# **Al Project**

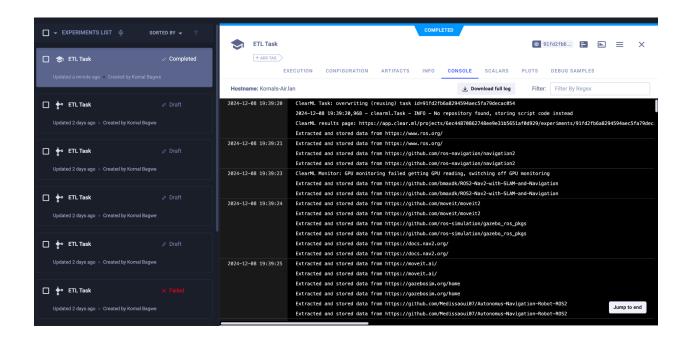
# **Configuration:**

## **ETL Pipeline:**

```
from bs4 import BeautifulSoup
from youtube_transcript_api import YouTubeTranscriptApi
import logging
mongo_client = MongoClient('mongodb://llm_engineering:llm_engineering@127.0.0.1:27017')
db = mongo_client['rag_system']
collection = db['raw_data']
logging.basicConfig(level=logging.INFO)
logger = logging.getLogger(__name__)
def etl_pipeline():
    # Initialize the task with the correct project name
task = Task.init(project_name="RAG_System", task_name="ETL Task")
    task_logger = task.get_logger()
    def extract_ros2_docs():
        ros2_docs_urls = [
         "https://www.ros.org/",
         "https://github.com/ros-navigation/navigation2",
        "https://github.com/bmaxdk/ROS2-Nav2-with-SLAM-and-Navigation",
        "https://github.com/moveit/moveit2",
"https://github.com/ros-simulation/gazebo_ros_pkgs",
         "https://docs.nav2.org/"
         "https://moveit.ai/",
         "https://gazebosim.org/home",
         "https://github.com/Medissaoui07/Autonomus-Navigation-Robot-ROS2",
         "https://github.com/ros-navigation/navigation2",
        "https://github.com/moveit/moveit2",
         "https://github.com/AndrejOrsula/ign_moveit2_examples/blob/master/docs/README.md",
         "https://github.com/oKermorgant/gz_moveit2_examples",
         "https://medium.com/schmiedeone/getting-started-with-ros2-part-1-d4c3b7335c71",
        "https://medium.com/@tetraengnrng/a-beginners-guide-to-ros2-29721dcf49c8",
    def extract_ros2_docs():
         "https://medium.com/robotics-zone/understanding-moveit2-and-ros-2-for-robotics-61d74832cd2a",
"https://medium.com/ros2-basics/building-your-own-ros2-navigation-stack-from-scratch-cc9c6b6a32e6",
         "https://medium.com/robotics-zone/introducing-ros2-and-gazebo-simulation-for-robotics-3e4054d4f8f8",
         "https://medium.com/robotics-zone/using-ros2-and-gazebo-to-simulate-robots-in-a-vibrant-world-34a6a4f35b28"
        for url in ros2_docs_urls:
                 response = requests.get(url)
                 soup = BeautifulSoup(response.content, 'html.parser')
                 text = soup.get_text()
                 collection.insert_one(data)
                 task_logger.report_text(f"Extracted and stored data from {url}")
                print(f"Extracted and stored data from {url}") # Use standard logger
                 task_logger.report_text(f"Failed to extract data from {url}: {e}")
                 print(f"Failed to extract data from {url}: {e}") # Use standard logger
    def extract_youtube_videos():
        youtube_video_ids = ['7TVW\ADXwRw', 'rtrGoGsMV\li&list=PLgG0XDQqJckkSJDPhXsFU_RIqEh08nG0V',"sVUKeHMBtpQ","Xbij9Tst-WA","jkoGkAd0GYk","QI
         for video_id in youtube_video_ids:
                 transcript = YouTubeTranscriptApi.get_transcript(video_id)
                 transcript_text = ' '.join([t['text'] for t in transcript])
                 data = {'source': 'youtube', 'url': f"https://www.youtube.com/watch?v={video_id}", 'content': transcript_text}
                 collection.insert_one(data)
                 task_logger.report_text(f"Extracted and stored data from YouTube video {video_id}")
                print(f"Extracted and stored data from YouTube video {video_id}")
             except Exception as e:
                 task_logger.report_text(f"Failed to extract data from YouTube video {video_id}: {e}")
                 print(f"Failed to extract data from YouTube video {video_id}: {e}")
```

```
def extract_ros2_docs():
    "https://medium.com/robotics-zone/understanding-moveit2-and-ros-2-for-robotics-61d74832cd2a",
    "https://medium.com/ros2-basics/building-your-own-ros2-navigation-stack-from-scratch-cc9c6b6a32e6",
    "https://medium.com/robotics-zone/introducing-ros2-and-gazebo-simulation-for-robotics-3e4054d4f8f8",
    "https://medium.com/robotics-zone/using-ros2-and-gazebo-to-simulate-robots-in-a-vibrant-world-34a6a4f35b28"
    for url in ros2_docs_urls:
           response = requests.get(url)
            soup = BeautifulSoup(response.content, 'html.parser')
            text = soup.get_text()
            data = {'source': 'github', 'url': url, 'content': text}
collection.insert_one(data)
            task_logger.report_text(f"Extracted and stored data from {url}")
            print(f"Extracted and stored data from {url}") # Use standard logger
        except Exception as e:
            task_logger.report_text(f"Failed to extract data from {url}: {e}")
            print(f"Failed to extract data from {url}: {e}") # Use standard logger
def extract voutube videos():
    youtube_video_ids = ['TTVWlADXwRw', 'rtrGoGsMVlI&list=PLgG0XDQqJckkSJDPhXsFU_RIqEh08nG0V',"sVUKeHMBtpQ","Xbij9Tst-WA","jkoGkAd0GYk","QI
    for video_id in youtube_video_ids:
            transcript = YouTubeTranscriptApi.get_transcript(video_id)
            transcript_text = ' '.join([t['text'] for t in transcript])
            data = {'source': 'youtube', 'url': f"https://www.youtube.com/watch?v={video_id}", 'content': transcript_text}
            collection.insert one(data)
            task_logger.report_text(f"Extracted and stored data from YouTube video {video_id}")
print(f"Extracted and stored data from YouTube video {video_id}")
        except Exception as e:
            task_logger.report_text(f"Failed to extract data from YouTube video {video_id}: {e}")
            print(f"Failed to extract data from YouTube video {video_id}: {e}")
```

```
etl_pipeline():
                                                                                                                 Aa <u>ab</u> * No results \uparrow \downarrow = \times
                                                                                        > fetch_similar_data
   def extract_youtube_videos():
               print(f"Failed to extract data from YouTube video {video_id}: {e}")
   def print ingested urls():
       print("Ingested URLs:")
       print("Ingested URLs:")
       for url in urls:
           logger.info(url)
           print(url)
   extract_ros2_docs()
   extract_youtube_videos()
   print_ingested_urls()
   task.close()
if __name__ == "__main__":
   etl_pipeline()
```



### **Feature Engineering:**

```
> fetch_similar_data
                                                                                                                 Aa <u>ab</u> \blacksquare* No results \uparrow \downarrow \equiv \times
import torch
from pymongo import MongoClient
from bson import ObjectId
from qdrant_client.models import VectorParams, Distance
mongo_client = MongoClient('mongodb://llm_engineering:llm_engineering@127.0.0.1:27017')
db = mongo_client['rag_system']
collection = db['raw_data']
qdrant = qdrant_client.QdrantClient("http://127.0.0.1:6333")
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
model = BertModel.from_pretrained('bert-base-uncased')
collection_name = "raw_data"
collections response = gdrant.get collections()
if not any(col.name == collection_name for col in collections_response.collections):
    qdrant.create_collection(
    collection name=collection name,
    vectors config=VectorParams(
       size=768,
        distance=Distance.COSINE
def extract_features(text):
    inputs = tokenizer(text, return_tensors="pt", padding=True, truncation=True)
    with torch.no_grad():
       outputs = model(**inputs)
    embeddings = outputs.last_hidden_state.mean(dim=1).squeeze().numpy()
    return embeddings
```

```
vectors_config=VectorParams(
                                                                                   > fetch_similar_data
       size=768,
def extract_features(text):
   inputs = tokenizer(text, return_tensors="pt", padding=True, truncation=True)
   with torch.no_grad():
      outputs = model(**inputs)
   embeddings = outputs.last_hidden_state.mean(dim=1).squeeze().numpy()
   return embeddings
def push_features_to_qdrant():
    for doc in collection.find():
       embedding = extract_features(content)
       point_id = str(uuid.uuid5(uuid.NAMESPACE_OID, str(doc['_id'])))
       payload = {
            "embedding": embedding.tolist(),
       point = models.PointStruct(id=point_id, vector=embedding.tolist(), payload=payload)
       qdrant.upsert(collection_name=collection_name, points=[point])
push_features_to_qdrant()
```



## FineTuning:

```
vectors_config=VectorParams(
                                                                                    > fetch_similar_data
        size=768,
        distance=Distance.COSINE
def extract_features(text):
   inputs = tokenizer(text, return_tensors="pt", padding=True, truncation=True)
       outputs = model(**inputs)
   embeddings = outputs.last_hidden_state.mean(dim=1).squeeze().numpy()
   return embeddings
def push_features_to_qdrant():
   for doc in collection.find():
       content = doc['content']
       embedding = extract_features(content)
       point_id = str(uuid.uuid5(uuid.NAMESPACE_OID, str(doc['_id'])))
            "embedding": embedding.tolist(),
           "metadata": {"source": doc['source'], "url": doc['url']}
       point = models.PointStruct(id=point_id, vector=embedding.tolist(), payload=payload)
       qdrant.upsert(collection_name=collection_name, points=[point])
       print(f"Features for {doc['url']} pushed to Qdrant.")
push_features_to_qdrant()
```

```
learning_rate=2e-5,
                                                                                          > fetch_similar_data
                                                                                                                    Aa ab * No results \uparrow \downarrow \equiv \times
    per_device_train_batch_size=1,
    gradient_accumulation_steps=16,
    num train epochs=3,
    weight_decay=0.01,
    save_strategy="epoch",
    logging_dir="./logs",
    logging_steps=10,
    fp16=False,
   report_to="clearml",
optim="adamw_torch",
    max_grad_norm=0.3,
torch.cuda.empty_cache()
model.gradient_checkpointing_enable()
trainer = Trainer(
    model=model.
    args=training_args,
    train dataset=tokenized dataset.
    data_collator=data_collator,
    model.save_pretrained("./fine_tuned_llama2")
    tokenizer.save_pretrained("./fine_tuned_llama2")
```

```
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

- Waiting for previous model to upload (2 pending, https://files.clear.ml/HuggingFace Transformers/Trainer.d6679e7829b744558 19e75b2b1ede8d8/models/checkpoint-3.zip)

- Waiting for previous model to upload (2 pending, https://files.clear.ml/HuggingFace Transformers/Trainer.d6679e7829b744558 19e75b2b1ede8d8/models/checkpoint-3.zip)

| 100% | 1424.25/1424.25 MB [09:34<00:00, 2.48MB/s]:
| 2024-12-08 23:01:25,021 - clearml.Task - INFO - Completed model upload to https://files.clear.ml/HuggingFace%20Transformers/Trainer.d6679e7829b74455819e75b2b1ede8d8/models/checkpoint-2.zip
| 2024-12-08 23:01:25,000 - clearml.storage - INFO - Uploading: 1424.25MB to /var/folders/66/q21rnv5143b58h0kj1w8rgmr0000gn/T/m odel_package.fibb8806.zip
| ('train_runtime': 1196.2704, 'train_samples_per_second': 0.008, 'train_steps_per_second': 0.003, 'train_loss': 0.425259470939 63623, 'epoch': 3.0}
| 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100
```