

[Open in app](#)

≡ Medium

Search

 Write



 Member-only story

The Ultimate OpenAI Alternative is Exploding, and It Runs on Your Laptop.

No GPU required. How to Replace the OpenAI API with a Free, Local-First Solution.



Bytefer

[Follow](#)

4 min read · 3 days ago

 11



...

This Free OpenAI Alternative Is Exploding

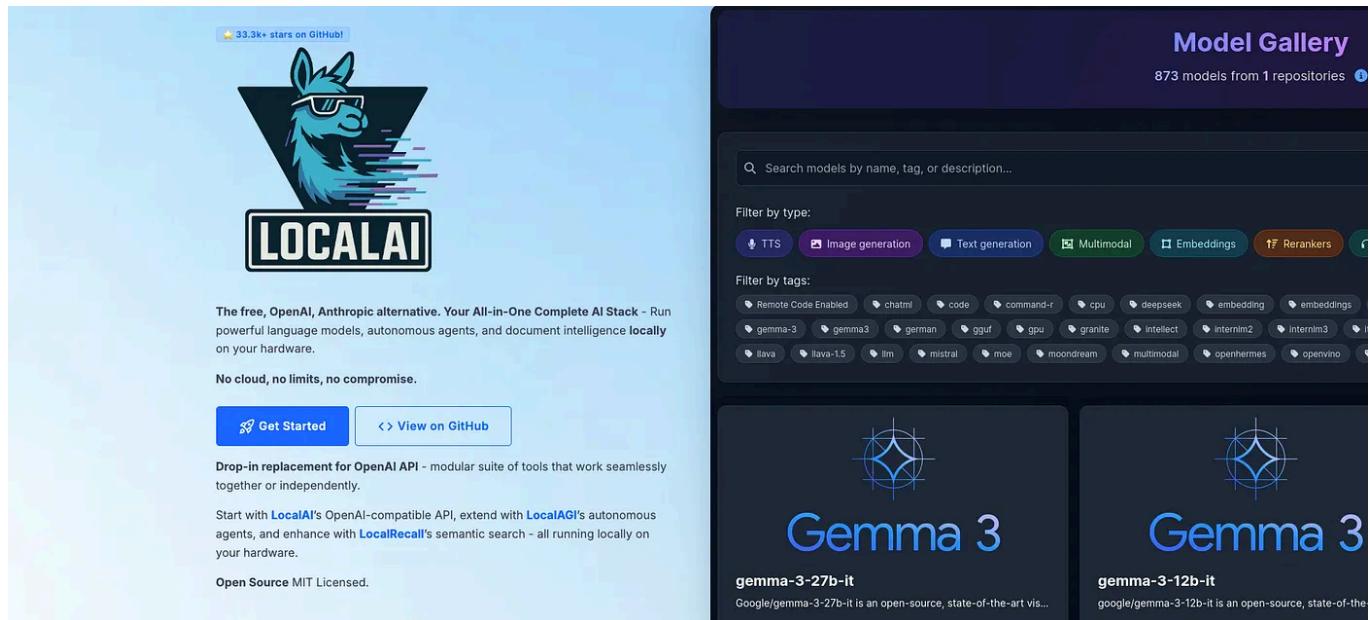


As developers, we all love to create something new with AI. But sometimes, before we can even start flexing our muscles, we're distracted by API Key requests, high fees, network latency, and data privacy issues. You've probably thought more than once: "If only I could simply run an OpenAI-like service on my own machine."

The open source project (**35.6K Starred**) introduced today could be the answer you're looking for.

What is LocalAI

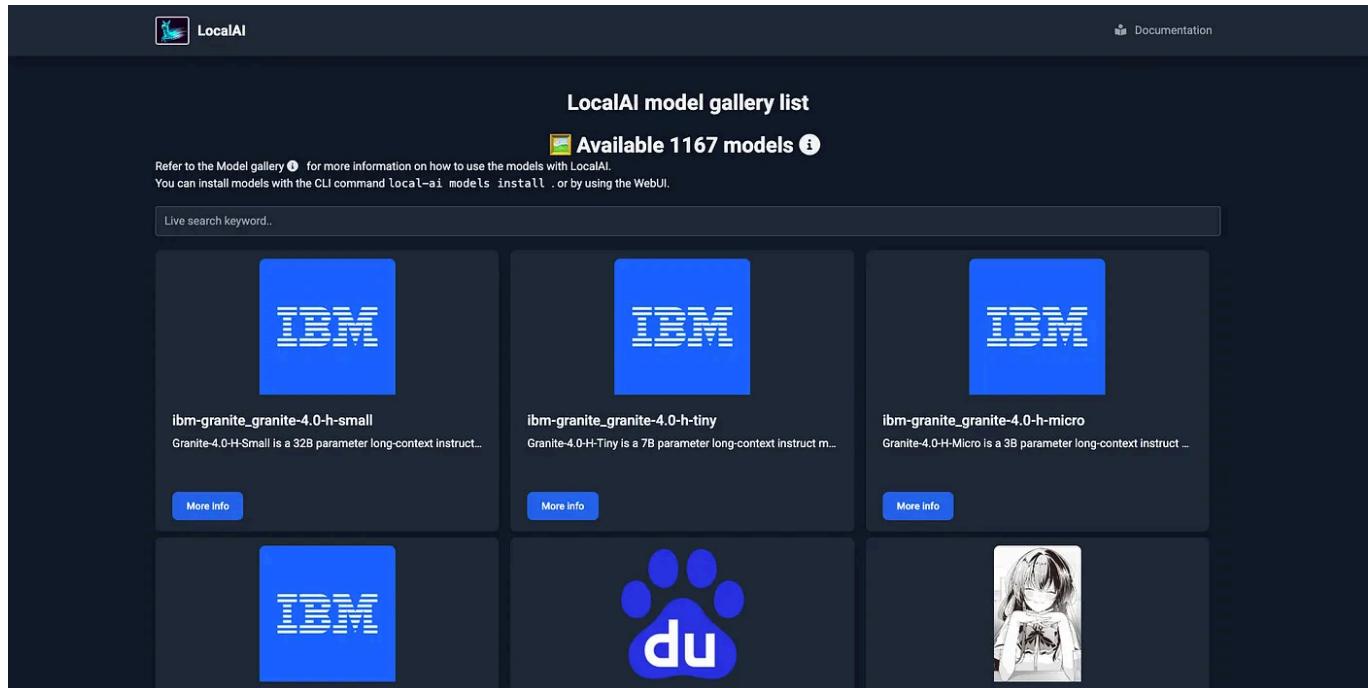
LocalAI is a free, open source alternative to OpenAI services. Its position is clear: to be a drop-in replacement REST API that lets you run Large Language Models (LLMs) and generate images and audio locally without having to send the data to any third-party servers.



<https://localai.io/>

Crucially, it's fully compatible with OpenAI's API standards. This means that your existing code that uses the OpenAI SDK or APIs will need little or no

modification — just point the API service address to your local LocalAI service address and everything will work as usual.



<https://localai.io/gallery.html>

This brings several obvious benefits to developers:

- Data Privacy and Security: All reasoning is done on your own hardware and the data never leaves your device. This is critical for applications that handle sensitive information or have compliance requirements.
- Manageable Costs: Has no token-based billing. Once deployed, you can call as many times as you want without worry.
- Runs offline: Your AI apps will continue to work even without an internet connection.

Quick Start

LocalAI officially offers several installation methods, the most convenient of which are the Bash installer script and Docker.

Using the Bash Installer

For most environments, a simple `curl` command is sufficient for a basic installation:

```
curl https://localai.io/install.sh | sh
```

Installing with Docker

If you are more comfortable with Docker, there are official images available for different hardware configurations.

```
docker run -ti --name local-ai -p 8080:8080 localai/localai:latest #CPU only image
```

If your machine has an NVIDIA graphics card and you have the appropriate drivers and CUDA toolkit installed, you can use GPU mirroring to get better performance.

```
# CUDA 12.0
docker run -ti --name local-ai -p 8080:8080 --gpus all localai/localai:latest-gpu

# CUDA 11.7
docker run -ti --name local-ai -p 8080:8080 --gpus all localai/localai:latest-gpu
```

If you are using an Intel graphics card, you can use the following command to install LocalAI:

```
docker run -ti --name local-ai -p 8080:8080 localai/localai:latest-gpu-intel
```

To install a model from the gallery, use the model name as the URI. For example, to run LocalAI with the Hermes model, execute:

```
local-ai run hermes-2-theta-llama-3-8b
```

To install only the model, use:

```
local-ai models install hermes-2-theta-llama-3-8b
```

Additionally, if you want to quickly experience the capabilities of LocalAI, you can install the AIO images:

```
# CPU version
docker run -ti --name local-ai -p 8080:8080 localai/localai:latest-aio-cpu

# NVIDIA CUDA 12 version
docker run -ti --name local-ai -p 8080:8080 --gpus all localai/localai:latest-ai

# NVIDIA CUDA 11 version
```

```
docker run -ti --name local-ai -p 8080:8080 --gpus all localai/localai:latest-ai

# Intel GPU version
docker run -ti --name local-ai -p 8080:8080 localai/localai:latest-aio-gpu-intel
```

The AIO images come pre-configured with the following features:

- Text to Speech (TTS)
- Speech to Text
- Function calling
- Large Language Models (LLM) for text generation
- Image generation
- Embedding server

Once the LocalAI service is started, you can send requests to <http://localhost:8080> as if you were calling the OpenAI API.

Using the LocalAI API

1.Text Generation

```
curl http://localhost:8080/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{ "model": "gpt-4", "messages": [{"role": "user", "content": "How are yo
```

2.Image Generation

```
curl http://localhost:8080/v1/images/generations \
-H "Content-Type: application/json" -d '{
  "prompt": "A cute baby sea otter",
  "size": "256x256"
}'
```

3.Text to speech

```
curl http://localhost:8080/v1/audio/speech \
-H "Content-Type: application/json" \
-d '{
  "model": "tts-1",
  "input": "The quick brown fox jumped over the lazy dog.",
  "voice": "alloy"
}' \
--output speech.mp3
```

4.Audio Transcription

```
curl http://localhost:8080/v1/audio/transcriptions \
-H "Content-Type: multipart/form-data" \
-F file="audio.wav" -F model="whisper-1"
```

If you want to learn how to call the Vision and Embeddings APIs, you can read the [official documentation](#).

Precautions

If you intend to deploy LocalAI to a public network environment, be sure to take good security precautions. The most straightforward way to do this is to

set the API_KEY environment variable when you start the service and add an access password to your API endpoint.

Summary

The [LocalAI](#) project provides a very useful option for developers. It allows us to build and test AI features in a more private, affordable, and controlled way. Whether you want to integrate LLM in your personal project or build a secure and controlled AI inference service within your organization, it's worth checking out.

AI

Python

Programming

Web Development

Generative Ai Tools



Written by Bytefer

3.7K followers · 18 following

Follow

Focus on web development and AI. <https://batchtool.com/>

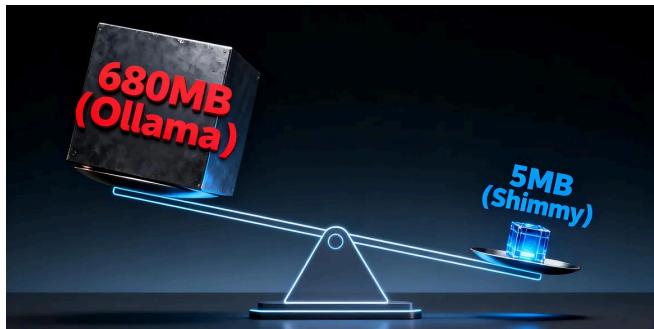
No responses yet



Bgerby

What are your thoughts?

More from Bytefer



 Bytefer

This Ollama alternative is under 5MB and compatible with the...

The Privacy-First Alternative to Ollama, Python-free Rust inference server—OpenAI...

 Sep 18  230  2

 + 

...



 Bytefer

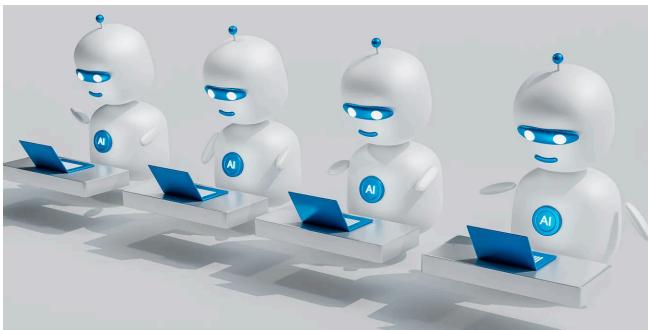
NVIDIA's Open-Source ASR Model Kills Whisper, With 1 Hour of Audi...

Open source and commercially available, super fast and high quality, the best choice f...

 Sep 30  40

 + 

...



 Bytefer

Top 8 Most Popular Open-Source AI Agent Framework

Don't Want To Just Tune APIs? These 8 Awesome Open Source AI Agent Framework...

 Sep 10  9  1



 Bytefer

The Open-Source ElevenLabs Alternative Is Here: Clone Voices i...

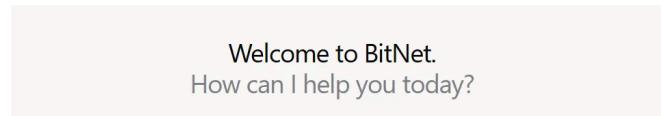
Stop choosing between quality and cost. Get production-grade, multilingual TTS with a...

 4d ago  41 



[See all from Bytefer](#)

Recommended from Medium



What do you want to know?

CPU  

By messaging BitNet, you agree to our [Terms](#) and [Privacy Policy](#).



In Data Science Collective by Erdogan T

Build Your Private Language Model: Local and Specialized For...

A complete step-by-step guide from setup to deployment of local language models, makin...

6d ago

644

4



...



Dr. Shouke Wei

Shimmy vs. Ollama: A Lightweight Alternative for Local LLM Serving

Why a 5 MB Rust Binary Could Change How We Run Local AI Models

Oct 1

106

3



...



Usman Writes

22 HTML Input Types That Will Make Your Forms 10x Better

The HTML <input> element is honestly one of the most versatile tags in web development....



In Coding Nexus by Algo Insights

Microsoft Just Declared War on the GPU Mafia: Meet bitnet.cpp

Bitnet.cpp will break the GPU Lock

4d ago

244

3



...



PY In Python in Plain English by Rizqi Mulki

The Hidden Python Framework That's Faster Than Node.js

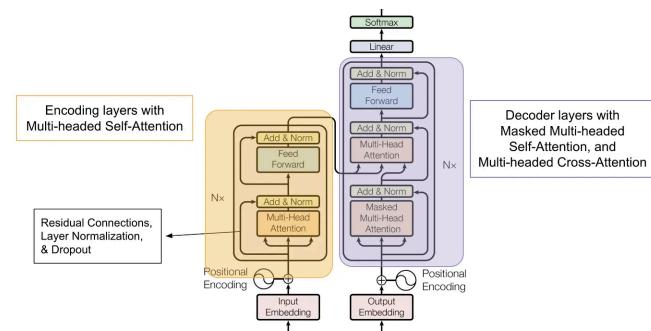
For years, developers have accepted a fundamental truth: Python is slow. When...

4d ago

11



...



In Towards AI by Ashish Abraham

No Libraries, No Shortcuts: LLM from Scratch with PyTorch

The no BS guide to build, train, and fine-tune a Transformer architecture from scratch

6d ago  188  6



•••



Oct 2

 646

 7



•••

[See more recommendations](#)