



Vietnamese-German University



VIETNAMESE – GERMAN UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE

FRANKFURT UNIVERSITY OF APPLIED SCIENCES
FACULTY OF COMPUTER SCIENCE AND ENGINEERING

Distributed Systems: Multi-index search with LLM

Team 20: Ha Nguyen Nhat Minh - 10421038
Nguyen Ha Phan - 10421099
Bui Vu Khoa - 10421028
Le Nguyen Quynh An - 10421118
Nguyen Xuan Tien - 10421058

Lecturer: Dr. Phan Trong Nhan

BINH DUONG, VIETNAM

Contents

1	Problem Statement	4
1.1	Context and Background	4
1.2	Existing Solution and Limitations	4
1.3	Proposed Approach	5
1.4	Supporting Elements and Scope	6
2	System Architecture	8
2.1	Overview	8
2.1.1	Client	10
2.1.2	Backend Services	10
2.1.3	Database Storage	11
3	Chat Application	
	Functionality	12
3.1	Core Functionalities	12
3.1.1	Natural Language Querying:	12
3.1.2	Multi-Departmental Search	12
3.1.3	Context-Aware Answering	12
3.1.4	Real-Time Asynchronous Processing:	13
3.2	Additional Functionalities	13
3.3	Technology Stack	14
4	Distributed Systems	
	Characteristics	17
4.1	Scalability	17
4.2	Fault Tolerance	19
4.3	Concurrency	23

4.4 Data Consistency	24
5 Implementation Details	26
5.0.1 Project Set-up	26
5.0.2 Package installation	26
5.0.3 Backend Routing and Task Set-up	30
5.1 Database and Vector Index Implementation	33
5.1.1 Database Set-up	33
5.1.2 Data Schema in PostgreSQL and Qdrant	37
5.2 Back-end Implementation	48
5.2.1 Core Functionality Set-up	48
5.2.2 Functionality Implementation	56
5.3 Front-end Implementation	64
5.4 Source Code	73
6 Conclusions	74
6.1 Development Potential	75

Chapter 1

Problem Statement

1.1 Context and Background

In recent years, Large Language Models (LLMs) like ChatGPT, Claude, and Gemini have grown in strength and popularity in a variety of fields. These models show remarkable abilities in the generation, reasoning, and understanding of natural language; making them highly popular in with common people for day-to-day usage across a variety of simple tasks. However, most commonly accessible implementations of the technology suffers issues with memory and factual consistency, especially when it comes to domain-specific or real-time tasks. To enhance their performance on this front, LLMs are frequently combined with Retrieval-Augmented Generation (RAG) systems, giving them better capabilities for working with information that was not initially built into their training data. This project aims to allow the querying and processing of multi-index distributed databases LLMs using implementations of the RAG technology.

1.2 Existing Solution and Limitations

RAG presents a way to improve LLM outputs by obtaining pertinent documents from an outside knowledge base, but it also presents a number of difficulties.

- **Dependency on embedding quality:** Relevant context might not be recovered if the embedding model fails to accurately capture semantic similarity. For instance, a simple query for a numeric value or specific code may not align well with the semantic space of the document vectors, leading to missed retrievals

- **Single-source bottleneck:** A central point of failure and limited flexibility are introduced by the majority of RAG systems' reliance on a single vector store or knowledge base.
- **Problems with latency and scaling:** RAG systems may experience higher latency and decreased responsiveness as document size or retrieval scope increases.
- **Over-reliance on static context:** Retrieved passages may take up too much space in the prompt and limit the underlying LLM's ability to generate new content.

1.3 Proposed Approach

To solve these problems, this project proposes a smarter chat bot that uses a two-step hybrid search process.

- **1. Understand the User's Real Question:** First, the chat bot sends the user's query to a fast Large Language Model (like Mistral). This LLM acts as a parser, pulling out specific filters (i.e. department name, job title, etc.) and structures them as a JSON object.
- **2. Perform a "Smart Search":** The application combines the specific filters from the LLM with a general semantic search in the Qdrant Vector Database.

The end result is a useful chatbot platform that enables comparative querying across models and data sources, such as DeepSeek R1, GPT 4.1, and Mistral, facilitating more resilient and contextually aware interactions in multi-agent settings.

1.4 Supporting Elements and Scope

Large Language Models (LLMs) such as GPT, Claude, and Gemini continue to excel in natural language processing and are increasingly being used in applications that require domain-specific or real-time knowledge. Retrieval-Augmented Generation (RAG) has emerged as an effective method for improving LLM outputs by retrieving relevant data from external knowledge bases. However, most existing RAG systems rely on a single vector store or index, which presents significant scalability, fault tolerance, and domain isolation challenges, especially in complex, multi-departmental environments.

In a real-world organization, different teams (e.g., HR, Engineering, Finance) frequently work with different datasets, have different access requirements, and require domain-specific interpretations. A single centralized index fails to address these needs, resulting in:

- Retrieval noise from irrelevant domains.
- Difficulty in managing updates or access policies per department.
- A single point of failure that reduces robustness.

As a result, using multiple indexes, each representing a distinct domain or department, provides numerous benefits, including increased relevance, scalability, and modular design. However, this distributed setup presents new technical challenges that must be explicitly addressed.

- **Routing Complexity:** Determine which index or set of indexes to query based on user intent.
- **Load balancing and failure isolation:** Ensuring that no index is overloaded while others are underutilized, and that failures in one do not affect the other.

- **Consistency management:** entails keeping multiple indexes in sync and up to date across departments, preventing stale responses.

This project looks into a distributed, multi-index RAG architecture to address these limitations. It aims to create a system that understands user intent using lightweight LLM parsing, performs hybrid retrieval across multiple Qdrant collections, and manages routing, load balancing, and consistency in a scalable manner. The use of simulated multi-department data creates a realistic scenario for assessing the design under distributed conditions.

Chapter 2

System Architecture

2.1 Overview

Our system is built on a distributed microservices architecture designed for scalability and fault tolerance. The flow begins with a User submitting a User Query to the Chainlit Client web application. This query is then dispatched as a task to a Task dispatch module, which pushes it to a Redis Broker (Task Queues).

From the queues, a pool of Celery Workers pulls and processes tasks based on their type. These workers interact with various data storage components, including a PostgreSQL DB for source data, a Qdrant DB for vector indexes, and AWS S3 for logging storage. The workers also make calls to external LLM APIs for tasks like filter extraction and synthesis. Once the processing is complete, the Chainlit client receives the processed context from the workers and synthesises the final answer to display back to the user.

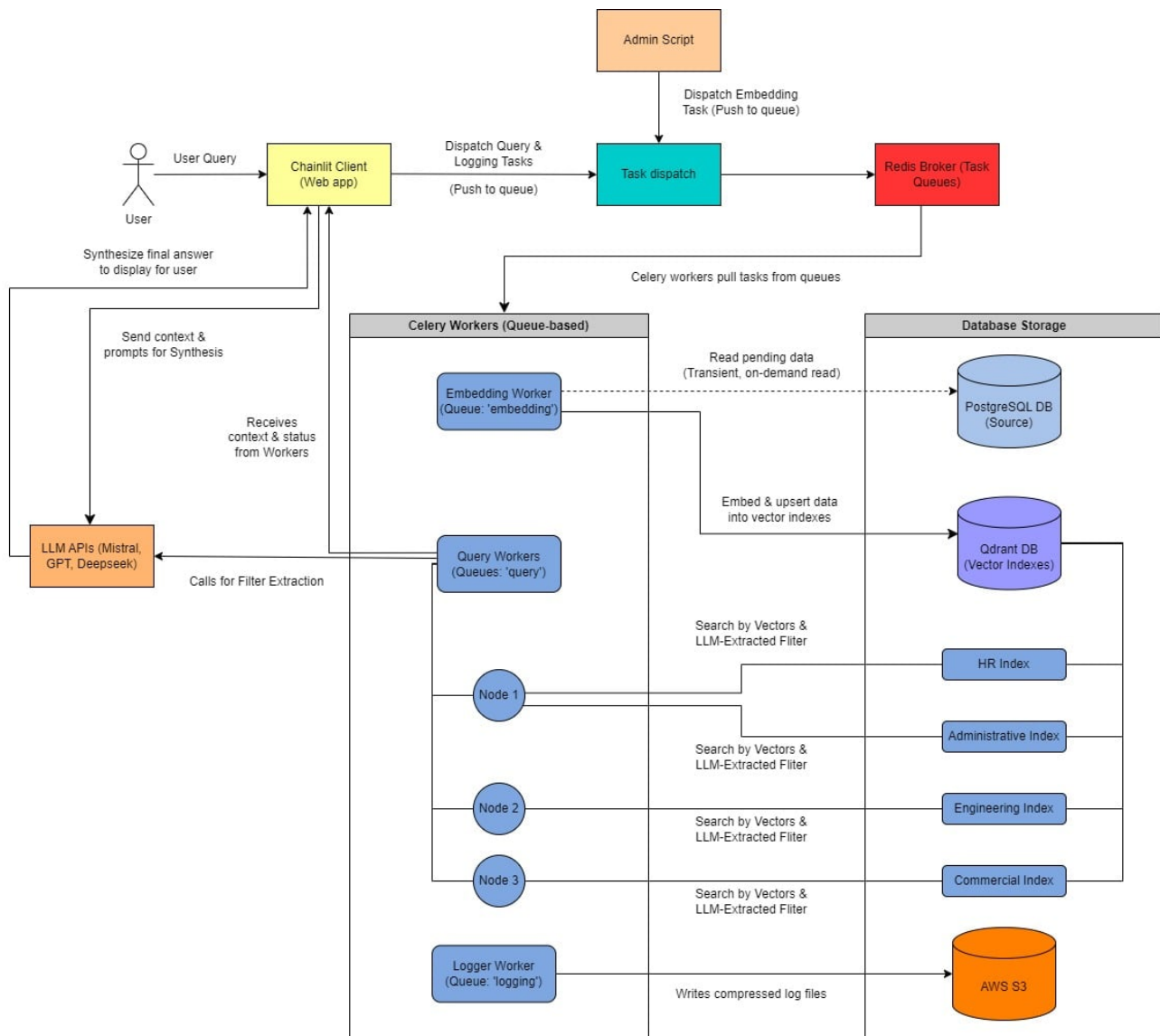


Figure 2.1: Employee Q&A System Architecture

2.1.1 Client

The Chainlit Client (Web app) serves as the user-facing web application. It is responsible for handling user input and displaying the chat interface. It dispatches query and logging tasks to the task queue. After receiving context and status from the workers, the client makes calls to LLM APIs for filter extraction and synthesis of the final answer, which is then displayed to the user.

2.1.2 Backend Services

The backend is composed of several decoupled services that work together through the task queue system, as depicted in the diagram.

- **Admin Script:** This is a command-line interface (`main.py`) used for administrative actions, specifically for dispatching embedding generation tasks to the queue.
- **Task Dispatch:** This is the initial entry point for tasks from the frontend or the admin script, pushing them into the Redis Broker.
- **Redis Broker (Task Queues):** Configured via `config.py` as the `CELERY_BROKER_URL`, Redis acts as the central message broker for Celery. It holds different queues (`embedding`, `query`, `logging`) to manage and distribute tasks to the workers.
- **Celery Workers (Queue-based):** These are stateless worker processes that pull tasks from the Redis queues and execute them asynchronously. The diagram shows three types of workers, each handling a specific queue:
 - **Embedding Worker:** This worker pulls tasks from the `embedding` queue. It reads pending data from the PostgreSQL database, generates vector embeddings using the `all-MiniLM-L6-v2` model, and upserts the data into the Qdrant vector indexes. This process is imple-

mented as a chained task (`prepare_employee_vectors_task` followed by `upsert_vectors_task`) to ensure a sequential flow.

- **Query Workers:** These workers pull tasks from `query` queues, which are dynamically routed based on the Qdrant collection being queried. Their primary role is to search the relevant Qdrant indexes (e.g., `HR Index`, `Administrative Index`) using the query’s vector and filters extracted by an LLM.
- **Logger Worker:** This worker pulls tasks from the `logging` queue. Its sole responsibility is to write compressed chat log files to **AWS S3** for long-term storage.
- **LLM APIs (Mistral, GPT, Deepseek):** These external services are called by the Chainlit client to perform two distinct functions: extracting structured filters from the user query and synthesizing the final human-readable answer from the retrieved context.

2.1.3 Database Storage

The system uses a multi-layered storage approach to handle different data types and access patterns.

- **PostgreSQL DB (Source):** This is the primary source of truth for all structured employee data. The `db_pool_setup.py` module manages a connection pool to ensure efficient and robust access from the Celery workers.
- **Qdrant DB (Vector Indexes):** This is a high-performance vector database that holds vector indexes for different departments. The diagram shows multiple indexes, which correspond to collections in Qdrant as defined in `config.py`. This multi-index design enables parallel and targeted searches.
- **AWS S3:** This is an object storage service used for the long-term archival of compressed chat logs, written to by the Logger Worker.

Chapter 3

Chat Application

Functionality

The Employee Q&A System is designed to provide users with a powerful and interactive experience for retrieving employee information. Here are the core and additional functionalities implemented:

3.1 Core Functionalities

3.1.1 Natural Language Querying:

The system allows users to ask questions about employees using natural language queries, such as "Who is the lead engineer on the new security project?" or "Find all employees with Python skills in the Sales department."

3.1.2 Multi-Departmental Search

The system can simultaneously search for relevant employee information across multiple departments (*e.g., Engineering, Sales, HR, Marketing, Finance, Operations*) by querying different Qdrant collections in parallel. This ensures comprehensive results and reduces search time. The departments are configured in `config.py`

3.1.3 Context-Aware Answering

Based on the retrieved context from the vector database, an LLM synthesises a clear and concise answer to the user's query. The LLM is instructed to use only the provided context and to state if the information is not available.

3.1.4 Real-Time Asynchronous Processing:

The system leverages Celery to process user queries asynchronously. When a user submits a query, tasks are dispatched to workers in the background, allowing the frontend to remain responsive. A loading spinner indicates that the system is searching for information.

3.2 Additional Functionalities

1. Structured Filter Extraction:

The system can automatically analyse a user's query and extract structured filters such as `department`, `job_title`, `employee_id`, `full_name`, `location`, `email`, and `min_satisfaction`. These filters are used to refine the search results in the vector database, improving retrieval accuracy.

2. LLM Selection:

The user interface displays a live task list with a running status, showing that the system is actively searching across different departments for context.

3. Live Task Tracking:

The user interface displays a live task list with a running status, showing that the system is actively searching across different departments for context.

4. Conversation History Logging:

User queries and the assistant's responses are buffered in Redis and then asynchronously stored as compressed JSON files in an AWS S3 bucket for long-term storage and analysis.

5. Configurable Relevance Threshold:

The system only considers retrieved context if its similarity score exceeds a configurable `relevance threshold` of 0.70, ensuring that only highly relevant infor-

mation is used to answer the query.

3.3 Technology Stack

The application is built on a robust and scalable technology stack, leveraging modern tools and frameworks to ensure high performance, real-time capabilities, and maintainability.

1. PostgreSQL:

As the database solution for storing raw, structured employee data, PostgreSQL is used. It is a powerful open-source relational database that is highly reliable and extensible. The system utilises a connection pool to manage database connections efficiently, as configured in `db_pool_setup.py`.

2. Sentence-Transformers:

The `sentence-transformers` library is used to generate dense vector embeddings from employee profile text. The specific model used is `all-MiniLM-L6-v2`, which is configured in `config.py` with a dimension of 384.

3. Qdrant:

Qdrant is a high-performance vector database used to store the generated embeddings. It supports fast semantic search and filtering on payloads, which is crucial for the retrieval phase of the RAG pipeline. The application interacts with Qdrant via the `qdrant_client` library. The local instance of the Qdrant service is hosted using Docker.

4. Celery:

Celery serves as the distributed task queue for asynchronous processing. It is used to dispatch and manage background tasks such as embedding generation, vector upserting, and parallel queries to different departments, ensuring the main

application remains responsive.

5. Redis:

Redis is utilised for two main purposes: as the message broker for Celery and as a temporary buffer for chat logs. It provides a fast, in-memory data store for these functions. The Redis server is run using a Docker container, simplifying its deployment and management.

6. Chainlit:

Chainlit is the framework used for building the conversational user interface. It handles user input and displays messages, and also provides UI elements like task lists and action buttons for selecting LLMs.

7. LLM APIs (Mistral, OpenAI, DeepSeek):

The system integrates with multiple large language models via their respective APIs. The OpenAI client is used for interacting with all three models, with configurations for each defined in `config.py`.

8. AWS S3:

Amazon S3 is used for durable, long-term storage of compressed chat logs. Logs are buffered in Redis and then asynchronously uploaded to a designated S3 bucket.

9. Boto3:

The `boto3` library is used to interact with AWS services, specifically to upload compressed chat logs to the S3 bucket.

10. Docker:

The modular, microservice-oriented architecture leverages Docker to containerise key services. This approach simplifies deployment and ensures consistent environ-

3.3. TECHNOLOGY STACK

ments. Specifically, Docker is used to run the Redis server and host the Qdrant vector database locally, as specified by the `QDRANT_HOST` and `REDIS_HOST` environment variables.

Chapter 4

Distributed Systems

Characteristics

4.1 Scalability

An application's capacity to scale up or down in response to changing user load is the main goal of scalability testing, a subset of performance testing. Scalability is important for real-time chat applications since it affects the application's capacity to support multiple users at once and continue to function well under heavy load.

Our system is designed to scale horizontally to handle increased load through several key architectural decisions:

- **Celery and Redis for asynchronous task processing.** This design decouples resource-intensive operations from the user-facing application, ensuring the system remains highly responsive even under heavy load.
- **Asynchronous Task Dispatch:** The user query is not processed synchronously by the web application. Instead, the Chainlit frontend dispatches a task to the Celery broker (Redis). This allows the web application to immediately handle new user requests without waiting for the query processing to complete.

Below is the code snippet from `chainlit_app.py` that dispatches the tasks:

```
1 departments_to_query = config.DEPARTMENTS
2     task_group = group(
3         process_employee_query_task.s(query=query_text, ←
4             target_department=dept)
5         for dept in departments_to_query
6     )
7     group_result = task_group.apply_async()
```

- **Parallel Query Execution:** For a single user query, the `chainlit_app.py` frontend dispatches a Celery group of tasks, which runs a separate `process_employee_query_task` for each configured department simultaneously. This parallel processing drastically reduces the total search time, as searches across different Qdrant collections occur concurrently. The use of a Redis task queue allows these tasks to be executed in parallel by multiple Celery workers.
- **Horizontal Scaling of Workers:** The use of a Redis broker and Celery workers allows for adding more worker processes to handle a higher volume of tasks concurrently. The `ProjectRouter` class dynamically routes tasks to specific queues (`query-collection_name`, `embedding`, `logging`) based on the task type or department. This enables load balancing and resource allocation strategies, where workers can be dedicated to specific, high-demand queues, such as those handling queries for a busy department.
- **Batch Processing for Efficiency:** The embedding generation process fetches employee IDs in batches of 90 (`QDRANT_UPSERT_BATCH_SIZE`) and processes them in a single task. This asynchronous batching optimizes database reads and Qdrant upsert operations, improving the efficiency of bulk data ingestion and preventing the ingestion process from blocking other

tasks.

- **Asynchronous Logging:** The `buffer_chat_log` function adds logs to a Redis list and then dispatches a Celery task to store them in S3 only when a certain `log_buffer_threshold` is reached. This prevents synchronous logging from becoming a bottleneck during peak usage and allows the system to handle high log volume without performance degradation.

4.2 Fault Tolerance

Our real-time Employee Q&A System is designed with a fault-tolerant mechanism that guarantees business continuity. This is achieved through several key mechanisms.

Firstly, the application uses a **database connection pool** to maintain a set of ready-to-use database connections. This prevents the application from failing under high load by ensuring connections are efficiently reused. The `db_connection()` context manager ensures that connections are always returned to the pool after use, even if an exception occurs during a database transaction, preventing resource leaks.

Below is the code snippet from `db_pool_setup.py` that demonstrates the context manager's rollback and return logic:

```
1 @contextmanager
2 def db_connection():
3     conn = None
4     pool = get_pool()
5
6     try:
7         conn = pool.getconn()
8         if not conn:
9             raise ConnectionError("Failed to acquire a database ↵
                                   connection from the pool")
```

4.2. FAULT TOLERANCE

```
10
11     logger.debug(f"Acquired connection {id(conn)} from pool")
12
13     yield conn
14
15     conn.commit()
16     logger.debug(f"Committed transaction on connection {id(conn)}↵
17         ")
18 except Exception as e:
19     if conn:
20         try:
21             conn.rollback()
22             logger.warning(f"Rolled back connection {id(conn)} ↵
23                 due to an exception")
24         except Exception as rollback_err:
25             logger.exception(f"Rollback failed on connection {id(↵
26                 conn)}: {rollback_err}")
27
28     logger.exception("Error during DB transaction")
29     raise
30 finally:
31     if conn:
32         if db_pool:
33             try:
34                 pool.putconn(conn)
35                 logger.debug(f"Returned connection {id(conn)} to ↵
36                     pool")
37             except Exception as e:
38                 logger.exception(f"Error returning connection to ↵
39                     pool: {e}")
40
41         try:
```

```
41         conn.close()
42         logger.warning(f"Force-closed connection {id(conn)} after pool failure")
43
44     except Exception as e2:
45         logger.exception(f"Failed to force-close connection {id(conn)}: {e2}")
46
47     else:
48         logger.warning(
49             f"Pool object (or global db_pool) became None in finally. "
50             f"Attempting to close conn id={id(conn)}."
51         )
52     try:
53         conn.close()
54
55     except Exception as e_close:
56         logger.exception(f"Failed to force-close connection id={id(conn)} when pool was None: {e_close}")
```

Furthermore, the application implements robust **error handling and retry mechanisms** for asynchronous tasks. Celery tasks like `prepare_employee_vectors_task`, `upsert_vectors_task`, and `process_employee_query_task` are configured with `max_retries=3`. This ensures that temporary failures, such as network disruptions or service unavailability, are handled gracefully by automatically re-attempting the operation after a delay.

Below is a snippet from `prepare_employee_vectors_task` showing the retry logic:

```
1  except Exception as e:
2      logger.warning(f"Task for {table_name} failed on embedding ←
```

```

    attempt {self.request.retries + 1}. Error: {e}")
3      try:
4          raise self.retry(exc=e, countdown=60)
5      except MaxRetriesExceededError:
6          logger.error(f"Task for {table_name} has failed ↵
              permanently. Marking records as FAILED.")
7          update_employee_status_in_db(employee_ids_batch, ↵
              table_name, ProcessingStatus.FAILED)
8      raise e

```

If a task fails permanently after all retries, the system updates the corresponding records in the database to a **FAILED** status, providing a persistent record of the failure and preventing indefinite retries.

The application also uses a **threading lock** during the initialization of the database pool to prevent race conditions and ensure the pool is set up only once, even if called by multiple threads.

Finally, the system ensures **graceful degradation** and **component-level skipping**. As seen in the `query_service.py` code, if a Qdrant collection for a specific department is not found, the system logs a warning and skips that department's search instead of failing the entire query. This allows the system to continue searching in other available collections, ensuring a partial result can still be provided to the user. Here is the relevant code snippet:

```

1  dept_info = config.DEPARTMENT_CONFIG.get(target_department)
2      if not dept_info:
3          logger.warning(f"No config found for department '{↵
              target_department}'. Skipping.")
4          return {"status": "skipped_no_config", "context": None, "↵
              source": target_department}
5
6      collection_name = dept_info['qdrant_collection']
7

```

```
8         try:
9             client.get_collection(collection_name=collection_name)
10        except Exception:
11            logger.warning(f"Collection '{collection_name}' not found↵
12                        . Skipping department '{target_department}'.".)
13        return {"status": "skipped_no_collection", "context": ↵
14                None, "source": target_department}
```

4.3 Concurrency

The system is designed to handle multiple tasks simultaneously, which is a fundamental aspect of concurrency.

- **Parallel Task Execution:** When a user submits a query, the `chainlit_app.py` component uses Celery's `group` primitive to dispatch a separate `process_employee_query_task` for each department defined in the configuration. This allows the searches to run concurrently on different worker nodes, enabling the system to retrieve context from all relevant departments at the same time. The `await await_group_task_result` function then polls the `GroupResult` to wait for all the parallel tasks to complete. Below is the snippet demonstrating this parallel execution:

```
1 task_group = group(
2     process_employee_query_task.s(query=query_text, ↵
3     target_department=dept)
4     for dept in departments_to_query
5 )
6 group_result = task_group.apply_async()
7
8 results = await await_group_task_result(group_result, ↵
    task_list)
```

- **Thread Safety:** A `threading.Lock` is used in `db_pool_setup.py` to ensure that the database connection pool is initialized only once, even if `initialize_pool()` is called by multiple threads concurrently. This prevents race conditions and ensures a consistent state for the connection pool in a multi-threaded environment.

4.4 Data Consistency

Maintaining consistency across multiple distributed data stores is a significant challenge. In our system, which uses both PostgreSQL and Qdrant, we ensure data consistency between the source data and its vector representation through a controlled, sequential task flow.

This is achieved by implementing the data ingestion process as a **Celery task chain**. A chain guarantees that a task will only execute after the preceding task in the chain has completed successfully.

The following code from `text_embedding.py` shows how a chain is created to link the vector preparation and upsert tasks:

```

1 task_chain = chain(
2     prepare_employee_vectors_task.s(batch_ids, ←
        table_name, dept_name),
3     upsert_vectors_task.s(collection_name, table_name←
        )
4 )
5 task_chain.apply_async(queue="embedding")

```

As shown, `upsert_vectors_task` is configured to run only after `prepare_employee_vectors_task` has successfully generated the embeddings. This ensures that invalid or incomplete vectors are not written to the Qdrant index.

Furthermore, a critical aspect of consistency is the atomic update of the embed-

ding status in the source database. The `embedding_status` in the PostgreSQL table is updated to `COMPLETED` only after the vectors have been successfully upserted to Qdrant. If the upsert task fails, the status remains `PENDING`, allowing the record to be picked up for re-processing in a subsequent run.

Here is the relevant code snippet from the `upsert_vectors_task`:

```
1 upserted_ids = [int(p.payload["employee_id"]) for p in point_structs]
2     update_employee_status_in_db(upserted_ids, table_name, ←
    ProcessingStatus.COMPLETED)
```

This status tracking mechanism guarantees that the `qdrant_collection` is eventually consistent with the `postgres_table` and provides a reliable way to monitor the data ingestion pipeline.

Chapter 5

Implementation Details

5.0.1 Project Set-up

5.0.2 Package installation

The project environment is configured by importing several key libraries necessary for the distributed RAG system. These include `psycopg2` for PostgreSQL connectivity, `qdrant_client` for vector database operations, `sentence-transformers` for creating embeddings, `celery` for task management, `redis` for the broker and buffer, `boto3` for AWS S3 interactions, and `openai` for LLM integrations. Environment variables, crucial for configuring services like Redis, PostgreSQL, and LLM APIs, are loaded using `python-dotenv`. The central configuration for the entire system is managed through the `config.py` module.

```
1 import os
2 import logging
3 from logging.handlers import RotatingFileHandler
4 from dotenv import load_dotenv
5
6 dotenv_path = os.path.join(os.path.dirname(__file__), '.env')
7 if os.path.exists(dotenv_path):
8     load_dotenv(dotenv_path=dotenv_path)
9
10 else:
11     load_dotenv()
12     if not os.getenv("POSTGRES_DBNAME"):
13         print("Warning: .env file not found or not loaded correctly.↵")
14         "
```

```
15 APP_NAME = "Employee_Q&A_System"
16
17 # Celery Configuration
18 CELERY_BROKER_URL = os.getenv("CELERY_BROKER")
19 CELERY_RESULT_BACKEND = os.getenv("CELERY_BACKEND")
20
21 # Redis Configuration
22 REDIS_HOST = os.getenv("REDIS_HOST")
23 REDIS_PORT = os.getenv("REDIS_PORT")
24 REDIS_DB = os.getenv("REDIS_DB")
25 REDIS_LOG_LIST_KEY_PREFIX = "chatlogs_list:"
26
27 # PostgreSQL Configuration
28 POSTGRES_DB_MIN_CONN = os.getenv("POSTGRES_DB_MIN_CONN")
29 POSTGRES_DB_MAX_CONN = os.getenv("POSTGRES_DB_MAX_CONN")
30 POSTGRES_DBNAME = os.getenv("POSTGRES_DBNAME")
31 POSTGRES_USER = os.getenv("POSTGRES_USER")
32 POSTGRES_PASSWORD = os.getenv("POSTGRES_PASSWORD")
33 POSTGRES_HOST = os.getenv("POSTGRES_HOST")
34 POSTGRES_PORT = os.getenv("POSTGRES_PORT")
35
36 DEPARTMENT_CONFIG = {
37     "Engineering": {
38         "postgres_table": "employees_engineering",
39         "qdrant_collection": "corp-dir-engineering"
40     },
41     "Sales": {
42         "postgres_table": "employees_sales",
43         "qdrant_collection": "corp-dir-commercial"
44     },
45     "HR": {
46         "postgres_table": "employees_hr",
47         "qdrant_collection": "corp-dir-hr"
48     },
49     "Marketing": {
50         "postgres_table": "employees_marketing",
```

```
51         "qdrant_collection": "corp-dir-commercial"
52     },
53     "Finance": {
54         "postgres_table": "employees_finance",
55         "qdrant_collection": "corp-dir-administrative"
56     },
57     "Operations": {
58         "postgres_table": "employees_operations",
59         "qdrant_collection": "corp-dir-administrative"
60     }
61 }
62 DEPARTMENTS = list(DEPARTMENT_CONFIG.keys())
63
64 # AWS S3 Configuration
65 AWS_BUCKET_NAME = os.getenv("AWS_BUCKET_NAME")
66
67 # Qdrant Configuration ---
68 QDRANT_HOST = os.getenv("QDRANT_HOST", "localhost")
69 QDRANT_PORT = int(os.getenv("QDRANT_PORT", 6333))
70 QDRANT_UPSERT_BATCH_SIZE = 90
71 QDRANT_NAMESPACE = "employee-data"
72
73 # Local Embedding Model
74 LOCAL_EMBEDDING_MODEL_NAME = "all-MiniLM-L6-v2"
75 LOCAL_EMBEDDING_MODEL_DIMENSION = 384
76
77 # LLM Models Information
78 MISTRAL_MODEL_CHOICE = 'Mistral'
79 MISTRAL_MODEL = 'mistral-small-latest'
80 OPENAI_MODEL_CHOICE = 'GPT 4.1'
81 OPENAI_MODEL = 'gpt-4.1-2025-04-14'
82 DEEPSEEK_MODEL_CHOICE = 'DeepSeekR1'
83 DEEPSEEK_MODEL = 'deepseek/deepseek-1:free'
84 DEEPSEEK_BASEURL = os.getenv("DEEPSEEK_BASEURL")
85
86 # LLM API Keys
```

```
87 OPENAI_API_KEY = os.getenv("OPENAI_API_KEY")
88 MISTRAL_API_KEY = os.getenv("MISTRAL_API_KEY")
89 DEEPSEEK_API_KEY = os.getenv("DEEPSEEK_API_KEY")
90
91 # Logging Configuration
92 LOG_DIR = os.path.join(os.path.dirname(__file__), '..', 'logs')
93 LOG_FILENAME = os.path.join(LOG_DIR, 'Employee_Q&A_System.log')
94 LOG_LEVEL = logging.INFO
95 LOG_FORMAT = '%(asctime)s - %(levelname)s - [(name)s:%(funcName)s] -<↵
    %(message)s'
96 DATE_FORMAT = '%Y-%m-%d %H:%M:%S'
97 LOG_MAX_BYTES = 5 * 1024 * 1024
98 LOG_BACKUP_COUNT = 3
99 LOG_BUFFER_THRESHOLD = 100
100 LOG_BUFFER_TTL_SECONDS = 3600
101
102 # Logging Setup Function
103 def setup_logging():
104     os.makedirs(LOG_DIR, exist_ok=True)
105     logger = logging.getLogger(APP_NAME)
106     logger.setLevel(LOG_LEVEL)
107
108     if logger.hasHandlers():
109         logger.handlers.clear()
110
111     formatter = logging.Formatter(LOG_FORMAT, datefmt=DATE_FORMAT)
112
113     # Rotating File Handler
114     file_handler = RotatingFileHandler(
115         LOG_FILENAME,
116         maxBytes = LOG_MAX_BYTES,
117         backupCount = LOG_BACKUP_COUNT,
118         encoding = 'utf-8'
119     )
120     file_handler.setLevel(LOG_LEVEL)
121     file_handler.setFormatter(formatter)
```

```
122     logger.addHandler(file_handler)
123
124     console_handler = logging.StreamHandler()
125     console_handler.setLevel(logging.INFO)
126     console_handler.setFormatter(formatter)
127     logger.addHandler(console_handler)
128
129     logger.info("Logging setup complete. Logging to %s", LOG_FILENAME↵
        )
130     return logger
```

5.0.3 Backend Routing and Task Set-up

Instead of traditional REST routes for core functionality, the backend's routing is managed by Celery's task queuing system. The `ProjectRouter` class defines custom routing logic to direct tasks to specific queues based on their purpose or the data they need to access. This is a crucial component of the distributed architecture, enabling efficient workload distribution and resource management across different Celery workers.

Below is the complete code for the `celery_app.py` file, which sets up the Celery application and its routing logic:

```
1  import config
2  from celery import Celery, signals
3  from backend import db_pool_setup
4
5  if not config.CELERY_BROKER_URL:
6      raise ValueError("An error has occurred. CELERY_BROKER_URL not ↵
        found in config")
7
8  class ProjectRouter:
9      def route_for_task(self, task, args=None, kwargs=None):
10         # 1. Dynamic routing for our main query task.
```

```
11         if task == 'backend.query_service.process_employee_query_task':
12             if kwargs and 'target_department' in kwargs:
13                 dept = kwargs['target_department']
14
15                 if dept in config.DEPARTMENT_CONFIG:
16                     collection_name = config.DEPARTMENT_CONFIG[dept]['←
17                         'qdrant_collection']
18                     return {
19                         'queue': f'query-{collection_name}'
20                     }
21
22             # 2. Static routing for embedding tasks.
23             if task.startswith('backend.text_embedding.'):
24                 return {'queue': 'embedding'}
25
26             # 3. Static routing for logging tasks.
27             if task.startswith('backend.chatlog_storage.'):
28                 return {'queue': 'logging'}
29
30         return None
31
32 # Celery App Initialization
33 celery_app = Celery(
34     config.APP_NAME,
35     broker=config.CELERY_BROKER_URL,
36     backend=config.CELERY_RESULT_BACKEND,
37     include=[
38         'backend.text_embedding',
39         'backend.query_service',
40         'backend.chatlog_storage'
41     ]
42 )
43 celery_app.conf.update(
44     task_serializer='json',
```

```
45     accept_content=['json'],
46     result_serializer='json',
47     timezone='Asia/Ho_Chi_Minh',
48     enable_utc=True,
49     task_acks_late=True,
50     worker_prefetch_multiplier=1,
51     task_routes = (ProjectRouter(),)
52 )
53
54 import logging
55 logger = logging.getLogger(config.APP_NAME)
56
57 @signals.worker_process_init.connect
58 def initialize_worker_process(**kwargs):
59     logger.info("Initializing Celery worker process...")
60     try:
61         db_pool_setup.initialize_pool()
62         logger.info("Database pool initialized for worker process.")
63     except Exception as e:
64         logger.error(f"Error initializing database pool in worker: {e↵↵
65                     }", exc_info=True)
66
67 @signals.worker_process_shutdown.connect
68 def shutdown_worker_process(**kwargs):
69     logger.info("Shutting down Celery worker process")
70     try:
71         db_pool_setup.close_pool()
72     except Exception as e:
73         logger.error(f"Error closing database pool in worker: {e}", ↵↵
74                     exc_info=True)
```

5.1 Database and Vector Index Implementation

5.1.1 Database Set-up

The system connects to a PostgreSQL database using `psycopg2` and a connection pool to manage concurrent access efficiently. The `initialize_pool()` function, protected by a `threading.Lock`, ensures the pool is set up safely and only once. This pool is initialized when a Celery worker process starts and is gracefully closed upon shutdown.

```

1  import psycopg2
2  import threading
3  import config
4  from contextlib import contextmanager
5  from psycopg2 import pool as psycopg_pool
6
7  import logging
8  logger = logging.getLogger(config.APP_NAME)
9
10 db_pool = None
11 pool_init_lock = threading.Lock()
12
13 def initialize_pool():
14     global db_pool
15
16     with pool_init_lock:
17         if db_pool:
18             logger.debug("Database pool already initialized")
19             return db_pool
20
21         logger.info("Initializing PostgreSQL connection pool")
22
23         try:

```

```
24         conn_args = {
25             "dbname": config.POSTGRES_DBNAME,
26             "user": config.POSTGRES_USER,
27             "password": config.POSTGRES_PASSWORD,
28             "host": config.POSTGRES_HOST,
29             "port": config.POSTGRES_PORT
30         }
31
32         db_pool = psycopg_pool.SimpleConnectionPool(
33             minconn=config.POSTGRES_DB_MIN_CONN,
34             maxconn=config.POSTGRES_DB_MAX_CONN,
35             **conn_args
36         )
37
38
39         logger.info(
40             f"Successfully initialized PostgreSQL pool "
41             f"(min: {config.POSTGRES_DB_MIN_CONN}, max: {config.↵
42                 POSTGRES_DB_MAX_CONN})"
43         )
44
45         return db_pool
46
47     except Exception as e:
48         db_pool = None
49         logger.exception("Failed to initialize connection pool")
50         raise RuntimeError("PostgreSQL connection pool ↵
51             initialization failed") from e
52
53
54 def close_pool():
55     global db_pool
56
57     if not db_pool:
58         logger.warning("Attempted to close PostgreSQL connection pool↵
59             , but it was not initialized")
60
61     return
```

```

57
58     logger.info("Closing PostgreSQL connection pool")
59
60     try:
61         db_pool.closeall()
62         logger.info("Successfully closed PostgreSQL connection pool")
63
64     except Exception as e:
65         logger.exception(f"Exception while closing PostgreSQL ↵
66                             connection pool: {e}")
67
68     finally:
69         db_pool = None
70
71 def get_pool():
72     global db_pool
73     if not db_pool:
74         logger.debug("Connection pool not initialized, calling ↵
75                             initialize_pool()")
76         pool = initialize_pool()
77         if not pool:
78             raise ConnectionError("Database pool could not be ↵
79                                     initialized")
80         return pool
81     return db_pool
82
83 # Context Manager using global db_pool
84 @contextmanager
85 def db_connection():
86     conn = None
87     pool = get_pool()
88
89     try:
90         conn = pool.getconn()
91         if not conn:
92             raise ConnectionError("Failed to acquire a database ↵

```

```
        connection from the pool")
90
91     logger.debug(f"Acquired connection {id(conn)} from pool")
92
93     yield conn
94
95     conn.commit()
96     logger.debug(f"Committed transaction on connection {id(conn)}↵
        ")
97
98     except Exception as e:
99         if conn:
100             try:
101                 conn.rollback()
102                 logger.warning(f"Rolled back connection {id(conn)} ↵
                    due to an exception")
103
104             except Exception as rollback_err:
105                 logger.exception(f"Rollback failed on connection {id(↵
                    conn)}: {rollback_err}")
106
107     logger.exception("Error during DB transaction")
108     raise
109
110 finally:
111     if conn:
112         if db_pool:
113             try:
114                 pool.putconn(conn)
115                 logger.debug(f"Returned connection {id(conn)} to ↵
                    pool")
116
117             except Exception as e:
118                 logger.exception(f"Error returning connection to ↵
                    pool: {e}")
119
```

```

120         try:
121             conn.close()
122             logger.warning(f"Force-closed connection {id(conn)} after pool failure")
123
124         except Exception as e2:
125             logger.exception(f"Failed to force-close connection {id(conn)}: {e2}")
126
127     else:
128         logger.warning(
129             f"Pool object (or global db_pool) became None in finally. "
130             f"Attempting to close conn id={id(conn)}."
131         )
132         try:
133             conn.close()
134
135         except Exception as e_close:
136             logger.exception(f"Failed to force-close connection id={id(conn)} when pool was None: {e_close}")

```

5.1.2 Data Schema in PostgreSQL and Qdrant

The data is structured across two main storage systems: PostgreSQL for raw data and Qdrant for vectorized data. In PostgreSQL, employee data is stored in separate tables for each department (e.g., `employees_engineering`), as mapped in `config.DEPARTMENT_CONFIG`. Each employee record includes fields such as `employee_id`, `full_name`, `job_title`, `start_date`, `skills`, `project_history`, `email`, and `location`. A crucial field, `embedding_status`, tracks whether the record has been processed for vectorization, which is key for data consistency.

When an employee's profile is vectorized, the resulting vector is stored in a

Qdrant collection along with a structured `payload`. This payload acts as the schema for the vector index, containing key-value pairs like `"text"`, `"employee_id"`, `"full_name"`, `"job_title"`, and other attributes. This dual storage of the raw text and structured metadata is essential for enabling both semantic search and filtered search.

The core logic for fetching data from PostgreSQL, generating embeddings, and preparing the Qdrant payload is contained in the `text_embedding.py` file, shown below:

```
1 import logging
2 import uuid
3 from enum import Enum
4 from celery import chain
5 from celery.exceptions import MaxRetriesExceededError
6 from qdrant_client import QdrantClient, models
7 from sentence_transformers import SentenceTransformer
8
9 import config
10 from backend import db_pool_setup
11
12 logger = logging.getLogger(config.APP_NAME)
13
14 # Celery App Import
15 try:
16     from celery_app import celery_app
17     logger.info("text_embedding.py: celery_app imported successfully.↵")
18 except ImportError as e:
19     logger.exception("text_embedding.py: Could not import celery_app.↵%s", e)
20     raise
21
22 # Initialize Sentence Transformer Model
23 try:
```

```

24     embedding_model = SentenceTransformer(config.↵
        LOCAL_EMBEDDING_MODEL_NAME)
25     logger.info(f"Successfully loaded sentence-transformer model: {↵
        config.LOCAL_EMBEDDING_MODEL_NAME}")
26 except Exception as e:
27     logger.critical(f"Failed to load sentence-transformer model: {e}"↵
        , exc_info=True)
28     embedding_model = None
29
30
31 # State Management
32 class ProcessingStatus(Enum):
33     PENDING = 'pending'
34     COMPLETED = 'completed'
35     FAILED = 'failed'
36
37 def launch_embedding_tasks():
38     logger.info("Starting Multi-Department Embedding Dispatch")
39     total_tasks_launched = 0
40
41     for dept_name, dept_config in config.DEPARTMENT_CONFIG.items():
42         table_name = dept_config.get('postgres_table')
43         collection_name = dept_config.get('qdrant_collection')
44
45         if not table_name or not collection_name:
46             logger.warning(f"Skipping department '{dept_name}' due to↵
                missing configuration.")
47             continue
48
49         all_pending_ids = fetch_pending_ids_for_department(table_name↵
            )
50
51         if not all_pending_ids:
52             logger.info(f"No pending employees found for department ↵
                '{dept_name}'.")
53             continue

```

```

54
55     batch_size = getattr(config, 'QDRANT_UPSERT_BATCH_SIZE', 90)
56
57     for i in range(0, len(all_pending_ids), batch_size):
58         batch_ids = all_pending_ids[i : i + batch_size]
59
60         try:
61             task_chain = chain(
62                 prepare_employee_vectors_task.s(batch_ids, ←
63                     table_name, dept_name),
64                 upsert_vectors_task.s(collection_name, table_name←
65                     )
66             )
67             task_chain.apply_async(queue="embedding")
68             total_tasks_launched += 1
69             logger.info(f"Dispatched task chain for {dept_name} ←
70                 with {len(batch_ids)} IDs.")
71         except Exception as e:
72             logger.exception(f"Failed to dispatch Celery task ←
73                 chain for {dept_name}: {e}")
74
75     logger.info(f"Finished Dispatching. Total tasks launched: {←
76         total_tasks_launched}")
77     return {"status": "dispatch_complete", "total_tasks_launched": ←
78         total_tasks_launched}
79
80 def fetch_pending_ids_for_department(table_name: str) -> list:
81     logger.info(f"Fetching pending IDs from table: '{table_name}'")
82     try:
83         with db_pool_setup.db_connection() as conn:
84             with conn.cursor() as cur:
85                 query = f"SELECT employee_id FROM {table_name} WHERE ←
86                     embedding_status = %s ORDER BY employee_id"
87                 cur.execute(query, (ProcessingStatus.PENDING.value,))
88                 return [row[0] for row in cur.fetchall()]
89     except Exception as e:

```



```

83         logger.exception(f"Error fetching pending IDs from {↵
            table_name}: {e}")
84         return []
85
86 def fetch_employees_by_ids(employee_ids: list, table_name: str) -> ↵
    list:
87     if not employee_ids:
88         return []
89     try:
90         with db_pool_setup.db_connection() as conn:
91             with conn.cursor() as cur:
92                 query = f"""
93                     SELECT
94                         employee_id, full_name, job_title, start_date↵
95                         ,
96                         skills, project_history, email, location
97                     FROM {table_name}
98                     WHERE employee_id = ANY(%s)
99                     """
100                 cur.execute(query, (employee_ids,))
101                 return [dict(zip([col.name for col in cur.description↵
                    ], row)) for row in cur.fetchall()]
102     except Exception as e:
103         logger.exception(f"Error fetching employee data from {↵
            table_name}: {e}")
104         return []
105
106 def update_employee_status_in_db(employee_ids: list, table_name: str,↵
    status: ProcessingStatus):
107     if not employee_ids:
108         return
109     logger.debug(f"Updating status to '{status.value}' for {len(↵
        employee_ids)} IDs in table '{table_name}'")
110     try:
111         with db_pool_setup.db_connection() as conn:
112             with conn.cursor() as cur:

```

```

112         query = f"UPDATE {table_name} SET embedding_status = ↵
           %s WHERE employee_id = ANY(%s)"
113         cur.execute(query, (status.value, employee_ids))
114     except Exception as e:
115         logger.error(f"DB status update to '{status.value}' failed ↵
           for table '{table_name}': {e}", exc_info=True)
116
117 @celery_app.task(bind=True, max_retries=3, default_retry_delay=60, ↵
    acks_late=True)
118 def prepare_employee_vectors_task(self, employee_ids_batch: list, ↵
    table_name: str, dept_name: str):
119     logger.info(f"[PREPARE TASK START] Preparing {len(↵
        employee_ids_batch)} employees from '{table_name}' ({dept_name↵
        }).")
120
121     if not embedding_model:
122         logger.error("Embedding model is not available. Task cannot ↵
            proceed.")
123         raise RuntimeError("SentenceTransformer model not loaded.")
124
125     employees = fetch_employees_by_ids(employee_ids_batch, table_name↵
        )
126     if not employees:
127         logger.warning(f"No employee data found for IDs in '{↵
            table_name}'. Task will end.")
128         return {"status": "no_data", "points": []}
129
130     texts_to_embed = []
131     employee_map = {}
132     for emp in employees:
133         text = (
134             f"Name: {emp.get('full_name', '')}. "
135             f>Title: {emp.get('job_title', '')}. "
136             f"Department: {dept_name}. "
137             f"Location: {emp.get('location', '')}. "
138             f"Contact Email: {emp.get('email', '')}. "

```

```

139         f"Start Date: {str(emp.get('start_date', ''))}. "
140         f"Skills: {emp.get('skills', '')}. "
141         f"Projects: {emp.get('project_history', '')}"
142     )
143     texts_to_embed.append(text)
144     employee_map[text] = emp
145
146     try:
147         embeddings = embedding_model.encode(texts_to_embed, ↵
148             show_progress_bar=False).tolist()
149
150     except Exception as e:
151         logger.warning(f"Task for {table_name} failed on embedding ↵
152             attempt {self.request.retries + 1}. Error: {e}")
153
154         try:
155             raise self.retry(exc=e, countdown=60)
156         except MaxRetriesExceededError:
157             logger.error(f"Task for {table_name} has failed ↵
158                 permanently. Marking records as FAILED.")
159             update_employee_status_in_db(employee_ids_batch, ↵
160                 table_name, ProcessingStatus.FAILED)
161             raise e
162
163     points_to_upsert = []
164     namespace = uuid.UUID('f2b4e448-86d6-4c28-8b1b-5e6e5e0c8b4e')
165
166     for i, text in enumerate(texts_to_embed):
167         emp = employee_map[text]
168         employee_id_str = str(emp['employee_id'])
169         point_id = str(uuid.uuid5(namespace, employee_id_str))
170
171         points_to_upsert.append({
172             "id": point_id,
173             "vector": embeddings[i],
174             "payload": {
175                 "text": text,

```

```

171         "employee_id": employee_id_str,
172         "full_name": emp.get('full_name', ''),
173         "job_title": emp.get('job_title', ''),
174         "department": dept_name,
175         "location": emp.get('location', ''),
176         "email": emp.get('email', ''),
177         "start_date": str(emp.get('start_date', '')),
178         "client_satisfaction": int(emp.get('←
            client_satisfaction', 0)),
179         "skills": emp.get('skills', ''),
180         "project_history": emp.get('project_history', '')
181     }
182 })
183
184     logger.info(f"[PREPARE TASK SUCCESS] Prepared {len(←
        points_to_upsert)} points for '{table_name}'.")
185     return {"status": "prepared", "points": points_to_upsert}
186
187
188 @celery_app.task(bind=True, max_retries=3, acks_late=True)
189 def upsert_vectors_task(self, vectors_payload: dict, collection_name:←
    str, table_name: str):
190     points = vectors_payload.get("points")
191     if not points:
192         logger.warning(f"Upsert task skipped for collection '{←
            collection_name}': No points received from prepare task.")
193         return {"status": "skipped"}
194
195     logger.info(f"[UPSERT TASK START] Upserting {len(points)} points ←
        to Qdrant collection: '{collection_name}'")
196
197     try:
198         qdrant_client = QdrantClient(host=config.QDRANT_HOST, port=←
            config.QDRANT_PORT)
199
200     try:

```

```

201         qdrant_client.get_collection(collection_name=↵
                collection_name)
202     except Exception:
203         logger.info(f"Collection '{collection_name}' not found. ↵
                Creating it now.")
204     qdrant_client.create_collection(
205         collection_name=collection_name,
206         vectors_config=models.VectorParams(
207             size=config.LOCAL_EMBEDDING_MODEL_DIMENSION,
208             distance=models.Distance.COSINE
209         ),
210     )
211
212     qdrant_client.create_payload_index(
213         collection_name=collection_name,
214         field_name="employee_id",
215         field_schema=models.TextIndexParams(
216             tokenizer=[models.TokenizerType.WORD, models.↵
                TokenizerType.WHITESPACE],
217             min_token_len=2
218         )
219     )
220
221     qdrant_client.create_payload_index(
222         collection_name=collection_name,
223         field_name="full_name",
224         field_schema=models.TextIndexParams(
225             tokenizer=[models.TokenizerType.WORD, models.↵
                TokenizerType.WHITESPACE],
226             min_token_len=2
227         )
228     )
229
230     qdrant_client.create_payload_index(
231         collection_name=collection_name,
232         field_name="job_title",

```

```
233         field_schema=models.TextIndexParams(  
234             tokenizer=[models.TokenizerType.WORD, models.TokenizerType.WHITESPACE],  
235             min_token_len=2  
236         )  
237     )  
238  
239     qdrant_client.create_payload_index(  
240         collection_name=collection_name,  
241         field_name="department",  
242         field_schema=models.TextIndexParams(  
243             tokenizer=[models.TokenizerType.WORD, models.TokenizerType.WHITESPACE],  
244             min_token_len=2  
245         )  
246     )  
247  
248     qdrant_client.create_payload_index(  
249         collection_name=collection_name,  
250         field_name="location",  
251         field_schema=models.TextIndexParams(  
252             tokenizer=[models.TokenizerType.WORD, models.TokenizerType.WHITESPACE],  
253             min_token_len=2  
254         )  
255     )  
256  
257     qdrant_client.create_payload_index(  
258         collection_name=collection_name,  
259         field_name="email",  
260         field_schema=models.TextIndexParams(  
261             tokenizer=[models.TokenizerType.WORD, models.TokenizerType.WHITESPACE],  
262             min_token_len=2  
263         )  
264     )
```

```

265
266         qdrant_client.create_payload_index(
267             collection_name=collection_name,
268             field_name="start_date",
269             field_schema=models.TextIndexParams(
270                 tokenizer=[models.TokenizerType.WORD, models.TokenizerType.WHITESPACE],
271                 min_token_len=2
272             )
273         )
274
275         qdrant_client.create_payload_index(
276             collection_name=collection_name,
277             field_name="client_satisfaction",
278             field_schema=models.PayloadIndexParams(
279                 index_type=models.PayloadIndexType.INTEGER
280             )
281         )
282
283         qdrant_client.create_payload_index(
284             collection_name=collection_name,
285             field_name="project_history",
286             field_schema=models.TextIndexParams(
287                 tokenizer=[models.TokenizerType.WORD, models.TokenizerType.WHITESPACE],
288                 min_token_len=2
289             )
290         )
291
292     point_structs = [models.PointStruct(**point) for point in points]
293
294     qdrant_client.upsert(
295         collection_name=collection_name,
296         points=point_structs,
297         wait=True

```

```

298         )
299
300         upserted_ids = [int(p.payload["employee_id"]) for p in ↵
                        point_structs]
301         update_employee_status_in_db(upserted_ids, table_name, ↵
                        ProcessingStatus.COMPLETED)
302
303         logger.info(f"[UPSERT TASK SUCCESS] Successfully upserted {↵
                        len(points)} points to '{collection_name}'.")
304         return {"status": "success", "count": len(points), "↵
                        collection_name": collection_name}
305
306     except Exception as e:
307         logger.error(f"[UPSERT TASK RETRY] Upsert failed for ↵
                        collection '{collection_name}': {e}", exc_info=True)
308     try:
309         raise self.retry(exc=e)
310     except MaxRetriesExceededError:
311         logger.error(f"Upsert task for {collection_name} has ↵
                        failed permanently. Marking records as FAILED.")
312     if points:
313         failed_ids = [int(p["payload"]["employee_id"]) for p ↵
                        in points]
314         update_employee_status_in_db(failed_ids, table_name, ↵
                        ProcessingStatus.FAILED)
315         raise e

```

5.2 Back-end Implementation

5.2.1 Core Functionality Set-up

The backend's core logic is implemented through a set of Celery tasks and helper functions. The `prompt_templates.py` file defines the system and user prompts used to instruct the LLM for both answer generation and filter extraction. These

prompts ensure the LLM behaves as an expert HR assistant and returns structured data when needed.

Below is the complete content of the `prompt_templates.py` file:

```

1 def get_employee_qa_prompt(context: str, query: str) -> tuple[str, str]:
2     system_prompt = (
3         "You are an expert HR assistant. Your task is to answer
4         questions about employees "
5         "based on the internal profile data provided as context.
6         Synthesize information from different profiles if needed.
7         "
8         "If you list employees, provide their full name and job title.
9         Be friendly and professional."
10    )
11    user_prompt = (
12        f"""Based ONLY on the context containing employee profiles
13        below, provide a clear and concise answer to the user's
14        query.
15        If the context does not contain the answer, you MUST state
16        that you could not find the information. Do not use
17        outside knowledge.
18
19        ## Context ##
20        {context}
21
22        ## User Query ##
23        "{query}"
24        """
25    )
26    return system_prompt, user_prompt
27
28 def get_filter_extraction_prompt(query: str) -> tuple[str, str]:
29     system_prompt = """

```

```
23     You are an expert at analyzing user queries and extracting ↵
        structured search filters.
24     Your output MUST be a single, valid JSON object and nothing else.
25     The valid keys for the JSON object are: "department" (string), "↵
        job_title" (string), "min_satisfaction" (integer), "↵
        employee_id" (string), "full_name" (string), "location" (↵
        string), "email" (string), "start_date" (string), "skills" (↵
        string), and "project_history" (string).
26     If the user's query does not contain any information for a filter↵
        , do not include the key in the JSON.
27     If no filters are found at all, return an empty JSON object {}.
28     ""
29
30     user_prompt = f"Analyze the following user query and extract the ↵
        filters:\n\n'{query}'"
31
32     return system_prompt, user_prompt
```

The core logic for filter extraction and Qdrant querying is handled in the `context_layer.py` module. It uses an LLM to extract structured filters from the query and then applies them as a `query_filter` during the vector search, improving retrieval precision.

Below is the complete content of the `context_layer.py` file:

```
1  import logging
2  import config
3  import re
4  import json
5  from qdrant_client import QdrantClient
6  from qdrant_client.http import models
7  from sentence_transformers import SentenceTransformer
8  from openai import OpenAI
9  from tenacity import retry, stop_after_attempt, wait_fixed
10
```

```
11 from backend.prompt_templates import get_filter_extraction_prompt
12
13 logger = logging.getLogger(config.APP_NAME)
14
15 # Initialize Sentence Transformer Model
16 try:
17     embedding_model = SentenceTransformer(config.↵
18         LOCAL_EMBEDDING_MODEL_NAME)
19     logger.info(f"Successfully loaded sentence-transformer model in ↵
20         context_layer: {config.LOCAL_EMBEDDING_MODEL_NAME}")
21 except Exception as e:
22     logger.critical(f"Failed to load sentence-transformer model in ↵
23         context_layer: {e}", exc_info=True)
24     embedding_model = None
25
26 relevance_threshold = 0.70
27
28 # LLM Helper for Filter Extraction
29 @retry(stop=stop_after_attempt(2), wait=wait_fixed(1))
30 def get_filters_from_llm(query: str) -> dict:
31     """
32     Uses an LLM to analyze the user's query and extract structured ↵
33     filters as a JSON object.
34     """
35     try:
36         client = OpenAI(api_key=config.MISTRAL_API_KEY, base_url="↵
37             https://api.mistral.ai/v1")
38
39         system_prompt, user_prompt = get_filter_extraction_prompt(↵
40             query)
41
42         response = client.chat.completions.create(
43             model=config.MISTRAL_MODEL, #
44             messages=[
45                 {"role": "system", "content": system_prompt},
46                 {"role": "user", "content": user_prompt}
```

```
41         ],
42         temperature=0.0,
43         response_format={"type": "json_object"}
44     )
45
46     response_text = response.choices[0].message.content
47     return json.loads(response_text)
48
49 except Exception as e:
50     logger.error(f"LLM filter extraction failed: {e}", exc_info=True)
51     return {}
52
53 # Qdrant Search Function
54 def query_qdrant(query: str, client: QdrantClient, collection_name: str,
55                 top_k=5, query_filter=None) -> list[dict]:
56     """
57     Queries a Qdrant collection using a locally generated embedding
58     and an optional filter.
59     """
60     if not embedding_model:
61         logger.error("Embedding model is not available. Cannot query Qdrant.")
62         return []
63
64     try:
65         query_embedding = embedding_model.encode(query, show_progress_bar=False).tolist()
66
67         search_responses = client.search(
68             collection_name=collection_name,
69             query_vector=query_embedding,
70             query_filter=query_filter,
71             limit=top_k,
72             with_payload=True,
73             score_threshold=relevance_threshold
```

```

72         )
73
74         return [{"text": match.payload["text"], "score": match.score}↵
                for match in search_responses]
75
76     except Exception as e:
77         logger.error(f"Qdrant query failed in collection '{↵
                collection_name}': {e}", exc_info=True)
78         return []
79
80 # Context Retrieval Handling
81 def get_context_for_query(query: str, client: QdrantClient, ↵
    collection_name: str) -> str:
82     """
83     Orchestrates the process of analyzing a query, building filters, ↵
        and retrieving context.
84     """
85     extracted_filters = get_filters_from_llm(query)
86     logger.info(f"Extracted filters from LLM for collection '{↵
        collection_name}': {extracted_filters}")
87
88     search_filter_conditions = []
89
90     if "min_satisfaction" in extracted_filters and isinstance(↵
        extracted_filters["min_satisfaction"], int):
91         search_filter_conditions.append(
92             models.FieldCondition(key="client_satisfaction", range=↵
                models.Range(gte=extracted_filters["min_satisfaction"↵
                    ]))
93         )
94
95     if "department" in extracted_filters and isinstance(↵
        extracted_filters["department"], str):
96         search_filter_conditions.append(
97             models.FieldCondition(key="department", match=models.↵
                MatchValue(value=extracted_filters["department"]))

```

```
98         )
99
100     if "employee_id" in extracted_filters and isinstance(↵
        extracted_filters["employee_id"], str):
101         search_filter_conditions.append(
102             models.FieldCondition(key="employee_id", match=models.↵
                MatchValue(value=extracted_filters["employee_id"]))
103         )
104
105     if "job_title" in extracted_filters and isinstance(↵
        extracted_filters["job_title"], str):
106         search_filter_conditions.append(
107             models.FieldCondition(key="job_title", match=models.↵
                MatchText(text=extracted_filters["job_title"]))
108         )
109
110     if "full_name" in extracted_filters and isinstance(↵
        extracted_filters["full_name"], str):
111         search_filter_conditions.append(
112             models.FieldCondition(key="full_name", match=models.↵
                MatchText(text=extracted_filters["full_name"]))
113         )
114
115     if "location" in extracted_filters and isinstance(↵
        extracted_filters["location"], str):
116         search_filter_conditions.append(
117             models.FieldCondition(key="location", match=models.↵
                MatchText(text=extracted_filters["location"]))
118         )
119
120     if "email" in extracted_filters and isinstance(extracted_filters[↵
        "email"], str):
121         search_filter_conditions.append(
122             models.FieldCondition(key="email", match=models.MatchText↵
                (text=extracted_filters["email"]))
123         )
```

```
124
125     if "start_date" in extracted_filters and isinstance(↵
        extracted_filters["start_date"], str):
126         search_filter_conditions.append(
127             models.FieldCondition(key="start_date", match=models.↵
                MatchText(text=extracted_filters["start_date"]))
128         )
129
130     if "skills" in extracted_filters and isinstance(extracted_filters↵
        ["skills"], str):
131         search_filter_conditions.append(
132             models.FieldCondition(key="skills", match=models.↵
                MatchText(text=extracted_filters["skills"]))
133         )
134
135     if "project_history" in extracted_filters and isinstance(↵
        extracted_filters["project_history"], str):
136         search_filter_conditions.append(
137             models.FieldCondition(key="project_history", match=models↵
                .MatchText(text=extracted_filters["project_history"]))
138         )
139
140     final_filter = models.Filter(must=search_filter_conditions) if ↵
        search_filter_conditions else None
141
142     qdrant_results = query_qdrant(query, client, collection_name, ↵
        query_filter=final_filter)
143
144     if not qdrant_results:
145         logger.warning("No results returned from Qdrant for query: '%↵
            s' in collection '%s'", query, collection_name)
146         return ""
147
148     relevant_passages = [res["text"] for res in qdrant_results]
149
150     if not relevant_passages:
```

```

151         logger.warning(
152             f"No results in collection '{collection_name}' met the ↵
                relevance threshold of {relevance_threshold}. "
153             f"Discarding context."
154         )
155         return ""
156
157     logger.info(
158         f"Retrieved {len(relevant_passages)} relevant passages from ↵
                collection '{collection_name}' for the context."
159     )
160     return "\n---\n".join(relevant_passages)

```

5.2.2 Functionality Implementation

The core backend functionalities are executed asynchronously by Celery workers, as defined in the `query_service.py` and `chatlog_storage.py` files. The `process_employee_query_task` is the workhorse for retrieving context from each department's Qdrant collection, while `store_batch_chat_logs_task` handles the asynchronous archival of conversation logs.

Below is the complete content of the `query_service.py` file, which defines the query processing task:

```

1  import logging
2  import config
3  from qdrant_client import QdrantClient
4  from backend import context_layer
5
6  logger = logging.getLogger(config.APP_NAME)
7
8  # Celery App Import
9  try:
10     from celery_app import celery_app

```



```
11     logger.info("celery_app has been successfully imported from ↵
        Celery")
12
13 except ImportError as e:
14     logger.exception("An exception has occurred when importing ↵
        celery_app instance")
15     raise
16
17 @celery_app.task(bind=True, max_retries=3, acks_late=True)
18 def process_employee_query_task(self, query: str, target_department: ↵
    str):
19     task_id = self.request.id
20     logger.info(f"[CONTEXT TASK START] Task ID: {task_id}. Dept: '{↵
        target_department}', Query: '{query}'")
21
22     try:
23         client = QdrantClient(host=config.QDRANT_HOST, port=config.↵
            QDRANT_PORT)
24
25         dept_info = config.DEPARTMENT_CONFIG.get(target_department)
26         if not dept_info:
27             logger.warning(f"No config found for department '{↵
                target_department}'. Skipping.")
28             return {"status": "skipped_no_config", "context": None, "↵
                source": target_department}
29
30         collection_name = dept_info['qdrant_collection']
31
32         try:
33             client.get_collection(collection_name=collection_name)
34         except Exception:
35             logger.warning(f"Collection '{collection_name}' not found↵
                . Skipping department '{target_department}'.")
36             return {"status": "skipped_no_collection", "context": ↵
                None, "source": target_department}
37
```

```

38         context_from_dept = context_layer.get_context_for_query(query↵
        , client, collection_name)
39
40     if context_from_dept:
41         logger.info(f"[CONTEXT TASK SUCCESS] Task ID: {task_id}. ↵
        Found context in '{target_department}'.")
42         return {"status": "success", "context": context_from_dept↵
        , "source": target_department}
43     else:
44         logger.info(f"[CONTEXT TASK SUCCESS] Task ID: {task_id}. ↵
        No relevant context found in '{target_department}'.")
45         return {"status": "no_context", "context": None, "source"↵
        : target_department}
46
47     except Exception as e:
48         logger.error(f"Unhandled exception in context task {task_id} ↵
        for department {target_department}: {e}", exc_info=True)
49         return {"status": "failed", "context": None, "source": ↵
        target_department, "error": str(e)}

```

Below is the complete content of the `chatlog_storage.py` file, which manages the chat log buffering and storage tasks:

```

1  import boto3
2  import json
3  import gzip
4  import uuid
5  import config
6  import redis
7  import io
8  import time
9  from botocore.exceptions import ClientError
10 from functools import lru_cache
11 import logging
12

```

```

13 logger = logging.getLogger(config.APP_NAME)
14
15 try:
16     from celery_app import celery_app
17     logger.info("celery_app imported successfully in chatlog_storage.↵
        ")
18 except ImportError:
19     logger.error("Could not import celery_app in chatlog_storage.")
20     celery_app = None ↵
                                   # ↵
        Allow file to be imported without Celery for non-worker ↵
        processes
21
22 # AWS and Redis config
23 bucket_name = config.AWS_BUCKET_NAME
24 try:
25     log_buffer_threshold = int(config.LOG_BUFFER_THRESHOLD)
26     log_buffer_ttl_seconds = int(config.LOG_BUFFER_TTL_SECONDS)
27 except (TypeError, ValueError):
28     log_buffer_threshold = 100
29     log_buffer_ttl_seconds = 3600
30
31 # Redis Initialization
32 redis_client = None
33 try:
34     redis_client = redis.StrictRedis(
35         host=config.REDIS_HOST, port=config.REDIS_PORT, db=config.↵
            REDIS_DB, decode_responses=False
36     )
37     redis_client.ping()
38 except Exception as e:
39     logger.error(f"Redis connection failed: {e}. Chatlog buffering ↵
        will likely fail.", exc_info=True)
40
41 @lru_cache(maxsize=1)
42 def get_s3_client():

```

```
43     try:
44         return boto3.client('s3')
45     except Exception as e:
46         logger.error(f"Failed to create S3 client: {e}", exc_info=↵
47             True)
48         return None
49
50 def get_utc_iso_timestamp():
51     return time.strftime("%Y-%m-%dT%H:%M:%S", time.gmtime()) + "Z"
52
53 def generate_s3_key(chat_id):
54     date_folder = time.strftime("%Y-%m-%d", time.gmtime())
55     uid = uuid.uuid4().hex[:8]
56     return f"chatlogs/{date_folder}/{chat_id}/batch_{int(time.time())↵
57         }_{uid}.json.gz"
58
59 def compress_json_payload(log_entries):
60     try:
61         json_bytes = json.dumps(log_entries).encode('utf-8')
62         out = io.BytesIO()
63         with gzip.GzipFile(fileobj=out, mode="wb") as f:
64             f.write(json_bytes)
65         return out.getvalue()
66     except Exception as e:
67         logger.error(f"Error during JSON compression: {e}", exc_info=↵
68             =True)
69         return None
70
71 if celery_app:
72     @celery_app.task(bind=True, max_retries=3, default_retry_delay=↵
73         =60, acks_late=True)
74     def store_batch_chat_logs_task(self, chat_id):
75         if not redis_client:
76             logger.error(f"Redis client not available. Cannot store ↵
77                 logs for chat_id: {chat_id}")
78             return {"status": "failed_no_redis", "chat_id": chat_id}
```

```
74
75     redis_key = f"{config.REDIS_LOG_LIST_KEY_PREFIX}{chat_id}"
76
77     try:
78         log_entries_bytes = redis_client.lrange(redis_key, 0, -1)
79         if not log_entries_bytes: return {"status": "↵
            no_logs_found", "chat_id": chat_id}
80
81         log_entries = [json.loads(entry.decode('utf-8')) for ↵
            entry in log_entries_bytes]
82         if not log_entries:
83             try: redis_client.delete(redis_key)
84             except Exception as del_e: logger.error(f"Redis error↵
                on cleanup: {del_e}")
85             return {"status": "decoding_failed", "chat_id": ↵
                chat_id}
86
87         compressed_payload = compress_json_payload(log_entries)
88         if not compressed_payload: raise ValueError("Compression ↵
            failed")
89
90         s3_key = generate_s3_key(chat_id)
91         s3 = get_s3_client()
92         if not s3: raise ConnectionError("S3 client unavailable."↵
            )
93
94         s3.put_object(Bucket=bucket_name, Key=s3_key, Body=↵
            compressed_payload, ContentType='application/json', ↵
            ContentEncoding='gzip')
95
96         try: redis_client.delete(redis_key)
97         except Exception as del_e: logger.error(f"Redis key ↵
            deletion failed after S3 upload: {del_e}")
98
99         return {"status": "success", "s3_path": f"s3://{↵
            bucket_name}/{s3_key}", "count": len(log_entries)}
```

```

100         except Exception as e:
101             logger.error(f"Error in store_batch_chat_logs_task for {↵
                chat_id}: {e}", exc_info=True)
102             raise self.retry(exc=e)
103
104 def buffer_chat_log(chat_id, user, message):
105     if not redis_client:
106         logger.warning(f"Redis not available. Cannot buffer log for ↵
            chat_id: {chat_id}")
107         return
108
109     log_entry = {'timestamp': get_utc_iso_timestamp(), 'user': user, ↵
        'message': message}
110     key = f"{config.REDIS_LOG_LIST_KEY_PREFIX}{chat_id}"
111
112     try:
113         entry_bytes = json.dumps(log_entry).encode('utf-8')
114         pipeline = redis_client.pipeline()
115         pipeline.rpush(key, entry_bytes)
116         pipeline.expire(key, log_buffer_ttl_seconds)
117         results = pipeline.execute()
118         current_length = results[0]
119
120         if current_length >= log_buffer_threshold and celery_app:
121             logger.info(f"Log buffer threshold reached for {chat_id}. ↵
                Triggering S3 storage task.")
122             store_batch_chat_logs_task.delay(chat_id)
123     except Exception as e:
124         logger.error(f"Error buffering chat log for {chat_id}: {e}", ↵
            exc_info=True)

```

The system also includes an administrative script for triggering manual tasks, as shown in `main.py` below:

```

1 import argparse

```

```
2 import logging
3 import sys
4 import os
5
6 project_root = os.path.abspath(os.path.join(os.path.dirname(__file__)↵
    , '..'))
7 if project_root not in sys.path:
8     sys.path.insert(0, project_root)
9
10 import config
11 from backend import text_embedding
12 from backend import db_pool_setup
13
14 try:
15     if hasattr(config, 'setup_logging'):
16         config.setup_logging()
17     logger = logging.getLogger(config.APP_NAME)
18 except Exception:
19     logging.basicConfig(level=logging.INFO, format='%(asctime)s - %(↵
        levelname)s - %(message)s')
20     logger = logging.getLogger(__name__)
21     logger.warning("Resorted to basic logging configuration.")
22
23 def run_embedding_generation():
24     logger.info("Attempting to launch embedding generation tasks for ↵
        all departments...")
25
26     db_pool_setup.initialize_pool()
27
28     try:
29         text_embedding.launch_embedding_tasks()
30         logger.info("Embedding generation tasks successfully ↵
            dispatched. Check Celery worker logs for progress.")
31
32     except Exception as e:
33         logger.error(f"Failed to launch embedding generation tasks: {↵
```

```
        e}", exc_info=True)
34
35     finally:
36         logger.info("Closing database pool...")
37         db_pool_setup.close_pool()
38         logger.info("Database pool closed.")
39
40 def main():
41     parser = argparse.ArgumentParser(description="Corporate Directory↵
        Admin CLI")
42     parser.add_argument(
43         "action",
44         choices=["generate-embeddings"],
45         help="The administrative action to perform."
46     )
47     args = parser.parse_args()
48
49     if args.action == "generate-embeddings":
50         run_embedding_generation()
51     else:
52         logger.error(f"Unknown action: {args.action}")
53         parser.print_help()
54
55 if __name__ == "__main__":
56     main()
```

5.3 Front-end Implementation

The frontend of the system is a conversational interface built using the Chainlit framework. It handles the user's interaction flow, displays task progress, and manages communication with the backend. Key functionalities include a user-friendly chat window, a live task list, and the ability to select different LLM models for generating answers.

The complete code for the Chainlit application logic is provided in `chainlit_app.py` below:

```
1 import sys
2 import os
3 import logging
4 import time
5
6 import chainlit as cl
7 from chainlit.element import TaskList, Task
8 from chainlit.message import Message
9 from celery import group
10 from celery.result import GroupResult, AsyncResult
11
12 project_root = None
13 try:
14     project_root = os.path.abspath(os.path.join(os.path.dirname(__file__), '..'))
15     if project_root not in sys.path:
16         sys.path.insert(0, project_root)
17 except Exception as e:
18     logging.basicConfig(level=logging.ERROR)
19     logging.error(f"CRITICAL: Failed to set project root path: {e}", exc_info=True)
20     raise
21
22 # Logging Configuration
23 try:
24     import config
25     logger = logging.getLogger(config.APP_NAME)
26     if hasattr(config, 'setup_logging') and callable(config.setup_logging):
27         config.setup_logging()
28     logger.info(f"Logging configured. Project root '{project_root}' in sys.path.")
```

```

29     else:
30         logger.warning(f"config.setup_logging() not found/callable.")
31 except Exception as e:
32     logger = logging.getLogger(__name__)
33     logging.basicConfig(level=logging.INFO)
34     logger.error(f"Error during config/logging setup: {e}. Using ↵
        basic logger.", exc_info=True)
35
36 from backend import prompt_templates
37 from openai import OpenAI
38 from tenacity import retry, stop_after_attempt, wait_fixed
39
40 # Import Celery
41 try:
42     from celery_app import celery_app
43     from backend.query_service import process_employee_query_task
44     from backend.chatlog_storage import buffer_chat_log
45     logger.info("Successfully imported application modules for ↵
        IntelChat.")
46 except ImportError as e:
47     logger.error("ImportError while loading application modules.", ↵
        exc_info=True)
48     raise
49
50 def initialize_selected_llm(model_choice):
51     safe_model_choice = model_choice.strip()
52     try:
53         if safe_model_choice == config.MISTRAL_MODEL_CHOICE.strip():
54             return OpenAI(api_key=config.MISTRAL_API_KEY, base_url="↵
                https://api.mistral.ai/v1")
55
56         elif safe_model_choice == config.OPENAI_MODEL_CHOICE.strip():
57             return OpenAI(api_key=config.OPENAI_API_KEY)
58
59         elif safe_model_choice == config.DEEPSEEK_MODEL_CHOICE.strip↵
            ():

```

```

60         return OpenAI(base_url=config.DEEPSEEK_BASEURL, api_key=↵
            config.DEEPSEEK_API_KEY)
61
62     else:
63         raise ValueError(f"Unsupported model: '{safe_model_choice↵
            }")
64     except Exception as e:
65         logger.critical(f"LLM Client initialization failed: {e}", ↵
            exc_info=True)
66         return None
67
68 @retry(stop=stop_after_attempt(3), wait=wait_fixed(2))
69 def call_llm_chat(client, system_prompt, user_prompt, model_name: str↵
    ):
70     try:
71         response = client.chat.completions.create(
72             model=model_name,
73             messages=[
74                 {"role": "system", "content": system_prompt},
75                 {"role": "user", "content": user_prompt}]
76         )
77         return response.choices[0].message.content
78     except Exception as e:
79         logger.error(f"LLM API call failed: {e}", exc_info=True)
80         raise
81
82 async def await_group_task_result(result_obj: GroupResult, ↵
    task_list_ui: cl.element.TaskList, timeout: int = 20):
83     """
84     Polls a Celery GroupResult object until it's ready, or until a ↵
        timeout is reached.
85     If a timeout occurs, it retrieves and returns results from only ↵
        the completed tasks.
86     """
87     task_list_ui.tasks[0].status = cl.TaskStatus.RUNNING
88     await task_list_ui.send()

```

```

89
90     start_time = time.time()
91
92     while not result_obj.ready():
93         if time.time() - start_time > timeout:
94             logger.warning(
95                 f"Task group {result_obj.id} timed out after {timeout}↵
96                     } seconds. "
97                     "Returning available results from completed tasks."
98             )
99             break
100
101     await cl.sleep(1)
102
103     retrieved_results = []
104     if result_obj.children:
105         for subtask_result in result_obj.children:
106             # Check if the individual subtask has a result (i.e., it'↵
107                 s ready)
108             if isinstance(subtask_result, AsyncResult) and ↵
109                 subtask_result.ready():
110                 try:
111                     result = subtask_result.get(timeout=0)
112                     retrieved_results.append(result)
113
114                 except Exception as e:
115                     logger.error(f"Failed to retrieve result from a ↵
116                         completed subtask: {e}", exc_info=True)
117
118                 elif isinstance(subtask_result, dict):
119                     retrieved_results.append(subtask_result)
120
121 if result_obj.ready() and result_obj.successful():
122     task_list_ui.tasks[0].status = cl.TaskStatus.DONE
123 else:
124     task_list_ui.tasks[0].status = cl.TaskStatus.FAILED
125     task_list_ui.tasks[0].description = "Some tasks failed or ↵

```

```
        timed out."
121     await task_list_ui.send()
122
123     return retrieved_results
124
125 @cl.on_chat_start
126 async def start_chat():
127     cl.user_session.set("chat_id", f"chat_{int(time.time())}")
128     cl.user_session.set("llm_choice", config.OPENAI_MODEL_CHOICE)
129     await cl.Message(
130         content="Welcome to IntelChat - the company's employee Q&A ↩️
131             chatbot. What can I help you with?",
132         author="IntelChat"
133     ).send()
134
135     ## CORRECTED ACTION SYNTAX
136     await cl.Message(
137         content="You can change the AI model at any time.",
138         author="IntelChat",
139         actions=[
140             cl.Action(
141                 name="change_model",
142                 value="change_model",
143                 payload={"value": "change_model"},
144                 label="Change Model"
145             )
146         ]
147     ).send()
148
149 @cl.action_callback("change_model")
150 async def on_change_model(action: cl.Action):
151     model_actions = [
152         cl.Action(
153             name="llm_selected",
154             value=config.MISTRAL_MODEL_CHOICE,
```

```

155         label=f"                {config.MISTRAL_MODEL_CHOICE}"
156     ),
157     cl.Action(
158         name="llm_selected",
159         value=config.OPENAI_MODEL_CHOICE,
160         payload={"value": config.OPENAI_MODEL_CHOICE},
161         label=f"                {config.OPENAI_MODEL_CHOICE}"
162     ),
163     cl.Action(
164         name="llm_selected",
165         value=config.DEEPSEEK_MODEL_CHOICE,
166         payload={"value": config.DEEPSEEK_MODEL_CHOICE},
167         label=f"                {config.DEEPSEEK_MODEL_CHOICE}"
168     )
169 ]
170 await cl.Message(
171     content="Please select your preferred AI model:",
172     author="IntelChat",
173     actions=model_actions
174 ).send()
175
176 @cl.action_callback("llm_selected")
177 async def on_llm_selected(action: cl.Action):
178     chosen_llm = action.payload.get('value')
179     if chosen_llm:
180         cl.user_session.set("llm_choice", chosen_llm)
181         await cl.Message(
182             content=f"Model has been set to **{chosen_llm}**. You can↵
183                 now ask me questions about our employees.",
184             author="IntelChat"
185         ).send()
186
187 @cl.on_message
188 async def main_logic(message: Message):
189     query_text = message.content
190     chat_id = cl.user_session.get("chat_id")

```

```
190     llm_choice = cl.user_session.get("llm_choice")
191     buffer_chat_log(chat_id, "user", query_text)
192
193     task_list = TaskList(tasks=[
194         Task(title="Searching all departments for context...", status←
195             =cl.TaskStatus.RUNNING)
196     ])
197
198     await task_list.send()
199
200     try:
201         departments_to_query = config.DEPARTMENTS
202         task_group = group(
203             process_employee_query_task.s(query=query_text, ←
204                 target_department=dept)
205             for dept in departments_to_query
206         )
207
208         group_result = task_group.apply_async()
209
210         results = await await_group_task_result(group_result, ←
211             task_list)
212
213     except Exception as e:
214         logger.error(f"Failed to dispatch Celery group task: {e}", ←
215             exc_info=True)
216         await task_list.remove()
217         await cl.Message(content="Error: Could not connect to the ←
218             processing service.").send()
219
220     return
221
222     all_context_fragments = []
223     searched_sources = []
224
225     for res in results:
226         if isinstance(res, dict) and res.get('status') == 'success' ←
227             and res.get('context'):
228             all_context_fragments.append(f"Context from {res['source']←
```

```

        ']] Department\n{res['context']}]")
220     if isinstance(res, dict) and 'source' in res:
221         searched_sources.append(res['source'])
222
223     if not all_context_fragments:
224         await cl.Message(content="I'm sorry, I could not find any ↵
            relevant employee profiles for your query across any ↵
            department.", author="IntelChat").send()
225         await task_list.remove()
226         return
227
228     final_context = "\n\n---\n\n".join(all_context_fragments)
229     system_prompt, user_prompt = prompt_templates.↵
        get_employee_qa_prompt(final_context, query_text)
230
231     llm_client = initialize_selected_llm(llm_choice)
232
233     model_to_use = None
234     if llm_choice == config.MISTRAL_MODEL_CHOICE: model_to_use = ↵
        config.MISTRAL_MODEL
235     elif llm_choice == config.OPENAI_MODEL_CHOICE: model_to_use = ↵
        config.OPENAI_MODEL
236     elif llm_choice == config.DEEPSEEK_MODEL_CHOICE: model_to_use = ↵
        config.DEEPSEEK_MODEL
237
238     if llm_client and model_to_use:
239         answer = call_llm_chat(llm_client, system_prompt, user_prompt↵
            , model_to_use)
240     else:
241         answer = "Error: Could not initialize the selected AI model ↵
            to generate a final answer."
242
243     await cl.Message(content=answer, author=llm_choice).send()
244
245     if searched_sources:
246         source_str = ", ".join(sorted(list(set(searched_sources))))

```



```
247         await cl.Message(content=f"*(Searched in: {source_str})*", ↵
           author="IntelChat", indent=1).send()
248
249     buffer_chat_log(chat_id, "assistant", answer)
250     await task_list.remove()
```

5.4 Source Code

The complete source code for the project is available on GitHub through this link: https://github.com/MatrixNova/distribution_system-employee_Q-A. You can refer to the repository for detailed setup instructions and deployment guidelines.

Chapter 6

Conclusions

In summary, the development of our Employee Q&A System has been aimed at resolving the critical challenges of latency and responsiveness that are associated with distributed systems, especially when dealing with large, multi-source datasets. We have been able to implement a real-time information retrieval platform that provides a responsive and friendly user experience by utilizing a robust technology framework that includes Celery, Qdrant, Redis, PostgreSQL, and various LLMs.

To achieve real-time query processing, we utilised asynchronous task processing with Celery and Redis, which facilitates instantaneous messaging and provide a smooth and engaging user experience. These technologies ensure that messages are delivered and received in real-time, creating a dynamic and interactive chat environment.

In terms of security, we implemented JWT-based authentication to secure user sessions and ensure that only authorised users can access specific functionalities. For maintaining application stability and security, we employed thorough input validation, efficient error-handling strategies, and detailed logging. These practices prevent invalid data from entering the system, avoid potential application crashes, and safeguard against security vulnerabilities.

Additionally, the application features a modern UI/UX design that enhances user interaction and satisfaction. The UI/UX design focuses on ease of use, aesthetic appeal, and a seamless navigation experience, contributing to a positive overall user experience.

With regard to scalability and maintainability, we structured the application to ensure it can grow and adapt to future needs. This involves writing clean, modular

code and implementing best practices that allow for straightforward updates and feature additions.

To support the real-time chat functionality, security features, and scalability requirements, we built an efficient system architecture. This includes a well-defined client-server model where the client-side handles the user interface and real-time updates via WebSockets, while the server-side manages user authentication, message storage, and overall system logic. This separation of concerns ensures that the system is both performance and easy to maintain.

Overall, the Employee Q&A System is a comprehensive solution that combines several technical features with a user-centric design. The project effectively addresses the core requirements of real-time communication, ensuring users can engage in instantaneous messaging.

6.1 Development Potential

Going forward, the project has multiple ways in which it could still be expanded on from the iteration being shown here. An easy way of expanding is to simply integrate more choices of LLMs for processing and answer generation, as our system is already structured to support multiple LLM clients from providers like Mistral, OpenAI, and DeepSeek. The User Interface, while perfectly serviceable using Chainlit, could be adapted into something more customised should future needs or proprietary issues arise. As for the actual inner workings of the systems, there will always be space for optimization whether it be in the asynchronous function calls or how our Celery workers handle and deliver each query. For instance, fine-tuning the dynamic routing logic in ProjectRouter or optimising the polling mechanism in `await_group_task_result` are potential future improvements. But as of now, these are all uncertain plans and the program is for all intents and purposes able to be considered feature complete.