# LLM Models

## What are LLMs? How do they work?

*Intuitive understanding for busy people*

**LLM Model**

Pre-Trained Model

**LLM Model**

Instruction Tuned Model

Generative AI is truly **revolutionary** technology. It is transforming the way we interact with technology. We are in a middle of a paradigm shift where for the first-time computers can understand humans via natural language and respond intelligently.
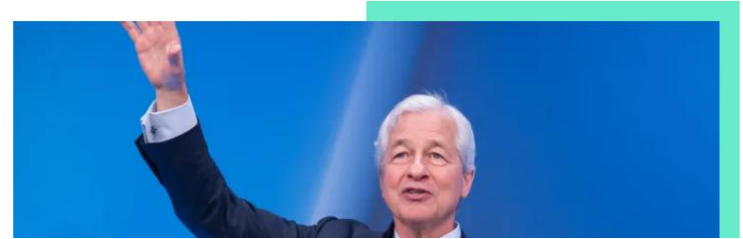
**NEXT GEN INVESTING**

## Jamie Dimon says AI could be as transformative as electricity or the internet—here's how to invest

Published Tue, Apr 9 2024·8:00 AM EDT

Cheyenne DeVon

SHARE

Source: CNBC

TECHNOLOGY | ARTIFICIAL INTELLIGENCE

## Amazon CEO Touts AI Revolution While Committing to Cost Cuts

In his letter to shareholders, Andy Jassy says generative AI could usher in the largest tech transformation since the Internet

By *Steven Russolillo* Follow *and Sebastian Herrera* Follow
Updated April 11, 2024 10:08 am ET

Resize          77          Gift unlocked article          Listen (6 min)
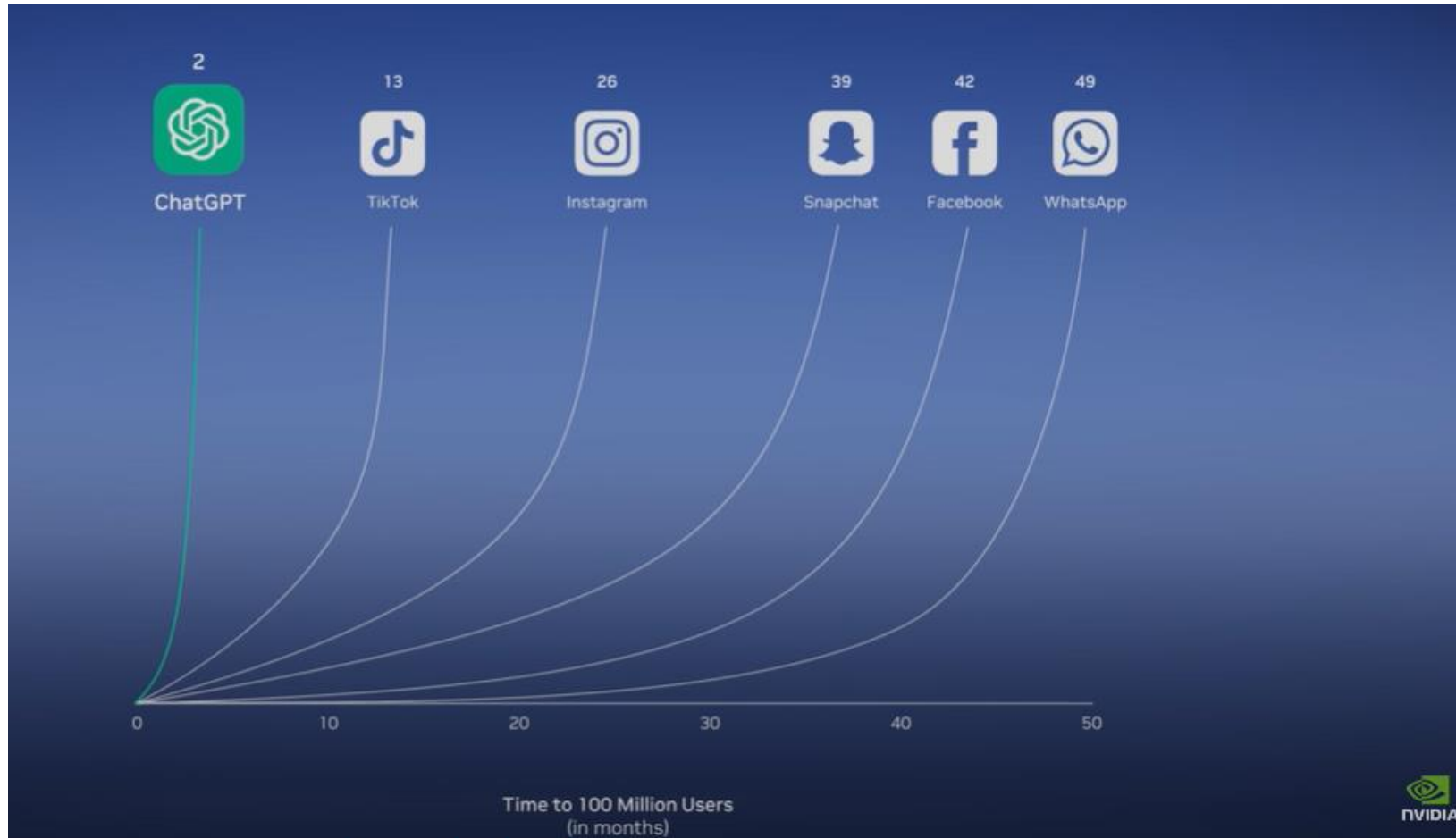
Source: WSJ

For the first time, we have a universal UI (User Interface). LLMs, can understand understand human natural language and can respond intelligently using natural language.

**User Interface**
Natural human language as input

Input (prompt) →

← output

**LLM Model**

# ChatGPT is the fastest growing application in human history.
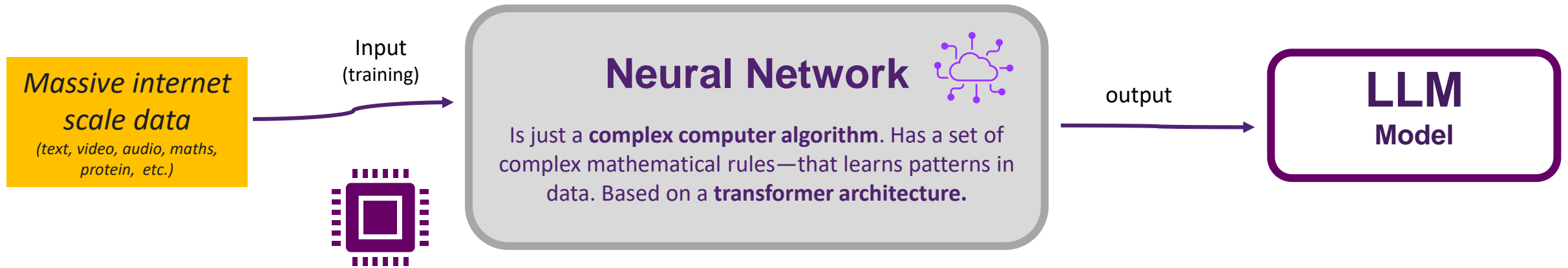That is because we use human natural language to interact with it.



Source: Nvidia

# Large Language Models

Given vast amount of data+compute, an algorithm can program itself to develop a deep understanding of patterns and meaning in the data.  This discipline is called **deep learning**. Once trained, LLMs can then generate this understanding or intelligence when **prompted** using natural language.

**Massive internet scale data**
*(text, video, audio, maths, protein,  etc.)*

Input
(training)

## Neural Network

Is just a **complex computer algorithm**. Has a set of complex mathematical rules—that learns patterns in data. Based on a **transformer architecture.**

output

## LLM
**Model**

**LLM**
**Model**

- Is a complex computer algorithm. Generally called a "neural network" within the technical community
- This neural network has an architecture called "transformer architecture" : or a set of defined complex mathematical rules—that has capability to learn patterns in data.
- Is a collection of few files. You can even download these files to your PC. The number & type vary based on framework (like TensorFlow or PyTorch). These include
  - Parameter files
  - Configuration/Setup files
  - Runtime files
- Has a "vocabulary size", which refers to the total number of unique tokens (words, characters, or subwords) that the model recognizes and uses to represent and process text
- It has a number of Layers. Layers can be though of steps in the process of transforming input into output.
- The context window of a LLM refers to the maximum amount of input text (in terms of tokens) that can be sent to the model when generating a response
- The parameters of a LLM can be thought of as variables. The parameter size of a LLM refers to the total number of learnable variables (weights and biases) within the model.  A larger parameter size generally means the model can capture more complex patterns and nuances in language, making it more powerful but also requiring more computational resources. For example, GPT-3 has 175 billion parameters, enabling it to generate highly sophisticated and human-like text.
- The process of invoking a LLM is called "inferencing".

Large language model

# Llama 2: open source, free for research and commercial use

We're unlocking the power of these large language models. Our latest version of Llama – Llama 2 – is now accessible to individuals, creators, researchers, and businesses so they can experiment, innovate, and scale their ideas responsibly.

**Download the model**

Llama 2 was trained on **40% more data** than Llama 1, and has double the context length.

## Llama 2

| MODEL SIZE (PARAMETERS) | PRETRAINED | FINE-TUNED FOR CHAT USE CASES |
|---|---|---|
| 7B | Model architecture: | Data collection for helpfulness and safety: |
| 13B | Pretraining Tokens: 2 Trillion | Supervised fine-tuning: Over 100,000 |
| 70B | Context Length: 4096 | Human Preferences: Over 1,000,000 |

With each model download you'll receive:

- Model code
- Model weights
- README (user guide)
- Responsible use guide
- License
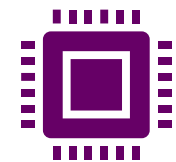- Acceptable use policy
- Model card

# What makes LLMs special

Large language models like GPT-4 or Llama 3 have state-of-the-art capabilities such as general **knowledge**, **steerability**, **advanced reasoning**, **math/science**, **tool use**, **data analysis**, **multilingual translation** and more.

Based on transformer architecture LLM models are giants and can learn to understand human knowledge without supervision & without labelled datasets.

A single LLM model can perform multiple tasks such as QA, summarization, content/code generation, data analysis, translation and more
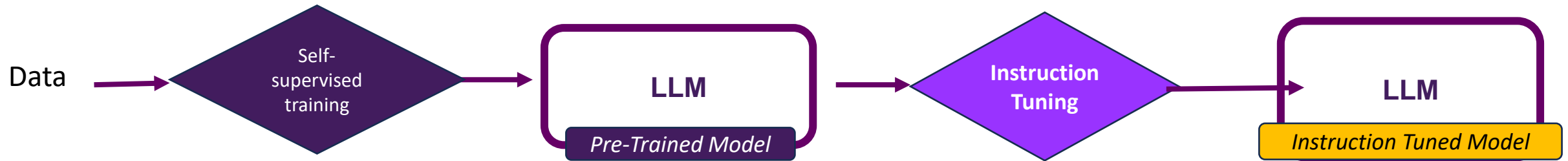
Models can be tuned to perform tasks for which they were never trained on.

LLM models can learn/understand patterns and representation of any sequence be it language, protein, biology, chemistry, etc.

LLMs are excellent few-shot learners. Using prompt engineering you can guide them to your request. LLMs can be multi-modal and so can be used in endless possible applications

# How are LLMs trained?

LLMs are very large <u>deep learning</u> models trained on huge amount of data. LLMs have a broad understanding of language, context, and world knowledge.



Both pre-trained and instruction-tuned models are foundation models. Because they are both built on a broad base of knowledge and are adaptable to a wide range of applications. The main difference is in the additional layer of training for instruction-tuned models, which is designed to enhance their ability to follow explicit instructions and perform tasks across different domains.
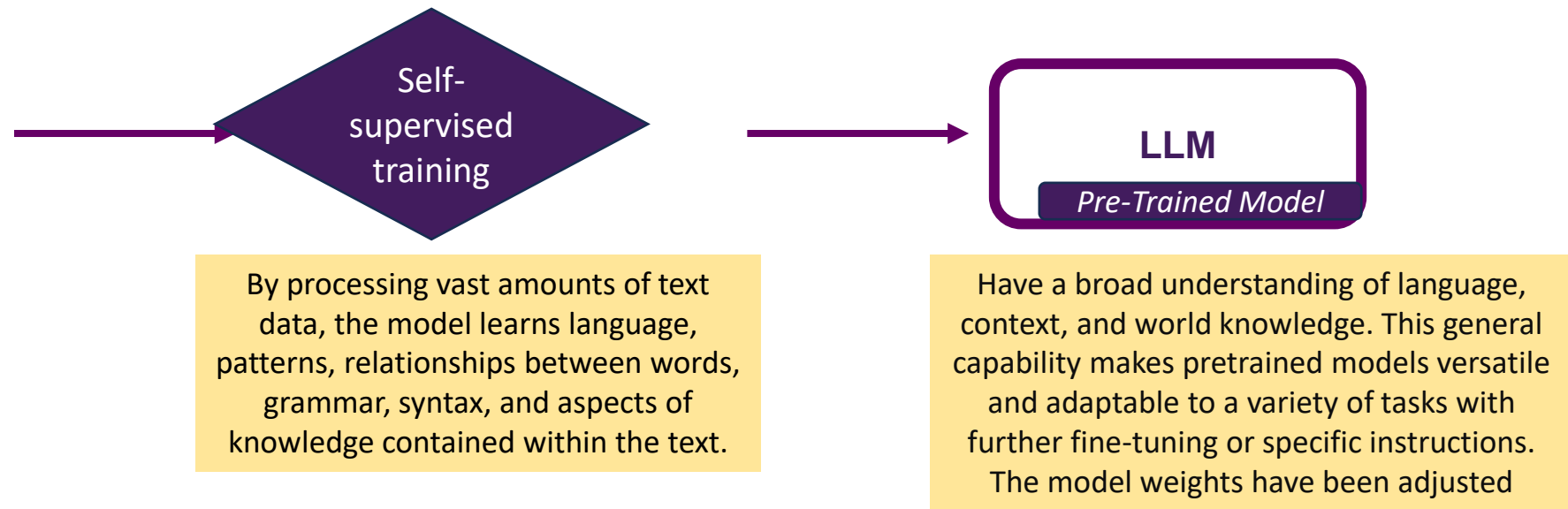
**LLMs are trained on massive corpus of internet data using GPUs**

| Dataset | Sampling prop. | Epochs | Disk size |
|---|---|---|---|
| CommonCrawl | 67.0% | 1.10 | 3.3 TB |
| C4 | 15.0% | 1.06 | 783 GB |
| Github | 4.5% | 0.64 | 328 GB |
| Wikipedia | 4.5% | 2.45 | 83 GB |
| Books | 4.5% | 2.23 | 85 GB |
| ArXiv | 2.5% | 1.06 | 92 GB |
| StackExchange | 2.0% | 1.03 | 78 GB |

Table 1: **Pre-training data.** Data mixtures used for pre-training, for each subset we list the sampling proportion, number of epochs performed on the subset when training on 1.4T tokens, and disk size. The pre-training runs on 1T tokens have the same sampling proportion.

Source (Paper on arxiv.org):
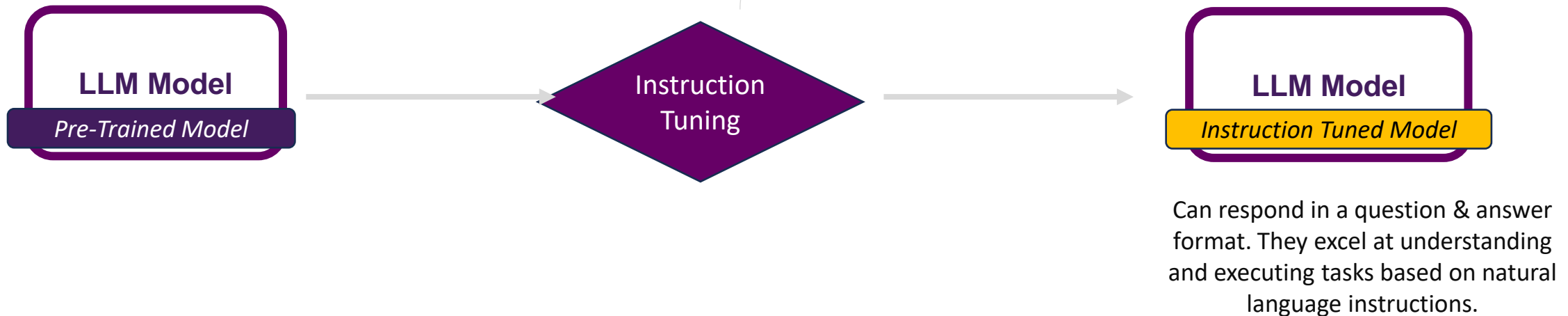LLaMA: Open and Efficient Foundation Language Models

**Self-supervised training**

By processing vast amounts of text data, the model learns language, patterns, relationships between words, grammar, syntax, and aspects of knowledge contained within the text.

**LLM**

*Pre-Trained Model*

Have a broad understanding of language, context, and world knowledge. This general capability makes pretrained models versatile and adaptable to a variety of tasks with further fine-tuning or specific instructions. The model weights have been adjusted

Supervised fine-tuning: Model is trained on low quantity/high quality labelled data such as Ideal Question/Response with human assistance

(+)

RLHF (Reinforcement Learning from Human Feedback)
Humans rank different responses generated by the model. This rating is captured in another model called the "reward model". Then the LLM model is trained with the help of the "reward model" to generate responses

**LLM Model**
*Pre-Trained Model*

Instruction Tuning

**LLM Model**
*Instruction Tuned Model*

Can respond in a question & answer format. They excel at understanding and executing tasks based on natural language instructions.

## LLM Model
### Pre-Trained Model

- **Definition:** These are models that have been initially trained on a large dataset to learn a wide range of patterns, knowledge, and language from that data. The process usually involves unsupervised learning, where the model learns to predict parts of the input (like the next word in a sentence) without explicit human-labeled instructions.
- **Purpose:** The main aim is to capture a broad understanding of language, context, and world knowledge. This general capability makes pretrained models versatile and adaptable to a variety of tasks with further fine-tuning or specific instructions.

## LLM Model
### Instruction Tuned Model

- **Definition:** These models start as pretrained models but undergo an additional phase of training (called instruction tuning or instruct-tuning) where they learn to follow human-like instructions or prompts more effectively. This stage involves supervised learning, typically using datasets where inputs are paired with instructions and desired outputs.
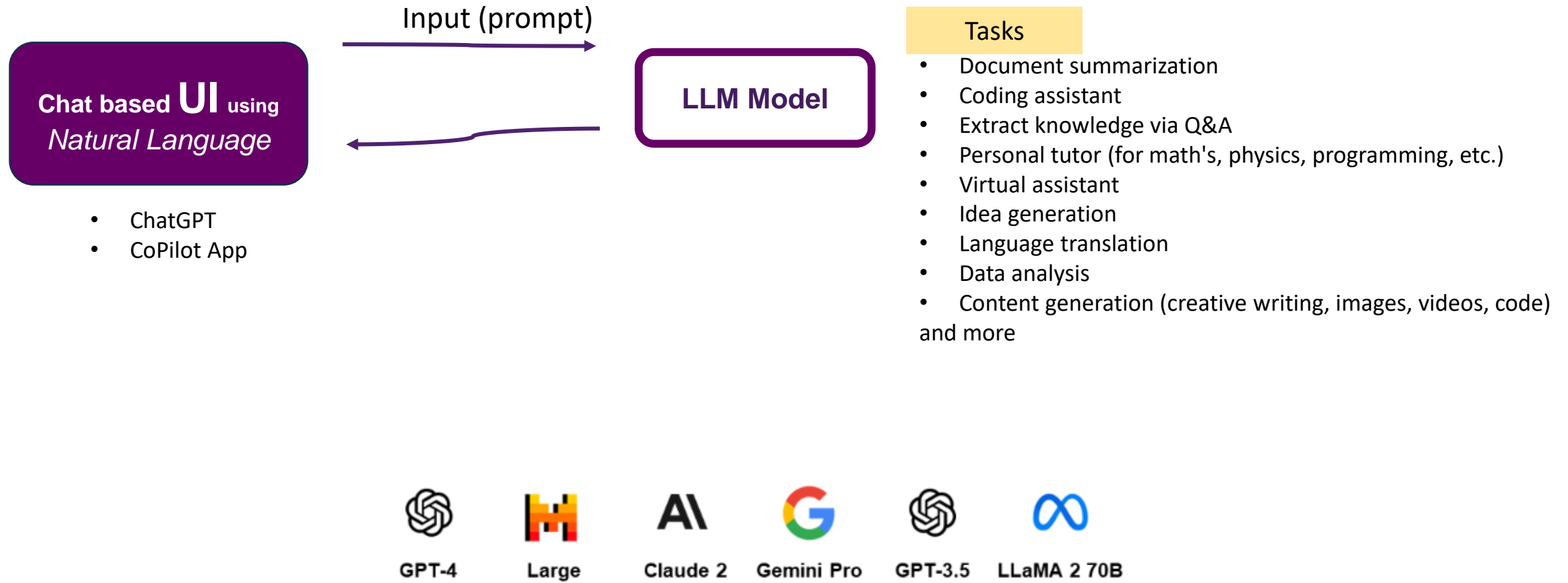- **Purpose:** The goal is to improve the model's ability to understand and execute complex instructions given in natural language, making it more user-friendly and effective for tasks specified by users through prompts.

GPT Assistant training pipeline

| Stage | Pretraining | Supervised Finetuning | Reward Modeling | Reinforcement Learning |
|---|---|---|---|---|
| Dataset | **Raw internet** text trillions of words low-quality, large quantity | **Demonstrations** Ideal Assistant responses, ~10-100K (prompt, response) written by contractors low quantity, high quality | **Comparisons** 100K –1M comparisons written by contractors low quantity, high quality | **Prompts** ~10K-100K prompts written by contractors low quantity, high quality |
| Algorithm | **Language modeling** predict the next token | **Language modeling** predict the next token | **Binary classification** predict rewards consistent w preferences | **Reinforcement Learning** generate tokens that maximize the reward |
| Model | **Base model** | init from → **SFT model** | init from → **RM model** | init from SFT use RM → **RL model** |
| Notes | 1000s of GPUs months of training ex: GPT, LLaMA, PaLM **can deploy this model** | 1-100 GPUs days of training ex: Vicuna-13B **can deploy this model** | 1-100 GPUs days of training | 1-100 GPUs days of training ex: ChatGPT, Claude **can deploy this model** |

Source:  https://www.youtube.com/watch?v=bZQun8Y4L2A&t=58s

# Large Language Models

LLMs have shown great promise as capable AI assistants for humans. LLMs can create new content, including text, images, videos, and music, that can resemble works made by humans. These AI systems are widely used for creativity, automation, and enhancing human work by providing novel ideas and solutions.

Input (prompt)

**Chat based UI using *Natural Language***

**LLM Model**

- ChatGPT
- CoPilot App

Tasks

- Document summarization
- Coding assistant
- Extract knowledge via Q&A
- Personal tutor (for math's, physics, programming, etc.)
- Virtual assistant
- Idea generation
- Language translation
- Data analysis
- Content generation (creative writing, images, videos, code) and more

GPT-4    Large    Claude 2    Gemini Pro    GPT-3.5    LLaMA 2 70B

# LLM inferencing

In the near future, AI will be infused into all applications. The process of invoking LLMs in applications is called **inferencing**.

LLMs can be used in apps via:

- **API**: Connect to LLM services online for easy access.
- **On-Premise**: Deploy on local servers for more control and privacy.
- **Edge Computing**: Run on local devices for low latency and offline use.

Each method balances performance, cost, and privacy differently. Small size models are more suitable for edge inferencing.

Gemma Open Models

| Parameters size | Input | Output | Tuned versions | Intended platforms |
|---|---|---|---|---|
| 2B | Text | Text | • Pretrained<br>• Instruction tuned | Mobile devices and laptops |
| 7B | Text | Text | • Pretrained<br>• Instruction tuned | Desktop computers and small servers |

Size & intended platform

# LLM Leaderboards

LLM leaderboards rank and compare language models based on performance metrics, track advancements, encourage innovation, and help users choose the best models for their needs.

Code to recreate leaderboard tables and plots in this notebook. You can contribute your vote at chat.lmsys.org!

| Category | Overall Questions |
|---|---|
| Overall ▼ | #models: **122 (100%)**  #votes: **1,559,385 (100%)** |

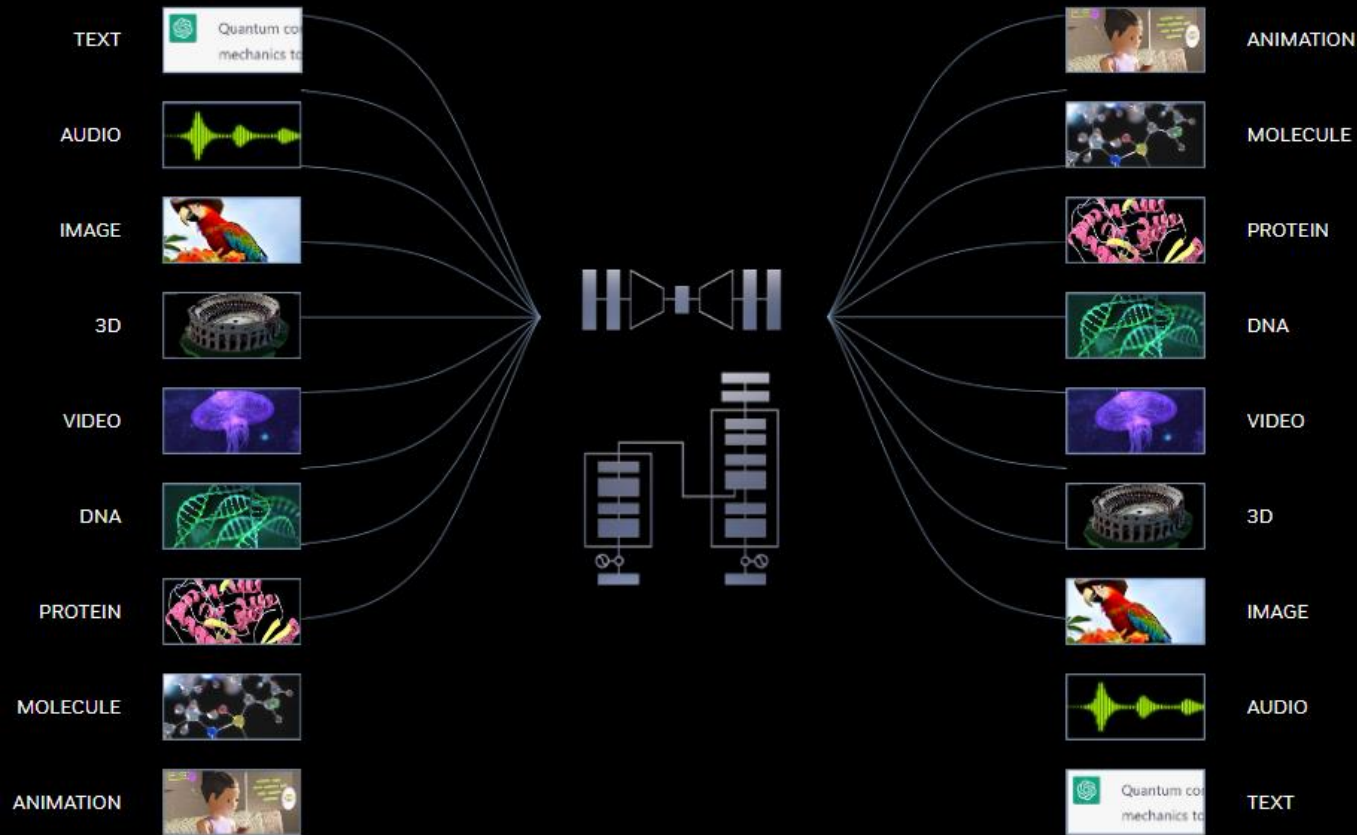| Rank* (UB) ▲ | Model ▲ | Arena Score ▲ | 95% CI ▲ | Votes ▲ | Organization ▲ | License ▲ | Knowledge Cutoff ▲ |
|---|---|---|---|---|---|---|---|
| 1 | GPT-4o-2024-05-13 | 1286 | +3/-3 | 68753 | OpenAI | Proprietary | 2023/10 |
| 1 | GPT-4o-mini-2024-07-18 | 1280 | +5/-6 | 11075 | OpenAI | Proprietary | 2023/10 |
| 2 | Claude 3.5 Sonnet | 1271 | +4/-2 | 38939 | Anthropic | Proprietary | 2024/4 |
| 3 | Gemini-Advanced-0514 | 1266 | +3/-4 | 50037 | Google | Proprietary | Online |
| 4 | Meta-Llama-3.1-405b-Instruct | 1262 | +7/-5 | 7322 | Meta | Llama 3.1 Community | 2023/12 |
| 4 | Gemini-1.5-Pro-API-0514 | 1261 | +3/-2 | 60928 | Google | Proprietary | 2023/11 |
| 5 | Gemini-1.5-Pro-API-0409-Preview | 1257 | +3/-3 | 55667 | Google | Proprietary | 2023/11 |
| 5 | GPT-4-Turbo-2024-04-09 | 1257 | +3/-3 | 78790 | OpenAI | Proprietary | 2023/12 |
| 9 | GPT-4-1106-preview | 1251 | +3/-3 | 89657 | OpenAI | Proprietary | 2023/4 |
| 9 | Claude 3 Opus | 1248 | +2/-3 | 150231 | Anthropic | Proprietary | 2023/8 |
| 9 | GPT-4-0125-preview | 1245 | +3/-3 | 82978 | OpenAI | Proprietary | 2023/12 |
| 9 | Athene-70b | 1245 | +7/-7 | 5137 | NexusFlow | CC-BY-NC-4.0 | 2024/7 |
| 9 | Meta-Llama-3.1-70b-Instruct | 1242 | +7/-7 | 3621 | Meta | Llama 3.1 Community | 2023/12 |
| 11 | Yi-Large-preview | 1240 | +3/-3 | 51499 | 01 AI | Proprietary | Unknown |
| 15 | Gemini-1.5-Flash-API-0514 | 1228 | +4/-3 | 50339 | Google | Proprietary | 2023/11 |
| 15 | Deepseek-v2-API-0628 | 1221 | +5/-5 | 10393 | DeepSeek AI | Proprietary | Unknown |

https://huggingface.co/spaces/lmsys/chatbot-arena-leaderboard

LLM benchmarking refers to the evaluation of large language models to assess their performance & efficiency across a variety of metrics.

| Capability | Benchmark<br>Higher is better | Description | Gemini 1.0 Ultra | GPT-4<br>API numbers calculated where reported numbers were missing |
|---|---|---|---|---|
| General | MMLU | Representation of questions in 57 subjects (incl. STEM, humanities, and others) | 90.0%<br>CoT@32* | 86.4%<br>5-shot** (reported) |
| Reasoning | Big-Bench Hard | Diverse set of challenging tasks requiring multi-step reasoning | 83.6%<br>3-shot | 83.1%<br>3-shot (API) |
| | DROP | Reading comprehension (F1 Score) | 82.4<br>Variable shots | 80.9<br>3-shot (reported) |
| | HellaSwag | Commonsense reasoning for everyday tasks | 87.8%<br>10-shot* | 95.3%<br>10-shot* (reported) |
| Math | GSM8K | Basic arithmetic manipulations (incl. Grade School math problems) | 94.4%<br>maj1@32 | 92.0%<br>5-shot CoT (reported) |
| | MATH | Challenging math problems (incl. algebra, geometry, pre-calculus, and others) | 53.2%<br>4-shot | 52.9%<br>4-shot (API) |
| Code | HumanEval | Python code generation | 74.4%<br>0-shot (IT)* | 67.0%<br>0-shot* (reported) |
| | Natural2Code | Python code generation. New held out dataset HumanEval-like, not leaked on the web | 74.9%<br>0-shot | 73.9%<br>0-shot (API) |

# Generative AI
## The most important computing platform of our generation

The era of generative AI has arrived, unlocking new opportunities for AI across many different applications

Generative AI is trained on large amounts of data to find patterns and relationships, learning the representation of almost anything with structure

It can then be prompted to generate text, images, video, code, or even proteins

For the very first time, computers can augment the human ability to generate information and create

1,600+ Generative AI companies are building on NVIDIA

Source: Nvidia

Source: blogs.nvidia.com/blog/llms-ai-horizon

## Importance of
## **Generative AI**

Improve productivity
Eliminate drudgery
Your reasoning engine
Increase innovation
Transform business
Personal Assistant

### CoPilots & Assistants

Empower humans in their line of work in business. Personal tutor.

### Universal UI

Natural language is the new interface for text, speech or video. Humans will learn & cocreate with AI using natural language

### AI Orchestrator: AI Agents

LLMs can function as AI orchestrators by coordinating the interaction between various systems & services.

# Gen AI introduces new risks

Gen AI offer great promise but comes with risks related to responsible AI. Gen AI systems can cause harm such as promote misinformation, hallucinate, etc. and lead to a wide range of other negative impacts..

## LLM models introduce new risks

**Bias & fairness**
LLMs can inherit and even amplify biases present in their training data. This can lead to outputs that are unfair or discriminatory, particularly in sensitive applications involving gender, race, or other personal characteristics.

**Security & Jailbreak**
refers to the potential vulnerabilities or threats that could lead to unauthorized access, data breaches, or misuse of the models. This includes concerns such as data leakage or manipulation, where sensitive information trained into the model might be inadvertently revealed through its responses.

**Hallucination**
instances where the model generates text that is factually incorrect, misleading, or entirely fabricated, despite being presented in a confident and plausible manner. This behavior can range from minor inaccuracies to completely erroneous statements.

**Offensive content**
LLM models may generate other types of inappropriate or offensive content, which may make it inappropriate to deploy for sensitive contexts without additional mitigations that are specific to the use case.