

# CSE 4/560 Milestone 2

## FilmVerse

Kajol  
UB ID: kajol  
kajol@buffalo.edu

Mahimitra Chirala  
UB ID: mahimitr  
mahimitr@buffalo.edu

Sagar Gopalasetti  
UB ID: sambasai  
sambasai@buffalo.edu

### I. INTRODUCTION

#### A. Problem Statement

Finding a good movie to watch can be a tedious and time-consuming task. With so many movies available to choose from, users often spend a lot of time searching the internet for ratings and reviews. To address this issue, we propose a movie information system that provides users with a comprehensive collection of movie data, including ratings, genres, actors, release year, and other relevant information. This system will enable users to search for and discover movies that meet their specific needs and preferences quickly and easily. By leveraging the system's advanced search and filtering capabilities, users can find top-rated films, movies in their preferred genres, films in a specific language, or movies featuring a favorite actor. This system aims to simplify the movie selection process and enhance the overall viewing experience for users.

### II. TARGET USER

#### A. The User

The IMDb dataset caters to a diverse audience, including researchers, data analysts, and movie enthusiasts, interested in discovering and studying movie and TV data. The intended audience comprises university researchers specializing in media and communication, marketing experts analyzing customer behavior and preferences, and film studios examining market trends. The dataset provides a valuable resource for these audiences to explore and gain insights into trends, patterns, and connections in the film industry. The comprehensive data available could enhance research efforts, improve marketing strategies, and drive innovation and success in the industry.

#### B. The Admin

The management of the IMDb dataset requires skilled professionals such as data administrators or data engineers. These professionals are responsible for ensuring that the data is correctly structured, maintained, and preserved. They must conduct data quality checks, make regular backups, and continuously work on improving database speed to ensure efficient and accurate data retrieval. Effective management of the database ensures that users have access to up-to-date, reliable information, which is essential for researchers, analysts, and industry professionals to make informed decisions. By employing experienced data professionals, the dataset can be

continuously improved and maintained to provide users with the most comprehensive and accurate movie data available.

#### C. Real-life scenario description

Real-world applications of the IMDb dataset extend to market research companies that aim to examine consumer preferences and behavior toward films and TV series. The dataset provides valuable insights into the most popular genres, actors and actresses that appear most frequently in popular films, and production companies that have the most commercial success. This information can be utilized to make informed decisions about movie distribution and development, as well as marketing and advertising tactics for streaming services and movie studios. The IMDb dataset is a comprehensive and reliable resource that market research companies can utilize to gain valuable insights into the film industry and make data-driven decisions.

### III. DATABASE

There are six tables in data namely:

- 1) Akas table  
The relation provides alternate titles for movies and TV shows in different regions and languages, which can be useful for international marketing and distribution.
- 2) Basics table  
The Basics table contains fundamental information about movie and TV show titles such as title, year of release, genre, movie type, etc.
- 3) Crew Table  
This table contains information about the individuals who contribute to a film or TV show's production, such as directors, and writers.
- 4) Episodes Table  
The Episodes table contains information about episodes, including the number of seasons, and episode number of a TV show.
- 5) Ratings Table  
The Ratings table contains information about the movie's average rating and number of votes.

### IV. BCNF

All the tables are in BCNF hence we did not have to decompose them. For all the dependencies of the form  $X \rightarrow Y$ .  $X$  is a super key and there are no transitive dependencies in any of the tables hence, they are in BCNF. To explain this,

Identify applicable funding agency here. If none, delete this.

dependencies for each table are listed as follows.

For the akas table,

akas(title\_id, ordering, title, region, language, types, attributes, is\_original\_title ) The Functional dependencies are:

titleId, ordering → title  
 titleId, ordering → region  
 titleId, ordering → language  
 titleId, ordering → types  
 titleId, ordering → attributes  
 titleId, ordering → is\_original\_title

For the basics table,

basics (tconst, title\_type, primary\_title, original\_title, is\_adult, start\_year, end\_year, runtime\_minutes, genres ) The

Functional dependencies are:

tconst → title\_type  
 tconst → primary\_title  
 tconst → original\_title  
 tconst → is\_adult  
 tconst → start\_year  
 tconst → end\_year  
 tconst → runtime\_minutes  
 tconst → genres

For crew table,

crew (tconst, directors, writers) The Functional dependencies are:

tconst → directors  
 tconst → writers

For episode table,

episode (tconst, parent\_tconst, season\_number, episode\_number) The Functional dependencies are:

tconst → parent\_tconst  
 tconst → season\_number  
 tconst → episode\_number

For ratings table,

ratings (tconst, average\_rating, num\_votes) The Functional dependencies are:

tconst → average\_rating  
 tconst → num\_votes

All the above tables are in BCNF as, in all the FDs of the form  $X \rightarrow Y$ .

## V. ENTITY-RELATIONSHIP DIAGRAM

ER Diagram stands for Entity Relationship Diagram, It displays the relationship of entity sets stored in a database. In other words, ER diagrams help to explain the logical structure of databases.

From Figure 1, we can observe that there are six entities, namely basics, akas, crew, episode, and ratings. All the entities are related to the basics entity.

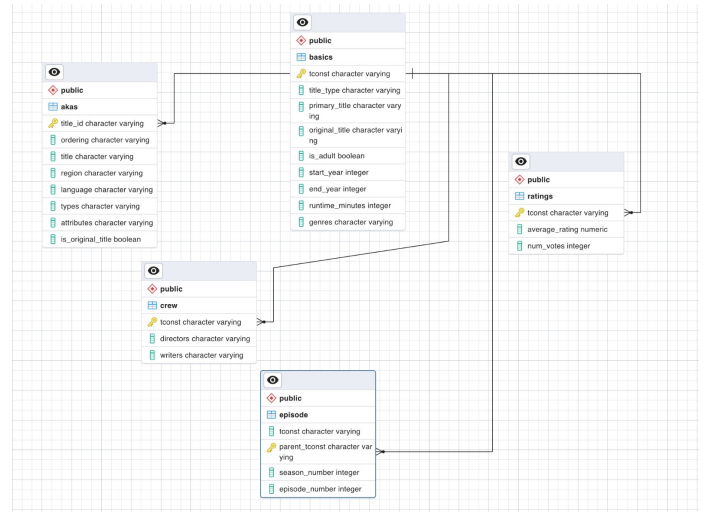


Fig. 1. ER-diagram.

## VI. CONSTRAINTS

The Constraints for our database are listed below as

For the akas table:

**Primary Key:** (title\_id, ordering), The combination of these two columns uniquely identifies a row in the table, which is necessary for data integrity and avoiding duplicate entries.

**Foreign Key:** The foreign key is title\_id, which refers to the tconst column in the basics table. This allows for linking alternative titles to their corresponding main titles.

For the basics table:

**Primary key:** tconst : The primary key is tconst. This is the unique identifier for each title and is necessary for data integrity and avoiding duplicate entries.

**Foreign key :** No foreign keys

For the crew table:

**Primary key:** tconst : The primary key is tconst. This is the unique identifier for each title and is necessary for data integrity and avoiding duplicate entries.

**Foreign key:** The foreign key is tconst, which refers to the tconst column in the basics table. This allows for linking crew information to the correct title.

For the episode table:

**Primary key:** tconst: The primary key is tconst. This is the unique identifier for each episode and is necessary for data integrity and avoiding duplicate entries.

**Foreign key:** The foreign key is parent\_const, which refers to the tconst column in the basics table. This allows

for linking episodes to their corresponding TV series.

For ratings table:

**Primary key:** The primary key is tconst. This is the unique identifier for each title and is necessary for data integrity and avoiding duplicate entries.

**Foreign key:** The foreign key is tconst, which refers to the tconst column in the basics table. This allows for linking rating information to the correct title.

## A. INSERTION

- 1) INSERT INTO public.basics (tconst, title\_type, primary\_title, original\_title, is\_adult, start\_year, end\_year, runtime\_minutes, genres) VALUES ('tt9999998', 'movie', 'New Movie', 'New Movie', false, 2023, NULL, 120, 'Action, Drama')

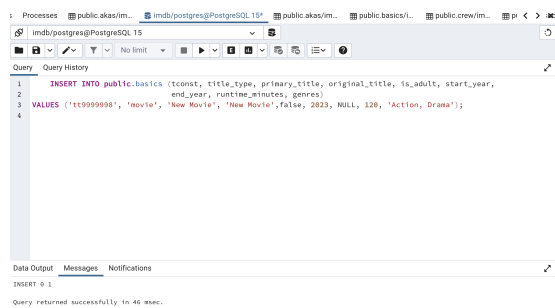


Fig. 2. Insert 1

- 2) INSERT INTO crew (tconst, directors, writers) VALUES ('tt9999998', 'nm9876543', 'nm1234567, nm2345678');

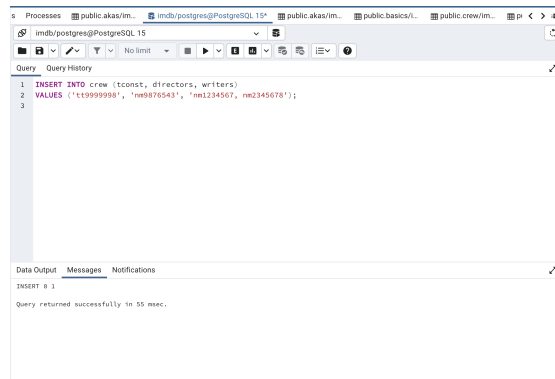


Fig. 3. Insert 2

## B. UPDATION

- 1) UPDATE akas SET title = 'The Enchanted Inn', region = 'US', types = 'working', is\_original\_title true WHERE title\_id = 'tt0000018';

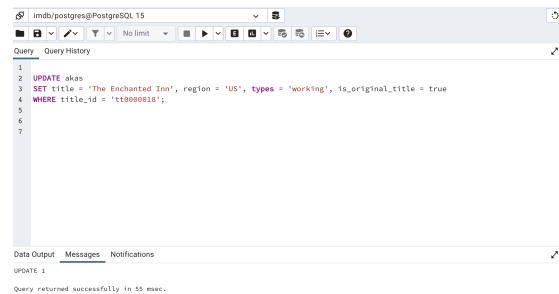


Fig. 4. Update 1

- 2) UPDATE basics SET title\_type = 'short', primary\_title = 'The Laughing Song', original\_title = 'Le rire', is\_adult = false, start\_year = 1896, end\_year = NULL, runtime\_minutes = 1, genres = 'Comedy' WHERE tconst = 'tt0000024';

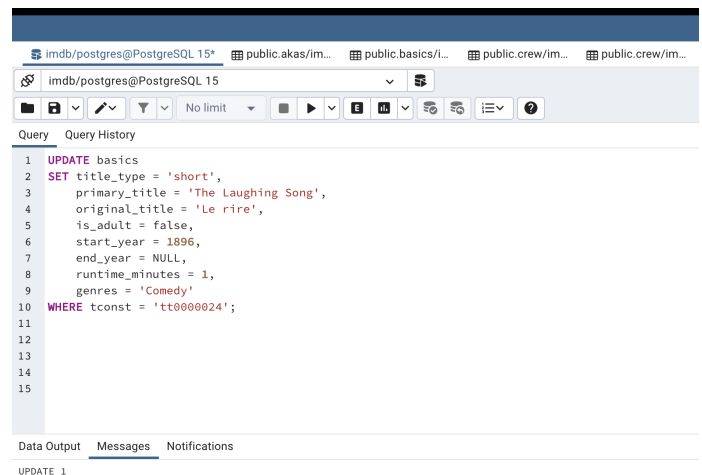


Fig. 5. Update 2

## C. DELETION

- 1) DELETE FROM akas WHERE title\_id = 'tt9999998' ;
- 2) DELETE FROM ratings WHERE tconst = 'tt9999998';

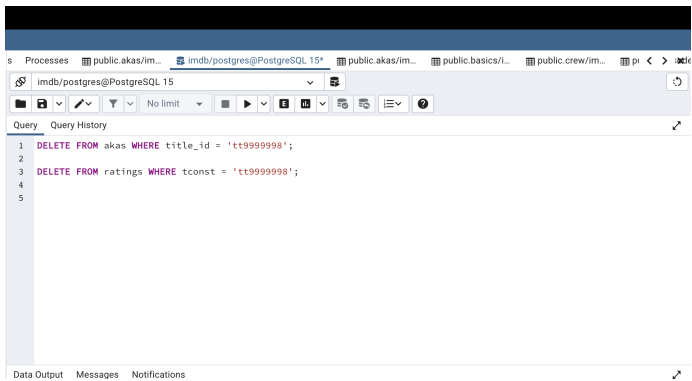


Fig. 6. Delete

#### D. SELECT QUERIES

Q1: Find the top 10 highest-rated movies.

Q2: Sort genres by the highest number of movies

Q3: Number of movies released in a year along with ratings

Q4: Episode and parent title information for a specific TV series, and sort by year and title

Q5: Give all comedy movies

tconst	primary_title	average_rating
tt000016	Barque sortant du port	8.4
tt000036	L'arroseur arrosé	8.3
tt000054	The Sign of the Cross	8.2
tt000030	Sea Bathing	8.1
tt000045	The Photographical Congress Arrives in Lyon	8.0
tt000024	The Laughing Song	7.9
tt000027	The Pillar of Fire	7.8
tt000057	Hawaiian Government Parade	7.7
tt000042	The Mystic Swing	7.6
tt000013	The Photographical Congress Arrives in Lyon	7.5

Fig. 7. Q1 : Top 10 highest rated movies

genre	num_movies
Short	32
Comedy	16
Documentary	14
Drama	8
Fantasy	8
Animation	6
Action	5
Sport	4
Romance	4
Horror	3
Drama	2
Crime	2
Short	1
short	1
History	1
Adventure	1

Fig. 8. Q2 : Genres sorted by highest number of movies

start_year	num_titles	avg_rating
1933	1	6.7000000000000000
1915	1	5.0000000000000000
1910	1	5.2000000000000000
1909	2	6.4500000000000000
1908	4	6.3000000000000000
1907	5	6.5400000000000000
1906	4	5.7500000000000000
1905	2	5.5500000000000000
1904	1	6.7000000000000000
1903	6	5.7166666666666667

Fig. 9. Q3 :Number of movies released in a year along with ratings

Fig 9. : This query joins the basics and ratings tables on their tconst columns and groups the results by the start\_year column of the basics table. It then calculates the count of titles and the average rating for each year using the COUNT() and AVG() aggregate functions, respectively. Finally, it orders the results by the start\_year column.

```

1 SELECT e.tconst AS id, 'episode' AS type, b.primary_title AS title, b.original_title AS original_title,
2 b.is_adult AS adult, b.start_year AS year, b.runtime_minutes AS runtime, b.genres AS genres
3 FROM public.episode e
4 JOIN public.basics b ON e.parent_tconst = b.tconst
5 WHERE e.parent_tconst = 'tt0000021'
6
7 UNION
8
9 SELECT b.tconst AS id, b.title_type AS type, b.primary_title AS title, b.original_title AS original_title,
10 b.is_adult AS adult, b.start_year AS year, b.runtime_minutes AS runtime, b.genres AS genres
11 FROM public.basics b
12 WHERE b.tconst = 'tt0000021'
13
14 ORDER BY year DESC, title ASC;
15

```

id	type	title	original_title	adult	year	runtime	genres
1	character varying	Short	Les forgerons	boolean	1895	1	Documentary
2	episode	Les forgerons	Les forgerons	false	1895	1	Documentary

Fig. 10. Q4 :Episode and parent title information for a specific TV series, and sort by year and title along with ratings

Fig 10. : This query retrieves information about all episodes and the parent title of the TV series with ID 'tt0000021' from the episodes and basics tables. It selects the relevant columns and renames them according to the specifications given earlier. It uses UNION to combine the results of two queries and sorts the result by year (descending) and title (ascending).

```

1 SELECT *
2 FROM basics
3 WHERE tconst IN (
4 SELECT DISTINCT e.parent_tconst
5 FROM episode e
6 JOIN basics b ON e.parent_tconst = b.tconst
7 WHERE b.genres LIKE '%Comedy%'
8 )
9 ORDER BY start_year DESC;

```

tconst	title_type	primary_title	original_title	is_adult	start_year	end_year	runtime_mir
1	Short	The Millinery Shop	The Millinery Shop	true	1909	[null]	[null]
2	Short	The Dancing Pig	The Dancing Pig	false	1907	[null]	[null]
3	Short	The Dancing Pig	Le cochoon danseur	true	1907	[null]	[null]
4	Short	Dream of a Rarebit Fiend	Dream of a Rarebit Fiend	false	1906	[null]	[null]
5	Short	Laughing Gas	Laughing Gas	false	1906	[null]	[null]
6	Short	The Mysterious Retort	The Mysterious Retort	false	1906	[null]	[null]
7	Short	The Porter's Tactics	Les réparateurs	true	1903	[null]	[null]
8	Short	The Messenger Boy's Mistake	The Messenger Boy's Mistake	false	1903	[null]	[null]
9	Short	The Comic Artist and His Model	The Comic Artist and His Model	false	1901	[null]	[null]

Fig. 11. Q5 :All comedy movies.

Fig 11. : This query will return all the TV series (parent titles) that have at least one episode and belong to the Comedy genre. The subquery first selects all the distinct parent\_tconst values from the episode table that match these criteria, and then the outer query returns all the corresponding rows from the basics table, sorted by start\_year in descending order

## VII. QUERY OPTIMISATION

In order to improve the performance of specific queries that may have taken longer to execute, indexing can be applied to them. For example, the execution time of a specific query can be reduced from 369ms to 142ms through indexing. For some queries, we followed different approaches to optimize them further. Refer to Fig. 12 and Fig. 13.

```

1 SELECT *
2 FROM public.ratings
3 WHERE tconst = 'tt0000091';
4

```

tconst	average_rating	num_votes
1	5.2	14

Total rows: 1 of 1 Query complete 00:00:00.369

Fig. 12. Before Indexing

```

1 SELECT *
2 FROM public.ratings
3 WHERE tconst = 'tt0000091';

```

tconst	average_rating	num_votes
1	5.7	1959

Total rows: 1 of 1 Query complete 00:00:00.142

Fig. 13. After Indexing

In this query, We replaced the IN clause with an EXISTS clause, which can be faster for large subqueries. We also removed the join with the basics table in the subquery since we only need to check if the parent\_tconst exists in the basics table. Refer to Fig. 14 and Fig. 15.

Query

```

1 SELECT *
2 FROM basics
3 WHERE tconst IN (
4   SELECT DISTINCT e.parent_tconst
5   FROM episode e
6   JOIN basics b ON e.parent_tconst = b.tconst
7   WHERE b.genres LIKE '%Comedy%'
8 )
9 ORDER BY start_year DESC;

```

Data Output

tconst	PKI	character varying	title_type	character varying	primary_title	character varying	original_title	character varying	is_adult	boolean	start_year	integer	end_year	integer	runtime_mnr	integer
1	tt000096		Short		The Millinery Shop		The Millinery Shop		true		1909					
2	tt000082		Short		The Dancing Pig		The Dancing Pig		false		1907					
3	tt000069		Short		The Dancing Pig		Le cochoon danseur		true		1907					
4	tt000072		Short		Dream of a Rarebit Fiend		Dream of a Rarebit Fiend		false		1906					
5	tt000078		Short		Laughing Gas		Laughing Gas		false		1906					
6	tt000041		Short		The Mysterious Retort		The Mysterious Retort		false		1906					
7	tt000028		Short		The Porter's Tactics		Les réparateurs		true		1903					
8	tt000095		Short		The Messenger Boy's Mistake		The Messenger Boy's Mistake		false		1903					
9	tt000080		Short		The Comic Artist and His Model		The Comic Artist and His Model		false		1901					

Total rows: 16 of 16 Query complete 00:00:00.205 Ln 9, Col 26

Fig. 14. Before Optimising

Query

```

1 SELECT *
2 FROM basics b
3 WHERE EXISTS (
4   SELECT 1
5   FROM episode e
6   WHERE e.parent_tconst = b.tconst
7 )
8 AND b.genres LIKE '%Comedy%'
9 ORDER BY b.start_year DESC;

```

Data Output

tconst	PKI	character varying	title_type	character varying	primary_title	character varying	original_title	character varying	is_adult	boolean	start_year	integer	end_year	integer	runtime_mnr	integer
1	tt000096		Short		The Millinery Shop		The Millinery Shop		true		1909					
2	tt000082		Short		The Dancing Pig		The Dancing Pig		false		1907					
3	tt000069		Short		The Dancing Pig		Le cochoon danseur		true		1907					
4	tt000072		Short		Laughing Gas		Laughing Gas		false		1906					
5	tt000078		Short		The Mysterious Retort		The Mysterious Retort		false		1906					
6	tt000041		Short		Dream of a Rarebit Fiend		Dream of a Rarebit Fiend		false		1906					
7	tt000028		Short		The Messenger Boy's Mistake		The Messenger Boy's Mistake		false		1903					
8	tt000095		Short		The Porter's Tactics		Les réparateurs		true		1903					
9	tt000080		Short		The Comic Artist and His Model		The Comic Artist and His Model		false		1901					

Total rows: 16 of 16 Query complete 00:00:00.082 Ln 9, Col 28

Fig. 15. After Optimising

In this query, We replaced the UNION operator with UNION ALL since we don't need to eliminate duplicates in this case. We also removed the unnecessary table aliases and used a single JOIN instead of two separate queries. Finally, I suggest creating an index on the parent\_tconst column of the episode table to speed up the join operation. Refer to Fig. 16 and Fig. 17.

Query

```

1 SELECT e.tconst AS id, 'episode' AS type, b.primary_title AS title, b.original_title AS original_title,
2 b.is_adult AS adult, b.start_year AS year, b.runtime_minutes AS runtime, b.genres AS genres
3 FROM public.episode e
4 JOIN public.basics b ON e.parent_tconst = b.tconst
5 WHERE e.parent_tconst = 'tt0000021'
6 UNION
7 SELECT b.tconst AS id, b.title_type AS type, b.primary_title AS title, b.original_title AS original_title,
8 b.is_adult AS adult, b.start_year AS year, b.runtime_minutes AS runtime, b.genres AS genres
9 FROM public.basics b
10 WHERE b.tconst = 'tt0000021'
11 ORDER BY year DESC, title ASC;

```

Data Output

id	type	title	original_title	adult	year	runtime	genres
1	tt0000021	Short	Les forgerons	Les forgerons	false	1895	1 Documentary
2	tt0001023	episode	Les forgerons	Les forgerons	false	1895	1 Documentary

Total rows: 2 of 2 Query complete 00:00:00.191 Ln 11, Col 31

Fig. 16. Before Optimising

Query

```

1 SELECT e.tconst AS id, 'episode' AS type, b.primary_title AS title, b.original_title AS original_title,
2 b.is_adult AS adult, b.start_year AS year, b.runtime_minutes AS runtime, b.genres AS genres
3 FROM public.episode e
4 JOIN public.basics b ON e.parent_tconst = b.tconst
5 WHERE e.parent_tconst = 'tt0000021'
6 UNION ALL
7 SELECT b.tconst AS id, b.title_type AS type, b.primary_title AS title, b.original_title AS original_title,
8 b.is_adult AS adult, b.start_year AS year, b.runtime_minutes AS runtime, b.genres AS genres
9 FROM public.basics b
10 WHERE b.tconst = 'tt0000021'
11 ORDER BY year DESC, title ASC;

```

Data Output

id	type	title	original_title	adult	year	runtime	genres
1	tt0001023	episode	Les forgerons	Les forgerons	false	1895	1 Documentary
2	tt0000021	Short	Les forgerons	Les forgerons	false	1895	1 Documentary

Total rows: 2 of 2 Query complete 00:00:00.074 Ln 12, Col 1

Fig. 17. After Optimising

## VIII. BONUS TASK

The user interface (UI) includes a home page that displays the top 10 movies at the current moment by default. Additionally, there is a dropdown menu that allows users to filter results by categories such as top 10 movies by different genres like comedy, short, etc. Below is the screenshot that displays our home page

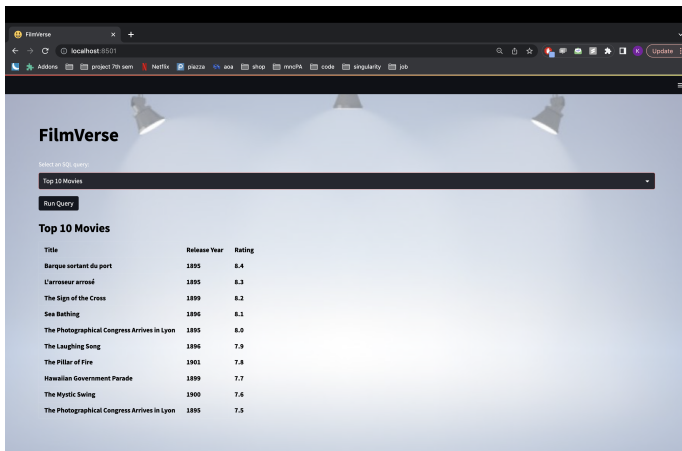


Fig. 18. UI Home Page

## IX. TEAM CONTRIBUTIONS

### Kajol

- Implemented Indexing
- Implemented Selection Queries
- Implemented frontend and backend for website
- Created Demo video

### Mahimitra Chirala

- Problem Statement Selection
- Normalisation check and dependencies in BCNF
- Writing Simple Queries
- Created Demo video/Report review/ Testing

### Sagar Gopalasetti

- Cleaning and modification of the dataset
- Creating Database
- ER diagram Implementation
- Report

## REFERENCES

- [1] The official IMDb datasets website: <https://www.imdb.com/interfaces/>
- [2] IMDbPy - a Python package useful for retrieving and managing IMDb data: <https://imdbpy.github.io/>
- [3] IMDb database schema: <https://www.imdb.com/interfaces/IMDb%20Schema>
- [4] SQL queries for IMDb datasets: <https://github.com/databasss/IMDb-SQL-Queries>
- [5] IMDb SQL database download: <https://www.kaggle.com/stefanoleone992/imdb-extensive-dataset>