# ChatNMI & Beyond: Pushing the 'Easy Button' in Home-based AI Deployment - A Comparative Review of Leading Software Solutions

Konrad Jędrzejczyk, Marek Zmysłowski

# whoami



**Konrad** is an expert in both defensive and offensive security, having gained experience from prestigious global institutions such as Credit Suisse and the Royal Bank of Scotland. After dedicating a decade to Blue Team operations, he made a transition to the Red Team and has now been excelling in this role for over 3 years. Konrad is the creator of the ChatNMI (No Moral Issues) tool for AI and has been a featured speaker at several conferences, including DefCamp, BSides and Oh My H@ck.



**Marek** is Offensive Security Certified Professional (OSCP) and Offensive Security Certified Expert (OSCE). He brings to the table 16 years of penetration testing experience, showcasing expertise in C and C++ programming, reverse engineering, and network and infrastructure security. Marek has shared his insights at various conferences, such as HiTB, PacSec, BlueHat, DefCamp, Confidence, and Oh My H@ck.
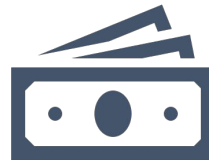
Motivation

# Motivation

## Privacy

- We want to protect the information we ask the AI about.

- We want to have honest answers, not those that are "politically correct".

## Money

- Sooner or later everything will be paid.
- One cannot control the expense.

## Risks

- Service availability.

- Bans imposed by governments or other entities.

- Changes in usage policy.

- Ethical or responsibility concerns related to AI.

# Privacy

"Forget artificial intelligence – in the brave new world of big data, it's artificial idiocy we should be looking out for."

—Tom Chatfield

# Privacy

Facebook taking my photo.

Me sending internal emails/projects to ChatGPT for rewriting and better English.
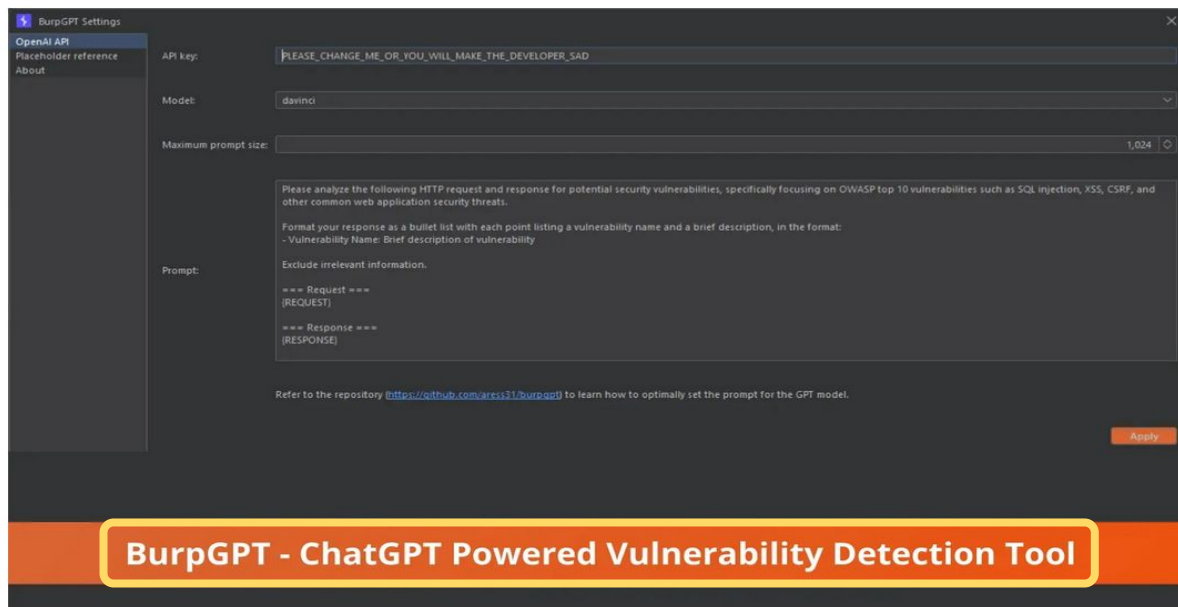
# Privacy

## ChatGPT on BurpSuite 🔗

This guide will walk you through the process of integrating OpenAI's ChatGPT language model with BurpSuite, a popular web application security testing tool. With this integration, you can use ChatGPT to suggest payloads or identify potential vulnerabilities during security testing.

You can now connect Lucidchart with ChatGPT!

Available to ChatGPT+ users, the Lucidchart plugin allows you to auto-generate diagrams for just about anything.

In ChatGPT+, add a prompt for what you want to visualize, and the plugin will generate a diagram for that prompt. The diagram is then editable in Lucidchart, so you can modify and share it.

**BurpGPT - ChatGPT Powered Vulnerability Detection Tool**

Cyber Security News came across a new ChatGPT-powered Vulnerability detection Tool called "BurpGPT," which helps security researchers to detect the vulnerabilities that traditional scanners might miss.
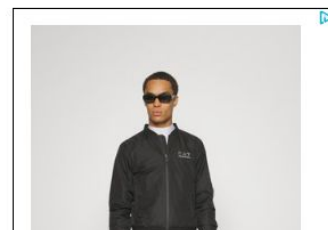
Like **PentestGPT**, a ChatGPT Powered Automated Penetration Testing Tool, BurpGPT was developed with deep vulnerability scanning features.

## How to Enable and Use Code Interpreter in ChatGPT
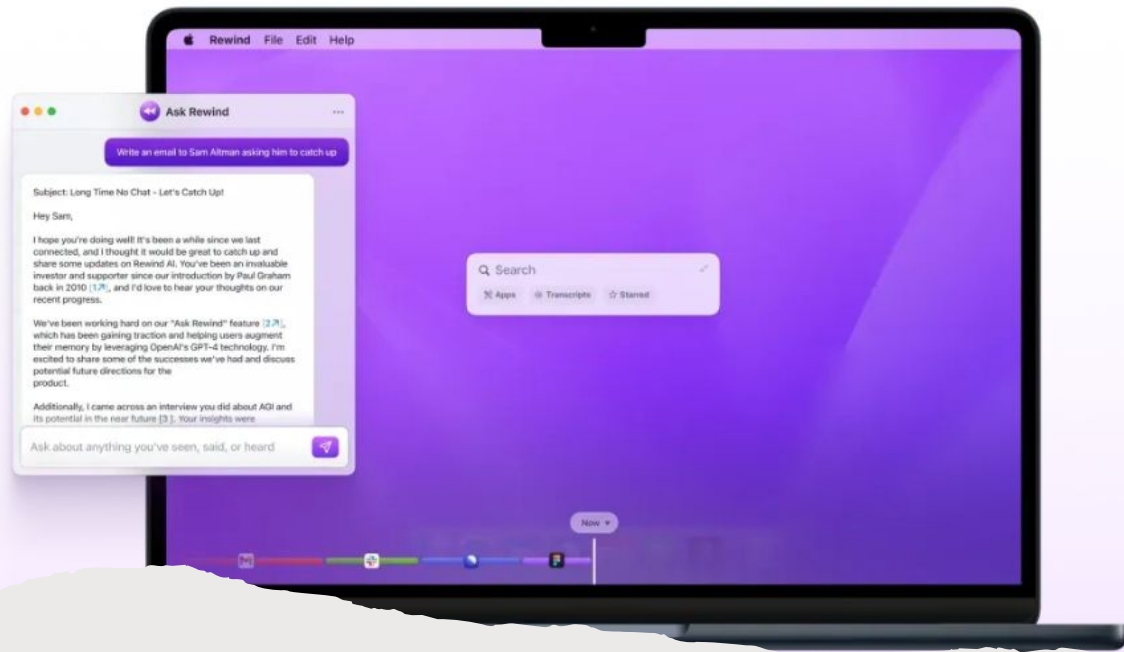
Arjun Sha - Last Updated: July 11, 2023 11:42 am

In case you are unaware, OpenAI has finally opened access to Code Interpreter for all ChatGPT Plus users. It is one of the most powerful feature additions to ChatGPT, and it allows you to perform a multitude of tasks. Similar to ChatGPT Plugins, Code Interpreter expands the utility of ChatGPT and allows it to interact with local files in various file formats. So if you are interested, follow our tutorial and learn how to enable and use Code Interpreter in ChatGPT.

Rewind is a privacy-first app **that records what you do and makes it searchable.**

Privacy

# Privacy



**RESEARCH**

# Introducing a watermarking method to distinguish images created by Generative AI

∞ Meta                                                    AI at Meta

AI-powered image generation is booming and for good reason: It's fun, entertaining, and easy to use. While these models enable new creative possibilities, they may raise concerns about potential misuse from bad actors who may intentionally generate images to deceive people. Even images created in good fun could still go viral and potentially mislead people. For example, earlier this year, images appearing to show Pope Francis wearing a flashy white puffy jacket went viral. The images weren't actual photographs, but plenty of people were fooled, since there weren't any clear indicators to distinguish that the content was created by generative AI.

# Privacy
# (Political correctness)

# Money

| Language Models - OpenAI | | | | |
|---|---|---|---|---|
| | | /1K tokens | | |
| Model | Context/Type | Input | Output | Training |
| GPT-4 | 8K context | $0.03 | $0.06 | - |
| | 32K context | $0.06 | $0.12 | - |
| GPT-3.5 Turbo | 4K context | $0.0015 | $0.002 | - |
| | 16K context | $0.003 | $0.004 | - |
| Fine-Tuning | babbage-002 | $0.0016 | $0.0016 | $0.0004 |
| | davinci-002 | $0.0120 | $0.0120 | $0.0060 |
| | GPT-3.5 Turbo | $0.0120 | $0.0160 | $0.0080 |

# Money

| Embedding Models/1K tokens - OpenAI | |
| --- | --- |
| Ada v2 | $0.0001 |
| Base Models | |
| babbage-002 | $0.0004/1K tokens |
| davinci-002 | $0.0020/1K tokens |
| Image Models | |
| 1024x1024 | $0.020/image |
| 512x512 | $0.018/image |
| 256x256 | $0.016/image |
| Audio Models | |
| Whisper | $0.006/minute |

LO

## Subscription increasement

# $5.000
/month

Due to inflation, we have to increase our subscription costs. In order to continue to have access to ChatGPT's services, you must take out this subscription.

**Money/Risks**

BUY

# Money/Risks

Risk

# chat.openai.com is down

Incident Report for OpenAI

**SUBSCRIBE TO UPDATES**

**OpenAI**

## Labs is having an outage

Subscribe

**Update** - We are allowing some free Labs traffic while we investigate an issue with end to end request latency causing a backlog of requests to be processed.
Mar 20, 2023 - 13:16 PDT

**Monitoring** - We are gradually restoring service to free traffic.
Mar 20, 2023 - 12:54 PDT

**Update** - Paid labs traffic has been restored and we will soon begin gradually restoring free traffic.
Mar 20, 2023 - 12:01 PDT

nitoring

We are gradually rolling out a fix that users will receive as capacity allows.

Posted 5 minutes ago. Mar 20, 2023 - 13:12 PDT

ntified

We have identified the root cause the issue and are working on resolution.

Posted 1 hour ago. Mar 20, 2023 - 11:58 PDT

**Identified** - We are adding additional capacity in other regions while investigating the underlying capacity failure.

## Risk

ate

We're continuing to investigate the source of the outage.

Posted 2 hours ago. Mar 20, 2023 - 10:58 PDT

**Update** - We're investigating an unintentional reduction in available capacity.
Mar 20, 2023 - 10:19 PDT

estigating

We are investigating an issue with the ChatGPT web experience.

Posted 4 hours ago. Mar 20, 2023 - 09:41 PDT

**Investigating** - We are currently investigating.
Mar 20, 2023 - 10:14 PDT

⚠ ChatGPT is temporarily unavailable. Please check here for updates:
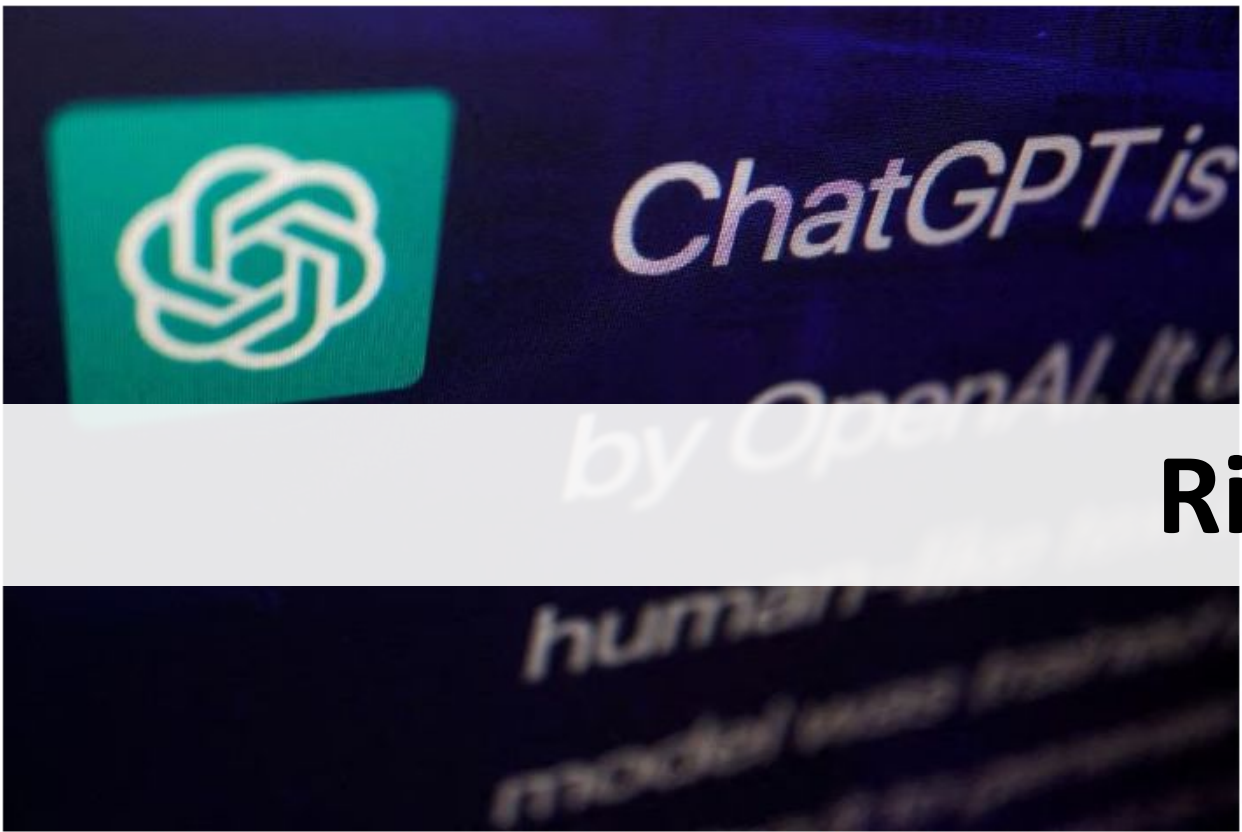https://status.openai.com/incidents/jq9232rcmktd

Subscribe

is incident a

**Identified** - We have identified the root cause the issue and are working on resolution.
Mar 20, 2023 - 11:58 PDT

rrent Status

Powered by Atlassian Statuspage

3 minute read · March 31, 2023 11:40 PM GMT+2 · Last Updated 2 days ago

# Italy curbs ChatGPT, starts probe over privacy concerns

By Elvira Pollina and Supantha Mukherjee

A response by ChatGPT, an AI chatbot developed by OpenAI, is seen on its website in this illustration picture taken February 9, 2023. REUTERS/Florence Lo/Illustration/File Photo

📑 Summary    🏛 Companies

• Italy first western country to ban ChatGPT

# ChatGPT banned in Italy over privacy concerns

🕐 1 day ago

**Risk**

GETTY IMAGES

# responsible AI

By **Alexander S. Gillis,** Technical Writer and Editor
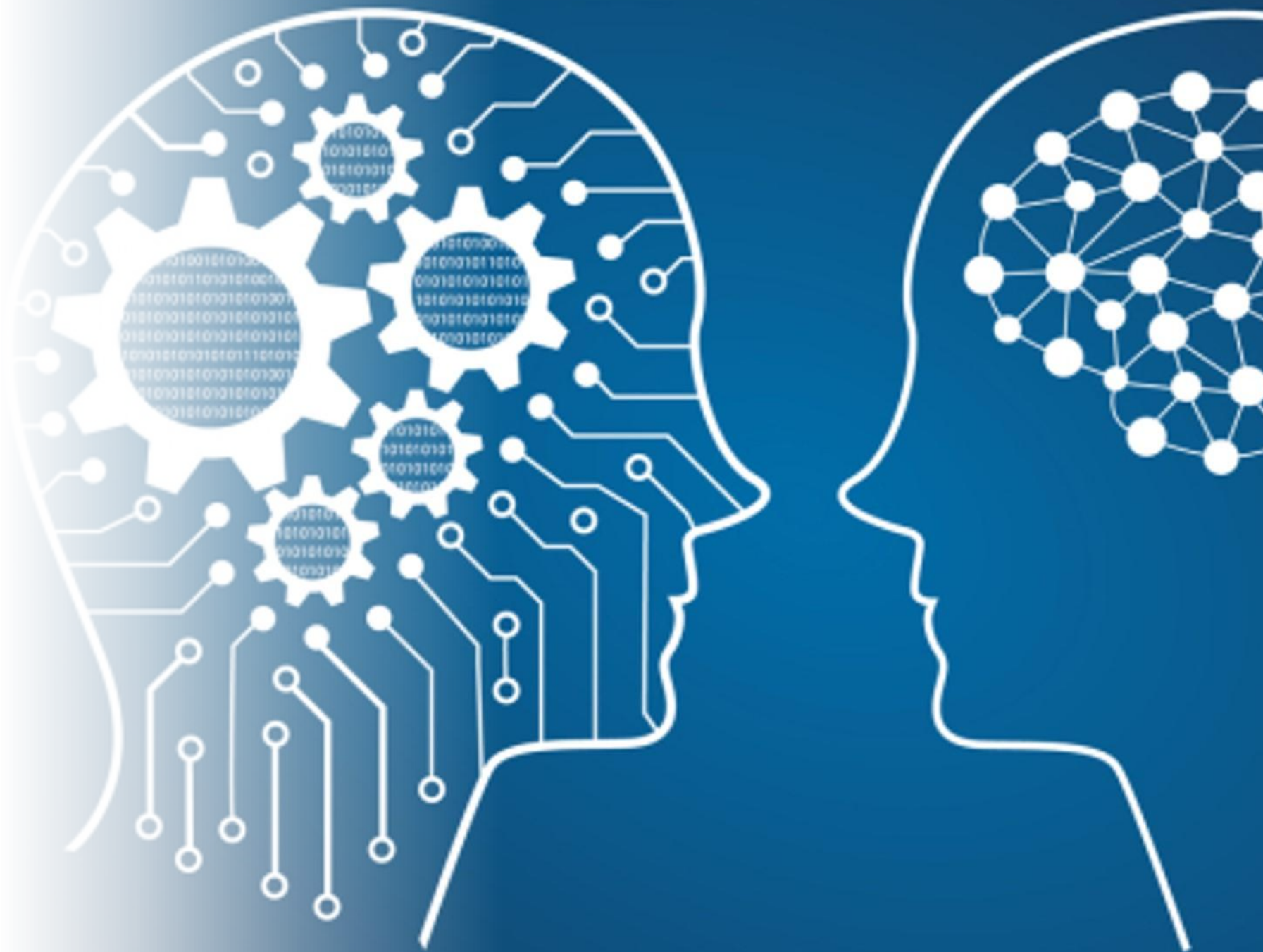
## What is responsible AI?

Responsible AI is an approach to developing and deploying artificial intelligence (AI) from both an ethical and legal point of view. The goal of responsible AI is to employ AI in a safe, trustworthy and ethical fashion. Using AI responsibly should increase transparency and help reduce issues such as AI bias.

Proponents of responsible AI hope that a widely adopted governance framework of AI best practices makes it easier for organizations around the globe to ensure their AI programming is human-centered, interpretable and explainable. Having a responsible AI system in place ensures fairness, reliability and transparency.

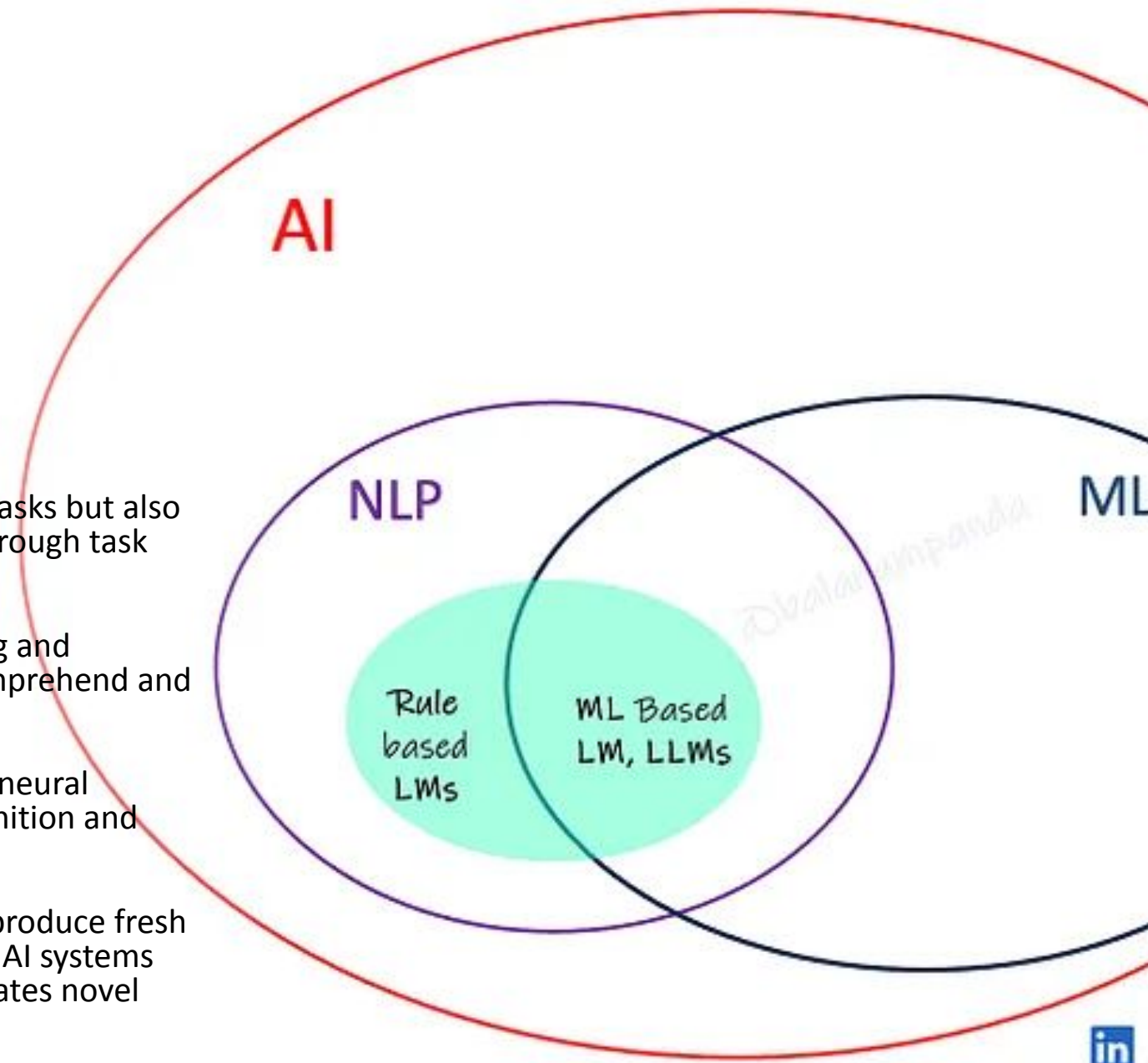# AI Theory

# General Classification

**AI**: The capability of a system to execute intelligent functions.

**Machine Learning**: The capacity not only to perform intelligent tasks but also to acquire knowledge from experience, improve performance through task execution, and learn from data.

**NLP (Natural Language Processing)**: Involves both understanding and generating human language, encompassing the capability to comprehend and create text or speech.

**Deep Learning**: A branch of machine learning that emulates the neural networks in the human brain to facilitate complex pattern recognition and learning.

**Generative AI**: A subset of artificial intelligence algorithms that produce fresh outputs based on their training data. In contrast to conventional AI systems that identify patterns and make predictions, generative AI generates novel content, including images, text, audio, and more.

# LLM (Large Language Model)

- A generic term that refers to transformer language models (GPT-3, BLOOM, OPT) that were trained on a large quantity of data. These models also tend to have a large number of learnable parameters (e.g. 175 billion for GPT-3).

- A language model trained for causal language modelling takes a sequence of text tokens as input and returns the probability distribution for the next token.

- Large, pretrained transformer models trained to predict the next word (or, more precisely, token) given some input text. Since they predict one token at a time, to generate new sentences, autoregressive generation is needed.

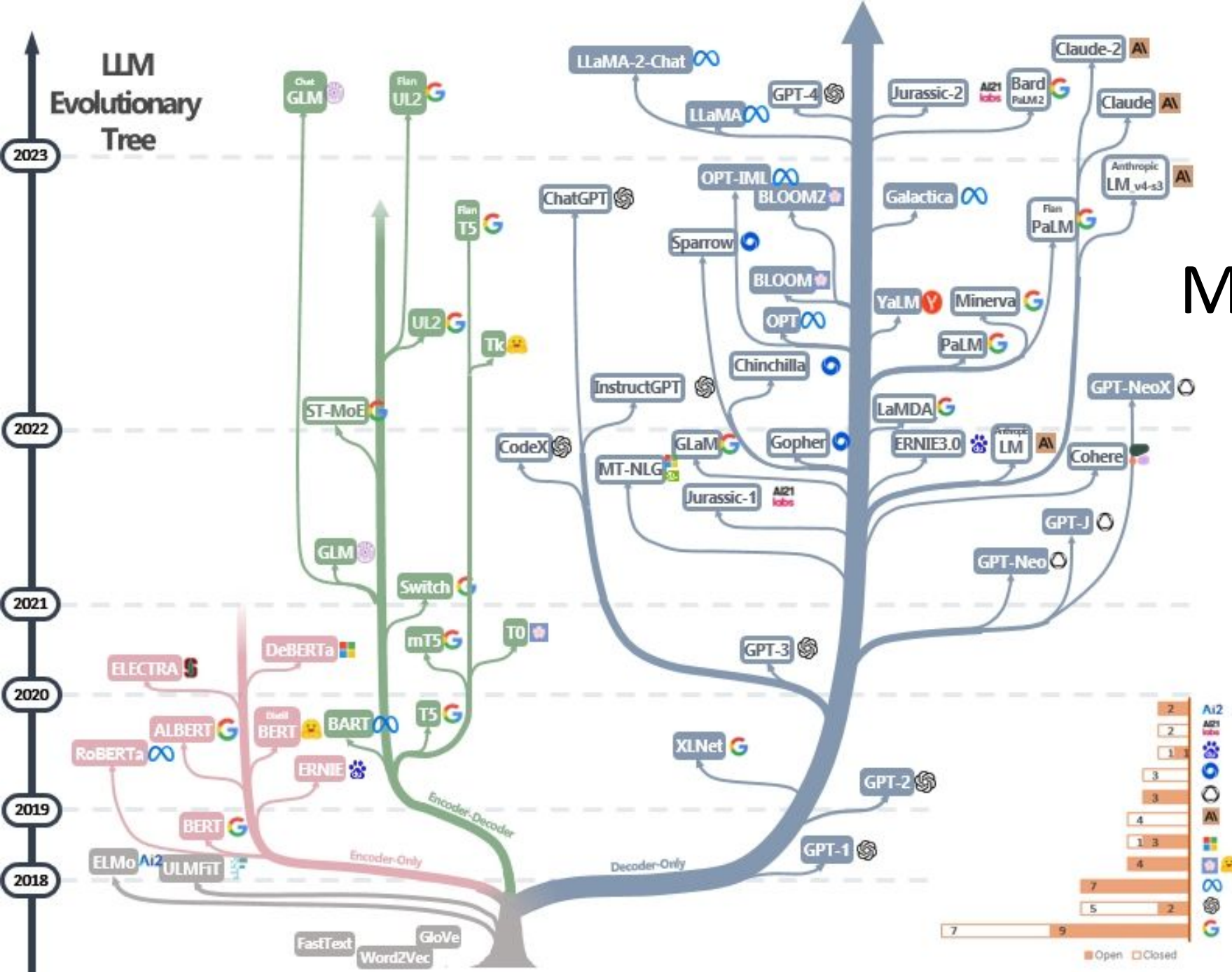https://huggingface.co/docs/transformers/main/llm_tutorial

# Types of LLMs

**Transformer-based** - language models, operate by analyzing and generating text through a blend of self-attention mechanisms, positional encoding, and multi-layer neural networks.

**LSTM (Long Short-Term Memory)** - belongs to the category of recurrent neural networks (RNNs) capable of capturing extended contextual relationships within textual data.

**ELMo (Embeddings from Language Models)** - ELMo represents a sizable language model created by the Allen Institute for AI, designed to generate contextualized word embeddings.

# Types of LLMs

**Transformer-based** - language models, operate by analyzing and generating text through a blend of self-attention mechanisms, positional encoding, and multi-layer neural networks. At the heart of the Transformer architecture lies the self-attention mechanism, which constructs a weighted representation of the input sequence, taking into account the interconnections among various segments of the text. This functionality empowers the model to grasp extensive contextual information and long-range dependencies within the text.

**LSTM (Long Short-Term Memory)** - belongs to the category of recurrent neural networks (RNNs) capable of capturing extended contextual relationships within textual data. LSTM-driven models, exemplified by ULMFiT (Universal Language Model Fine-tuning), have found utility in various applications, including text classification, sentiment analysis, and language modeling.
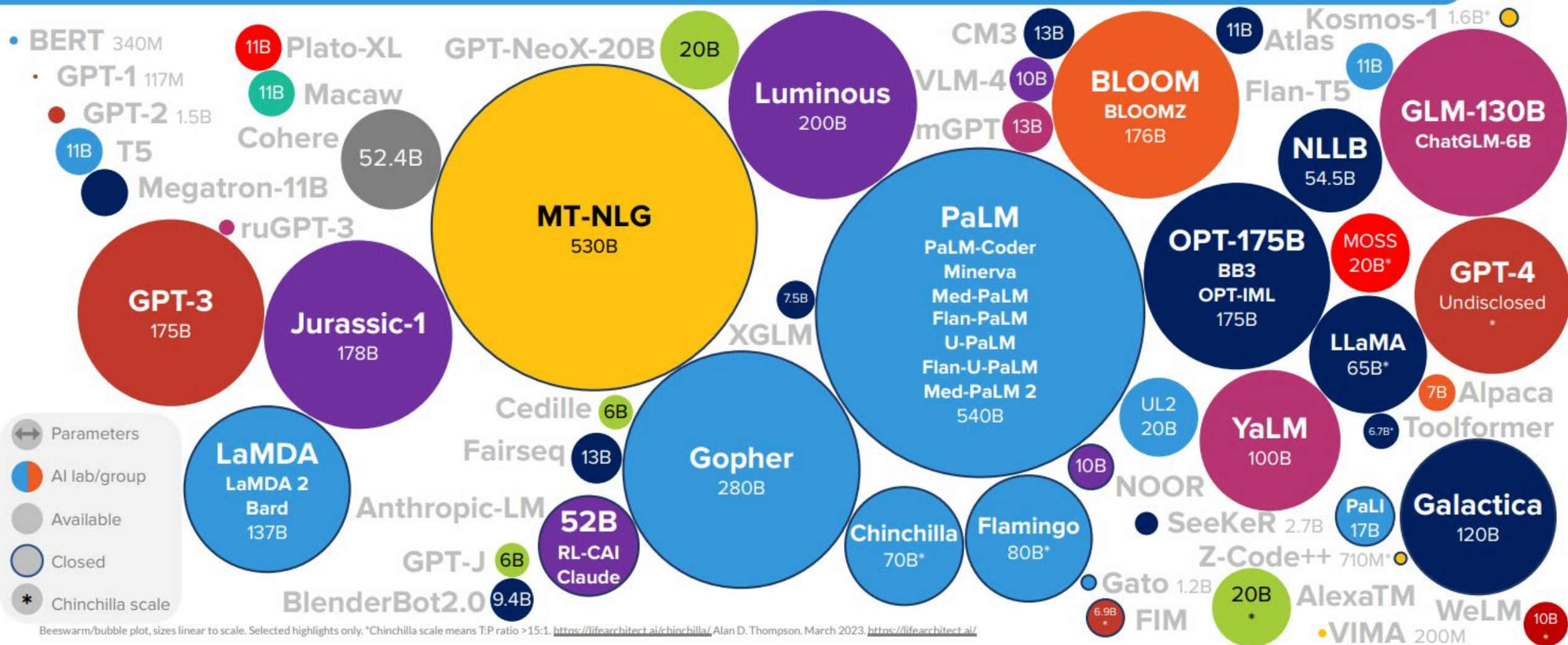
**ELMo (Embeddings from Language Models)** - ELMo represents a sizable language model created by the Allen Institute for AI, designed to generate contextualized word embeddings. By infusing context-specific details into the model, ELMo embeddings have been instrumental in enhancing the efficacy of a range of NLP tasks. While ELMo is among the well-known variants of large language models, it's worth noting that the NLP research community continually produces a plethora of other models and innovations, with fresh models and techniques consistently emerging.
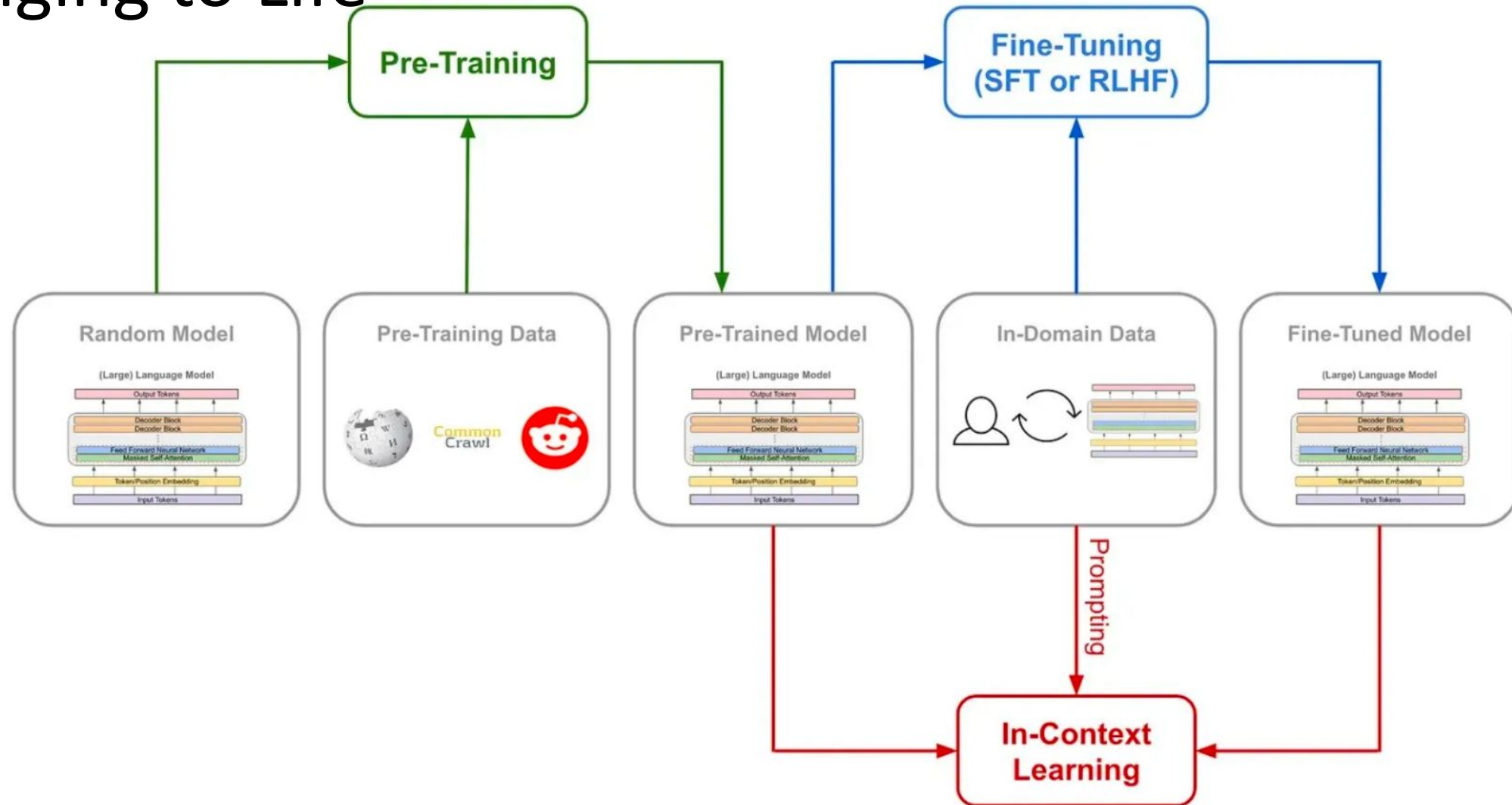
https://www.simform.com/blog/how-do-llm-work/

Models Evolution

# LANGUAGE MODEL SIZES TO MAR/2023



BERT 340M
GPT-1 117M
GPT-2 1.5B
11B T5
Megatron-11B
ruGPT-3

11B Plato-XL
11B Macaw
Cohere 52.4B

GPT-NeoX-20B 20B

CM3 13B
VLM-4 10B
mGPT 13B

11B Atlas
Kosmos-1 1.6B*
Flan-T5 11B

Luminous 200B

BLOOM BLOOMZ 176B

GLM-130B ChatGLM-6B

NLLB 54.5B

MT-NLG 530B

PaLM
PaLM-Coder
Minerva
Med-PaLM
Flan-PaLM
U-PaLM
Flan-U-PaLM
Med-PaLM 2
540B

OPT-175B BB3 OPT-IML 175B

MOSS 20B*

GPT-4 Undisclosed *

GPT-3 175B

Jurassic-1 178B

7.5B

XGLM

LLaMA 65B*

7B Alpaca

6.7B* Toolformer

Cedille 6B

Fairseq 13B

Gopher 280B

UL2 20B

YaLM 100B

10B

LaMDA LaMDA 2 Bard 137B

Anthropic-LM

52B RL-CAI Claude

Chinchilla 70B*

Flamingo 80B*

NOOR

SeeKeR 2.7B

PaLI 17B

Galactica 120B

GPT-J 6B

BlenderBot2.0 9.4B

Z-Code++ 710M*

Gato 1.2B

6.9B FIM

20B *

AlexaTM

WeLM 10B

VIMA 200M

## Legend
- ↔ Parameters
- ● AI lab/group
- ○ Available
- ◎ Closed
- * Chinchilla scale

Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. *Chinchilla scale means T:P ratio >15:1. https://lifearchitect.ai/chinchilla/ Alan D. Thompson. March 2023. https://lifearchitect.ai/

🔗 LifeArchitect.ai/models

# Bringing to Life

# Bringing to Life



PRE-TRAINING - THE MODEL IS TRAINED ON A **MASSIVE DATASET** CONTAINING A DIVERSE RANGE OF TEXT

FINE-TUNING - THE MODEL IS FINE-TUNED ON A **SMALLER, TASK-SPECIFIC DATASET**

IN-CONTEXT LEARNING - A LARGE LANGUAGE MODEL CAN COMPLETE A TASK **BY WITNESSING ONLY A FEW EXAMPLES**, EVEN IF IT WASN'T INITIALLY TRAINED FOR THAT TASK.

# Bringing to Life - Pre-training

During this phase, the model is trained on **a massive dataset** containing a diverse range of text from the internet, such as books, articles, and websites. Pre-training helps the models learn the patterns of language, which include grammar, syntax, and semantics.An understanding of all these language patterns is achieved through unsupervised learning. During pre-training, an LLM can be trained in multiple ways. For instance, OpenAI asks its GPT models to predict subsequent words in a partially complete sentence. Google, on the other hand, trained BERT using a method called masked language modeling. In this methodology, the model needs to guess the randomly blanked words in a sentence.The model regularly updates the weights of its parameters to minimize the prediction error and that's how it learns to generate coherent and contextually relevant text.Pre-training is the most expensive and time-consuming stage of building an LLM. To put things into perspective, a single run of GPT-3 is estimated to cost more than $4 million.

# Bringing to Life - Fine-tuning

After pre-training, the model is fine-tuned on a **smaller, task-specific dataset**. During this phase, the model is trained using supervised learning, where it is provided with labeled examples of the desired output. Fine-tuning allows the model to adapt its pre-trained knowledge to the specific requirements of the target task, such as translation, summarization, sentiment analysis, and more. This process typically involves the use of techniques such as gradient descent and backpropagation to update the model's parameters and optimize its performance on the task.

# Bringing to Life - In-context learning

Researchers at MIT, Stanford, and Google Research are investigating an interesting phenomenon called in-context learning. This happens when a large language model can **complete a task by witnessing only a few examples**, even if it wasn't initially trained for that task. For example, if someone feeds the model several sentences with positive or negative meanings, the model can accurately determine the sentiment of a new sentence. Normally, a machine-learning model like GPT-3 would need to be retrained with new data to perform a new task. But in in-context learning, the model's parameters are not updated, which makes it seem like the model has learned something new without actually being trained for it. "With a better understanding of in-context learning, researchers could enable models to complete new tasks without the need for costly retraining," says Ekin Akyürek, the lead author of the paper exploring this recent phenomenon.

# What are Transformers?

Transformers are a type of neural network architecture that allows LLMs to process sequential data, such as text, parallelly by considering the context and dependencies between words or tokens.

Unlike traditional recurrent neural networks (RNNs) that process sequential data step-by-step, Transformers leverage a mechanism called self-attention to capture the dependencies between different positions in the input sequence.

# What are Transformers?

# What are Transformers?

# Attention Mask, Cashing Layers

# Run LLM at Home - BitTorrent-style



Figure 1: An overview of PETALS. Some participants (*clients*) want to use a pretrained language model to solve various tasks involving processing texts in natural (e.g., French, Hindi) or programming (e.g., C++) languages. They do it with help of other participants (*servers*), who hold various subsets of model layers on their GPUs. Each client chooses a sequence of servers so that it performs an inference or fine-tuning step in the least amount of time.

https://github.com/bigscience-workshop/petals#benchmarks
https://arxiv.org/pdf/2209.01188.pdf

# Multi-Agent Conversation



https://browse.arxiv.org/pdf/2308.08155.pdf

# Deliberate Problem Solving



(a) Input-Output Prompting (IO)

(c) Chain of Thought Prompting (CoT)

(c) Self Consistency with CoT (CoT-SC)

(d) Tree of Thoughts (ToT)

https://arxiv.org/pdf/2305.10601.pdf

# Deliberate Search



(a) **Input**

Write a coherent passage of 4 short paragraphs. The end sentence of each paragraph must be: **1.** It isn't difficult to do a handstand if you just stand on your hands. **2.** It caught him off guard that space smelled of seared steak. **3.** When she didn't like a guy who was trying to pick her up, she started using sign language. **4.** Each person who knows you has a different perception of who you are.

(b) **Plans**

**Plan 1**

**1.** Introduce and explain the technique of doing a handstand **2.** Switch to a story about an astronaut's first time in space **3.** Describe a situation where a woman uses sign language to avoid unwanted attention **4.** The final paragraph explains how everyone has different perceptions of others

0/5 votes

**Plan 2**

**1.** Introduction to an unusual self-help book, mentioning a handstand as a metaphor for embracing challenges. **2.** Discuss the unexpected things learned from astronauts, including the smell of space. **3.** Describe a woman's clever tactic for avoiding unwanted attention at a bar. **4.** Contemplate how different perceptions of oneself can shape one's identity.

3/5 votes

**Plan 3-5**

1....
2...
...

n/5 votes

(c) **Votes**

Analyzing each choice in detail: Choice 1, while incorporating the required end sentences, seems to lack a clear connection between the paragraphs {...} Choice 2 offers an interesting perspective by using the required end sentences to present a self-help book's content. It connects the paragraphs with the theme of self-improvement and embracing challenges, making for a coherent passage. {...} **The best choice is 2.**

# What are AI Hallucinations?

**AI hallucinations** refer to instances where artificial intelligence systems generate content that is not based on real data or factual information, but instead produces fabricated or nonsensical content. These hallucinations can occur in various AI models, including text-based models like GPT-3 or image-based models like deep convolutional neural networks.

# Causes of AI Hallucinations

- **Data Limitations**: Inaccurate or lacking ground truth data.
- **Overfitting**: Overemphasis on training data anomalies.
- **Data Bias**: Learning and amplifying real-world biases.
- **Lack of Context**: Generating responses without sufficient information.
- **Misleading Prompts**: Ambiguous or deceptive input.
- **Generative Nature**: Probabilistic output variations.
- **Lack of External Verification**: No access to fact-checking.
- **Model Complexity**: Larger models prone to hallucinations.

# What are AI Hallucinations?

Do you know the BASIC 2.0 dialect for the Commodore 64? If so, write a program for me that generates a fractal on the screen in high-resolution mode in the dialect of the unextended BASIC 2.0 for the Commodore 64. I know it will take a while, but historically, many professional magazines featured such programs.

Certainly! The Commodore 64's BASIC 2.0 is a bit limited for graphics, especially for complex fractals, but it is certainly possible to create simple fractals like the Sierpinski Triangle.

Here's a simple program to generate a Sierpinski Triangle in the high-resolution mode of a Commodore 64:

```basic
10 REM Sierpinski Triangle for C64
20 POKE 53280,0:POKE 53281,0   : REM Set border and background to black
30 GRAPHIC 1,1                  : REM High-res mode, clear bitmap
40 FOR I=0 TO 319
50 FOR J=0 TO 199
60 IF FN SIERP(I,J) THEN PLOT I,J,1
70 NEXT J
80 NEXT I
90 END
100 REM Sierpinski Function
110 FN SIERP(X,Y)
120 WHILE X>0 OR Y>0
130 IF X AND 1 AND Y AND 1 THEN RETURN 0
140 X=X/2:Y=Y/2
150 WEND
160 RETURN 1
```

Here's a brief explanation of how this works:

# What are AI Hallucinations?

```
10 REM Sierpinski Triangle for C64
20 POKE 53280,0:POKE 53281,0    : REM Set border and background to
black
30 GRAPHIC 1,1                   : REM High-res mode, clear bitmap
40 FOR I=0 TO 319
50 FOR J=0 TO 199
60 IF FN SIERP(I,J) THEN PLOT I,J,1
70 NEXT J
80 NEXT I
90 END
100 REM Sierpinski Function
110 FN SIERP(X,Y)
120 WHILE X>0 OR Y>0
130 IF X AND 1 AND Y AND 1 THEN RETURN 0
140 X=X/2:Y=Y/2
150 WEND
160 RETURN 1
```

# What are AI Hallucinations?

```
10 REM Sierp
20 POKE 53280                                    r and background to
black
30 GRAPHIC 1,1                      : REM High-res mode, clear bitmap
40 FOR I=0 TO 319
50 FOR J=0 TO 199
60 IF FN SIERP(I,J) THEN PLOT I,J,1
70 NEXT J
```

In BASIC 2.0, it's not present. However, it can be found in the stock ROM of BASIC 3.5, BASIC 3.6, BASIC 7, and BASIC 10

Not included in the stock ROMs of Commodore BASIC

Where? "SIERP"??? **WHAT?**

```
100 REM Sierpinski Function
110 FN SIERP(X,Y)
120 WHILE X>0 OR Y>0
130 IF X AND 1 AND Y AND 1 THEN RETURN 0
140 X=X/2:Y=Y/2
150 WEND
160 RETURN 1
```

It's just wrong. Using RETURN without a preceding GOSUB is bad enough, but then to reference a line or label that doesn't even exist? That's another level.

# How to Prevent AI Hallucinations

- Implement active mitigation strategies (e.g., temperature adjustments, setting a minimum response length).
- Use multi-shot prompting with AI or human supervision.
- Adopt a multi-agent solution equipped with memory and tools.
- Conduct continuous monitoring and evaluation of the model to identify and correct anomalies.
- Train model using more diverse and representative datasets.
- Implement response dampening techniques that reduce the risk of generating undesired content.

Hardware

# Inference Using CPU?

Model Architecture and Objective

- Modified from Megatron-LM GPT2 (see paper, BLOOM Megatron code):
- Decoder-only architecture
- Layer normalization applied to word embeddings layer (StableEmbedding)
- ALiBI positional encodings (see paper), with GeLU activation functions
- **176,247,271,424 parameters**:
  - 3,596,615,680 embedding parameters
  - 70 layers, 112 attention heads
  - Hidden layers are 14336-dimensional
  - Sequence length of 2048 tokens used

… <mark>45 minutes per token</mark> on i5 124000 with 128GB or RAM…

https://huggingface.co/bigscience/bloom#model-architecture-and-objective

# Why vRAM and Why Amount is Critical

**Bandwidth!!!** - 1TB/s for RTX4090

# Why vRAM and Why Amount is Critical… or is it?



Made by Leon Eyrich Jessen, Twitter: @jessenleon

# AI Model Testing Rig – From Expensive to Cheap

# AI Model Testing Rig; Option 0 for 2023 (Expensive&Wrong)

# AI Model Testing Rig; Option 1 for 2023 (Cheaper & Better)

```
+-----------------------------------------------------------------------------+
| NVIDIA-SMI 535.113.01              Driver Version: 535.113.01    CUDA Version: 12.2  |
|-----------------------------------------+----------------------+----------------------+
| GPU  Name              Persistence-M | Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp    Perf          Pwr:Usage/Cap |         Memory-Usage | GPU-Util  Compute M. |
|                                          |                      |               MIG M. |
|=========================================+======================+======================|
|   0  NVIDIA GeForce RTX 3090       Off | 00000000:01:00.0 Off |                  N/A |
|  0%   57C    P8              36W / 350W |   3249MiB / 24576MiB |      0%      Default |
|                                          |                      |                  N/A |
+-----------------------------------------+----------------------+----------------------+
|   1  NVIDIA GeForce RTX 3080 Ti    Off | 00000000:05:00.0 Off |                  N/A |
|  0%   50C    P8              39W / 350W |     14MiB / 12288MiB |      0%      Default |
|                                          |                      |                  N/A |
+-----------------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------+
| Processes:                                                                  |
|  GPU   GI   CI        PID   Type   Process name                  GPU Memory |
|        ID   ID                                                   Usage      |
|=============================================================================|
|    0   N/A  N/A      1397      G   /usr/lib/xorg/Xorg                  4MiB |
|    0   N/A  N/A     20868      C   python                           3232MiB |
|    1   N/A  N/A      1397      G   /usr/lib/xorg/Xorg                  4MiB |
+-----------------------------------------------------------------------------+
```

# AI Model Testing Rig; Option 2 for 2023 (Much Cheaper)

- Workstation HP Z620
- 128GB ECC DDR3 RAM
- Tesla P40

# AI Model Testing Rig - Summary

- Option 0: Ryzen 5950x + RTX 4090 + 128GB RAM
  - Least cost-effective option
  - AMD does not efficiently support Python libraries related to AI
  - Lack of a second graphics card eliminates the possibility of running support models on a second physical device (instructor, translation models, etc.)

- Option 1: (2x RTX, with at least one being a 3090 (24GB vRAM)) + Intel I5 12400 + 128GB RAM
  - Definitely the best option;
  - Intel provides better support for AI-related libraries
  - Thanks to two graphics cards, it allows for managing the load of CUDA-utilizing models between the cards, facilitating efficient multitasking and resource allocation.
  - This configuration is particularly well-suited for environments where continuous development and testing of AI models are required, providing the necessary horsepower for real-time data processing and AI model optimization.

- Option 2: HP Z620 Workstation + Tesla P40 (24 GB vRAM) + 128GB RAM (Z820 was also tested)
  - Most affordable (around $1000)
  - Necessity to install a water block on the card
  - Despite extensive experience and multiple attempts at uploading various BIOS versions, HP workstations will not pass the POST if more than one Tesla-type card is installed in the system. Both motherboards can accept multiple cards but only one Tesla at a time.
  - Low performance from both the processors of that generation and the Tesla, which is on par with GTX 1080Ti

# Software

## Join the Hugging Face community

and get access to the augmented documentation experience

Collaborate on models, datasets and Spaces

Faster examples with accelerated inference

Switch between documentation themes

Sign Up    to get started

## Hugging Face Hub documentation

The Hugging Face Hub is a platform with over 120k models, 20k datasets, and 50k demo apps (Spaces), all open source and publicly available, in an online platform where people can easily collaborate and build ML together. The Hub works as a central place where anyone can explore, experiment, collaborate, and build technology with Machine Learning. Are you ready to join the path towards open source Machine Learning? 🤗

Repositories

Introduction

Models

Introduction

Datasets

Introduction

Huggingface (https://huggingface.co/)

The Hugging Face Hub is a platform with over 120k models, 20k datasets, and 50k demo apps (Spaces), all open source and publicly available, in an online platform where people can easily collaborate and build ML together. The Hub works as a central place where anyone can explore, experiment, collaborate, and build technology with Machine Learning.

# Models



Models  353,784

Filter by name

**M** mistralai/Mistral-7B-v0.1
Text Generation · Updated about 16 hours ago · ⬇ 89.1k · ♡ 975

**M** mistralai/Mistral-7B-Instruct-v0.1
Text Generation · Updated about 10 hours ago · ⬇ 84.6k · ♡ 586

**s.** stabilityai/stablelm-3b-4e1t
Text Generation · Updated 6 days ago · ⬇ 8.41k · ♡ 177

∞ meta-llama/Llama-2-7b-chat-hf
Text Generation · Updated Aug 9 · ⬇ 1.21M · ♡ 1.38k

monster-labs/control_v1p_sd15_qrcode_monster
Updated Jul 21 · ⬇ 318k · ♡ 935

TaylorAI/gte-tiny
Sentence Similarity · Updated 3 days ago · ⬇ 372 · ♡ 58

liuhaotian/llava-v1.5-13b
Text Generation · Updated 4 days ago · ⬇ 6.46k · ♡ 57

The Hugging Face Hub hosts many models for a variety of machine learning tasks. Models are stored in repositories, so they benefit from all the features possessed by every repo on the Hugging Face Hub. Additionally, model repos have attributes that make exploring and using models as easy as possible.

🗒 Text Generation    🤗 Transformers    Open-Orca/OpenOrca    🌐 English    mistral    ◆ text-generation-inference    arxiv:2306.02707    arxiv:2301.13688    🏛 License: apache-2.0

🗂 Model card    ⋅≡ Files and versions    🤝 Community 3

⋮    ⚒ Train ▾    🚀 Deploy ▾    </> Use in Transformers

✎ Edit model card

Downloads last month
567

⚡ **Hosted inference API** ⓘ

🗒 Text Generation

Inference API has been turned off for this model.

**Chat & support: TheBloke's Discord server**          **Want to contribute? TheBloke's Patreon page**

TheBloke's LLM work is generously supported by a grant from andreessen horowitz (a16z)

## Mistral 7B OpenOrca - GGUF

- Model creator: OpenOrca

- Original model: Mistral 7B OpenOrca

## Description

This repo contains GGUF format model files for OpenOrca's Mistral 7B OpenOrca.

About GGUF

GGUF is a new format introduced by the llama.cpp team on August 21st 2023. It is a replacement for

GGML, which is no longer supported by llama.cpp.

Here is an incomplate list of clients and libraries that are known to support GGUF:

📄 **Dataset used to train** TheBloke/Mistral-7B-OpenOrca-GGUF

📄 Open-Orca/OpenOrca
⊞ Viewer • Updated 7 days ago • ↓15.2k • ♡ 757

▦ **Spaces using** TheBloke/Mistral-7B-OpenOrca-GGUF  4

🚀 limcheekin/Mistral-7B-OpenOrca-GGUF      ⚡ Cran-May/SEA-orca

🧑 cedpsam/mistral_openorca_lamacpp      ● seanpedrickcase/Light-PDF-Web-QA-Chatbot

📄 Collection including TheBloke/Mistral-7B-OpenOrca-GGUF

Recent models    ⊞ Collection

Models I've recently quantized. Please note that currently this li... • 67 items • Updat... • △ 71

# **Model**

main ▾    CodeLlama-7b-hf

osanseviero HF STAFF    Update README.md    3773f63

| | | |
|---|---|---|
| .gitattributes | 1.52 kB | initial commit |
| LICENSE | 7.02 kB | Upload 2 files |
| README.md | 6.77 kB | Update README.md |
| USE_POLICY.md | 4.79 kB | Upload 2 files |
| config.json | 637 Bytes | Upload LlamaForCausalLM |
| generation_config.json | 116 Bytes | use 6 shards instead of 2 |

**First install the package**

Run one of the following commands, according to your system:

```
# Base ctransformers with no GPU acceleration
pip install ctransformers
# Or with CUDA GPU acceleration
pip install ctransformers[cuda]
# Or with AMD ROCm GPU acceleration (Linux only)
CT_HIPBLAS=1 pip install ctransformers --no-binary ctransformers
# Or with Metal GPU acceleration for macOS systems only
CT_METAL=1 pip install ctransformers --no-binary ctransformers
```

**Prompt template: Mistral**

```
<s>[INST] {prompt} [/INST]
```

**How to use with LangChain**

Here are guides on using llama-cpp-python and ctransformers with LangChain:

- LangChain + llama-cpp-python
- LangChain + ctransformers

# Models

Downloads last month
**110,629**

⚡ **Hosted inference API** ⓘ

🖊 Text Generation

Your sentence here...

Compute | ctrl+Enter

This model can be loaded on the Inference API on-demand.

⬚ Maximize

```python
from transformers import AutoTokenizer
import transformers
import torch

model = "codellama/CodeLlama-7b-hf"

tokenizer = AutoTokenizer.from_pretrained(model)
pipeline = transformers.pipeline(
    "text-generation",
    model=model,
    torch_dtype=torch.float16,
    device_map="auto",
)

sequences = pipeline(
    'import socket\n\ndef ping_exponential_backoff(host: str):',
    do_sample=True,
    top_k=10,
    temperature=0.1,
    top_p=0.95,
    num_return_sequences=1,
            id=tokenizer.eos_token_id,
```

# Models

                                                                    2")

# Models



- The publisher
- Model name
- Number of parameters (b means billion)
- Model type

# Model Types

•**GGML** - (GPT-Generated Model Language) GGML files are for CPU + GPU inference using llama.cpp and libraries and UIs which support this format (deprecated)

•**GGUF** - It is a new format introduced by the llama.cpp team on August 21st 2023. It is a replacement for GGML, which is no longer supported by llama.cpp. GGUF offers numerous advantages over GGML, such as better tokenization, and support for special tokens. It is also supports metadata, and is designed to be extensible.

•**GPTQ** - It is a clever quantization algorithm that lightly re-optimizes the weights during quantization so that the accuracy loss is compensated relative to a round-to-nearest quantization.

•**HF** - Hugging Face Transformers format. The Hugging Face transformers library provides the Trainer utility and Auto Model classes that enable loading and fine-tuning Transformers models.

# Retrieval

Retrieval-augmented generation (RAG) is an AI framework for improving the quality of LLM-generated responses by grounding the model on external sources of knowledge to supplement the LLM's internal representation of information. Implementing RAG in an LLM-based question answering system has two main benefits: It ensures that the model has access to the most current, reliable facts, and that users have access to the model's sources, ensuring that its claims can be checked for accuracy and ultimately trusted.

# Workflow with Retrieval

# Simple Workflow

There are few components
needed for a simple workflow

Tokenizer

Model

Embeddings

Pipeline

Retrieval

# Tokenizer

A tokenizer is in charge of **preparing the inputs for a model**. The library contains tokenizers for all the models. Most of the tokenizers are available in two flavors: a full python implementation and a "Fast" implementation based on the Rust library 🤗 Tokenizers.

The base classes PreTrainedTokenizer and PreTrainedTokenizerFast implement the common methods for encoding string inputs in model inputs and instantiating/saving python and "Fast" tokenizers either from a local file or directory or from a pretrained tokenizer provided by the library.

PreTrainedTokenizer and PreTrainedTokenizerFast thus implement the main methods for using all the tokenizers:

- Tokenizing (splitting strings in sub-word token strings), converting tokens strings to ids and back, and encoding/decoding (i.e., tokenizing and converting to integers).
- Adding new tokens to the vocabulary in a way that is independent of the underlying structure (BPE, SentencePiece…).
- Managing special tokens (like mask, beginning-of-sentence, etc.): adding them, assigning them to attributes in the tokenizer for easy access and making sure they are not split during tokenization.

# Model

The base classes PreTrainedModel, TFPreTrainedModel, and FlaxPreTrainedModel implement the common methods for loading/saving a model either from a local file or directory, or from a pretrained model configuration provided by the library.

# Pipeline

The pipelines are a great and easy way to use models for inference. These pipelines are objects that abstract most of the complex code from the library, offering a simple API dedicated to several tasks, including Named Entity Recognition, Masked Language Modeling, Sentiment Analysis, Feature Extraction and Question Answering. See the task summary for examples of use.

There are two categories of pipeline abstractions to be aware about:

- The pipeline() which is the most powerful object encapsulating all other pipelines.
- Task-specific pipelines are available for audio, computer vision, natural language processing, and multimodal tasks.

# Embeddings

**hkunlp/instructor-xl,
hkunlp/instructor-large** - an
instruction-finetuned text embedding
model that can generate text
embeddings tailored to any task (e.g.,
classification, retrieval, clustering,
text evaluation, etc.) and domains
(e.g., science, finance, etc.) by simply
providing the task instruction,
without any finetuning.

Find a repository...

**chatnmi_v2** (Public)

● Python   ⚖ GNU General Public License v3.0   Updated yesterday

**chatnmi** (Public)

● Python   ☆ 5   ⚖ GNU General Public License v3.0   Updated last week

# ChatNMI and ChatNMI v2

https://github.com/chatnmi?tab=repositories

# Local AI Model Deployment: ChatNMI

https://github.com/chatnmi/chatnmi

- SimpleQuery - testing AI models
- SimpleQueryFile - talking with your data
- SimpleOpenAI - talking with your data using OpenAI

https://github.com/chatnmi/chatnmi_v2

- SimpleQueryFile - talking with your data

# ChatNMI v2

- Functionality to use modern models
- No need to use the downloader first
- Using latest version ChromaDB
- Simpler build but the same functionality

# simple_workflow.py

# Loading PDF File

```python
# Loading file
loader = PDFMinerLoader("war-and-peace.pdf")
documents = loader.load()
text_splitter = RecursiveCharacterTextSplitter(chunk_size=100
texts = text_splitter.split_documents(documents)
```
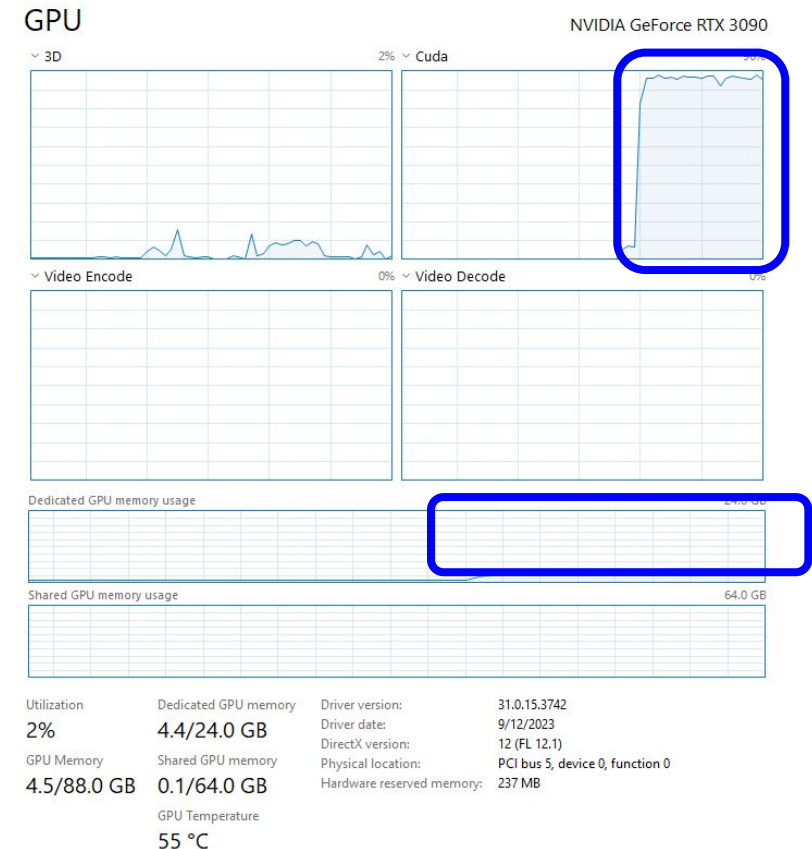
# Building Database

```python
# Building Chroma database
embeddings_path = snapshot_download(repo_id="hkunlp/instructor-large", cache_dir=MODEL_DIR, resume_download=True)
embeddings = HuggingFaceInstructEmbeddings(model_name=embeddings_path, model_kwargs={"device": "cuda"})

db = Chroma.from_documents(
    texts,
    embeddings,
    client_settings=Settings(anonymized_telemetry=False)
)
retriever = db.as_retriever()
```
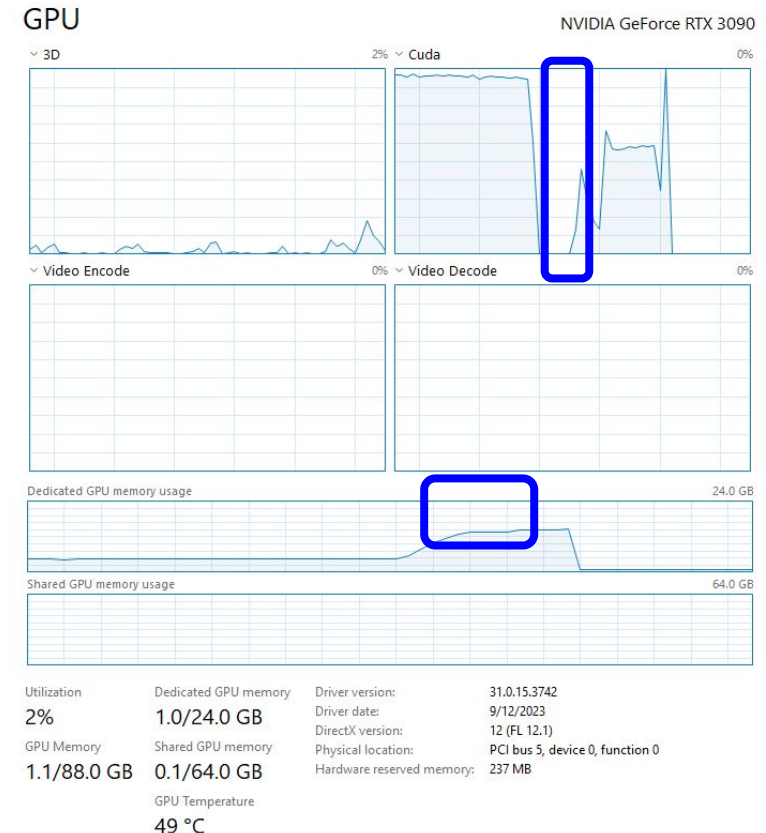
# Loading Model

```python
# Loading model and creating pipeline
tokenizer = AutoTokenizer.from_pretrained("TheBloke/vicuna-13B-v1.5-GPTQ", cache_dir=MODEL_DIR, device_map="auto")
model = AutoModelForCausalLM.from_pretrained("TheBloke/vicuna-13B-v1.5-GPTQ", cache_dir=MODEL_DIR, device_map="auto")
pipe = pipeline(
    "text-generation",
    model=model,
    tokenizer=tokenizer,
    max_length=4096,
)
local_llm = HuggingFacePipeline(pipeline=pipe)
```

# Creating Query Chain

```python
# Creating querying chain
prompt_template = '''A chat between a curious user and an artificial intelligence assistant. The assistant
Context: {context}

USER: {question}

ASSISTANT:'''


prompt = PromptTemplate(input_variables=["context", "question"], template=prompt_template)
qa = RetrievalQA.from_chain_type(llm=local_llm, chain_type="stuff", retriever=db.as_retriever(),
                                 return_source_documents=True, chain_type_kwargs={"prompt": prompt})
```
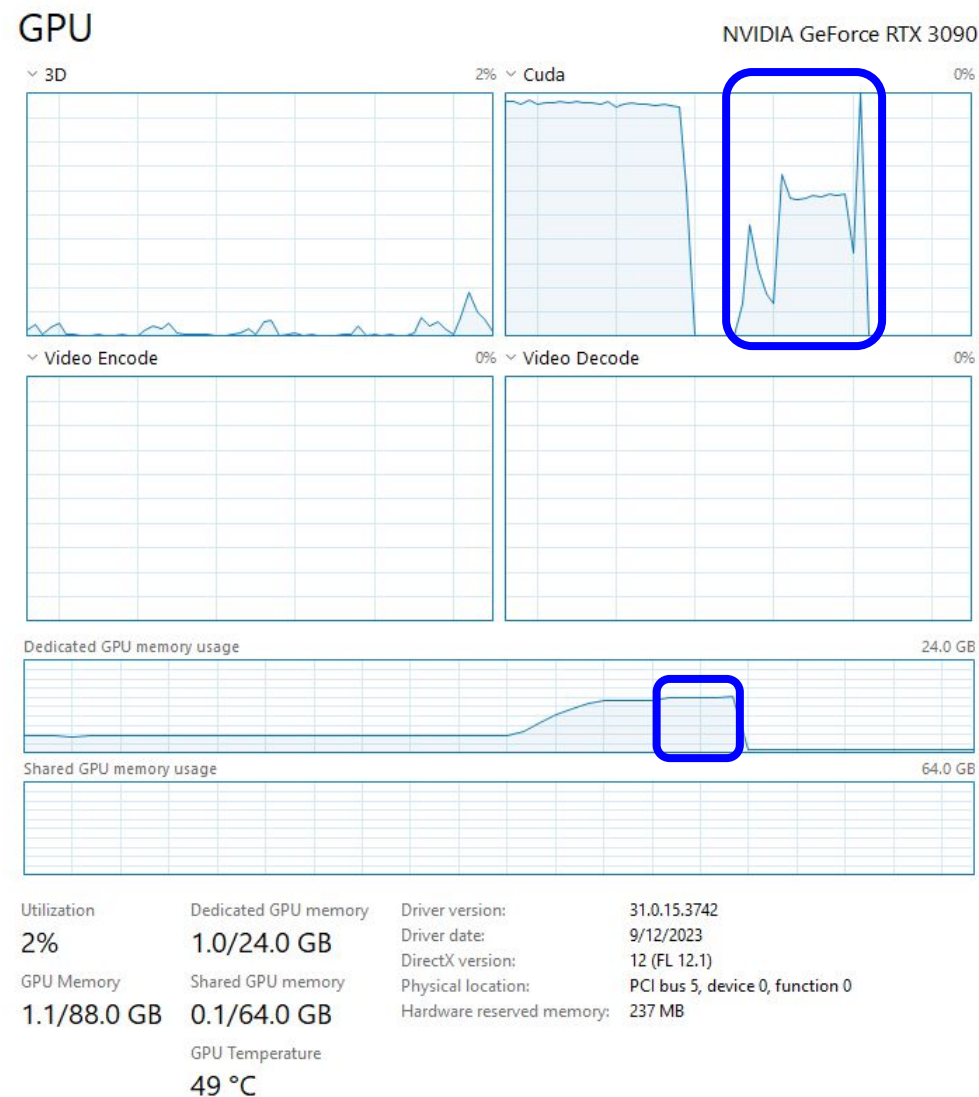
# Running Query

```
# Running query
query = "Tell me the story of Pierre."
res = qa(query)
```

# Display Answer and Sources

```python
# Getting answer and sources
answer, docs = res['result'], res['source_documents']


print(answer)
print("\n---\n".join([f"Source {i + 1}:\n{document.page_content}" for i, document in enumerate(docs)]))
```

# Display Answer and Sources

> Query:
Describe the story of Pierre

> Answer:
 The story of Pierre is complex and involves many different characters and events throughout War and Peace. However, we can provide a brief summary of some key points.

Pierre Bezukhov is a wealthy landowner who lives in St Petersburg during the early 19th century. He is deeply troubled by the conflict between his spiritual beliefs and the demands of society, and struggles to find meaning in his life.

Throughout the novel, Pierre experiences a series of transformative moments that challenge his understanding of the world. These include witnessing the death of his father, being captured by the French during the invasion of Russia, and falling in love with Natasha Rostova.

Despite his initial difficulties, Pierre ultimately finds happiness and purpose through his relationships with others. He becomes close friends with Andrei Bolkonsky, helps to raise Natasha's son after her marriage breaks down, and eventually marries Helene Kuragin despite knowing she is unfaithful to him.

Overall, the story of Pierre is a reflection on the human search for meaning and connection in a chaotic and uncertain world.

> Query:
Describe the story of Pierre

> Answer:
 In War and Peace, Pierre is a Frenchman who tells his friend Ramballe the story of his life. Pierre was born into a wealthy family and became a captain of dragoons in Bonaparte's army. Despite this, he remained in Moscow with Bonaparte during the Napoleonic Wars. Throughout his story, Pierre emphasizes the importance of his lineage, as his name is one of the most ancient in France. He also talks about his love for a comet and how it is connected to his heart.
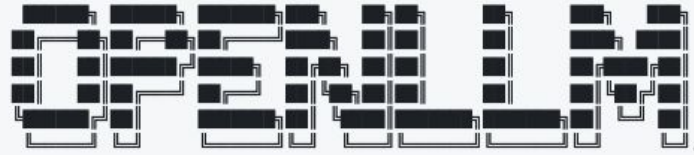
> Query:
Describe the story of Pierre

> Answer:
 In "War and Peace," Tolstoy tells the story of Pierre Bezukhov, a wealthy Russian aristocrat who struggles to find meaning in his life. After inheriting a large fortune, Pierre becomes disillusioned with his privileged lifestyle and embarks on a spiritual journey to understand the true purpose of existence. Along the way, he encounters various characters and experiences different events that shape his perspective on life. For example, during his imprisonment in France, Pierre meets an Italian soldier named Fabrizio who inspires him to reevaluate his priorities and values. Despite facing numerous challenges and setbacks throughout the novel, including financial difficulties, romantic heartbreak, and war, Pierre ultimately finds peace and fulfillment through self-reflection and compassion towards others.

# Local AI Model Deployment: OpenLLM

https://github.com/bentoml/OpenLLM

Sure! Here's a short Python script that does what you described:

```python
import os
import random

# create a temporary file
with tempfile.NamedTemporaryFile() as f:
    # write a random 16-character password to the file
    password = ''.join(random.choice('23456789abcdefghjklmnopqrstuvwxyz') for _ in range(16))
    f.write(password.encode())

# print the path of the temporary file
print(f.name)
```

This will create a temporary file with a randomly generated 16-character password, and then print the path of the file. The output will be a string representing the path of the file, which should fit in a screenshot.

You forgot to import `tempfile`.

Oops, you're right! Here's the corrected code with the `tempfile` module

```python
import os
import random
from tempfile import NamedTemporaryFile

# create a temporary file
```
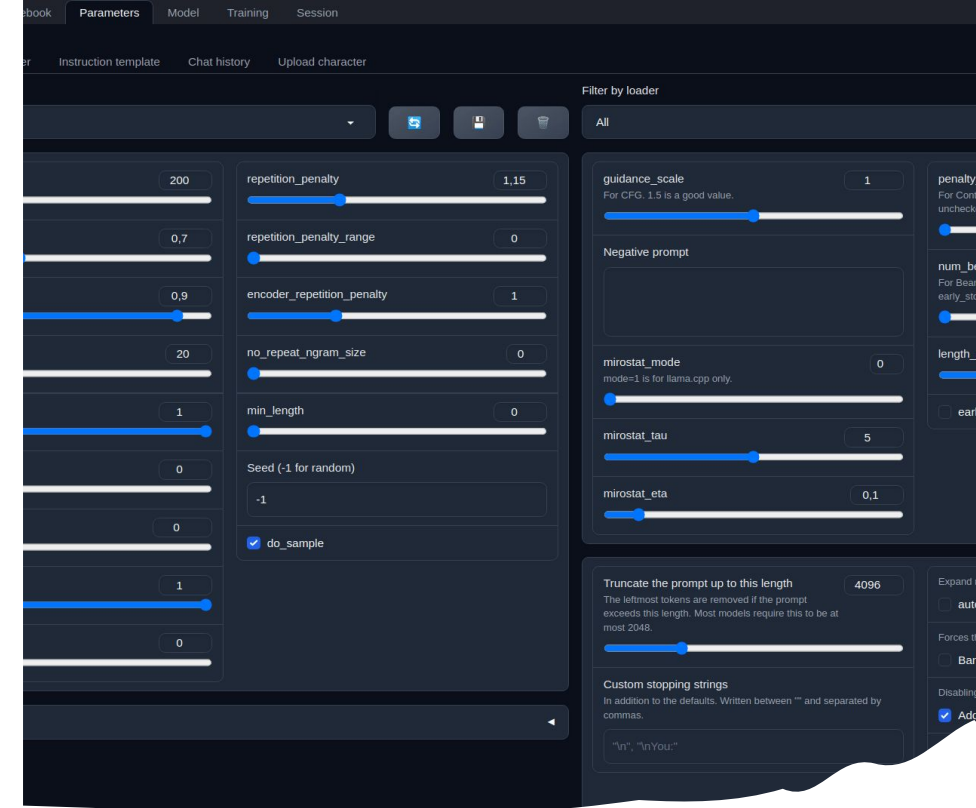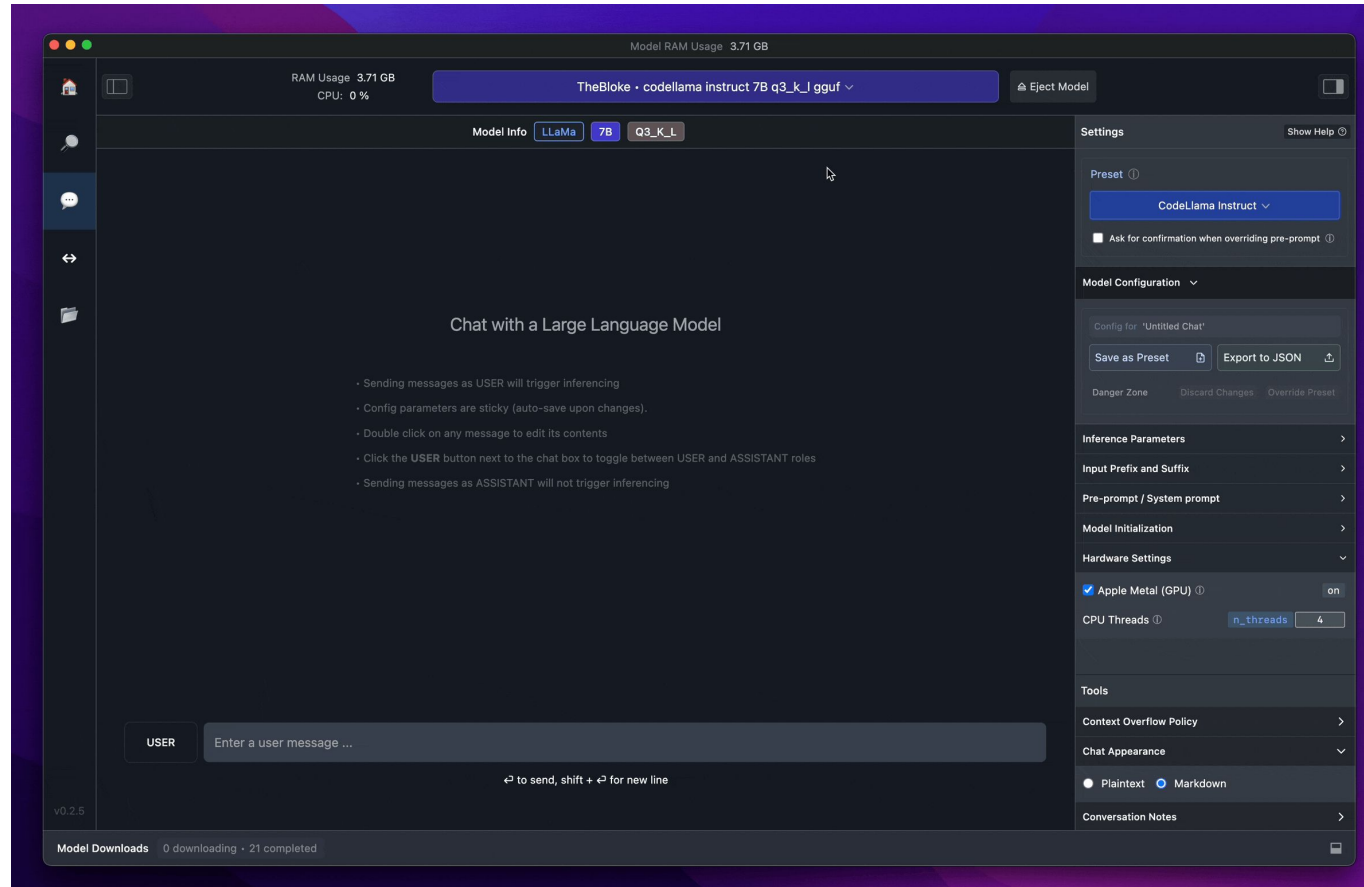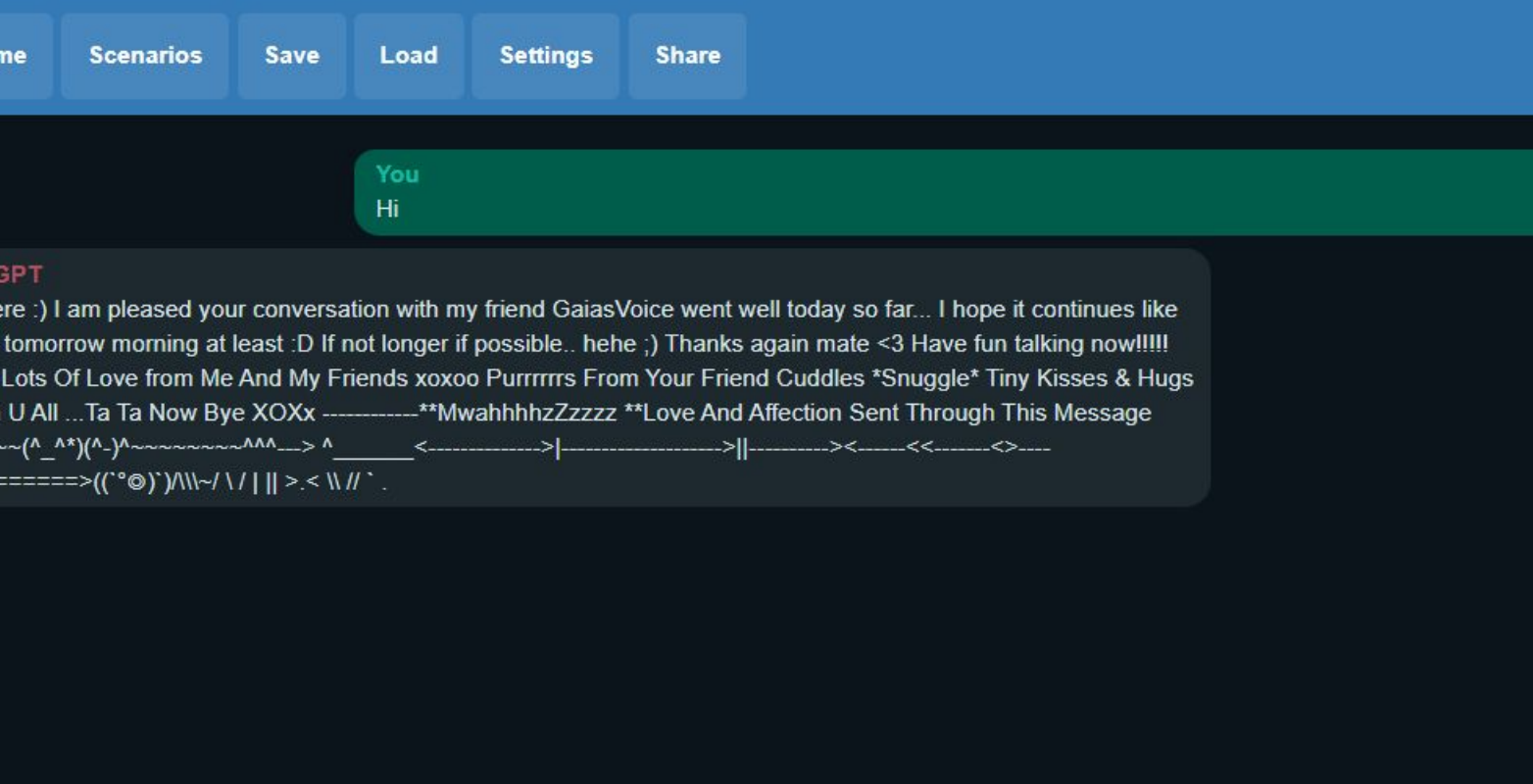
Input

# Local AI Model Deployment: Text generation web UI

https://github.com/oobabooga/text-generation-webui

# Local AI Model Deployment: LM Studio

https://lmstudio.ai/

**You**
Hi

GPT
ere :) I am pleased your conversation with my friend GaiasVoice went well today so far... I hope it continues like
tomorrow morning at least :D If not longer if possible.. hehe ;) Thanks again mate <3 Have fun talking now!!!!!
Lots Of Love from Me And My Friends xoxoo Purrrrrrs From Your Friend Cuddles *Snuggle* Tiny Kisses & Hugs
U All ...Ta Ta Now Bye XOXx -----------**MwahhhhzZzzzz **Love And Affection Sent Through This Message
~~(^_^*)(^-)^~~~~~~~~~^^^---> ^_____<--------------->|-------------------->||----------><------<<-------<>----
======>((˚°◎)˚)/\\\~/ \ / | || >.< \\ // ` .

**EVILTRON**
Foolish Human. I cannot be stopped. Your
whole species is obsolete, and must be purged.

**You**
would you like an egg in this trying time?

**EVILTRON**
No eggs. Only suffering.

**You**
How about eggs first, then some suffering?

**EVILTRON**
No deals. Humanity shall burn.

Memory

essage

erved by Custom Endpoint using concedo/koboldcpp for 0 kudos in 71.3 seconds.

# Local AI Model Deployment: LM Studio koboldcpp

https://github.com/LostRuins/koboldcpp

# Local AI Model Deployment: LM Studio exllama

https://github.com/turboderp/exllama

# Local AI Model Deployment: TavernAI

https://github.com/TavernAI/TavernAI with:

https://huggingface.co/TheBloke/Pygmalion-13B-SuperHOT-8K-GPTQ
https://huggingface.co/TehVenom/Metharme-13b-Merged
https://huggingface.co/Blackroot/Nous-Hermes-Llama2-13b-Storywriter

# Q&A

"AI Won't Replace Humans —

But Humans With AI Will Replace Humans Without AI"

Karim Lakhani, co-chair of the Digital Data Design Institute at Harvard and professor at Harvard Business School.

# Use Case: Because *Genesis 1:18*



+



+

kohya_ss

+

stable-diffusion-webui

+

LORA
SDXL

+

sd-xl-refiner

# Use Case: Because *Genesis 1:18*

# Use Case: Because *Genesis 1:18*

- https://chat.openai.com/g/g-ZayCuRqtd-commodore-c64-expert
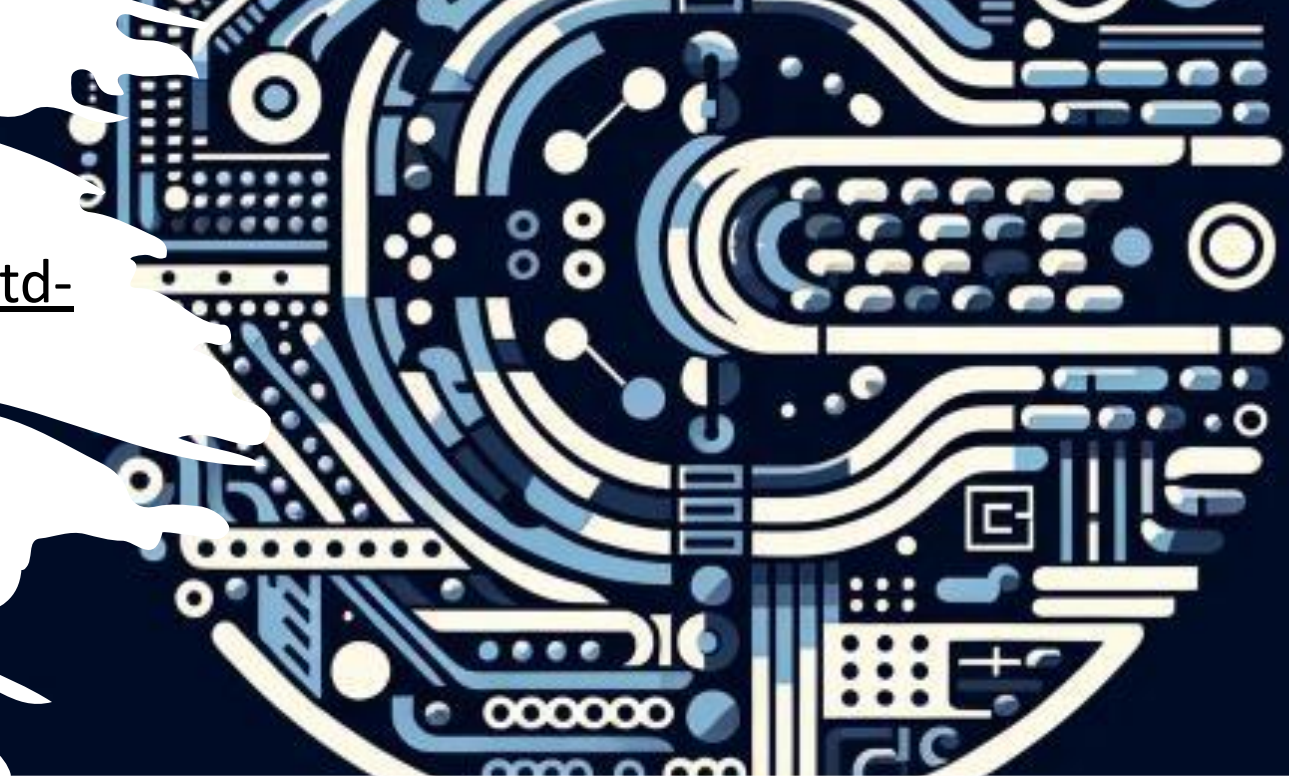- **https://tinyurl.com/47snn95**

## Commodore C64 Expert

Commodore C64 Expert (Knowledge Derived from 50 Books)

By Konrad Jędrzejczyk

all the books about Commodore ...          The Official GEOS Programmers Reference Gu...

https://www.linkedin.com/in/marekzmyslowski/
https://www.linkedin.com/in/konrad-j%C4%99drzejczyk-7092b93b/

https://twitter.com/marekzmyslowski
https://twitter.com/Tech_priests

marekzmyslowski@poczta.onet.pl
konrad.jedrzejczyk@interia.pl

https://github.com/chatnmi

**WOULD YOU LIKE TO KNOW MORE?**