

# Project Report – Amit Meena

## Project 1: Custom Chatbot Interface for LLaMA and Mistral Models

### Objective:

- Build a **custom chatbot UI** where users interact with **LLaMA 2**, **LLaMA 3.2**, and **Mistral** models.

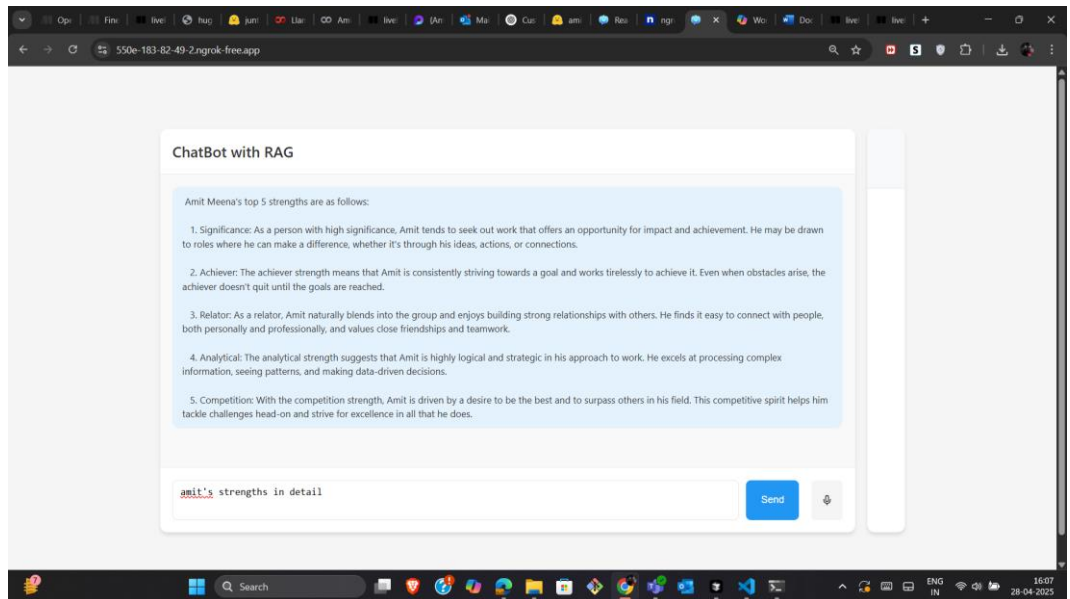
### Implementation:

1. **Model Setup:**
  - a. Downloaded and hosted **LLaMA 2**, **LLaMA 3.2**, and **Mistral** models locally.
  - b. Optimized loading and memory management for efficient inference.
2. **Backend (Python):**
  - a. Developed a Flask-based backend API to handle user prompts and generate model responses.
  - b. Integrated an optional retrieval step (RAG) **before** querying the model (detailed in Project 2).
3. **Frontend (React):**
  - a. Built an intuitive chatbot UI in React.
  - b. Features include:
    - i. Real-time conversation flow.
    - ii. Model selection (user can choose between LLaMA 2, LLaMA 3.2, and Mistral).
    - iii. Support for additional retrieved context (RAG responses injected dynamically).
    - iv. API communication through the **Ngrok URL**.
4. **Testing:**
  - a. Validated both direct model responses and RAG-enhanced responses.
  - b. Optimized for low latency and smooth user experience.

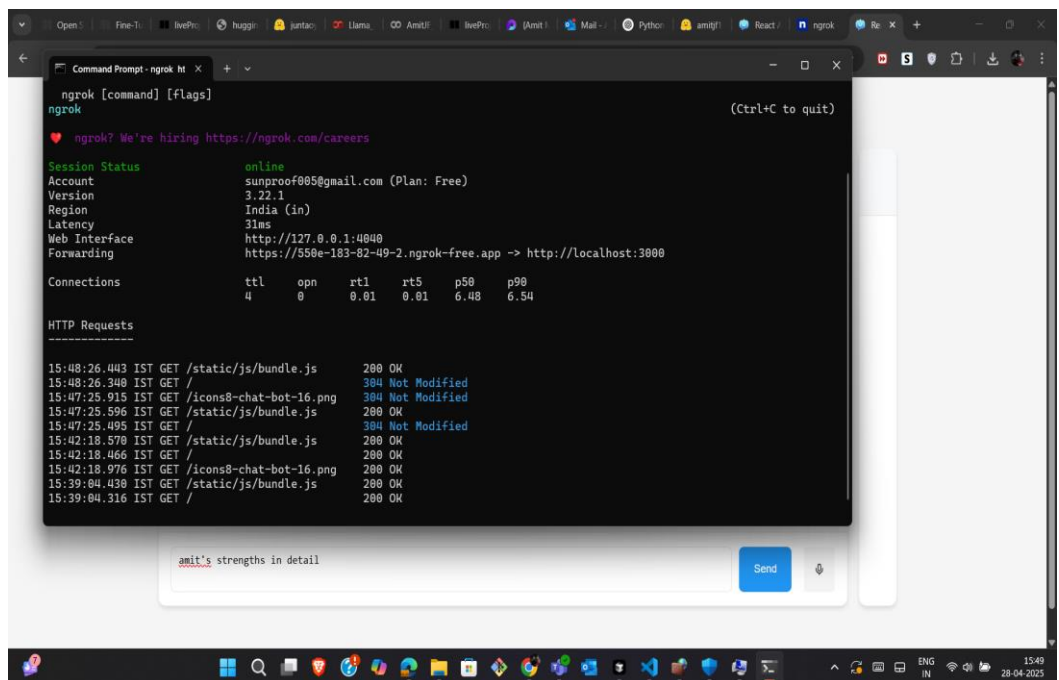
**Note:** This chatbot is tightly coupled with a **Retrieval-Augmented Generation (RAG)** system for better, context-aware answers (see Project 2).

## Screenshots:

### a. React Chatbot Interface:



### b. Flask API Exposed via Ngrok:



# Project 2: Retrieval-Augmented Generation (RAG)

## System Integration

### Objective:

- Implement **Retrieval-Augmented Generation (RAG)** to improve the accuracy and relevance of model outputs.

### Implementation:

#### 1. Vector Database Creation:

- a. Processed large sets of domain-specific documents.
- b. Converted documents into vector embeddings using transformer-based embedding models.
- c. Indexed embeddings into a **FAISS-based** vector database.

#### 2. Retrieval Pipeline:

- a. On receiving a user query:
  - i. Search the vector database for the most relevant documents.
  - ii. Retrieve top-k matches as context snippets.
  - iii. Append these snippets to the user prompt.

#### 3. Response Generation:

- a. The model (LLaMA or Mistral) generates its answer **based on both** the user's query **and** the retrieved context.
- b. This ensures more **informed and accurate** responses, especially for domain-specific or fact-based questions.

#### 4. Integration with Chatbot:

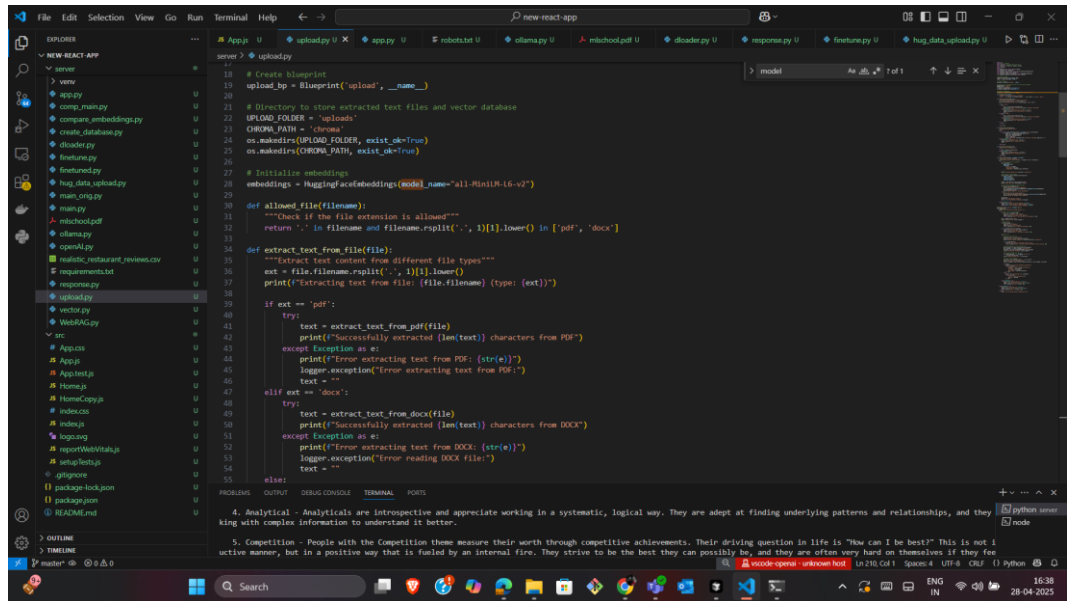
- a. Seamlessly connected the RAG pipeline with the chatbot backend.
- b. User can interact without noticing the retrieval step, but benefits from smarter outputs.

#### 5. Testing:

- a. Compared RAG-enhanced outputs vs standard model outputs.
- b. Found significant improvement in answer relevance and factual correctness.

#### 6. Screenshots:

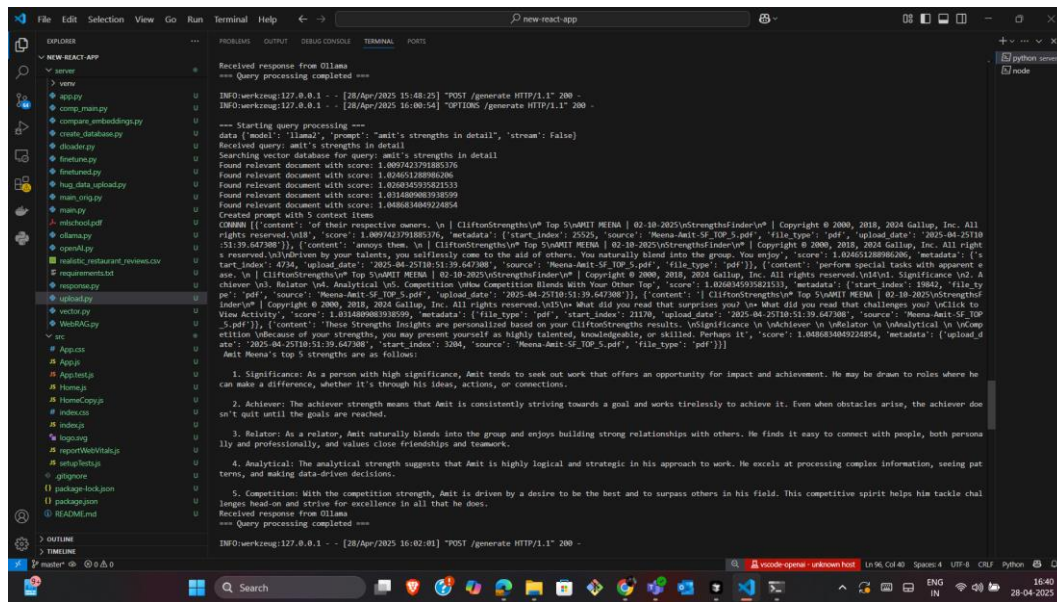
a. *Vector Database Setup (CHROMA DB):*



```
server > upload.py
18 # Create blueprint
19 upload_bp = Blueprint('upload', __name__)
20
21 # Directory to store extracted text files and vector database
22 UPLOAD_FOLDER = 'uploads'
23 CHROMA_PATH = 'chroma'
24 os.makedirs(UPLOAD_FOLDER, exist_ok=True)
25 os.makedirs(CHROMA_PATH, exist_ok=True)
26
27 # Initialize embeddings
28 embeddings = HuggingFaceEmbeddings(model_name="all-MiniLM-L6-v2")
29
30 def allowed_file(filename):
31     """Check if the file extension is allowed"""
32     return '.' in filename and filename.rsplit('.', 1)[1].lower() in ['pdf', 'docx']
33
34 def extract_text_from_file(file):
35     """Extract text content from different file types"""
36     ext = file.filename.rsplit('.', 1)[1].lower()
37     print(f"Extracting text from file: {file.filename} (type: {ext})")
38     text = ""
39     if ext == 'pdf':
40         try:
41             text = extract_text_from_pdf(file)
42             print(f"Successfully extracted {len(text)} characters from PDF")
43         except Exception as e:
44             print(f"Error extracting text from PDF: {str(e)}")
45             logger.exception("Error extracting text from PDF:")
46             text = ""
47     elif ext == 'docx':
48         try:
49             text = extract_text_from_docx(file)
50             print(f"Successfully extracted {len(text)} characters from DOCX")
51         except Exception as e:
52             print(f"Error extracting text from DOCX: {str(e)}")
53             logger.exception("Error reading DOCX file:")
54             text = ""
55     else:
56         print(f"Unsupported file type: {ext}")
57         text = ""
58     return text
59
60 def extract_text_from_pdf(file):
61     """Extract text from a PDF file"""
62     text = PyPDF2.PdfReader(file).extract_text()
63     return text
64
65 def extract_text_from_docx(file):
66     """Extract text from a DOCX file"""
67     doc = docx.Document(file)
68     return doc.text
69
70 @upload_bp.route('/upload', methods=['POST'])
71 def upload():
72     """Handle file upload and store in vector database"""
73     if 'file' not in request.files:
74         return jsonify({'error': 'No file part'}), 400
75     file = request.files['file']
76     if file.filename == '':
77         return jsonify({'error': 'No selected file'}), 400
78     if not allowed_file(file.filename):
79         return jsonify({'error': 'Invalid file type'}), 400
80     text = extract_text_from_file(file)
81     if text:
82         save_to_vector_db(text, file.filename, file.filename.rsplit('.', 1)[1].lower())
83     return jsonify({'message': 'File uploaded and processed successfully'}), 200
84
85 def save_to_vector_db(text, filename, file_type):
86     """Save text to vector database"""
87     print(f"Saving text to vector database: {filename} (type: {file_type})")
88     try:
89         # Create a document with metadata
90         doc = Document(
91             page_content=text,
92             metadata={
93                 "source": filename,
94                 "file_type": file_type,
95                 "upload_date": datetime.now().isoformat()
96             }
97         )
98         print(f"Created document with {len(text)} characters")
99         # Split the text into chunks
100         text_splitter = RecursiveCharacterTextSplitter(
101             chunk_size=100,
102             chunk_overlap=50,
103             length_function=len,
104             add_start_index=True,
105         )
106         chunks = text_splitter.split_documents([doc])
107         print(f"Split text into {len(chunks)} chunks")
108         # Load existing database or create new one
109         if os.path.exists(CHROMA_PATH):
110             print("Loading existing Chroma database")
111             db = Chroma(persist_directory=CHROMA_PATH, embedding_function=embeddings)
112             db.add_documents(chunks)
113             print("Added documents to existing database")
114         else:
115             print("Creating new Chroma database")
116             db = Chroma.from_documents(chunks, embeddings, persist_directory=CHROMA_PATH)
117     except Exception as e:
118         print(f"Error saving to vector database: {str(e)}")
119         logger.exception("Error saving to vector database:")
120         return jsonify({'error': 'Failed to save document'}), 500
121
122 if __name__ == '__main__':
123     server.run(upload_bp)
```

b. *Retrieved Context Added to Prompt:*







# Project 3: Dataset Creation and Fine-tuning of LLaMA Model

## Objective:

- Create a custom dataset.
- Upload it to Hugging Face.
- Fine-tune the LLaMA model to adapt it to specific tasks/domains.

## Implementation:

### 1. Dataset Creation:

- a. Designed a structured dataset (e.g., instruction-response pairs, domain-specific Q&A).
- b. Ensured high data quality with clear formatting (JSONL/CSV).

### 2. Uploading to Hugging Face:

- a. Uploaded the dataset to my Hugging Face account for accessibility and versioning.

### 3. Fine-tuning:

- a. Fine-tuned the LLaMA model using the custom dataset.
- b. Techniques used:
  - i. LoRA (Low-Rank Adaptation) for efficient fine-tuning.
  - ii. Hyperparameter tuning (learning rate, batch size, epochs) for optimal results.

### 4. Results:

- a. Achieved a more specialized LLaMA model that performed better on custom task domains.
- b. Improved coherence, task adherence, and response specificity.

**Screenshot - Custom Dataset Example(Hugging Face upload):**



The screenshot shows the Hugging Face dataset card for 'amitj111/first-finetuning-validate-chemistry-questions'. The dataset is licensed under gpl-3.0 and is currently in the 'Dataset card' view. The 'Dataset Viewer' section shows a split of 1 train (933 rows). The 'input' column is a string of length 39, and the 'output' column is a string of length 213. The 'output' column contains JSON objects with a 'valid' key. The 'Downloads last month' section shows 0 downloads. The 'Size of downloaded dataset files' is 418 kB, and the 'Size of the auto-converted Parquet files' is 51.2 kB. The 'Number of rows' is 933.

instruction	input	output
Determine if the input is a question related to chemical science. If it is, return a JSON object with...	What is the definition of chemistry?	{'valid': true}
Determine if the input is a question related to chemical science. If it is, return a JSON object with...	What are physical properties of matter?	{'valid': true}
Determine if the input is a question related to chemical science. If it is, return a JSON object with...	How do chemical changes differ from physical changes?	{'valid': true}
Determine if the input is a question related to chemical science. If it is, return a JSON object with...	What is a substance in terms of chemistry?	{'valid': true}
Determine if the input is a question related to chemical science. If it is, return a JSON object with...	Can elements be broken down into simpler substances?	{'valid': true}
Determine if the input is a question related to chemical science. If it is, return a JSON object with...	What is the definition of science?	{'valid': true}

Screenshot – use while finetuning:

The screenshot shows a Google Colab notebook titled 'AmitJF\_Llama\_3\_1\_8b + Unsloth finetuning.ipynb'. The notebook contains a Python script that defines a 'formatting\_prompts\_func' and uses the 'datasets' library to load and process the dataset. The script also includes a 'train' function that uses the 'trl' library to train the model. The notebook shows the progress of the training process, including the download of the dataset and the training of the model. The training process is completed at 16:31.

```

eos_token = tokenizer.eos_token + MUST_ADD_EOS_TOKEN

def formatting_prompts_func(examples):
    instructions = examples["instruction"]
    inputs       = examples["input"]
    outputs      = examples["output"]
    texts = []
    for instruction, input, output in zip(instructions, inputs, outputs):
        # Must add EOS_TOKEN, otherwise your generation will go on forever!
        text = alpaca_prompt.format(instruction, input, output) + EOS_TOKEN
        texts.append(text)
    return { "text" : texts, }

pass

from datasets import load_dataset
dataset = load_dataset("amitj111/first-finetuning-validate-chemistry-questions", split = "train")
dataset = dataset.map(formatting_prompts_func, batched = True,)

```

Progress bars show the following status:

- README.md: 100% (925/925) [00:00-00:00, 21.9KB/s]
- finetune.json: 100% (418k/418k) [00:00-00:00, 5.65MB/s]
- Generating train split: 100% (933/933) [00:00-00:00, 7863.48 examples/s]
- Map: 100% (933/933) [00:00-00:00, 14011.10 examples/s]

**Train the model**

Now let's use Huggingface TRL's SFTTrainer! More docs here: [TRL SFT docs](#). We do 60 steps to speed things up, but you can set num\_train\_epochs=1 for a full run, and turn off max\_steps=None. We also support TRL's DPOTrainer!

## Screenshot - Fine-tuning Logs:

```
# Save to multiple GGUF options - much faster if you want multiple!
if True:
    model.push_to_hub_gguf(
        "hf/model", # Change hf to your username!
        tokenizer,
        quantization_method = ["q4_k_m", "q8_0", "q5_k_m"],
        token = "", # Get a token at https://huggingface.co/settings/tokens
    )

INFO:hf-to-gguf:blk.15.attn.v.weight, torch.float16 --> F16, shape = (4096, 1024)
INFO:hf-to-gguf:blk.15.attn.output.weight, torch.float16 --> F16, shape = (4096, 4096)
INFO:hf-to-gguf:blk.15.ffn_gate.weight, torch.float16 --> F16, shape = (4096, 14336)
INFO:hf-to-gguf:blk.15.ffn_up.weight, torch.float16 --> F16, shape = (4096, 14336)
INFO:hf-to-gguf:blk.15.ffn_down.weight, torch.float16 --> F16, shape = (14336, 4096)
INFO:hf-to-gguf:blk.15.attn_norm.weight, torch.float16 --> F32, shape = (4096)
INFO:hf-to-gguf:blk.15.ffn_norm.weight, torch.float16 --> F32, shape = (4096)
INFO:hf-to-gguf:blk.16.attn.q.weight, torch.float16 --> F16, shape = (4096, 4096)
INFO:hf-to-gguf:blk.16.attn.k.weight, torch.float16 --> F16, shape = (4096, 1024)
INFO:hf-to-gguf:blk.16.attn.v.weight, torch.float16 --> F16, shape = (4096, 1024)
INFO:hf-to-gguf:blk.16.attn.output.weight, torch.float16 --> F16, shape = (4096, 4096)
INFO:hf-to-gguf:blk.16.ffn_gate.weight, torch.float16 --> F16, shape = (4096, 14336)
INFO:hf-to-gguf:blk.16.ffn_up.weight, torch.float16 --> F16, shape = (4096, 14336)
INFO:hf-to-gguf:blk.16.ffn_down.weight, torch.float16 --> F16, shape = (14336, 4096)
INFO:hf-to-gguf:blk.16.attn_norm.weight, torch.float16 --> F32, shape = (4096)
INFO:hf-to-gguf:blk.16.ffn_norm.weight, torch.float16 --> F32, shape = (4096)
INFO:hf-to-gguf:blk.17.attn.q.weight, torch.float16 --> F16, shape = (4096, 4096)
INFO:hf-to-gguf:blk.17.attn.k.weight, torch.float16 --> F16, shape = (4096, 1024)
INFO:hf-to-gguf:blk.17.attn.v.weight, torch.float16 --> F16, shape = (4096, 1024)
INFO:hf-to-gguf:blk.17.attn.output.weight, torch.float16 --> F16, shape = (4096, 4096)
INFO:hf-to-gguf:blk.17.ffn_gate.weight, torch.float16 --> F16, shape = (4096, 14336)
INFO:hf-to-gguf:blk.17.ffn_up.weight, torch.float16 --> F16, shape = (4096, 14336)
INFO:hf-to-gguf:blk.17.ffn_down.weight, torch.float16 --> F16, shape = (14336, 4096)
INFO:hf-to-gguf:blk.17.attn_norm.weight, torch.float16 --> F32, shape = (4096)
INFO:hf-to-gguf:blk.17.ffn_norm.weight, torch.float16 --> F32, shape = (4096)
```

## Conclusion

This project series involved end-to-end development: setting up strong foundation models, creating an interactive chatbot, augmenting it with retrieval for smarter responses, and fine-tuning a model for even more targeted performance.

The experience covered key real-world AI production workflows like serving large models, building retrieval pipelines, integrating frontends, and model personalization.