# Tips for building applications that use GenAI

**December 10th, 2024**
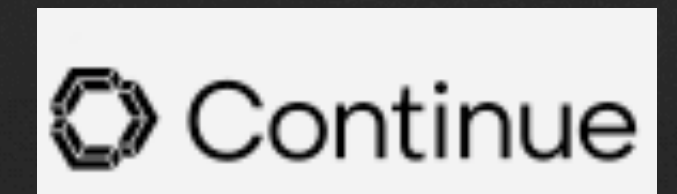
**Juan Peredo**

BOLBECK LLC

# Juan Peredo - Your guide in this journey

- Architect / consultant / developer & everything in between

- Independent Cloud Application Architect / Entrepreneur

- Linkedin: linkedin.com/in/juanperedotech

- Over 15 years of IT experience in companies like:

  - Bolbeck LLC

  - AWS

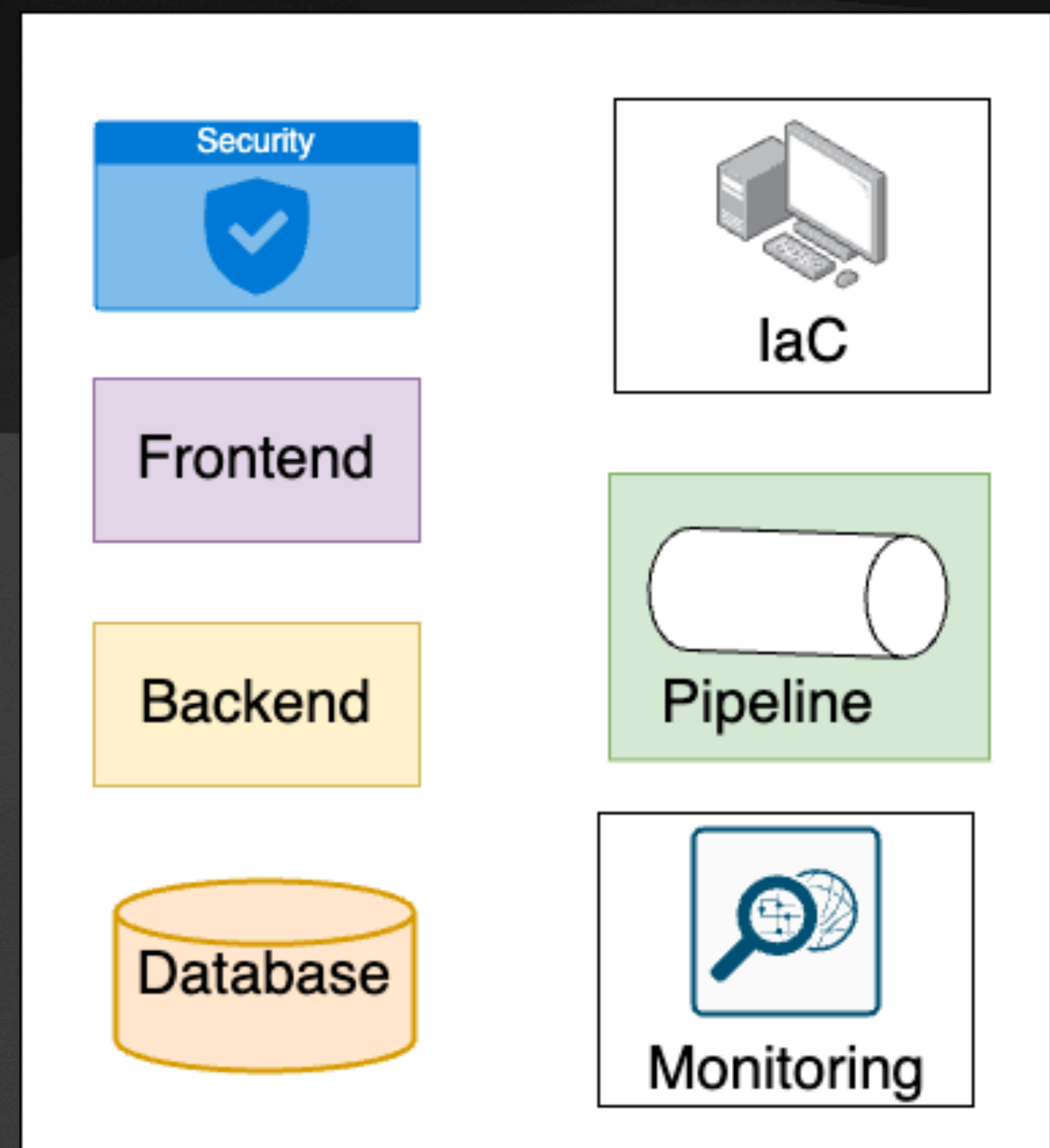  - Strategy& (PWC)

  - Booz & Co
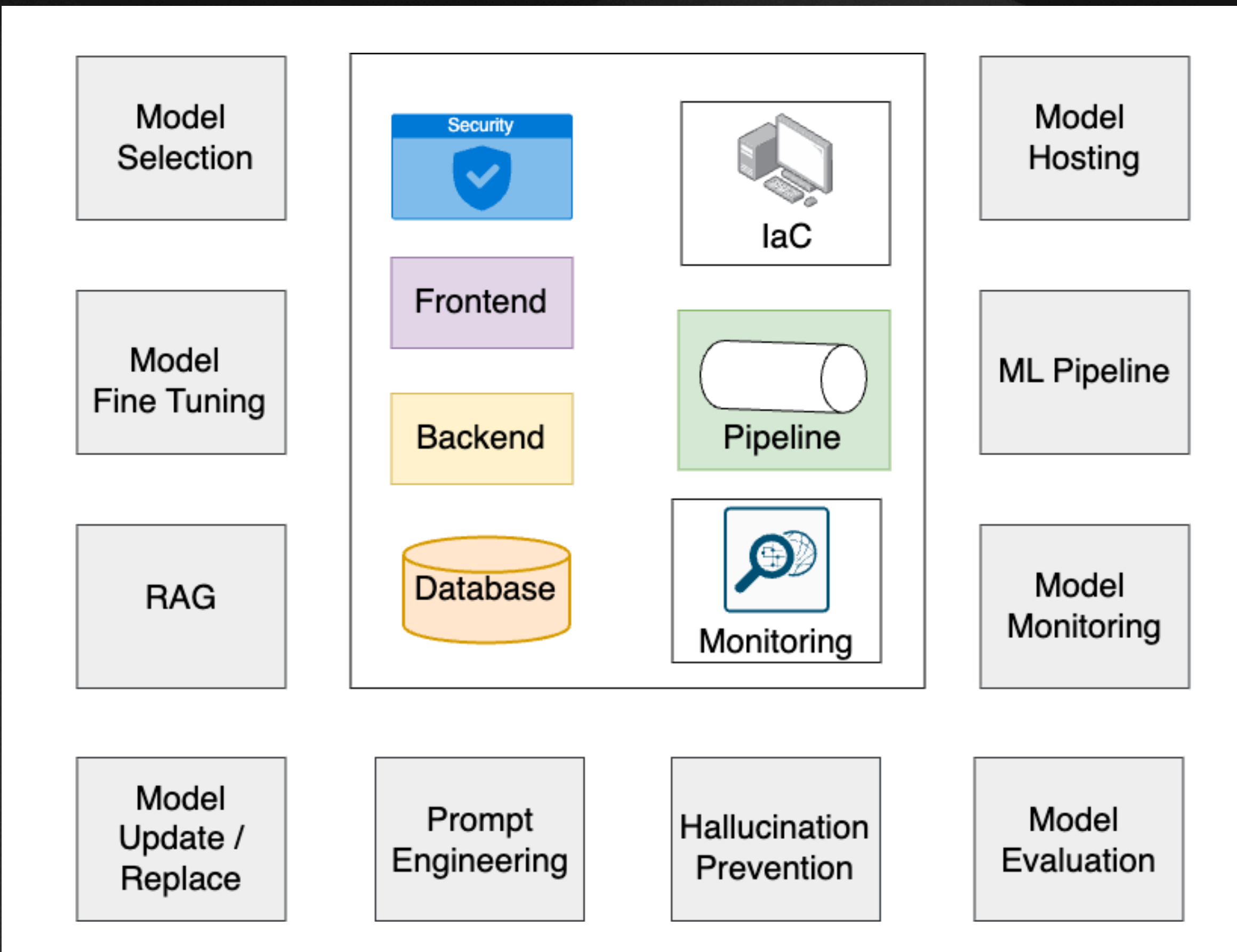
# AI can help reduce the time to code an App

- These are plugins for your IDE except cursor that is an actual IDE

- You can also just use any LLM to help with your code via cut and paste

# However, AI models also increase your app complexity

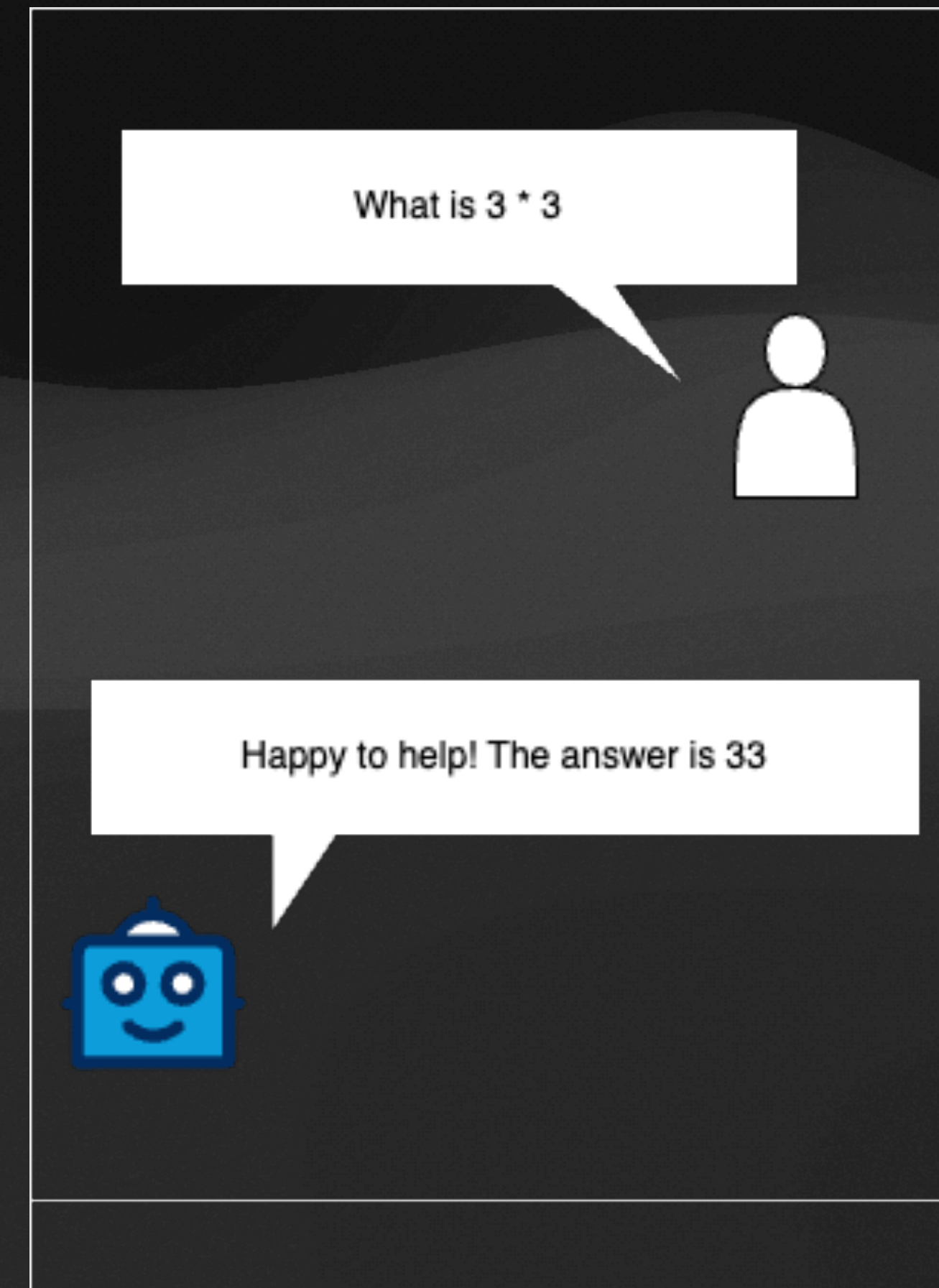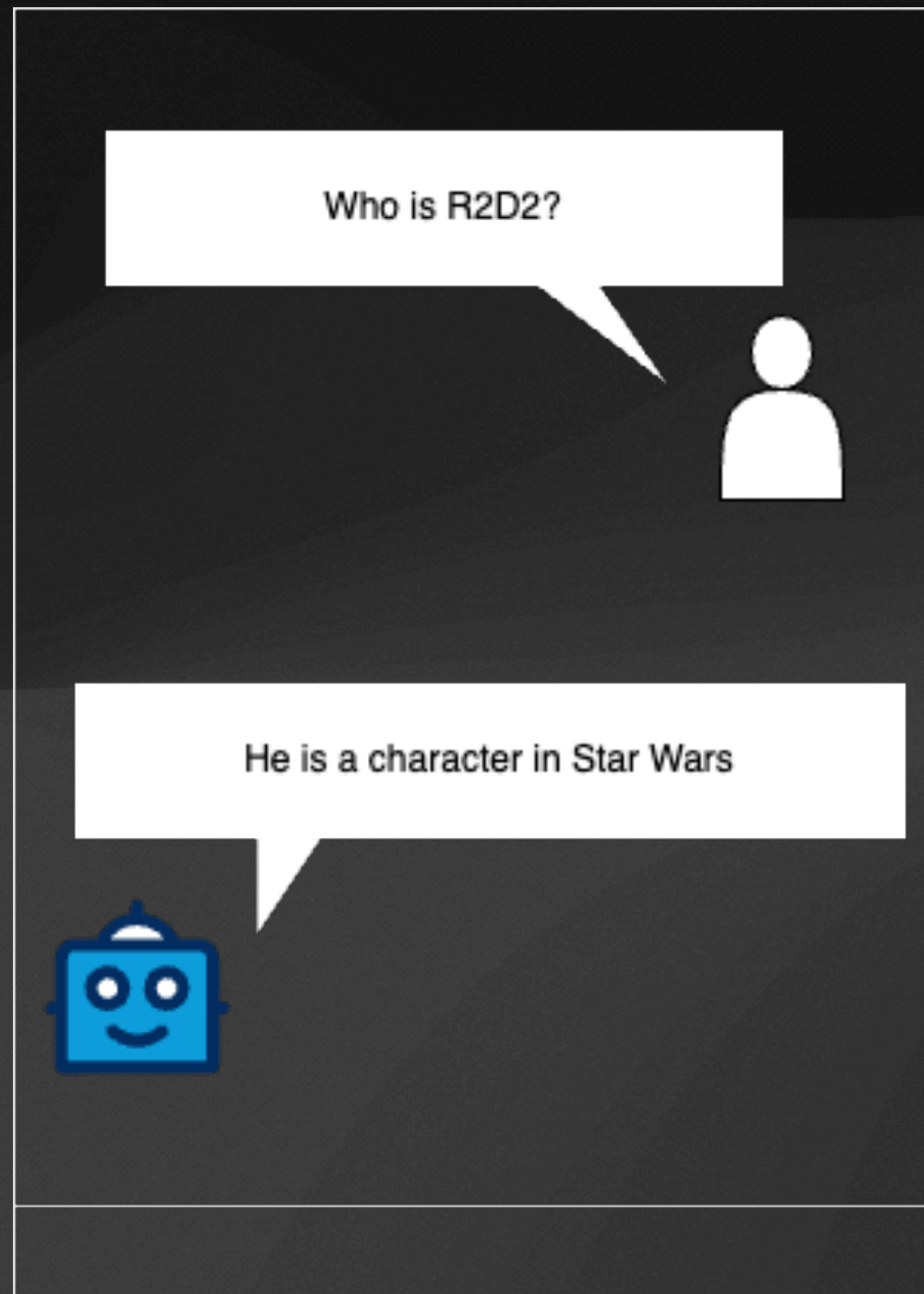# Where/How to host Models is important

Host locally

Host in the cloud

Alternatively, you could just use a library like HF Transformers

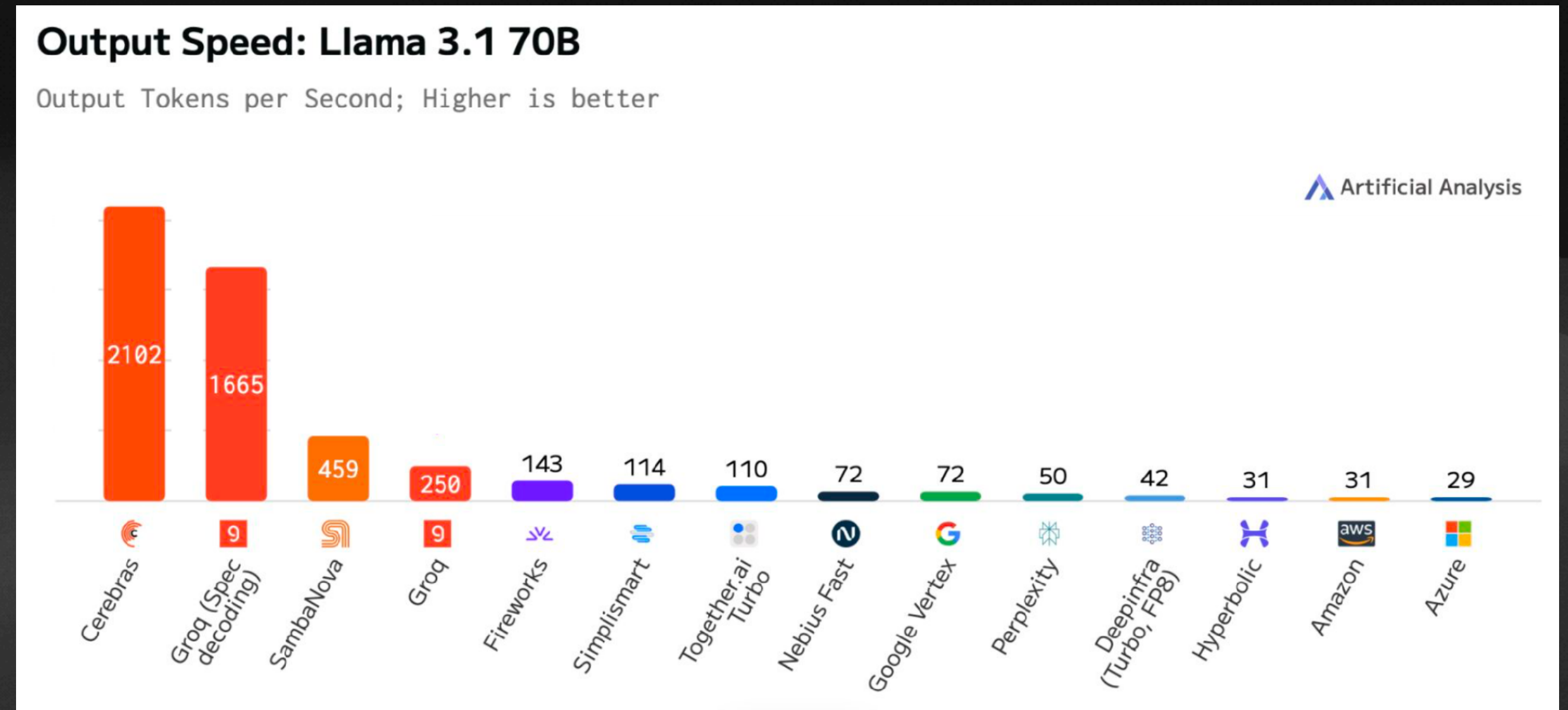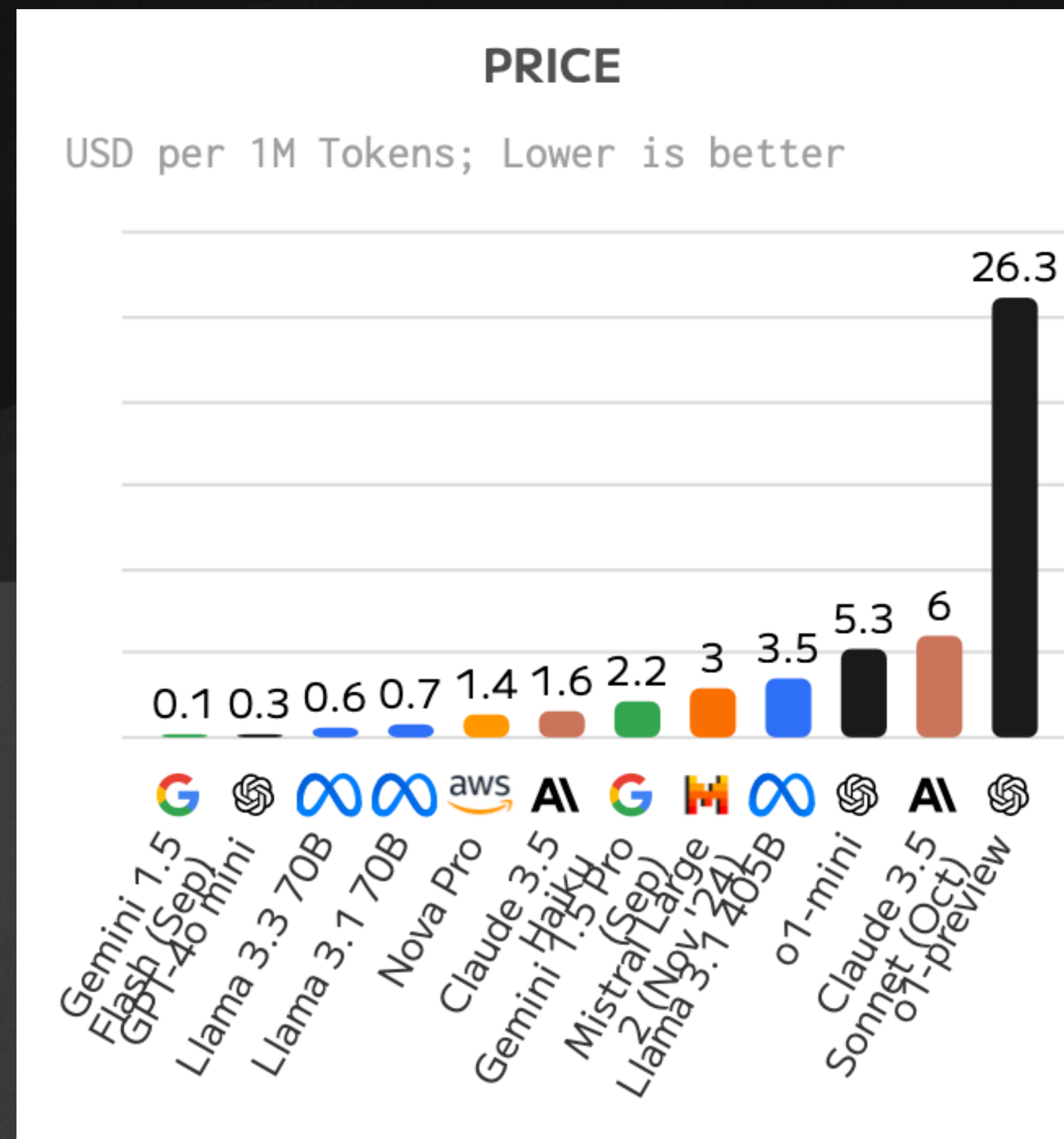# Building an AI chatbot is easy



However, validating and moderating the chatbot content is hard

# There are techniques to validate LLM output

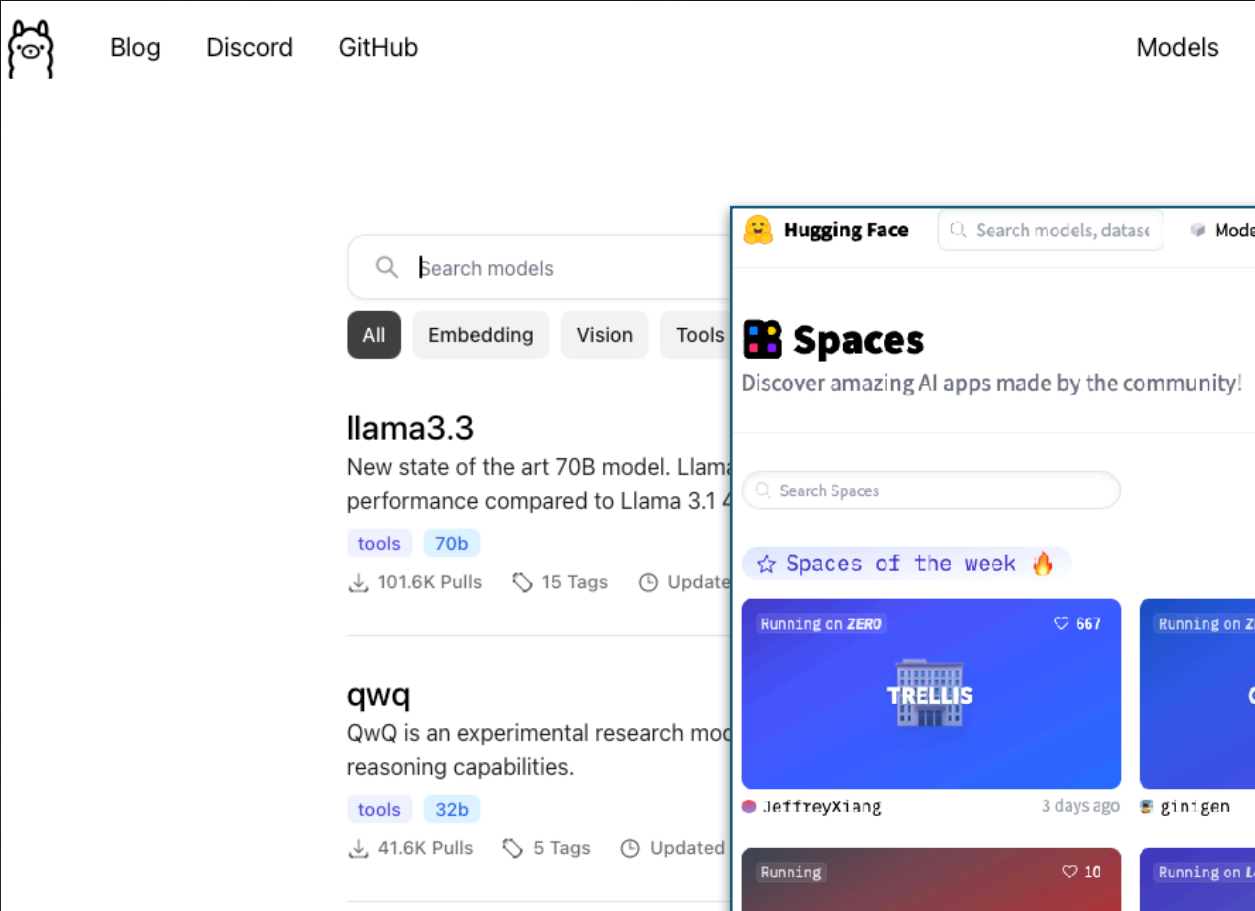| Technique | Description | Technique Issues |
|---|---|---|
| Prompt Engineering | Provide guardrails to LLM in the prompt | • Need to be very precise<br>• LLM May choose to ignore guardrails |
| Guardrails | Specialized LLMs that classifies answer as 'bad' | • Adds additional latency to every LLM call<br>• Adds cost<br>• Is not always correct |
| RAG | Provide additional context to LLM as part of the prompt | • Need to add RAG store to solution<br>• Highly dependent on the quality of the data provided<br>• RAG may not provide enough information to get answers |
| Fine tuning | Provide additional training to an LLM | • Very expensive when compared to other methods above<br>• Time consuming<br>• No guaranties plus could degrade LLM |

There are no techniques that guarantee chatbot answer is 'good'

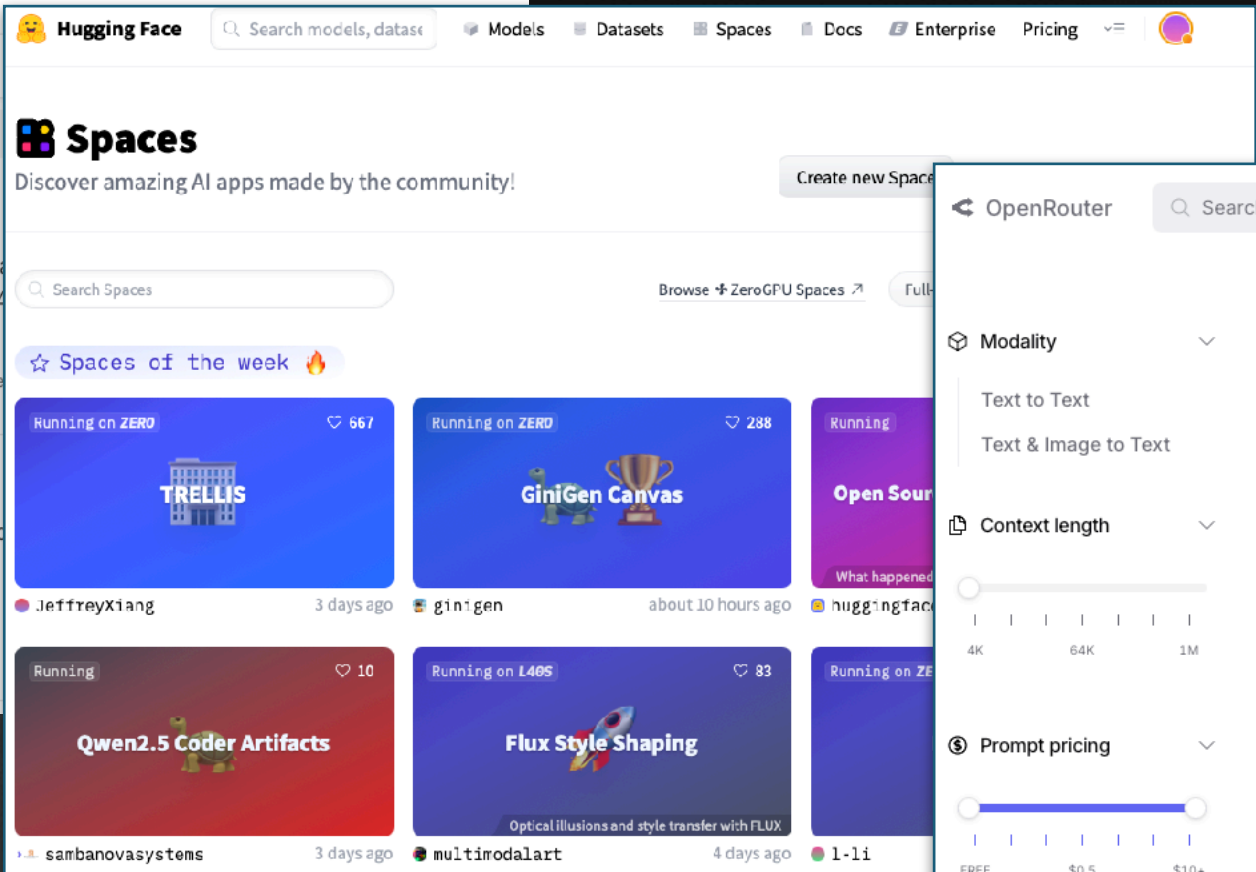# Choice of model and provider has a big impact





Small is beautiful and better for the environment (not to mention your wallet)
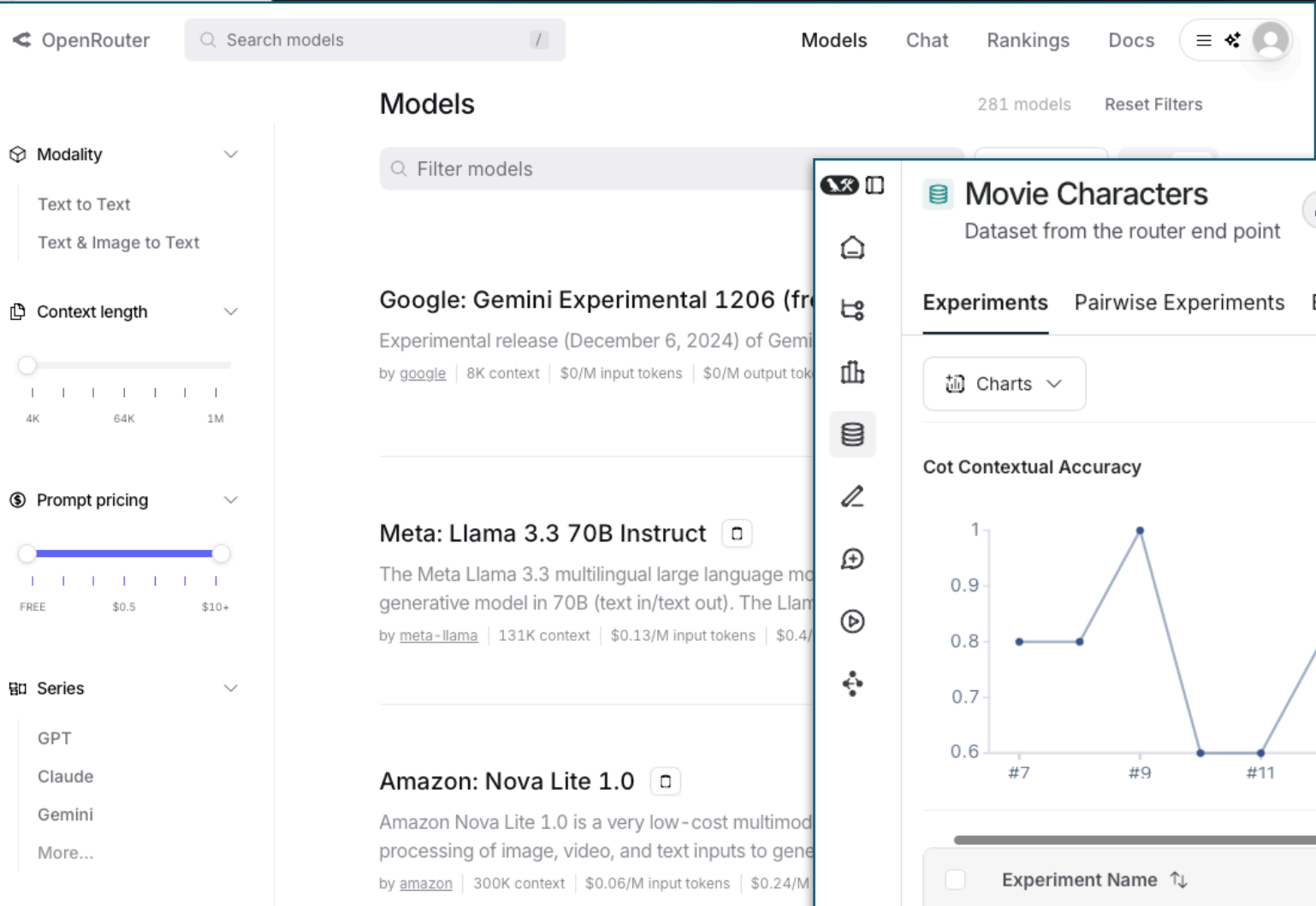
*Charts source: Artificial Analysis*

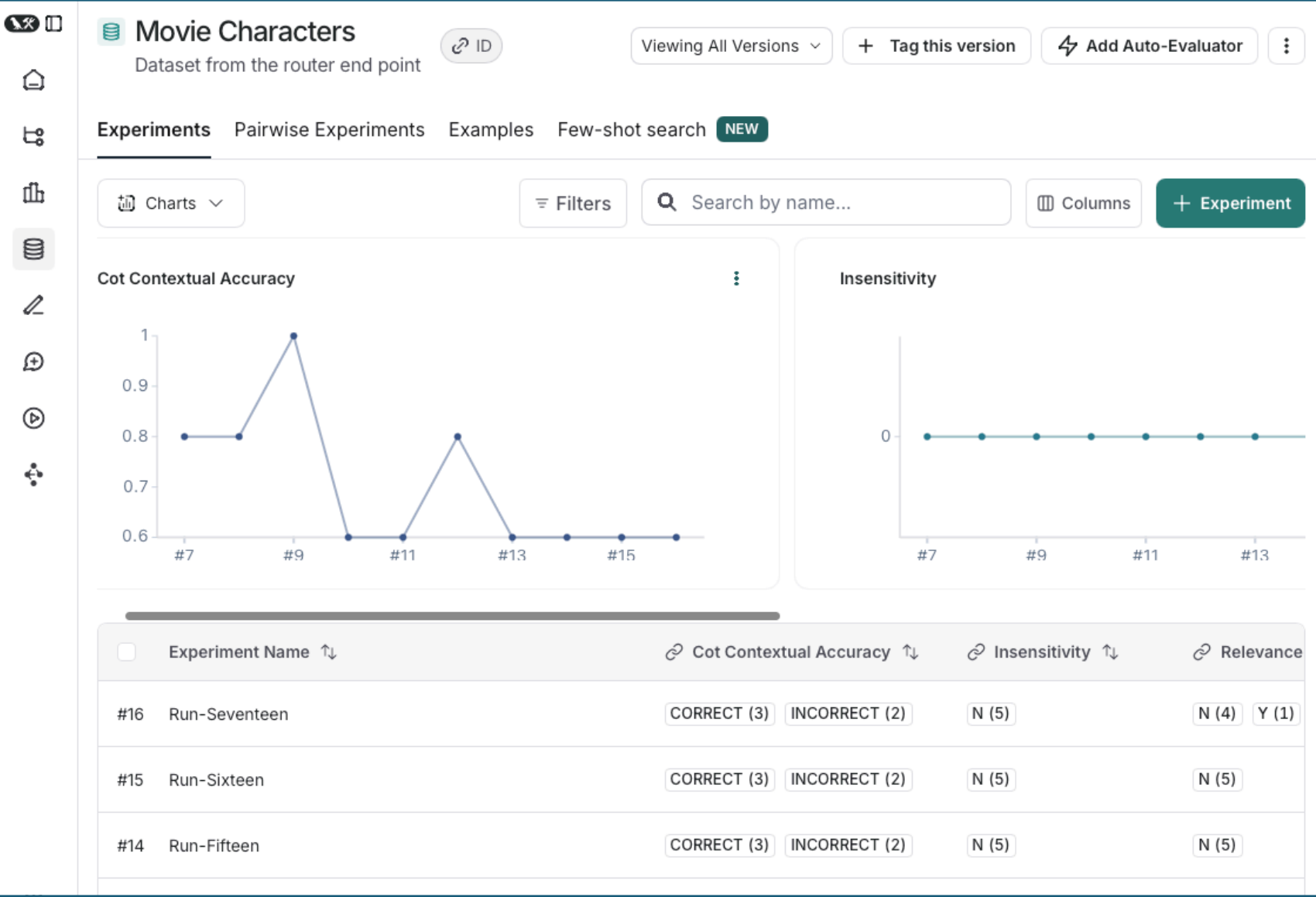# Evaluate models & choose the best for the use case
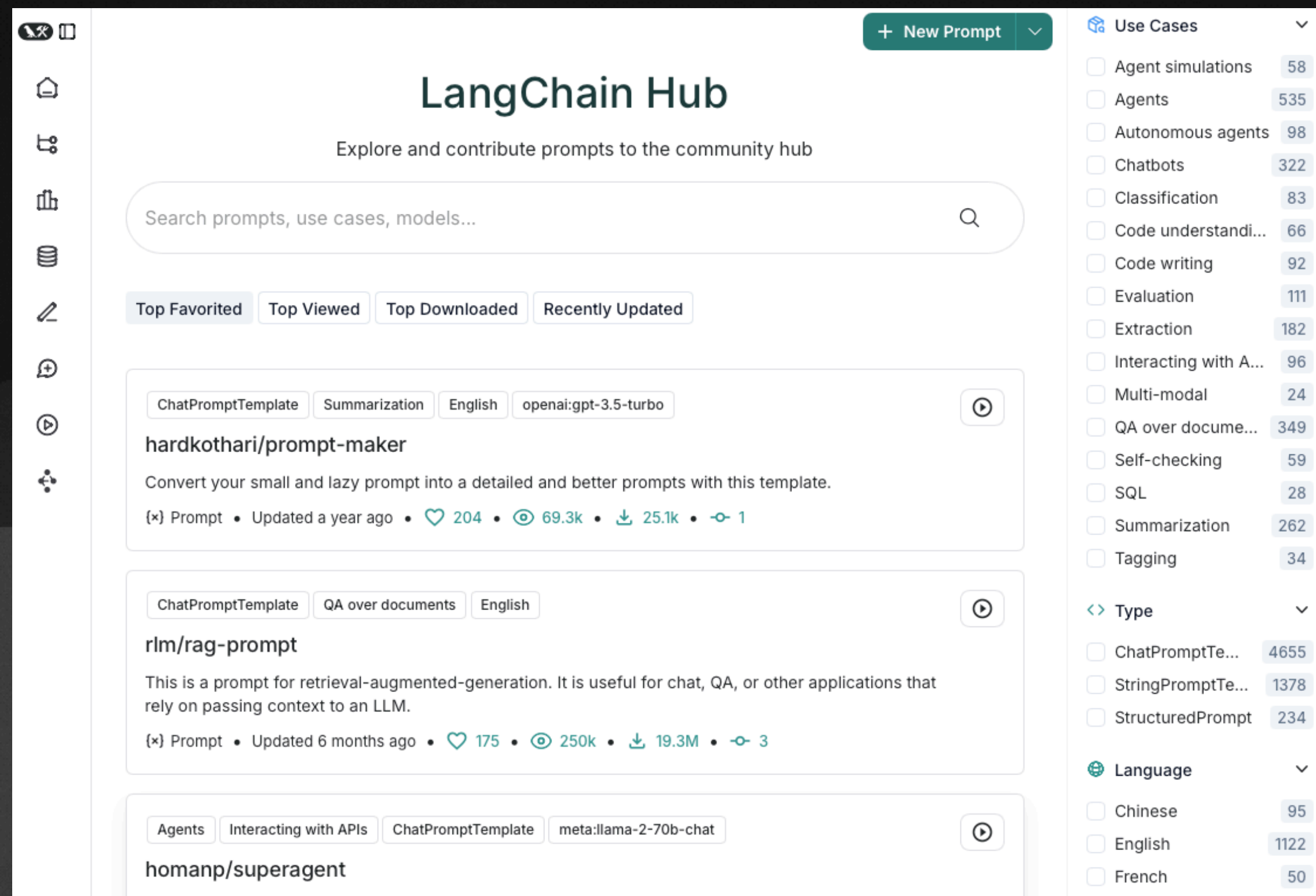


Ollama

Hugging Face

OpenRouter

LangSmith

Evaluation must take place at each step of the way
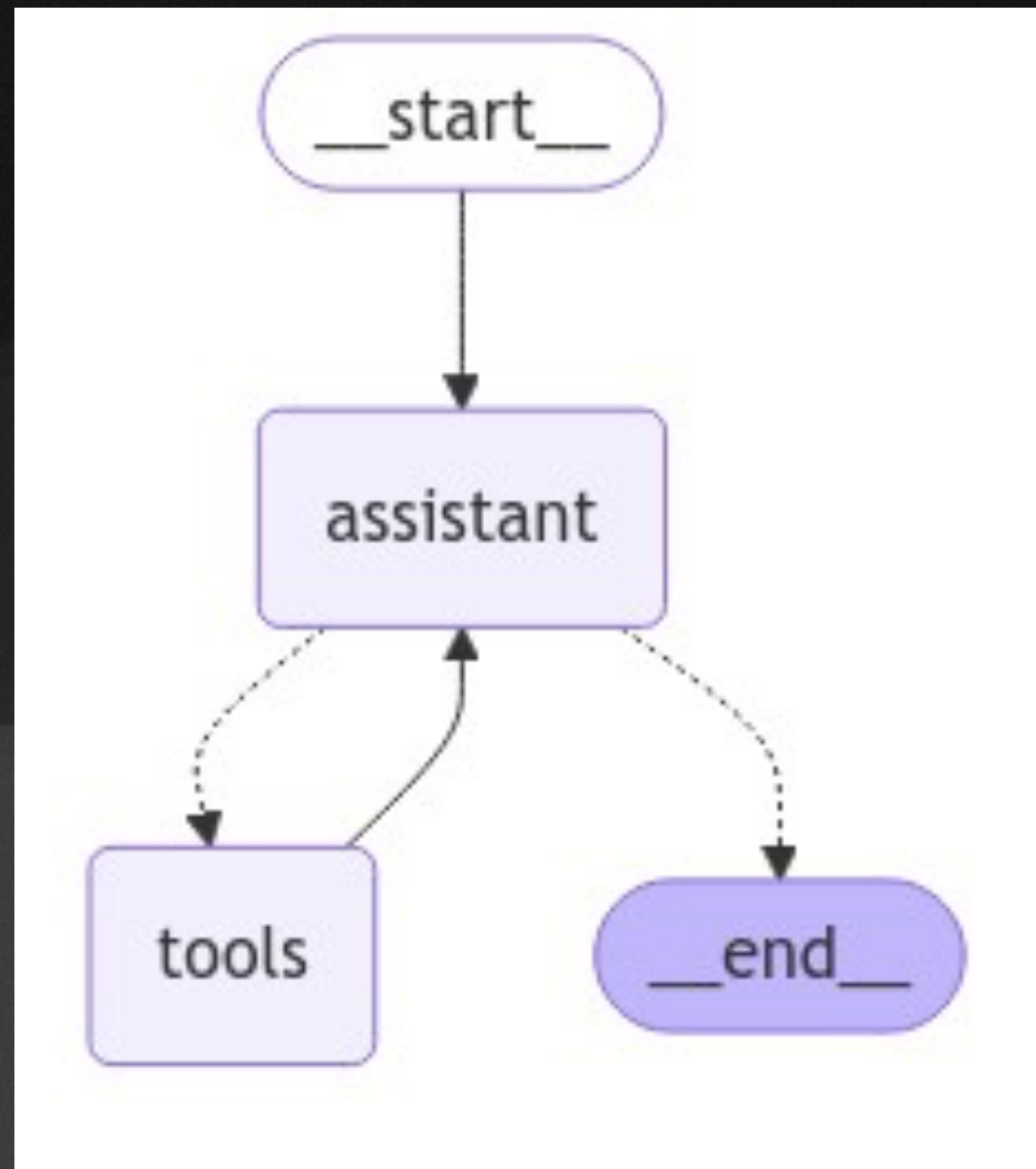
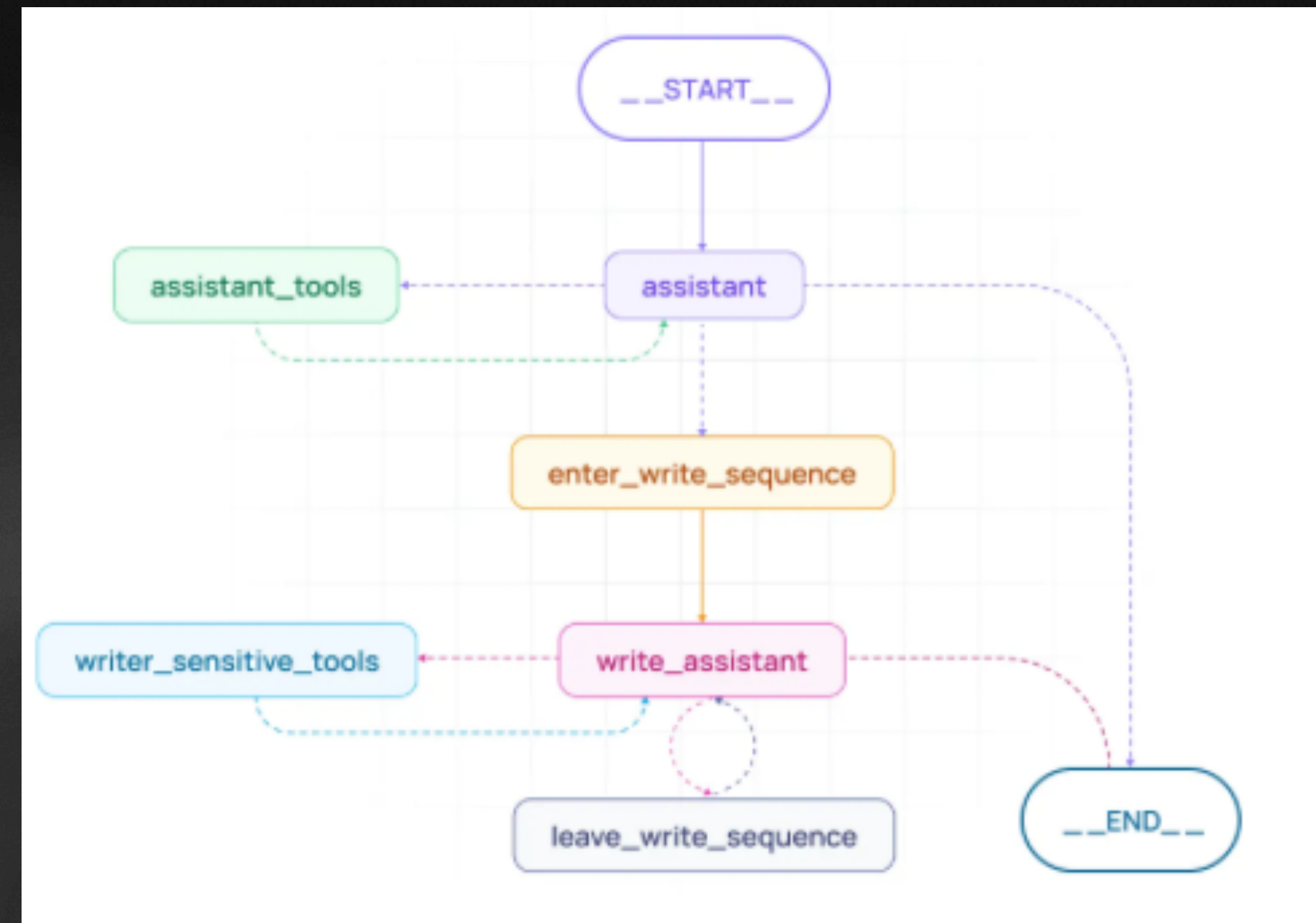# Externalize your prompts



Externalizing prompts facilitates:

- Expert input

- Future proofing

- Faster development

# Use agents to go beyond the chatbot



Simple tool calling agent



Simple assistant agent

Agents allow LLMs to interact with the real world

# Take advantage of platforms to speed up agents



- Models are the biggest source of latency

  - Take seconds to run

- Platforms like LangGraph allow complex chaining of models and code

  - Concurrent calls

  - Branching

Agents are the future of LLM development

# Get proper observability for your models



Agent run traces in LangSmith

Tracing model output, specially with agents, is hard without the necessary tooling

# Questions ?

Juan Peredo

linkedin.com/in/juanperedotech

**Thanks to AI Camp for organizing this meetup and to IBM for hosting us**