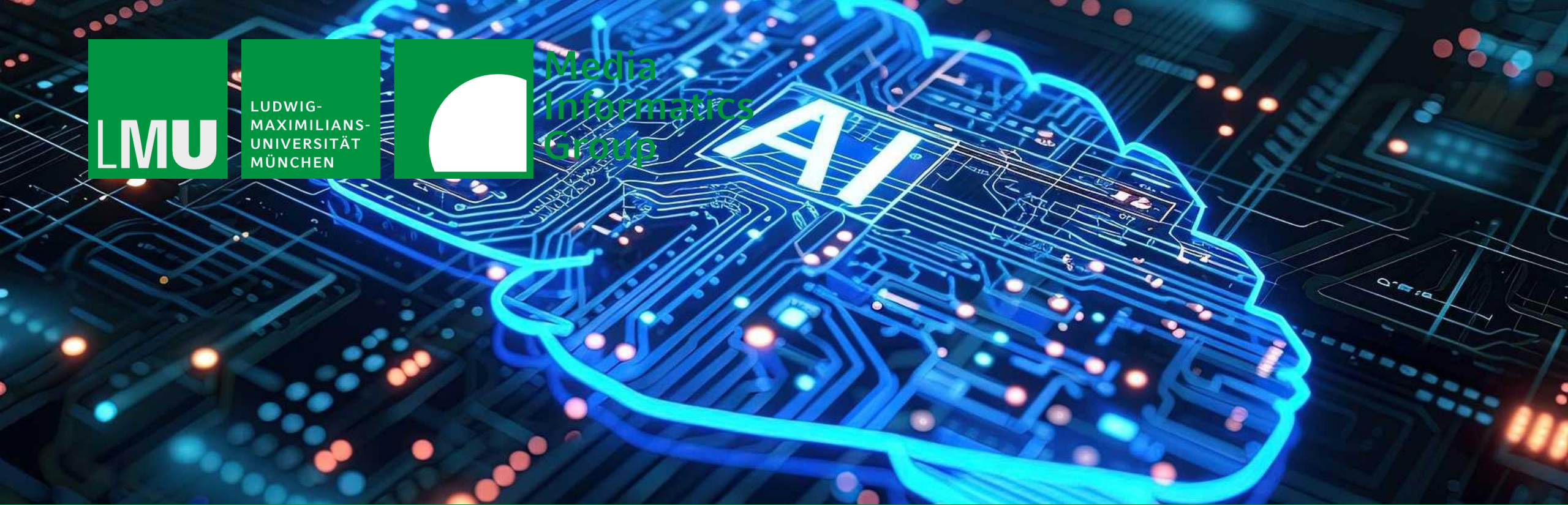




LUDWIG-  
MAXIMILIANS-  
UNIVERSITÄT  
MÜNCHEN



Media  
Informatics  
Group



# Local AI Models with LM Studio

Philipp Thalhammer

# Why are local LLMs a Good Idea?

- Keep control of your data
- No (less) censorship
- Works offline
- Cheap
- Great for prototyping

# But...

- Not as fast
- Worse performance



# What You (of Course) Already Did

- Download LM Studio
  - <https://lmstudio.ai/>

# What You (of Course) Already Did

- Download a *small* model
  - E.G., *meta-llma-3.1-8b-instruct*

LM Studio - 0.3.14 (Build 5)

My Models

Models DirectoryC:\Users\Philipp\lmstudio\models

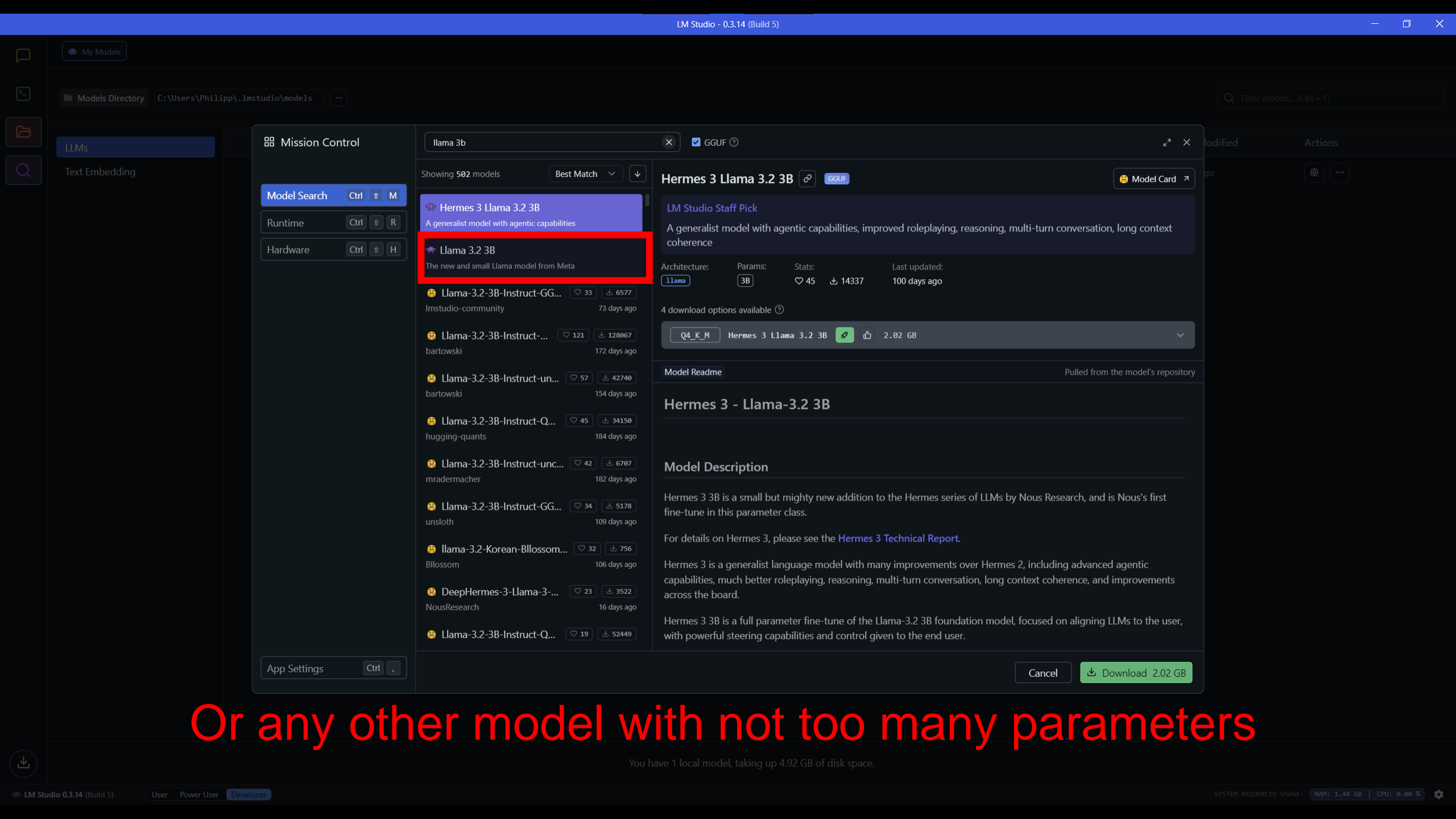
Filter models... (Ctrl + F)

LLMs

Arch	Params	Publisher	Model	Quant	Size	Date Modified	Actions
llama	8B	lmstudio-community	meta-llama-3.1-8b-instruct	Q4_K_M	4.92 GB	2 days ago	

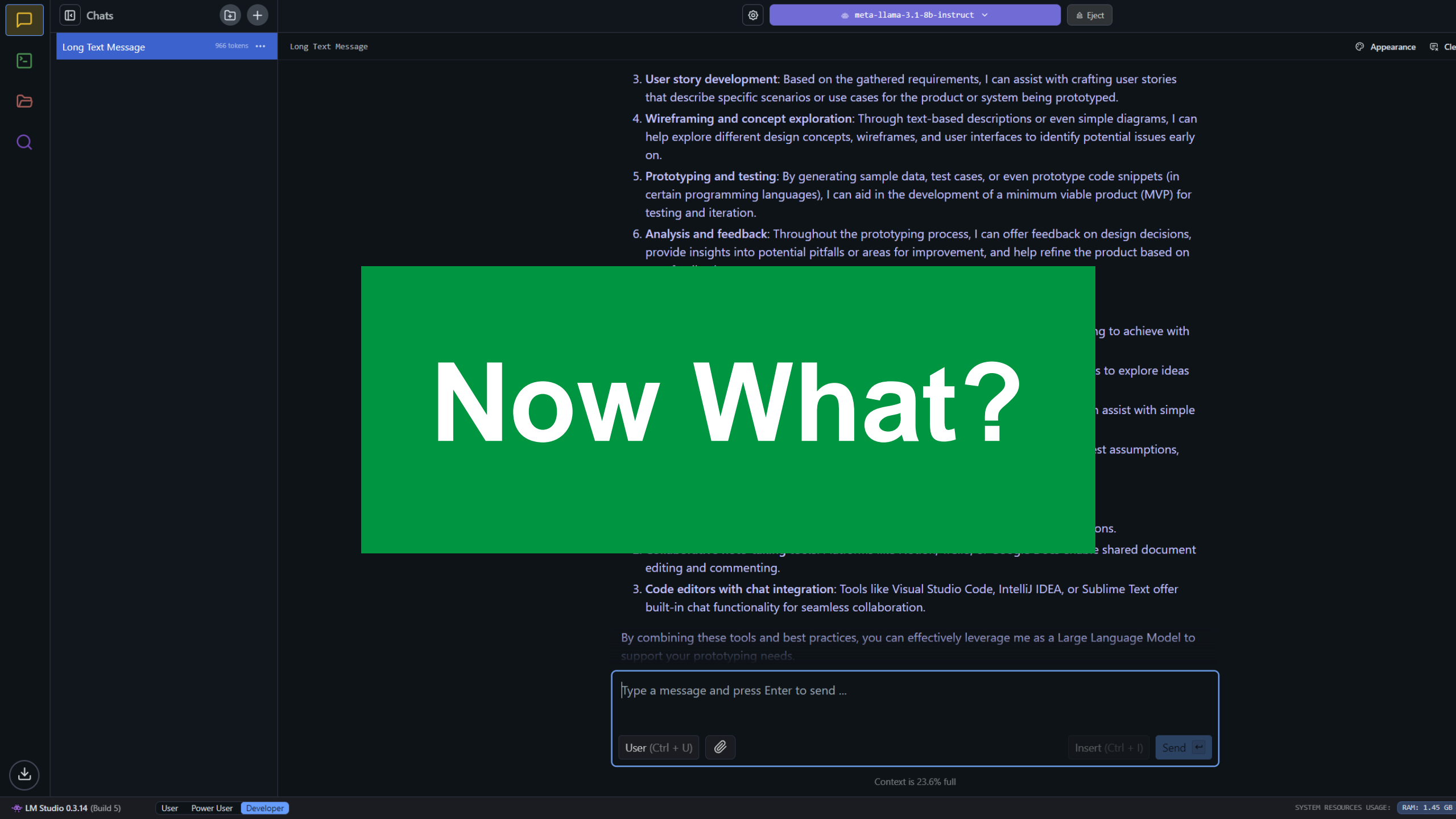
You have 1 local model, taking up 4.92 GB of disk space.

LM Studio 0.3.14 (Build 5)UserPower UserDeveloperSYSTEM RESOURCES USAGE: RAM: 1.44 GB | CPU: 0.00 %



Or any other model with not too many parameters





Now What?

**Let's set up  
a server!**

Chats

Unnamed Chat

0 tokens ...

Unnamed Chat

Appearance Clear All Duplicate

## Mission Control

Model Search Ctrl ↑ M

Runtime Ctrl ↑ R

Hardware Ctrl ↑ H

App Settings

Ctrl ,

## App Settings

☐ Use ↑ + ← to send message☒ Use Ctrl + R to regenerate the last message in chat

## Chat AI Naming

Auto

Decides whether to create names based on generation speed

## Local LLM Service (headless)

Use LM Studio's LLM server without having to keep the LM Studio application open

☒ Enable Local LLM Service ⓘ

## Developer

☒ JIT models auto-evict: ensure at most 1 model is loaded via JIT at any given time (unloads previous model)☐ Show debug info blocks in chat☐ Enable model load configuration support in presets ⓘ Experimental

## LM Studio Extension Packs Download Channel

Stable

☐ When applicable, separate reasoning\_content and content in API responses ⓘ Experimental

## Onboarding Hints

Dismissed Onboarding Hints

User (Ctrl + U)



Insert (Ctrl + I)

Send ↵

Context is 0.0% full





Status: Running ⏻

Settings

Reachable at: <http://192.168.178.30:1234> 🔗

READY

meta-llama-3.1-8b

Server Port ⓘ 1234

Enable CORS ⓘ ⏻

Serve on Local Network ⓘ ⏻

Just-in-Time Model Loading ⓘ ⏻

Auto unload unused JIT loaded models ⓘ ⏻

Max idle TTL 60 minutes

Only Keep Last JIT Loaded Model ⓘ ⏻

Size 4.92 GB

Eject

Supported endpoints (OpenAI-like) ▾

- GET /v1/models ⓘ
- POST /v1/chat/completions ⓘ
- POST /v1/completions ⓘ
- POST /v1/embeddings ⓘ

Developer Logs



```
    "finish_reason": "length",
    "message": {
      "role": "assistant",
      "content": "A man walked into a library and asked the librarian, \"Do you have any books on Pavlov's dogs and Schrödinger's cat?\" The librarian replied, \"It rings a bell, but I'm not sure if it's\"
    }
  ],
  "usage": {
    "prompt_tokens": 40,
    "completion_tokens": 49,
    "total_tokens": 89
  },
  "stats": {},
  "system_fingerprint": "meta-llama-3.1-8b-instruct"
}
```

2025-03-28 15:03:06 [INFO] [LM STUDIO SERVER] Client disconnected. Stopping generation... (If the model is busy processing the prompt, it will finish first.)

Info Inference Load

Model Information ▾

Model  
lmstudio-community/Meta-Llama-3.1-8B-Instruct-G...

File  
Meta-Llama-3.1-8B-Instruct-Q4\_K\_M.gguf

Format  
GGUF

Quantization  
Q4\_K\_M

Arch  
llama Trained for Tool Use

Domain  
llm

Size on disk  
4.92 GB

API Usage ▾

This model's API identifier  
meta-llama-3.1-8b-instruct 🔗

☒ The local server is reachable at this address  
<http://192.168.178.30:1234> 🔗

Tool Use ▾

This model is detected to have been trained for tool use

[Open quick docs](#) for more information

Minimize to Tray

LM Studio Server: Not Running

No Models Loaded

2 → Start Server on Port 1234...

1 → Load Model >

Quit LM Studio

Strg + Q

# HTML Example

## Local AI Chat

Type your prompt here...

Send





**Why Does It Talk  
Like That?**





**Try It**  
**Yourself!**

# Resources

<https://github.com/Phlipinator/Local-LLM-Examples>



# Conclusion

# Questions?

## Contact

Philipp Thalhammer  
LMU Munich  
[philipp.thalhammer@ifi.lmu.de](mailto:philipp.thalhammer@ifi.lmu.de)