# 14

# BEST

# LARGE LANGUAGE MODELS (LLM)

# IN 2024

Swipe →

# 1. GPT-4 O (OMNI)

OpenAI is launching GPT-4o, an iteration of the GPT-4 model that powers its hallmark product, ChatGPT. The updated model "is much faster" and improves "capabilities across text, vision, and audio," OpenAI CTO Mira Murati said in a livestream announcement. It'll be free for all users, and paid users will continue to "have up to five times the capacity limits" of free users, Murati added.

```python
from datetime import datetime
import matplotlib.pyplot as plt
from meteostat import Point, Daily


def foo(x, y):
    x['Average Temperature'] = x['Average Temperature'].rolling(window=y).mean()
    x['Minimum Temperature'] = x['Minimum Temperature'].rolling(window=y).mean()
    x['Maximum Temperature'] = x['Maximum Temperature'].rolling(window=y).mean()
    return x


def bar(ax, events):
    for date, label in events.items():
        ax.annotate(label, xy=(date, data.loc[date, 'Maximum Temperature']),
                    xytext=(date, data.loc[date, 'Maximum Temperature'] + 5),
                    arrowprops=dict(facecolor='black', arrowstyle='->'))


start = datetime(2018, 1, 1)
end = datetime(2018, 12, 31)

location = Point(49.2497, -123.1193, 70)

data = Daily(location, start, end)
data = data.fetch()

data = data.rename(columns={
    'tavg': 'Average Temperature',
```
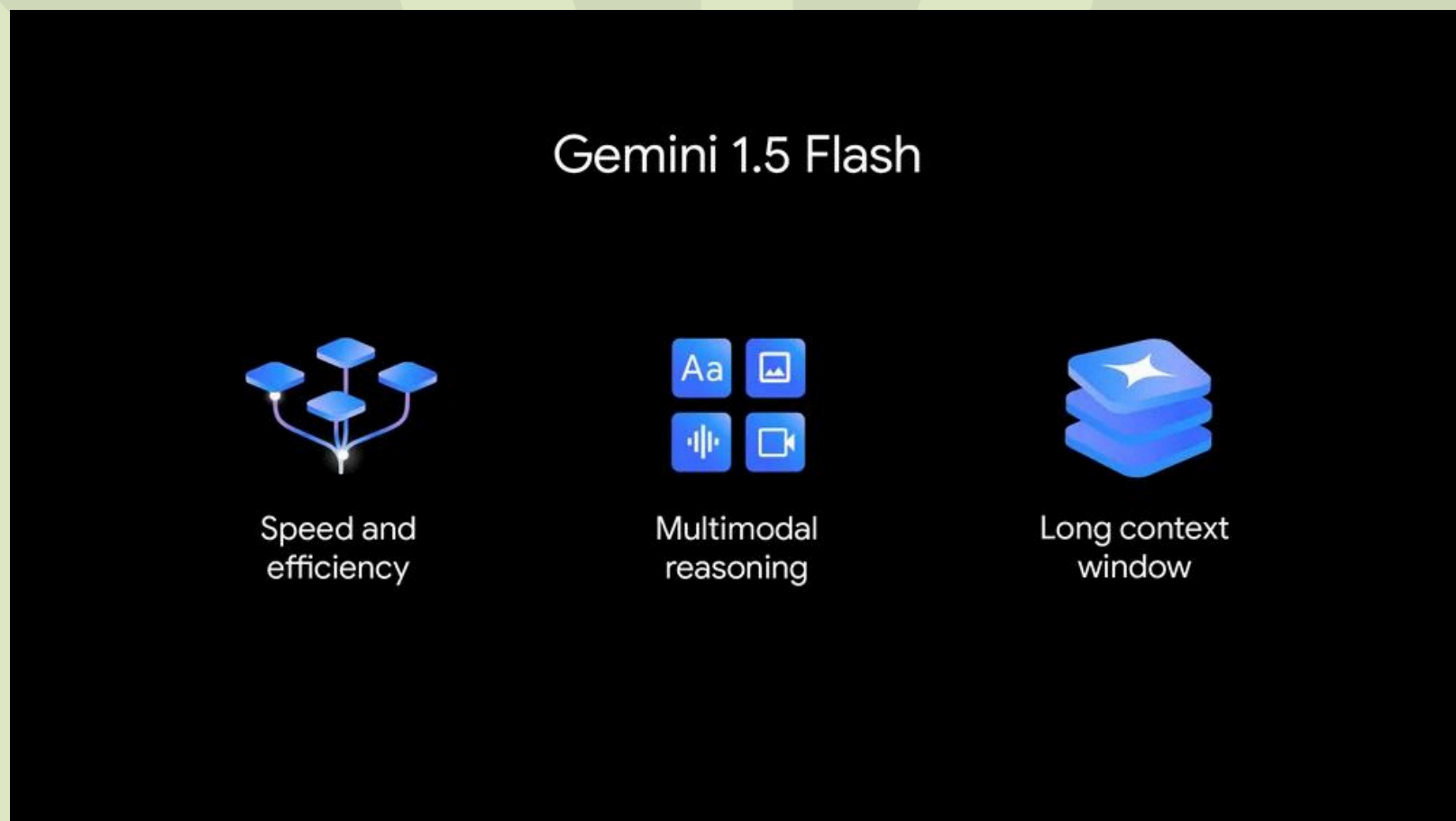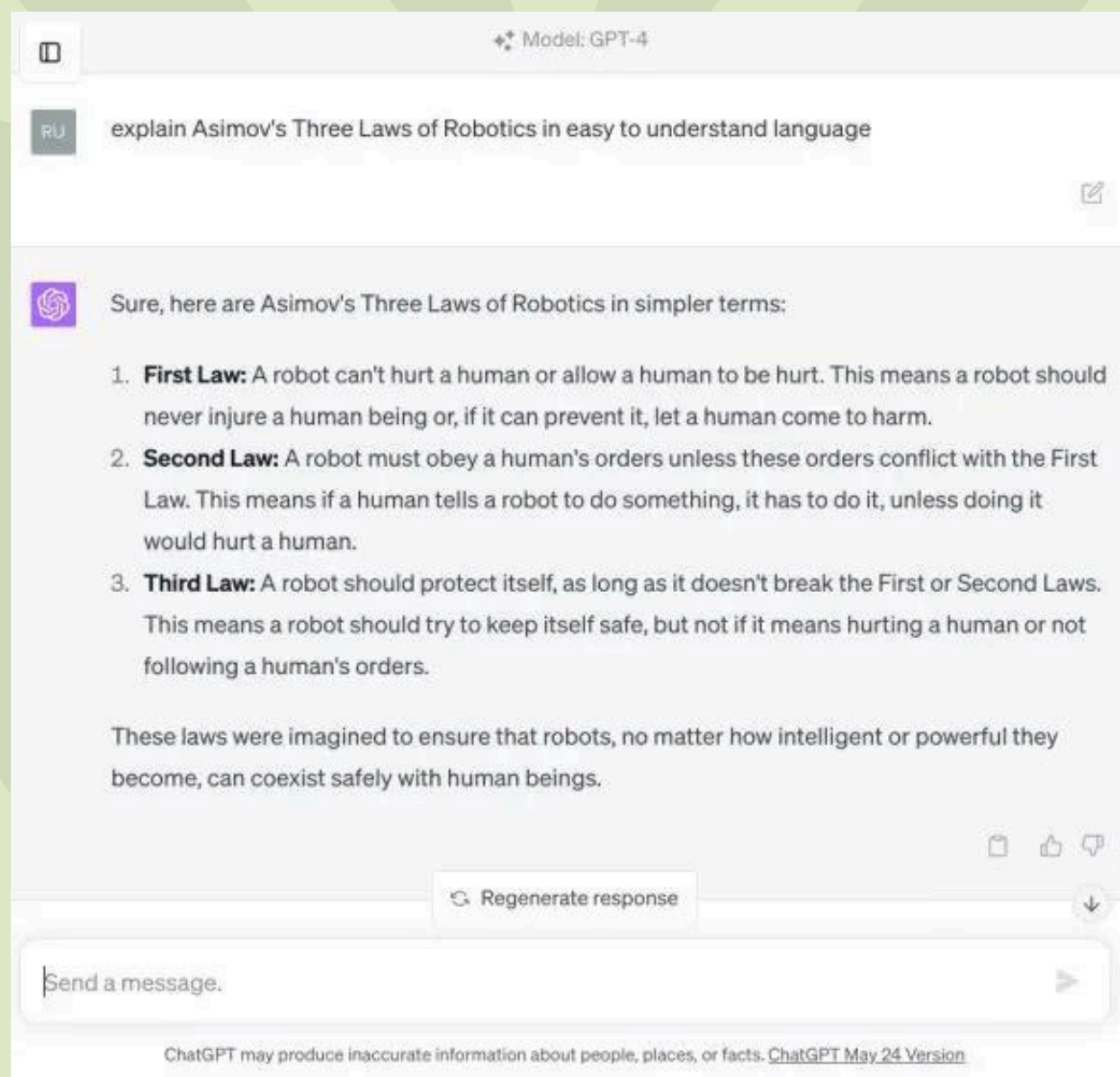
## 2. GEMINI FLASH 1.5

Google has also introduced a new addition to the family of Gemini AI models. The new AI model, dubbed Gemini 1.5 Flash is a light-weight model that is designed to be faster, more responsive, and cost-efficient. The tech giant said that it has worked on improving its latency to improve its speed. While solving complex tasks would not be its strength, it can do tasks such as summarisation, chat applications, image and video captioning, data extraction from long documents and tables, and more.

Gemini 1.5 Flash

Speed and efficiency

Multimodal reasoning

Long context window

# 3. GPT-4

The GPT-4 model by OpenAI is the best AI large language model (LLM) available in 2024. Released in March 2023, the GPT-4 model has showcased tremendous capabilities with complex reasoning understanding, advanced coding capability, proficiency in multiple academic exams, skills that exhibit human-level performance, and much more
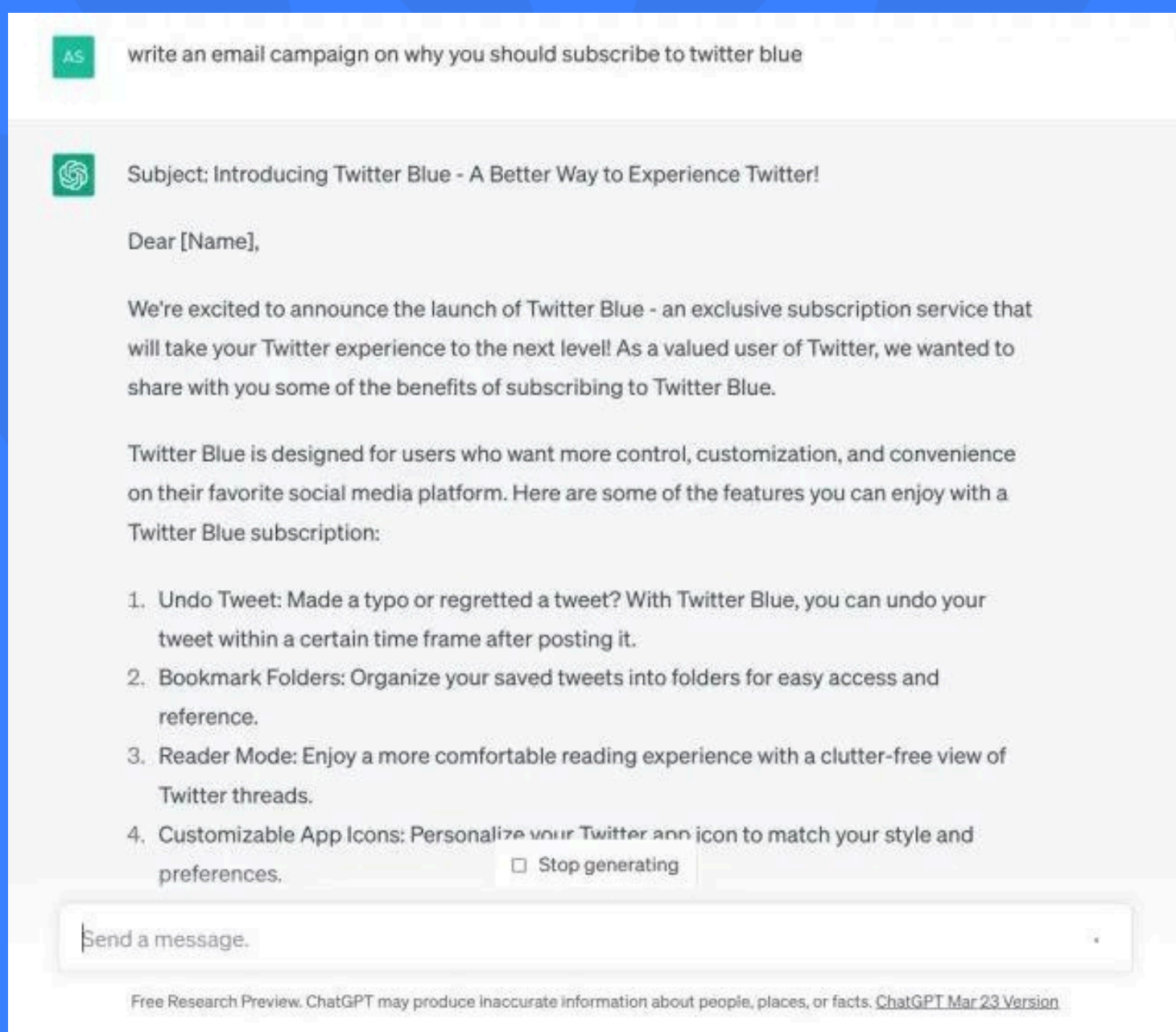
# 4. GPT-3.5

After GPT 4, OpenAI takes the second spot again with GPT-3.5. It's a general-purpose LLM similar to GPT-4 but lacks expertise in specific domains. Talking about the pros first, it's an incredibly fast model and generates a complete response within seconds.
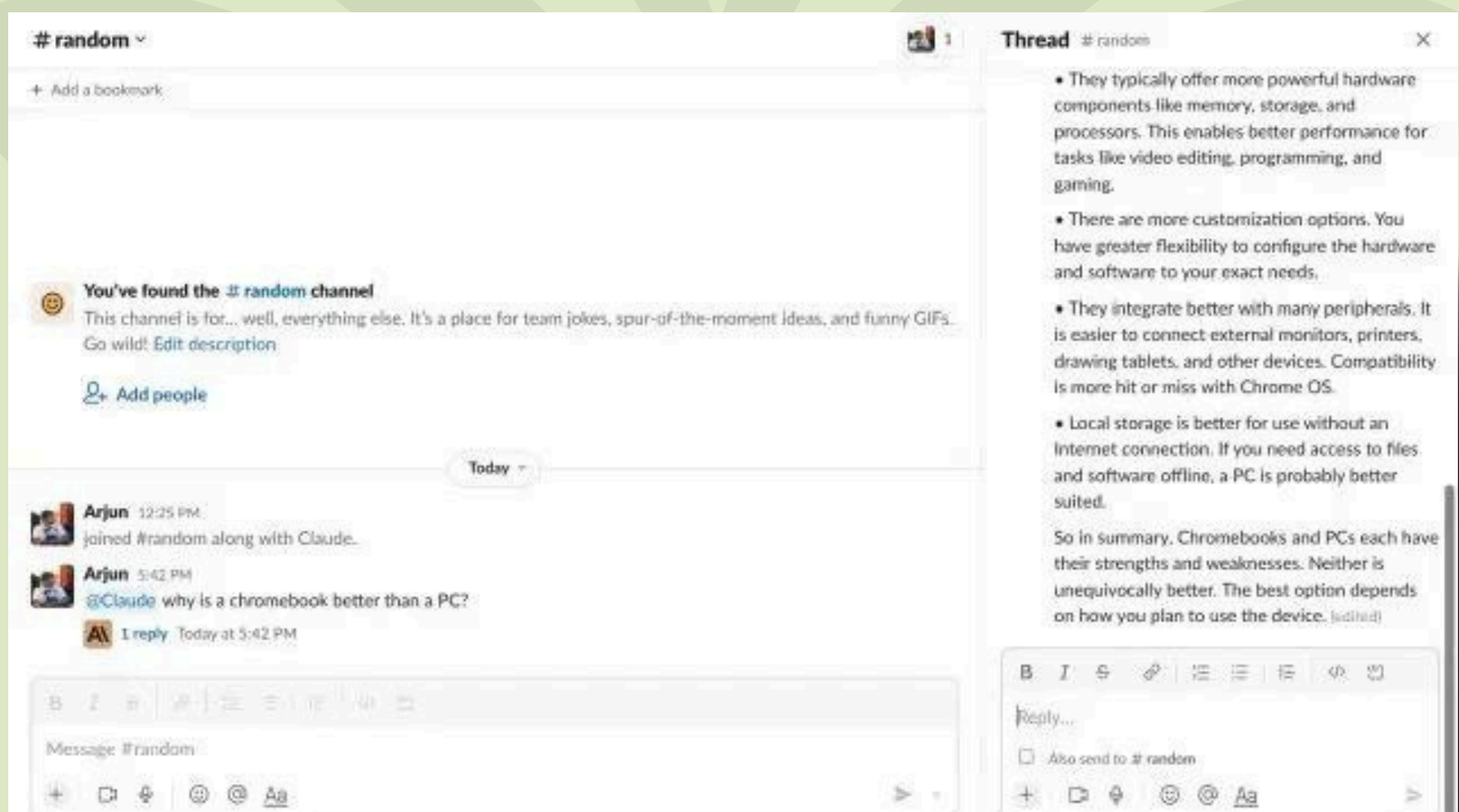
# 5. PALM 2 (BISON-001)

Next, we have the PaLM 2 AI model from Google, which is ranked among the best large language models of 2024. Google has focused on commonsense reasoning, formal logic, mathematics, and advanced coding in 20+ languages on the PaLM 2 model. It's being said that the largest PaLM 2 model has been trained on 540 billion parameters and has a maximum context length of 4096 tokens.
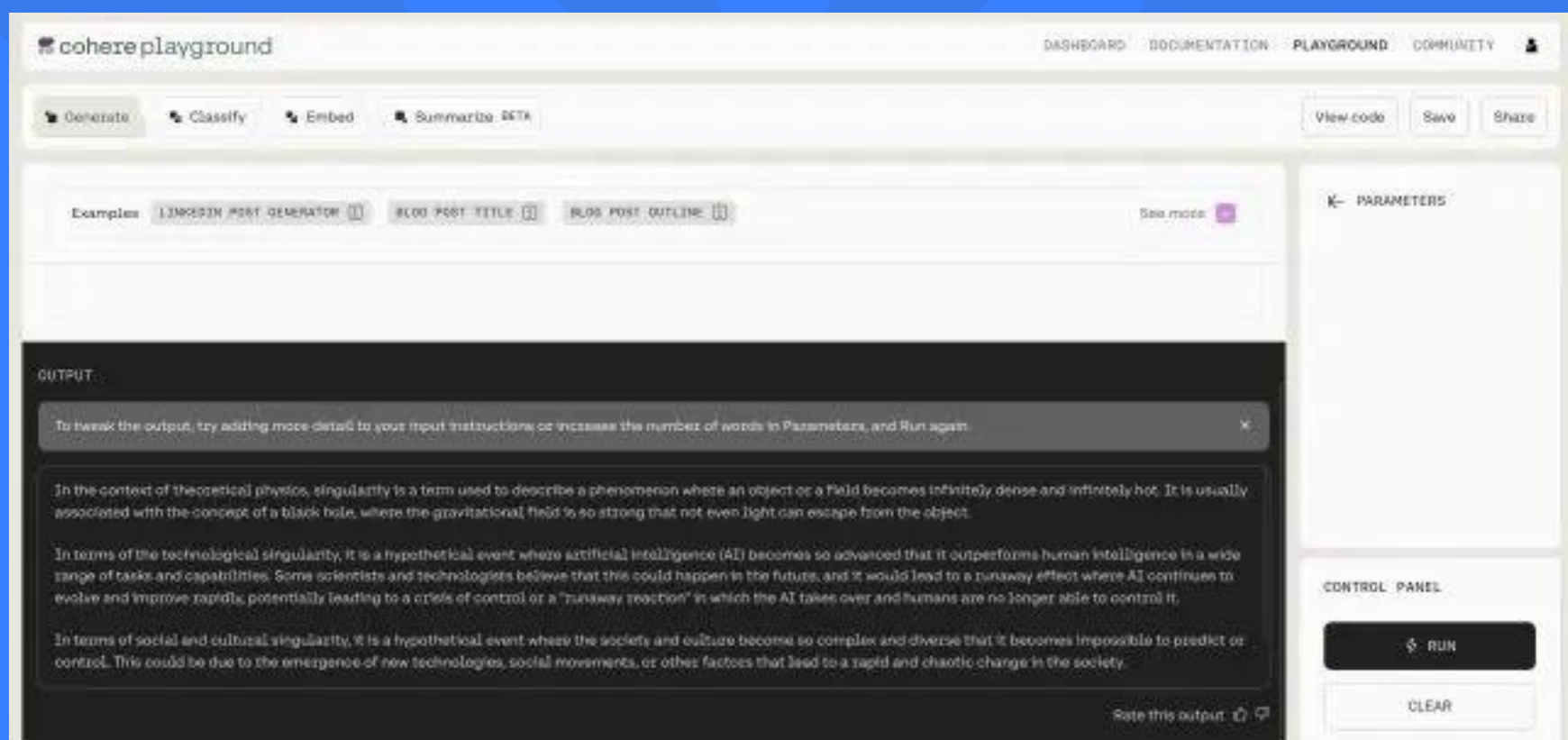
# 6. CLAUDE V1

In case you are unaware, Claude is a powerful LLM developed by Anthropic, which has been backed by Google. It has been co-founded by former OpenAI employees and its approach is to build AI assistants which are helpful, honest, and harmless. In multiple benchmark tests, Anthropic's Claude v1 and Claude Instant models have shown great promise. In fact, Claude v1 performs better than PaLM 2 in MMLU and MT-Bench tests.

# 7. COHERE

Cohere is an AI startup founded by former Google employees who worked on the Google Brain team. One of its co-founders, Aidan Gomez was part of the "Attention is all you Need" paper that introduced the Transformer architecture. Unlike other AI companies, Cohere is here for enterprises and solving generative AI use cases for corporations. Cohere has a number of models from small to large — having just 6B parameters to large models trained on 52B parameters.
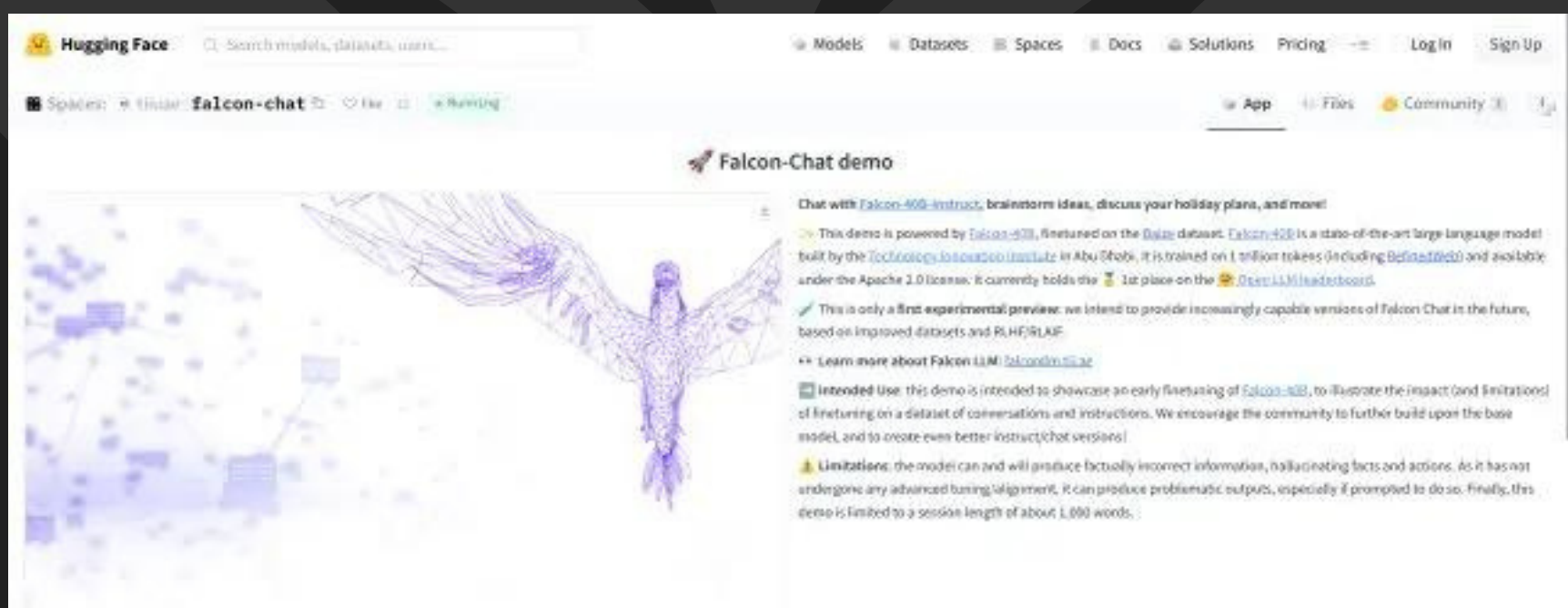
# 8. FALCON

Falcon is the first open-source large language model on this list, and it has outranked all the open-source models released so far, including LLaMA, StableLM, MPT, and more. It has been developed by the Technology Innovation Institute (TII), UAE. The best thing about Falcon is that it has been open-sourced with Apache 2.0 license, which means you can use the model for commercial purposes. There are no royalties or restrictions either.
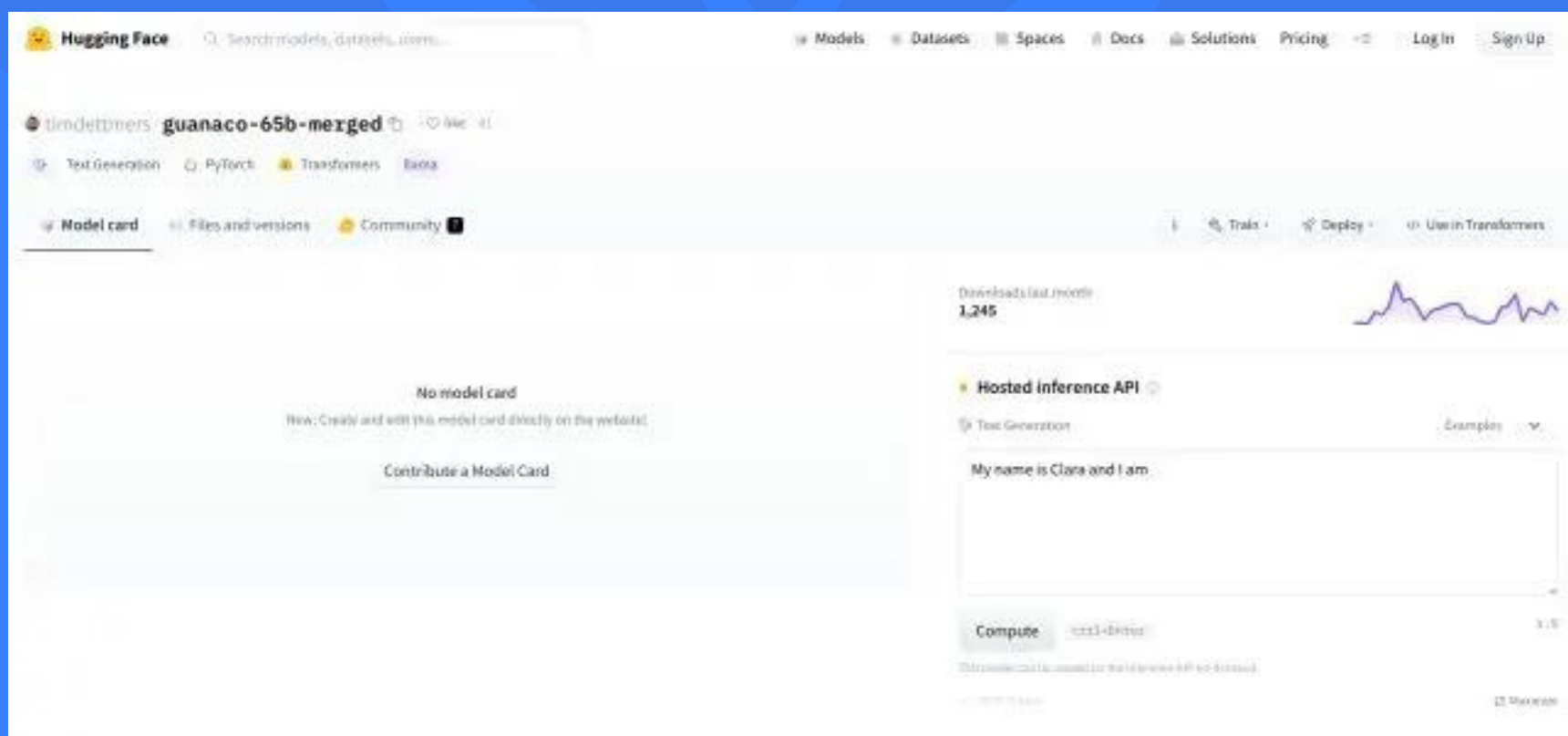
# 9. LLAMA

Ever since LLaMA models leaked online, Meta has gone all-in on open-source. It officially released LLaMA models in various sizes, from 7 billion parameters to 65 billion parameters. According to Meta, its LLaMA-13B model outperforms the GPT-3 model from OpenAI which has been trained on 175 billion parameters. Many developers are using LLaMA to fine-tune and create some of the best open-source models out there. Having said that, do keep in mind, LLaMA has been released for research only and can't be used commercially unlike the Falcon model by the TII.



**Meta AI**

Research

## Introducing LLaMA: A foundational, 65-billion-parameter large language model

February 24, 2023

As part of Meta's commitment to open science, today we are publicly releasing LLaMA (Large Language Model Meta AI), a state-of-the-art foundational large language model designed to help researchers advance their work in this subfield of AI. Smaller, more performant models such as LLaMA enable others in the research community who don't have access to large amounts of infrastructure to study these models, further democratizing access in this important, fast-changing field.

Training smaller foundation models like LLaMA is desirable in the large language model space because it requires far less computing power and resources to test new approaches, validate others' work, and explore new use cases. Foundation models train on a large set of unlabeled data, which makes them ideal for fine-tuning for a variety of tasks. We are making LLaMA available at several sizes (7B, 13B, 33B, and 65B
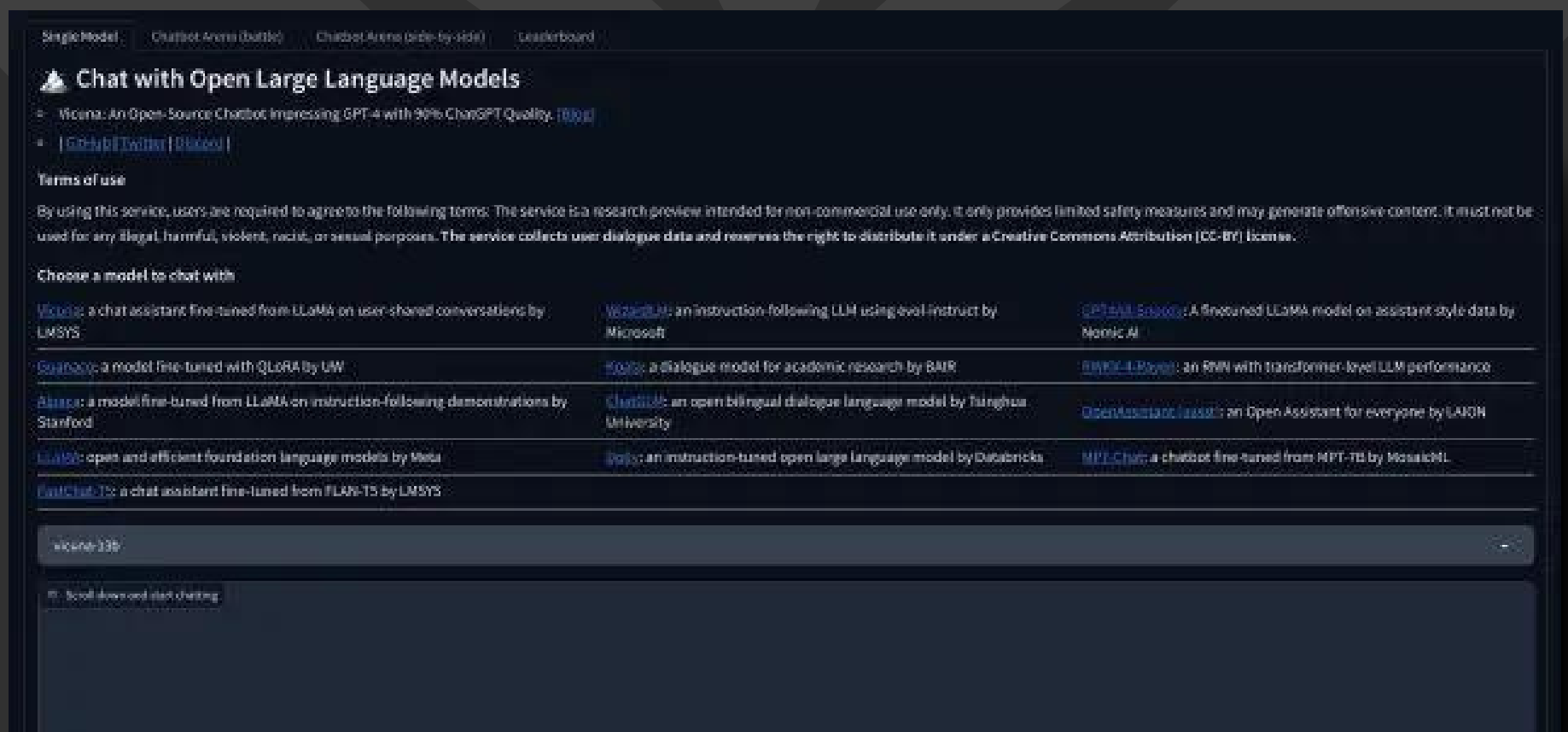
# 10. GUANACO-65B

Among the several LLaMA-derived models, Guanaco-65B has turned out to be the best open-source LLM, just after the Falcon model. In the MMLU test, it scored 52.7 whereas the Falcon model scored 54.1. Similarly, in the TruthfulQA evaluation, Guanaco came up with a 51.3 score and Falcon was a notch higher at 52.5. There are four flavors of Guanaco: 7B, 13B, 33B, and 65B models. All of the models have been fine-tuned on the OASST1 dataset by Tim Dettmers and other researchers.

# 11. VICUNA 33B

Vicuna is another powerful open-source LLM that has been developed by LMSYS. It has been derived from LLaMA like many other open-source models. It has been fine-tuned using supervised instruction and the training data has been collected from sharegpt.com, a portal where users share their incredible ChatGPT conversations. It's an auto-regressive large language model and is trained on 33 billion parameters.

## 12. MPT-30B

MPT-30B is another open-source LLM that competes against LLaMA-derived models. It has been developed by Mosaic ML and fine-tuned on a large corpus of data from different sources. It uses datasets from ShareGPT-Vicuna, Camel-AI, GPTeacher, Guanaco, Baize, and other sources. The best part about this open-source model is that it has a context length of 8K tokens.

# 13. 30B-LAZARUS



MPT-30B is another open-source LLM that competes against LLaMA-derived models. It has been developed by Mosaic ML and fine-tuned on a large corpus of data from different sources. It uses datasets from ShareGPT-Vicuna, Camel-AI, GPTeacher, Guanaco, Baize, and other sources. The best part about this open-source model is that it has a context length of 8K tokens.

# 14. WIZARDLM

WizardLM is our next open-source large language model that is built to follow complex instructions. A team of AI researchers has come up with an Evol-instruct approach to rewrite the initial set of instructions into more complex instructions. And the generated instruction data is used to fine-tune the LLaMA model.

# Dr. Martha Boeckenfeld

TOP 100 WOMEN OF THE FUTURE | WEB3 ADVISOR & INVESTOR | UN
PEACE AMBASSADOR | FOUNDER MARTHAVERSE
– GUIDING EXECUTIVES & TEAMS FROM WEB2 TO WEB3 WITH
STRATEGY, CREATION & EDUCATION

# FOUND THIS POST
# HELPFUL ?

# FOLLOW ME
## FOR MORE INSIGHTFUL POSTS

CHECKOUT MORE ABOUT ME

https://linktr.ee/marthaverse

CLICK HERE