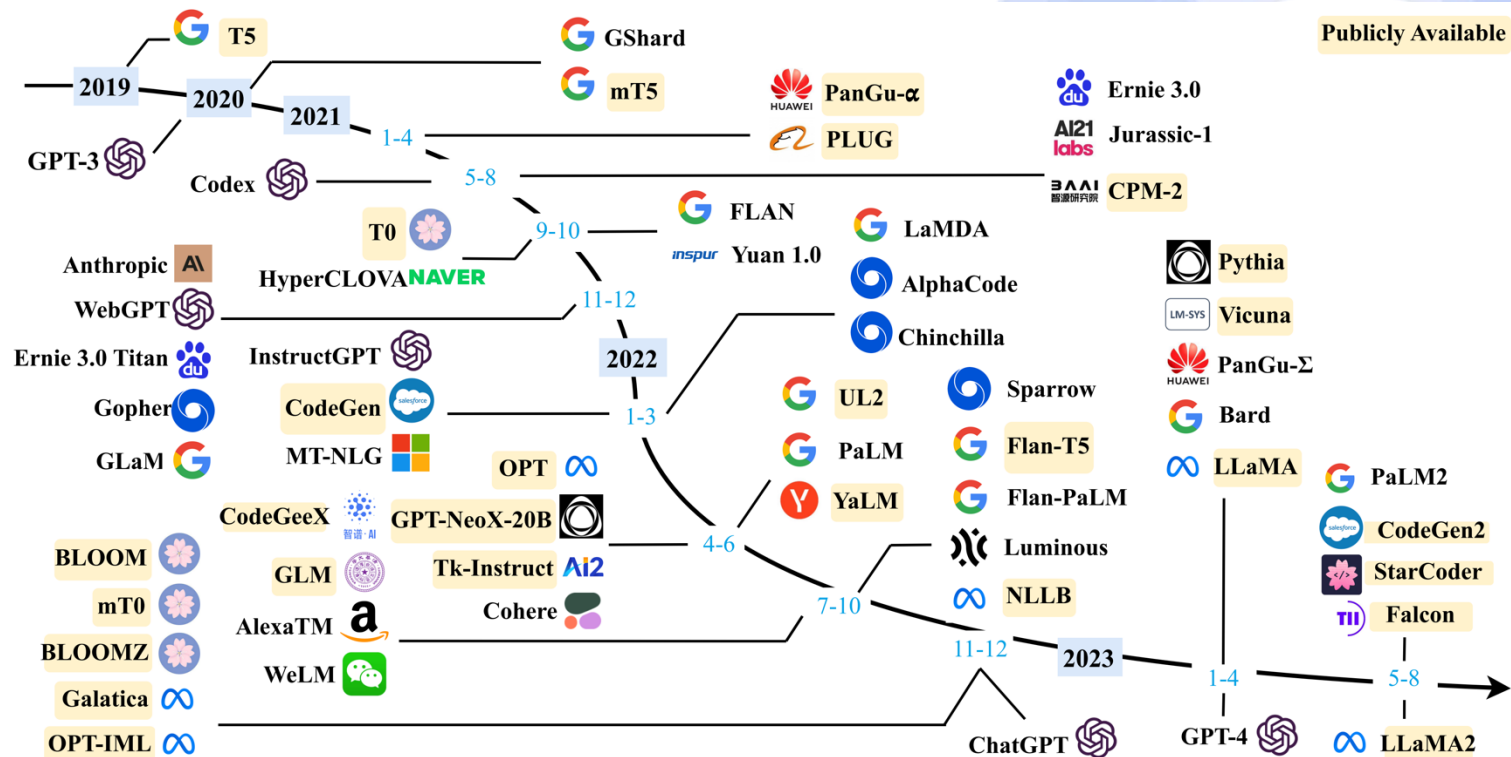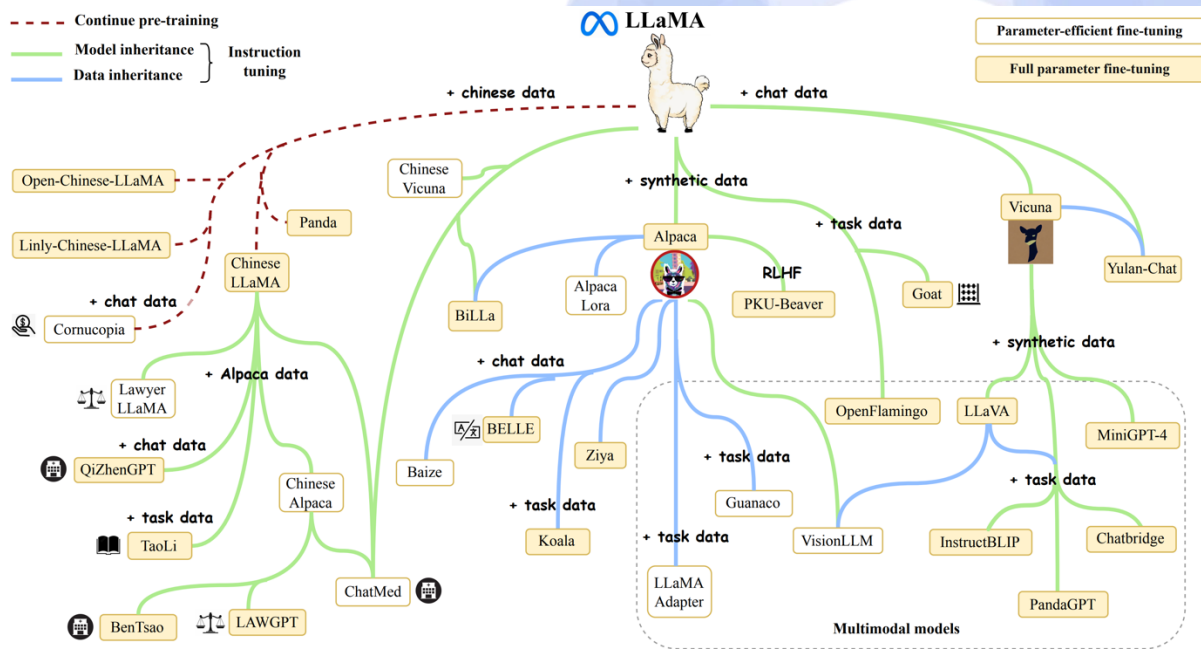# What's out there?

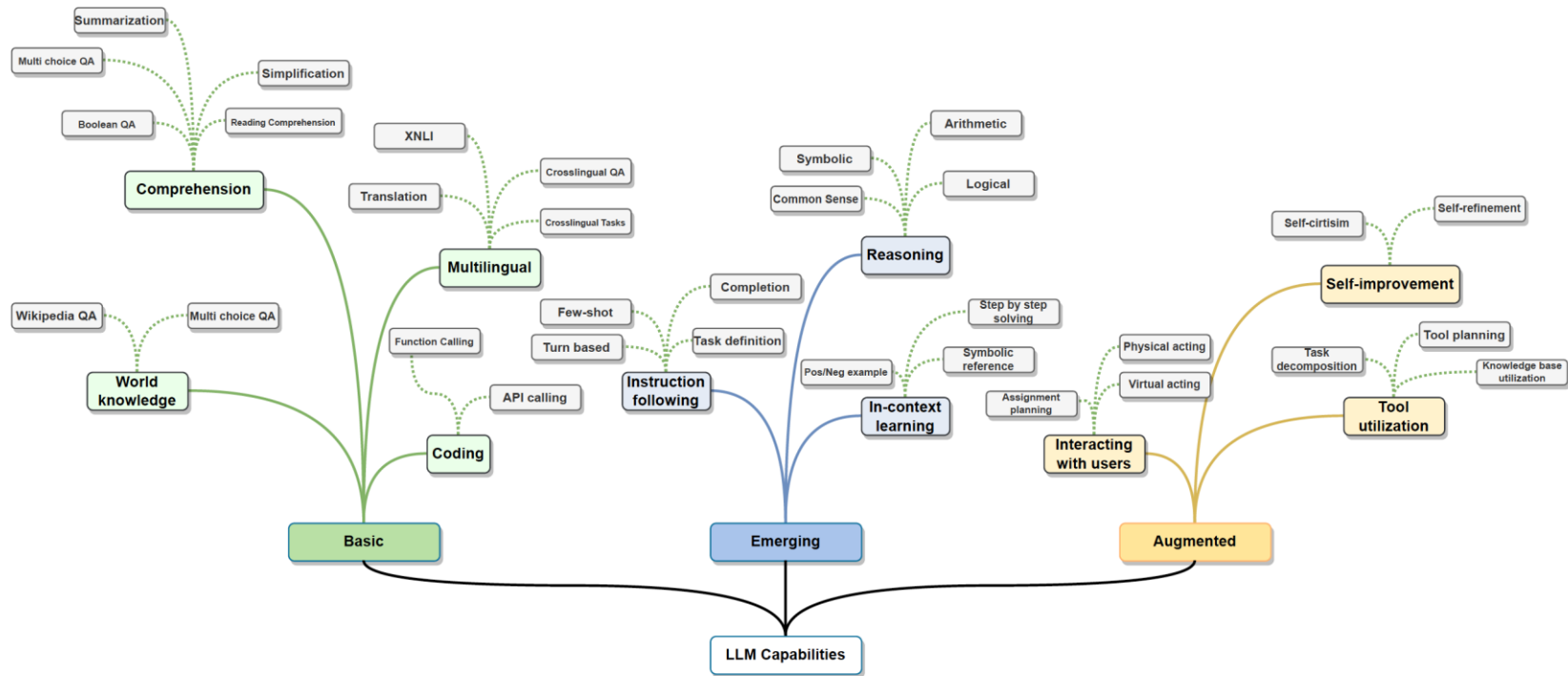# Timeline of >10B parameter models

# How can we categorize LLMs?

- Architecture
- Objective or use-case
- Scale
- Modalities
- Language
- Availability

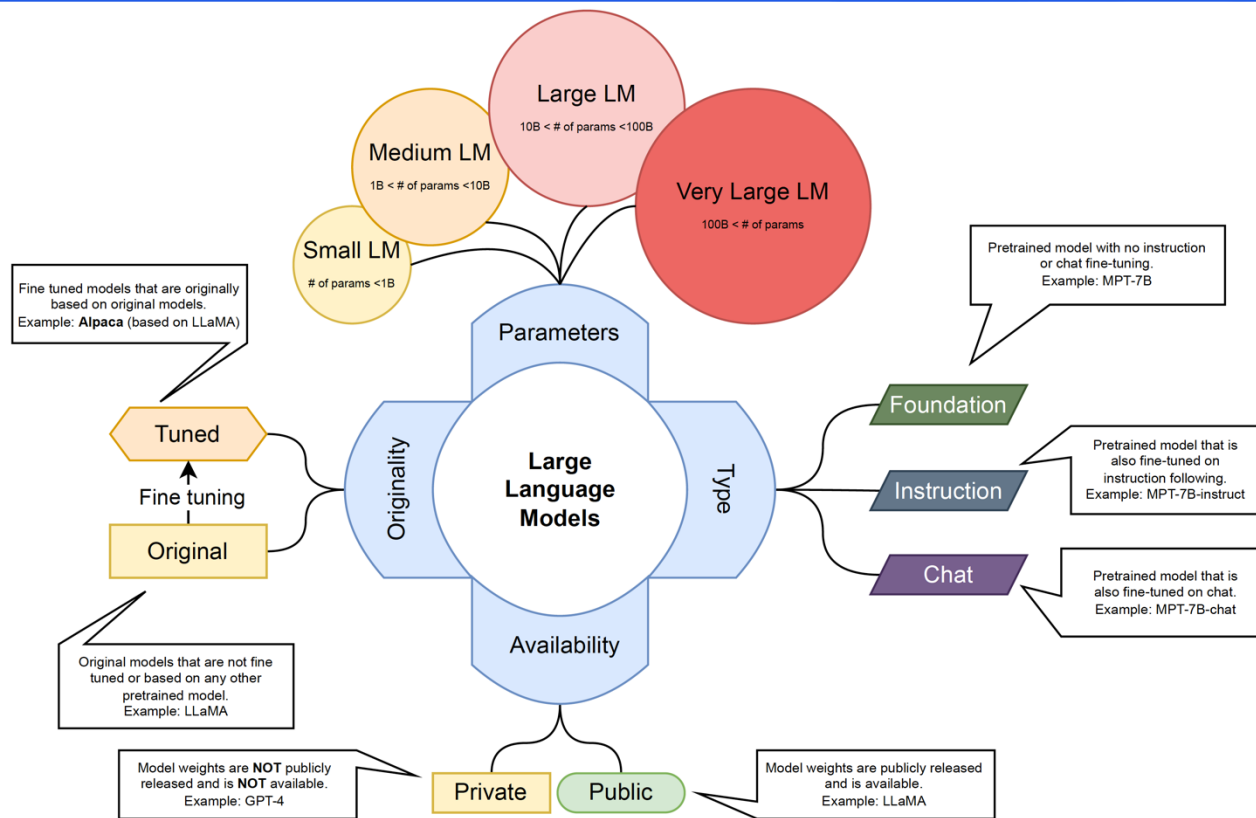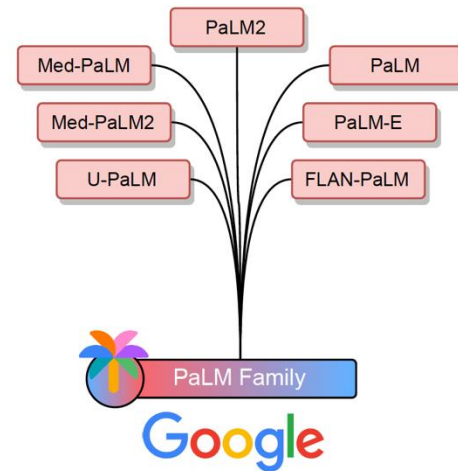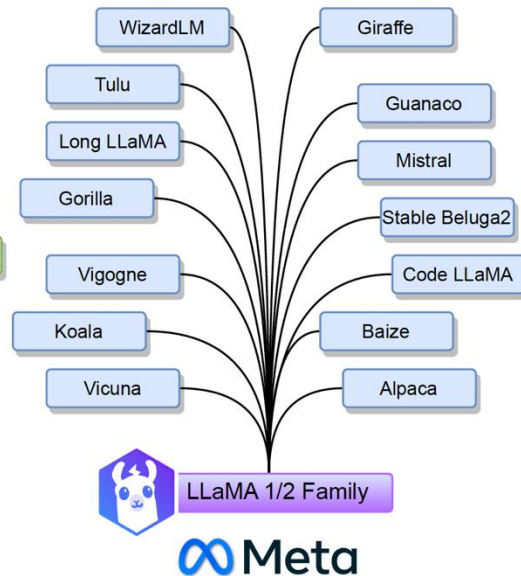Sometimes, one base model can evolve into all of the above!

# How can we categorize LLMs?

# How can we categorize LLMs?

# Popular LLM Families

# BERT

## Bidirectional Encoder Representations from Transformers

| '[CLS]' | 101 |
| 'the' | 'the' | 1996 |
| 'quick' | 'quick' | 4248 |
| 'brown' | 'brown' | 2829 |
| 'fox' | 'fox' | 4419 |
| 'jumps' | 'jumps' | 14523 |
| '[SEP]' | 102 |

$$X \in \mathbb{R}^{(2C+1) \times d}$$

| 0.17 | 0.01 | | 0.04 |
| 0.33 | 0.54 | | 0.27 |
| 0.71 | 0.57 | $e = X_w \in \mathbb{R}^{1 \times d}$ | 0.49 |
| 0.65 | 0.91 | $\circ \circ \circ$ | 0.11 |
| 0.31 | 0.45 | | 0.22 |
| 0.43 | 0.71 | | 0.42 |
| 0.96 | 0.19 | | 0.69 |

BROWN

Brown things

Surname

Brown hair/eyes

Brown in the negative sense

| 0.32 | 0.51 | $\circ \circ \circ$ | 0.99 |
| 0.05 | 0.13 | | 0.40 |
| 0.65 | 0.91 | $\circ \circ \circ$ | 0.11 |
| | | | |
| 0.65 | 0.01 | $\circ \circ \circ$ | 0.01 |

$e$

$$E \in \mathbb{R}^{n \times d}$$

# BERT

BERT is a very versatile family

DistilBERT

CamemBERT

ERNIE

XML-RoBERTa

RoBERTa

ALBERT

BART

Q-BERT

# CLIP

## Contrastive Language-Image Pre-training

○ *Learning Transferable Visual Models for Natural Language Supervision*, Radford et al, 2021, ~17k citations

○ Takes in images and texts and connects them together



Radford, Alec, et al. "Learning transferable visual models from natural language supervision." International conference on machine learning. PMLR, 2021.

# What does OpenAI offer?

Browser-based service:

o ChatGPT and ChatGPT Plus

o DALL-E

o GPTs

API-based service:

o Embeddings

o Fine-tuning

o Image generation

o Vision

o Text-to-speech (and speech-to-text)

# Open source LLMs



**Audio**
- 🔊 Text-to-Speech
- 🔊 Text-to-Audio
- 👤 Automatic Speech Recognition
- 🎚 Audio-to-Audio
- 🎵 Audio Classification
- 🗣 Voice Activity Detection

**Tabular**
- 🗒 Tabular Classification
- 📈 Tabular Regression

**Reinforcement Learning**
- 🎮 Reinforcement Learning
- 🤖 Robotics

**Multimodal**
- 🔲 Feature Extraction
- 📝 Text-to-Image
- 🖼 Image-to-Text
- 🗂 Text-to-Video
- 🔳 Visual Question Answering
- 📄 Document Question Answering
- 🔗 Graph Machine Learning

**Computer Vision**
- Depth Estimation
- 🖼 Image Classification
- Object Detection
- ⬜ Image Segmentation
- 🖼 Image-to-Image
- 🖼 Unconditional Image Generation
- 📹 Video Classification
- 🔳 Zero-Shot Image Classification

**Natural Language Processing**
- ⠿ Text Classification
- 🔲 Token Classification
- 🔲 Table Question Answering
- 🔲 Question Answering
- ✳ Zero-Shot Classification
- 🔤 Translation
- 📑 Summarization
- 💬 Conversational
- 📝 Text Generation
- 🔲 Text2Text Generation
- 🔳 Fill-Mask
- 🔲 Sentence Similarity

# But is it really open source…?

Here are some things to look out for:

- Weights
- Training data (and RL data)
- Training code
- License
- Architecture
- Preprint/paper
- Modelcard
- Package
- API

| Project | Availability | | | | | | Documentation | | | | | | Access | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Open code | LLM data | LLM weights | RL data | RL weights | License | Code | Architecture | Preprint | Paper | Modelcard | Datasheet | Package | API |
| OLMo 7B Instruct | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ~ |
| BLOOMZ | ✓ | ✓ | ✓ | ✓ | ~ | ~ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ |
| AmberChat | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ~ | ~ | ✓ | ✗ | ~ | ~ | ✗ | ✓ |
| Open Assistant | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ~ | ✗ | ✗ | ✗ | ✓ | ✓ |
| OpenChat 3.5 7B | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ~ | ✗ | ✗ | ~ |
| Pythia-Chat-Base-7... | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ~ | ~ | ✗ | ~ | ~ | ✓ | ✗ |
| Cerebras GPT 111... | ~ | ✓ | ✓ | ✓ | ✓ | ~ | ✗ | ✓ | ~ | ✗ | ✗ | ✓ | ✗ | ✓ |
| RedPajama-INCITE... | ~ | ✓ | ✓ | ✓ | ✓ | ✓ | ~ | ~ | ✗ | ✗ | ✓ | ✓ | ✗ | ~ |
| dolly | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ~ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Tulu V2 DPO 70B | ✓ | ✗ | ~ | ✓ | ✓ | ✓ | ~ | ~ | ~ | ✓ | ✗ | ~ | ✗ | ✓ |
| MPT-30B Instruct | ✓ | ~ | ✓ | ~ | ✗ | ✓ | ✓ | ~ | ✗ | ✗ | ~ | ✗ | ✓ | ~ |
| MPT-7B Instruct | ✓ | ~ | ✓ | ~ | ✗ | ✓ | ✓ | ~ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| trlx | ✓ | ✓ | ✓ | ~ | ✗ | ✓ | ✓ | ~ | ✗ | ✗ | ✗ | ✗ | ~ | ✓ |
| Vicuna 13B v 1.3 | ✓ | ~ | ✓ | ✗ | ✗ | ~ | ✓ | ✗ | ✓ | ✗ | ~ | ✗ | ✗ | ~ |
| minChatGPT | ✓ | ~ | ✓ | ~ | ✗ | ✓ | ✓ | ✓ | ~ | ✗ | ✗ | ✗ | ✗ | ✓ |
| ChatRWKV | ✓ | ~ | ✓ | ✗ | ✗ | ✓ | ~ | ~ | ~ | ✗ | ✗ | ✗ | ✓ | ~ |
| BELLE | ✓ | ~ | ~ | ~ | ~ | ✗ | ~ | ✓ | ✓ | ✗ | ✗ | ~ | ✗ | ✗ |
| WizardLM 13B v1.2 | ~ | ✗ | ~ | ✓ | ✓ | ~ | ~ | ✓ | ~ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Airoboros L2 70B G... | ~ | ✗ | ~ | ✓ | ✓ | ~ | ~ | ~ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| ChatGLM-6B | ~ | ~ | ✓ | ✗ | ✗ | ✓ | ~ | ~ | ~ | ✗ | ~ | ✗ | ✗ | ✓ |
| Mistral 7B-Instruct | ~ | ✗ | ✓ | ✗ | ~ | ✓ | ✗ | ~ | ~ | ✗ | ✗ | ✗ | ~ | ✓ |

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **WizardLM-7B** | ~ | ~ | ✗ | ✓ | ~ | ~ | | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **Qwen 1.5** | ~ | ✗ | ✓ | ✗ | ✓ | ✗ | ~ | ~ | ✗ | ✗ | ✗ | ✗ | ~ | ✓ |
| **StableVicuna-13B** | ~ | ✗ | ~ | ~ | ~ | ~ | ~ | | ✗ | ~ | ✗ | ✗ | ~ |
| **Falcon-40B-instruct** | ✗ | ~ | ✓ | ~ | ✗ | ✓ | ✗ | | ✗ | ~ | ✗ | ✗ | ✗ |
| **UltraLM** | ✗ | ✗ | ~ | ✓ | ~ | ✗ | ✗ | ✓ | ~ | ~ | ✗ | ✗ |
| **Yi 34B Chat** | ~ | ✗ | ✓ | ✗ | ✓ | ~ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ~ |
| **Koala 13B** | ✓ | ~ | ~ | ~ | ✗ | ~ | | ✗ | ✗ | ✗ | ✗ | ✗ |
| **Mixtral 8x7B Instruct** | ✗ | ✗ | ✓ | ✗ | ~ | ✓ | ✗ | ~ | ✗ | ✗ | ✗ | ~ | ✗ |
| **Stable Beluga 2** | ✗ | ✗ | ~ | ✗ | ✓ | ~ | ✗ | ~ | ✗ | ✗ | ✗ | ✗ | ~ |
| **Stanford Alpaca** | ✓ | ✗ | ~ | ~ | ~ | ✗ | ~ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **Falcon-180B-chat** | ✗ | ~ | ~ | ~ | ✗ | ✗ | ~ | | ✗ | ~ | ✗ | ✗ | ✗ |
| **Orca 2** | ✗ | ✗ | ~ | ✗ | ✓ | ✗ | ~ | | ✗ | ~ | ✗ | ✗ | ~ |
| **Command R+** | ✗ | ✗ | ✗ | ✓ | ✓ | ~ | ✗ | ✗ | ✗ | ✗ | ~ | ✗ | ✗ | ✗ |
| **Gemma 7B Instruct** | ~ | ✗ | ~ | ✗ | ~ | ✗ | ✗ | ~ | ✗ | ✓ | ✗ | ✗ | ✗ |
| **LLaMA2 Chat** | ✗ | ✗ | ~ | ✗ | ~ | ✗ | ✗ | ~ | ✗ | ✗ | ~ | ✗ | ~ |
| **Nanbeige2-Chat** | ✓ | ✗ | ✗ | ✗ | ✓ | ~ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ~ |
| **Llama 3 Instruct** | ✗ | ✗ | ~ | ✗ | ~ | ✗ | ✗ | ~ | ✗ | ✗ | ~ | ✗ | ~ |
| **Solar 70B** | ✗ | ✗ | ~ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **Xwin-LM** | ✗ | ✗ | ~ | ✗ | ✗ | ~ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ~ |
| **ChatGPT** | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ~ | ✗ | ✗ | ✗ | ✗ | ✗ |

Liesenfeld, Andreas, and Mark Dingemanse. "Rethinking open source generative AI: open washing and the EU AI Act." The 2024 ACM Conference on Fairness, Accountability, and Transparency. 2024.

# How can you run open source models?

There are two main ways to access open source models:

○ API
  - Huggingface Hub API
  - Individual company APIs
  - No compute required but restrictions on requests

○ Cloud
  - Google Colab
  - Commercial cloud providers are expensive
  - Steep learning curve

○ Locally
  - Code
  - No code
  - Requires compute…

# No code

# No Code

Some no code options for running LLMs locally

o LMStudio
- Not open source
- Desktop client

o GPT4All
- Open source desktop client
- Upload documents for question answering
- Bindings for Python and NodeJS, and LangChain itegration

o Textgen-webui
- Open source
- Runs in the browser which means you can spin it up on a remote GPU
- Requires some effort to install and run on Linux machines.
- Ability to fine-tune models

# Resources

LMStudio

GPT4All

Textgen-webui

OpenAI API Usage, Pricing, and Policy

HuggingFace Tutorials

# How to keep up?

- o LinkedIn
    - – OpenAI
    - – Meta and AI at Meta
    - – Hugging Face
    - – Microsoft and Microsoft Research
    - – Apple
    - – PyTorch
- o YouTube
    - – [Fireship](Fireship)
    - – [Two Minute Papers](Two%20Minute%20Papers)
    - – [Joma Tech](Joma%20Tech)
- o Reddit
    - – r/LocalLLaMA
    - – r/MachineLearning and r/learnmachinelearning
    - – r/OpenAI
    - – r/homelabs
    - – r/LLMDevs
    - – r/StableDiffusion