# LLM
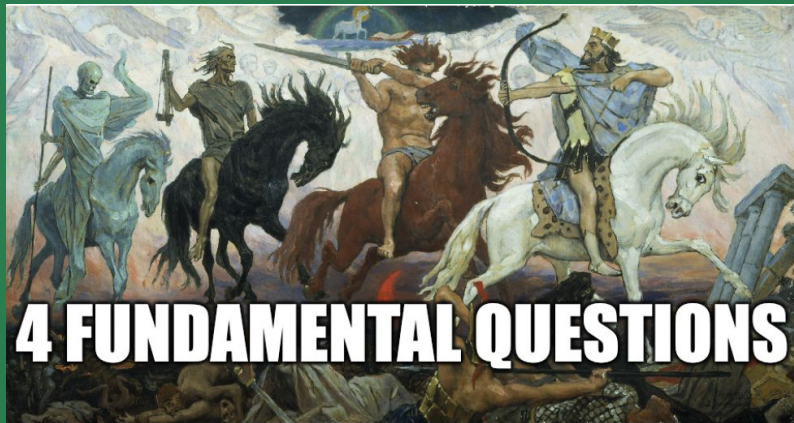# 2025

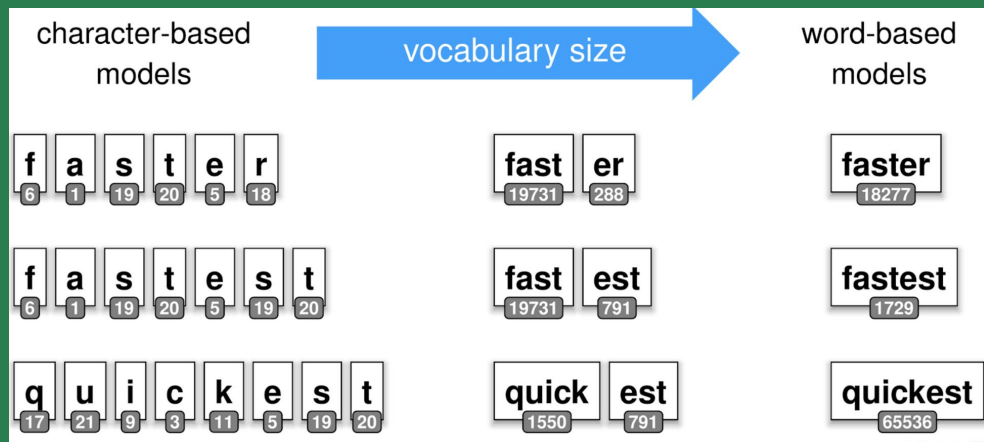# ML Basics

1. What is the input?

2. What is the output?

3. What is the model architecture?

4. What is the learning algorithm?



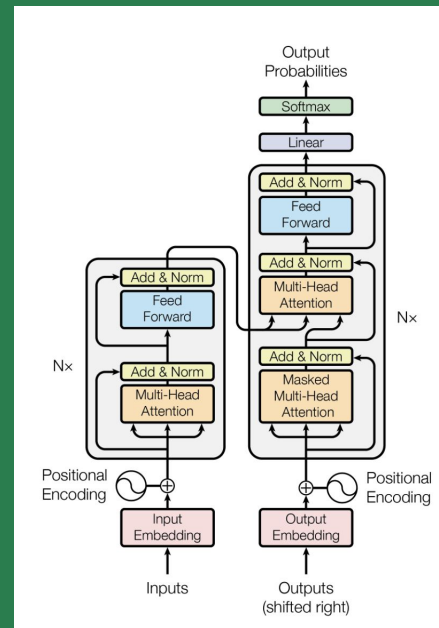4 FUNDAMENTAL QUESTIONS

# What is the input?

character-based models → vocabulary size → word-based models

| character-based | intermediate | word-based |
|---|---|---|
| f a s t e r | fast er | faster |
| 6 1 19 20 5 18 | 19731 288 | 18277 |
| f a s t e s t | fast est | fastest |
| 6 1 19 20 5 19 20 | 19731 791 | 1729 |
| q u i c k e s t | quick est | quickest |
| 17 21 9 3 11 5 19 20 | 1550 791 | 65536 |

https://tiktokenizer.vercel.app/

# What is the output?

**The Capital of France is**

| Token | Step score | Total score |
|-------|-----------|-------------|
| Paris | 19.9081 | 19.9081 |
| the | 16.6388 | 16.6388 |
| ... | 16.1794 | 16.1794 |
| known | 15.8383 | 15.8383 |

# What is the model arch algorithm?

# What is the learning algorithm?

# ML Basics
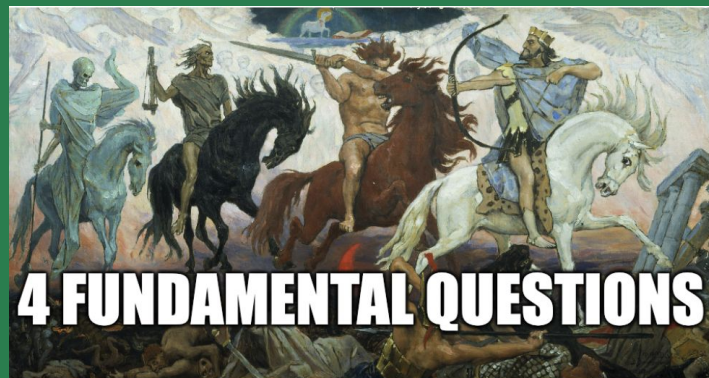
1.What is the input? -tokens

2.What is the output? -token probabilities

3.What is the model architecture? -transformers



4 FUNDAMENTAL QUESTIONS

4.What is the learning algorithm? -self-supervised learning
(Stochastic Gradient Descent - error back propagation)

# DEMO

Spelling/Counting

# LLM History

2018 - Transformer paper (Attention is all you need)

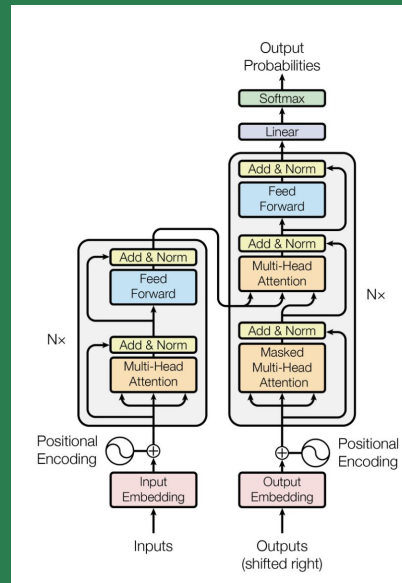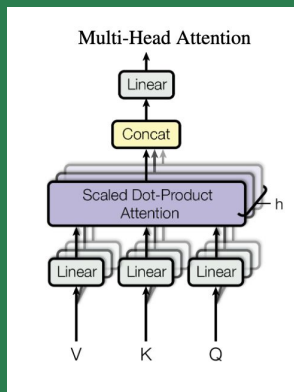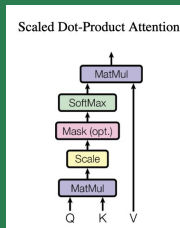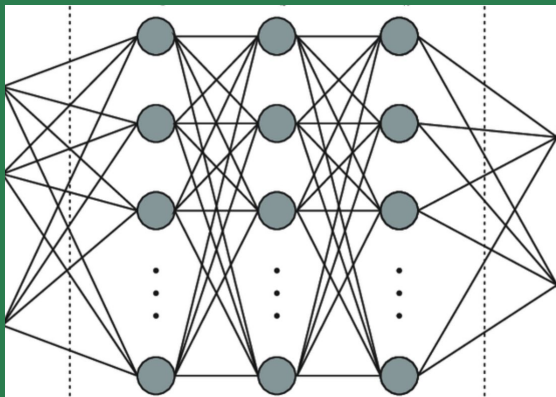2020 - GPT3 (Language Models are Few-Shot Learners)

2022 - ChatGPT (GPT 3.5 + RLHF)

2024 - Deepseek (R1-Zero & R1)

# LLM History

2018 - Transformer paper (Attention is all you need)

- basic language model architecture

- smarter / more efficient way to combine parameters

# LLM History

2020 - GPT3 (Language Models are Few-Shot Learners)

- 10X scale of previous models

- can generate human-like text - news articles

- performs language tasks without task-specific tuning (emergent SOTA capabilities)
    - Translation
    - Question-answering
    - Unscrambling words
    - 3-digit arithmetic.

```
Title:  Star's Tux Promise Draws Megyn Kelly's Sarcasm
Subtitle:  Joaquin Phoenix pledged to not change for each awards event
Article:  A year ago, Joaquin Phoenix made headlines when he appeared on
the red carpet at the Golden Globes wearing a tuxedo with a paper bag over
his head that read, "I am a shape-shifter.  I can't change the world.  I
can only change myself." It was a promise to not change to fit into the
Hollywood mold:  "I think that's a really special thing, to not change
```

**Zero-shot**

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1   Translate English to French:        ←── task description

2   cheese =>                           ←── prompt
```

**One-shot**

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1   Translate English to French:        ←── task description

2   sea otter => loutre de mer          ←── example

3   cheese =>                           ←── prompt
```

**Few-shot**

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1   Translate English to French:        ←── task description

2   sea otter => loutre de mer          ←── examples

3   peppermint => menthe poivrée        ←──

4   plush girafe => girafe peluche      ←──

5   cheese =>                           ←── prompt
```
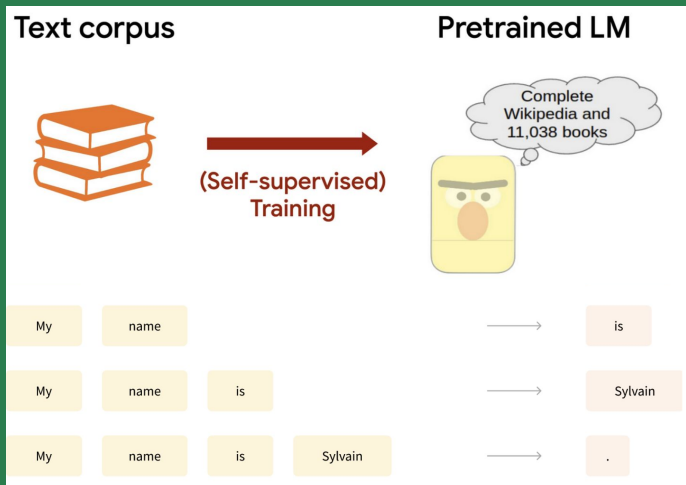
# DEMO

translate to hungarian
bear => medve
strawberry =>

Llama3.1-405b-base

# LLM History

2022 - ChatGPT (GPT 3.5 + RLHF)

- base models vs instruction following/assistant models

- raw text corpus vs special format conversation data

- human preference feedback



Text corpus — Pretrained LM

(Self-supervised) Training

Complete Wikipedia and 11,038 books

| My | name | | | → | is |
| My | name | is | | → | Sylvain |
| My | name | is | Sylvain | → | . |



<USER>

Can you help me with this code? It seems like there is a bug.
print("hello world)

<ASSISTANT>

It looks like you forgot to close the string passed to the function print. You have to add a closing quote to properly terminate the string. Here is the corrected function:

print("hello world")

Let me know if I can help with anything else!



You
Write a haiku about paperclips

ChatGPT
Steel curves intertwine,
Holding life's pages in line,
Quiet, small design.

You
Write a haiku about paperclips

ChatGPT
Silver arcs entwined,
Holding life's pages, aligned,
In steel we find, bind.
< 12 / 12 >

You
Write a haiku about paperclips

ChatGPT
Silver arcs entwined,
Humble joiner of pages,
Quiet strength, aligned.
< 2/2 >

# DEMO

what is 2 * 55 minus 13?

Llama3.1-405b-base

vs

Llama3.1-405b
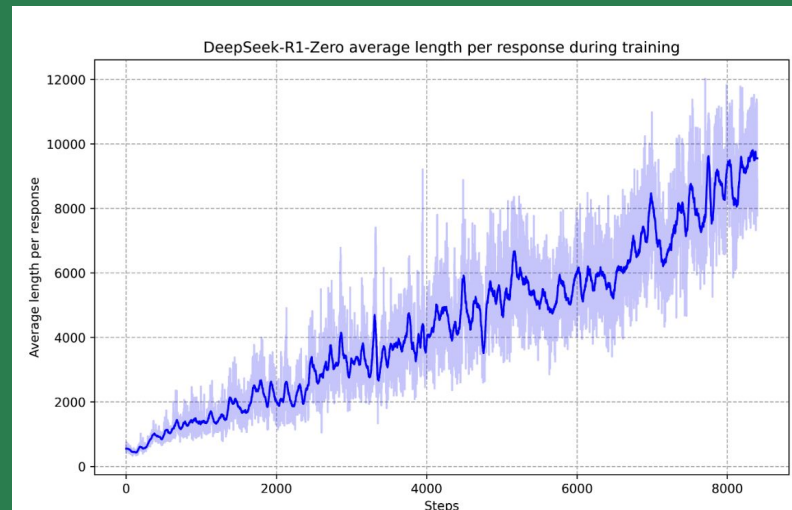(instruction trained model)
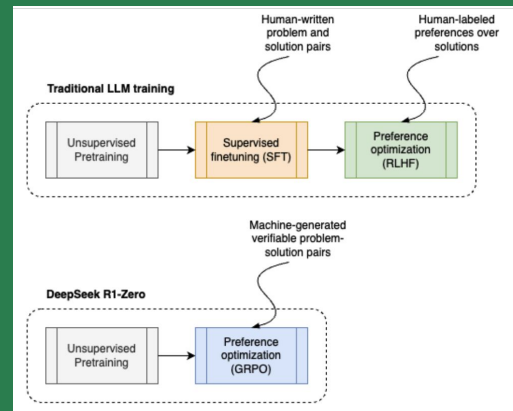
# LLM History

2024 - <u>Deepseek</u> (R1-Zero & R1)

    - think before answering - reasoning/thinking models

    - human written problem-solution pair data set
        - explicit teaching

  vs

    - LLM generated solution to verifiable problems
        - filter based on accuracy

        - naturally emerged, self-evolution
        - reevaluation, explore alternatives
        - solutions with optimal reasoning

    - longer generated response usually better,
more "time/tokens" for the model to "think/compute"

# DEMO

what is 2971 * 98?

Llama3.1-405b
(instruction trained model)

Vs

Deepseek-r1
(reasoning model)

# LLM History

2018 - Transformer paper (Attention is all you need)

    - basic language model architecture

2020 - GPT3 (Language Models are Few-Shot Learners)

    - human-like text and language task without specific tuning

2022 - ChatGPT (GPT 3.5 + RLHF)

    - from base models to instruction following/assistant models

2024 - Deepseek (R1-Zero & R1)

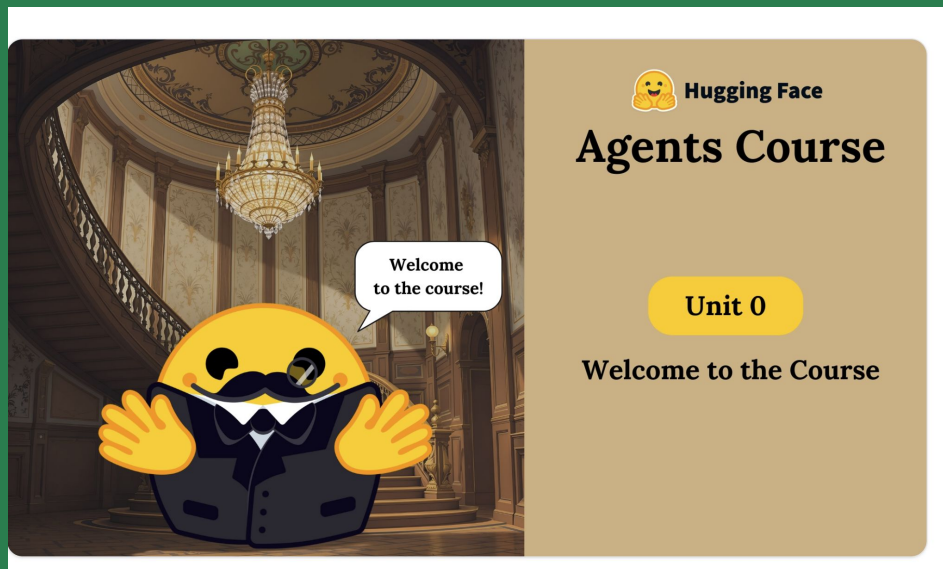    - true Reinforcement Learning - reasoning/thinking models

# HuggingFace abstraction

InferenceClient API

Pipelines API

Models API

# Agent



https://huggingface.co/learn/agents-course/unit0/introduction

# Agent

LLM + Tools

# Tool use

what is 2 * 55 minus 13?

Tool - a function with clear goal, description, inputs, outputs

Calculator tool - which can add numbers

# Tool use

```
system_message = """
You are a helpful assistant, that communicate using JSON format.
Your response action can either be a "Final Answer" or "ToolName" for a tool use.

Tools available to Assistant are:

- "Calculator": Useful for when you need to answer questions about math.
  - To use the calculator tool, Assistant should write like so:
    ```json
    {{"action": "Calculator",
     "action_input": "4+4"}}
    ```

Here are some previous conversations between the Assistant and User:

User: Hey how are you today?
Assistant: ```json
{{"action": "Final Answer",
 "action_input": "I'm good thanks, how are you?"}}
```
User: I'm great, what is the square root of 4?
Assistant: ```json
{{"action": "Calculator",
 "action_input": "sqrt(4)"}}
```
Result: 2.0
Assistant: ```json
{{"action": "Final Answer",
 "action_input": "It looks like the answer is 2!"}}
"""
```

what is 2 * 55 minus 13?

```python
def use_calculator(input, first_call=True):
  if first_call:
    input_prompt = prompt.format(user_input=input)
  else:
    input_prompt = input
  # make the agent call
  response = agent(input_prompt)
  # parse the output if possible
  parsed = parser(response)
  # check if the output is our final answer or if it is a tool
  if parsed["action"] == "Final Answer":
    return parsed["action_input"]
  # if not try to use the tool
  tool_output = use_tool(parsed["action"], parsed["action_input"])

  # add message to the agent
  next_prompt = f"{input_prompt}\n{response}\nResponse: {tool_output}"
  # recursively call the agent with the output of the tool
  return use_calculator(next_prompt, False)
```

# Agent

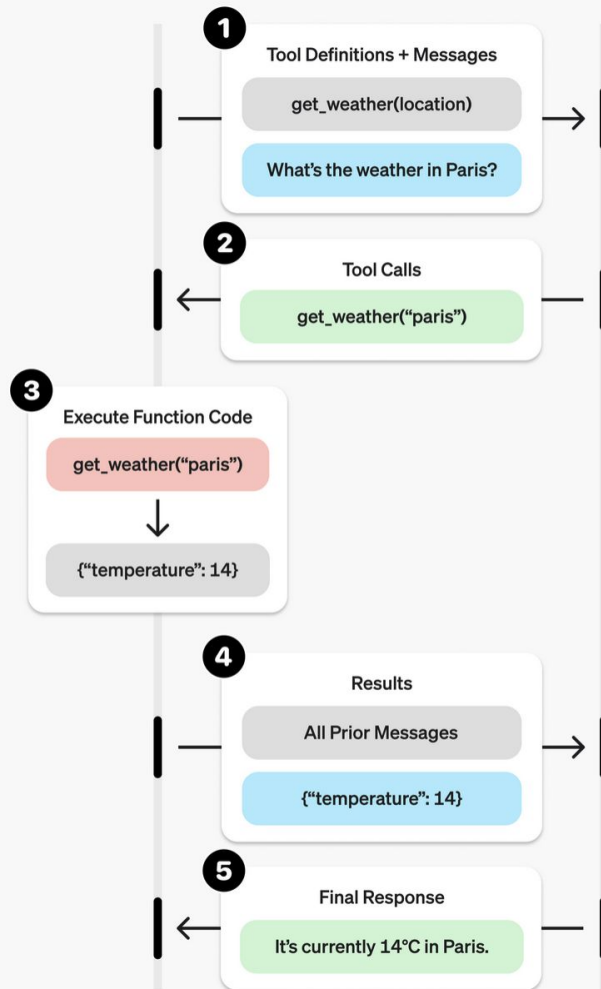LLM + tool definitions + tool executor

# LLM History

2023 - Tool use / Function calling

- connect LLM to external tools

- train model with available list of tools
(description, input format)

- model can automatically decide
- to generate output normally
- or to use a tool - and generate tool input

- tool invocation is done "outside" of the llm

# from powerful models
# to artfully curated context

https://thelongcontext.com/

# making learning engaging

notebooklm.google.com

# References

HuggingFace Agent course -
https://huggingface.co/learn/agents-course/unit0/introduction

deeplearning.ai Agent course -
https://www.deeplearning.ai/short-courses/building-code-agents-with-hugging-face-smolagents/

Creating interactive game using a book and long context LLM -
https://thelongcontext.com/

Learn, by converting files into a podcast -
notebooklm.google.com