# The Rise of Open Source LLMs

Exploring the open-source Large Language Model (LLM) ecosystem on the code.siemens.com infrastructure

Fabio Huser & Dr. Ercan Uçan - on behalf of the code team

**SIEMENS**

# Who we are

Quick overview on code.siemens.com,
serving thousands of developers,
enabling and building communities

SIEMENS

# code.siemens.com - For developers, from developers

**DevOps and developer enablement for Open Source, Inner Source and Closed Source**
**A vibrant community - users get answers in less than an hour**

## Core features (GitLab CE)

- **Collaboration** on Code (version control, merge requests, issue tracker, wiki, social coding)
- **CI/CD** with shared Linux Windows and Mac runners, or bring your own
- **Pages** hosting (e.g., docs)
- Container and Package **Registry**

## Ecosystem

- Collaboration with **business partners** possible
- Continuous **security scans**
- **Core APIs** available for everyone, e.g. llm
- **Open & Inner Source** management
- Inner Source **Tax Transfers**, Export Control
- docs-theme (1900 re-uses)

- **High confidentiality** rating
- **> 70k users, 79 countries**
- **> 300k builds / day**
- Pure Open Source

## Ongoing Topics
- ☐ **Security** guidance
- ☐ Open Source **contributions**
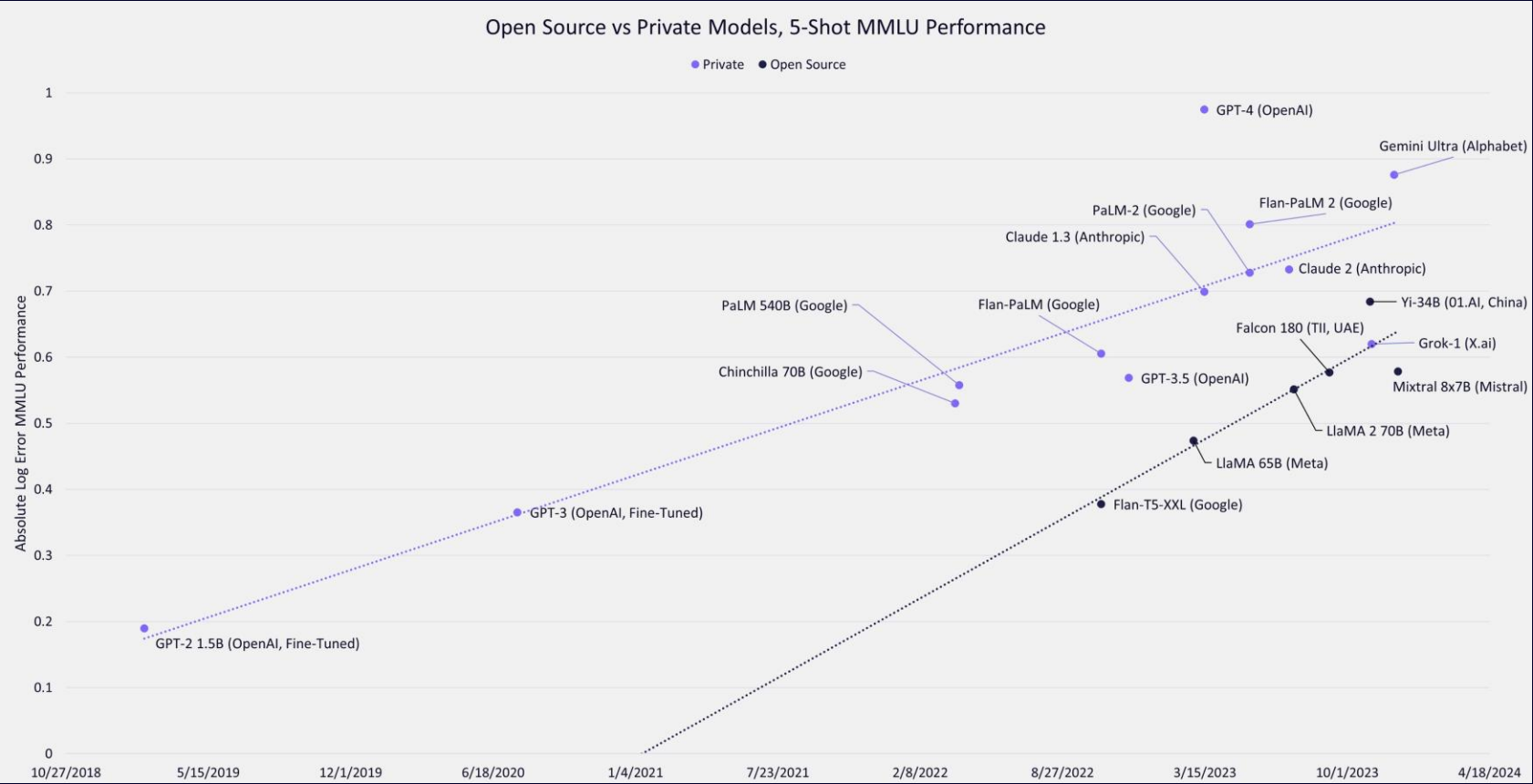- ☐ **AI assistants** (platform, IDE)

**SIEMENS**

# Why?

**SIEMENS**

# Advantages of open source LLMs in-house

- All **data stays internal**, highest confidentiality rating LLM.

- **Cost-effective** (no license/subscription fees).

- **Customizability** (in-house code usage) and control.

- **Context-based** answers about code and documentation.

- The same **ecosystem** code.siemens.com developers are used to.

- Open source LLMs are **catching up.**

**SIEMENS**

# Open Source LLMs are catching up



Open Source vs Private Models, 5-Shot MMLU Performance

• Private • Open Source

Less than 50% of
the evaluated models
are **Proprietary**
on the Leaderboard
of the LMSYS Chatbot Arena.

*March 7th, 2024*

The hyper-scalers are
nowadays also offering
Open Source models such as
Mistral.

Source: Cathy Wood, ARK Invest  https://twitter.com/CathieDWood/status/1735136182651027952

*"Open Source LLMs are about 6 months behind the commercial LLMs"* **,** *Arthur Mensch (CEO Mistral)*

# In-house LLM APIs

SIEMENS

# Current setup - all based on Open Source and Inner Source



- vLLM: https://github.com/vllm-project/vllm

- LiteLLM: https://github.com/BerriAI/litellm

- Kong: https://github.com/Kong/kong

- **OpenAI compatible** API endpoints

IDE plugin continue.dev
(like GitHub Copilot)

mention
@code-ai

self developed
apps and tools

code.siemens.com

Chat bot

Kong

@code-ai
bot

LLM Proxy

to be developed

worker

worker

worker

Siemens code context

LLM

LLM

LLM

**code.siemens.com AI enabled services on AWS**

**Sustainable AI (solar power, lake cooling) at HQ Zug**

**SIEMENS**

# Easy Self-Service, Get API key, Configure IDE, use it



- **Users get started in less than 5 minutes.**

- **All data stays within Siemens.**

- **Open Source IDE plugin continue.dev (Apache-2.0).**

## Next

- **Tab completion.**

- **Siemens context provider within IDE.**

https://blog.siemens.com/2024/04/open-source-llms-for-everyone/

**SIEMENS**

# Usage via VSCode IDE integration (continue.dev)

**SIEMENS**

# Continue.dev – open-source Visual Studio Code (VS Code) integration

**SIEMENS**

# Usage via local UI (Chatbox)

**SIEMENS**

# Chatbox



- Easily configured for usage with code.siemens.com LLM offerings.

- You can select models and personas (i.e., profiles).

- Open Source on GitHub:

  https://github.com/Bin-Huang/chatbox

**SIEMENS**

# Usage with code-ai bot

SIEMENS

# Code-AI bot



## Current Functionality:

- **summary**: summarize changes in Merge Requests.
- **ci-failure**: identifies root cause for failed pipeline.
- **clean-code**: clean code suggestions (community contribution).
- **discussion-summary**: summarize the discussions in the issue or merge request.

## Next Topics: Unit test creation, Secure coding analysis.

**SIEMENS**

# Usage via RAG with context

**SIEMENS**
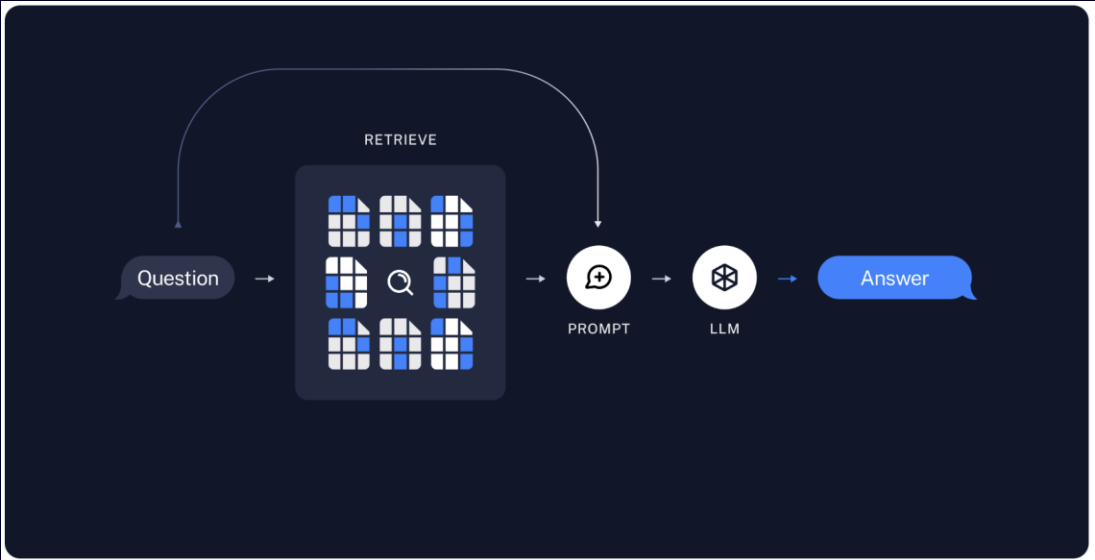
# RAG: Retrieval Augment Generation with context (code.siemens.com)

Ingest

Retrieve



**LangChain:** https://python.langchain.com/docs/use_cases/question_answering/

- Ingested with code.siemens.com documentation, various inner source project code samples.

**SIEMENS**

# RAG: Chat bot with Siemens context



## Retrieval Augmented Generation (RAG)



Initial implementation in a hackathon

Now in production and within testing phase

## Upcoming features:

- Projects on code.siemens.com use it with their context.

- Selection of context based on permissions.

**SIEMENS**

# Legal Aspects

Copyright, The Stack

**SIEMENS**

# How can we avoid code from a dataset ending up within our code?



The Stack v2 (67.5 TB) was released February 28th

**input**

**LLM based on The Stack Data Set**

**Machine generated content**

**Within Stack?**

As we know the dataset used to train the LLM, we can easily check for copies from The Stack and plagiarism.
E.g., scan automatically and notify user.

StarCoder2 is based on The Stack v2

**SIEMENS**

# Sustainable AI

On-prem Hardware
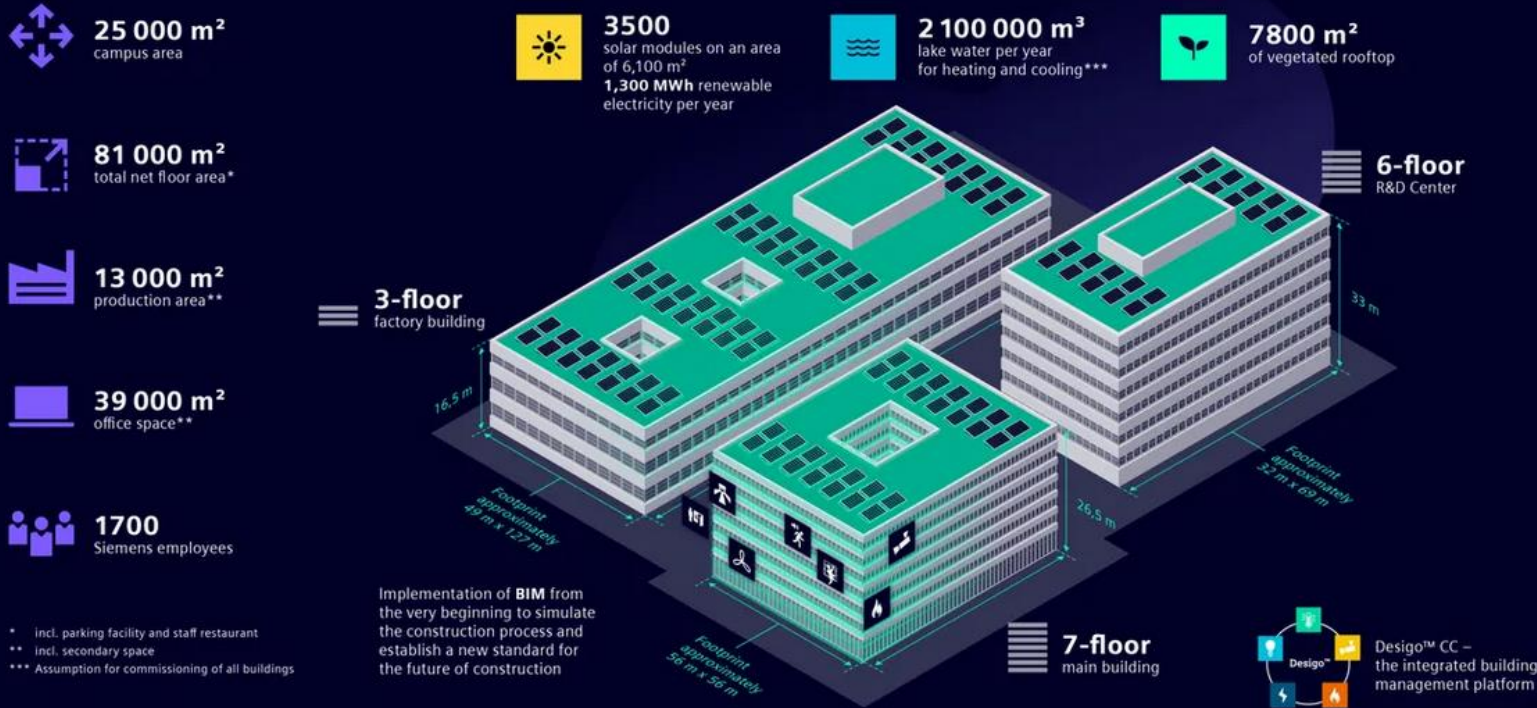(Future work)

**SIEMENS**

# Sustainability, in-house knowledge, climate-neutral AI in the long run

- Typical AI offerings charge per in- and outbound tokens -> **expensive**
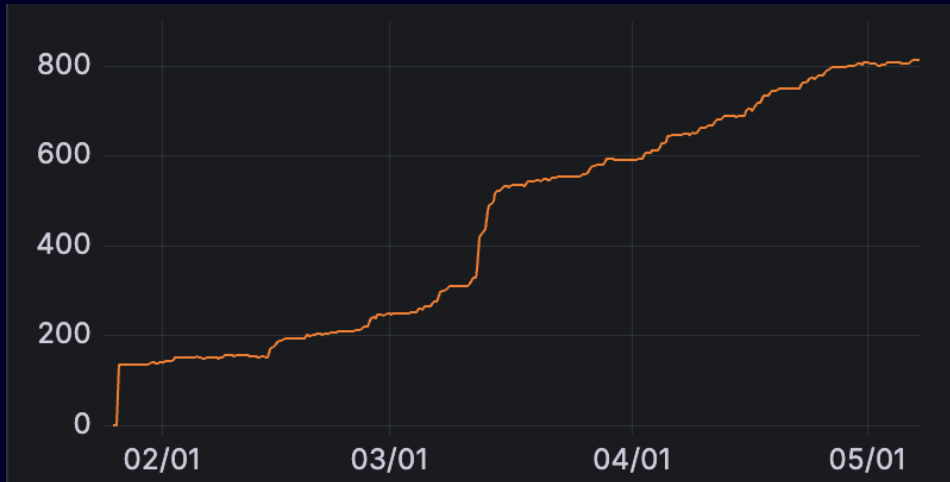- **Amortization** of a physical machine vs. AWS g5 ~ **3-4 months**



**Facts & Figures about the climate-neutral Campus Zug**

25 000 m² campus area

3500 solar modules on an area of 6,100 m²
1,300 MWh renewable electricity per year

2 100 000 m³ lake water per year for heating and cooling***

7800 m² of vegetated rooftop

81 000 m² total net floor area*

6-floor R&D Center

13 000 m² production area**

3-floor factory building

39 000 m² office space**

1700 Siemens employees

33 m

16.5 m

Footprint approximately 49 m x 127 m

Footprint approximately 32 m x 69 m

26.5 m

Footprint approximately 56 m x 56 m

7-floor main building

Desigo™ CC – the integrated building management platform

* incl. parking facility and staff restaurant
** incl. secondary space
*** Assumption for commissioning of all buildings

Implementation of **BIM** from the very beginning to simulate the construction process and establish a new standard for the future of construction



**Server room within the R&D center**

*"Companies with serious plans for AI have the GPUs and the knowledge in-house."*

**SIEMENS**

# Feedback from the community



Keep up the pace and your great work that you are doing. I really like it :)

👍 UNLIKE    ↩ REPLY    ↗ SHARE    ⋯

You, Wittig, Max (IT APS DEV COD) and Huser, Fabio (IT APS DEV COD) reacted to this

Huser, Fabio (IT APS DEV COD) I really like the usability of this, couldn't be any easier for us users. Keep up the great work!

**Feedback**

I tried the PoC with `continue` in VS Code according to https://code.siemens.io/general/ai/#large-language-model-llm-api-usage:

I used it on multiple Python projects we are developing so far:

- Sometimes continue creates loops with endless new lines.
- Indexing on large projects (with 500+ files) seems to not work as expected (for smaller projects this works fine). I guess more specific configuration might help.
- When providing two or more files as context for writing doc strings, the results are quite promising and only little changes are required.
- Sometimes response speed is quite fast, but sometimes a quite big delay is visible.
- Useful for fixing typos and grammar
- Having access to a chatbot directly in VS Code is helpful and speeds up development.

Key takeaways so far:

- This technology can already be helpful in early stages. Setup was quite smooth but for big projects more instructions are required.
- Thanks for bringing this to life!
- This feedback was also supported by continue 🙂

💚 5    😊

Great work colleagues. Leading the way for open source adoption even in this new and niche field.

Really love the work that you're doing here. I'd love to get involved on the Healthineers side of things.

you guys rock! cool stuff

This is really nice work, we'll be trying it in a project pretty much immediately :) Thanks a lot to you all!

**SIEMENS**

# Contact



- opensource.siemens.com

- Fabio Huser: [fabio.huser@siemens.com](mailto:fabio.huser@siemens.com)

- Dr. Ercan Uçan: [ercan.ucan@siemens.com](mailto:ercan.ucan@siemens.com)

**SIEMENS**

# Backup slides

**SIEMENS**

# Integration into internal crash reporting

**SIEMENS**

# API usage for all possible applications – demo with Bruno (API Client)

**SIEMENS**

# Sustainability, in house knowledge
# climate-neutral AI in the long run

- Typical offerings charge per inbound and outbound tokens, which is more expensive, and we have no influence if we reach limits.

**Example cost calculation**
- AWS g5.12xlarge (4 x 24 GB GPU)
  - 192 GB GPU = ~ **$ 122.640 / year**
- Physical AIME machine
  - 192 GB GPU = ~ **€ 42.149 / one time**

- **Amortization** for a physical machine is **~ 3/4 months**

**Estimate:** GitHub Copilot IDE plugin vs. GPU on prem
- Serving 1000 developers with a 192GB machine for 3 years
- **~ €1.38 / month**

*"Companies with serious plans for AI have the GPUs and the knowledge in-house."*



**3.2 kW max**



**6.5 kW max**

**SIEMENS**