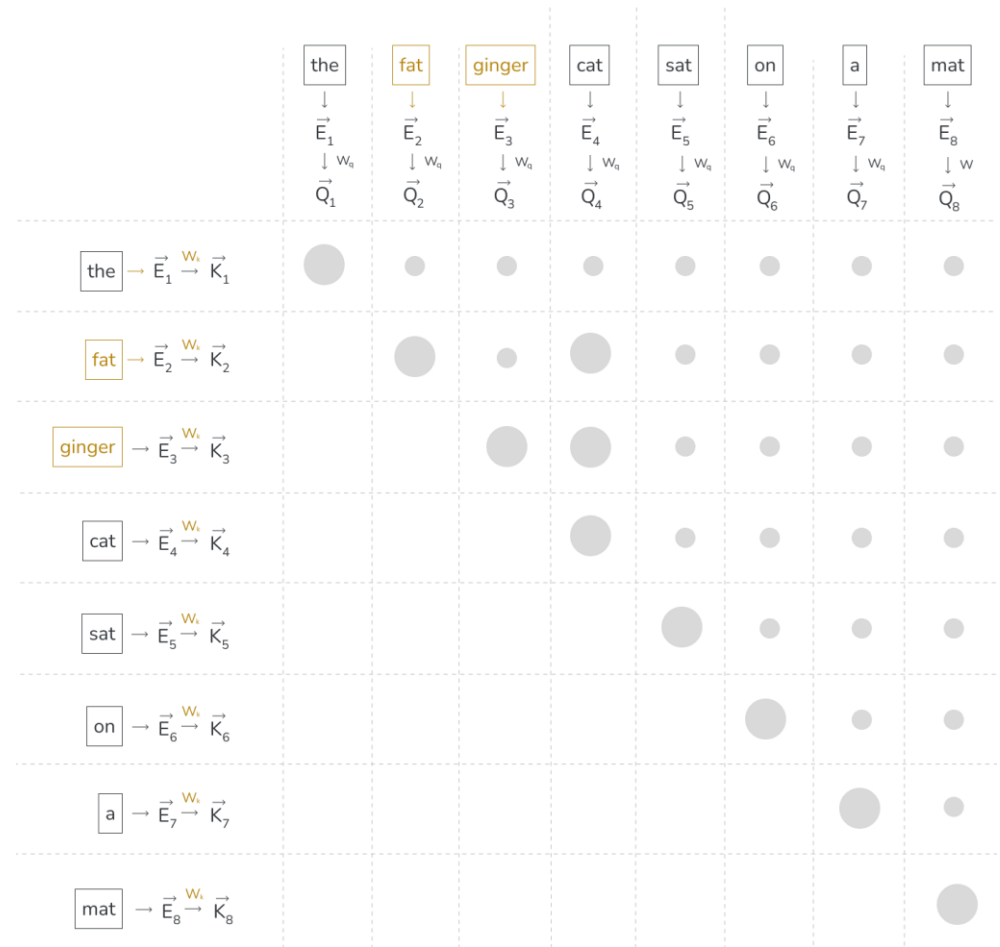
## Query – Key

Compute the dot product with each query-key pair, to determine how well the key matches the query. Where the queries and keys align, the dot product is larger.

**Imagine it like this:**
With our previous example, the dot product of the key vectors of „fat" and „ginger" with the vector of „cat" yields the largest result. The embeddings of „fat" and „ginger" attend to „cat".
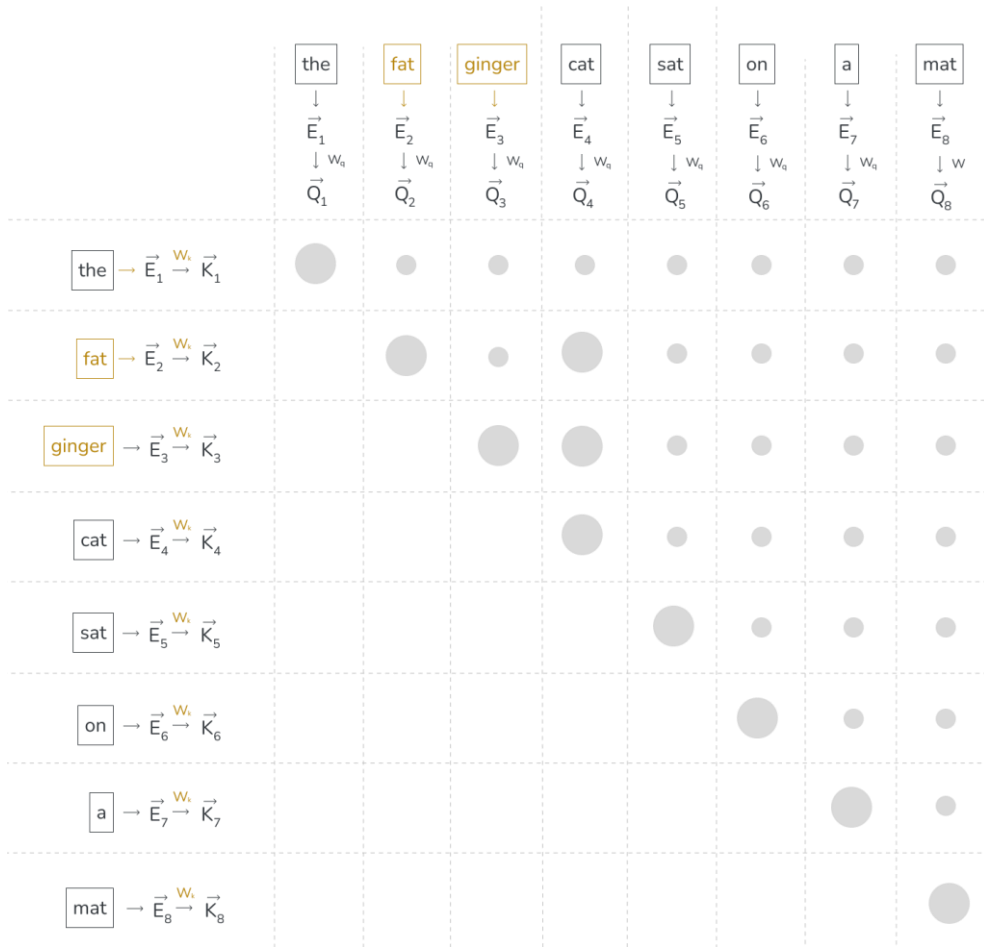
# Attention! Queries, Keys and Values

## Attention Pattern

Lower left dot products are masked, as we want the model to predict every next word during training. To prevent data leakage, future words should not influence previous ones.

The size of the attention pattern is the context size squared.
This is why the context size can be a substantial bottleneck.

# Attention! Queries, Keys and Values

## Attention

So far, we have determined which word is relevant to which other word (dot product of query and key).

We would like to use this as a score for how relevant every word is to update the meaning of other words.

For numerical stability the dot product is devided by the square root of the dimension of the query-key space.

To normalize the numbers to be between 0 and 1 we apply softmax.

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

# Attention! Queries, Keys and Values

## Value

.The value matrix multiplied by the embedding of the preceding word results in the value vector.

$$\vec{V_i} = W_K \vec{E_i}$$

| | the | fat | ginger | cat | sat | on | a | mat |
|---|---|---|---|---|---|---|---|---|
| | $\vec{E_1}$ | $\vec{E_2}$ | $\vec{E_3}$ | $\vec{E_4}$ | $\vec{E_5}$ | $\vec{E_6}$ | $\vec{E_7}$ | $\vec{E_8}$ |
| $\vec{E_1} \xrightarrow{W_v} \vec{V_1}$ | $1.0\ \vec{V_1}$ | $0.0\ \vec{V_1}$ | $0.0\ \vec{V_1}$ | $0.0\ \vec{V_1}$ | $0.0\ \vec{V_1}$ | $0.0\ \vec{V_1}$ | $0.0\ \vec{V_1}$ | $0.0\ \vec{V_1}$ |
| | + | + | + | + | + | + | + | + |
| $\vec{E_2} \xrightarrow{W_v} \vec{V_2}$ | $0.0\ \vec{V_2}$ | $1.0\ \vec{V_2}$ | $0.0\ \vec{V_2}$ | $0.42\ \vec{V_2}$ | $0.0\ \vec{V_2}$ | $0.0\ \vec{V_2}$ | $0.0\ \vec{V_2}$ | $0.0\ \vec{V_2}$ |
| | + | + | + | + | + | + | + | + |
| $\vec{E_3} \xrightarrow{W_v} \vec{V_3}$ | $0.0\ \vec{V_3}$ | $0.0\ \vec{V_3}$ | $1.0\ \vec{V_3}$ | $0.58\ \vec{V_3}$ | $0.0\ \vec{V_3}$ | $0.0\ \vec{V_3}$ | $0.0\ \vec{V_3}$ | $0.0\ \vec{V_3}$ |
| | + | + | + | + | + | + | + | + |
| $\vec{E_4} \xrightarrow{W_v} \vec{V_4}$ | $0.0\ \vec{V_4}$ | $0.0\ \vec{V_4}$ | $0.0\ \vec{V_4}$ | $0.0\ \vec{V_4}$ | $0.0\ \vec{V_4}$ | $0.0\ \vec{V_4}$ | $0.0\ \vec{V_4}$ | $0.0\ \vec{V_4}$ |
| | + | + | + | + | + | + | + | + |
| $\vec{E_5} \xrightarrow{W_v} \vec{V_5}$ | $0.0\ \vec{V_5}$ | $0.0\ \vec{V_5}$ | $0.0\ \vec{V_5}$ | $0.0\ \vec{V_5}$ | $0.0\ \vec{V_5}$ | $0.0\ \vec{V_5}$ | $0.0\ \vec{V_5}$ | $0.0\ \vec{V_5}$ |
| | + | + | + | + | + | + | + | + |
| $\vec{E_6} \xrightarrow{W_v} \vec{V_6}$ | $0.0\ \vec{V_6}$ | $0.0\ \vec{V_6}$ | $0.0\ \vec{V_6}$ | $0.0\ \vec{V_6}$ | $0.99\ \vec{V_6}$ | $1.0\ \vec{V_6}$ | $0.0\ \vec{V_6}$ | $0.0\ \vec{V_6}$ |
| | + | + | + | + | + | + | + | + |
| $\vec{E_7} \xrightarrow{W_v} \vec{V_7}$ | $0.0\ \vec{V_7}$ | $0.0\ \vec{V_7}$ | $0.0\ \vec{V_7}$ | $0.0\ \vec{V_7}$ | $0.0\ \vec{V_7}$ | $0.0\ \vec{V_7}$ | $1.0\ \vec{V_7}$ | $1.0\ \vec{V_7}$ |
| | + | + | + | + | + | + | + | + |
| $\vec{E_8} \xrightarrow{W_v} \vec{V_8}$ | $0.0\ \vec{V_8}$ | $0.0\ \vec{V_8}$ | $0.0\ \vec{V_8}$ | $0.0\ \vec{V_8}$ | $0.0\ \vec{V_8}$ | $0.0\ \vec{V_8}$ | $0.0\ \vec{V_8}$ | $0.0\ \vec{V_8}$ |
| | ‖ | ‖ | ‖ | ‖ | ‖ | ‖ | ‖ | ‖ |
| | $\Delta\vec{E_1}$ | $\Delta\vec{E_2}$ | $\Delta\vec{E_3}$ | $\Delta\vec{E_4}$ | $\Delta\vec{E_5}$ | $\Delta\vec{E_6}$ | $\Delta\vec{E_7}$ | $\Delta\vec{E_8}$ |

# Attention! Queries, Keys and Values

## Value

Eeach value vector is then multiplied by the corresponding weight for this word and the results summed up.

The result $\Delta\vec{E_i}$ is the change, that needs to be added to the origial embedding to get an updated, richer meaning.

Since this happens to every word in the sequence, we end up with a set of more refined embeddings.

cat
$\downarrow$
$\vec{E}_4$

the $\rightarrow \vec{E}_1 \xrightarrow{W_v} \vec{V}_1$    $0.0\,\vec{V}_1$

$+$

fat $\rightarrow \vec{E}_2 \xrightarrow{W_v} \vec{V}_2$    $0.42\,\vec{V}_2$

$+$

ginger $\rightarrow \vec{E}_3 \xrightarrow{W_v} \vec{V}_3$    $0.58\,\vec{V}_3$

$+$

cat $\rightarrow \vec{E}_4 \xrightarrow{W_v} \vec{V}_4$    $0.0\,\vec{V}_4$

$+$

sat $\rightarrow \vec{E}_5 \xrightarrow{W_v} \vec{V}_5$    $0.0\,\vec{V}_5$

$+$

on $\rightarrow \vec{E}_6 \xrightarrow{W_v} \vec{V}_6$    $0.0\,\vec{V}_6$

$+$

a $\rightarrow \vec{E}_7 \xrightarrow{W_v} \vec{V}_7$    $0.0\,\vec{V}_7$

$+$

mat $\rightarrow \vec{E}_8 \xrightarrow{W_v} \vec{V}_8$    $0.0\,\vec{V}_8$

cat
$\downarrow$
$\vec{E}_4$
$+$
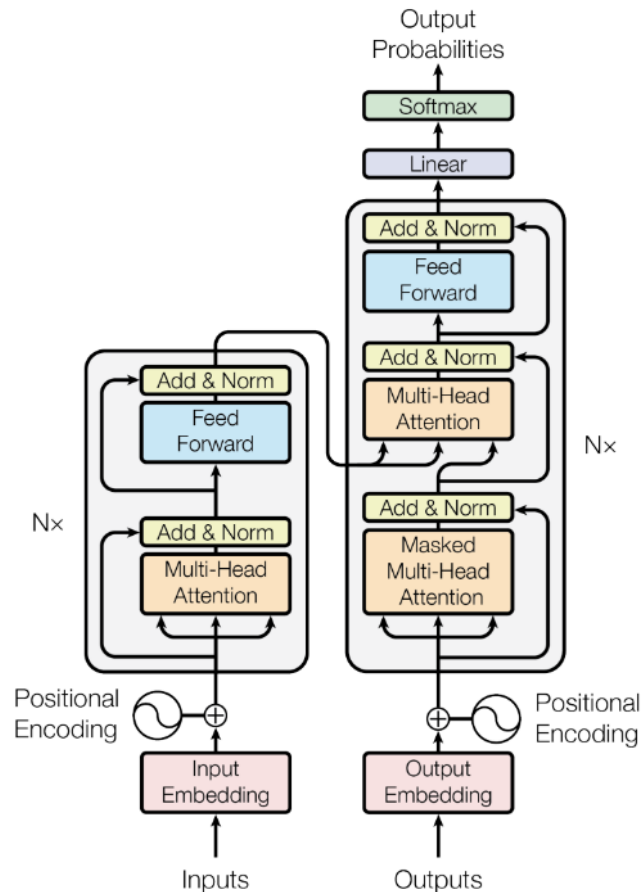$\Delta\vec{E}_4$
$\|$
$\vec{E}_4'$

# Attention

## Multi-Head Attention

The attention mechanism we just looked at, is done several times in parallel.

Each attention head is focussing on different features of the embeddings and computes ist own $\Delta\overrightarrow{E_i}^{(j)}$

The resulting $\Delta\overrightarrow{E_i}^{(j)}$ vectors, each suggesting the necessary change, are added to the orginal embedding.

This can be done in parallel on GPUs.



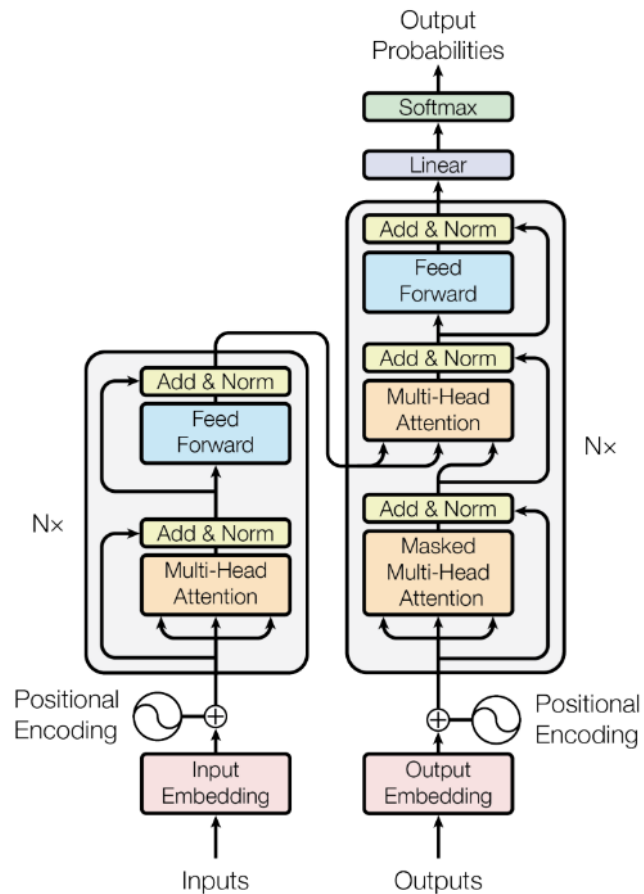Source: "Attention Is All You Need", Vaswani et al.

# Attention

## Cross Attention

Cross attention is almost the same as self attention.

Difference:

- query and key maps act on different data sets (e.g. 2 different languages in machine translation)
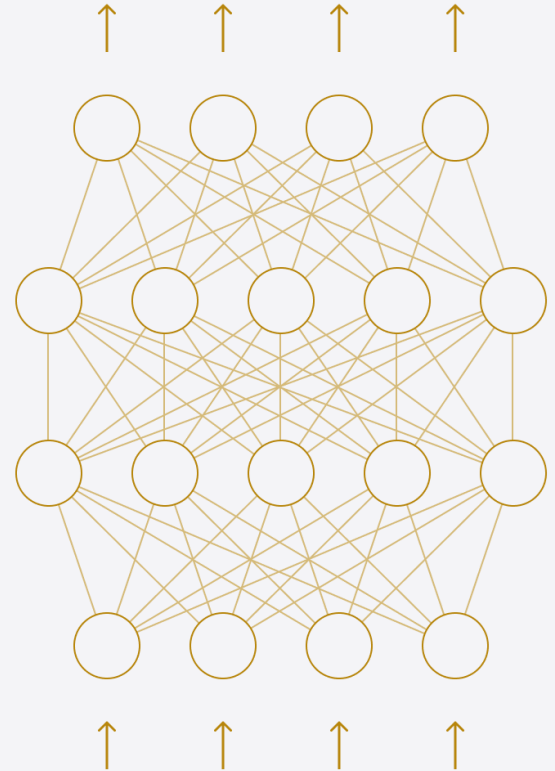- no masking, since there is no issue of later words affecting earlier ones.

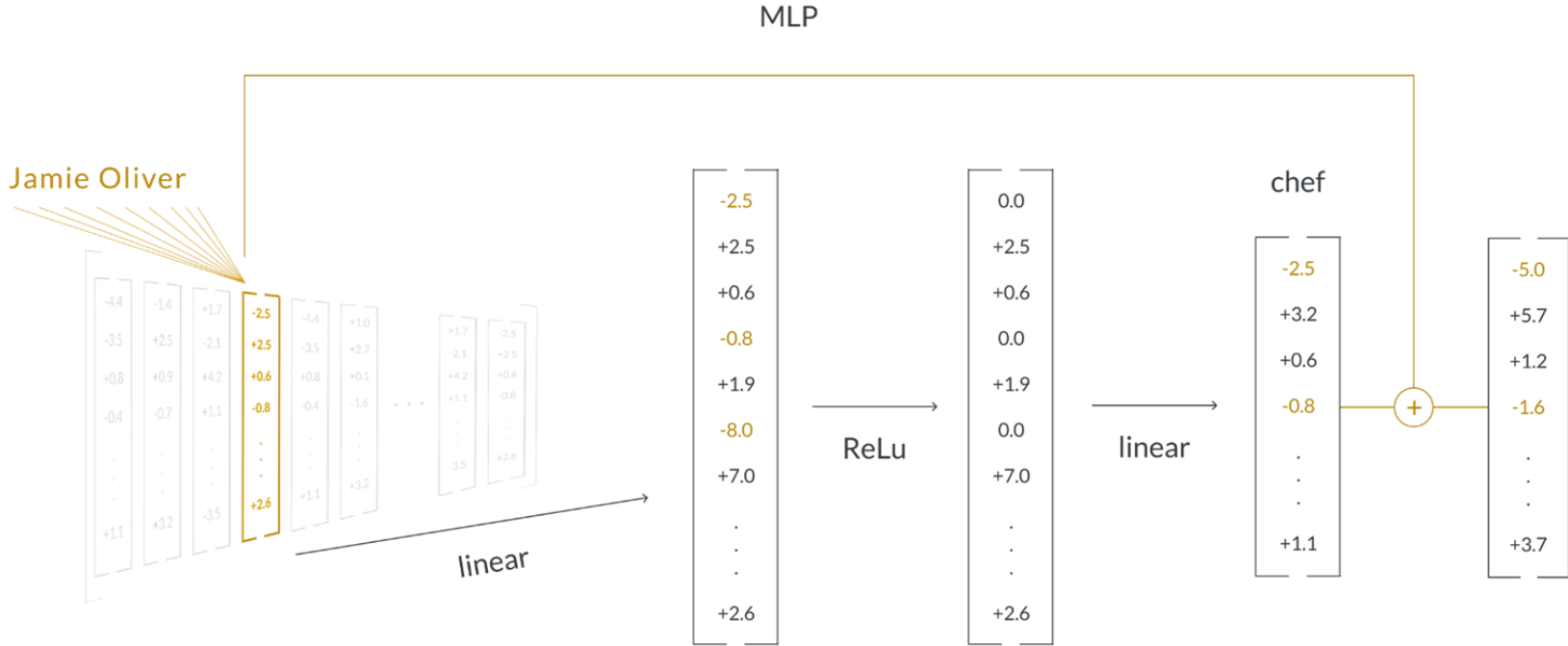Source: "Attention Is All You Need", Vaswani et al.

# Feed Forward Network

## Multi Layer Perceptron (MLP)

- Home to approx. 2/3s of all parameters

- This is where „knowledge" is baked in

- Source of halluzinations

- Facts that are associated with input embeddings are added to the input embeddings

- Each embedding vector can be processed independently -> parallelization!



EURO
AUSTRIA

# Feed Forward Network

# Tokenization & Embeddings

Turning words into numbers

Speaker: Simeon Harrison
Trainer at EuroCC Austria

EURO
AUSTRIA

# From Text to Tokens

## Character Tokenization

- Only needs 256 integers

- Good to tokenize rare words

- Helps with misspellings

- Looses linguistic structure

- Rarely used in practice

Largelanguagem
odelsonsupercom
puters

EURO
AUSTRIA

# From Text to Tokens

## Word Tokenization

- Model does not have to learn words from characters

- Each word has specific ID

- Size of vocabulary explodes

- Model needs to learn different tokens for e.g. singular and plural

Large lanugage models on supercomputers

# From Text to Tokens

## Subword Tokenization

- Best of both worlds

- Deal with complex words easily

- Frequent words remain as one token

Large lanugage models
on super com put ers

EURO
AUSTRIA

# From Tokens to Vectors

## Embeddings

Tokens are mapped to unique integers accoding to the vocabulary size of the tokenizer.

Now, the tokens need to be embedded, which means turned into a vector representation.

This is done by an embedding layer of a model. The model takes each token ID and looks it up in an embedding matrix. The embedding matrix is a learned set of weights that maps each token ID to a corresponding high-dimensional vector (embedding).

Woman ⟶
$$\begin{bmatrix} 1.2 \\ 0.4 \\ 3.6 \\ 5.0 \\ 3.9 \\ 8.5 \\ 2.1 \\ 0.6 \\ 7.0 \\ 8.2 \\ 4.2 \\ 9.8 \\ 3.0 \\ \vdots \\ 1.1 \\ 4.3 \\ 2.7 \\ 3.3 \\ 0.0 \end{bmatrix}$$

EURO
AUSTRIA

# From Tokens to Vectors

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

## Embeddings

Each word/token has a unique direction in the embedding space

Similar words point in a similar direction.

E (uncle)

E (man)

E (aunt)

E (woman)

EURO
AUSTRIA

# THANK YOU