# Setup

- Download LM Studio

  - https://lmstudio.ai/

- Download a *small* model

  - E.G., *meta-llma-3.1-8b-instruct*

My Models

Models Directory    C:\Users\Philipp\.lmstudio\models    ...

Filter models... (Ctrl + F)

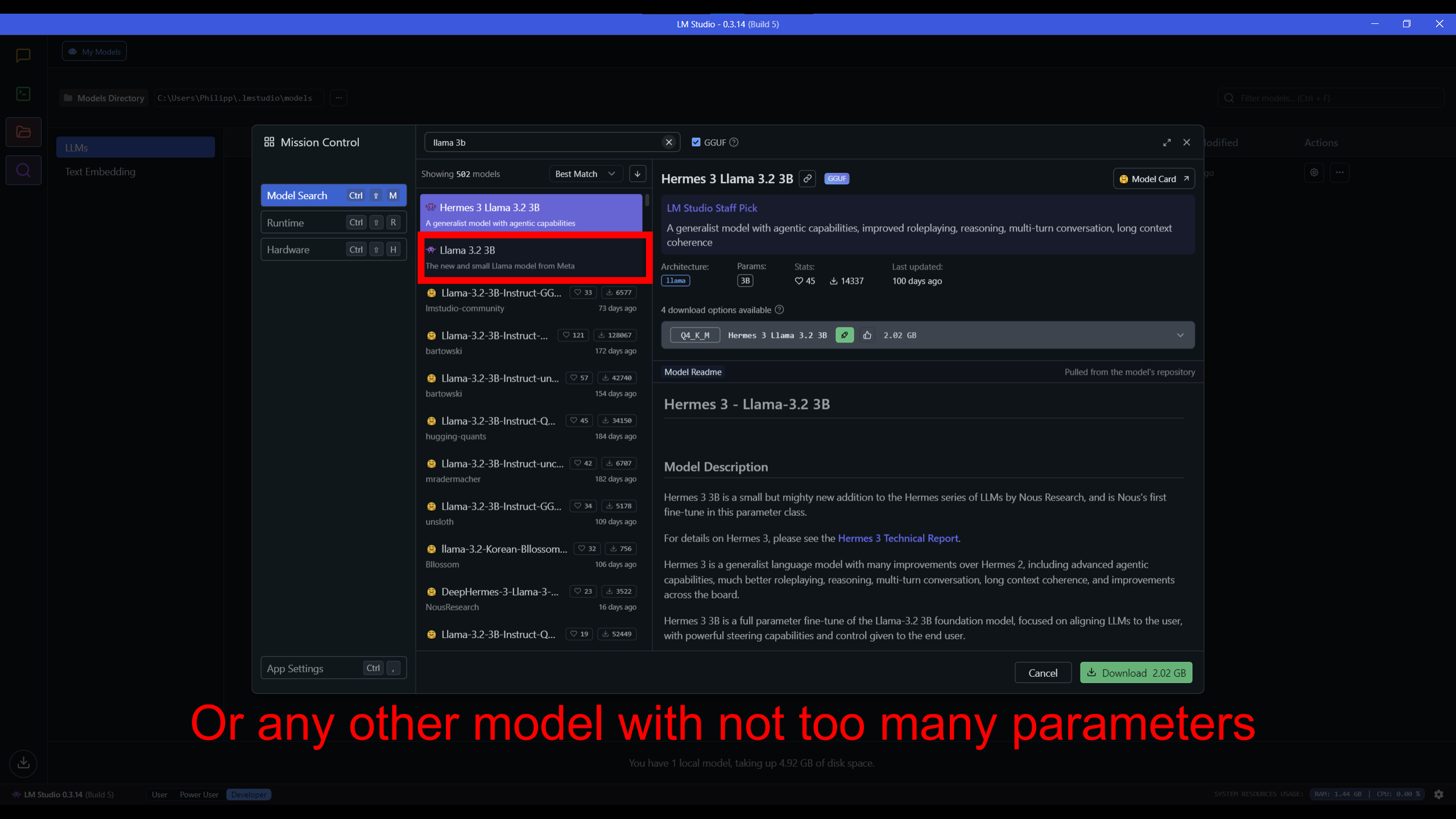| | Arch | Params | Publisher | Model | Quant | Size | Date Modified | Actions |
|---|---|---|---|---|---|---|---|---|
| LLMs | llama | 8B | lmstudio-community | meta-llama-3.1-8b-instruct | Q4_K_M | 4.92 GB | 2 days ago | |

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

Or any other model with not too many parameters

# Set up a
# server!

Chats

Unnamed Chat                    0 tokens

Unnamed Chat                    Appearance    Clear All    Duplicate

meta-llama-3.1-8b-instruct    Eject

**Mission Control**                    **App Settings**

Model Search    Ctrl ⇧ M

Runtime    Ctrl ⇧ R

Hardware    Ctrl ⇧ H

☐ 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  ⓘ    ← **2**

## 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**

App Settings    Ctrl ,

User (Ctrl + U)                    Insert (Ctrl + I)    Send ↵

Context is 0.0% full

**1**

LM Studio 0.3.14 (Build 5)    User    Power User    Developer    SYSTEM RESOURCES USAGE:  RAM: 1.44 GB    CPU: 0.00 %

LM Studio - 0.3.14 (Build 5)

LM Studio    LM Runtimes

Select a model to load  (Ctrl +L)

Community    Quick Docs

Status: Running    Settings

Reachable at:  http://192.168.178.30:1234

Info    Inference    Load

Server Port    1234

READY

**1**

Enable CORS    ⬅ **2**

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

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

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.)

LM Studio 0.3.14 (Build 5)    User  Power User  Developer

SYSTEM RESOURCES USAGE:  RAM: 1.45 GB  CPU: 0.00 %

# HTML Example



**Local AI Chat**

Type your prompt here...

Send

Why Does It Talk Like That?

Welcome    webinterface.html ✕    micropython.py 7

LLM STUDIO TEST
  micropython.py                7
  webinterface.html

```html
  2     <html lang="en">
  3     <head>
  7         <style>
 46             }
 47         </style>
 48     </head>
 49     <body>
 50         <div class="container">
 51             <h1>Local AI Chat</h1>
 52             <textarea id="userInput" placeholder="Type your prompt here..."></
                    textarea>
 53             <button onclick="sendRequest()">Send</button>
 54             <div id="response" class="response"></div>
 55         </div>
 56
 57         <script>
 58             async function sendRequest() {
 59                 const userInput = document.getElementById('userInput').value;
 60                 const responseDiv = document.getElementById('response');
 61
 62                 responseDiv.textContent = 'Processing...';
 63
 64                 try {
 65                     const response = await fetch('http://localhost:1234/v1/chat/
                        completions', {
 66                         method: 'POST',
 67                         headers: { 'Content-Type': 'application/json' },
 68                         body: JSON.stringify({
 69                             model: "local-model",
 70                             messages: [{ role: "user", content: userInput },
 71                                 { role: "system", content: "Only answer in
                                    rhymes" }],
 72                             max_tokens: 200
 73                         })
 74                     });
 75
 76                     if (!response.ok) {
 77                         throw new Error(`HTTP error! status: ${response.status}`);
 78                     }
 79
 80                     const data = await response.json();
 81                     responseDiv.textContent = data.choices[0].message.content.trim();
 82                 } catch (error) {
 83                     responseDiv.textContent = `Error: ${error.message}`;
 84                 }
 85             }
 86         </script>
 87     </body>
 88 </html>
 89
```

OUTLINE
TIMELINE
SVN
VS CODE PETS

⊗ 3  ⚠ 4  ⓘ 12    Live Share                                    UTF-8   HTML         Go Live

# Now, let's try it with a microcontroller!

LM Studio    LM Runtimes

Select a model to load  (Ctrl +L)

Community    Quick Docs

Status: Running    Settings

Reachable at: http://192.168.178.30:1234

Info    Inference    Load

READY

llm meta-llama-3.1-8b-instruct

cURL

Size 4.92 GB    Eject

**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

Supported endpoints (OpenAI-like)

GET /v1/models ?

POST /v1/chat/completions ?

POST /v1/completions ?

POST /v1/embeddings ?

Developer Logs

```
                    "system_fingerprint": "meta-llama-3.1-8b-instruct"
                }
2025-03-31 13:27:56  [INFO] [LM STUDIO SERVER] Client disconnected. Stopping generation... (If the model is busy processing the prompt, it will finish first.)
2025-03-31 14:13:36  [INFO] Server stopped.
2025-03-31 14:42:24  [INFO] [LM STUDIO SERVER] Success! HTTP server listening on port 1234
2025-03-31 14:42:24  [WARN] [LM STUDIO SERVER] Server accepting connections from the local network. Only use this if you know what you are doing!
2025-03-31 14:42:24  [INFO]
2025-03-31 14:42:24  [INFO] [LM STUDIO SERVER] Supported endpoints:
2025-03-31 14:42:24  [INFO] [LM STUDIO SERVER] -> GET http://192.168.30:1234/v1/models
2025-03-31 14:42:24  [INFO] [LM STUDIO SERVER] -> POST http://192.168.30:1234/v1/chat/completions
2025-03-31 14:42:24  [INFO] [LM STUDIO SERVER] -> POST http://192.168.30:1234/v1/completions
2025-03-31 14:42:24  [INFO] [LM STUDIO SERVER] -> POST http://192.168.30:1234/v1/embeddings
2025-03-31 14:42:24  [INFO]
2025-03-31 14:42:24  [INFO] [LM STUDIO SERVER] Logs are saved into C:\Users\Philipp\.lmstudio\server-logs
2025-03-31 14:42:24  [INFO] Server started.
2025-03-31 14:42:24  [INFO] Just-in-time model loading active.
```

LM Studio 0.3.14 (Build 5)

User    Power User    Developer

SYSTEM RESOURCES USAGE:    RAM: 1.36 GB    CPU: 0.00 %

basic.py ×

```python
1   import time
2
3   import network
4   import ujson
5   import urequests
6
7   # ====== USER CONFIGURABLE VARIABLES ======
8   SSID = ""  # WiFi SSID
9   PASSWORD = ""  # WiFi Password
10
11  # API endpoint URL
12  API_URL = "http:XXX.XXX.XXX.XX:1234/v1/chat/completions"  # Replace with your local IP
13
14  # Request data
15  REQUEST_DATA = {
16      "model": "XXXX",  # Replace with your model name
17      "messages": [{"role": "user", "content": "Tell me a joke!"}], # For Reference: https://platform.openai.com/docs/api-reference/responses/create
18      "max_tokens": 50
19  }
20
21  # =========================================
22
23  def connect_wifi():
24      wlan = network.WLAN(network.STA_IF)
25      wlan.active(True)
26      wlan.connect(SSID, PASSWORD)
27      while not wlan.isconnected():
28          time.sleep(1)
29      print("Connected to WiFi:", wlan.ifconfig())
30
31  def send_request():
32
33
34      headers = {
35          "Content-Type": "application/json"
36      }
37
38      try:
39          print("Sending request to:", APT URL)
```

Shell ×

```
Python 3.10.11 (/Users/philipp/Applications/Thonny.app/Contents/Frameworks/Python.framework/Versions/3.10/bin/python3.10)
>>> %cd /Users/philipp/Documents
>>> %cd '/Users/philipp/Documents/Git Repos'
>>> %cd '/Users/philipp/Documents/Git Repos/Local-LLM-Examples'
```