

# Llm presentation

Transformer architecture, activation...

Tilian Bourachot

June 16, 2024

# Table des matières

- 1 Introduction
- 2 Fonctionnement des LLM
- 3 The Transformer model

# Large Language Models (LLMs)

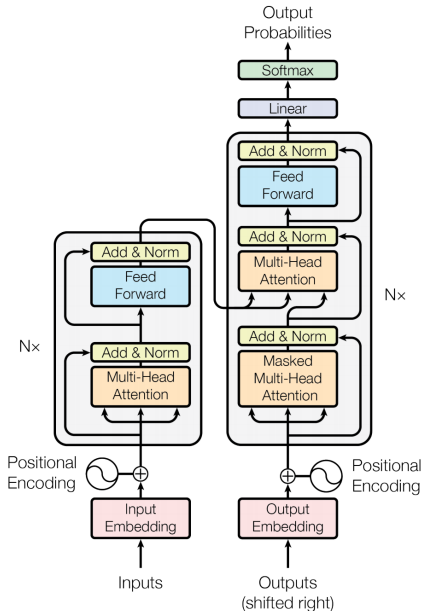
- **Large Language Models (LLMs)** use advanced deep learning techniques to intelligently understand and generate natural language. Their goal is to generate coherent and plausible text based on a set of training data. Their functionality is based on transformer architectures.
- **Language models (LM)**: A model trained to predict the next word in a sequence, given the words that come before it. It is a probabilistic model capable of generating text that follows a certain linguistic style or pattern (e.g., GPT, BERT).

# Example Models

- **n-gram models**
- **Hidden Markov Models (HMMs)**
- **Recurrent Neural Networks (RNNs):**
  - RNNs process sequences using internal state (memory), making them suitable for language modeling but struggling with long-range dependencies due to the vanishing gradient problem.
- **Long Short-Term Memory (LSTM):**
  - A specialized RNN designed to capture long-term dependencies using gates that control information flow in and out of the network's memory state.
- **Gated Recurrent Unit (GRU):**
  - A simplified version of LSTM with different gating mechanisms, often faster to train but performance varies by task.

- LLM as GPT3 or GPT4 are simply LMs that are trained on a very large amount of text and have a very large number of parameters.
- Important advantage : we can train them unsupervised on a large corpus of data and then fine-tune with a limited amount of data for different tasks.

# Introduction to Transformers



# Introduction to Transformers

- ① Enables machines to understand the context of the words in sentences more effectively.
- ② Can handle words in parallel (unlike RNNs, which are sequential).
- ③ Uses attention mechanisms to understand word context within sentences.

# Introduction to Transformers

- Pre-trained models (BERT, GPT) are trained on vast amounts of text data.
- They can be fine-tuned for specific tasks, such as:
  - Sentiment analysis
  - Language translation
  - Text summarization



# Introduction to Transformers

- Specific kind of neural network
- Uses self-attention mechanisms to weigh the relevance of different words in the input when making predictions.
- Suitable for parallel computing (in RNN (LSTM, GRU), the sequence of tokens must be processed sequentially because each one depends on the previous tokens of the sequence)
- Self attention process allows each token to be computed as a weighted sum of all tokens in the sequence, with the weight determined by the attention mechanism.
- Speed ++

# Introduction to Transformers

- The core of ChatGpt is a transformer.
- ChatGPT is based on the GPT version of the transformer. The GPT models are trained to predict the next word in a sequence of words, given all the previous words. They process text from left to right, which makes them well-suited for text generation tasks.
- GPT-3: 175 billion parameters.
- These models are trained using self-supervised learning on a large corpus of text, which means they predict the next word in a sentence (GPT) or a word based on surrounding words (BERT).
- Because they are exposed to a large amount of texts, these models learn grammar, facts about the world, reasoning abilities, and also biases in the data they are trained on.

# How does Transformer work ?

- **First step : Tokenization**



- The input is broken up into little pieces (tokens), words for text, pixels for picture.

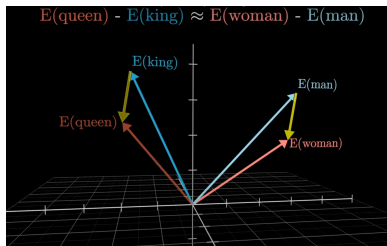
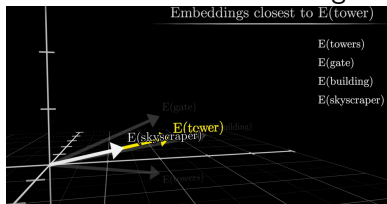
# Tokenization

- Embedding matrix has a single columns for each words, wich are represented as vector. Initially, they are set randomly.
- In GPT3, each words have 12,288 coordinates, and total tokens : 50,257 Each words coordinate are called weight. As long the model is trained, the weights are changing, and so the embedding matrix.

All words, ~ 50k																						
aah	aardvark	aardwolf	aargh	ab	aback	abacterial	abacus	abalone	abandon	...	zygoid	zygomatic	zygomorphic	zygosis	zygote	zygotic	zyne	zymogen	zymosis	zzz		
+1.0	+4.3	+2.0	+0.9	-1.5	+2.9	-1.2	+7.8	+9.2	-2.3	...	+0.6	+1.3	+8.4	-8.5	-8.2	-9.5	+6.6	+5.5	+7.3	+9.5		
+5.9	-0.8	+5.6	-7.6	+2.8	-7.1	+8.8	+0.4	-1.7	-4.7	...	-0.9	+1.4	-9.5	+2.3	+2.2	+2.3	+8.8	+3.6	-2.8	-1.2		
+3.9	-8.7	+3.3	+3.4	-5.7	-7.3	-3.7	-2.7	+1.4	-1.2	...	-7.9	-5.8	-6.7	+3.0	-4.9	-0.7	-5.1	-6.8	-7.7	+3.1		
-7.2	-6.0	-2.6	+6.4	-8.0	+6.7	-8.0	+9.4	-0.6	+9.4	...	+4.7	-9.1	-4.3	-7.5	-4.0	-7.5	-3.6	-1.7	-8.6	+3.8		
+1.3	-4.6	+0.5	-8.0	+1.5	+8.5	-3.6	+3.3	-7.3	+4.3	...	-6.3	+1.7	-9.5	+6.5	-9.8	+3.5	-4.6	+4.7	+9.2	-5.0		
+1.5	+1.8	+1.4	-5.5	+9.0	-1.0	+6.9	+3.9	-4.0	+6.2	...	+7.5	+1.6	+7.6	+3.8	+4.5	+0.0	+9.0	+2.9	-1.5	+2.1		
-9.5	-3.9	+3.2	-4.2	+2.3	-1.4	-7.2	-4.0	+1.4	+1.8	...	+3.0	+3.0	-1.4	+7.9	-2.6	-1.3	+7.8	+6.1	+4.0	-7.9		
+8.3	+4.2	+9.9	-6.9	+7.3	-6.7	+2.3	-7.4	+6.9	+6.1	...	-1.8	-8.5	+3.9	-0.9	+4.4	+7.3	+9.4	+7.0	-9.7	-2.8		
:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:		
-3.7	-2.0	-5.7	-6.2	+8.8	+4.7	-0.2	-5.4	-4.9	-8.8	...	-3.7	+3.9	-2.4	-6.3	-9.4	-8.6	+3.6	-0.9	+0.7	+7.9		

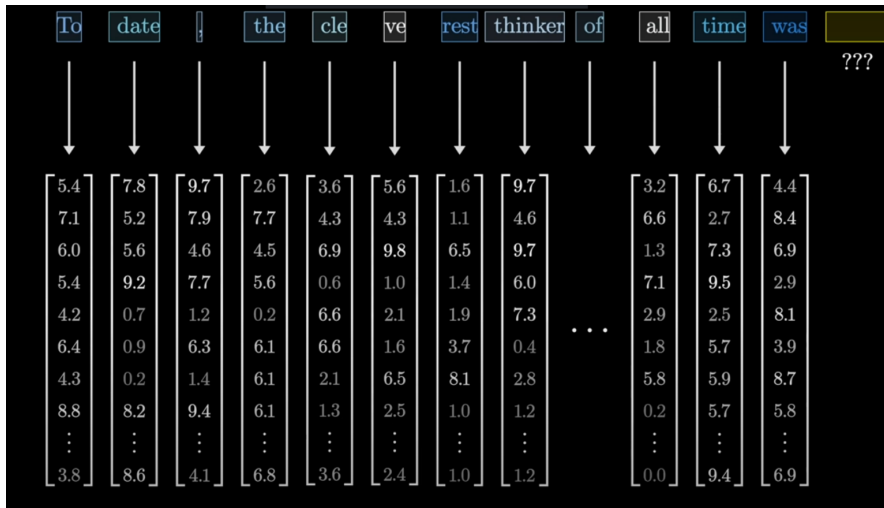
# Tokenization

- Words with closest meaning have closest vector



# Tokenization

- After several steps of attention, which we will detail next, the sentence is represented by a sequence of vectors.



# Tokenization

- To determine the next word, we use the transpose of the embedding matrix, referred to as the unembedding matrix. We compute the matrix product between this matrix and the last vector of the sequence. Applying softmax to the result yields a vector of size 1x50,257 containing the probabilities for each word.



# How does Transformer work ?

- The sequence of vectors of the sentence passes through into an operation call attention and who will able the vectors to talk to each other and to pass information to uptade their weights according to the sentence.
- For example : a machin Learning model and a fashion model, model has not the same meaning in both sentences.



# How does Attention work ?

- At the beginning, when a sentence is input, each token is represented by a vector that only corresponds to the word's meaning without context and its position in the sentence. Now, we are going to name each word's vector or embedding  $\vec{E}_i$ . The goal is to transform each embedding into new ones that inject the meanings from corresponding adjectives, for example.

# How does Attention work ?

- Each embedding also has a query vector  $\vec{Q}_i$  and a key vector  $\vec{K}_i$ , which are determined as:

$$\vec{Q}_i = W_q \cdot \vec{E}_i \quad \text{where} \quad W_q \text{ is the query matrix of size } 12,288 \times 128$$

Similarly, the key vector  $\vec{K}_i$  is determined using a similar matrix  $W_k$ .

The diagram illustrates the calculation of the query vector  $\vec{Q}_4$  from the embedding  $\vec{E}_4$  using the query matrix  $W_Q$ . The matrix  $W_Q$  is a  $12 \times 128$  matrix, and the embedding  $\vec{E}_4$  is a  $128 \times 1$  vector. The resulting query vector  $\vec{Q}_4$  is a  $12 \times 1$  vector.

The matrix  $W_Q$  is shown as a grid of values, with the first 12 rows and the first 12 columns highlighted in red. The values are:

+7.5	-3.2	+9.1	-5.3	+8.9	+8.7	+5.9	+2.6	+7.4	-4.1	...	+2.3
-9.6	-3.0	-7.0	+9.5	-0.4	-0.1	+2.8	-2.6	-7.2	+6.4	...	+0.2
-5.5	-8.0	+7.2	+9.4	+9.1	+8.0	+5.4	-3.3	-8.3	-1.8	...	-7.3
-8.8	+4.5	-9.7	+5.4	-7.0	-8.3	-8.1	+3.4	-5.0	-1.6	...	+7.1
+4.5	-4.5	-7.3	-8.8	-3.9	-4.7	-0.9	+3.6	+3.9	-4.3	...	-6.3
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
-9.0	+5.9	-8.4	+0.4	-3.8	+1.5	+9.1	+2.9	-9.2	-1.4	...	+0.7

The embedding  $\vec{E}_4$  is shown as a vertical vector of values:

2.9
2.4
1.0
0.2
9.2
6.6
7.8
2.8
5.8
0.6
⋮
9.7

The resulting query vector  $\vec{Q}_4$  is shown as a vertical vector of values:

⋮
⋮
⋮
⋮
⋮
⋮
⋮
⋮
⋮
⋮
⋮
⋮

# How does Attention work ?

- The entries of  $W_q$  and  $W_k$  are parameters of the model, learned from data.
- Difficult to analyse, somekind encode information into lower dimension space.

# How does Attention work ?

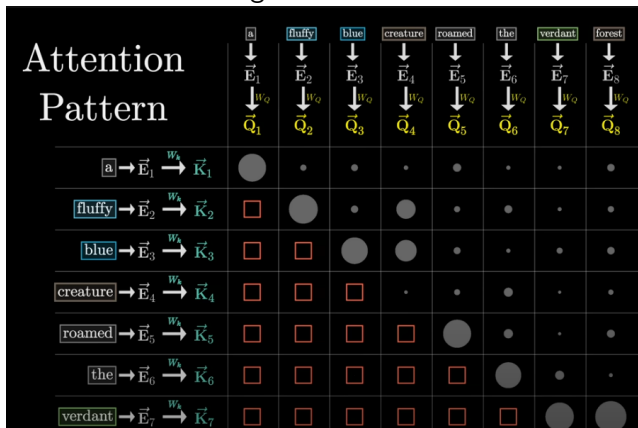
- We can consider that keys are answering queries (for instance  $\vec{Q}_i$ : "any adjectives in front of me?").
- Keys are matching the queries when  $\vec{K}_i$  and  $\vec{Q}_j$  are closely aligned, meaning  $\vec{E}_i$  has an impact/link with  $\vec{E}_j$  (example: adjective).
- To determine if they are matching, we use a dot product:  $\vec{K}_i \cdot \vec{Q}_j$  for all  $i, j$ .
- If the dot product has a small value, it indicates that they are unrelated.

# How does Attention work ?



# How does Attention work ?

- Attention: We can't allow later words to influence earlier words; otherwise, they can kind of give away the answer of what comes next. So all the values in red become  $-\infty$  because with softmax, these elements will become 0 and the column will stay normalized.
- This process is called masking



# Value Matrix

- Now, we want to update the embeddings, allowing words to pass information to other words they are relevant to. We use a value matrix  $W_v$  that we multiply by the embedding of the first word that is relevant (e.g., adjective). This gives a value vector that we add to the embedding of the second word (e.g., noun).

$$\vec{V}_i = W_v \cdot \vec{E}_i$$

	a	fluffy	blue	creature	roamed	the	verdant	forest
	$\vec{E}_1$	$\vec{E}_2$	$\vec{E}_3$	$\vec{E}_4$	$\vec{E}_5$	$\vec{E}_6$	$\vec{E}_7$	$\vec{E}_8$
a $\rightarrow \vec{E}_1 \xrightarrow{W_v} \vec{V}_1$	1.00 $\vec{V}_1$	0.00 $\vec{V}_1$	0.00 $\vec{V}_1$	0.00 $\vec{V}_1$	0.00 $\vec{V}_1$	0.00 $\vec{V}_1$	0.00 $\vec{V}_1$	0.00 $\vec{V}_1$
fluffy $\rightarrow \vec{E}_2 \xrightarrow{W_v} \vec{V}_2$	0.00 $\vec{V}_2$	1.00 $\vec{V}_2$	0.00 $\vec{V}_2$	0.42 $\vec{V}_2$	0.00 $\vec{V}_2$	0.00 $\vec{V}_2$	0.00 $\vec{V}_2$	0.00 $\vec{V}_2$
blue $\rightarrow \vec{E}_3 \xrightarrow{W_v} \vec{V}_3$	0.00 $\vec{V}_3$	0.00 $\vec{V}_3$	1.00 $\vec{V}_3$	0.58 $\vec{V}_3$	0.00 $\vec{V}_3$	0.00 $\vec{V}_3$	0.00 $\vec{V}_3$	0.00 $\vec{V}_3$
creature $\rightarrow \vec{E}_4 \xrightarrow{W_v} \vec{V}_4$	0.00 $\vec{V}_4$	0.00 $\vec{V}_4$	0.00 $\vec{V}_4$	0.00 $\vec{V}_4$	0.00 $\vec{V}_4$	0.00 $\vec{V}_4$	0.00 $\vec{V}_4$	0.00 $\vec{V}_4$
roamed $\rightarrow \vec{E}_5 \xrightarrow{W_v} \vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$	0.01 $\vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$	0.00 $\vec{V}_5$
the $\rightarrow \vec{E}_6 \xrightarrow{W_v} \vec{V}_6$	0.00 $\vec{V}_6$	0.00 $\vec{V}_6$	0.00 $\vec{V}_6$	0.00 $\vec{V}_6$	0.99 $\vec{V}_6$	1.00 $\vec{V}_6$	0.00 $\vec{V}_6$	0.00 $\vec{V}_6$

- and then,  $\vec{E}_3 + \Delta\vec{E}_3 = \vec{E}'_3$ : a more refined vector encoding a richer contextual meaning. And we have  $\vec{E}'_i$  for all  $i$ .



# Multi head attention

- This is only one head attention. For instance, one head can focus on adjectives, grammar, verbs, etc.
- Multi-head attention is when you run many of these operations in parallel, each with its own key, query, and value matrices:  $W_q(i)$ ,  $W_v(i)$ ,  $W_k(i)$ .
- GPT-3 uses 96 attention heads.
- This gives the model the capacity to learn many distinct ways that context changes the meaning.

# Back to the first draw

Core of the Transformer is Self Attention as in the original paper: “**Attention** is all you need”.

[arxiv.org/pdf/1706.03762.pdf](https://arxiv.org/pdf/1706.03762.pdf)

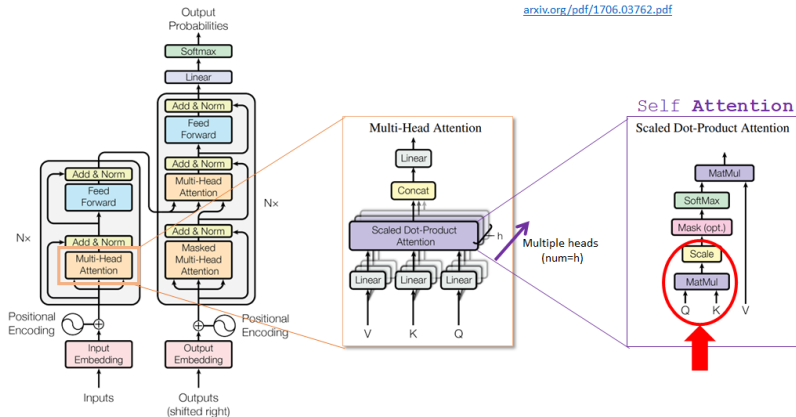
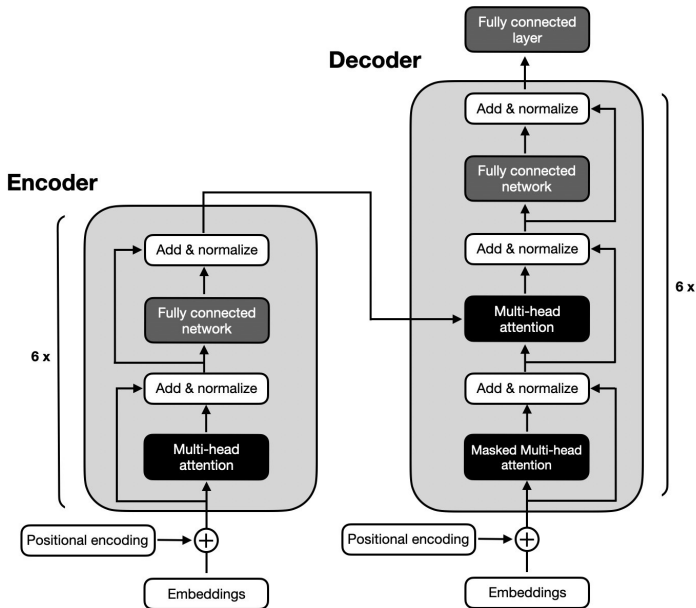


Figure 1: The Transformer - model architecture.

# Encoder/ Decoder



- **Encoder:** Responsible for understanding and extracting the relevant information from the input text. It then outputs a continuous representation (embedding) of the input text that is passed to the decoder.
  - Good for classification, sequence tagging, sentiment analysis
  - Exemple : Bert...
- The masked language (masking random word tokens in the input sequence and then training the model to predict the original masked tokens based on the surrounding context) and next sentence pretraining objectives allow BERT to learn rich contextual representations of the input texts.
- Typically need to be fine-tuned for various downstream tasks like sentiment analysis or specific tasks
- Cannot generate text (only understand text)

- **Encoder/Decoder:** Use encoder to understand and use decoder to generate output sequence
  - Transformer
  - Good for translation and summarization

- **The decoder** generates output word by word, and the predictions for position  $i$  can depend only on the known outputs at positions less than  $i$ .
  - Language Models
  - GPT2,3,4, ChatGPT
  - Phi3, Llama3
- Controlled by prompting (instruction)
- Few-shots > Zero- shots (shots=examples)