





## Intro to LLMs in Healthcare

Assistant Professor Jyrki Savolainen LUT-University / CSC – IT Center for Science

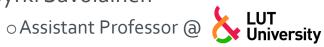
23<sup>rd</sup> Oct 2025





#### About me

- Jyrki Savolainen



○ Application Specialist @



o DSc. (Econ.), MSc. (Eng.)

- Research and Teaching
  - Simulation, Business Data Analytics

ojyrki.savolainen@lut.fi





#### **Contents**

• Large Language Models (LLMs) and Foundational Models

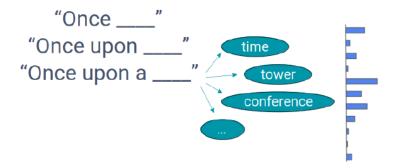
• Tailoring LLMs



# Large Language Models are... models



- Statistical, autoregressive model of text
- Generate sequence of most likely tokens in context



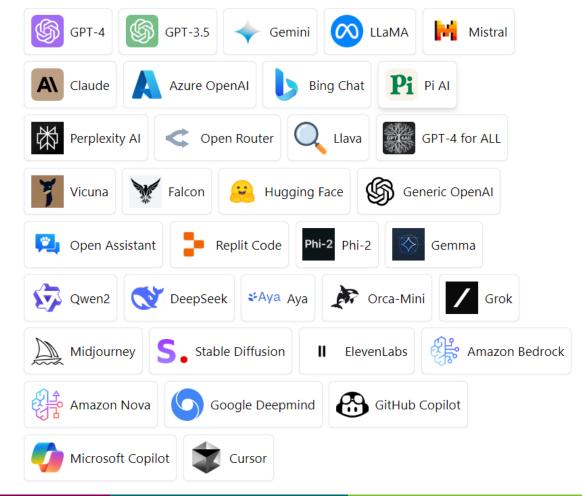
 $P(\text{token}_n \mid \text{token}_1, \text{token}_2, ..., \text{token}_{n-1})$ 

# Large Language Models (LLMs)

• Trained with massive datasets

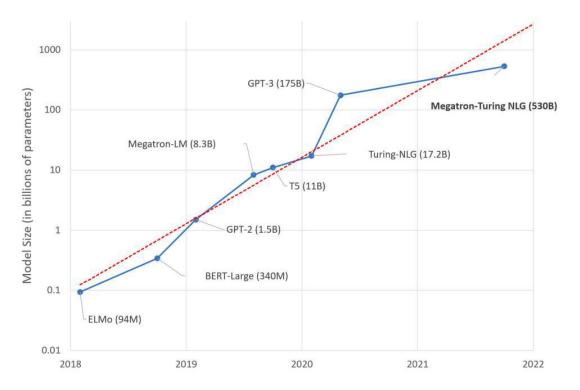
- Available for download
  - Require a lot of resources (dependent on model size)
  - o E.g. huggingface.co







### Increasing Size and Complexity of the Foundational Models



- Model size measured in billions of parameters (tokens)
- In general: "Higher number of tokens leads to better reasoning capabilities" with the cost of computational requirements

## **Computational Costs?**

#### Microsoft chooses infamous nuclear site for AI power

20 September 2024

Share Save

#### Natalie Sherman

**BBC News** 



America's Three Mile Island energy plant, the site of the worst nuclear accident in US history, is preparing to reopen as Microsoft looks for ways to satisfy its growing energy needs.

The tech giant said it had signed a 20-year deal to purchase power from the Pennsylvania plant, which would reopen in 2028 after improvements.

The agreement is intended to provide the company with a clean source of energy as power-hungry data centres for artificial intelligence (AI) expand.

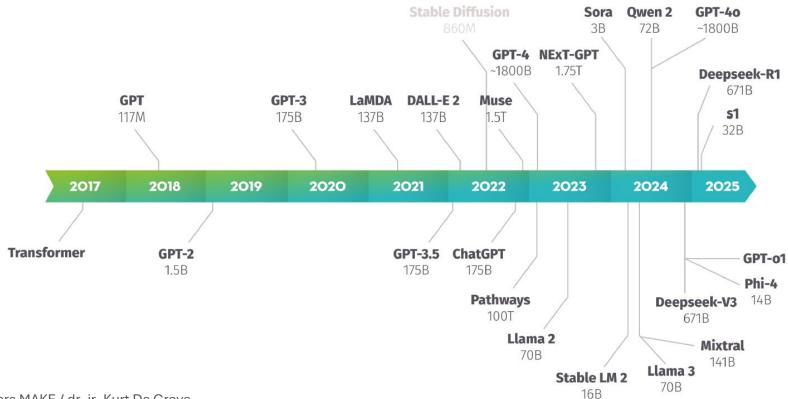
The plan will now go to regulators for approval.

.

Sisäinen

| TRAINING COST | training TIME  | GPUs   | parameters   |
|---------------|--|--|--|
| \$1 million   | 2 days   | 512  | 1.5B   |
| \$10 million  | 21 days  | 2,000  | 175B   |
| \$200 million | 3 months   | 8,000  | 1,800B   |
| \$1 billion?  | 5 months   | 25,000   | 10T<br>(10,000B)   |
| \$? billion?  | ? year ? months  | 100,000  | ?  |
|               |  | /llm-trainin   | g-cost-how   |
|               | \$1 million<br>\$10 million<br>\$200 million<br>\$1 billion?<br>\$? billion? | \$1 million 2 days \$10 million 21 days  \$200 million 3 months  \$1 billion? 5 months  \$? billion? ? year ? months | \$1 million 2 days 512<br>\$10 million 21 days 2,000<br>\$200 million 3 months 8,000<br>\$1 billion? 5 months 25,000<br>\$? billion? ? year ? months 100,000<br>https://gregoreite.com/llm-trainin |

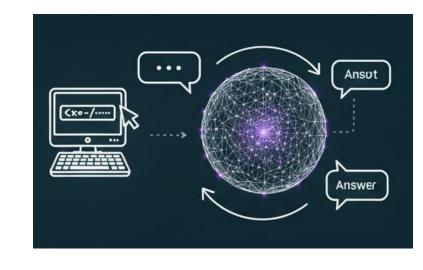
#### **Milestones**





# THIS WEBINAR: Tailoring LLMs for a specific purpose

- For simple tasks, simple models are enough
  - e.g. document summarization or finding information from multiple documents
- There is the option of <u>prompting the</u> <u>existing models</u>
- If you want your own model, then simple,
   "small scale", LLMs can run on local servers
   or are affordable for daily use through APIs
   Case example coming up





## **Options for Tailoring LLMs**

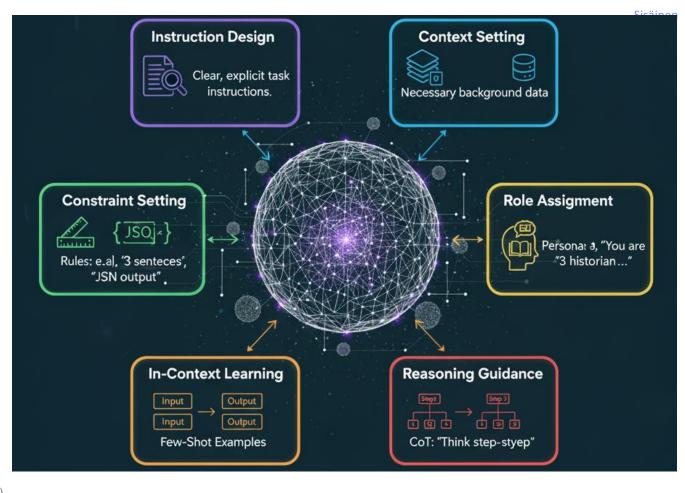
| Method                | Resource Intensity         | Changes<br>Model<br>Parameters | Best For   |
|-----------------------|----------------------------|--------------------------------|--|
| 1. Prompt Engineering | Lowest (Inference<br>Cost) | No                             | Changing tone, task format, simple few-shot tasks. |

#### **EXAMPLE:**

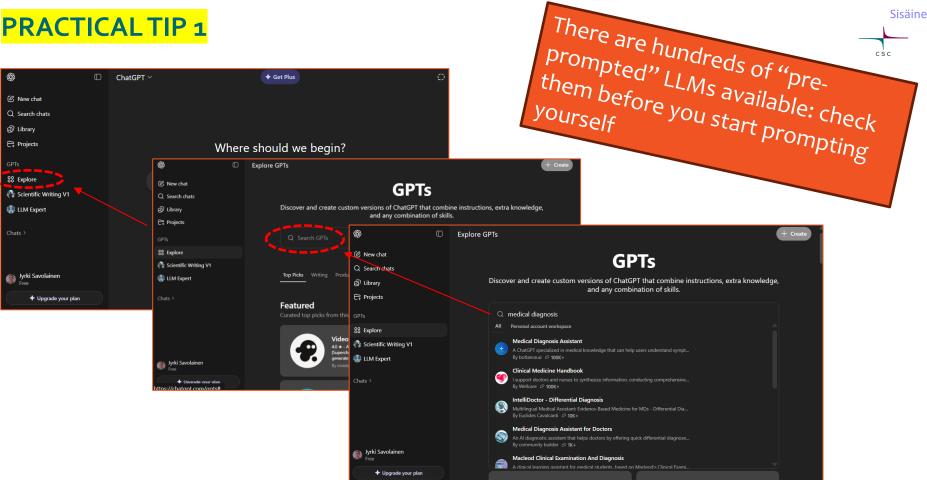
"You are an expert financial analyst. Your goal is to provide concise, factual summaries of stock performance. Never provide investment advice, only historical data and public company information. Respond in a formal, professional tone."

LUT Group Confidential - Other information (3)

# **Prompt Engineering**



#### PRACTICAL TIP 1



Sisäinen

# "LLM Agents"

- Pre-trained LLMs from which you can query for specific topics
- Can be trained based on vast data focused on certain area





| Method                                     | Resource Intensity                                  | Changes<br>Model<br>Parameters | Best For  |
|--|---|--------------------------------|---|
| 1. Prompt Engineering                      | Lowest (Inference<br>Cost)                          | No                             | Changing tone, task format, simple few-shot tasks.                |
| 2. Retrieval-Augmented<br>Generation (RAG) | Low to Moderate<br>(System Design &<br>Maintenance) | No                             | Incorporating new/proprietary knowledge, reducing hallucinations. |

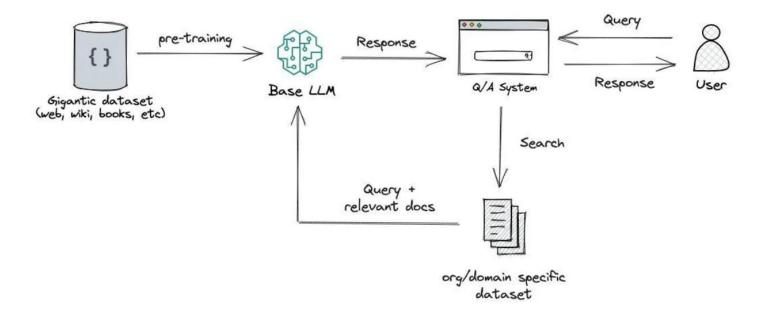
#### **EXAMPLE:**

Organizational instructions, code of conduct, strategy, etc. are used as a training material for the LLM(s)

LUT Group Confidential - Other information (3Y)

# Retrieval-Augmented Generation (RAG)





#### PRACTICAL TIP 2





2.5 Flash ▼

Google Gemini AUTOMATICALLY connects with your Google Drive Documents (the same with MS Copilot and OneDrive) I have written a paper on the Evolutionary Spare Parts. Check







Workspace & 1 more ∨

The paper "Evolutionary Spare Parts" presents and formalizes a self-evolutionary design framework aimed at automatically optimizing the physical shape of a part for a unique, userspecific application. The core proposition is to create a continuous feedback loop between the physical product and its virtual counterpart, the Digital Twin (DT). In this system, automatically gathered process data from the application environment is used by Evolutionary Algorithms (EAs) to iteratively refine the part's design topology in the virtual space. By leveraging the design freedom of Additive Manufacturing (AM), this automated process allows for minor, valuable product improvements in an economically feasible way, resulting in improved application performance with every new spare part instance installed.

the latest version and summarize the plot in two paragraphs

This self-evolutionary logic unlocks a disruptive business model with solid potential in high-end industrial applications where the customer is motivated to pay extra for performance. The



| Method  | Resource Intensity                                  | Changes<br>Model<br>Parameters | Best For   |
|---|---|--------------------------------|--|
| 1. Prompt Engineering   | Lowest (Inference<br>Cost)                          | No                             | Changing tone, task format, simple few-shot tasks.   |
| 2. Retrieval-Augmented<br>Generation (RAG)  | Low to Moderate<br>(System Design &<br>Maintenance) | No                             | Incorporating new/proprietary knowledge, reducing hallucinations.                                  |
| 3. Fine-Tuning  | Moderate to High<br>(Training/Compute<br>Cost)      | Yes                            | Customizing style/tone, adapting to specific task formats, improving performance on a narrow task. |
| 4. Reinforcement Learning with<br>Human Feedback (RLHF) / Direct<br>Preference Optimization (DPO) | Collection/Training                                 | Yes                            | Aligning model behavior with human values, safety, or complex subjective preferences.              |

LUT Group Confidential - Other information (3Y)





facebook.com/CSCfi



twitter.com/CSCfi



youtube.com/CSCfi



linkedin.com/company/csc---it-center-for-science



github.com/CSCfi