# CSE 460/560 - DMQL Project Final Report

## I. 1. PROJECT DETAILS

Name of your project, your team, and all team members, everyone's UB id(not the UB number);

**Project Name: Preview before Review Movies**
**Group members:**
**Ryan O'Sullivan - rposulli**
**Sakshi Singhal - ss666**
**Ismail Ajaz - ismailaj**

## II. 2. PROBLEM STATEMENT [5 POINTS]:

Form a title and problem statement that clearly state the problem and questions you are trying to answer. Why do you need a database instead of an Excel file?

The title of the project is "Preview before Review Movies" and the problem statement is "Lack of prior knowledge about a movie exposes viewers with unexpected and even disappointing entertainment options." The questions that we can answer by using this dataset would be "I am feeling low. Which is the best movie to watch? I don't wanna waste my time like I did for the XYZ movie. It's been a long time since I haven't watched an action movie, let's find something." Because of its scalability, concurrent access, data integrity, security, and capacity for managing data relationships, databases will be a much better option than Excel files. They can handle massive volumes of data, preserve data consistency, and ensure data security because they have advanced querying, backup, and recovery capabilities. Additionally, they would reduce the redundancy that we would otherwise face if we combined several entity sets into a single table. Due to the limitations of Excel files in these areas, databases are a more reliable option for complicated data management requirements.

### A. Discuss the background of the problem leading to your objectives. Why is it a significant problem?

At the moment, people rely heavily on online reviews and suggestions, which frequently leads to surprising and potentially disappointing movie-watching experiences. Moviegoing consumers want to make sure that they are watching something that is worth their time and money in the theater, and without resources like movie databases, it will be much harder for them to make an informed decision about which movies to attend. This is a significant problem because the absence of resources like ours will worsen the moviegoer's average experience, and it may convince them to visit the theater less often. This will result in less revenue for the theater and less enjoyment from movie fans.

### B. Explain the potential of your project to contribute to your problem domain. Discuss why this contribution is crucial.

The project will provide moviegoers with a way of attaining information much more efficiently (as opposed to endless scrolling through reviews/suggestions on a variety of sites) and with more quality. Scrolling through random reviews means the user will have to ingest many unbiased and unqualified opinions of many different people who may expect and desire different things in a movie. Our contribution can resolve these issues because all of our data will be in one place (the database) and the data itself is significantly more unbiased because the data consists of reviews and ratings that come from regular users and actual critics, whose opinions are much more valuable than some random person who gave a bad rating because they were upset that they did not get to see a particular scene that they were expecting.

## III. TARGET USER [5 POINTS] :

Who will use your database? Who will administer the database? You are encouraged to give a real-life scenario;
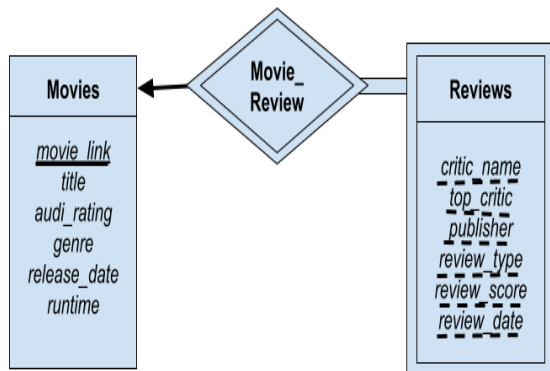
Anyone who wants to get a recommendation on a movie will be using our database. Suppose you are not someone who regularly sees movies, and you buy tickets to the Barbie movie without ever seeing a single review, just because of the branding of the Title. Imagine if the movie completely flopped because it was another corporate cash grab on a recognizable name. Imagine how disappointed you'd feel walking out of the theater. Wouldn't it have been so helpful to see the negative reviews about the movie? That way you could have made a more informed decision, like buying tickets to the Oppenheimer movie which was well-received. The people administering our database would be the popular theater businesses like Regal theaters, or alternatively, this could be presented by streaming services like Netflix or Hulu which have similar programs. They would get increased engagement because viewers would feel comfortable if they knew they were getting their money's worth before they subscribe, or continue their subscription to the service.

## IV. E/R DIAGRAM [10 POINTS]:

Draw an E/R diagram for your database and briefly describe the relationships between different tables. (Do not draw the figure by hand, you may use any tools to design or generate your E/R diagram

Movies is a strong entity set wherein each tuple corresponds to one movie, and other relevant information about that particular movie. It is identified by the primary key movie_link which is guaranteed to be unique for each movie on the site. Reviews is a weak entity set which depends on identifying set

Movies vía Movie_Review relationship, because each tuple is a movie review which logically depends on the existence of the movie being reviewed. This is also why all attributes of Reviews must be the partial key. One movie can receive many reviews, but each review can only be for one movie, so there is a one-to-many relationship. Additionally, every review must participate in Movie_Review because it always has a parent movie, but not every movie must have a review, hence the partial participation arrow.



V. TASKS (3 AND 4) SHOULD BE COMPLETED [10 POINTS]:

(Task 3) Acquire the large "production" dataset, either by downloading it from a real data source or generating it using a program. Make sure the dataset fits your schema. You might need to write programs/scripts to transform them into a suitable form for loading into a database for real datasets. For program-generated datasets, make sure they contain interesting enough "links" across rows of different tables to show the results of different Advanced SQL queries learned in class.

**Done. See figures 1 and 2 at near the bottom of this report.**

(Task 4) You are required to make sure all of your relations are in Boyce- Codd Normal Form. Provide a list of dependencies for each relation. Decompose them if the tables are not in BCNF. If you decide to keep it in 3NF instead of BCNF, justify the decision for a particular relation. Your report for this milestone should contain a separate section with the details of the transformation from the initial schema to the final schema where the relations are in BCNF

The only noteworthy FDs are: $movie\_link \rightarrow movie\_link, title, audi\_rating, genre, release\_date, runtime,$

and (movie_link, critic_name, top_critic, publisher, review_type, review_score, review_date) $\rightarrow$ (movie_link, critic_name, top_critic, publisher, review_type, review_score, review_date).

For the Reviews table, because all attributes are needed to uniquely identify a tuple, there are no nontrivial functional dependencies here. Only the 2nd FD mentioned, the trivial one, holds. Thus, Reviews are in BCNF.

For the Movies table, the only related FD is that movie_link -¿ movie_link, title, audi_rating, genre, release_date, runtime. movie_link is the primary key/superkey, and none of the remaining attributes are in the LHS of an FD, so Movies is also in BCNF

*A. 4.1 Define a list of relations and their attributes.*

The relation schemas for each relation is as follows: Movies(movie_link, title, audi_rating, genre, release_date, runtime) Reviews(movie_link, critic_name, top_critic, publisher, review_type, review_score, review_date ) Indicate the primary key and foreign keys (if any) for each relation. Justify your choice; For Movies, the primary key is movie_link because each movie is hosted at a different URL on the site, and thus is guaranteed to be a unique identifier for each movie. For Reviews, a weak entity set, there is no single attribute primary key, or even an obvious single-attribute partial key. This is because reviews can have the same critic name, review the same movie, suggest the same score on the same date, etc. Because there are no true identifiers in the weak entity set, we must use ALL attributes as a primary key for reviews. This is common in weak entity sets without an obvious partial key. (movie_link, critic_name, top_critic, publisher, review_type, review_score, review_date)

*B. 4.2:*

Write a detailed description of each attribute (for each table), its purpose, and datatype;

| Movies() | | |
|---|---|---|
| ATTRIBUTE NAME | PURPOSE | DATATYPE |
| movie_link | URL extension to RT movie's webpage. Unique for each movie. | Varchar(MAX) |
| title | Title of the movie. | Varchar(MAX) |
| audi_rating | Audience score of the movie, from 0 to 100. | Integer |
| genre | Genre of the movie (Science Fiction, Action, Fantasy). | Varchar(MAX) |
| release_date | Release date of the movie in format yyyy-mm-dd. | Date |
| runtime | Movie runtime in minutes. | Integer |

| Reviews() | | |
|---|---|---|
| ATTRIBUTE NAME | PURPOSE | DATATYPE |
| movie_link | URL extension to RT movie' webpage. Unique for each movie. | Varchar(MAX) |
| critic_name | Name of the critic. | Varchar(MAX) |
| top_critic | Boolean "True" or "False" indicating the critic's authority. (Would be mapped 1 or 0 for BIT) | Varchar(5) or BIT |
| publisher | Publishing company/organization of the critique. | Varchar(MAX) |
| review_type | Either "Fresh" or "Rotten", indicates insult or appraisal in the review. | Varchar(MAX) |
| review_score | Score fraction from 0 to 1. Often represented out of 5 but easily translatable to decimal. | Decimal(4,4) |
| review_date | Release date of the review in format yyyy-mm-dd. | Date |

| Movies() | |
|---|---|
| ATTRIBUTE NAME | DEFAULT VALUE OR NULL |
| movie_link | No default value. NULL is prohibited because this is a primary key. |
| title | No default value. Every movie in the set has, and must have, a title. |
| audi_rating | NULL; 0 cannot be default here because it would bring down average movie ratings unfairly. |
| genre | NULL; "" is also acceptable but NULL is preferred |
| release_date | NULL; |
| runtime | NULL; |

| Reviews() | |
|---|---|
| ATTRIBUTE NAME | DEFAULT VALUE OR NULL |
| movie_link | No default value. NULL is prohibited because this is the primary key of the identifying set. |
| critic_name | NULL; "" is acceptable but NULL is used in practice. |
| top_critic | "False". If not specified or new, assume the critic is not a top critic. |
| publisher | NULL; **In practice this doesn't because all of the critics in the set have publishers |
| review_type | NULL; |
| review_score | NULL; |
| review_date | NULL; |

## C. 4.3

Indicate each attribute's default value (if any) or if the attribute can be set to 'null';

## D. 4.4

Explain the actions taken on any foreign key when the primary key (that the foreign key refer to) is deleted (e.g., no action, delete cascade, set null, set default)

When a particular movie_link is deleted from Movies, you would have to remove the entire tuple because that attribute is the primary key of Movies. Thus, the only occurrence of this would be if the entire tuple were deleted from the identifying entity set. Even then, we would take "NO ACTION". Just because a movie is deleted from Movies, doesn't mean that the movie review information is useless. We want to keep the correlated tuples in Reviews for further analysis. Additionally if the related movie_link value is inserted again later in Movies, we will still have the review data that goes along with it.

**This concludes Phase 2 of the Final Report.**

Fig. 1. Query sample from Reviews() table.



Fig. 2. Query sample from Movies() table.

## VI. TASK 5: PROBLEMS AND INDEXING CONCEPTS

"Do you specifically run into any problems while handling the larger dataset? Did you try to adopt some indexing concepts to resolve this? Briefly describe the questions you faced and how you solved them."

Even on our larger dataset, we did not run into any problems handling it using PostGreSQL. However, to demonstrate our understanding, we will discuss some common indexing methods used to handle large datasets. On ordered files like our Movies() table, which is ordered by movie_link, it would be implemented as a primary index, and NOT a clustered index. This is because records close in the index would also be close in data. Within this primary index, we can choose either a dense index model or a sparse index model. The main difference is that dense indices contain a record for every search-key value in the data file (in our case, an index record for every data record). However, sparse indexes have an advantage because they only include one index key per data block. Because of this and the fact that the data file is only sorted by the search key, we would prefer a sparse index for Movies(). The Reviews() table is a weak entity set and while it is ordered with respect to movie_link(), it is UNordered with respect to critic_name, which is a common search key for movie reviews. Because it's sorted on one of these dimensions, a secondary index would not be necessary. In this case, we would want to use a dense clustered index (variation of a primary index) in order to manage this large dataset. Dense clustered indexes are ideal for this case because the Reviews() table has many possible search keys, including movie_link, critic_name and publisher.

## VII. TASK 6: TESTING THE DATABASE ON SQL QUERIES

"Test your database with more than 10 SQL queries. You are supposed to design 1 or 2 queries for each inserting, deleting, and updating operation in your dataset."

For demonstration, we have performed over ten sample queries on the database, and we have displayed their results below. Figure 3 inserts a new movie "Kuch Kuch Hota Hai" into the Movies() table, Figure 4 updates the genre and runtime of the newly inserted tuple, and Figure 5 deletes this movie from the Movies() table entirely. Similarly, Figures 6, 7 and 8 are an insertion, update and deletion of a review for the movie "Kuch Kuch Hota Hai." In this case, the update changed the review_score and publisher name to 'Buffalo News'. In figures Figure 9 and 10, we performed the inner join on the Movies() and Reviews() table.

In figures 11 we first order the movies on their titles from Z-A (descending order). Then in figure 12, we ordered the movies in ascending order by genre.

We performed some similar queries for the Reviews() table. In figure 13, we ordered the reviews by the last name of their critics, from Z to A. Then in reference 14, we reordered the table alphabetically by release date.

Finally, we performed three grouping queries. In figure 15, we count the number of movies with the genre "Comedy" and find it to be 1263 total movies. Figure 16 is more general and shows us the movie count for every genre. Lastly, in figure 17, we group by movie_links to get the total number of reviews for each movie.



Fig. 3. Insertion of a tuple in the Movies() table.



Fig. 4. Updating in the Movies() table

Fig. 5. Deletion of a tuple from the Movies() table



Fig. 6. Insertion of a tuple in the Reviews() table



Fig. 7. Updating in Reviews() table



Fig. 8. Deletion of a tuple from the Review() table



Fig. 9. Joining the Movies() and Reviews() tables.



Fig. 10. Joining the Movies() and Reviews() tables, continued.



Fig. 11. Ordering in Movies() table in descending order.

## VIII. TASK 7: PROBLEMATIC QUERIES AND PERFORMANCE IMPROVEMENTS

"Query execution analysis: identify three problematic queries (show their cost), where the performance can be improved. Provide a detailed execution plan (you may use EXPLAIN in PostgreSQL) on how you plan to improve these queries."

Our first problematic query, shown in figure 18, is a simple query with a glaring inefficiency: We join Reviews() and Tables(), but we only look for attributes from Reviews()! In a database with a small schema, this may be obvious, but in larger schemas it is important to only join tables which would provide you with the desired attributes. This strategy of "only necessary joins" is how we would improve

Fig. 12. Ordering in Movies() table in ascending order.



Fig. 13. Ordering in Reviews() table in descending order.



Fig. 14. Ordering in Reviews() table in ascending order.



Fig. 15. Performing GroupBy operation to get the total count of Comedy Movies in the Movies() table.



Fig. 16. Performing GroupBy operation in the Movies() table, grouping movies by genre.

on queries like this in large databases. Otherwise, as shown in the EXPLAIN function, you will incur extra time costs from the join operation.

The second problematic query we present is shown in figure 19. In this case, we join Reviews() and Movies(), only to later apply a condition on the release date of the movie in question. This is problematic because the joined table will have more total tuples than Movies(), and will require a longer scan than if we had just applied the release date condition to the Movies() table BEFORE joining it, and this inefficiency is reflected in the JOIN cost in the figure. In bigger databases it is important to apply conditions before joins when possible, in order to shrink the tables and lower the cost of the join. Applying conditions before joins is another useful data strategy to manage queries like this one.

Our third and final query, presented in figure 20, has a number of problematic components that increase the total execution time of the process. For one thing, including the phrase "HAVING COUNT ¿ 0" is redundant because if a movie exists at all, it will have a title since we're grouping by movie link. Additionally, we scan on two conditions in a nested select statement instead of just applying them together. This is cumbersome and leads to many of the inefficiencies you see in the explanation table associated with the figure. In this case, the query is problematic because it includes redundant clauses and unnecessary nesting which bloats the runtime of the query. To combat this, we would modify the query to put both WHERE clauses in a single SELECT statement, and

Fig. 17. Performing GroupBy operation using to acquire the total number of movie links for each movie.

remove the HAVING clause entirely. In general you want to simplify overcomplicated queries to prevent running excessive scans, and this is especially true for larger databases where each scan takes much longer.

**This concludes Phase 2 of the Final Report.**



Fig. 18. Problem Query 1: It is inefficient to join several tables when you're only searching for values in one of them.

**References :- Link to Rotten Tomatoes Dataset**



Fig. 19. Problem Query 2: It is inefficient to join before applying a condition; We should apply the condition to lower the amount of tuples before joining.



Fig. 20. Problem Query 3: There are several inefficiencies here; The conditions are separated and the addition of a COUNT call should be optimized out.