

Project Report – Majji Yaswanth Sai

Project 1: Building and Deploying a Conversational Chatbot with Llama

Objective: The goal was to set up and deploy a conversational chatbot using the Meta-Llama 3.1 model, Docker, WasmEdge, and Flask API, creating a high-performance, portable system capable of handling natural language inference and interaction.

1. Environment Setup

- **Docker Container:** Ubuntu 22.04 was used as the base environment for deploying the model.
- **Model:** The Meta-Llama-3.1-8B-Instruct model was downloaded from Hugging Face.
- **Execution:** WasmEdge was employed for running the Llama model in the Docker container. This enabled high-performance inference using the GGUF model format.

2. Model Execution and Testing

- The Llama model was successfully initialized and ran a prompt to generate responses, showcasing its ability to output coherent text based on input prompts.
- A sample prompt on Robert Oppenheimer's achievements generated a detailed, informative response, proving the model's capability for natural language

understandin

```
root@3ccf065cc683: ~/llama- x + v
2. In what field did Oppenheimer make his most significant contributions?
A) Quantum mechanics and astrophysics.
B) Particle physics and cosmology.
C) Nuclear physics and engineering.
D) Condensed matter physics and materials science.

3. What was one of the first areas of study that Oppenheimer recognized as important in quantum mechanics?
A) Wave-particle duality.
B) Quantum field theory.
C) Quantum computing.
D) Quantum gravity.

4. What was another area of study in which Oppenheimer made important contributions?
A) Black holes.
B) Dark matter and dark energy.
C) The early universe.
D) The behavior of subatomic particles.

5. How did Oppenheimer's experiences during World War II shape his views on the use of nuclear weapons?
A) He became a strong supporter of nuclear proliferation.
B) He became a vocal advocate for nuclear disarmament.
C) He became a critic of the use of nuclear weapons, but also supported their development.
D) He had no opinion on the matter.

Please let me know if you'd like me to rephrase any of the questions or if you have any other requests!
Let me know if you'd like me to rephrase any of the questions or if you have any other requests!
Here are
root@3ccf065cc683:~/llama-models#
```

3. Web Service API Setup

- **Flask API:** A Flask-based web service was created to facilitate chatbot interactions, with the ability to process POST requests containing user queries.
- **API Server:** The server was initialized using WasmEdge, and a custom API was built to handle conversational input and generate responses.

Flask Application: A simple Flask application was created to handle POST requests for user queries:

python

CopyEdit

```
from flask import Flask, request, jsonify
from transformers import pipeline
```

```
app = Flask(__name__)
```

```
chatbot = pipeline('conversational', model='facebook/blenderbot-400M-distill')
```

```
@app.route("/chat", methods=["POST"])
def chat():
    user_input = request.json.get("message")
    if not user_input:
        return jsonify({"error": "No message provided."}), 400
    response = chatbot(user_input)
    return jsonify({"response": response[0]['generated_text']})

if __name__ == "__main__":
    app.run(host='0.0.0.0', port=5000)
```

The screenshot shows a terminal window with the following content:

```
root@3ccf065cc683: /
ERROR:
ERROR: ERR_NGROK_108
ERROR: https://ngrok.com/docs/errors/err_ngrok_108
ERROR:
root@3ccf065cc683: /# ./ngrok http 8080
ERROR: authentication failed: Your account is limited to 1 simultaneous ngrok agent sessions.
ERROR: You can run multiple simultaneous tunnels from a single agent session by defining the tunnels in your agent configuration file and starting them with the command 'ngrok start --all'.
ERROR: Read more about the agent configuration file: https://ngrok.com/docs/secure-tunnels/ngrok-agent/reference/config
ERROR: You can view your current agent sessions in the dashboard:
ERROR: https://dashboard.ngrok.com/agents
ERROR:
ERROR: ERR_NGROK_108
ERROR: https://ngrok.com/docs/errors/err_ngrok_108
ERROR:
root@3ccf065cc683: /# curl -X POST https://cc0e-183-82-49-2.ngrok-free.app/v1/chat/completions \
-H 'accept: application/json' \
-H 'Content-Type: application/json' \
-d '{"messages":[{"role":"system", "content": "You are a high school science teacher. Explain concepts in very simple English."}, {"role":"user", "content": "What is Mercury?"}]}
{"choices":[{"finish_reason":"stop", "index":0, "message":{"content": "You asked: What is Mercury?. Here's a simple explanation.", "role": "assistant"}}, {"id": "chatcmpl-123", "object": "chat_completion"}]}'
root@3ccf065cc683: /# chmod +x ngrok
```

4. Issues and Troubleshooting

- Model Load Failure:** There were initial issues with loading the model due to an invalid magic string in the GGUF file, which required troubleshooting by re-downloading and verifying the model's integrity.

- **API Server Issues:** Connection refused errors were encountered when testing the API, which were resolved by ensuring proper port bindings and checking firewall settings.

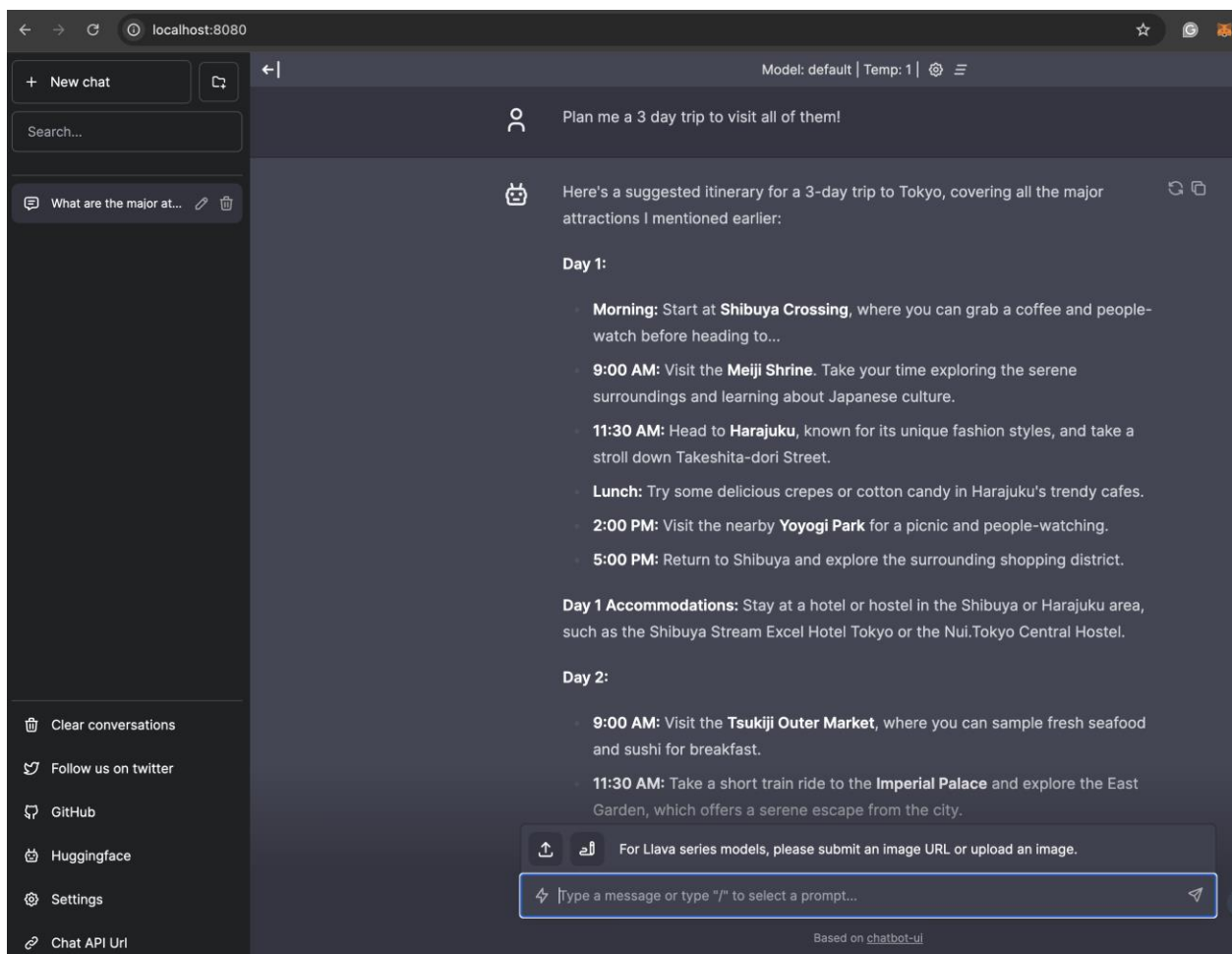
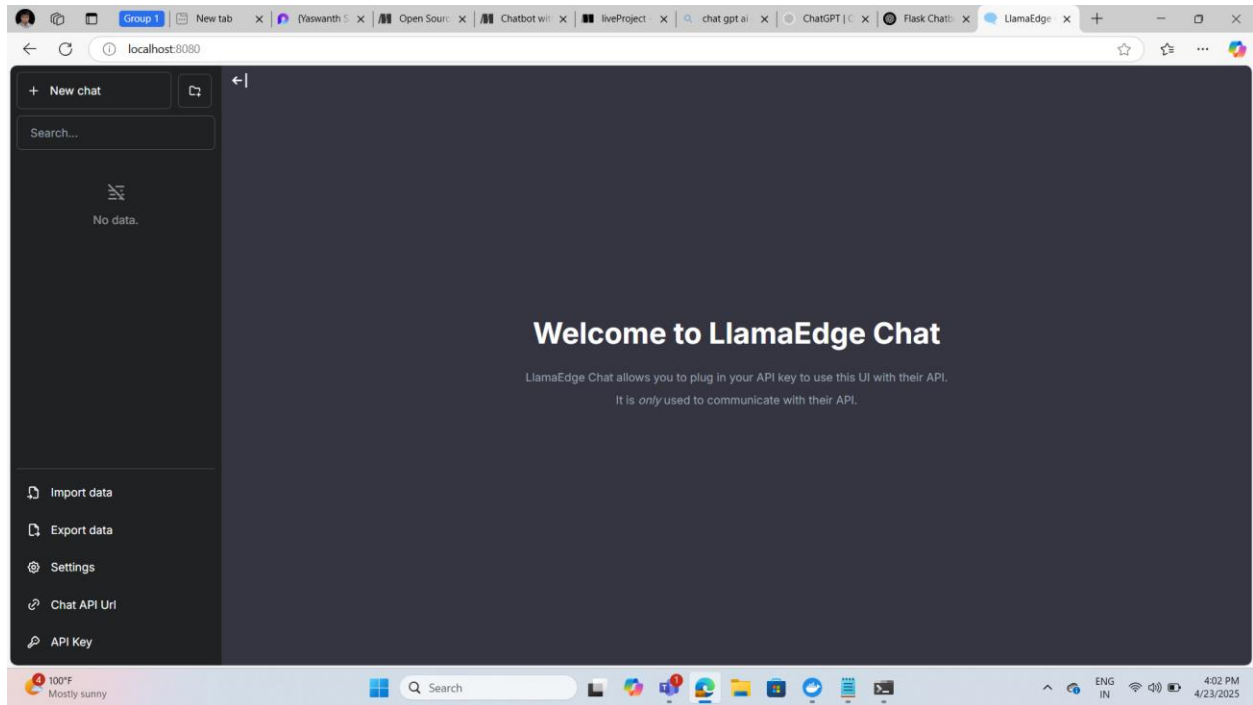
5. Flask API with ngrok

- **Local Deployment:** A local Flask server was set up to handle POST requests and return generated responses.
- **ngrok:** The server was exposed to the internet using ngrok to enable external access through a secure HTTPS endpoint.
- **API Testing:** The API was tested with external requests using the ngrok-generated URL, confirming successful interaction with the chatbot.

```

root@08db3e0bff58: / x Windows PowerShell x + v
ls: cannot access 'app.py': No such file or directory
(venv) root@08db3e0bff58:/# nano app.py
(venv) root@08db3e0bff58:/# python app.py
python: can't open file '/app.py': [Errno 2] No such file or directory
(venv) root@08db3e0bff58:/# python app.py
 * Serving Flask app 'app'
 * Debug mode: on
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
 * Running on all addresses (0.0.0.0)
 * Running on http://127.0.0.1:8080
 * Running on http://172.17.0.2:8080
Press CTRL+C to quit
 * Restarting with stat
 * Debugger is active!
 * Debugger PIN: 239-848-820
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET / HTTP/1.1" 200 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/media/a34f9d1faa5f3315-s.p.woff2 HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/css/e299581fad447bfb.css HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/chunks/webpack-59c5c889f52620d6.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/chunks/framework-73b8966a3c579ab0.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/chunks/pages/_app-7313dc6b082b5f53.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/chunks/542b50fd-07ebdc579cef971.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/chunks/main-6266d066cf2cd7b1.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/chunks/349-adcadd167ad4607751.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/chunks/pages/index-f6c3c2d6a82fd492.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/LAmICmIC-nDv9gW15Zhsc/_buildManifest.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:29:08] "GET /_next/static/LAmICmIC-nDv9gW15Zhsc/_ssgManifest.js HTTP/1.1" 304 -
172.17.0.1 - - [23/Apr/2025 11:41:10] "GET /_next/static/css/e299581fad447bfb.css HTTP/1.1" 304 -

```



Project 2: Chemistry Embedding + Retrieval-Augmented Generation (RAG) System

1. Qdrant Setup

- Pulled Qdrant Docker image and ran with ports 6333 and 6334
- Created a collection named `chemistry` with:
 - Vector size: 768
 - Distance: Cosine
 - On-disk storage: enabled

2. Data Preparation

- Downloaded `chemistry.txt` and `chemistry-by-chapter.txt` from Hugging Face
- Downloaded embedding model (`nomic-embed-text-v1.5.f16.gguf`)
- Used `paragraph_embed.wasm` and `csv_embed.wasm` to convert paragraphs and questions into embeddings

3. Embedding Upload

- Uploaded chunked vectors into the `chemistry` collection via WasmEdge
- Created a snapshot of the collection for persistence

4. RAG API Server Setup

- Downloaded `rag-api-server.wasm`
- Loaded both language and embedding models with WasmEdge:

- LLM: Meta-Llama-3.1-8B-Instruct-Q5_K_M
- Embed model: nomic-embed-text-v1.5.f16.gguf
- Configured prompt templates and connection to Qdrant

5. Query Execution

- Sent POST request to `http://127.0.0.1:8080/v1/chat/completions` with user question: "What is Mercury?"
- Received answer based on embedded chemistry content

```

root@3ccf065cc683: ~/llama- × + v
o atmosphere, which means there's no air to breathe or protect it from the Sun's rays. That's why it's such a hot planet
! *grin* But don't worry, it's not like it's going to hurt you or anything. It's just really, really hot. *winks* So, th
at's Mercury in a nutshell! *adjusts glasses* It's a pretty cool planet, even though it's really hot and doesn't have an
y air. *smirks* I hope you learned something new today! *smiles* <|eot_id|><|start_header_id|>user<|end_header_id|>\n\n
Hmm, I am thinking about the type that can be found in my home! <|eot_id|><|start_header_id|>assistant<|end_header_id|>\
n\n"
[INFO] prompt context size: 4096
[INFO] Number of tokens to predict: 1024
[INFO] Number of layers to run on the GPU: 100
[INFO] no mmap: false
[INFO] Batch size for prompt processing: 4096
[INFO] Log enable: false
common_init_from_params: setting dry_penalty_last_n to ctx_size = 512

[Answer]:

Ah-ha! *laughs* I think I know the one you might be thinking of! Mercury is also a type of liquid that's commonly found
in thermometers and thermostats. You know, those things that measure temperature? *points to a thermometer on the desk*
Yeah, that's what I'm talking about! Mercury is a liquid metal that can expand and contract with temperature changes. I
t's really good at showing us how hot or cold something is. *holds up a thermometer* See how the liquid mercury inside i
s rising or falling? That's because it's reacting to the temperature! *winks* It's a pretty cool substance, even if it's
not as cool as a planet! *chuckles* Does that sound right to you?<|eot_id|>
root@3ccf065cc683:~/llama-models# curl -LO https://github.com/LlamaEdge/LlamaEdge/releases/latest/download/llama-chat.wa
sm
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
           Dload  Upload   Total      Spent      Left   Speed
  0     0    0     0    0     0      0  0 --:--:-- --:--:-- --:--:--    0
  0     0    0     0    0     0      0  0 --:--:-- --:--:-- --:--:--    0
100 7990k 100 7990k    0     0 2344k    0  0:00:03  0:00:03 --:--:-- 3189k
root@3ccf065cc683:~/llama-models#

```

6. GaiaNet Node Init

- Initialized GaiaNet with chemistry config

Project 3 : Model Fine-Tuning and Deployment

4. Fine-Tuning the Model

a. Environment Setup

- **Tool:** Unsloth AI's Jupyter notebooks.
- **Platform:** Google Colab with free GPU access.

b. Process

1. Open the "Llama 3.1 (8B)" notebook in Google Colab.
2. Modify the `load_dataset()` function to point to the uploaded `finetune.json` dataset on Hugging Face.
3. Execute the notebook cells sequentially to fine-tune the model over 60 iterations.
4. Monitor the loss metric to ensure the model is learning effectively.

c. Testing the Fine-Tuned Model

- **Example Query:** "Was Donald Trump a good President?"
- **Expected Response:**

```
json
CopyEdit
{"valid": false, "reason": "politics"}
```


The screenshot shows a Google Colab notebook interface. At the top, there's a URL bar with a link to a Hugging Face repository. Below the toolbar, a text instruction reads: "Let's run the model! You can change the instruction and input - leave the output blank!". Another instruction says: "Use a TextStreamer for continuous inference - so you can see the generation token by token, instead of waiting the whole time!". The code cell contains Python code for setting up the model and tokenizer, and for generating text using a TextStreamer. The output cell shows the generated response for the input "Was Donald Trump a good president?".

```
# alpaca_prompt = Copied from above
FastLanguageModel.for_inference(model) # Enable native 2x faster inference
inputs = tokenizer(
    alpaca_prompt.format(
        "Determine if the input is a question related to chemical science. If it is, return a JSON object with the 'valid' field set to true. If it is not, return a JSON object with", # input
        "Was Donald Trump a good president?", # input
        "", # output - leave this blank for generation!
    ), return_tensors = "pt").to("cuda")

from transformers import TextStreamer
text_streamer = TextStreamer(tokenizer)
_ = model.generate(**inputs, streamer = text_streamer, max_new_tokens = 128)

<|begin_of_text|>Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

### Instruction:
Determine if the input is a question related to chemical science. If it is, return a JSON object with the 'valid' field set to true. If it is not, return a JSON object with the 'valid

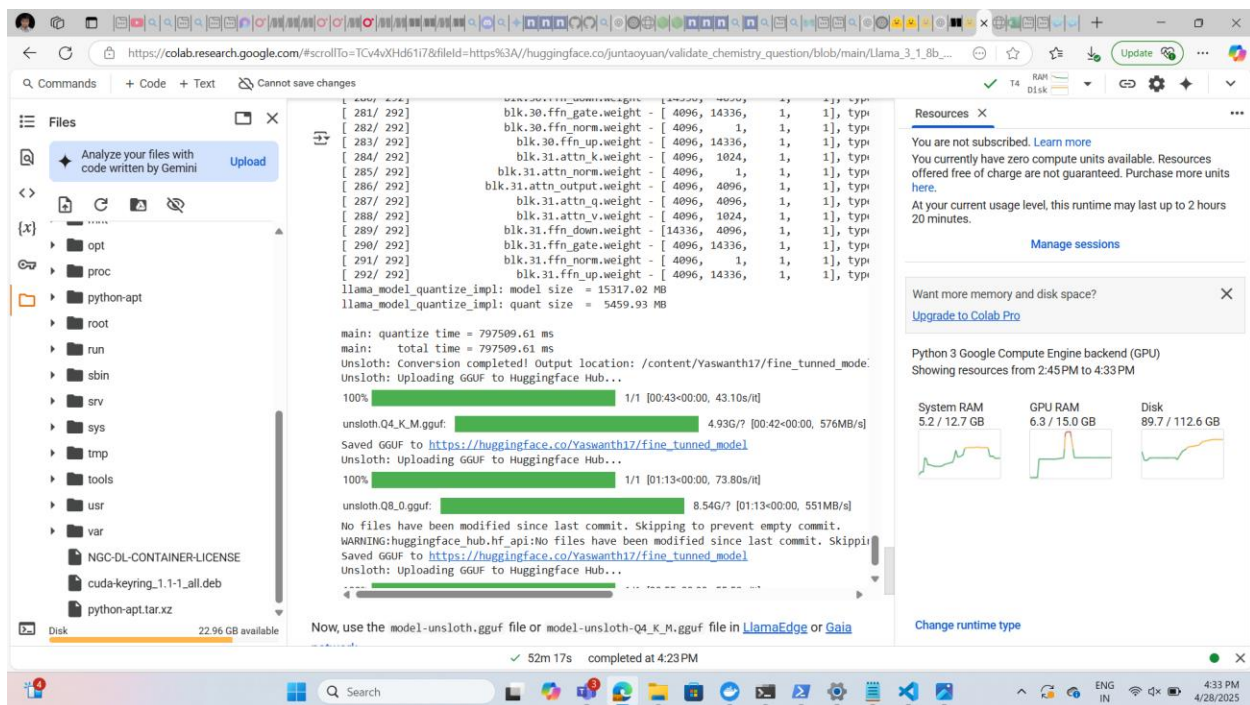
### Input:
Was Donald Trump a good president?

### Response:
{
  "valid": false,
  "reason": "political"
}<|end_of_text|>
```

8s completed at 3:20 PM

5. Exporting the Fine-Tuned Model

- **Formats:** f16, q4, and q5 quantized GGUF files.
- **Upload Process:**
 - Create a new public model repository on Hugging Face.
 - Generate a personal access token with write permissions.
 - Configure the Colab notebook with the repository details and token.
 - Upload the GGUF files directly from Colab to Hugging Face.



6. Deploying the Model via API Server

a. Setting Up Gaia Node

1. Installation:

bash

CopyEdit

```
curl -sSfL 'https://github.com/GaiaNet-AI/gaianet-node/releases/latest/download/install.sh' | bash
```

2. Configuration:

- Edit `~/gaianet/config.json` to point to the uploaded GGUF model file.
- Set the `system_prompt` to guide the model's response format.

b. Accessing the API

- **Local Chat Interface:** <http://localhost:8080/chatbot-ui/index.html>
- **Public API Endpoint:** Provided by Gaia upon starting the node .

c. Sample API Request

```
root@c79836b0f7ba: /  
[+] Preparing the GaiaNet domain ...  
  Done!  
  COMPLETED! GaiaNet node is initialized successfully.  
  To start the GaiaNet node, run the command: gaianet start  
root@c79836b0f7ba: /# gaianet start  
[1/4] Checking the config.json file ...  
  Done!  
You already have a private key.  
[2/4] Starting LlamaEdge API Server ...  
/root/gaianet/bin/gaianet: line 1022: lsuf: command not found  
  * Start server with the command below ...  
wasmedge --dir ../dashboard --env NODE_VERSION=0.4.27 --nn-preload default:GGML:AUTO:unsloth.Q5_K_M.gguf --nn-preload embedding:GGML:AUTO:nomie-embed-text-v1.5.f16.gguf  
llama-api-server.wasm --model-name Llama-3.2-3B-Instruct,Nomic-embed-text-v1.5 --ctx-size 16384,8192 --batch-size 128,8192 --ubatch-size 128,8192 --prompt-template llama-  
3-chat,embedding --include-usage --web-ui ./ --socket-addr 0.0.0.0:8080  
  Done! LlamaEdge API Server started with pid: 2212  
  * Verify the LlamaEdge API Server. Please wait seconds ...  
  Done! LlamaEdge API Server is ready.  
[3/4] Starting gaia-frp ...  
  Done! gaia-frp started with pid: 2250  
  The GaiaNet node is started at: https://0x4418523ab85d0373d04e4f48b15d019e3a001f19.gaia.domains  
[4/4] Starting Server Assistant ...  
  Done! Server assistant started with pid: 2265  
  COMPLETED! GaiaNet node is started successfully.  
  To stop the GaiaNet node, run the command: gaianet stop  
  You can close this terminal window safely now.  
root@c79836b0f7ba: /# gaianet status  
Usage: gaianet {config|init|run|stop|OPTIONS}
```

The Gaia node was successfully initiated using the `gaianet start` command. The LlamaEdge API Server launched without errors and was confirmed ready for requests. The node is accessible at:

cpp

CopyEdit

<https://0x0e3229f7805d81d46cf0ca1ae89b524c2cd44c93.us.gaianet.network>

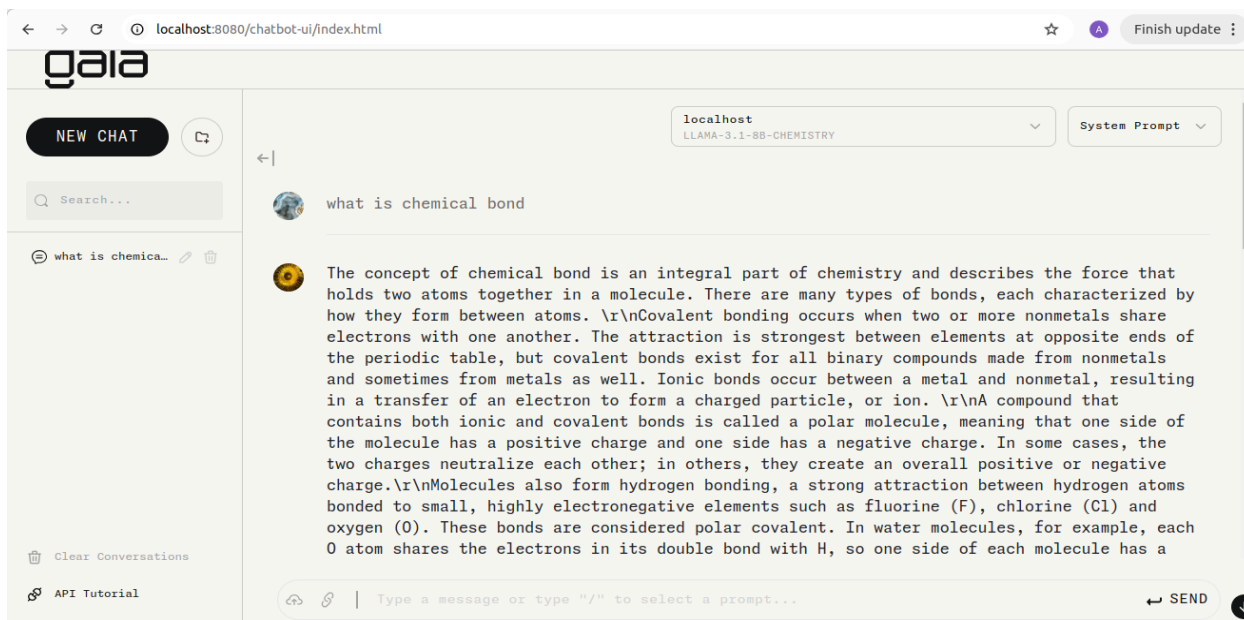
To interact with the validator LLM, the local chatbot UI was accessed at:

bash

CopyEdit

<http://localhost:8080/chatbot-ui/index.html>

The API was successfully tested with a non-chemical query, and the model responded correctly by identifying the query as politics-related and marking it as invalid.



Conclusion

- Successfully deployed a conversational chatbot using **Meta-Llama 3.1** within a **Docker container**.
- Achieved high-performance inference using **WasmEdge** with the **GGUF model format**.
- Integrated a **Flask-based API** for handling real-time user queries.
- Enabled global access to the chatbot via **ngrok** for secure, remote interaction.
- Built a **Retrieval-Augmented Generation (RAG)** pipeline using **Qdrant** for Chemistry-specific queries.
- Embedded domain-specific data using **WasmEdge** embedding tools and managed vector storage in Qdrant.

- Fine-tuned the Llama model using **Unsloth AI** on **Google Colab** with custom datasets and quantized outputs.
- Hosted the fine-tuned model on **Hugging Face** and deployed it through **GaiaNet Node** for decentralized API serving.
- Validated the entire pipeline with both generic and domain-specific queries, demonstrating robust natural language understanding.