

UNIVERSITÀ DI PISA



Department of Information Engineering(DII)

UMDb - Movie Review

Tesfaye Yimam Mohammad

Arsalen Bidani

Aida Himmiche

Semester Final Project of Large Scale and
Multistructure Database course.

Winter Session
Academic Year 2021/22

Contents

List of Figures	ii
List of Tables	iii
1 Introduction	1
2 Design	1
2.1 The Application Requirements	1
2.1.1 Actors	1
2.1.2 Functional Requirements	2
2.1.3 Non-Functional Requirements	3
2.2 Use case Diagrams	4
2.3 Class Analysis Diagram	6
2.4 Dataset	6
2.5 Database Design	8
2.5.1 Data Models	8
2.5.2 Graph Database	12
2.5.3 Indexes	15
2.6 Software Architecture	19
2.6.1 Modules	19
2.6.2 User Interface	20
2.7 Database Deployment	21
2.7.1 Cluster replicas Topology and Consistency management	21
2.7.2 Cross Database consistency management	23
3 Sharding	25
4 Implementation	25
4.0.1 MongoDB Query Implementations And Analysis	25
4.1 Main Queries in the GraphDB	34
4.1.1 User nodes	34
4.1.2 Movies Nodes	35
4.1.3 Watchlist nodes	36
4.1.4 Recommendations/Suggestions in the GraphDB	36
4.1.5 Summary of Queries analysis for the GraphDB	41
References	41

List of Figures

Figure 1	Class Diagram of UMDB	6
Figure 2	A document inside the Movie collection	11
Figure 3	A document inside the User collection	11
Figure 4	A document inside the Review collection	12
Figure 5	The GraphDB Model of UMDB	14
Figure 6	The Software architecture of UMDB	19
Figure 7	Home Page	21
Figure 8	Replica topology for the document database	22
Figure 9	Cross Database Coherence for ADD USER operation .	23
Figure 10	Cross Database Coherence for ADD MOVIE operation	24
Figure 11	Cross Database Coherence for DELETE USER operation	24
Figure 12	Cross Database Coherence for DELETE MOVIE operation	25

List of Tables

Table 1	Read Operations on MongoDB	26
Table 2	Write Operations on MongoDB	27
Table 3	Read Operations in the GraphDB	41
Table 4	Write Operations in the GraphDB	41

Listings

1	Popularity Simple Indexing type	15
2	MongoDB: Year Genre Component Index type	16
3	Neo4j Indexing: User Id	16
4	Neo4j Indexing:Full text: One	16
5	Neo4j Indexing:Full text:Two	16
6	Deleting a user node	34
7	Creating User FOLLOW User relation	34
8	Deleting a user FOLLOW user relation	35
9	Creating a movie node	35
10	Delete Movie node	35
11	Creating a watchlist node	36
12	Updating a watchlist node	36
13	Recommended Users	36
14	Recommended Users Express JS Method	36
15	Recommended Movies	37
16	Recommended Movies Express JS Method	37
17	Recommended Watchlists Express JS Method	38

1 Introduction

UMDB is a web based application that aims to provide a platform for UniPi students where they can find information about a large collection of movies, give some reviews on them, and interact with other users as well. UMDB stands for “UNiPi Movie Database”, therefore it stores movie-related information including ratings and reviews. Users can view movies, rate them, write reviews on them, and/or add them to their watch-lists. Users can also follow other users and view other users’ profiles in addition to their own, moreover, movies and user profiles will be suggested based on their interactions on the platform. While User records may be manipulated by the users themselves, movie data manipulation will be controlled solely by the administrators.

2 Design

2.1 The Application Requirements

2.1.1 Actors

The major actors in the application are grouped in three categories: Anonymous users, Registered users and Administrators. Every registered user has their own privilege in the system based on their credentials given to them during a registration time.

- **Anonymous User:** A user with no account in the system has the most restricted access in the application. As a result, this user type is invited to create an account to exploit the extensive functionalities of our application.
- **Registered User:** This is the type of user which has login credentials to access their own session in the application. Once logged in he/she

can post reviews, create watch-lists, follow watch-lists, follow other users, and more.

- **Administrator:** it is the third type of users which has access to all previously mentioned functionalities, in addition to adding/removing movies and/or users from the database(s).

2.1.2 Functional Requirements

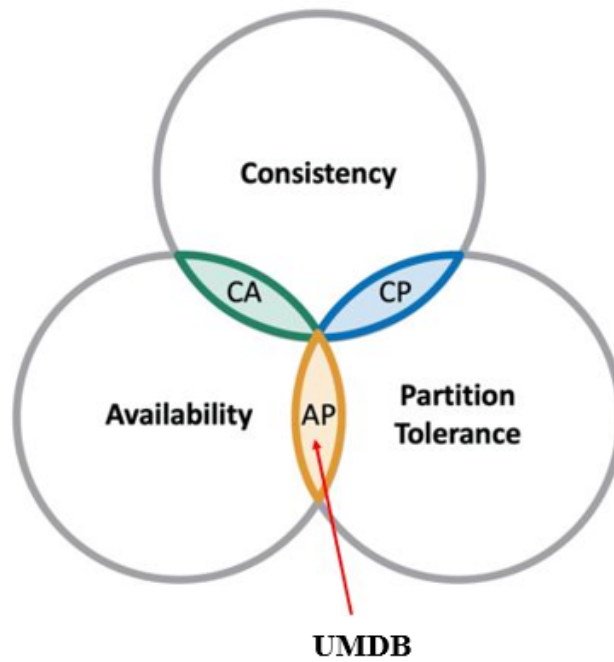
Here is a list organized by actors, describing the functionalities available for each:

- **Anonymous User:**
 - Create a personal account in the application.
 - Log into a personal account using credentials.
 - Search for movies and read user reviews.
 - Search for registered users on the platform.
 - Find general analytics about movies, users, and reviews on the platform.
- **Registered User:** All previously mentioned in addition to:
 - View/edit personal information on the user-dedicated page.
 - Post a review about a movie.
 - Edit/delete personal previously-posted reviews.
 - Create a watch-list.
 - Edit/delete personal watch-lists.
 - Add/remove a movie from a personal watch-list.
 - View other users' watch-lists.
 - Follow/unfollow other users' watch-lists.
 - View list of followed watch-lists.
 - Follow/unfollow other users.

- View list of followed users.
- Receive user profile suggestions.
- Receive watch-list suggestions.
- Receive movie suggestions.
- Logout of personal account.
- Delete personal account.
- **Administrator:** All previously mentioned in addition to:
 - Add/delete a movie.
 - Delete any user account of choice.
 - Delete any user review of choice.
 - Give administrator privileges to users of choice.

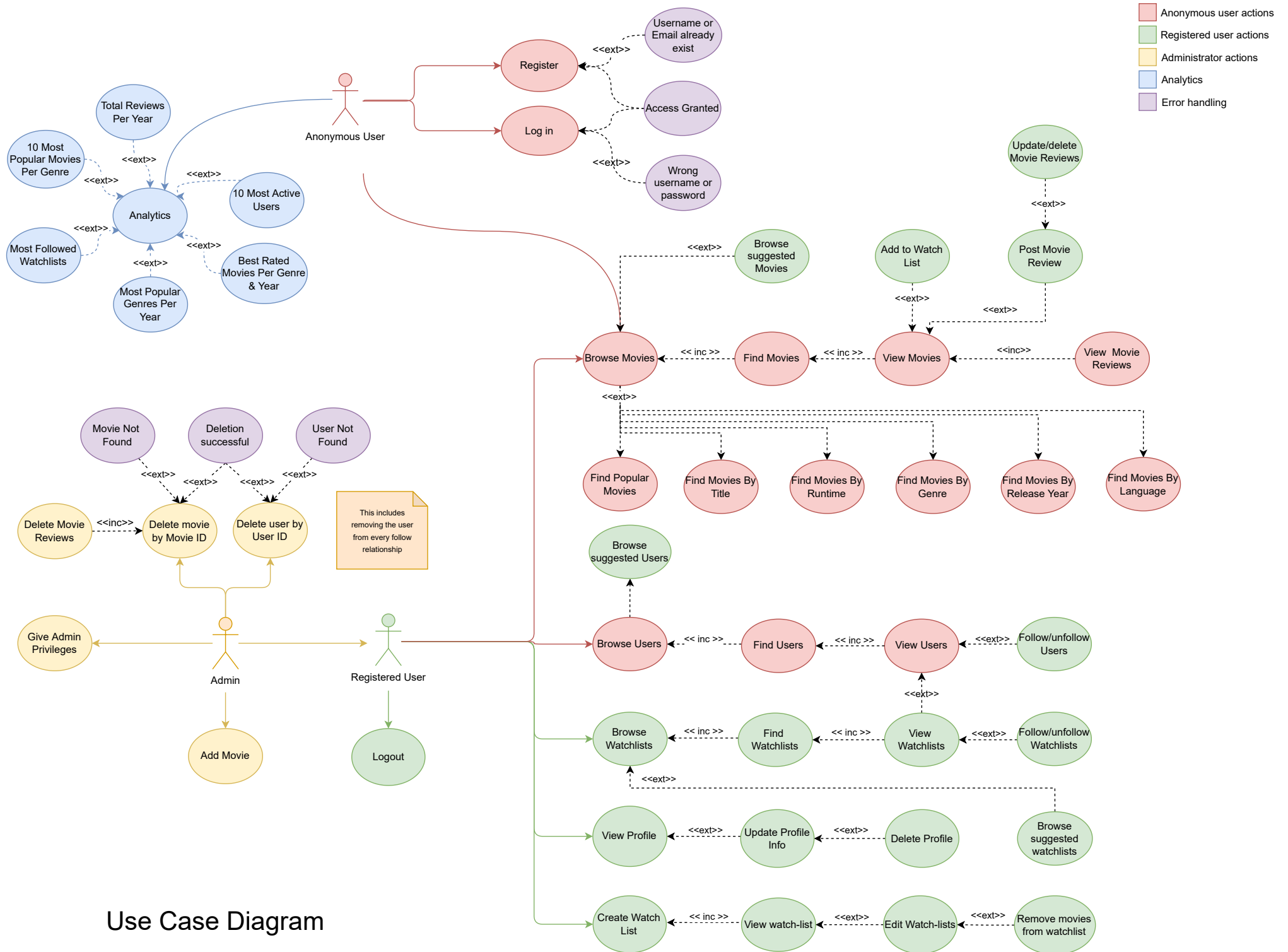
2.1.3 Non-Functional Requirements

- Usability: The user interface should be as comfortable as possible to the end user.
- Users must be able to exploit all the functionalities of the application being supported by high availability, a low latency and a good level of tolerance with respect to the possibility that there might be network partitions. Therefore, the application has to be designed in order to cover the AP (Availability / Partition protection) edge of the CAP triangle. More detail is found at fig.2.1.3.



- Maintainability: The application's code must be easy to read and to maintain.
- Security: Users' passwords must be stored in an encrypted form within the database, to guarantee an high level of security the SHA-256 function will be adopted.

2.2 Use case Diagrams



2.3 Class Analysis Diagram

Fig.1 shows the Class diagram of the application.

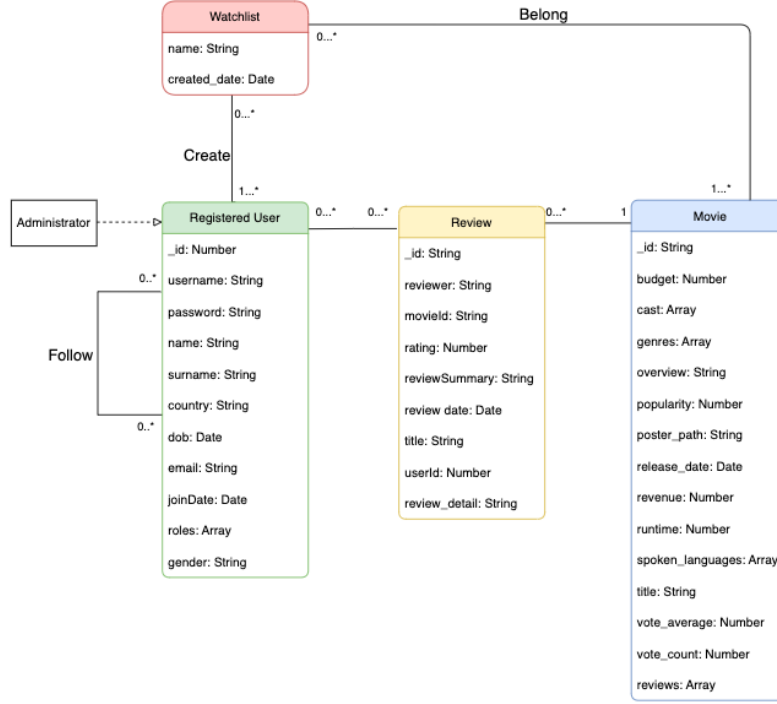


Figure 1: Class Diagram of UMDB

2.4 Dataset

Description To simulate a large-scale database for this application, we consulted different sources to amass data.

We chose a movie dataset provided on the Kaggle platform under the name *The Movies Dataset*. It contains a total of 45k movies, with information including the title, budget, cast, release date, genres, runtime, overview, languages spoken, revenue, average rating, and number of votes.

The user reviews were acquired from a different dataset on the same platform with the name *IMDb Review Dataset*, it grouped a total of over 190k reviews with information on the reviewer, movie title, rating, review summary, review detail, and review date. The way we matched the reviews with the movies was using a common field referring to the IMDb ID of the movie.

User information was not provided by any dataset, instead we used the open-source API *Random User Generator* to randomly generate 78k profiles with personal details such as the E-mail, password, name, surname, gender, country, date of birth, and date of account creation. To connect this data to the reviews, we randomly assigned each user profile one of the unique usernames in the reviews.

Volume and Format Here are the details about the volume and format of the datasets after complete manipulation:

- **The Movies Dataset:** Size : 323.6MB ; Format : CSV.
- **IMDb Review Dataset:** Size : 258.2MB, Format: JSON.
- **Users Dataset:** 13.3MB; Format : CSV.

Variety The data we used is multi-sourced and multi-formatted in order to achieve the aspect of variety.

Velocity/Variability The administrator can populate the database with new movies to achieve velocity.

2.5 Database Design

To store all the data of the application we have decided to use a Document Database (MongoDb) and a Graph database (Neo4j). The designed data model ensures that each read operation requires a single call either to the Graph Database or to the Document database.

2.5.1 Data Models

The **MongoDb** has been used in order to maintain data for the following collections: *users*, *movies*, *reviews*. Here we present the detailed information about the attributes found in the documents of each.

- Users
 - `_id`: Number
 - `username`: String
 - `email`: String
 - `password`: String
 - `gender`: String
 - `name`: String
 - `surname`: String
 - `country`: String
 - `dob`: Date
 - `roles`: Array[String]
 - `joinDate`: Date
- Reviews
 - `_id`: String
 - `userId`: Number
 - `reviewer`: String

- movieId: String
- rating: Number
- title: String
- review_summary: String
- review_detail: String
- review_date: Date
- Movies
 - _id: String
 - title: String
 - budget: Number
 - cast: Array[Object]
 - genres: Array[String]
 - overview: String
 - popularity: Number
 - poster_path: String
 - release_date: Date
 - revenue: Number
 - runtime: Number
 - spoken_languages: Array[String]
 - vote_average: Number
 - vote_count: Number
 - reviews: Array[Object]

NB: The 'roles' entry in the users collection differentiates the normal registered users from those who have administrative privileges in the system.

In the upcoming paragraphs we are going to present the skeleton of the document database and the reasoning behind it.

Movie Collection: As shown in fig.2 the "movies" collection comprises data about the movies imported from the appropriate aforementioned CSV file.

Our first design of this collection stated that the movie document included its reviews as embedded documents in the form of an array, since the reviews would only be retrieved when a Read operation on the movie is executed. However, considering the nature and context of our application (being Read-heavy and prone to accumulating large numbers of reviews per movie) we thought about the size of the documents and the speed of the Read operations on the collection.

The maximum size a document can have in a MongoDB collection is 16Mb. In our "movies" collection, the average size of a movie document excluding the reviews is 4.5kB, and the average size of a review document is 1.5kB. The maximum number of reviews per movie in our dataset is more or less 2880. At first glance, it seems as though the final size of the document with embedded reviews won't be a problem because of the small size of the review dataset at hand. However, looking at a case where a low number of 15000 registered users wanted to post one review each on a popular movie, the size would easily exceed 16Mb, not only making the Read operations very slow, but also making the Write operations impossible.

Having ruled out the possibility of embedding the reviews into the "movies" collection, we thought of keeping them in a collection of their own, with a reference to the movie they are posted for. However this would mean that whenever a movie is to be displayed, two collections need to be accessed to get the entire data, and that beats the purpose of using NoSQL.

Since neither extremes are suitable for our application, we decided to opt for the middle ground. We created our "movies" collection with an embedded array of only the 20 most recent reviews for each movie, and also kept the reviews in another collection in the case of the user wanting to load more. This allows us to avoid a join operation right from the start, but only resort to it when required.

```

    _id: "tt0114709"
    budget: 30000000
  > cast: Array
  > genres: Array
  overview: "Led by Woody, Andy's toys live happily in his room until Andy's birthd..."
  popularity: 21.946943
  poster_path: "/uXDfjJbdP4ijW5hWSBrPrLKpxab.jpg"
  release_date: 1995-10-30T00:00:00.000+00:00
  revenue: 373554033
  runtime: 81
  > spoken_languages: Array
  title: "Toy Story"
  vote_average: 7.7
  vote_count: 5415
  > reviews: Array
    > 0: Object
      _id: "rw0968444"
      review_summary: "The greatest animation ever made."
      userId: "105"
      reviewer: "nickfurze"
      review_date: "2004-11-27"
      review_detail: "This film amazed me it has one of the greatest scripts ever made it's ..."
      rating: "10"
    > 1: Object

```

Figure 2: A document inside the Movie collection

User Collection: As depicted in fig:3 The documents in this collection provide the all the necessary information about the Users. The queries on this collection are independent from the other two and comply with the NoSQL requirements.

```

    _id: 0
    username: "ahim"
    email: "aida@email.com"
    password: "$2b$10$Sv97s4Y/Ko8wRqjPGtZ4AucF2dzoqS8pVu3oXNZmMQw0v7rgwxGdG"
    gender: "f"
    name: "aida"
    surname: "him"
    country: "Morroco"
    dob: 1967-11-01T00:00:00.000+00:00
  > roles: Array
    > 0: "User"
    > 1: "Admin"
  joinDate: 2015-04-22T00:00:00.000+00:00

```

Figure 3: A document inside the User collection

Review Collection: This collection groups all the reviews imported from the appropriate dataset, including the ones embedded in the movie documents. Our reason behind allowing repetition in this case, is because without

it, it would be hard to keep track of the reviews that are being added as they would be pushed directly in the embedded array, moreover, aggregations on the reviews would require a lookup operation with the access of both the "reviews" and "movies" collection.

Since we would prefer to avoid join operations as much as possible, we made our decision keeping in mind that no inconsistencies between the duplicated reviews shall arise. For that, we made sure of the atomicity of Create-Update-Delete actions on said duplicates.

NB: The "userId" attribute stored in the review documents is not used in any join/lookup operations. It is stored as such for purposes related to front-end navigation and authentication. The fact that the field never undergoes any sort of modification makes it so that no consistency management is needed.

```
_id: "rw1986516"
movieId: "tt0114709"
rating: 10
review_date: 2008-07-12T00:00:00.000+00:00
review_detail: "The movie was totally fun and cool. The toys in Andy's Room are superb..."
review_summary: "It was totally Fun and cool"
reviewer: "manikhero"
title: "Toy Story"
userId: 1
```

Figure 4: A document inside the Review collection

2.5.2 Graph Database

The graph database has been used in order to move in the most efficient way through the FOLLOW, CREATE and BELONGS relationships between different users, watch-lists and the movies nodes. By exploiting this NoSQL architecture, some queries can be efficiently performed to find, having a certain user as a starting point, a list of suggested users to follow or a list of recommended watch-lists based on the the number of followers they have and the movie recommendations based on their occurrence in a number of watch-lists to mention a few. **GraphDB Nodes:** As mentioned above there are three main nodes in the GraphDB and below we have presented the attributes of each node.

Taking into account these objectives, it was not necessary to store in the Graph Database too much information concerning the instances of each and every entity. In particular, for each of these ones, only the following information are present in the graph database:

- User
 - username :String
 - id: Int
- Watchlist
 - created_date :Date
 - name :String
- Movie
 - title: String
 - movie_id: String
 - overview :String
 - poster_path :String
 - vote_average :Int

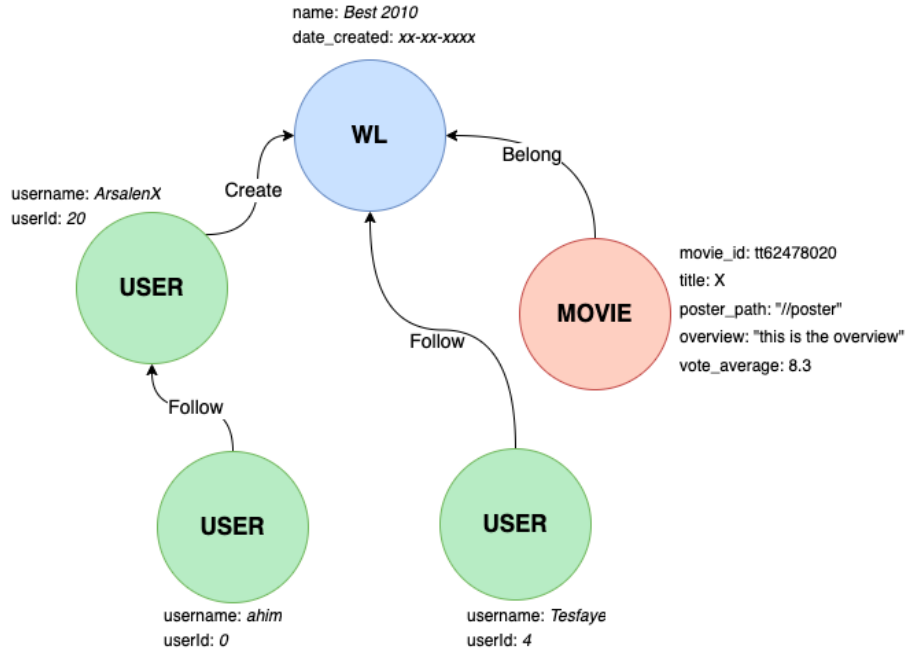


Figure 5: The GraphDB Model of UMDB

Relation types in the Neo4j As mentioned above the three nodes are connected in the neo4j database as follows.

- User-User : $(u1:User)-[:FOLLOW]\rightarrow(u2:User)$ It connects a user to an other user. It has no attributes
- User-Watchlist : $(u:User)-[:FOLLOW]\rightarrow(w:Watchlist)$. This relation connects user to watchlist nodes he/she is following. It has no attributes.
- Movie-Watchlist : $(m:Movie)-[:BELONGS]\rightarrow(w:Watchlist)$. This relation connects a movie node with a watchlist node. It has no attributes
- User-Watchlist : $(u1:User)-[:CREATE]\rightarrow(u1:User)$. It connect the user with the watchlists he/she has created. It has no attributes

2.5.3 Indexes

One of the components of the database design is indexes, as it makes it easier to query the large documents. We stated that one of the requirements is the speed handling of the web application.

We opted for 4 indexes on **MongoDB** as well as 4 indexes in **Neo4j**.

- **MongoDB**
 - **Popularity** Simple Index Type
 - **Year,Genre** Compound Index Type
 - **_id** Hashed Index Type (Default)
 - **Username** Unique Index Type
- **Neo4j**
 - **User_id** BTREE Index Type
 - **username** Unique Index Type
 - **watchlist_id** BTREE Index (Default)
 - **Username** FULLTEXT Index Type

Indexes in the performance as they optimise multi-key lookups but we need to be careful of the costs of indexes on update,inserts and deletes.

Popularity Simple Indexing type

```
1 db.movies.aggregate([ { $sort: { popularity: -1 } }, { $project: { title:  
    1, poster_path: 1, overview: 1, vote_average: 1, release_date: 1, },  
    }, ])
```

Listing 1: Popularity Simple Indexing type

Year Genre Component Index type

```

1 db.movies.aggregate([ { $match: { $and: [ { release_date: { $gte: new
    ISODate("2012-01-12T20:15:31Z"), $lt: new ISODate("2020-01-12T20:15:31
    Z") } }, { genres: { $regex: "Comedy", $options: "i" } }, ], }, }, {
    $sort: { vote_average: -1 } }, { $project: { title: 1, poster_path: 1,
    vote_average: 1, release_date: 1, }, }, ]).explain("executionStats")

```

Listing 2: MongoDB: Year Genre Component Index type

The indexing for the Neo4j part

For the Neo4j we have done the User Id and the Full text indexing. Below are the listings for both of them, according to their order.

Indexing for User Id

```

1 match(u:User)
2     where u.user_id=${userId}
3     with *
4     optional match(ux:User)-[followers:FOLLOW]->(u)
5     optional match(u)-[following:FOLLOW]->(Uy:User)
6     return {user_id:u.user_id,username:u.username,following:count(
    distinct following), followers:count(distinct followers) }

```

Listing 3: Neo4j Indexing: User Id

FULLTEXT with index:

```

1 CALL db.index.fulltext.queryNodes("userSearch", "a")
2     YIELD node,score
3     return {user_name:node.username,user_id:node.user_id}
4     ORDER BY score
5     SKIP ${skip}
6     LIMIT 7

```

Listing 4: Neo4j Indexing:Full text: One

FullTEXT with out index:

```

1 MATCH (u:User)
2     WHERE toUpper(u.username) CONTAINS toUpper("${searchString}")
3     with *

```

```

4 optional match(ux:User)-[follower:FOLLOW]->(u)
5 with u,count(follower) as num
6 where num in range(${min},${max})
7 return {user_name:u.username,user_id:u.user_id}
8 ORDER BY u.username
9 SKIP ${skip}
10 LIMIT 7

```

Listing 5: Neo4j Indexing:Full text:Two

With Out indexing for Compound:

```

1 nReturned: 568,
2   executionTimeMillis: 92,
3   totalKeysExamined: 0,
4   totalDocsExamined: 45145,

```

With Indexing for Compound:

```

1 nReturned: 568,
2   executionTimeMillis: 8,
3   totalKeysExamined: 3417,
4   totalDocsExamined: 568,
5   executionStages:

```

Without index for the Genre:

```

1 nReturned: 568,
2   executionTimeMillis: 92,
3   totalKeysExamined: 0,
4   totalDocsExamined: 45145,

```

With indexing for the Genre:

```

1 nReturned: 568,
2   executionTimeMillis: 8,
3   totalKeysExamined: 3417,
4   totalDocsExamined: 568,
5   executionStages:

```

Neo4j User Id: without index

```
1 without index
2 | Time | DbHits | Rows | Memory (Bytes) |
3 | 105  | 157306 | 1    | 1512
```

Neo4j User Id: with index

```
1 | Time | DbHits | Rows | Memory (Bytes) |
2 +-----+
3 | 10   | 57     | 1    | 1432           |
4 +-----+
```

Neo4j FULLTEXT: without indexing

```
1 | Time | DbHits | Rows | Memory (Bytes) |
2 | 105  | 157306 | 1    | 1512
```

Neo4j FULLTEXT: with indexing

```
1 time | DbHits | Rows | Memory (Bytes) |
2 | 10   | 226    | 113  | 0
```

2.6 Software Architecture

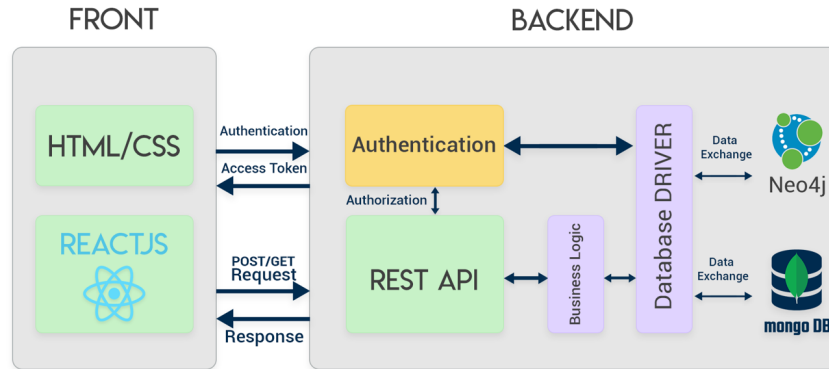


Figure 6: The Software architecture of UMDB

We used the well known Model-View-Controller (MVC) software development pattern to design our project. The model contains the data logic of the movies, reviews and users collections. The controller contains all the control logic for all the documents and graph DB nodes.

2.6.1 Modules

As depicted in fig.6 the UMDB application has the following organization, according to the MVC architecture.

- **Backend:-** This is a folder that contains all the back end related files. All the controllers, models, routes and middle-ware of the application are grouped there.
- **Public:** This folder includes the index.html file and images referenced throughout the application.

- src: it comprises the API, hooks and all the JavaScript files that play a role in the rendering of the application.

We used the MERN stack plus the Neo4j graph database. Below are the list of the software tools used to accomplish this work.

- HTML: to design the basic structure of the web application.
- CSS: for manipulating the aesthetic of the front-end pages.
- React: to design the UI of our application.
- POSTMAN: for testing the routes to our business logic.
- Neo4j: The nodes of the graph DB are stored and queried from the Neo4j database.
- MongoDB: The collections of the document database are stored and queried from the MongoDB database.

2.6.2 User Interface

For a better interaction between the end user and application under the hood of the databases, it has been decided to use the the React framework for the front end UI. As said previously from the MERN stack React is used to build all the UI part so that the end user can have one of the finest experiences to the state-of-the-art web technology. The Home page of the application is shown in fig.7. It is very easy and user friendly application so that the end user can easily click on the provided navigation buttons to surf our website. Moreover, in this page the most popular movies are depicted.

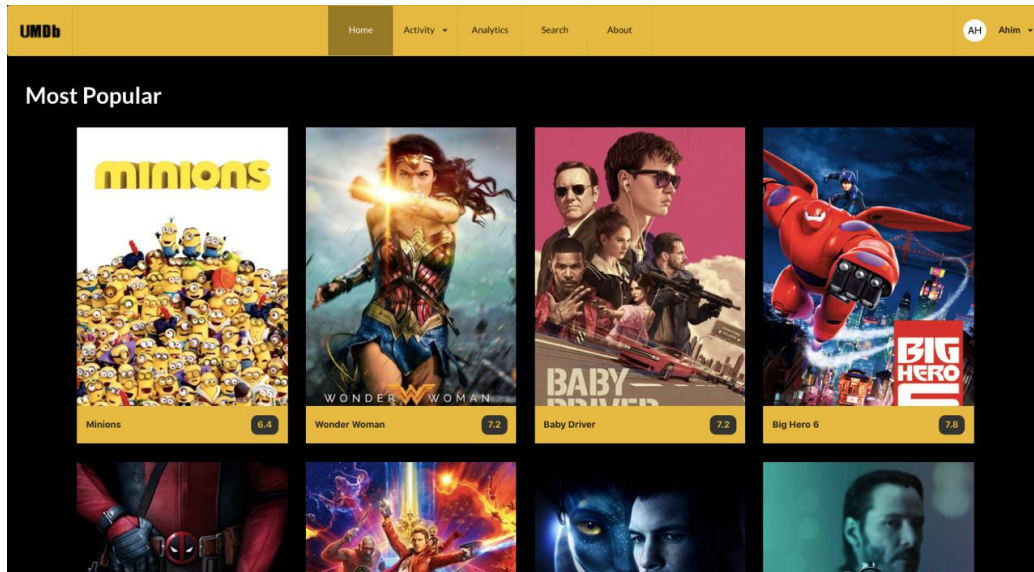


Figure 7: Home Page

As seen in fig.7 all the buttons in the navigation bar helps the user to surf the site. In the right top corner there is a Sing in and Login option. As can be seen the home page has got a suggestion of the most popular movies based on the user activities.

2.7 Database Deployment

2.7.1 Cluster replicas Topology and Consistency management

The document database has been deployed on a cluster of three virtual machines by which their IP address as depicted in fig.8

- The primary replica is at the *172.16.4.60*
- The Secondary replica is at the *172.16.4.64*
- The last replica is at the *172.16.4.65*

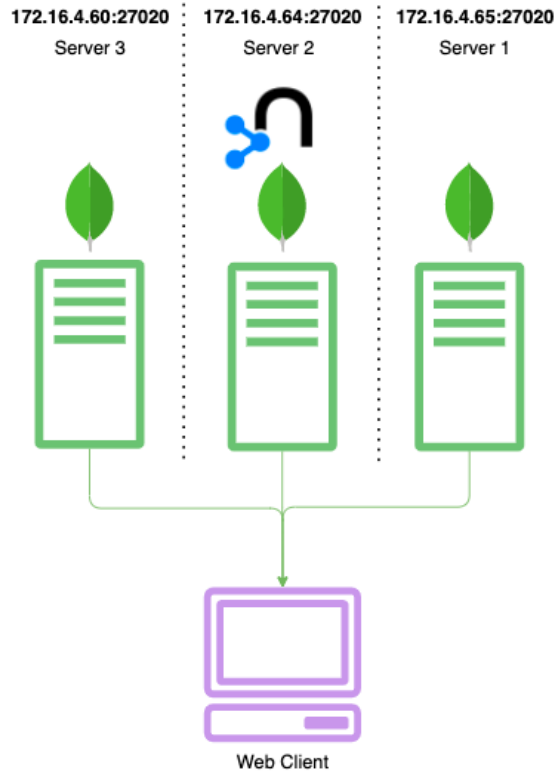


Figure 8: Replica topology for the document database

The data consistency between the different collections of the document database is completely handled by the application code.

A higher availability level is going to be provided, since the databases won't get stuck communicating in order to synchronize every modification occurred; at the same time, plus to this the usage of partition replica, removes the possibility of a denial of service situation, since the eventual failure of the primary server won't stop the application's runtime. Hence, the combination of the two properties from the CAP triangle locates our UMDB application into the (A-P) edge of the triangle as seen from fig.2.1.3.

2.7.2 Cross Database consistency management

Consistency between MongoDB and GraphDB

So as to maintain the consistency between the two databases, for each write operation that took place in the DocumentDB, the GraphDB is updated accordingly before the control is returned to the application. In fig.9 the logic to maintain the consistency between the databases for adding a user data. This same logic is applied to add a movie to the database. Generally, speaking the deletion also follows a similar reasoning except the fact that the two operations are opposite.

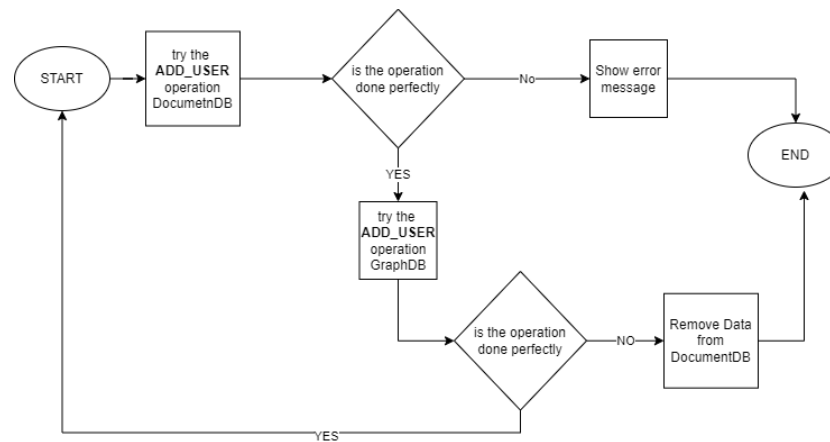


Figure 9: Cross Database Coherence for ADD USER operation

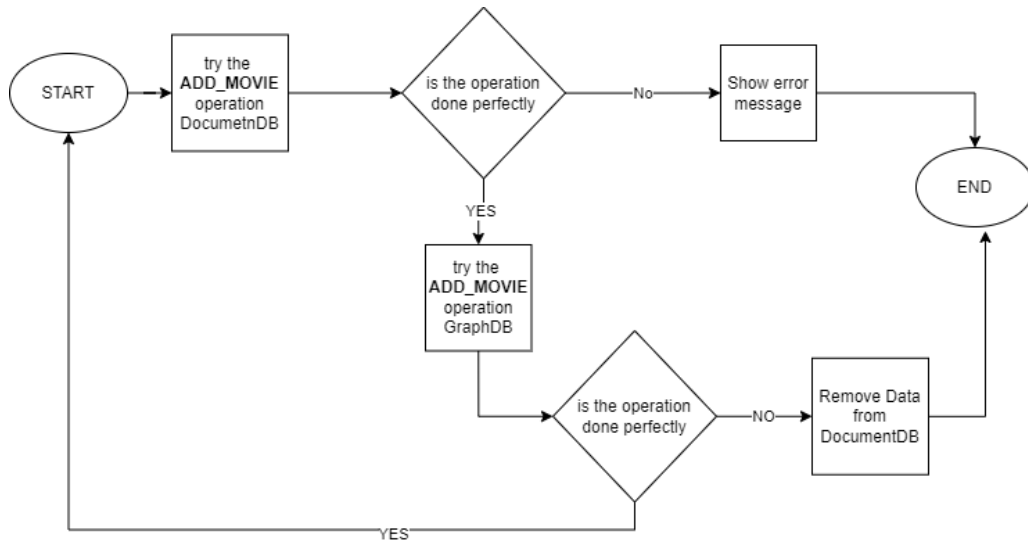


Figure 10: Cross Database Coherence for ADD MOVIE operation

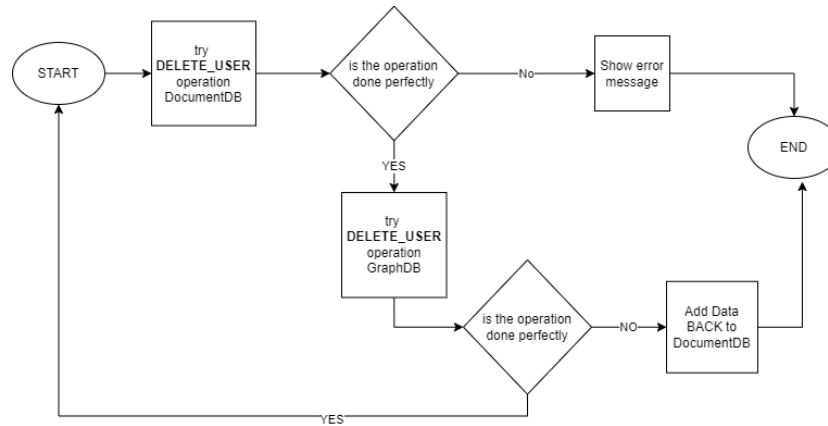


Figure 11: Cross Database Coherence for DELETE USER operation

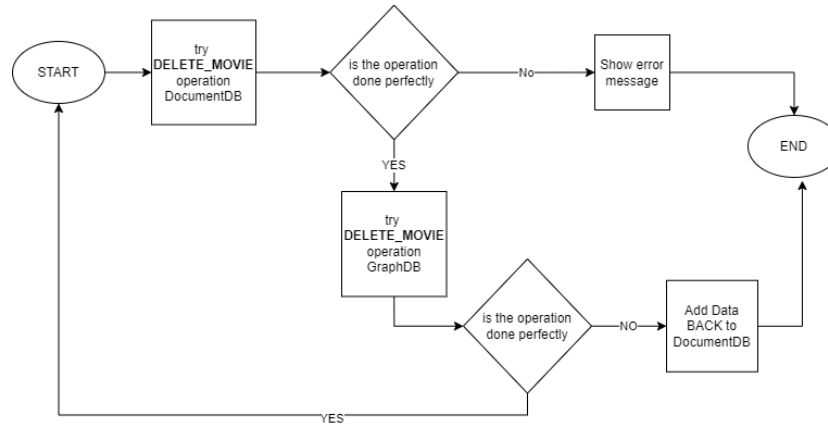


Figure 12: Cross Database Coherence for DELETE MOVIE operation

3 Sharding

A further improvement on the availability could be achieved by deploying the document database on a sharded architecture. A possible sharding key capable of distribute in a relatively could be the `_id` field inserted by the DMBS itself. This field, in fact, is computed by means of an hash algorithm on the document itself. Hash algorithms are able per design to produce random values uniformly distributed on a quite wide spectrum of values(given by the digest length itself).

4 Implementation

4.0.1 MongoDB Query Implementations And Analysis

The document part of our application is implemented using MongoDB, a source-available cross-platform document-oriented database program. Classified as a NoSQL database program, MongoDB uses JSON-like documents with optional schemas. [1].

Read Operations		
Operation	Expected Frequency	Cost
View movie information with reviews	Very High	Low(One Read)
Get popular movies	High	High (Aggregation)
Get popular movies by genre	Average	High (Aggregation)
Get top movies by year	Average	High (Aggregation)
Get top movies by year and genre	Average	Very high (Complex Aggregation)
Get top movies by year	Average	High (Aggregation)
Get the most popular genre per year	Average	Very High (Complex Aggregation)
Movie search with dynamic filtering	High	Average (multiple Reads)
Get the reviews of a movie	Very High	Low (One Read)
Get "more" reviews of a movie	Average	High (Aggregation)
Get the total reviews added per year	Average	Very High(Complex Aggregation)

Table 1: Read Operations on MongoDB

Write Operations		
Operation	Expected Frequency	Cost
Add user	High	Average (Document add)
Edit user	Average	Average (Document update)
Delete user	Average	Average (Document remove)
Add movie	Low	Average (Document add)
Delete movie	Low	Document remove + Documents (reviews) remove
Add review	High	Document add + Document update
Edit review	Average	Document update + Document update
Delete review	Average	Document Delete + Document update
Update a movie's vote information	High	Low(1 or 2 attribute writes)

Table 2: Write Operations on MongoDB

From this Query analysis we can say that the application is equally Read-Write heavy, because of the complex aggregations in the read operations, and also the 2-part atomic write operations related to reviews, as they are duplicated in two collections. **Read Operations:** Here is a selection of Read operation implementations in MongoDB:

GET POPULAR MOVIES

```
1 db.collection("movies").aggregate([
2   { $sort: { popularity: -1 }},
3   { $project: {
4     title: 1, poster_path: 1, overview: 1, vote_average: 1,
5     release_date: 1}}
6 ]).skip(skipped).limit(28).toArray();
```

With this aggregation we get the popular movies by sorting based on the "popularity" attribute and projecting only the fields used in the front-end.

GET TOP MOVIES BY YEAR AND GENRE

```
1 db.collection("movies").aggregate([
2   { $match: {
3     $and: [
4       { release_date: { $gte: reqYear, $lt: yearLimit }},
5       { genres: { $regex: genre, $options: "i" }}
6     ]}},
7   { $sort: { vote_average: -1 } },
8   { $project: {
9     title: 1, poster_path: 1, vote_average: 1, release_date: 1}},
10  ]).skip(skipped).limit(30).toArray();
```

This aggregation takes a genre and year as inputs and finds the movies that match, then gets the top movies by sorting based on the "average vote" attribute and projecting only the fields used in the front-end.

GET THE MOST POPULAR GENRE PER YEAR

```
1 db.collection("movies").aggregate([
2   { $unwind: { path: "$genres" }},
3   { $project: {
4     year: { $year: { $toDate: "$release_date" }}, genres: 1,
5     popularity: 1}},
6   { $group: {
7     _id: { genre: "$genres", year: "$year" },
8     average_popularity: { $avg: "$popularity" }},
9   { $group: {
10    _id: "$_id.year",
11    most_popular_genre: {
```

```

11         $max: "$average_popularity",
12     },
13     data: { $push: "$$ROOT" },}},
14     { $sort: { "data._id": -1 } },
15 ]).toArray();

```

To find the most popular genres in each year, there needs to be a grouping based on the genres and year, but given as our data for the movies contained an array of "genres" and not individual genre Strings, we need to unwind the array and create a document with the same movie but for each genre in the array. Then we project the release data to only retrieve the year, and genre and popularity. After that we start grouping by year and genre so that we can calculate the average popularity of that genre in that year. Then we need to do a comparison between the genres of each year to find the most popular one.

Since "group" aggregations don't allow to keep data other than the accumulator, we needed to store the members of the new year-based groups, so we use "push: ROOT" and end up with an array of objects containing the previously generated documents.

Since the final object structure is complex, it undergoes manipulation in the function before returning it.

GET "MORE" REVIEWS OF A MOVIE

```

1 db.collection("reviews").aggregate([
2     { $match: { movieId: movieId } },
3     {$project: {
4         review_date: { $toDate: "$review_date" },movieId: 1,title: 1,
5         reviewer: 1,rating: 1,review_detail: 1,review_summary: 1,userId: 1,
6         },
7     { $sort: { date: -1 } },
8     ]).skip(20).toArray();

```

Normally, retrieving a movie's review is done in a simple find(), since we keep 20 most recent review embedded in the movie document. However, a user may request to load more reviews and therefore, we need to fetch the "next"

20 most recent reviews from the reviews collection.

We match by movie ID in the reviews collection and project the fields we need. Then we sort the documents by date and skip the first 20 as they are already embedded.

DYNAMIC FILTERING FOR MOVIES

```
1   var query = { $and: [] };
2
3   if (title) { query.$and.push({ title: { $regex: title, $options: "i" }
4     });}
5   if (genre) { query.$and.push({ genres: { $regex: genre, $options: "i"
6     });}
7   if (minRuntime) { query.$and.push({ runtime: { $gte: minRuntime } });}
8   if (maxRuntime) { query.$and.push({ runtime: { $lte: maxRuntime } });}
9   if (minDate) { query.$and.push({ release_date: { $gte: new Date(
10     minDate) } });}
11  if (maxDate) { query.$and.push({ release_date: { $lte: new Date(
12     maxDate) } });}
13  if (language) { query.$and.push({spoken_languages: { $regex: language,
14     $options: "i" },,});}
15
16  let db = await mongoDriver.mongo();
17  const movies = await db.collection("movies").find(query, {
18    projection: { title: 1, poster_path: 1, vote_average: 1, overview:
19      1 },
20  }).toArray();
```

This query allows us to filter the search for movies based on the content of a form input sent by the front end, where the user can use many combinations of filters. We check if each attribute is not empty, then push it into the find() query, then send back an array of matching movies.

Write Operations: Here is a selection of Write operation implementations in MongoDB. In this section we will only include Creating a user, Deleting a movie, and Creating a review for the sake of conciseness.

CREATE USER/REGISTER: ANONYMOUS USERS

We start the query by checking if the user exists in the users collection, with

the assumption that if it is found in MongoDB it will also be in Neo4j.

```
1 // Check if user exists already by username and email
2 if (exists) {
3   return res.status(409).json({ message: "User already exist." });
4 } else {
5   // Hash password
6   // Create user document on Mongo
7   // Calculate a new _id for the user
8   // Create user object
9   try {
10    var newUser = {
11      ...set the attributes
12    };
13
14    //Consistency management
15    try {
16      await db.collection("users").insertOne(newUser);
17    } catch (e) {
18      return res.status(500).json({ message: "User Document not
19      Created" });
20    }
21    //Create the user node in Neo4j
22    const node = await userController.createUser(username, parseInt(
23    newId));
24
25    //No error was thrown
26    return res
27      .status(201)
28      .json({ User_Document: newUser, User_Node: node });
29    } catch (err) {
30      //Delete the user that was only inserted in MongoDB and not Neo4j
31      await userController.backtrackDelete(newId);
32      return res.status(500).json({ message: err.message });
33    }
34  }
35 }
```

DELETE A MOVIE: ADMIN PRIVILEGES

```
1 try {
2   let db = await mongoDriver.mongo();
3   //Store the movie to be deleted for backtracking
4   movie_to_be_deleted = await db.collection("movies").findOne({ _id:
5   movieId });
6 }
```

```

6      try {
7          await db.collection("movies").deleteOne({ _id: movieId });
8      } catch (err) {
9          return res.status(400).json({ message: "Mongo: Movie Deletion
Failed" });
10     }
11     //Delete the movie node in Neo4j
12     await deleteMovie(movieId);
13     //No errors were caught => delete the movie reviews from the
review collection
14     await reviewController.deleteAllMovieReviews(movieId);
15     res.status(200).json({ message: "Task executed successfully" });
16 } catch (err) {
17     try {
18         //If any of the previous steps fail, reinsert the movie document
in Mongo DB
19         await backtrackInsert(movie_to_be_deleted);
20     } catch (err) {
21         res.status(400).json({ message: err.message });}}

```

CREATE A REVIEW: REGISTERED USERS

```

1      // Calculate new ID for the review
2      newReview = {
3          ... set attributes of the review object};
4
5      let movie = await db.collection("movies").findOne({ _id: movieId });
6
7      // If the embedded array is at max length => push review at the top
and shift(delete from array copy) the least recent review
8      if (movie["reviews"].length >= 20) {
9          new_reviews = movie["reviews"]
10         new_reviews.push(newReview)
11         new_reviews.shift();
12         console.log(new_reviews)
13     } else {
14         new_reviews = movie["reviews"]
15         new_reviews.push(newReview);
16     }
17
18     // Insert the review in MongoDB
19     try {
20         await db.collection("reviews").insertOne(newReview);
21     } catch (e){

```

```

22     return res.status(400).json({ message: "Review Insertion
23     unsuccessful" });
24 }
25
26 // Calculate the new vote_average and vote_count of the movie
27 const new_vote = await updateMovieRatingCreation(movie, newReview["
28 rating"]);
29
30 try {
31     //update the movie node on Neo4j with the new average_vote
32     await updateVoteMovieNode(movieId, (new_vote.vote_average).toFixed
33     (1));
34     //No errors caught => update the movie document with the new
35     embedded array, the vote_count and the vote_average
36     await db.collection("movies").updateOne({_id: movieId}, {
37         $set: {
38             reviews: new_reviews,
39             vote_average: new_vote.vote_average,
40             vote_count: new_vote.vote_count}}});
41 } catch(e){
42     throw new Error ("Update movie info unsuccessful")
43 }
44 res.status(201).json({ review: newReview, message: "Review Added
45 successfully" });
46 } catch (err) {
47     //Delete the new review from the reviews collection and log error
48     let db = await mongoDriver.mongo();
49     await db.collection("reviews").deleteOne({_id: newReview._id});
50     res.status(400).json({ message: err.message });
51 }

```

The grapDB part is implemented using the neo4j grapDB. It is well known for its index-free adjacency technique which is cheaper and more efficient than doing the same task with foreign keys as in RDBMS or indexes, because query times are proportional to the amount of the graph searched, rather than increasing with the overall size of the data [2]. *Cypher*, which is a declarative language, is the de facto quering language for the neo4j GrapDB.

4.1 Main Queries in the GraphDB

In the neo4j database there are three entities/nodes : **User**, **Watchlist** and **Movie**

4.1.1 User nodes

Creating a user node. Upon the the time the user registers to the the system a user node is added to the graph database. Here is a code snippet from the neo4j Desktop.

```
1 match(u:User)
2 with u ORDER BY u.user_id DESC LIMIT 1
3 create(ux:User {user_id:u.user_id+1,username:"${username}"})
4 return {user_id:ux.user_id,username:ux.username}‘
```

Delete a user node: When a user node is removed from a database all his/her watchlists will also be removed from the Graph database.

```
1 MATCH(u:User{user_id:${user_id}})-[:CREATE]->(w:Watchlist)
2 DETACH DELETE u,w
```

Listing 6: Deleting a user node

Adding a follow relation: As we are going to build a movie review social application the user has the all the capabilities to follow other users in the social network. When the user wants to do so the FOLLOW relation is created as below. In our application a user can not follow him/her self. In the following listing a sample code that depicts the user-follow-user relation in the GraphDB.

```
1 match(u:User{user_id:${userId}})
2 match(ux:User{user_id:${followedUserId}})
3 where u.user_id<>ux.user_id and NOT (u)-[:FOLLOW]-(ux)
4 create (u)-[r:FOLLOW]->(ux)
5 return 'relationship created'
```

Listing 7: Creating User FOLLOW User relation

Deleting the follow relation: from above we have seen that a user can follow an other user, on the other hand a user can decide to stop following an other user. In line number 3 of the listing 7 before a relation is created the, a boolean check is done if there is a follows relation between those users.

```

1 MATCH(u:User{user_id:${user_id}}-[f:FOLLOW]->(ux:User{user_id:${user_id}}))
2 WHERE u.user_id<>ux.user_id
3 DELETE f
4 RETURN 'follows relationship deleted'

```

Listing 8: Deleting a user FOLLOW user relation

4.1.2 Movies Nodes

Create movie node: The administrator of the system can add a movie node to the database. The below snippet of code creates a movie node in the noe4j database.

```

1 MERGE (m:Movie{
2     movie_id:"${movie_id}",
3     title:"${title}",
4     poster_path:"${poster_path}",
5     vote_average:${vote_average},
6     overview:"${overview}")
7 return m

```

Listing 9: Creating a movie node

Delete Movie Node: When a movie is deleted from the watchlist it belongs the relation that connect the movie with the watchlist and the node itself will be deleted from the database. Here is a snippet of a code that depicts the deletion of a movie node from the GraphDB.

```

1 MATCH(m:Movie{movie_id:"${movie_id}"})
2 DETACH DELETE m

```

Listing 10: Delete Movie node

4.1.3 Watchlist nodes

Create Watchlist

```
1 MATCH (u:User{user_id:${user_id}})
2 CREATE(w:Watchlist{name:"${name}",created_date:date()})
3 CREATE (u)-[:CREATE]->(w)
4 return {created_by:u.username,user_id:u.user_id,watchlist_name:w.name,id:
  ID(w),created_date:w.created_date}
```

Listing 11: Creating a watchlist node

Update Watchlist

```
1 match(w:Watchlist)<-[:CREATE]-(u:User)
2   where ID(w)=${watchlist_ID} AND u.user_id=${user_id}
3   SET w.name="${name}"
4   return {name:w.name,id:ID(w),created_date:w.created_date}
```

Listing 12: Updating a watchlist node

Delete watchlist

```
1 Match(w:Watchlist)<-[:CREATE]-(u:User)
2   where ID(w)=${watchlist_ID} and (u.user_id=${user_id} or true)
3   detach delete w
```

4.1.4 Recommendations/Suggestions in the GraphDB

Recommended Users

```
1 MATCH (u1:User{user_id:${userId}})-[f1:FOLLOW]->(w1:Watchlist)<-[f2:
  FOLLOW]-(ux:User)
2   where NOT (u1)-[:FOLLOW]->(ux)
3   return {user_name:ux.username,user_id:ux.user_id} LIMIT 7
```

Listing 13: Recommended Users

```
1 try {
```

```

2   let session = neo4jdbconnection.session();
3   const suggestedMovies = await session.run(
4     'MATCH(u:User{user_id:${userId}})-[:FOLLOW]->(ux:User)
5     with *
6     optional match(ux)-[:CREATE]->(w:Watchlist)<-[:BELONGS]-(m:Movie)
7     with m,count(m) as num
8     where num>=2
9     return {title:m.title,poster_path:m.poster_path,overview:m.overview,
vote_average:m.vote_average,id:m.movie_id}
10    limit 7'
11  );
12  session.close();
13  if (!suggestedMovies.records)
14    return res.status(400).json({ message: 'an error occurred' });
15  const result = suggestedMovies.records.map((e) => {
16    return {
17      ...e["_fields"][0],
18      vote_average: e["_fields"][0]["vote_average"],
19    };
20  });
21
22  res.status(200).json(result);
23 } catch (err) {
24   res.status(500).json({ message: err.message });
25 }

```

Listing 14: Recommended Users Express JS Method

Recommended Movies

```

1  MATCH(u:User{user_id:${userId}})-[:FOLLOW]->(ux:User)
2    with *
3    optional match(ux)-[:CREATE]->(w:Watchlist)<-[:BELONGS]-(m:Movie)
4    with m,count(m) as num
5    where num>=2
6    return {title:m.title,poster_path:m.poster_path,overview:m.overview,
vote_average:m.vote_average,id:m.movie_id}
7    limit 7

```

Listing 15: Recommended Movies

```

1  try {
2    let session = neo4jdbconnection.session();
3    const suggestedMovies = await session.run(

```

```

4      'MATCH(u:User{user_id:${userId}})-[:FOLLOW]->(ux:User)
5      with *
6      optional match(ux)-[:CREATE]->(w:Watchlist)<-[:BELONGS]-(m:Movie)
7      with m,count(m) as num
8      where num>=2
9      return {title:m.title,poster_path:m.poster_path,overview:m.overview,
vote_average:m.vote_average,id:m.movie_id}
10     limit 7'
11 );
12 session.close();
13 if (!suggestedMovies.records)
14     return res.status(400).json({ message: 'an error occurred' });
15 const result = suggestedMovies.records.map((e) => {
16     return {
17         ...e["_fields"][0],
18         vote_average: e["_fields"][0]["vote_average"],
19     };
20 });
21
22 res.status(200).json(result);
23 } catch (err) {
24     res.status(500).json({ message: err.message });
25 }

```

Listing 16: Recommended Movies Express JS Method

View suggested Watchlists that are from the users you follow.

```

1 \begin{lstlisting}[caption=Recommended Watchlists,]
2 MATCH (u1:User{user_id:${userId}})-[f1:FOLLOW]->(ux:User)-[f1:FOLLOW]->(
3     wl:Watchlist)
4     where NOT (u1)-[:FOLLOW]->(wl) AND NOT (u1)-[:CREATE]->(wl)
5     with wl
6     match (ux:User)-[:CREATE]->(wl)
7     return {name:wl.name,id:ID(wl),owner:ux.username,owner_id:ux.user_id,
created_date:wl.created_date}

1 try {
2     let session = neo4jdbconnection.session();
3     const suggestedWatchlists = await session.run('
4     MATCH (u1:User{user_id:${userId}})-[f1:FOLLOW]->(ux:User)-[f1:FOLLOW
5     ]->(wl:Watchlist)
6     where NOT (u1)-[:FOLLOW]->(wl) AND NOT (u1)-[:CREATE]->(wl)

```

```

6   with wl
7   match (ux:User)-[:CREATE]->(wl)
8   return {name:wl.name,id:ID(wl),owner:ux.username,owner_id:ux.user_id,
9   created_date:wl.created_date}');
10  session.close();
11  if (!suggestedWatchlists.records)
12    return res.status(400).json({ message: 'an error occurred' });
13  const result = suggestedWatchlists.records.map((e) => {
14    return {
15      ...e["_fields"][0],
16      id: e["_fields"][0].id.low,
17      owner_id: e["_fields"][0].owner_id.low,
18    };
19  });
20  res.status(200).json(result);
21 } catch (err) {
22   res.status(500).json({ message: err.message });
23 }

```

Listing 17: Recommended Watchlists Express JS Method

View suggested user profiles that follow the same Watchlists you do.

```

1  //View suggested user profiles that follow the same Watch Lists you do.
2  //parameters:user_id
3  MATCH (u1:User{user_id:${user_id}})-[f1:FOLLOW]->(wl:Watchlist)<-[f2:
4    FOLLOW]-(ux:User)
5  where NOT (u1)-[:FOLLOW]->(ux)
6  return {user_name:ux.username,user_id:ux.user_id} LIMIT 7

```

View suggested movies based on their occurrence on the watchlists of users you follow.

```

1  //View suggested movies based on their occurrence on the watchlists of
2  users you follow.
3  // parameters : user_id
4  MATCH(u:User{user_id:${user_id}})-[:FOLLOW]->(ux:User)
5  with *
6  optional match(ux)-[:CREATE]->(w:Watchlist)<-[:BELONGS]-(m:Movie)
7  with m,count(m) as num
8  where num>=2

```

```

8 return {title:m.title,poster_path:m.poster_path,overview:m.overview,
          vote_average:m.vote_average,id:m.movie_id}
9 limit 7
10

```

Retrieve a user's own watchlist.

```

1  try {
2    let session = neo4jdbconnection.session();
3    const suggestedUsers = await session.run(
4      'MATCH (u1:User{user_id:${userId}})-[f1:FOLLOW]->(w1:Watchlist)<-[f2:
      FOLLOW]-(ux:User)
5    where NOT (u1)-[:FOLLOW]->(ux)
6    return {user_name:ux.username,user_id:ux.user_id} LIMIT 7'
7    );
8    if (!suggestedUsers.records)
9      return res.status(400).json({ message: 'an error occurred' });
10   const result = suggestedUsers.records.map((e) => {
11     return {
12       ...e["_fields"][0],
13       user_id: e["_fields"][0].user_id.low,
14     };
15   });
16
17   res.status(200).json(result);
18 } catch (err) {
19   res.status(500).json({ message: err.message });
20 }

```

Find the most active users based on the number of followers of their watchlists.

```

1  MATCH (u1:User)-[f1]->(w:Watchlist)<-[f2]-(u2:User)
2  with *
3  WHERE
4  type (f1) in ["CREATE"]AND
5  type (f2) in ["FOLLOW"]AND
6  NOT (u1)-[:FOLLOW]-(w) AND
7  NOT (u2)-[:CREATE]-(w)
8  RETURN u1.user_name AS YOU, w.watchlist_name AS TitleOfYourWatchlist,
9  COLLECT(u2.user_name)AS WatchlistFollowers, count(f2) AS NumFollowers
10 ORDER BY NumFollowers

```

```

11 DESC
12 LIMIT 10

```

4.1.5 Summary of Queries analysis for the GraphDB

Read Operations		
Operation	Expected Frequency	Cost
Most followed Watchlists	Medium	Medium (Multiple reads)
Most Active Users based on their number of followers of their watchlist	Low	High
Suggested Movies based on their occurrence on the watchlists of users you follow	Low	High
Suggested Watchlists from the users that you follow	Average	Medium
Suggested Users	Average	Medium (Multiple reads)
Retrieve a user's own watchlists	High	Low

Table 3: Read Operations in the GraphDB

Write Operations		
Operation	Expected Frequency	Cost
Insert/update/Remove a Watchlist	Low	Low
Update User information	Low	Low

Table 4: Write Operations in the GraphDB

References

- [1] <https://www.mongodb.com/docs/manual/introduction/>, “Mongo db documentation,” 2022.

- [2] N. I. Joy Chao, “Getting started guide for neo4j,” April, 2022.