# Introduction

Our team proposes the development of **FilmFinder**, aimed at providing users with an efficient way to search for movies based on various criteria. The platform will utilize the information contained in the TMDB dataset to offer users a customizable movie search experience.

## 1 Application System

### 1a System Description:

The platform will serve as a catalog of movies and related information to facilitate users to search for movies using various criteria, including *Genre, Release Date, Language, Popularity*. A search option would be provided where users can enter their search criteria and receive appropriate movie recommendations. Additionally, users will have the option to view detailed information about each movie, including its title, poster, overview, rating, and other metadata.

### Limitations:

The analysis is limited to the data available in the database, which may not cover all the information. The system does not support real time data analysis, as the data set is updated periodically. It supports only basic analysis, it does not support complex analysis.

### 1b Database Necessity

A database is essential for the Movie Search and Discovery Platform for the following reasons:

- Storage : The application requires an organized repository to efficiently store and manage movie data and its attributes.
- Retrieval : A database facilitates efficient retrieval of information, allowing for quick and accurate search results.
- Data Integrity : A database management system ensures data integrity, uniformity and consistency.
- Scalability : A DBMS provides a decent scalability, to handle a growing volume of movie data and user requests effectively.

### 1c Justification for Migrating to AWS:

Migrating the database to AWS offers several benefits :

- Cost savings: We can reduce infrastructure costs and pay only for the resources they use, leading to significant cost savings.
- Scalability and flexibility: AWS provides scalable and flexible infrastructure, allowing us to easily adjust resources based on demand.
- Disaster recovery: We can create backup copies of the data across different regions, minimizing the risk of data loss in case of a disaster.
- Accessibility: The global accessibility allows users to access the database from anywhere with an internet connection.
- Agility: AWS enables quick deployment and scaling of new applications.

# 2 Solution Requirements

## 2a Requirements of the solution:

- Data Storage and Management: The system should have a very strong data storage to securely store and manage the dataset efficiently.
- Performance: The system should deliver efficient query processing.
- Scalability:  The system should be scalable to handle growing datasets and high loads by maintaining the performance.
- Security: Security measures should be implemented to protect the data.
- Highly available and Disaster Recovery: The system should ensure high availability and disaster recovery to reduce the downtime and data loss.
- Data source: The initial source of information for the application is the TMDB dataset.
- Data preprocessing: The raw data needs to be cleaned and normalized in order to ensure query-ability.
- User Interface: The system should provide a user friendly interface to interact with the system and to perform analysis.

## 2b System's Functionalities and Limitations:

### Functionalities:

Users can explore the TMDB dataset by browsing movie metadata, genres, ratings, cast and crew, languages, production companies and keywords. Users can execute custom queries to filter data on a particular criteria. As the application grows, the users can also perform various analytics tasks such as trend analysis,  popularity analysis, rating analysis, genre analysis to know more.

Here are a few functionalities and requirements:

**Search Functionality:** Users should be able to search for movies based on criteria such as genre, release date, language, and popularity.
**Filtering Options:** Provide users with filtering options to refine their search results, including genre selection, release date range, and language preferences.

**Detailed Movie Information:** Display detailed information about each movie, including title, poster, overview, rating, and other metadata.
**Feedback Mechanism:** Provide users with the ability to provide feedback on search results and suggest improvements to the platform.

Limitations:

**Data Availability -** The service may be limited by the availability and quality of data in the TMDB Analysis database. Inaccurate or incomplete data may lead to poor search results.
**Scalability -** As the user base / database size increases, we may encounter scalability challenges in terms of handling large numbers of service (data retrieval) requests.
**Performance -** The performance may be affected by service availability, load, and querying efficiency.

## 2c User Interaction:

Users will interact with the system through a web-based interface. The interface will feature search and filtering options prominently, allowing anyone, even *unregistered users*, to have *read only* access. Users can  access detailed information and ratings on individual movies.

For registered users, the system will provide features such as user accounts and *search and view history* to enhance the user experience. Feedback mechanisms, such as rating movies, will be available for registered users.

# 3 Conceptual Database Design

## 3a Database Requirements

For an application utilizing the TMDB dataset, the database requirements involve:

- **MySQL on RDS:** For scalable and reliable database management.
- **Schema Design:** For efficient storage of all the information on movies, collections, genres, etc., with appropriate keys for data integrity. (Refer ER Diagram in the next section)
- **Views:** To simplify complex queries and to enable restricted access to registered users.
- **Stored Procedures:** For common tasks like fetching recommendations or calculating statistics.
- **Triggers:** For data integrity checks, validation, or automation based on database events.
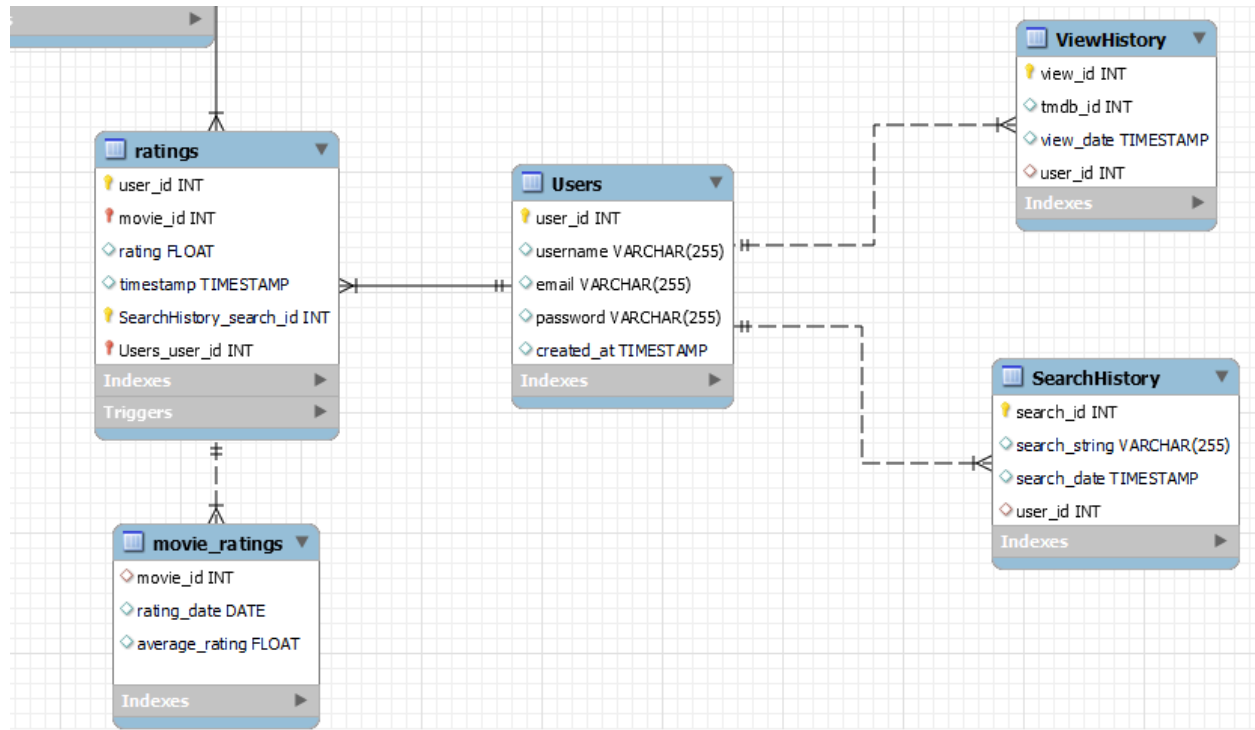
# 3b Entity Relationship Diagram

**casts**
- credit_id VARCHAR(255)
- cast_id INT
- id INT
- character_name VARCHAR(500)
- gender INT
- name VARCHAR(255)
- order_num INT
- profile_path VARCHAR(255)
- Indexes

**keywords**
- keywords_id VARCHAR(255)
- name VARCHAR(255)
- Indexes

**movie_genres**
- tmdb_id INT
- genres_id VARCHAR(255)
- Indexes

**movie_spoken_languages**
- tmdb_id INT
- spoken_languages_id VARCHAR(255)
- Indexes

**movie_casts**
- tmdb_id INT
- credit_id VARCHAR(255)
- Indexes

**genres**
- genres_id VARCHAR(255)
- name VARCHAR(255)
- Indexes

**movie_keywords**
- tmdb_id INT
- keywords_id VARCHAR(255)
- Indexes

**spoken_languages**
- spoken_languages_id VARCHAR(255)
- name VARCHAR(255)
- Indexes

**movie_metadata**
- tmdb_id INT
- title VARCHAR(255)
- adult TINYINT(1)
- budget BIGINT
- homepage VARCHAR(255)
- imdb_id INT
- original_language VARCHAR(255)
- original_title VARCHAR(255)
- overview VARCHAR(1000)
- popularity FLOAT
- poster_path VARCHAR(255)
- release_date VARCHAR(255)
- revenue BIGINT
- runtime INT
- status VARCHAR(255)
- tagline VARCHAR(500)
- video TINYINT(1)
- vote_average FLOAT
- vote_count INT
- Indexes
- Triggers

**ratings**
- user_id INT
- movie_id INT
- rating FLOAT
- timestamp TIMESTAMP
- Indexes
- Triggers

**links**
- movie_id INT
- imdb_id INT
- tmdb_id INT
- Indexes

**movie_production_countries**
- tmdb_id INT
- production_countries_id VARCHAR(255)
- Indexes

**production_countries**
- production_countries_id VARCHAR(255)
- name VARCHAR(255)
- Indexes

**belongs_to_collection**
- belongs_to_collection_id VARCHAR(255)
- name VARCHAR(255)
- poster_path VARCHAR(255)
- backdrop_path VARCHAR(255)
- Indexes

**movie_crews**
- tmdb_id INT
- credit_id VARCHAR(255)
- Indexes

**production_companies**
- production_companies_id VARCHAR(255)
- name VARCHAR(255)
- Indexes

**movie_production_companies**
- tmdb_id INT
- production_companies_id VARCHAR(255)
- Indexes

**movie_belongs_to_collection**
- tmdb_id INT
- belongs_to_collection_id VARCHAR(255)
- movie_metadata_tmdb_id INT
- movie_metadata_tmdb_id1 INT
- Indexes

**crews**
- credit_id VARCHAR(255)
- id INT
- department VARCHAR(255)
- gender INT
- name VARCHAR(255)
- job VARCHAR(255)
- profile_path VARCHAR(255)
- Indexes

The schema is centered around *movie_metadata* for the most part with tables like *genres* and *spoken_languages* having many-to-many relationships with *movie_metadata*.

Here, *belongs_to_collection* and *movie_metadata* is a one-to-many relationship, but we have a middle table *movie_belongs_to_collection* to avoid too many rows having NULL values in the collection column.

Additional tables:

To enable user-specific functionalities in the application, we have introduced 4 more tables with user data for managing their ratings, average ratings and search / view history.



# 3c Normalization process:

Normalization is a crucial process in database design which is used to organize data efficiently reducing the data redundancy. Normalization is done by breaking down large tables into smaller, related tables and establishing relationships between them. In the TMDB dataset, normalization can help organize movie-related data into structured tables, eliminating duplication and improving data management.

Normalization in TMDB dataset:

1. First Normal Form (1NF): Each table needs to have atomic values. For instance, in movies_metadata.csv, columns collections, genres, spoken languages, production companies and production countries don't have atomic values. Splitting the one big table of *movie_metadata* into separate tables, *movie_metadata, belongs_to_collection, genres, production_companies, production_countries* and *spoken_languages* with their links preserved in a different set of tables - *movie_genres, movie_belongs_to_collection, movie_production_companies, movie_production_countries* and *movie_spoken_languages* respectively to achieve 1NF atomicity.

2. Second Normal Form (2NF): Non-key attributes are fully functionally dependent on the primary key, that is every attribute depends on the entire primary key, not just part of it. For instance, in the movies_metadata table, attributes such as title, release_date depend on the entire primary key, which is *tmdb_id*.

3. Third Normal Form (3NF): No transitive dependencies between non-key attributes, that is non-key attribute doesn't depend on other non-key attributes. In the TMDB dataset, no non-key attribute depends on a non-key attribute.

Primary keys in the TMDB dataset represent unique identifiers for each movie or related entity. For example, in the *movie_metadata* table, the *tmdb_id* column serves as the primary key, uniquely identifying each Movie record.

Foreign keys are established to represent relationships between different tables. For example, in the *ratings* table, the *movie_id* column may serve as a foreign key referencing the *movie_id* column in the *links* table. This establishes a relationship between the ratings table and the links table, indicating which movie each rating entry is associated with.

In summary, normalization in the TMDB dataset involves structuring data to eliminate redundancy, ensure data integrity, and establish relationships between related entities. This is achieved by breaking down tables, ensuring atomic values, and establishing primary and foreign key relationships between tables.

# 4 Functional Analysis

## 4a Functional components of the application:

## Data preparation

Data preprocessing and preparation are essential steps in database design to ensure that the data is well-organized, consistent, and optimized for efficient storage and retrieval. In the case of the TMDB dataset, which contains information about movies and related entities, such as collections, genres, production companies, and spoken languages, several preprocessing steps are necessary. Initially, the data is read from CSV files, which contains discrepancies and redundancies. To address this, the data is transformed into a structured format, adhering to the principles of normalization. This involves breaking down large tables into smaller, related tables and establishing relationships between them.

For example, the one big table of movie metadata is split into separate tables, such as movie metadata, belongs_to_collection, genres, production companies, production countries, and spoken languages. Each table is designed to have atomic values and to eliminate transitive dependencies, ensuring that non-key attributes are fully functionally dependent on the primary key. Once the data is in a normalized and structured format, it is connected to a local database, such as MySQL, and populated into the appropriate

tables. This preprocessing and preparation process ensures that the data is well-organized, consistent, and ready for efficient storage and retrieval in the database system. All this is accomplished using Python and the corresponding code can be found [here](#).

Before we dive into discussing the functional DB components of the application, let us discuss the access distribution.

One section of the dataset has a variety of attributes centered around the movie_metadata table. The other section of the database consists of user and user related data including ratings and search / view history. Access control is achieved by many ways, one of which is by maintaining different users with different access privileges at the MySQL level.
- **Unregistered users** have access solely to the *movie_details* view, to search for movies based on specific criteria. A separate MySQL user ***viewro*** is created with appropriate read-only access to the *movie_details* view.
- **Registered users** are free to rate movies and have access to recommendations based on their usage history. Queries from registered users are executed from the MySQL user named ***ltdrw*** with limited read-write privileges. This user has read-only access to all the tables in tmdb and read-write access to *ratings, Users, ViewHistory* and *SearchHistory*
- **Administrators** have RW access to all tables and they access the tmdb database via the user ***admin***

The application can be divided into 4 main functional components:
1. **Search and Retrieval:** This component allows users to search for movies based on various criteria. It includes functionalities for search, filtering, and displaying movie details.
2. **Recommendation system based on feedback mechanism:** Allowing for analysis of search patterns, the logs can be used to personalize user experiences by recommending similar or trending movies, also utilizing the user ratings and recommendations.
3. **Admin Ops:** Administrators have access to manage movie data, user ratings, and system settings. They can add, edit, or remove movies, genres, languages, and other metadata.
4. **Logging and Monitoring:** The system logs user interactions, database queries, and system events for monitoring, analysis, and troubleshooting purposes.
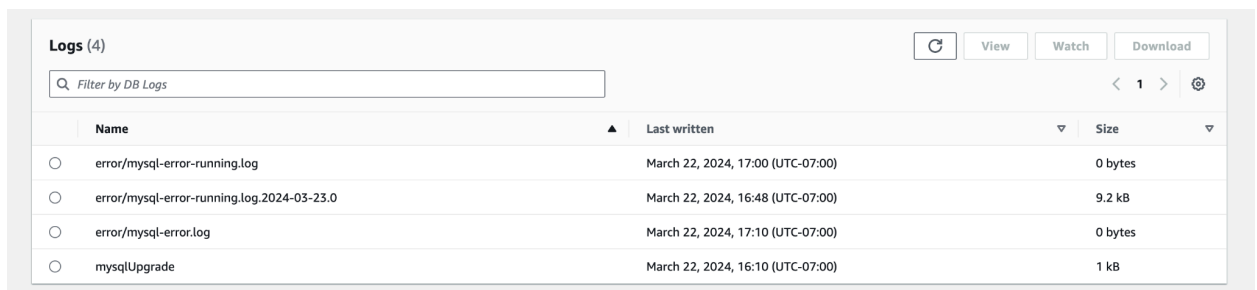
## 4b Database interactions for each Functional components:

1. <u>Search and Retrieval</u>

   This component allows users to search for movies based on various criteria. It includes functionalities for search, filtering, and displaying movie details.
   *This can be split into two sections:*
   a. *Open to all - Provides access only to **movie_details** view*
   b. *For registered users - Provides restricted access to all the tables*

## Table structure:

Users with the **open to all** view have permission only to the movie_details VIEW containing most of the information from movie_metadata.

The **registered user** view would provide querying access to all the tables.

For the full table structure, please refer to section 3b.

## Access privileges:

As stated above, we have 3 users at MySQL level - *viewro, ltdrw* and *admin*.
- Users with the **open to all** view would be able to query the database through the ***viewro*** user (MySQL) that has access only to movie_details VIEW.
- The **registered users** with MySQL user **ltdrw** would have read-access to all the tables including write access to a subset of them viz., ratings, Users, ViewHistory and SearchHistory.
- **Administrators** with MySQL user **admin** are the maintainers of the TMDB database, they have write access to all the tables in the schema.

*Please refer to the grant access queries [here](#) .*

## SQL code:

1. Find all comedy movies.
2. Get the cast and crew of Ariel movie.
3. Fetch the movie production company and country.
4. Fetch the collection that has more than 10 movies.
5. Top 10 movies based on revenue generated
6. Number of Movies in each genre
7. Top 5 movies based on Popularity score in each genre
8. Top 15 Keywords
9. Weighted Rating based on Number of Votes and Vote Average

## Stored Procedures:

1. MovieSearch()
2. CalculateMovieStatistics()
3. CountOfMovies_InYear()
4. display_top_collaborations()
5. filter_movie()

## Logging Mechanism:

Logging mechanisms for Administrators and to maintain Search History of registered users are described at the end of this section.

**2. Recommendation system based on feedback mechanism:**

Allowing for analysis of search patterns, the logs can be used to personalize user experiences by recommending similar or trending movies, also utilizing the user ratings and recommendations.

Table structure:

The tables - *ratings, ViewHistory* and *SearchHistory* serve to provide insights into user behavior on the platform, allowing for analysis of search patterns and the popularity of viewed movies. They can be used to personalize user experiences by recommending similar or trending movies based on search and viewing history.

Access privileges:

*Registered users (MySQL user - **ltdrw**) can rate movies and have access to their search / view history. This requires RW access to user-related tables. The users have to be registered to enjoy maximum benefits of the application and to receive movie recommendations.*

SQL code:

1. Get the average rating of movies and order by highly rated movies.

Triggers:

1. update_average_rating_trigger
2. update_average_rating_log
3. Prevent_negative_rating

Stored Procedures:

1. GenerateMovieRecommendations()
2. GetTopRatedMovies( )

3. Admin Ops -

Administrators have access to manage movie data, user ratings, and system settings. They can add, edit, or remove movies, genres, languages, and other metadata.

Table structure:

All the tables represented in the ER diagram are accessible to admin users. They are the maintainers of the TMDB database who are the only ones to have insert / update / delete permissions to all the Movie tables.

Access privileges:

*All other tables including movie_metada are RW-accessible only to administrators (MySQL user - **admin**).*

1. Introduce new collections
2. Update status after a movie release

1. cascade_delete_movie_records
2. enforce_budget_constraint
3. enforce_unique_title_year
4. check_genre_limit

1. check_if_belongs_to_collection()

4. Logging and Monitoring -

The system logs user interactions, database queries, and system events for monitoring, analysis, and troubleshooting purposes.

In-built:

MySQL has in-built logging mechanisms for various purposes, viz., Error Logs, General Query Logs, Binary Logs and Slow Query Logs. To get a log of all queries that manipulated the data, we use the binary logs. These logs are also used to backup and recover the database.



In RDS, there is a Logs and Events tab in which we will be able to access recent events on the DB server and the error logs as well.

We can also use Amazon CloudWatch Logs to monitor, store, and access your log files from RDS sources. CloudWatch Logs enables us to centralize the logs in a single, highly scalable service.

**Log exports**

Select the log types to publish to Amazon CloudWatch Logs

- ☑ Audit log
- ☑ Error log
- ☑ General log
- ☐ Slow query log

User-defined:

We can also establish a user-defined logging mechanism like how we implemented search and view history. This is more customized as per the requirement, user-friendly and also easily accessible. The "SearchHistory" table maintains the log of search queries, while the "ViewHistory" table records viewed movies corresponding to users. Both tables serve to provide insights into user behavior on the platform, allowing for analysis of search patterns and the popularity of viewed movies, thereby contributing to the recommendation system.

Table structure:

**SearchHistory** and **ViewHistory** *maintain a log of movies accessed by registered users. Users and movie_metadata tables also play a part in maintaining the history.*

Access privileges:

*Open to registered users with **ltdrw** permission. Only registered users have access to their Search and View history.*

SQL code:

1. Retrieve all search history records for a specific user
2. Retrieve all view history records for a specific user
3. Delete all search history records older than 6 months
4. Retrieve the top 10 users with the highest number of view history records

Triggers:

1. log_popularity_update
2. rating_deletion_audit

Stored Procedures:

1. LogSearchHistory()
2. LogViewHistory()

# 5 Database Migration to AWS

## 5a Step-by-Step guide:

**Source Database Setup:**
The dataset imported from multiple csv files was thoroughly examined to identify and rectify any inconsistencies, inaccuracies, or missing values, ensuring data cleanliness and integrity using python scripts.

Subsequently, the dataset was grouped into multiple data frames, each representing a distinct entity or aspect of the data as described in the ER diagram of TMDB. This partitioning facilitated the adherence to the principles of normalization, ensuring that data was organized efficiently and redundancies were minimized. Furthermore, by applying referential integrity principles the relationships between entities were accurately represented and maintained.

Finally, a database was created in the local instance of MySQL and all the preprocessed data was inserted into the tables. The data preprocessing scripts developed in python ensured that the data was structured optimally for seamless integration into MySQL tables, fostering a robust and reliable database foundation for the migration process.

Code Execution Results:

**Target Database Setup:**

A Relational Database Service (RDS) instance was created to facilitate the storage and management of data. A new user was created in the AWS Management Console with access to RDS service. The necessary parameters for the RDS instance, including the database engine type, instance class, storage type, and allocated storage, were configured based on project requirements. Security settings such as network access permissions, encryption options, and database credentials were established to ensure data security and integrity. Once all configurations were finalized, the RDS instance was provisioned, and the database engine was deployed. Create an empty TMDB database in the new AWS RDS instance.



**Migration Process Execution:**

From the MySQL workbench local instance, export the TMDB Database to a self contained .sql file using Server->Data Export option. While exporting the database select all the tables and

views. Also, select the options to dump stored procedures, dump triggers and include create schema.





From MySQL workbench connect to the AWS RDS instance created in the target db setup using the hostname, port, user and password. Verify the connection using the Test Connection option.

Using Server->Import option in MySQL RDS instance, import the dump file exported in the previous step.

After importing the dump file, refresh the schemas tab and verify that all the tables, stored procedures, triggers and views are available under the TMDB database.

## 5b Connectivity to AWS using Python:

Python MySQL Connector module was used to establish connectivity with AWS RDS instance by following these steps:

1. **Install MySQL Connector using pip:**

```
pip install mysql-connector-python
```

2. **Get RDS Endpoint and Database Credentials:** Necessary information to connect to AWS RDS instance like endpoint (hostname), port number, database username, and password were gathered from AWS RDS Dashboard.
3. **Connect to RDS Instance:** Using the MySQL Connector a connection was established with the RDS instance. Using the following code snippet:

```python
import mysql.connector
from mysql.connector import Error

def get_db_connection(host, user, password, database=None):
        try:
                if (database == None):
                        connection = mysql.connector.connect(host=host, user=user, password=password)
                else:
                        connection = mysql.connector.connect(host=host, user=user, password=password,
database=database)
                print(f"Sucessfully connected to AWS RDS as {user} user.")
                print("AWS RDS Connection Details:")
                print(f"Host: {connection.server_host}")
                print(f"Port: {connection.server_port}")
                print(f"User: {connection.user}")
                print(f"Database: {connection.database}")
                print_executed = True
        except Error as e:
                print(f"Error while connecting to AWS RDS as {user} user.", e)
                quit()
        return connection
```

4. **Verify Connectivity:** To ensure connectivity with AWS RDS instance, details of the connection were printed on the console and a "SELECT VERSION();" query was executed using embedded sql in python

```
● (test_env) soumiya@Sunaadas-MacBook-Air DATA225-Lab1 % python AWS/main.py  dblab1.cva0cgg2qoot.us-east-1.rds.amazonaws.com admin
  Sucessfully connected to AWS RDS as admin user.
  AWS RDS Connection Details:
  Host: dblab1.cva0cgg2qoot.us-east-1.rds.amazonaws.com
  Port: 3306
  User: admin
  Database: tmdb
  MySQL Server version: ('8.0.36',)
○ (test_env) soumiya@Sunaadas-MacBook-Air DATA225-Lab1 % ▊
```

## 5c Migration Validation

A SELECT query was executed using embedded sql in python to validate migration. Also, all the stored procedures and triggers in the TMDB database migrated to the AWS RDS instance were listed to ensure that the migration was successful.

Functions to execute the queries and print the results:

```python
def execute_query_with_params(host, user, password, database, query, param_values):
        try:
        cnx = get_db_connection(host,user,password,database)
        cursor = cnx.cursor()
        cursor.execute(query, param_values)
        rows = cursor.fetchall()
        column_names = [i[0] for i in cursor.description]
        df = pd.DataFrame(rows, columns=column_names)
        cursor.close()
        cnx.close()
        return df

        except mysql.connector.Error as err:
        print(f"Error: {err}")
        return None

def exexute_and_print_query_results(host,user,password, query, param_values, header):
        result_df = execute_query_with_params(host,user,password,constants.DATABASE_NAME, query,
param_values)

        if result_df is not None:
                print(header)
                print(result_df)
        else:
                print("Query execution failed.")
```

Sample query to fetch top 10 movies by revenue:

```python
def fetch_top_movie_by_revenue(host,user,password):
        query = """SELECT
                    original_title AS 'Movie Title',
                    FORMAT(budget / 1000000, 2) AS 'Budget in Millions',
                    FORMAT(revenue / 1000000, 2) AS 'Revenue in Millions',
                    FORMAT((revenue - budget) / 1000000, 2) AS 'Profit in Millions'
        FROM
                    movie_metadata
        ORDER BY revenue DESC
        LIMIT 10;
        """

        param_values = ()
        header = "Top 10 Movies based on revenue generated"
        exexute_and_print_query_results(host,user,password, query, param_values, header)
```

Function to list all the stored procedures in the database:

```python
def list_stored_procedures(host,user,password):
        query="SELECT routine_name FROM information_schema.routines WHERE routine_type =
'PROCEDURE' AND routine_schema = %s"
        # Execute the query
        param_values = ((constants.DATABASE_NAME,))
        header = "List of Stored Procedures"
        exexute_and_print_query_results(host,user,password, query, param_values, header)
```

Function to list all the triggers and their associated tables in the database:

```python
def list_triggers(host,user,password):
        query="SELECT trigger_name, event_object_table FROM information_schema.triggers WHERE
trigger_schema = %s"
        # Execute the query
        param_values = ((constants.DATABASE_NAME,))
        header = "List of Triggers"
        exexute_and_print_query_results(host,user,password, query, param_values, header)
```

Screenshot of the results:



# 6 SQL Performance

[Screenshot for the commands](#)

| S.NO | Command | DML/DDL | Local Instance- Time Taken in seconds | AWS RDS- Time Taken in seconds |
|------|---------|---------|----------------------------------------|--------------------------------|
| 1 | CREATE | DDL | 0.031 | 0.156 |
| 2 | DROP | DDL | 0.031 | 0.125 |
| 3 | INSERT | DML | 0.000 | 0.078 |

| 4 | SELECT | DML | 0.047 | 0.109 |
|---|--------|-----|-------|-------|

By analyzing the time taken, we see that the performance in local is faster than RDS for both DDL and DML commands. RDS is a little slow compared to Local MySQL because of the Network Latency. But AWS RDS has other advantages like scalability, manageability, and reliability.

RDS Performance Insights

**Performance Insights** expands on existing Amazon RDS monitoring features to illustrate and help you analyze your database performance. With the Performance Insights dashboard, you can visualize the database load on your Amazon RDS DB instance load and filter the load by waits, SQL statements, hosts, or users. RDS Performance Insights is not available on the small and micro instance sizes but it's an efficient way to monitor and enhance performance of the application.



# 7 Summary

In conclusion, the migration of the FilmFinder database to AWS RDS offers numerous benefits such as cost savings, scalability, disaster recovery, and accessibility. Despite slight differences in performance compared to a local MySQL instance, AWS RDS provides advantages in terms of manageability and reliability. With the implementation of Performance Insights, the application's performance can be further monitored and optimized. Overall, the migration to AWS RDS ensures a robust and efficient database management solution for the FilmFinder application.

# 8 References

GitHub URL :
 https://github.com/soumiyarao/DATA225-Lab1

Queries and result screenshots used in the application :
https://docs.google.com/document/d/1bF2yBIiPd2gXxN0zh5KZ2pEOvNcZubqzbdRTZV3UfFQ

Performance analysis screenshots :
https://acrobat.adobe.com/link/review?uri=urn%3Aaaid%3Ascds%3AUS%3Acb971d94-4787-3d2f-b987-6b7f5e01006d