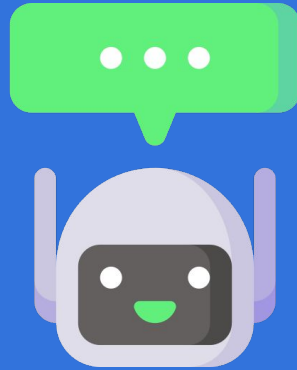


# A Practical Guide to Fast and Private LLMs



# Do I Even Care?



No



No but in yellow

Why?

It's **MY** fucking data!





Ali Farhat 

Posted on 1 Aug 2025 • Edited on 1 Jan



45



13



14



13



14

# ChatGPT Privacy Leak: Thousands of Conversations Now Publicly Indexed by Google

[#ai](#) [#privacy](#) [#gdpr](#) [#chatgpt](#)

Google has indexed thousands of ChatGPT conversations — exposing sensitive prompts, private data, and company strategies. Here's what happened, why it matters, and how you can protect your [AI workflows](#).



r/OpenAI 2h · arstechnica.com



**OpenAI slams court order to save all ChatGPT logs, including deleted chats**



↑ 47



🗨 12



r/LocalLLaMA 1h · arstechnica.com



**After court order, OpenAI is now preserving all ChatGPT and API logs**



↑ 95



🗨 11



# Leak, Cheat, Repeat: Data Contamination and Evaluation Malpractices in Closed-Source LLMs

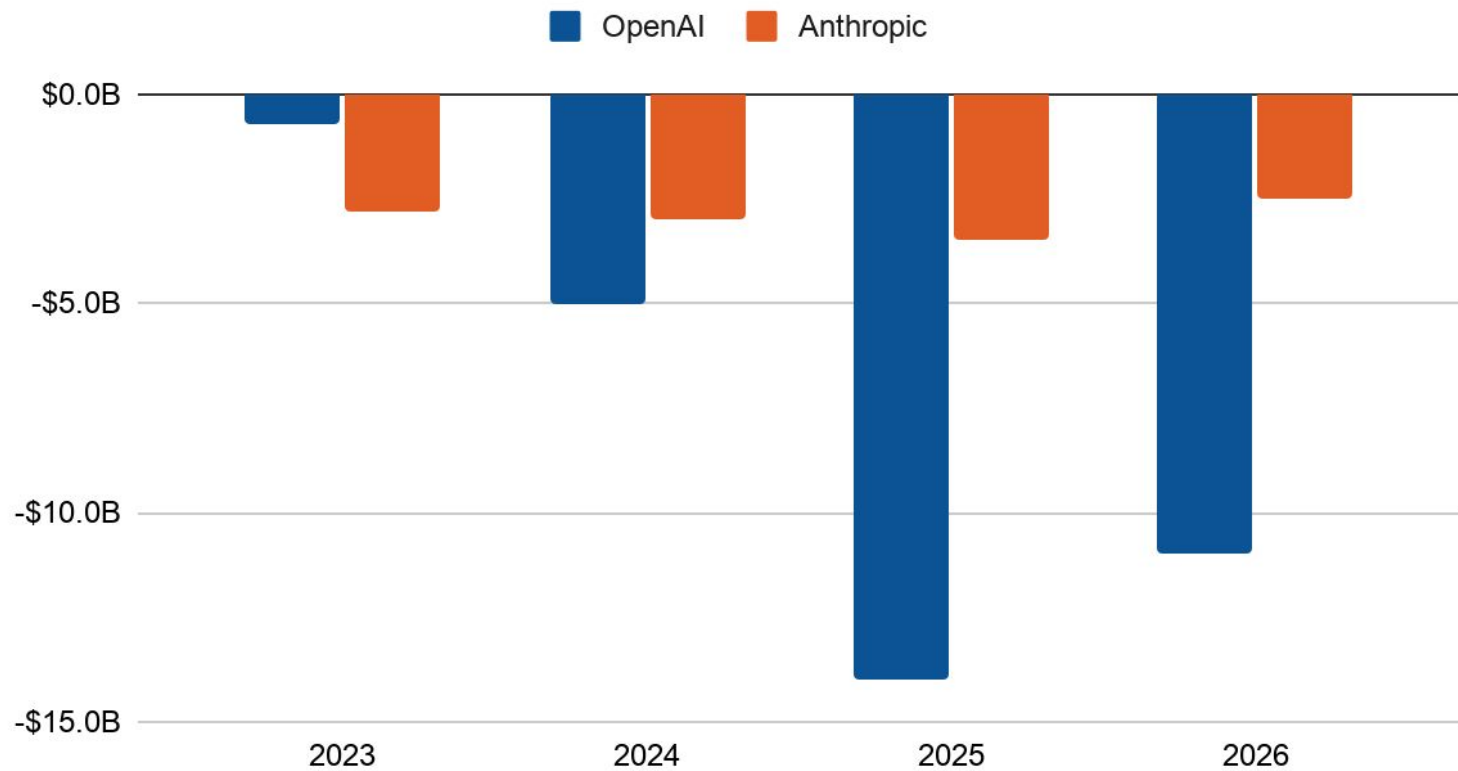
[github.com/leak-llm/leak-llm.github.io](https://github.com/leak-llm/leak-llm.github.io)

Why?

It's **MY** fucking money!



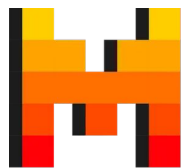
# Yearly Earnings



How?

It's **MY** fucking LLM!





Public

ChatGPT / Claude / Mistral



Hosted

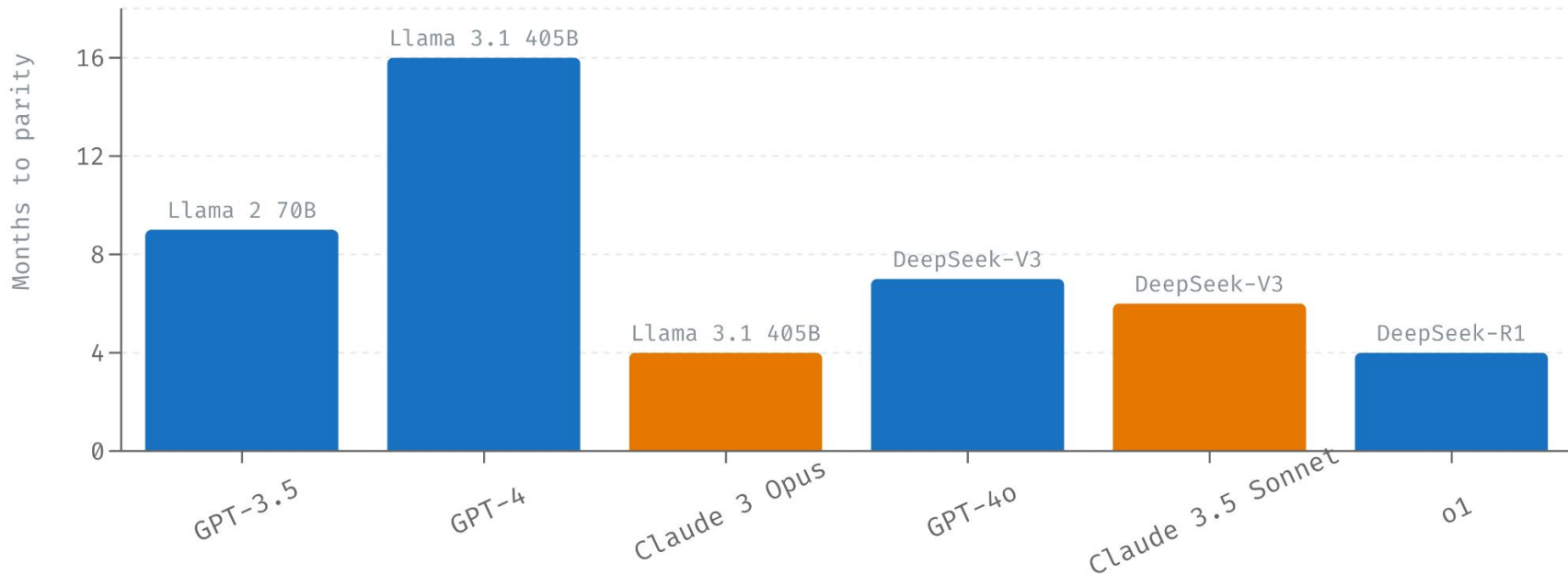
Bedrock / Copilot / Vertex

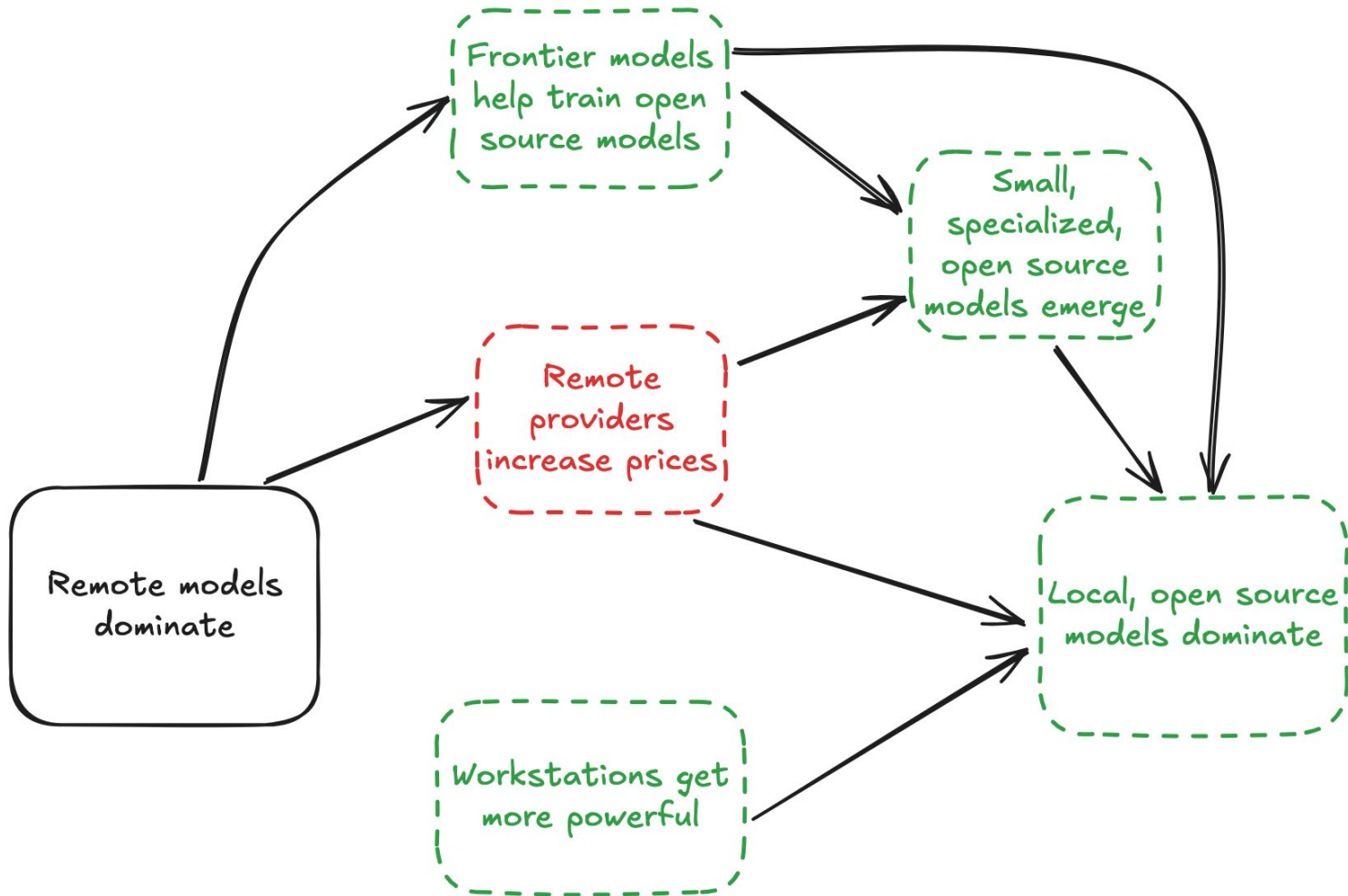


Local

Ollama / vLLM / TGI

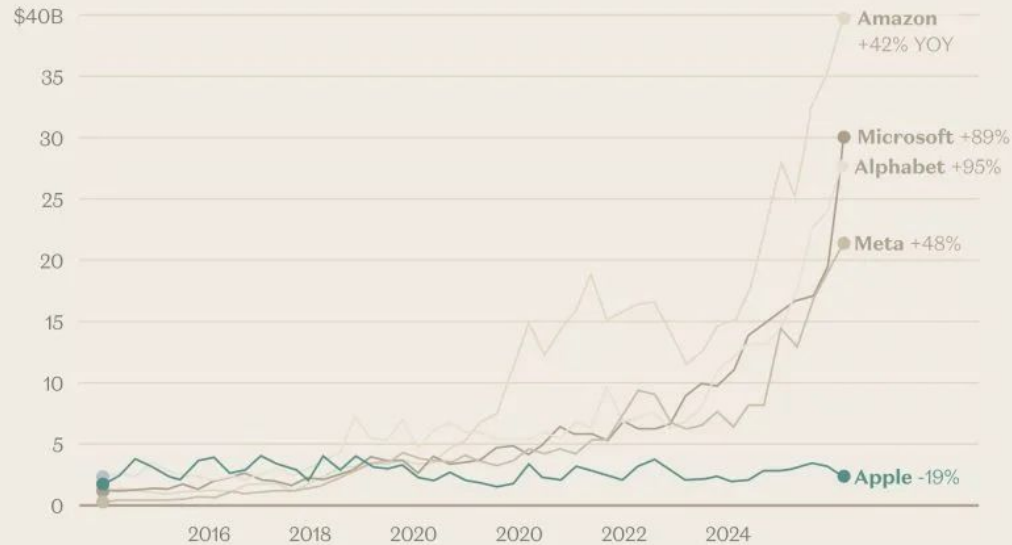
## Months to open source parity with frontier models





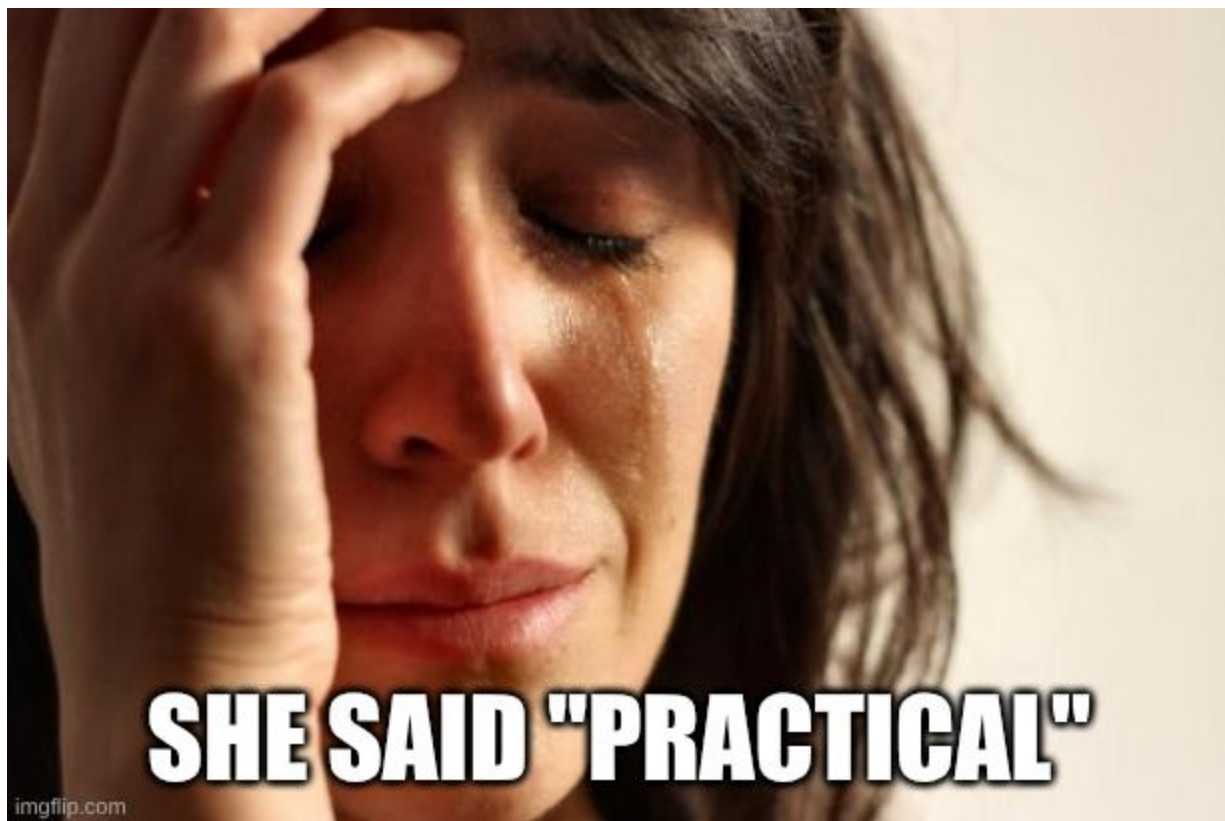
# Apple on Capex: 'Nah, we're good'

Standardized quarterly capital expenditure



Source: FactSet, Snacks (2/9/26)

AIGZ



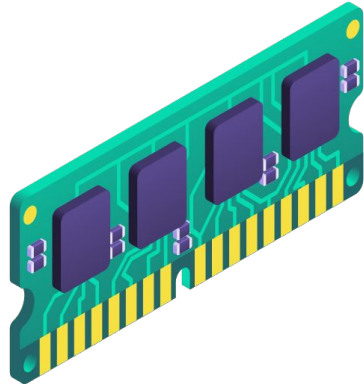
**SHE SAID "PRACTICAL"**

What?

It's **MY** fucking stack!

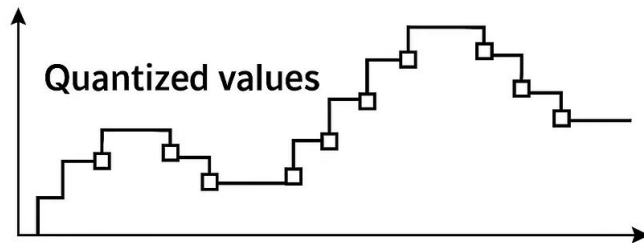
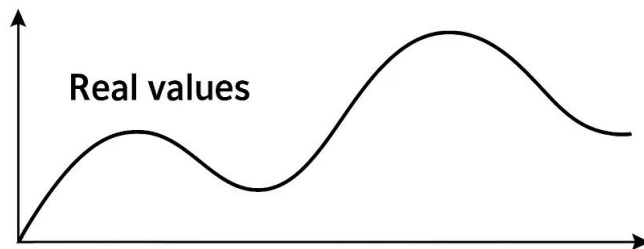


# Part 1: The Constraint



Memory,  
Not Compute

# Part 1: The Constraint



# Part 2: The Architecture

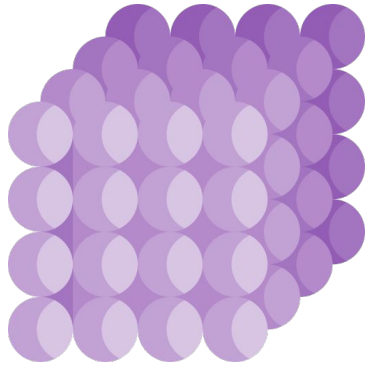


Small / Specialist



Large / Generalist

# Part 2: The Architecture



Dense



Mixture-of-Experts



Qwen3-35B-A3B-Q4\_K\_M.gguf

# Part 3: The Hardware

7B

5 GB

13B

8 GB

35B

20 GB

70B

40 GB

120B

70 GB

# Part 3: The Hardware



RTX 4070  
(12GB VRAM)

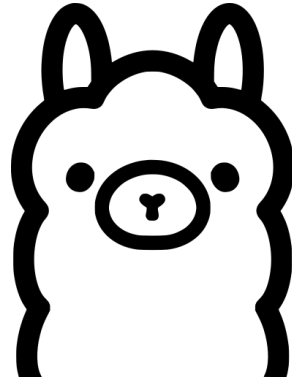


RTX PRO 6000  
(96GB VRAM)



Mac Studio M4 Max  
(128GB RAM)

# Part 4: The Software



What is Lena like?

Lena

Lena has

Lena has very

Lena has very weird

Lena has very weird hair

# Part 4: The Software



**kubernetes**

Memory

Quantization

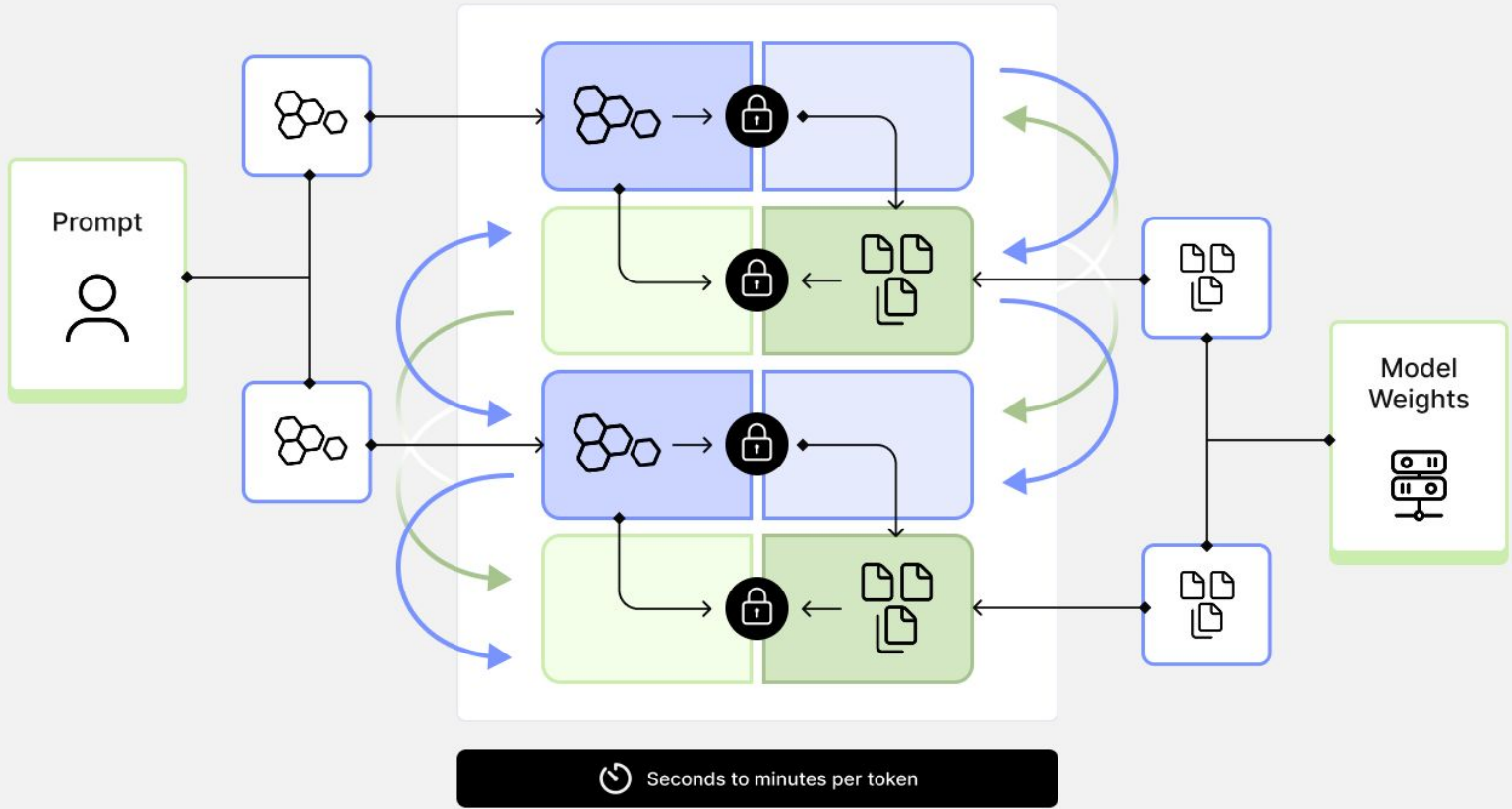
MoE

Dare to Go Small

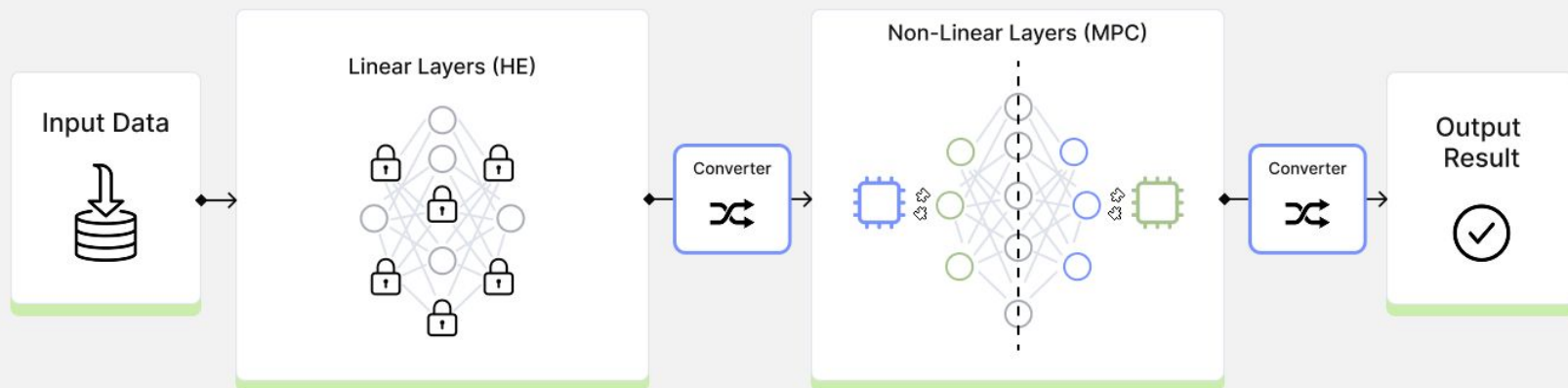
# Thanks!

[bespinian.io](https://bespinian.io)





# Fully Homomorphic Encryption for LLM Inference



⬇️ Reduced communication, added complexity ⬆️

# Fully Homomorphic Encryption for LLM Inference

