# GenAl-Powered Data Pipeline Project Report

Date: March 31, 2025

Student: Vignesh Jetty Ravi(Assumed from screenshots)

Student ID: 018286970

### 1. Project Overview

This project focuses on building a **cloud-based GenAI-powered data pipeline** using **Langflow** and **AstraDB**.

# **Key Components:**

- **Data Ingestion**: Unstructured FAQ dataset is processed and ingested.
- **Vectorization**: The dataset is transformed using **OpenAI embeddings**.
- **Retrieval-Augmented Generation (RAG)**: The system enables intelligent query responses by leveraging vectorized data.

This architecture enhances the ability to provide accurate, context-aware answers based on the ingested FAQs.

### 2. Architecture Overview

The system is structured into two primary workflows in **Langflow**:

#### • Load Data Flow:

- Ingests CSV data
- o Splits text into manageable chunks
- Generates embeddings using **OpenAI**
- Stores the vectorized data in **AstraDB**

#### • Retriever Flow:

• Processes user queries

- Embeds the input question
- Retrieves relevant text chunks from **AstraDB**
- Constructs a prompt and generates a response using **GPT-3.5**

This streamlined approach ensures efficient data handling and accurate retrieval-based answering.

# 3. Tools & Technologies

- Langflow (Cloud Version) Orchestrates the data pipeline
- AstraDB Stores vectorized data with vector search enabled
- **OpenAI Embeddings** Utilizes **text-embedding-3-small** for vectorization
- **OpenAI GPT-3.5 Turbo** Powers the Retrieval-Augmented Generation (RAG) process
- Dataset faq\_bulk\_150\_records.csv, an unstructured FAQ dataset

### 4. Key Design Decisions

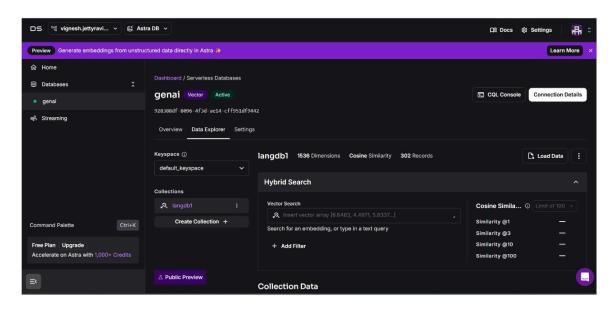
- **Text Chunking**: Set to **100 tokens** with a **20-token overlap** to enhance context retrieval.
- **Embedding Model**: Used **text-embedding-3-small** for faster and efficient vectorization.
- **Retrieval Optimization**: Limited results to **5 chunks** to minimize LLM context size and improve response speed.
- LLM Selection: Chose GPT-3.5 and 40 mini for its balance of speed and response quality.
- Dynamic Prompting: Implemented prompt templating with variables for optimized RAG-based queries.

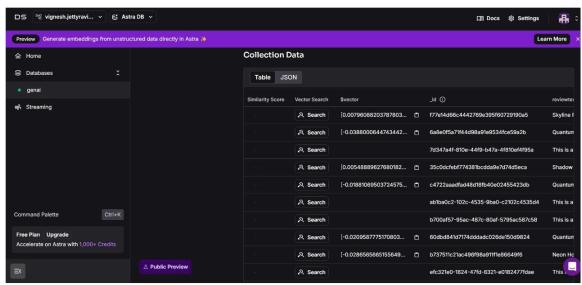
# 5. Challenges & Solutions

- -ERROR building Component Astra DB: Error adding documents to AstraDBVectorStore: Length of vector parameter different from declared '\$vector' dimension: root cause = (InvalidQueryException) Failed to insert document with \_id 979313aa79e74367a333c98102224431: Unexpected 2048 extraneous bytes after vector<float, 1024> . I had to adjust the vector dimension from 1024 to 1536.
- -ERROR during langflow run: solved via careful step-by-step testing and visual flow debugging and refreshing the page for re-run.

#### 6. Results & Screenshots

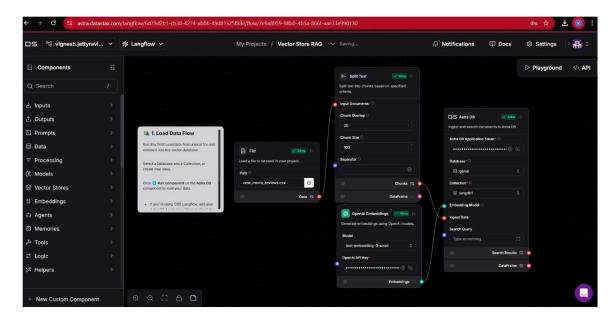
#### Astra DB



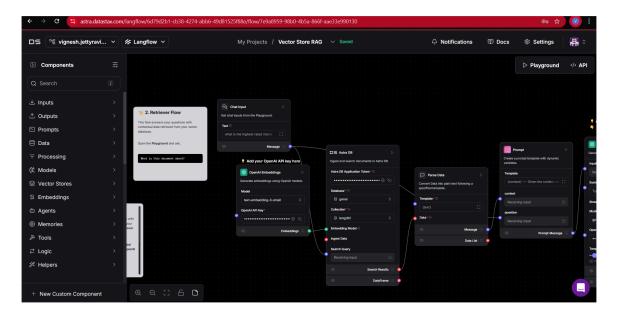


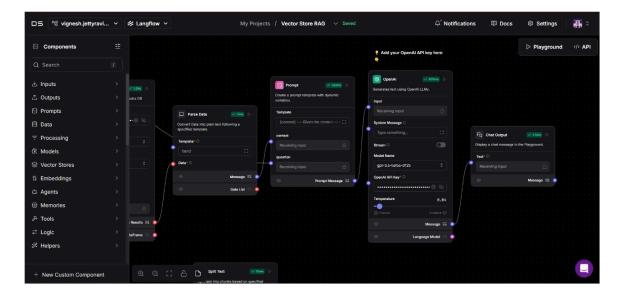
# Langflow

# 1.load data flow

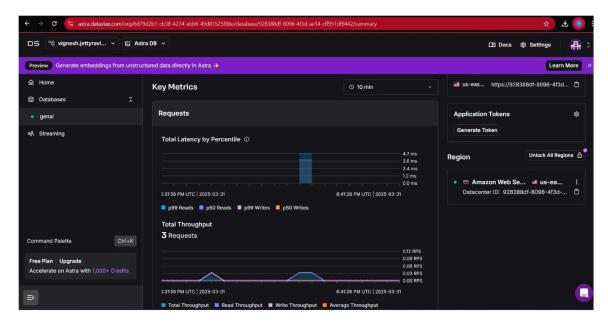


# 2. retriever flow

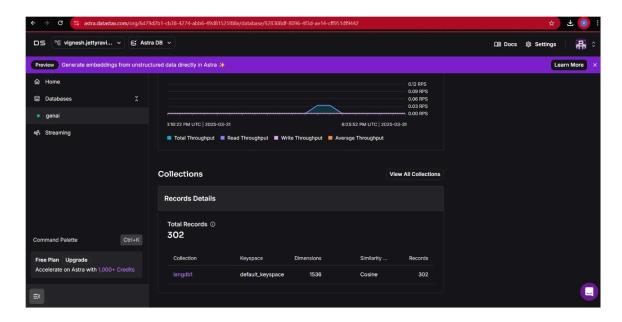




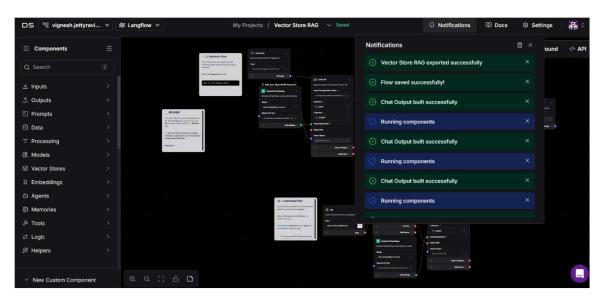
### 3.key metrics



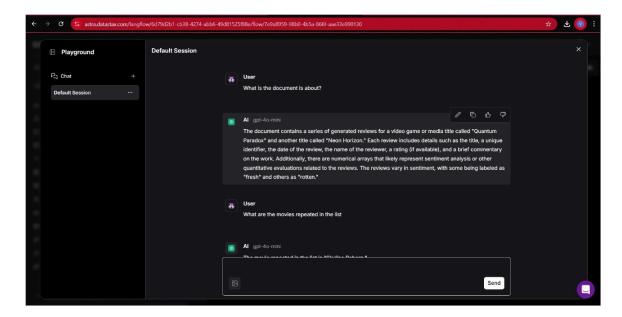
### 4. collections

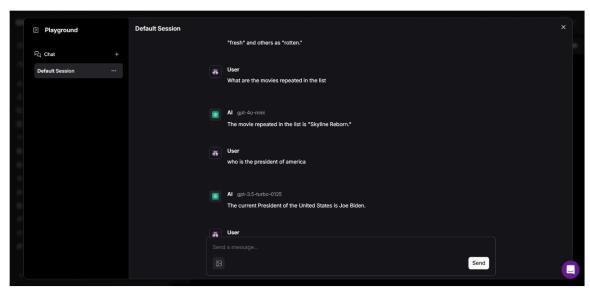


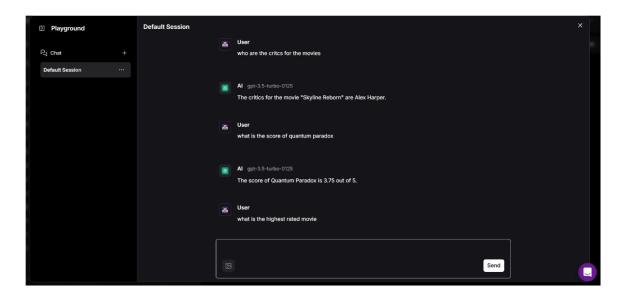
5. notifications of completed prompts

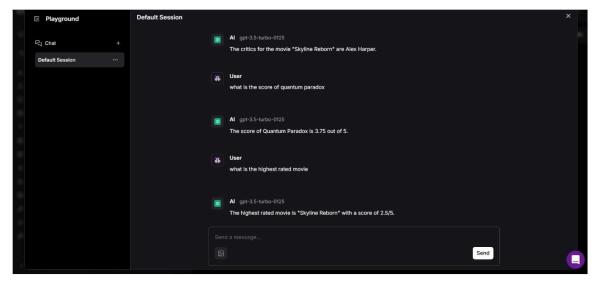


6.prompts and session results

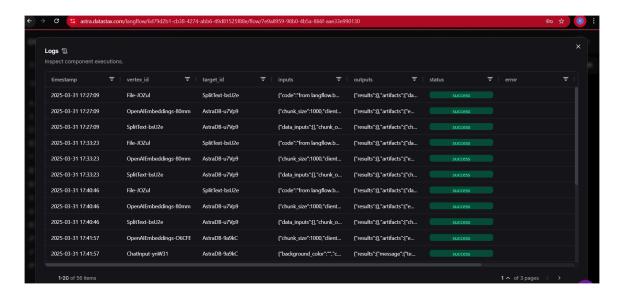




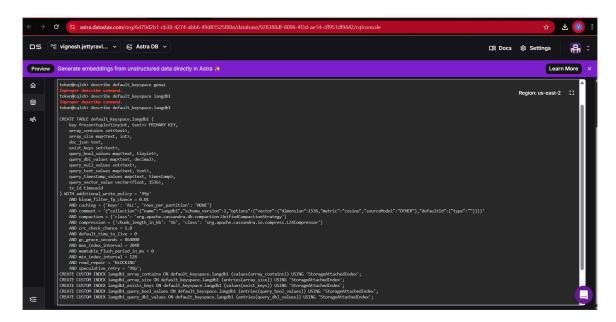




### 7. Logs



#### 8.dbschema

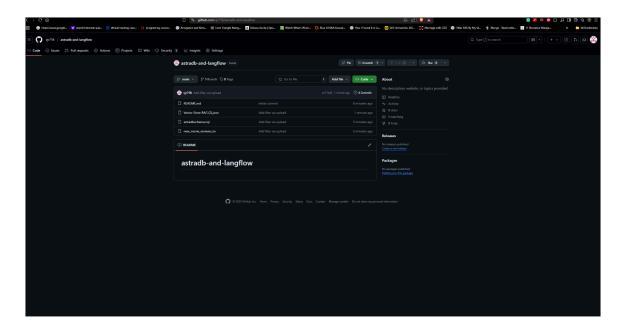


### 7. Future Enhancements

- Implement a feedback loop to refine and improve answers over time.
- Deploy as a public chatbot using Streamlit or Gradio for wider accessibility.
- Integrate LangSmith to enhance observability and monitoring.
- Leverage a fine-tuned LLM for improved domain-specific responses.
- **Explore Hugging Face open-source LLMs** as a cost-effective alternative.

Github repository for the assignment

https://github.com/vjr718/astradb-and-langflow



->Due to the recent update we can't share the deployment model of langflow based genai pipeline.

## so I have given the link here

https://astra.datastax.com/org/6d79d2b1-cb38-4274-abb6-49d81525f88e/database/928388df-8096-4f3d-ae14-cff951df9442/data-explorer

https://astra.datastax.com/langflow/6d79d2b1-cb38-4274-abb6-49d81525f88e/flow/7e9a8959-98b0-4b5a-866f-aae33e990130

------END------