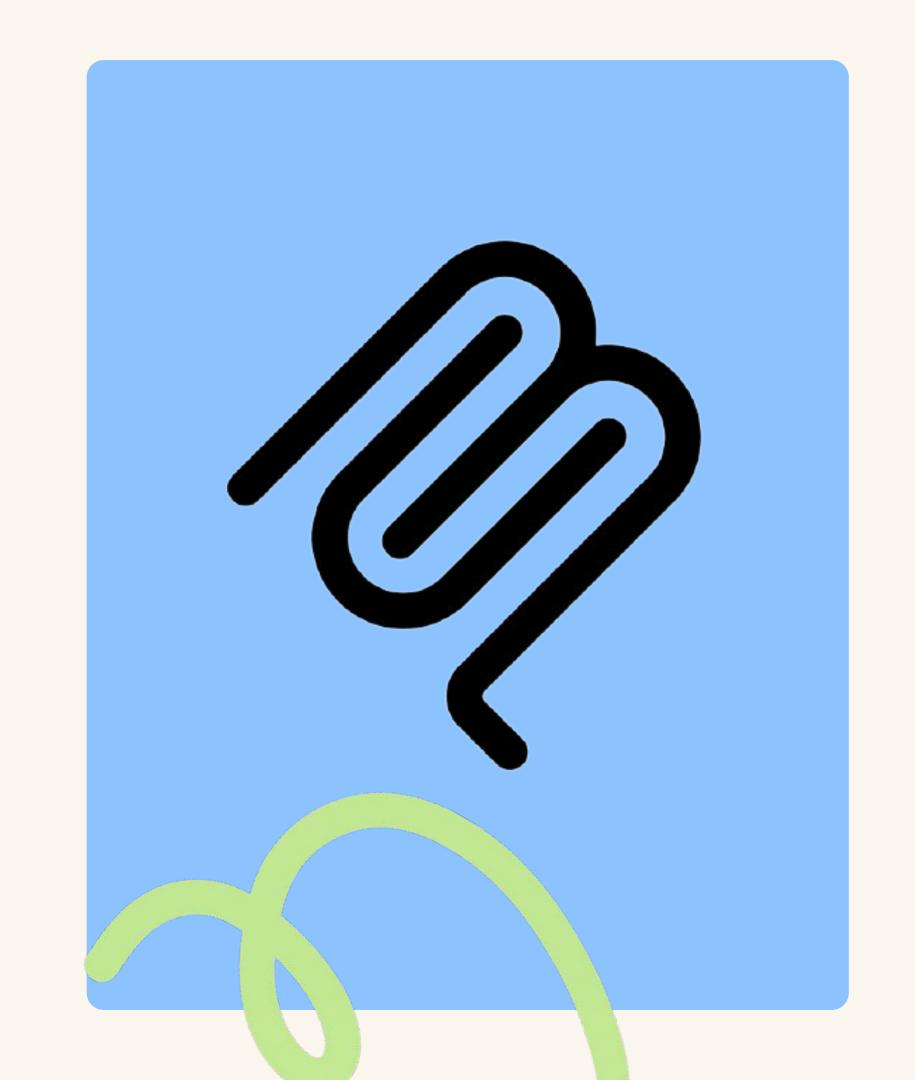
Model Context Protocol

Local

Real world interaction

A full walkthrough



What you'll learn

Local setup

Ollama & Open WebUI

MCP

Getting to know the protocol.

Own MCP server:)

Creating tools to interact with the world how we want!



Setup!

Wonder how to set the foundation?

We need:

- Tool to run models
- Nice interface



Ol Open WebUl





Ollama

CLI Tool

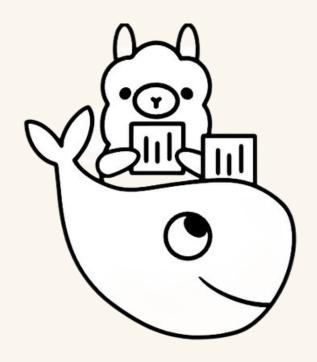
Server for local LLMs

What is Ollama?

- wraps the llama.cpp engine
- Go-based server and CLI
- exposing a simple HTTP API

Why this tool?

- usage
- setup
- community & frequent updates



Easiest way to run it

https://ollama.com/

docker run -d --gpus=all -v ollama:/root/.ollama -p 11434:11434 --name ollama ollama/ollama



Ol Open WebUl

Web Tool

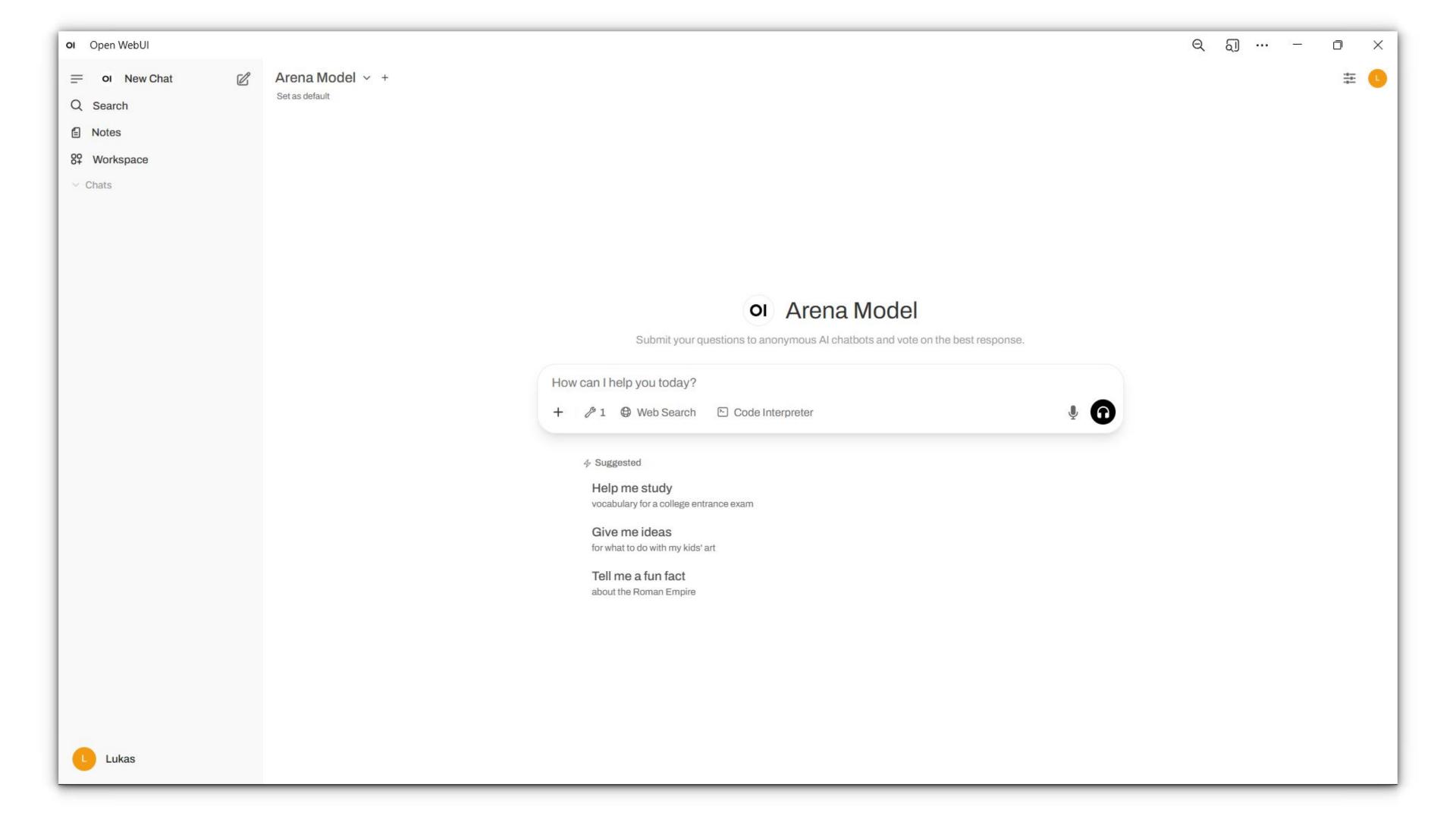
Interface for LLMs

What is Open WebUI?

- Modern Chat tool for interacting with LLMs
- Interests local & cloud LLMs

Why this tool?

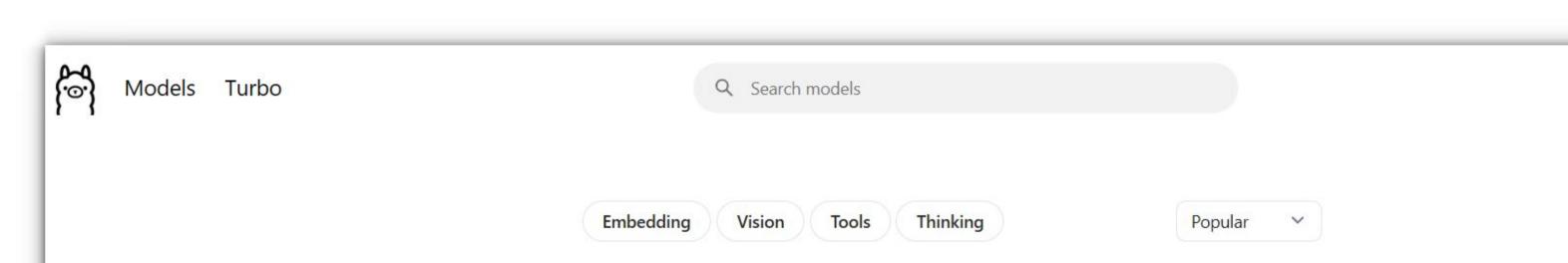
- Ease of use
- many great features
- detects Ollama directly
- Supports MCP tols:)



What is the task? The hardware? speed/accuracy tradeoff?

LMArena

Ollama Search



gpt-oss

OpenAl's open-weight models designed for powerful reasoning, agentic tasks, and versatile developer use cases.

Download

Sign in



deepseek-r1

DeepSeek-R1 is a family of open reasoning models with performance approaching that of leading models, such as O3 and Gemini 2.5 Pro.

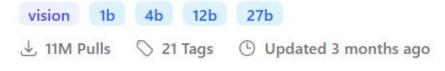
```
tools thinking 1.5b 7b 8b 14b 32b 70b 671b

1.5b 7b 8b 14b 32b 70b 671b

1.5c Topic Topic
```

gemma3

The current, most capable model that runs on a single GPU.



qwen3

☑ New Chat ☑ Leaderboard Experimental Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord ☐ New Chat	☑ New Chat ☑ Leaderboard Experimental Introducing Video Arena Generate your own videos and vote with our Discord community					
Experimental Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord 2	Experimental Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord 7		LMAr	ena v		0
Experimental Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord 3	Experimental Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord Try on Discord	Ø	New Ch	nat		
Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord 7	Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord	蛩	Leader	board		
Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord 7	Introducing Video Arena Generate your own videos and vote with our Discord community Try on Discord			(2)		
Generate your own videos and vote with our Discord community Try on Discord	Generate your own videos and vote with our Discord community Try on Discord 7	Exp	erimental			
with our Discord community Try on Discord	with our Discord community Try on Discord 7					
		-				
₩ Hide this	⊗ Hide this		nerate you	ur own vi	deos and	d vote
			nerate you n our Disc	ur own vi cord com	deos and munity	d vote

Send Feedback

nttps://lmarena.ai/leaderboard/text okies

郑 Report Bugs

22	100	1000	3	
l ead	erbo	ard O	vervie	W

Text

WebDev

Vision

Overview

See how leading models stack up across text, image, vision, and beyond. This page gives you a snapshot of each Arena, you can explore deeper insights in their dedicated tabs. Learn more about it here.

Text-to-Image

Image Edit

Search

恒 Text		O	3 days ago
Rank (UB) ↑	Model ↑↓	Score ↑↓	Votes ↑↓
1	G gemini-2.5-pro	1459	26,137
1	⊚ o3-2025-04-16	1452	32,151
2		1442	30,875
3		1438	15,271
3	p qwen3-235b-a22b-instruct-2507	1430	4,264
4	x grok-4-0709	1428	13,120
5	A\ claude-opus-4-20250514-think	1420	18,063
5	kimi-k2-0711-preview	1420	11,809

WebDev ○ View →			∵View →
Rank (UB) ↑	Model ↑↓	Score ↑↓	Votes ↑↓
1	G Gemini-2.5-Pro	1403	6,813
1	▼ DeepSeek-R1-0528	1389	4,441
2	A\ Claude Opus 4 (20250514)	1380	8,708
2	Z GLM-4.5	1360	722
3	ℤ GLM-4.5-Air	1347	752
4	☆ Qwen3-Coder	1361	6,119
4	A\ Claude Sonnet 4 (20250514)	1359	7,913
4	A\ Claude 3.7 Sonnet (20250219)	1358	7,460

Image-to-Video

Text-to-Video

Copilot

Start Voting

Rules of Thumb

Model-size VRAM

~7 B parameters → ~8 GB VRAM ~13 B parameters → ~16 GB VRAM

Quantization q8

Size: roughly 50 % Speed: ~1.56 × faster Quality loss: minimal

Quantization q4

Size: roughly 25 %

Speed: up to 2.4 × faster

Quality loss: typically negligible

Testing!

At the end, it really comes down to test out the best options.





Ol Open WebUl







Model Context Protocol

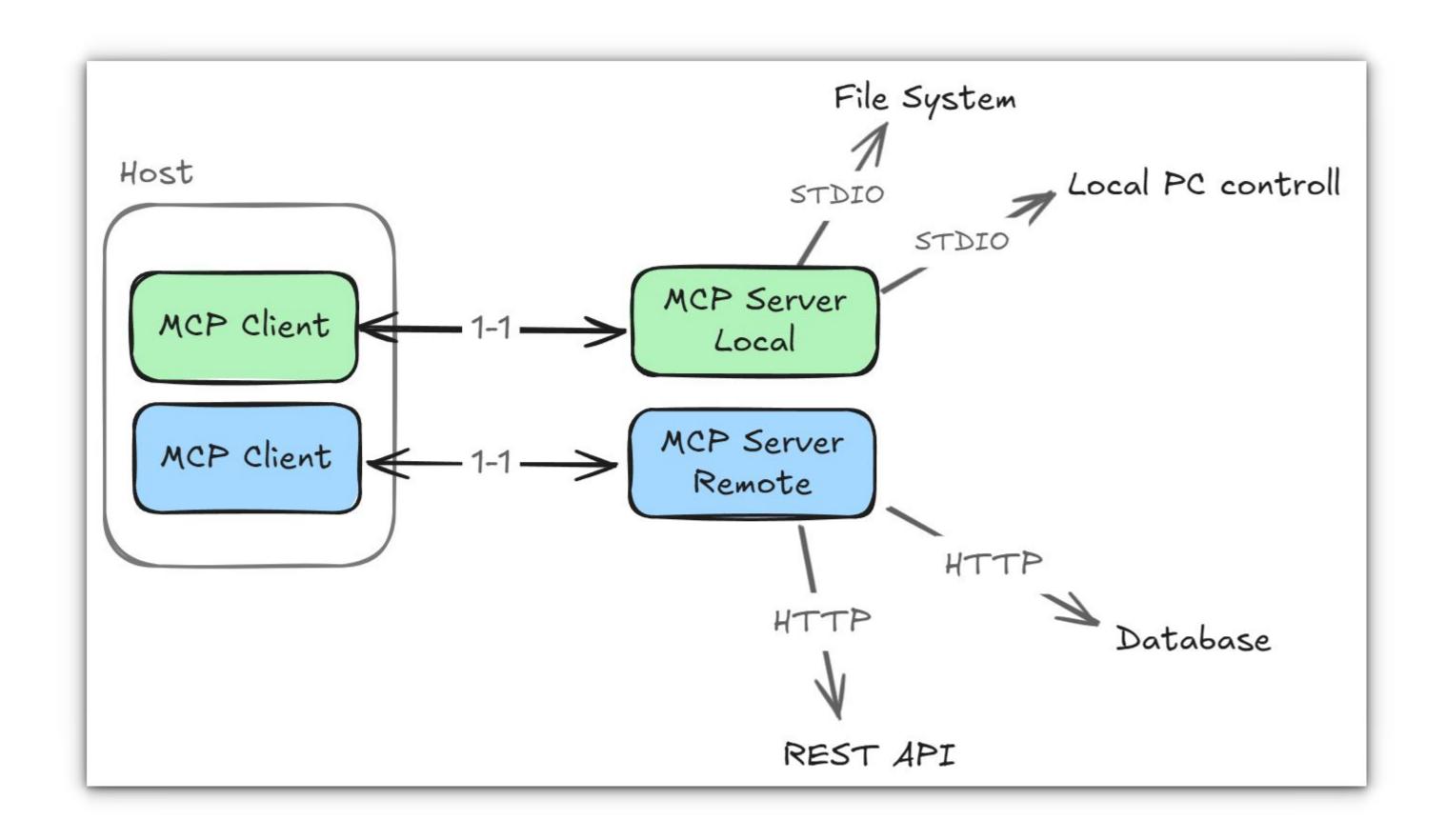
Wonder how why we need it?

- Open interoperability standard
- Client-server architecture
- Rich primitives with human-in-the-loop



Client-server architecture

Simple design where both Server and Client can implement functionalities and communicate bidirectional.



Primitives

The MCP defines 3 Primitives, with all of them serving different purposes.

Primitives

Tools

- Executable functions
- Provided to the model
- model controlled

Resources

- Data sources that provide contextual information
- Client controlled

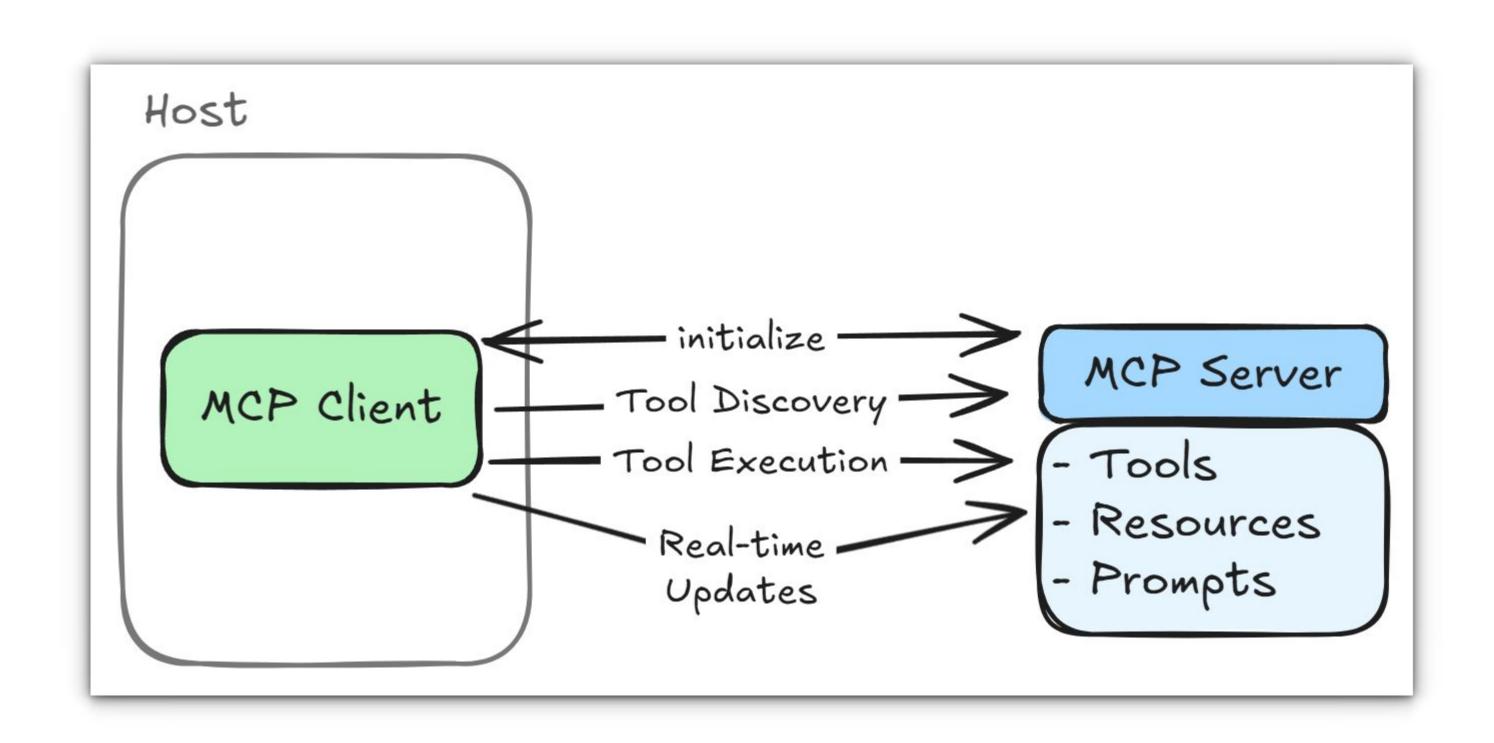
Prompts

- Reusable templates
- system prompts
- few-shot examples



Primitives

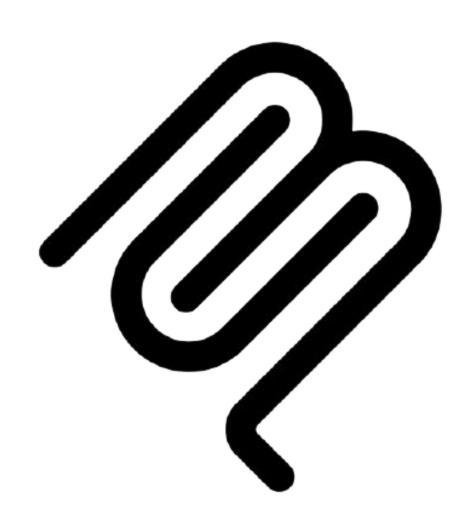
The MCP defines 3 Primitives, with all of them serving different purposes.

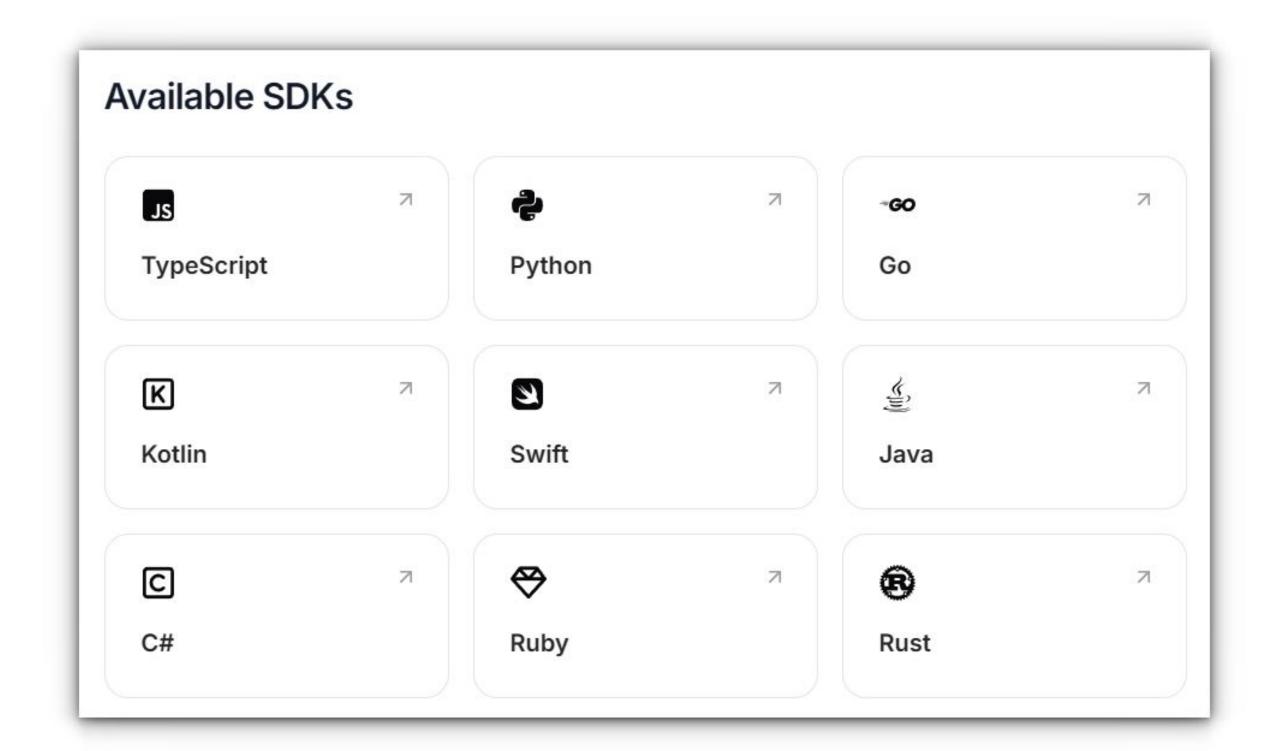


Own MCP

Wonder how why we need it?

- Open interoperability standard
- Client-server architecture
- Rich primitives with human-in-the-loop





```
from mcp.server.fastmcp import FastMCP
from pynput.keyboard import Key, Controller
import time
# Initialize the FastMCP server
mcp = FastMCP("LazyPresentator")
# Create a keyboard controller instance
keyboard = Controller()
# Placeholder :)
if __name__ == "__main__":
    mcp.run()
```

```
@mcp.tool()
def press_n_key(delay_seconds: int = 5) -> str:
   For lazy presentators: Automatically presses the 'n' key to advance to the next slide
    after a specified delay.
    This is perfect for presentators who are too lazy to press the 'n' key themselves.
    The default delay gives you time to finish speaking before advancing.
    Args:
        delay_seconds: Number of seconds to wait before pressing the 'n' key (default: 2)
   Returns:
        A confirmation message about the slide advancement.
    11 11 11
    try:
        # Press and release the 'n' key
        time.sleep(delay_seconds)
        keyboard.press('n')
        keyboard.release('n')
        # Add a small delay to ensure the keypress is registered
        time.sleep(0.1)
        return "Successfully pressed the 'n' key. Screen should be blanked/unblanked if in
presentation mode."
    except Exception as e:
        return f"Error pressing 'n' key: {str(e)}"
```

mcpo --port 8000 -- python .\lazy_presentator.py



Altook over!

Demo

Now we are ready to get into the demo how easy it is to interact with the **Open Data Hub** via our new MCP Server skills!

OPEN DATA HUB

Thank you! Happy hacking!

Article

Model Context Protocol | Lukas's ὑπόμνημα



