

GenAI

Session 1 : What's GenAI ?

2023-2024

em
lyon
business
school

 early makers
since 1872



Table of contents of session

01 Course Overview

02 What's going on with LLMs ?

03 Transformers : how does it works ?



You

- What do you want to learn ?
- What do you expect from this class ?

Me

- Revenue Operations @Side
- MSc in Statistics / ML
- Using & building AI automatisations/tools to optimize processes





Course overview

The world

Prompt Engineer

The Role

- Work with cross-functional teams to discuss product development
- Identify uses of AI tools
- Design, develop and refine AI-generated text prompts

Skills



- Additional certifications recommended

Salary

Junior: \$ 280,000
Average: \$ 327,000
Senior: \$ 375,000



Mamba

Mamba MOE

mambaByte MOE

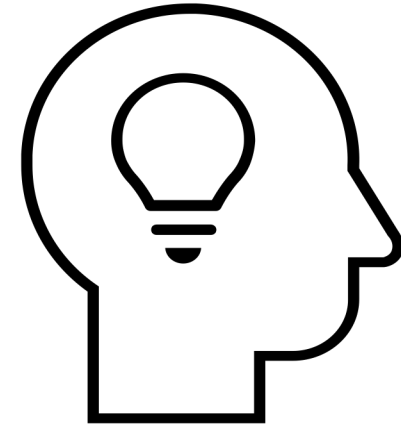
MambaByte MOE + self-reward fine-tuning

Multimodal MambaByte MOE + self-reward fine-tuning + cascade speculative drafting + LASER + DRpGS + AQLM Quantized by The Bloke



Objectives

- Learning the basics
- Testing/Prompting the new tools on the market
- Implementing small usecases



Structure

- 1. Introduction to GenAi through LLMs :** "Exploring Large Language Models"
- 2. Practical Applications: & Tuning a Model to a Specific Need :** “Leveraging LLMs in Business”
- 3. How to Make a Movie with AI?** " Tools, Techniques, and Boundaries "

What do think of this program/titles ?

-

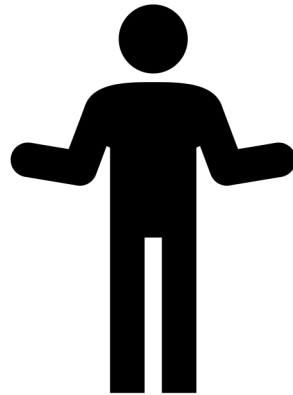




What is going on with LLMs ?

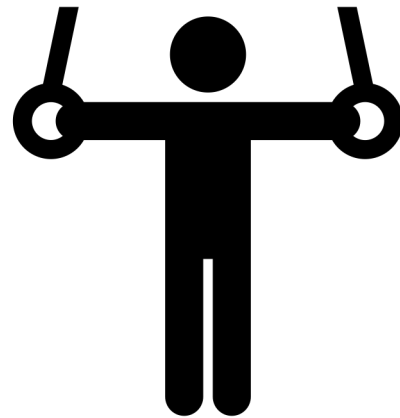
Context

But... the class is on GenAI, not LLMs ?



Context

Ok, so what's the difference between LLM and GenAI ?



What is Generative AI?

- GenAI is a type of Artificial Intelligence that creates new content based on what it has learned from existing content.
- The process of learning from existing content is called training and results in the creation of a statistical model.
- When given a prompt, GenAI uses this statistical model to predict what an expected response might be—and this generates new content.

Models in GenAI

**Generative Adversarial
Networks (GANs)**

Transformer Models

**Variational Autoencoders
(VAEs)**

**Restricted Boltzmann
Machines (RBMs)**

Diffusion Models

and « text » was, for a long time, a type a content we had a hard time to deal with (NLP tasks, in general)

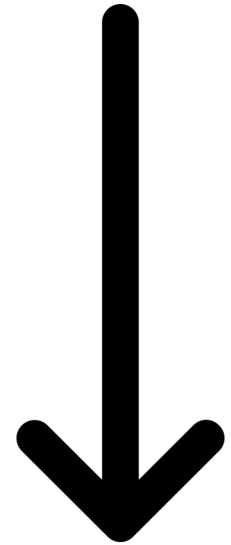
•



Definition

ChatGPT = **Generative Pre-trained Transformer**

which is a type of **Large Language Model**,
which is a type of **General-purpose Transformer**,
which is a type of **Artificial Neural Network**,
which is a type of **Machine Learning**,
which is a type of **Artificial Intelligence**.



What do you know about each item ?

Generative Pre-trained Transformer (GPT)

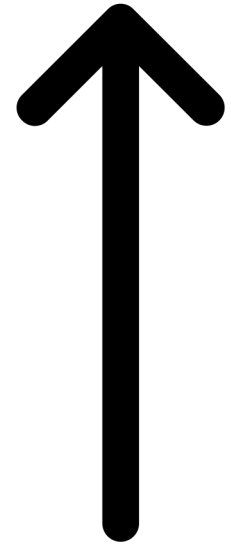
which is a type of **Large Language Model**,

which is a type of **General-purpose Transformer**,

which is a type of **Artificial Neural Network**,

which is a type of **Machine Learning**,

which is a type of **Artificial Intelligence**.



Transformers

Designed versus limitations of « classic » DL architectures (RNN, LSTM, CNN) :

- sequential processing
- difficulty with long-term dependencies
- less computational efficiency

'Attention is All You Need'.

Defining an architecture with 3 new aspects :

- Encoder-Decoder structure
- Positional encoding
- Self-Attention

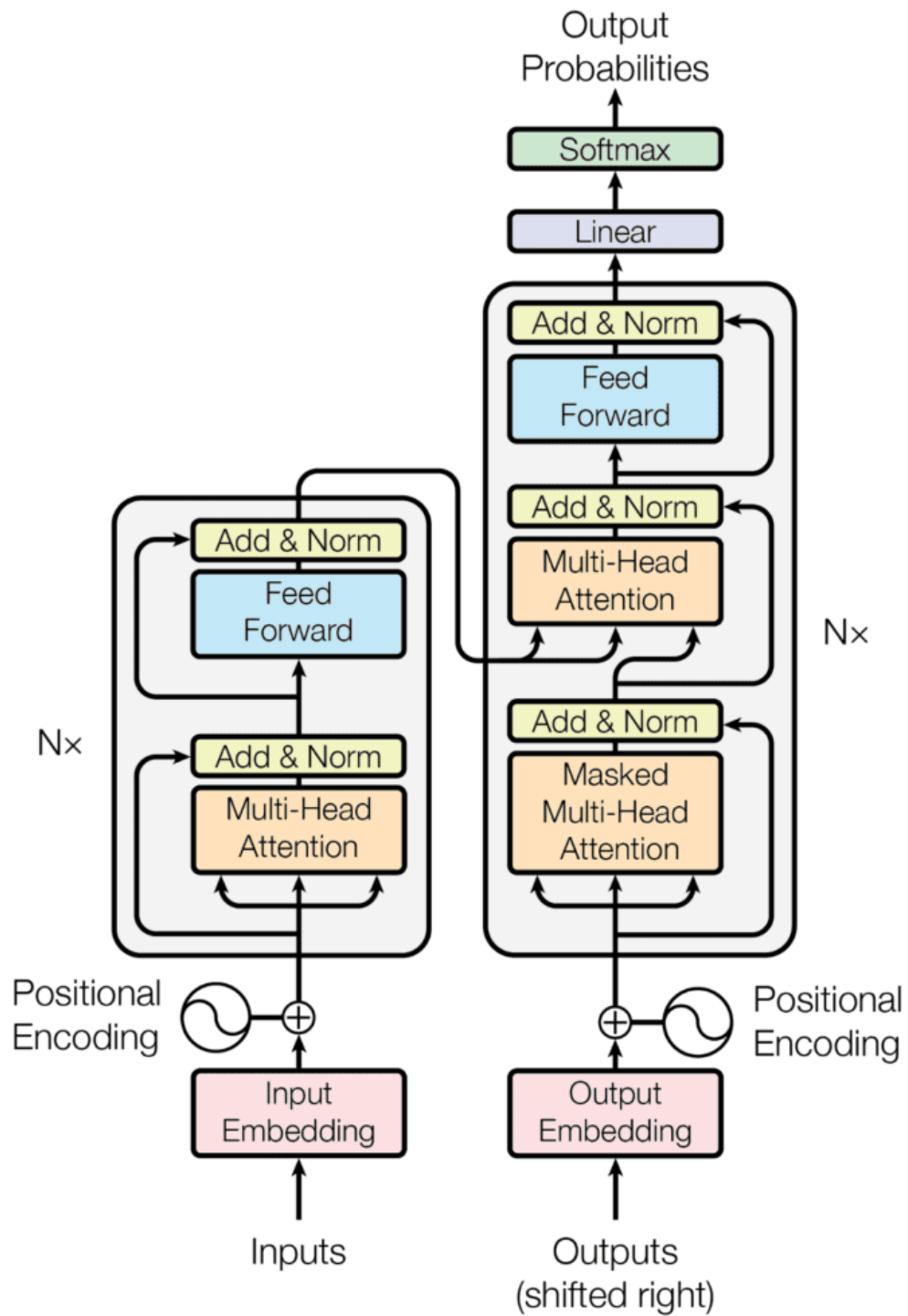
Parallel processing + Handling long-term dependencies
=
Scalability & Performance

Let's use them a bit...



Transformers : how does it work ?

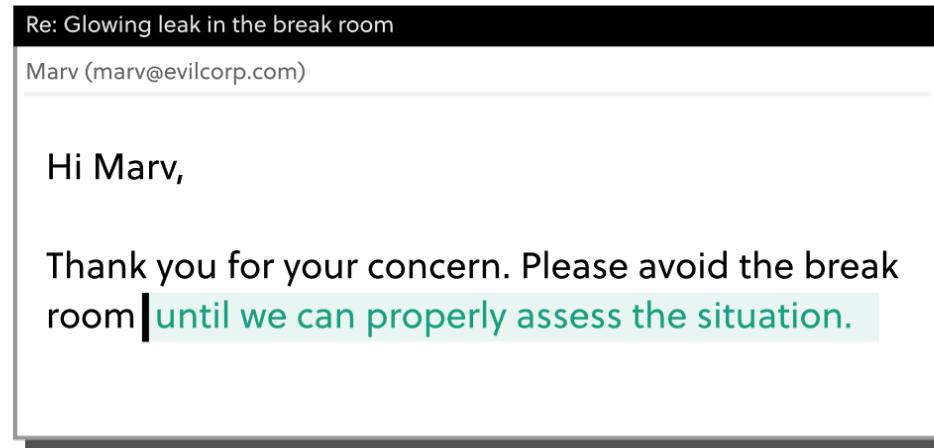
Transformers



Transformers



Transformers



Modern email programs try to predict the next word in a sentence. How would you guess they do this?

Predicting the next word

N-gram language models, while not the same as large language models (LLMs), are surprisingly simple and offer a good model to understand how this works.

Let's begin, Bigram models answer the question "given this first word, what is a likely next word?"

The last word you gave it provides the **context** for the next word.

The last word would be guessed thanks to the probability in the **corpus** that you used for the training.

Predicting the next word

If you gave the bigram model a long prompt like “To be or not to be, that is the question,” what is the bigram model using as context?

Predicting the next word

If we could use more than one word of context, would it improve the predictive power of our model?

Predicting the next word

You could make your context window, or the n in n -gram, as large as you want.

- 1) How large would the context window of an 8-gram model be?
- 2) Do you think there could ever be a problem with using a very large n ?

Predicting the next word

Which uses a smaller n , plagiarism detection or text autocorrect?

Calculating probabilities

Let's create "our model" from this only input :

"In the heart of the serene valley, in the golden light of dawn, the quick brown fox elegantly leaps over the indifferent, sleepy dog, while in the gentle caress of the morning breeze, whispers weave through the vibrant green leaves."

- 1) Uni-gram
- 2) Bi-gram

Calculating probabilities

Imagine a prompt that ends with the word “egg.” This table shows all the words that appear after “egg,” and how many times each bigram appears in the book.

Word	Count	Percentage (%)	<i>What percent of the time does “Bravo” appear after “egg”?</i> <ul style="list-style-type: none">• 100%• 23,43%• Not enough information to tell• Never ?
Alpha	52	13.10	
Bravo	93	23.43	
Charlie	15	3.78	
Delta	72	18.14	
Echo	61	15.37	
Foxtrot	21	5.29	
Golf	83	20.91	

Limits of n-gram

N-gram models are simple models that are good at predicting the next word, but not at more complicated tasks.

Consider this sentence:

I bought it for the cake she'll make with Julia tomorrow.

*What would a bigram model **not** be able to do?*

Limits of n-gram

N-gram models cannot connect pieces of information that are separated by a lot of words.

=> Neural networks

Needs **processing & tokenization**

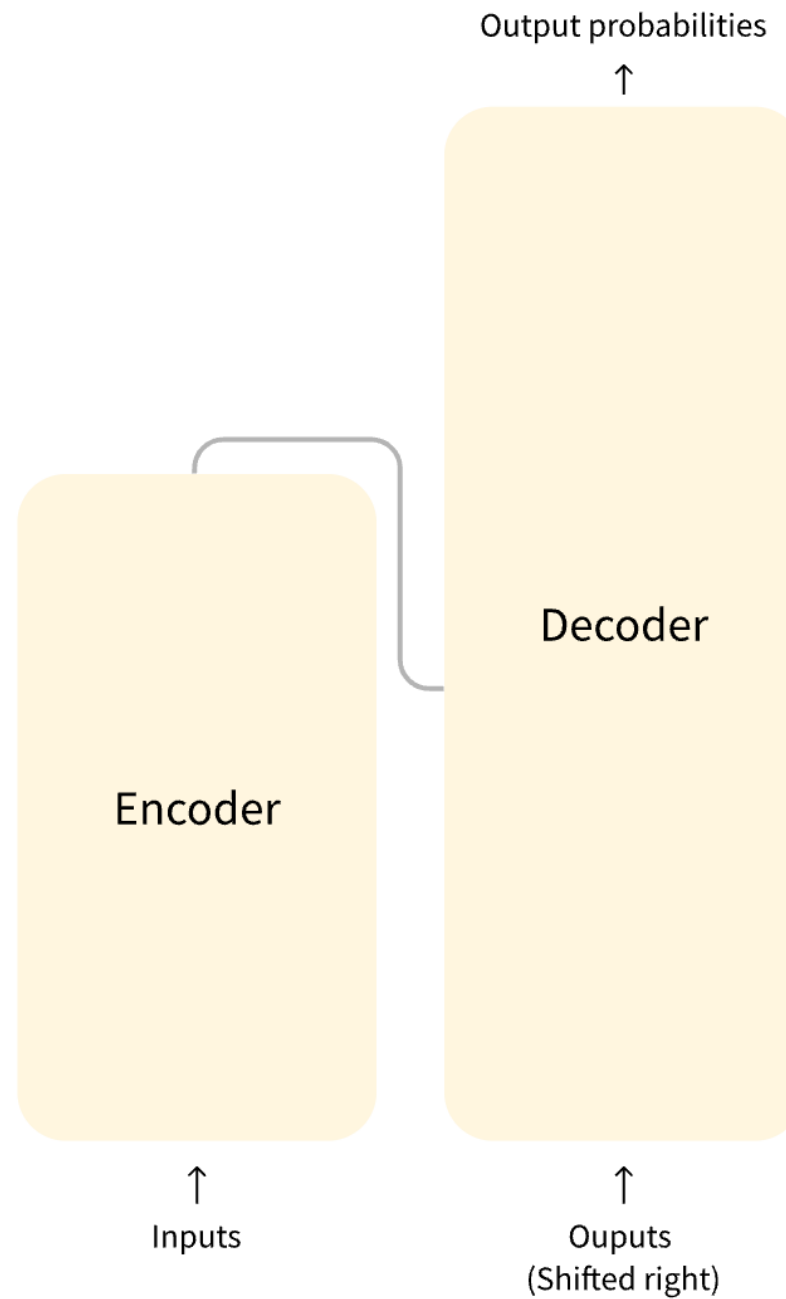
Limits of n-gram

N-gram models cannot connect pieces of information that are separated by a lot of words.

=> Neural networks
Needs **processing & tokenization**

⇒ Transformers
Adds encoding

Architecture



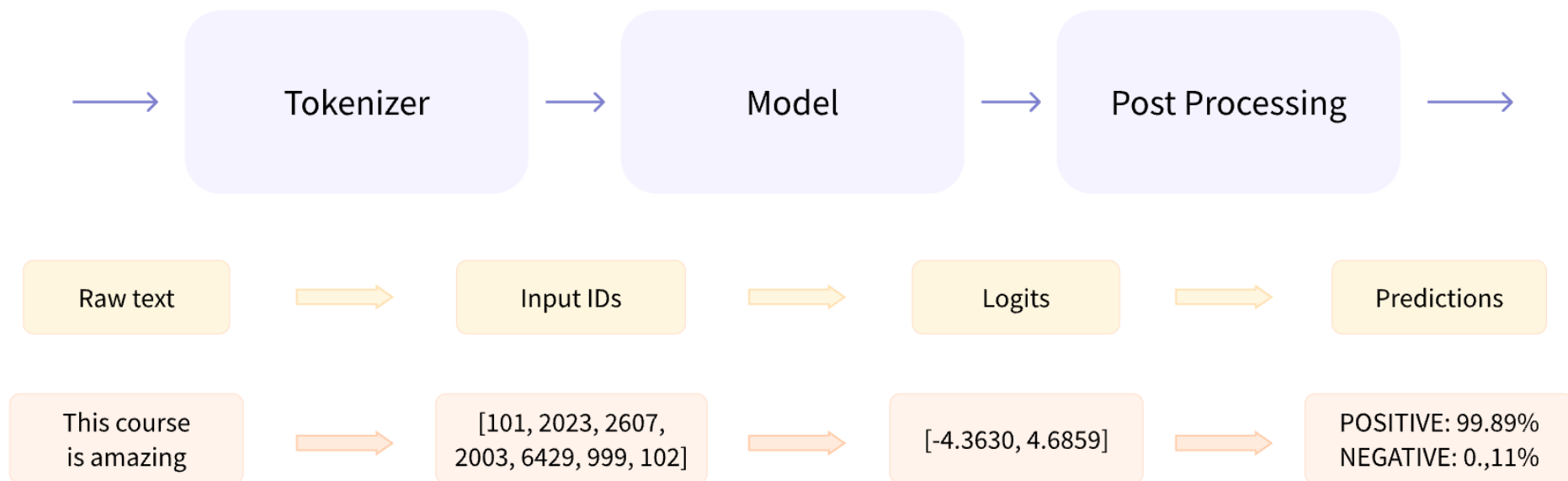
Transformers

1) The Self-Attention Mechanism

=> Example : “You like this course”.

2) Positional Encoding: Capturing Sequence Order

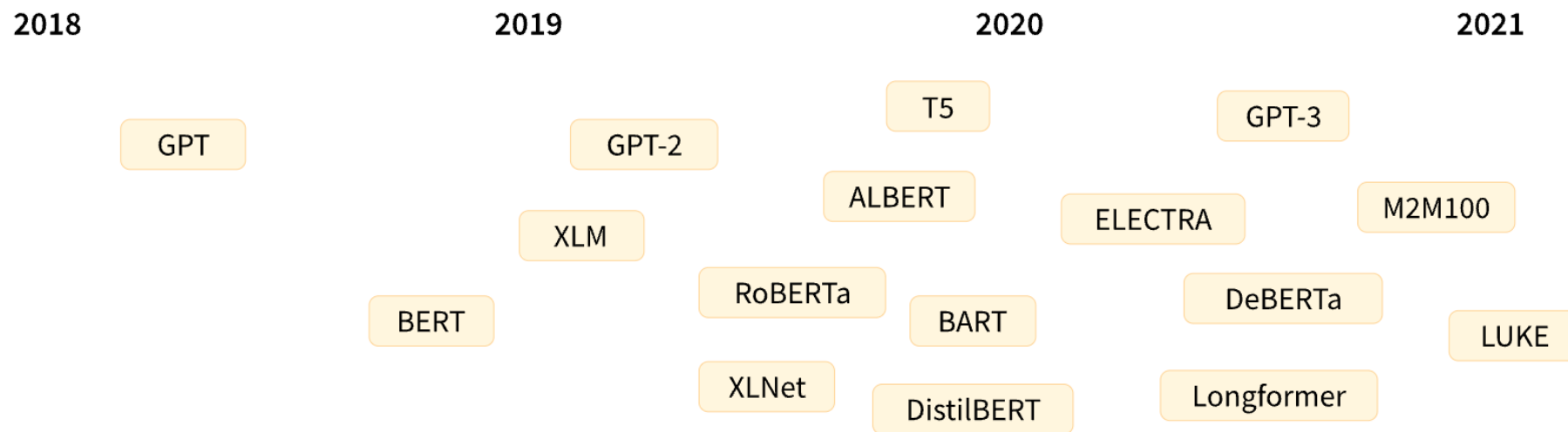
Transformers



Transformers

Second part of notebook : let's check this

Transformers



Model	Examples	Tasks
Encoder	ALBERT, BERT, DistilBERT, ELECTRA, RoBERTa	Sentence classification, named entity recognition, extractive question answering
Decoder	CTRL, GPT, GPT-2, Transformer XL	Text generation
Encoder-decoder	BART, T5, Marian, mBART	Summarization, translation, generative question answering

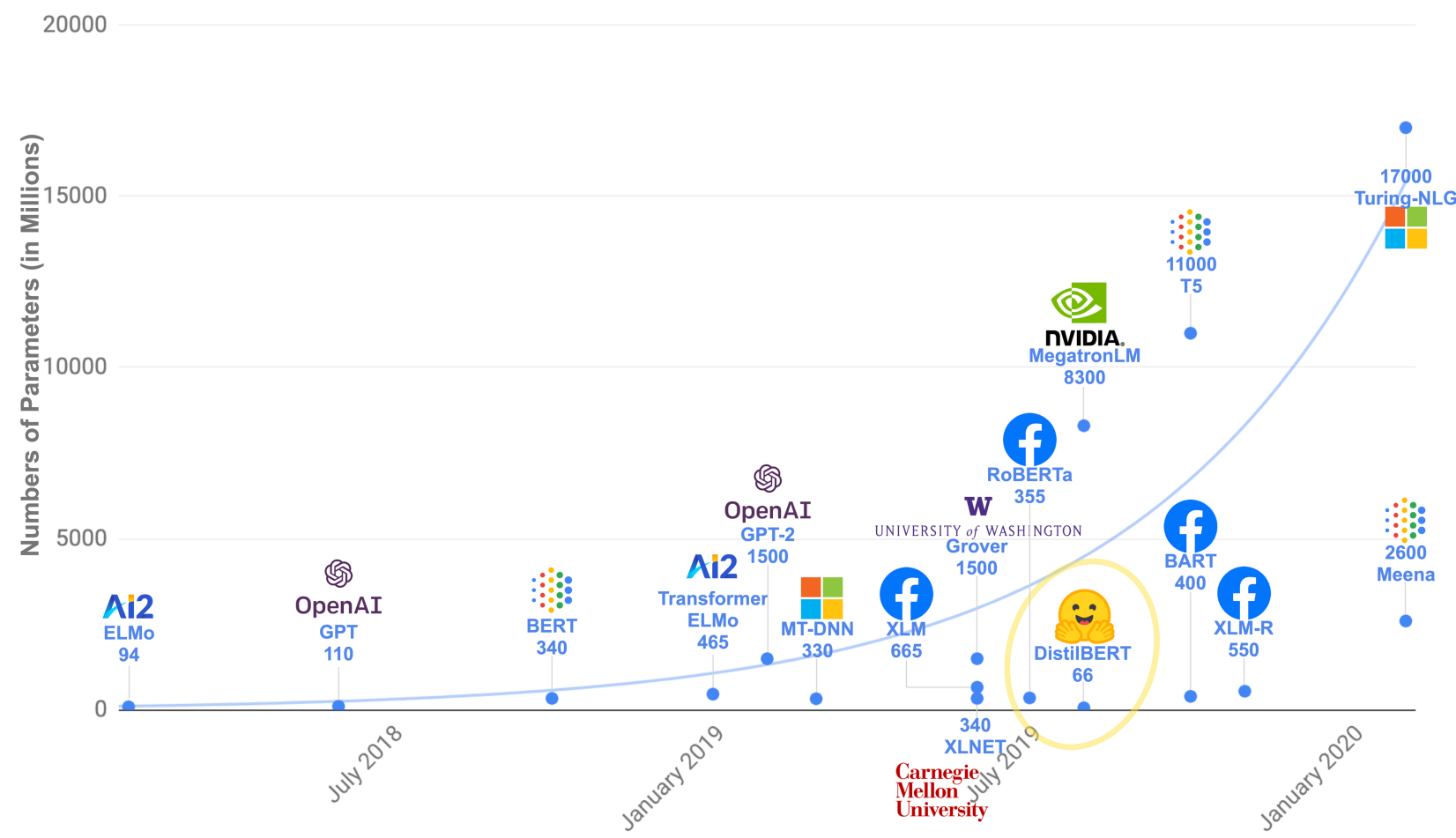
LLM

“Large” Language Models, what is large ?

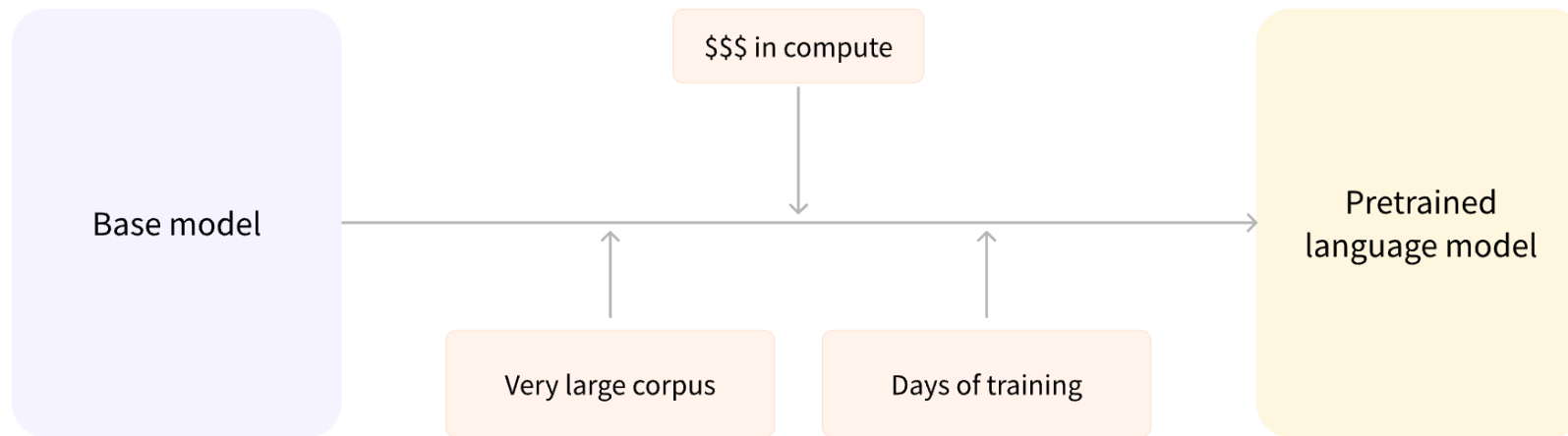
What's the use ?



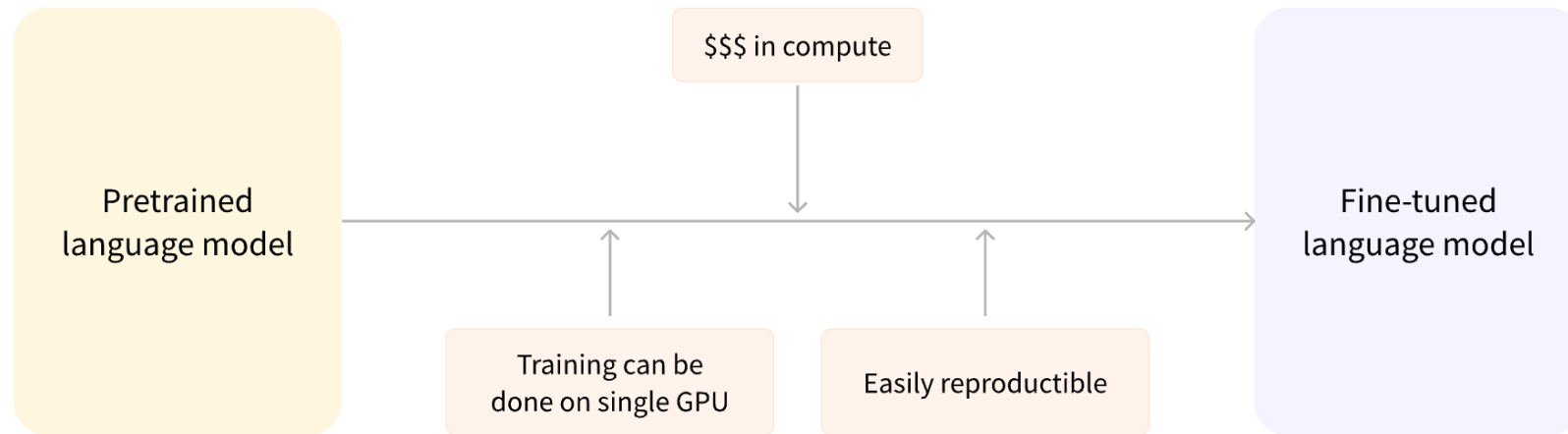
Transformers



Recap : This session



Recap : Next session



Thank you for your attention

If you have any questions,
feel free to get in touch:

simon.abad@ext.emlyon.com

