# Finetuned RAG Systems Engineering Report

## Step 1: Infrastucture

- **Goal**: Set up your environment, including Qdrant for vector storage, MongoDB for raw data storage, Gradio App , ClearML.

- **Completed Tasks**:

  - Infrastructure appears ready, including Qdrant and MongoDB.

## Step 2:  ETL Pipeline

- **Goal**: Collect high-quality Q&A data, store it in MongoDB, and prepare it for feature engineering.

- **Completed Tasks**:

  - Blogs and other data sources like Medium articles have been ingested.

  - Data is formatted into a Q&A structure.

Added question-and-answer style blog posts to Medium to create a Q&A dataset

https://abdulrazzaq0902.medium.com/introduction-to-ros2-robot-operating-system-9fa9e9367263

## Step 3: Feature Engineering

EXPLORER

RAG_FINETUNING
- results
- docker-compose.yml
- gradio_test.py
- pipeline_ETL.py
- pipeline_FeatureEngineering.py
- pipeline_FineTune.py
- pipeline_Inference.py
- poetry.lock
- pyproject.toml
- qdrant_connection.py

pipeline_FeatureEngineering.py > feature_engineering_and_instruct_task

```python
from sentence_transformers import SentenceTransformer
from clearml import Task, PipelineController
from loguru import logger
import re
import numpy as np

# Configuration
MONGO_URI = "mongodb://llm_engineering:llm_engineering@127.0.0.1:27017"
DB_NAME = "etl_database_n"
COLLECTION_NAME = "platform_scraped_data"
QDRANT_HOST = "http://localhost:6333"


# Feature engineering and storage task
def feature_engineering_and_instruct_task(mongo_uri: str, db_name: str, collection_name: str, qdrant_host: str):
    QDRANT_COLLECTION_NAME = "instruct_dataset_clearml"
    VECTOR_DIMENSION = 384
    DISTANCE_METRIC = Distance.COSINE
    BATCH_SIZE = 64
    # Normalize vector for cosine similarity (optional, improves scores)
    def normalize_vector(vector):
        return vector / np.linalg.norm(vector)

    # Split text into smaller chunks for better context alignment
    def split_into_chunks(text: str, max_chunk_size: int = 150) -> list:
        sentences = re.split(r'(?<=[.!?])\s+', text)  # Split by sentence endings
        chunks = []
        current_chunk = []

        for sentence in sentences:
            if sum(len(s) for s in current_chunk) + len(sentence) <= max_chunk_size:
```

PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL  PORTS

TERMINAL

```
2024-12-08 18:40:06.590 | INFO     | __main__:feature_engineering_and_instruct_task:120 - Processed batch 5/6
2024-12-08 18:40:06.746 | INFO     | __main__:feature_engineering_and_instruct_task:120 - Processed batch 6/6
2024-12-08 18:40:06.746 | INFO     | __main__:feature_engineering_and_instruct_task:121 - Generated embeddings for 332 chunks. Embedding sh
ape: 384
2024-12-08 18:40:07.036 | INFO     | __main__:feature_engineering_and_instruct_task:130 - Upserted 64 records to Qdrant.
2024-12-08 18:40:07.097 | INFO     | __main__:feature_engineering_and_instruct_task:130 - Upserted 64 records to Qdrant.
2024-12-08 18:40:07.143 | INFO     | __main__:feature_engineering_and_instruct_task:130 - Upserted 64 records to Qdrant.
2024-12-08 18:40:07.202 | INFO     | __main__:feature_engineering_and_instruct_task:130 - Upserted 64 records to Qdrant.
2024-12-08 18:40:07.233 | INFO     | __main__:feature_engineering_and_instruct_task:130 - Upserted 64 records to Qdrant.
2024-12-08 18:40:07.240 | INFO     | __main__:feature_engineering_and_instruct_task:130 - Upserted 12 records to Qdrant.
2024-12-08 18:40:07.240 | INFO     | __main__:feature_engineering_and_instruct_task:131 - All records successfully upserted to Qdrant.
```

---

localhost:6333/dashboard#/collections/instruct_dataset_clearml

Finish update

(2) WhatsApp  New Tab  How to use Enhan...  Posit Cloud  Part Time Work Fr...  Surveys - Surveys  Gmail  YouTube  Maps  SRK Traders  All Bookmarks

drant

**Point 3**

Payload:

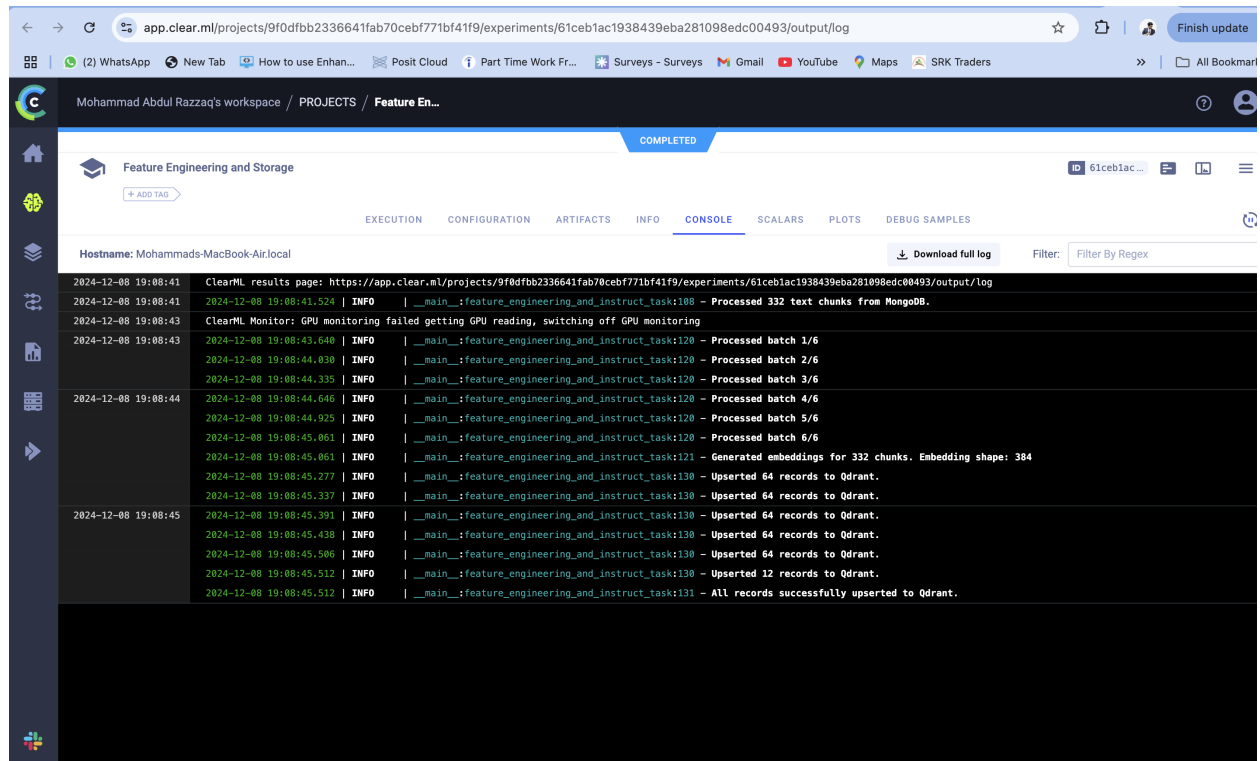| text | The work outlined here not only fills this void but also offers a practical guide for beginners and interested parties looking to replicate similar multi-robot simulations in Gazebo. The UR5 robotic arm, a versatile and widely-used model in robotics, serves as the core of this demonstration. By focusing on this particular model, the tutorial ensures a broad relevance to a multitude of potential robotic applications. |
|---|---|
| link | https://medium.com/@arshad.mehmood/enabling-multi-robot-arm-in-gazebo-for-ros2-dc18981c03c6 |

metadata

```
{ 4 Items
  "_id": "67562d8eb40f6a8778373e96"
  "platform": "medium"
  "content": "Title: Enabling Multi-Robot ARM in Gazebo for ROS2 …"
  "link": "https://medium.com/@arshad.mehmood/enabling-multi- …"
}
```
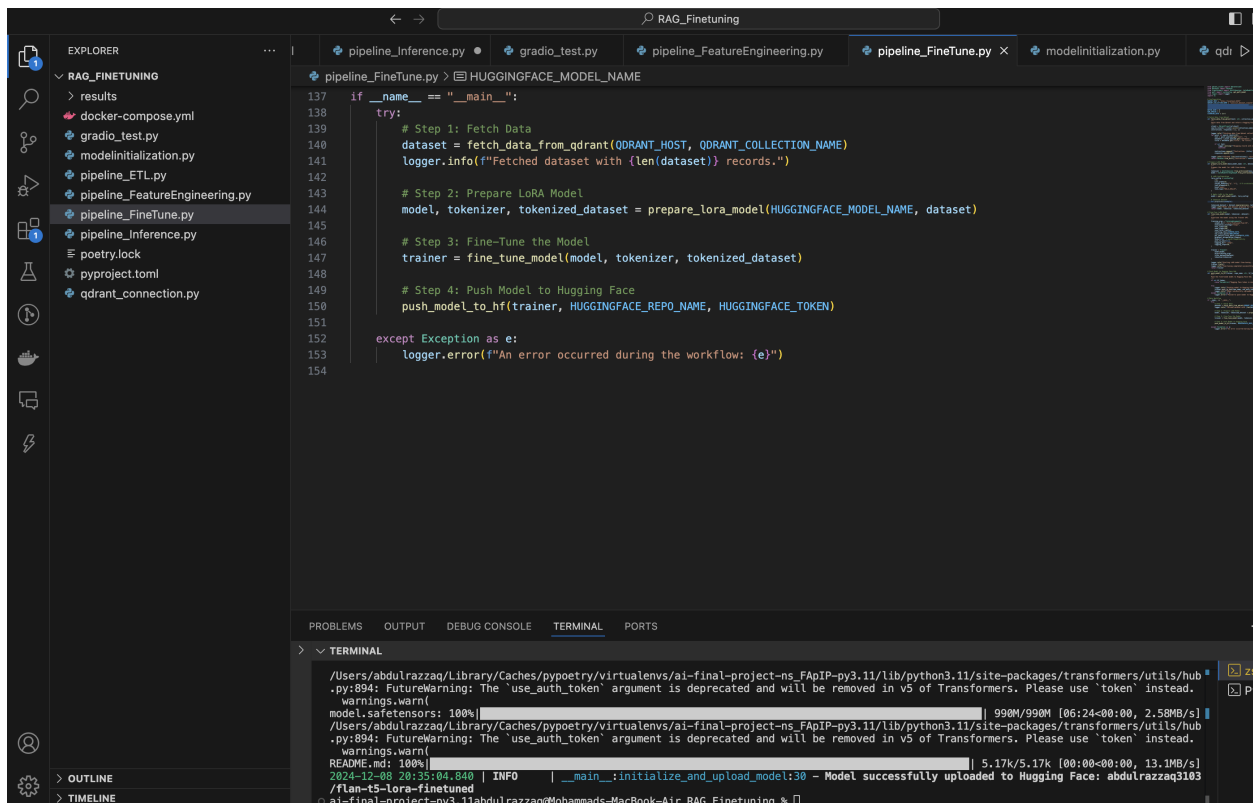
Vectors:

v1.12.4

Vectors:

- **Goal**: Extract and preprocess features, such as embeddings, for use in Qdrant.

- **Completed Tasks**:

    - Sentence embeddings generated using `all-MiniLM-L6-v2`.

    - Data has been chunked and upserted into Qdrant.

## Step 4: Fine-Tuning

- **Goal**: Fine-tune the `Flan-T5` model using LoRA to adapt it to your specific Q&A dataset.

- **Completed Tasks**:

  - Initial fine-tuning setup with LoRA and Hugging Face libraries.

## Step 5: Contextual Retrieval and Inference

- **Goal**: Implement the end-to-end RAG pipeline for context-based response generation.

- **Completed Tasks**:

  - Query embeddings are generated, and Qdrant retrieves relevant context.

- Prompt structure integrates context with the user query.