

Big Data Assignment 1

CRGMAT002 GRDDAN017 KFWJOR001 MRCGAB004 WHLJOS001

March 1, 2024



Contents

1	Find or Create a Suitable Data Set	3
1.1	Data Set Explanation	3
1.2	Data Pre-Processing	3
2	Design a MongoDB Database	5
2.1	Collection 1 - books	5
2.2	Collection 2 - users	8
2.3	Explanation and Justification	11
2.3.1	Books	11
2.3.2	Users	11
3	Create and Load This MongoDB Database	13
3.1	Load the Database	13
3.1.1	Books Collection	13
3.1.2	Users Collection	13
3.2	Testing	13
4	Discuss the Relative Benefits and Disadvantages of MongoDB	15
4.1	MongoDB (Document Store)	15
4.2	Graph DB	15
4.3	Key-value Store	16
4.4	Relational	16
4.5	Column Family	17
4.6	Hierarchical	17
5	Query and Updating the Database	19
5.1	GRDDAN017	19
5.1.1	1	19
5.1.2	2	20
5.1.3	3	22
5.1.4	4	23
5.2	MRCGAB004	24
5.2.1	1	24
5.2.2	2	25
5.2.3	3	26
5.2.4	4	27
5.3	CRGMAT002	28
5.3.1	1	28
5.3.2	2	30
5.3.3	3	32
5.3.4	4	33
5.4	WHLJOS001	34
5.4.1	1	34

5.4.2	2	35
5.4.3	3	36
5.4.4	4	37
5.5	KFWJOR001	38
5.5.1	1	38
5.5.2	2	40
5.5.3	3	41
5.5.4	4	43
6	Link the Database to a Program	47
7	Contribution Statement	51

1 Find or Create a Suitable Data Set

1.1 Data Set Explanation

Link to the dataset: <https://github.com/zygmuntz/goodbooks-10k>

The dataset includes data from an online book review platform (<https://goodreads.com/>). It includes information about the books, users, book tags, and book rating scores. The dataset initially contained multiple csv files. This dataset was chosen as its ideal for a MongoDB database due to its semi-structured nature and nested data, which is particularly useful for storing ratings and book tags.

Dataset Content:

- **books.csv**: Each entry represents a book with a unique `book_id`. There are multiple data fields for a book:
 - **book_id, goodreads_book_id, best_book_id, work_id**: Unique id's representing a book, each with a different purpose. We only used `book_id` and `goodreads_book_id` as they're used to link books to user `ratings` and user `to_read` lists.
 - **ratings_1, ratings_2, ...**: Number of user ratings by rating value. eg. `ratings_1` represents the number of 1-star ratings given to that book.
 - The rest of the fields are self-explanatory but include info relating to authors, title, release date, and isbn number.
- **ratings.csv**: Each entry is a `user_id` to `book_id` mapping with a rating.
 - **book_tags.csv**: Each entry is a `book_id` to `tag_id` mapping.
 - **tags.csv**: Each entry is a `tag_id` to `tag_name` mapping.
 - **to_read.csv** : Each entry is a `user_id` to `goodreads_book_id` mapping which represents a user adding a book to their `to_read` list.

1.2 Data Pre-Processing

The data was processed such that the data was represented in JSON format with evidence of nested objects so that we could demonstrate the capabilities of MongoDB.

Here is a quick outline of how we processed the data to create JSON files:

Libraries used: `Pandas`, `PyArrow`, `Faker`

`Pandas` was used to load the csv files into dataframes where we merged data and applied `group by` aggregate functions to obtain lists of data objects per a unique entry id. This was useful, for example, when we obtained a list of tags per `book_id`.

`Faker` was used to generate random usernames for each id that were then written to `user_data.csv`. The dataframes were then converted into JSON files.

All data pre-processing code is in the data-processing directory but the output JSON files are included in the final submission.

2 Design a MongoDB Database

Both Collection Schemas were designed by creating hand-made JSON example objects. Each of these objects shows what a document in the DB would look like. Underneath each JSON example, we have included a diagram which represents the example's nesting visually.

2.1 Collection 1 - books

JSON example

```
{
  "book_id": "98",
  "isbn": "1401359604",
  "isbn13": "9781401359610.0",
  "authors": [
    "Plum Sykes"
  ],
  "original_publication_year": 2004,
  "title": "Bergdorf Blondes",
  "language_code": "en-US",
  "average_rating": 3.26,
  "total_ratings": 23795,
  "rating_counts": {
    "ratings_1": 2020,
    "ratings_2": 4428,
    "ratings_3": 8669,
    "ratings_4": 6144,
    "ratings_5": 4561
  },
  "image_url": "https://s.gr-assets.com/assets/nophoto/book/111x148-bcc042a9c91a29c1d680899eff700a03.png",
  "tags": [
    {
      "tag_id": 11743,
      "tag_name": "fiction"
    },
    {
      "tag_id": 8717,
      "tag_name": "currently-reading"
    },
    {
      "tag_id": 8055,
      "tag_name": "contemporary"
    }
  ],
  "ratings": [
    {
      "user": {
        "user_id": 237,
        "user_name": "David Smith"
      },
      "rating": 1
    },
  ],
}
```

```
{
  "user": {
    "user_id": 364,
    "user_name": "Christina Calderon"
  },
  "rating": 1
},
{
  "user": {
    "user_id": 399,
    "user_name": "Stephen Pugh"
  },
  "rating": 2
}
]
```

book_id: "98"
isbn: "1401359604"
isbn13: "9781401359610.0"
original_publication_year: 2004
title: "Bergdorf Blondes"
language_code: "en-US"
average_rating: 3.26
total_ratings: 23795
image_url: "https://s.gr-assets.com/assets/nophoto/book/111x148-bcc042a9c91a29c1d680899eff7..."

authors (1)

Plum Sykes

rating_counts (1)

ratings_1: 2020
ratings_2: 4428
ratings_3: 8669
ratings_4: 6144
ratings_5: 4561

tags (3)

tag_id: 11743
tag_name: "fiction"

tag_id: 8717
tag_name: "currently-reading"

tag_id: 8055
tag_name: "contemporary"

ratings (6)

user (1)
user_id: 237
user_name: "David Smith"

rating: 1

user (1)
user_id: 364
user_name: "Christina Calderon"

rating: 1

user (1)
user_id: 399
user_name: "Stephen Pugh"

rating: 2

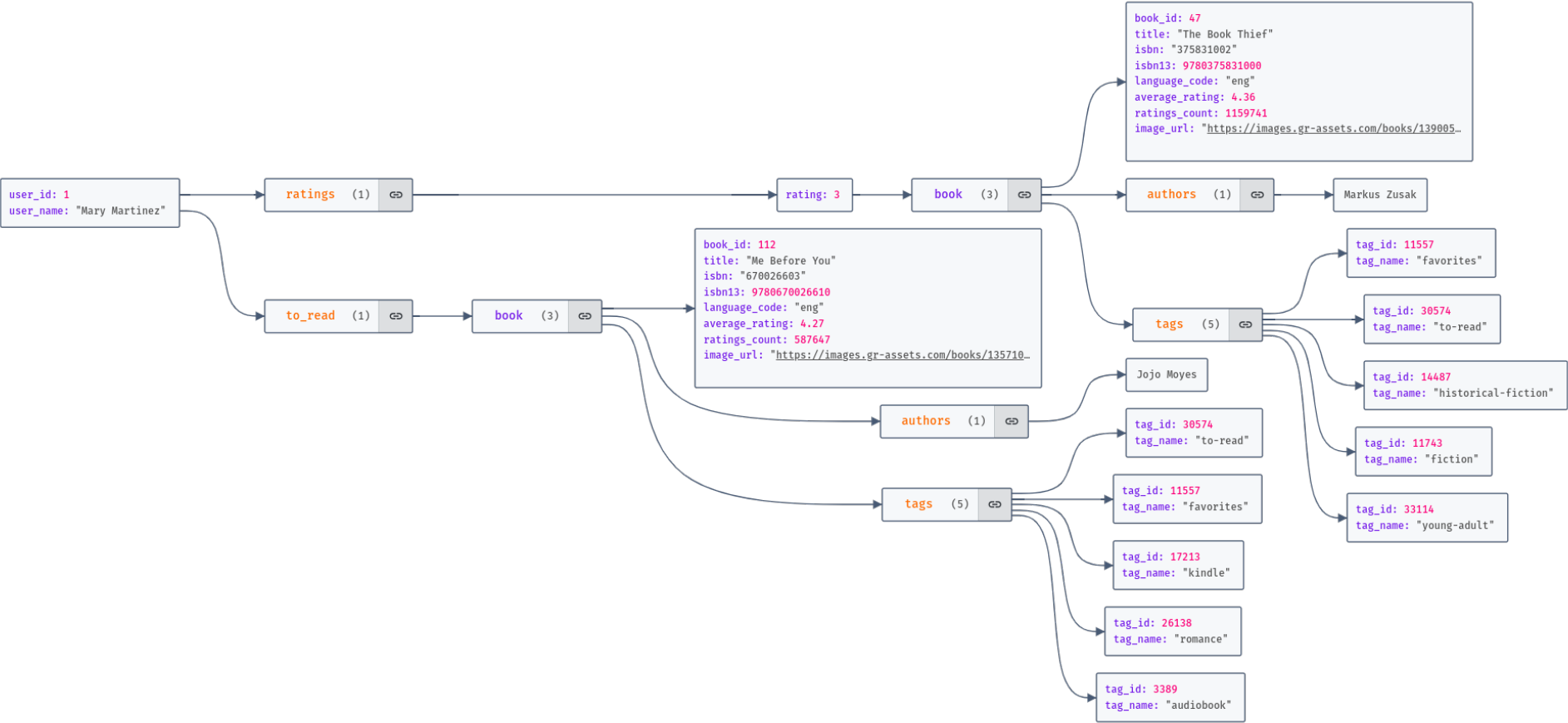
2.2 Collection 2 - users

```
{
  "user_id": 1,
  "user_name": "Mary Martinez",
  "ratings": [
    {
      "book": {
        "book_id": 47,
        "authors": [
          "Markus Zusak"
        ],
        "title": "The Book Thief",
        "isbn": "375831002",
        "isbn13": 9780375831000.0,
        "language_code": "eng",
        "average_rating": 4.36,
        "ratings_count": 1159741,
        "image_url": "https://images.gr-assets.com/books/1390053681m/19063.jpg",
        "tags": [
          {
            "tag_id": 11557,
            "tag_name": "favorites"
          },
          {
            "tag_id": 30574,
            "tag_name": "to-read"
          },
          {
            "tag_id": 14487,
            "tag_name": "historical-fiction"
          },
          {
            "tag_id": 11743,
            "tag_name": "fiction"
          },
          {
            "tag_id": 33114,
            "tag_name": "young-adult"
          }
        ]
      },
      "rating": 3
    }
  ],
  "to_read": [
```

```

{
  "book": {
    "book_id": 112,
    "authors": [
      "Jojo Moyes"
    ],
    "title": "Me Before You",
    "isbn": "670026603",
    "isbn13": 9780670026610.0,
    "language_code": "eng",
    "average_rating": 4.27,
    "ratings_count": 587647,
    "image_url": "https://images.gr-assets.com/books/1357108762m/15507958.jpg",
    "tags": [
      {
        "tag_id": 30574,
        "tag_name": "to-read"
      },
      {
        "tag_id": 11557,
        "tag_name": "favorites"
      },
      {
        "tag_id": 17213,
        "tag_name": "kindle"
      },
      {
        "tag_id": 26138,
        "tag_name": "romance"
      },
      {
        "tag_id": 3389,
        "tag_name": "audiobook"
      }
    ]
  }
}
]
}

```



2.3 Explanation and Justification

The data in CSV format emulates the functionality of a relational database. Many of the fields have foreign keys that point to elements in the other CSV files. Document store database favour efficiency over consistency, thus, we have nested a copy of the relevant object where the value would have otherwise been a foreign key. The process by how this was achieved was highlighted in “[Data Pre-Processing](#)”.

The data was seeded into the following 2 collections:

- books
- users

2.3.1 Books

The books collection roughly followed the format of books.csv (outlined in “[Data Set Explanation](#)”) with some modifications. We omitted unnecessary information and altered the names of some of the properties to make their semantic meaning clearer. As a result each book document includes:

- General information about the book (title, author, original_publication_year, etc.).
- Aggregated rating values (average_rating, total_ratings, rating_counts).
- A ratings list
 - Each element represents a user’s review of that book. It includes the user’s basic information and the assigned rating score.
- A tags list. Each tag represents a genre or category the book belongs to.

The most significant element of our design was the aforementioned nesting. Rather than store a separate tags collection, all the tags associated with the book are stored as a list of objects. The same is true for ratings, which is a list of rating objects.

Use case examples:

- Collection of books: Querying to see the average rating of a book and the distribution of ratings of a particular book

2.3.2 Users

Unlike books, the users collection does not directly correspond to a csv file. Instead, users was created by combining data from ratings.csv, to_read.csv, and books.csv. We chose to create the user collection in order to demonstrate the importance of collection design with regards to query efficiency. While the users contains a lot of duplicate data from books, it does so in a way that places information about the users at the top of the nesting hierarchy. This means that data about individual users can be obtained without performing expensive joins. The chosen collection design allows one to access, store and perform analytics from the perspective of the user. Each document in the users collection includes the following:

- A `user_id` and `user_name` (randomly generated, as explained in “[Data Pre-Processing](#)”).
- A `ratings` list:
 - Each element represents a score that the user has given to a book.
 - The element includes a `book` object and a given `rating` score.
 - The nested `book` object includes all high-level data about that book, as would be found in the `books` collection.
- A `to_read` list.
 - Each element represents a book that the user has added to their `to_read` list (ie. plans to read that book).
 - The element is represented as a `book` object that is identical in structure to those represented in the `ratings` list.

Use case examples:

Collection of users: what books user X wants to read.

3 Create and Load This MongoDB Database

3.1 Load the Database

The process of creating and loading the database was significantly simplified on account of the fact that we had already [Pre-Processed](#) the data into JSON format. For each collection, a shell command is called to seed the JSON data into their associated collections. The creation of the database is implicit ie. inserting data into the non-existent database leads to its creation. The commands use `mongoimport`, a CLI tool designed for extracting data from plain-text formats (eg. JSON) and inserting them into MongoDB databases.

3.1.1 Books Collection

```
mongoimport --db goodbooks --collection books --file mongo-seed/books.json --jsonArray
```

3.1.2 Users Collection

```
mongoimport --db goodbooks --collection users --file mongo-seed/users.json --jsonArray
```

3.2 Testing

The shell commands were tested and the output has been annotated on the following page:

```
8DAssignment1 on / main [ $? ] via @ v18.17.1 via v3.11.5
} mongosh
Current Mongosh Log ID: 65e02d3d4ff0cb22198e19a8
Connecting to: mongod://127.0.0.1:27017/?directConnection=true&serverSelectionTimeoutMS=2000&appName=mongosh+2.1.5
Using MongoDB: 7.0.5
Using Mongosh: 2.1.5

For mongosh info see: https://docs.mongodb.com/mongod-shell/

-----
The server generated these startup warnings when booting
2024-02-28T15:34:02.442+00:00: Access control is not enabled for the database. Read and write access to data and configuration is unrestricted
2024-02-28T15:34:02.442+00:00: vm.max_map_count is too low
-----

test> show dbs
admin      40.00 KiB
config     108.00 KiB
local      72.00 KiB
test>

8DAssignment1 on / main [ $? ] via @ v18.17.1 via v3.11.5 took 28s
} mongoimport --db bookstore --collection users --file mongo-seed/users.json --jsonArray
} mongoimport --db bookstore --collection books --file mongo-seed/books.json --jsonArray
2024-02-29T09:08:31.249+0200 connected to: mongod://localhost/
2024-02-29T09:08:31.589+0200 255 document(s) imported successfully. 0 document(s) failed to import.
2024-02-29T09:08:31.602+0200 connected to: mongod://localhost/
2024-02-29T09:08:31.824+0200 2184 document(s) imported successfully. 0 document(s) failed to import.

8DAssignment1 on / main [ $? ] via @ v18.17.1 via v3.11.5
} mongosh
Current Mongosh Log ID: 65e02d938dc73afd1fd4d2df
Connecting to: mongod://127.0.0.1:27017/?directConnection=true&serverSelectionTimeoutMS=2000&appName=mongosh+2.1.5
Using MongoDB: 7.0.5
Using Mongosh: 2.1.5

For mongosh info see: https://docs.mongodb.com/mongod-shell/

-----
The server generated these startup warnings when booting
2024-02-28T15:34:02.442+00:00: Access control is not enabled for the database. Read and write access to data and configuration is unrestricted
2024-02-28T15:34:02.442+00:00: vm.max_map_count is too low
-----

test> use bookstore
switched to db bookstore
bookstore>

8DAssignment1 on / main [ $? ] via @ v18.17.1 via v3.11.5 took 1m
} mongosh
Current Mongosh Log ID: 65e02dd18ea8219511814a69
Connecting to: mongod://127.0.0.1:27017/?directConnection=true&serverSelectionTimeoutMS=2000&appName=mongosh+2.1.5
Using MongoDB: 7.0.5
Using Mongosh: 2.1.5

For mongosh info see: https://docs.mongodb.com/mongod-shell/

-----
The server generated these startup warnings when booting
2024-02-28T15:34:02.442+00:00: Access control is not enabled for the database. Read and write access to data and configuration is unrestricted
2024-02-28T15:34:02.442+00:00: vm.max_map_count is too low
-----

test> show dbs
admin      40.00 KiB
bookstore   2.28 MiB
config     108.00 KiB
local      72.00 KiB
test> use bookstore
switched to db bookstore
bookstore> db.books.find( { tags: { $elemMatch: { tag_name: "mystery" } } }, { _id: 0, title: 1, authors: { $slice: 1 }, average_rating: 1 }).sort( { average_rating: -1 }).limit(5)
[
  {
    authors: [ 'Arthur Conan Doyle' ],
    title: 'SherLock Holmes: The Complete Novels and Stories, Volume I',
    average_rating: 4.45
  }
]
```

Show that db "bookstore" does not exist.

Import JSON

use bookstore

Bookstore is now 2.28MiB

Confirm with query.
Couldn't fit it in fully.

NB: the mongoDB Database was renamed from "bookstore" to "goodbooks".

4 Discuss the Relative Benefits and Disadvantages of MongoDB

4.1 MongoDB (Document Store)

MongoDB is a document store database. Document store databases favour low retrieval latency and efficiency over consistency and functionality (expressiveness) when compared with relational databases. MongoDB is flexible: it facilitates incremental schema modification as the schema is not rigidly defined. This can improve development velocity and account for changing stakeholder requirements.

MongoDB is highly scalable and will be efficient with the ever increasing addition of books being published. MongoDB meets a good middle-ground when compared with other NoSQL databases. It facilitates nested data and objects and allows querying on those objects without being overly strict or complex. Each Book and User is able to store all necessary information while allow for performant data retrieval without the overhead of complex joins.

The disadvantages of using a MongoDB are that there is a high degree of duplication and redundant data. This means that storing and transferring all of this data is computationally expensive. Secondly, inserting and updating data can be much more expensive than a relational database on account of the necessity to make updates in multiple places (due to the duplication). However, data about books is unlikely to change frequently (as the book has already been published) and thus this downside is heavily mitigated in the chosen example.

- **Benefits** (When compared with relational): Scalable, schema-free design, data-access performance
- **Disadvantages** (When compared with relational): Duplication, high storage use, slow modification, poor query expressiveness (functionality)
- **Role in polyglot persistence:**
 - A future polyglot design could store data that is high in volume but does not change frequently inside a document store db.
 - Books do not change frequently so the information about the books could be stored in the document store while users and ratings could be stored in other databases.
 - User telemetry could also be stored as in a document store database as it is high volume with very infrequent modification

4.2 Graph DB

Graph databases allow for expressive and performant relationship representation and would be ideal for relationship-related queries. Each relationship can be given properties that qualitatively describe that relationship. This is ideal for data where the relationships to other data is the primary provider of utility (eg. a knowledge graph), as opposed to the structured categorisation (as is the case in MongoDB and relational databases).

The chosen dataset is highly structured without a large number of annotated relationships between

data and it thus is a poor choice for graph DBs. The dataset would not leverage the benefits of a graph DB and would miss out on the well-fitted benefits of MongoDB, namely schema flexibility, performance, scalability, and simplicity.

- **Benefits:** Efficient modelling of highly inter-linked data, Query expressiveness (functionality), Great relationship representation (for which it is scalable and performant)
- **Disadvantages:** Unnecessary complexity, poor modelling of structured data
- **Role in polyglot persistence:**
 - A future polyglot design could include a graph database that models the relationship books have with each other.
 - Each book could have an outgoing edge to related books with a descriptive tag that describes their relationship strength and type.
 - This would prove highly useful for the site’s recommendation engine.

4.3 Key-value Store

Key-Value stores are extremely simple. They are very similar to a document store but with fewer capabilities in data representation and querying. They offer many of the same benefits of schema flexibility, scalability, and performance. Their key benefit and disadvantage over document stores is their simplicity. The simplicity is limiting but ideal for data that can be modelled simply. The chosen dataset is not suited for a key-value store as it includes more complex nested objects.

- **Benefits:** High performance, Scalability, Simplicity, Flexibility
- **Disadvantages:** Limited query capabilities and Does not support complex data structures.
- **Role in polyglot persistence:**
 - It can store the current session data of users. Easy to keep track of users’ sessions such as how long they were on the app.
 - The key-value store’s flexibility would be very important as it would allow us to store more complex session information about the user which can be used for processing later on to improve the user’s experience on the application.

4.4 Relational

Relational databases are great for keeping data consistent and for effectively maintaining the relationships of the data. Their ability to store relationships (via foreign keys) between data ensures the following benefits: data consistency, minimal duplication, and expressive queries. Relational databases falter when needing to retrieve large quantities of data from different tables. Joins become very expensive as the size of a dataset grows. Data analysis and operations are impacted substantially. It becomes inefficient with very large (big) data.

The goodbooks dataset at present would be well suited to a relational database, however, if we were to massively increase the scale of the data stored the efficiency gains of MongoDB would quickly become apparent. The data is not updated frequently enough for critical operations and thus the consistency benefits provided by a relational database would be negligible.

- **Benefits:** Query Expressiveness (functionality), consistency, reliability, low storage use (avoids duplication), accuracy, data integrity
- **Disadvantages:** Decrease in performance at scale, low flexibility (requires up-front schema design)
- **Role in polyglot persistence:**
 - A future polyglot design could include a relational database for data that would benefit from reliability and consistency.
 - Information about users, including their profiles, and security information could be stored in this format.
 - User data could benefit from the expressive queries provided by relational databases.

4.5 Column Family

Column Family DBs are similar to relational databases, however, they bundle groups of frequently accessed columns together. This reduces the performance issues that may occur in relational databases. It does, however, still require up-front schema design without obtaining the expressiveness benefits provided by relational databases.

At present, the `goodbooks` database is not well suited to being effectively grouped. The only instance in which this would be useful would be a bundling of ratings data in the books table. However, the benefits are not extensive enough to utilise this over MongoDB for this dataset.

- **Benefits:** Compression of data, Scalable, fast to load queries, good for simple analytics.
- **Disadvantages:** Limited querying capabilities, limited data modelling capabilities
- **Role in polyglot persistence:**
 - A future polyglot design could include a column family for calculating metrics and performance analytics, such as:
 - * How many books are users reading and the total number of stars people rated individual books
 - * How long people stay on the website, etc.
 - Each set of metrics relating to a particular type of user data could be grouped into a column family for that user.

4.6 Hierarchical

Hierarchical databases are optimised for data that can be stored as deeply nested hierarchies. Document Store DBs are also well suited for hierarchical nesting, however, their access performance greatly suffers when the depth of those hierarchies becomes excessive. Thus, hierarchical databases are great for data with inherently deep hierarchical structure. Hierarchical databases would be effective for fast lookup and to effectively store this data.

While the `goodbooks` dataset does include hierarchies, the hierarchies seldom extend beyond a depth of

5. This means the data is not inherently hierarchical and does not suffer from performance detriments when traversing the nested data. On account of this, a hierarchical database would be a poor choice as it would miss out on the key benefits of MongoDB (flexibility in particular).

- **Benefits:** Fast retrieval of nested data, efficient storage of data, predictable data structure
- **Disadvantages:** Limited flexibility, hard to maintain and update, limited interoperability, reduced expressiveness
- **Role in polyglot persistence:**
 - In the future, the platform may wish to hold a representation of books, their predecessors, and their successors.
 - A hierarchical relationship could be modelled such that each book belongs to a subgenre, which belongs to a series of parent genres. These books would have parent and descendant books that are influenced by them.
 - This hierarchy could be stored separately from the primary book database to avoid affecting its performance and size.

5 Query and Updating the Database

5.1 GRDDAN017

5.1.1 1

Description

Find the book title and author list of the top 3 highest rated books.

Query

```
db.books.find({}, {
  _id: false, book_id: true, title: true, authors: true, average_rating: true
}).sort({average_rating:-1}).limit(3);
```

Output

```
bookstore> db.books.find({}, {_id: false, book_id: true, title: true, authors: true, average_rating: true}).sort({average_rating:-1}).limit(3);
[
  {
    book_id: '25599',
    authors: [ 'Bill Watterson' ],
    title: 'The Complete Calvin and Hobbes',
    average_rating: 4.82
  },
  {
    book_id: '25601',
    authors: [ 'Bill Watterson' ],
    title: "It's a Magical World: A Calvin and Hobbes Collection",
    average_rating: 4.75
  },
  {
    book_id: '170846',
    authors: [ 'Bill Watterson' ],
    title: "There's Treasure Everywhere: A Calvin and Hobbes Collection",
    average_rating: 4.74
  }
]
bookstore>
```

5.1.2 2

Description

Find 2 books with the “fantasy” tag with a high average rating.

Query

```
db.books.find({
  $and: [{ "tags.tag_name": "fantasy" }, { "average_rating": { $gt: 4 } }]
}).limit(2)
```

Output

```

mongosh mongodb://127.0.0.1:27017/?directConnection=true&serverSelectionTimeoutMS=2000
bookstore> db.books.find({$and: [{ "tags.tag_name": "fantasy" }, { "average_rating": { $gt: 4 } }]}).limit(2)
[
  {
    _id: ObjectId('65dceab927a93ef44b0559e4'),
    book_id: '852',
    isbn: '786931582',
    isbn13: '9780786931580.0',
    authors: [ 'Margaret Weis', 'Tracy Hickman' ],
    original_publication_year: 1986,
    title: 'Time of the Twins',
    language_code: 'eng',
    average_rating: 4.13,
    total_ratings: 24671,
    rating_counts: {
      ratings_1: 142,
      ratings_2: 937,
      ratings_3: 5249,
      ratings_4: 9846,
      ratings_5: 11090
    },
    image_url: 'https://images.gr-assets.com/books/1390667790m/29187.jpg',
    tags: [
      { tag_id: 30574, tag_name: 'to-read' },
      { tag_id: 11305, tag_name: 'fantasy' },
      { tag_id: 9866, tag_name: 'dragonlance' },
      { tag_id: 11743, tag_name: 'fiction' },
      { tag_id: 11557, tag_name: 'favorites' }
    ],
    ratings: [
      {
        user: { user_id: 149, user_name: 'Anthony Palmer' },
        rating: 5
      }
    ]
  },
  {
    _id: ObjectId('65dceab927a93ef44b0559e5'),
    book_id: '1096',
    isbn: '671742515',
    isbn13: '9780671742520.0',
    authors: [ 'Douglas Adams' ],
    original_publication_year: 1988,
    title: 'The Long Dark Tea-Time of the Soul',
    language_code: 'en-US',
    average_rating: 4.05,
    total_ratings: 59276,
    rating_counts: {
      ratings_1: 393,
      ratings_2: 2464,
      ratings_3: 13614,
      ratings_4: 24400,
      ratings_5: 22898
    },
    image_url: 'https://images.gr-assets.com/books/1388257271m/357.jpg',
    tags: [
      { tag_id: 30574, tag_name: 'to-read' },
      { tag_id: 11743, tag_name: 'fiction' },
      { tag_id: 11305, tag_name: 'fantasy' },
      { tag_id: 26837, tag_name: 'science-fiction' },
      { tag_id: 26771, tag_name: 'sci-fi' }
    ],
    ratings: [
      {
        user: { user_id: 149, user_name: 'Anthony Palmer' },
        rating: 5
      }
    ]
  },
  {
    _id: ObjectId('65dceab927a93ef44b0559e5'),
    book_id: '1096',
    isbn: '671742515',
    isbn13: '9780671742520.0',
    authors: [ 'Douglas Adams' ],
    original_publication_year: 1988,
    title: 'The Long Dark Tea-Time of the Soul',
    language_code: 'en-US',
    average_rating: 4.05,
    total_ratings: 59276,
    rating_counts: {
      ratings_1: 393,
      ratings_2: 2464,
      ratings_3: 13614,
      ratings_4: 24400,
      ratings_5: 22898
    },
    image_url: 'https://images.gr-assets.com/books/1388257271m/357.jpg',
    tags: [
      { tag_id: 30574, tag_name: 'to-read' },
      { tag_id: 11743, tag_name: 'fiction' },
      { tag_id: 11305, tag_name: 'fantasy' },
      { tag_id: 26837, tag_name: 'science-fiction' },
      { tag_id: 26771, tag_name: 'sci-fi' }
    ],
    ratings: [
      {
        user: { user_id: 149, user_name: 'Anthony Palmer' },
        rating: 5
      }
    ]
  }
]
bookstore>

```

5.1.3 3

Description

Find users who have rated a book but have an empty `to_read` list.

Query

```
db.users.find({
  ratings: { $exists: true, $not: { $size: 0 } },
  to_read: { $size: 0 }
}, {
  _id: false, ratings: false
});
```

Output

```
bookstore> db.users.find((ratings: { $exists: true, $not: { $size: 0 } }, to_read: { $size: 0 }}, {_id: false, ratings: false});
[
  { user_id: 4, user_name: 'Tracy Howard', to_read: [] },
  { user_id: 18, user_name: 'Dale Peterson', to_read: [] },
  { user_id: 21, user_name: 'Alan Dalton', to_read: [] },
  { user_id: 25, user_name: 'Cheryl Ross', to_read: [] },
  { user_id: 38, user_name: 'Felicia Turner', to_read: [] },
  { user_id: 42, user_name: 'Courtney Jones', to_read: [] },
  { user_id: 53, user_name: 'Robert Sullivan', to_read: [] },
  { user_id: 62, user_name: 'Mrs. Chelsea May', to_read: [] },
  { user_id: 63, user_name: 'Xavier Mora', to_read: [] },
  { user_id: 68, user_name: 'Hannah Hardy', to_read: [] },
  { user_id: 74, user_name: 'David Morgan', to_read: [] },
  { user_id: 77, user_name: 'Julia Rocha', to_read: [] },
  { user_id: 80, user_name: 'Mr. Jerome Rose', to_read: [] },
  { user_id: 85, user_name: 'Adam Hall MD', to_read: [] },
  { user_id: 86, user_name: 'Rachel Hernandez', to_read: [] },
  { user_id: 100, user_name: 'Christie Rivera', to_read: [] },
  { user_id: 110, user_name: 'Terry Smith', to_read: [] },
  { user_id: 114, user_name: 'James Hansen', to_read: [] },
  { user_id: 115, user_name: 'William Potter', to_read: [] },
  { user_id: 117, user_name: 'Joann Price', to_read: [] }
]
Type "it" for more
bookstore>
```

5.1.4 4

Description

Add a book to a user's to_read list.

Query

```
db.users.updateOne({
  "user_id" : 4
}, {
  $push: {
    "to_read": {
      book: { book_id: 0, authors: ["Daniel Gordon"], title: My Story,
        isbn: "0", isbn13: 0, average_rating: -1, tags: ["trash"]}
    }
  }
})
```

Output

```
bookstore> db.users.find((user_id: 4),(user_id: true, user_name: true, to_read: true))
{
  _id: ObjectId('65dceab027a03ef44b0558e0'),
  user_id: 4,
  user_name: 'Tracy Howard',
  to_read: []
}
Before
bookstore> db.users.updateOne({ "user_id": 4 }, { $push: { "to_read": { book: { book_id: 0, authors: ["Daniel Gordon"], title: "My Story", isbn: "0", isbn13: 0, average_rating: -1, tags: ["trash"]} } } })
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
Query
bookstore> db.users.find((user_id: 4),(user_id: true, user_name: true, to_read: true))
{
  _id: ObjectId('65dceab027a03ef44b0558e0'),
  user_id: 4,
  user_name: 'Tracy Howard',
  to_read: [
    {
      book: {
        book_id: 0,
        authors: [ 'Daniel Gordon' ],
        title: 'My Story',
        isbn: '0',
        isbn13: 0,
        average_rating: -1,
        tags: [ 'trash' ]
      }
    }
  ]
}
After
bookstore>
```


5.2 MRCGAB004

5.2.1 1

Description

Find and return a list of users' names who want to read "The book Thief".

Query

```
db.users.find( {  
  "to_read.book.title":"The Book Thief"  
}, {  
  "_id":0,"user_name":1  
})
```

Output

```
bookstore> db.users.find( {"to_read.book.title":"The Book Thief"}, {"_id":0,"user_name":1})  
[  
  { user_name: 'Mrs. Cynthia Martinez DVM' },  
  { user_name: 'Stephanie Pham' }  
]
```

5.2.2 2

Description

Find a list of books published after 2004 and return the title and year it was published.

Query

```
db.books.find({
  "original_publication_year":{ "$gt": 2004 }
},{
  "title":1, "_id":0,"original_publication_year":1
})
```

Output

```
bookstore> db.books.find({"original_publication_year":{"$gt": 2004 }},{ "title":1, "_id":0,"original_publication_year":1})
[
  { original_publication_year: 2005, title: 'Bergdorf Blondes' },
  {
    original_publication_year: 2006,
    title: 'World War Z: An Oral History of the Zombie War'
  },
  { original_publication_year: 2006, title: 'Forest Mage' },
  { original_publication_year: 2005, title: 'The Undomestic Goddess' },
  { original_publication_year: 2005, title: 'Shadow of the Giant' },
  { original_publication_year: 2006, title: 'Mistral's Kiss' },
  { original_publication_year: 2006, title: 'Hunters of Dune' },
  {
    original_publication_year: 2005,
    title: 'Freakonomics: A Rogue Economist Explores the Hidden Side of Everything'
  },
  {
    original_publication_year: 2006,
    title: 'I Like You: Hospitality Under the Influence'
  },
  { original_publication_year: 2006, title: 'Everyman' },
  {
    original_publication_year: 2005,
    title: 'Marley & Me: Life and Love with the World's Worst Dog'
  },
  {
    original_publication_year: 2005,
    title: 'The Complete Calvin and Hobbes'
  },
  { original_publication_year: 2006, title: 'Kushiel's Scion' },
  { original_publication_year: 2006, title: 'Hannibal Rising' },
  {
    original_publication_year: 2005,
    title: 'Misquoting Jesus: The Story Behind Who Changed the Bible and Why'
  },
  {
    original_publication_year: 2006,
    title: 'Y: The Last Man Vol. 7: Paper Dolls'
  },
  {
    original_publication_year: 2005,
    title: 'One Night @ The Call Center'
  },
  { original_publication_year: 2006, title: 'American Born Chinese' },
  { original_publication_year: 2006, title: 'The Night Watch' },
  { original_publication_year: 2006, title: 'Stumbling on Happiness' }
]
Type "it" for more
```

5.2.3 3

Description

Update the book “Bergdorf Blondes” with a new publishing date of 2005.

Query

```
db.books.updateOne({  
  "title": "Bergdorf Blondes"  
}, {  
  "$set": { "original_publication_year": 2005 }  
})
```

Output

```
bookstore> db.books.find({ "title": "Bergdorf Blondes" }, { "original_publication_year": 1, "title": 1, "_id": 0 })  
[ { original_publication_year: 2004, title: 'Bergdorf Blondes' } ]  
bookstore>  
  
bookstore> db.books.updateOne({ "title": "Bergdorf Blondes" }, { "$set": { "original_publication_year": 2005 } })  
{  
  acknowledged: true,  
  insertedId: null,  
  matchedCount: 1,  
  modifiedCount: 1,  
  upsertedCount: 0  
}  
bookstore>  
  
bookstore> db.books.find({ "title": "Bergdorf Blondes" }, { "original_publication_year": 1, "title": 1, "_id": 0 })  
[ { original_publication_year: 2005, title: 'Bergdorf Blondes' } ]
```

5.2.4 4

Description

Delete user with the user_id 32.

Query

```
db.users.deleteOne({"user_id":32})
```

Output

```
bookstore> db.users.find({"user_id":32},{ "user_id":1, "user_name":1, "_id":0})
[ { user_id: 32, user_name: 'Kaitlyn Walsh' } ]
bookstore>

bookstore> db.users.deleteOne({"user_id":32})
{ acknowledged: true, deletedCount: 1 }
bookstore>

bookstore> db.users.find({"user_id":32},{ "user_id":1, "user_name":1, "_id":0})
bookstore>
```

5.3 CRGMAT002

5.3.1 1

Description

Find the top 5 books (by `average_rating`) with `tag_name` “sci-fi”. Show only a single author for each book.

Query

```
db.books.find({
  tags: {$elemMatch: {tag_name: "sci-fi"}}
}, {
  _id: 0, title: 1, authors: {$slice: 1}, average_rating: 1,
}).sort({average_rating : -1 }).limit(5)
```

Output

```

bookstore> db.books.find({
...   tags: {
...     $elemMatch: {
...       tag_name: "sci-fi"
...     }
...   }
... }, {
...   _id: 0,
...   title: 1,
...   authors: {$slice: 1},
...   average_rating: 1,
... }).sort({average_rating : -1 }).limit(5)
[
  {
    authors: [ 'Isaac Asimov' ],
    title: 'The Foundation Trilogy',
    average_rating: 4.39
  },
  {
    authors: [ 'Douglas Adams' ],
    title: "The Ultimate Hitchhiker's Guide: Five Complete Novels and One Story",
    average_rating: 4.37
  },
  {
    authors: [ 'Warren Ellis' ],
    title: 'Transmetropolitan, Vol. 2: Lust for Life',
    average_rating: 4.35
  },
  {
    authors: [ 'Isaac Asimov' ],
    title: 'The Complete Robot',
    average_rating: 4.34
  },
  {
    authors: [ 'Vernor Vinge' ],
    title: 'A Deepness in the Sky',
    average_rating: 4.32
  }
]

```

5.3.2 2

Description

Aggregate `user_ids` alongside the total number of ratings that user has submitted (size of ratings array). The result is sorted by the number of ratings they have submitted and the top 10 are shown.

Query

```
db.users.aggregate([ {  
  $project: { _id: 0, user_id: 1, numRatings: { $size: "$ratings" } }  
}, {  
  $sort: { numRatings: -1 }  
}, {  
  $limit: 10  
}])
```

Output

```
bookstore> db.users.aggregate([ {
...     $project: {
...         _id: 0,
...         user_id: 1,
...         numRatings: { $size: "$ratings" }
...     }
... }, { $sort: { numRatings: -1 }},
... { $limit: 10 }
... ])
[
  { user_id: 158, numRatings: 100 },
  { user_id: 370, numRatings: 100 },
  { user_id: 395, numRatings: 100 },
  { user_id: 137, numRatings: 100 },
  { user_id: 247, numRatings: 100 },
  { user_id: 207, numRatings: 100 },
  { user_id: 446, numRatings: 100 },
  { user_id: 324, numRatings: 100 },
  { user_id: 116, numRatings: 100 },
  { user_id: 301, numRatings: 99 }
]
```


5.3.3 3

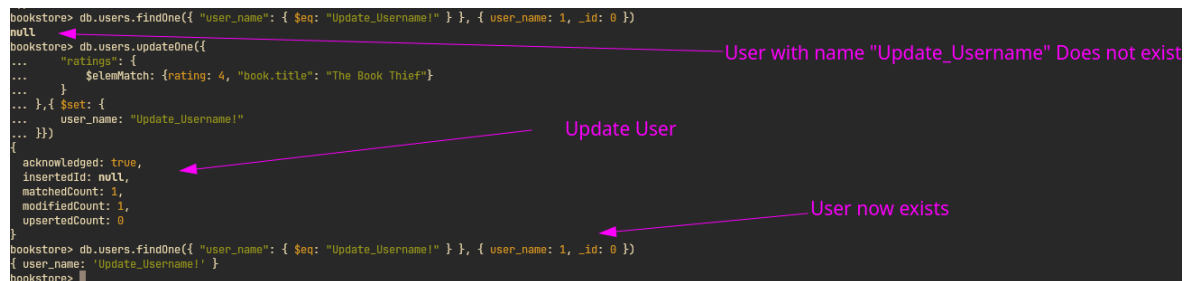
Description

Update a user that has rated “The Book Thief” with 4 stars. Change their name to: “Update_User!”

Query

```
db.users.updateOne({
  "ratings": {$elemMatch: {rating: 4, "book.title": "The Book Thief"}}
},{
  $set: {user_name: "Update_Username!"}
})
```

Output



```
bookstore> db.users.findOne({ "user_name": { $eq: "Update_Username!" } }, { user_name: 1, _id: 0 })
null
bookstore> db.users.updateOne({
...   "ratings": {
...     $elemMatch: {rating: 4, "book.title": "The Book Thief"}
...   },
...   $set: {
...     user_name: "Update_Username!"
...   })
{
  acknowledged: true,
  insertedId: null,
  matchedCount: 1,
  modifiedCount: 1,
  upsertedCount: 0
}
bookstore> db.users.findOne({ "user_name": { $eq: "Update_Username!" } }, { user_name: 1, _id: 0 })
{ user_name: 'Update_Username!' }
bookstore>
```

Annotations in the image:

- “User with name “Update_Username” Does not exist” points to the null result of the first findOne query.
- “Update User” points to the updateOne query.
- “User now exists” points to the result of the second findOne query.

5.3.4 4

Description

Delete all book's written by Roald Dahl.

Query

```
db.books.deleteMany({
  authors: "Roald Dahl"
},{
  title: 1
})
```

Output



```
bookstore> db.books.find({
...   authors: "Roald Dahl"
... }, {title: 1})
[
  {
    id: ObjectId('65e07ed87b18979e01983247'),
    title: 'Fantastic Mr. Fox'
  },
  {
    id: ObjectId('65e07ed87b18979e01983292'), title: 'The Witches' },
  {
    id: ObjectId('65e07ed87b18979e01983298'),
    title: 'Boy: Tales of Childhood'
  },
  {
    id: ObjectId('65e07ed87b18979e0198332c'),
    title: 'The Giraffe and the Pelly and Me'
  },
  {
    id: ObjectId('65e07ed87b18979e01983427'), title: 'The Twits' },
    id: ObjectId('65e07ed87b18979e019834df'), title: 'Matilda' },
    id: ObjectId('65e07ed87b18979e019835bd'), title: 'The BFG' },
    id: ObjectId('65e07ed87b18979e01983614'), title: 'Going Solo' },
  {
    id: ObjectId('65e07ed87b18979e019836c7'),
    title: 'George's Marvellous Medicine'
  },
  {
    id: ObjectId('65e07ed87b18979e019837cf'),
    title: 'James and the Giant Peach'
  },
  {
    id: ObjectId('65e07ed87b18979e01983856'),
    title: 'Charlie and the Chocolate Factory'
  },
  {
    id: ObjectId('65e07ed87b18979e01983905'),
    title: 'The Wonderful Story Of Henry Sugar And Six More'
  }
]
bookstore> db.books.deleteMany({
...   authors: "Roald Dahl"
... }, {title: 1})
{ acknowledged: true, deletedCount: 12 }
bookstore> db.books.find({ authors: "Roald Dahl" }, { title: 1 })
[]
bookstore>
```

Find for author "Roald Dahl" returns many books

Delete operation

Find for author "Roald Dahl" returns no books

5.4 WHLJOS001

5.4.1 1

Description

Find the names and ratings of the top 50 fiction books with at least 1000 ratings. The result would be too large to show in a single screen but several are shown in the image below.

Query

```
db.books.aggregate([
  {
    $match: { "total_ratings": { $gte: 1000 }, "tags.tag_name": "fiction" }
  }, {
    $sort: { "average_rating": -1 }
  }, {
    $limit: 50
  }, {
    $project: { _id: 0, title: 1, average_rating: 1 }
  } ] )
```

Output

```
bookstore> db.books.aggregate([ { $match: { "total_ratings": { $gte: 1000 }, "tags.tag_name": "fiction" } }, { $sort: { "average_rating": -1 } }, { $limit: 50 }, { $project: { _id: 0, title: 1, average_rating:
  } } ] )
[
  {
    title: "It's a Magical World: A Calvin and Hobbes Collection",
    average_rating: 4.75
  },
  {
    title: "It's a Magical World: A Calvin and Hobbes Collection",
    average_rating: 4.75
  },
  {
    title: "There's Treasure Everywhere: A Calvin and Hobbes Collection",
    average_rating: 4.74
  },
  {
    title: "There's Treasure Everywhere: A Calvin and Hobbes Collection",
    average_rating: 4.74
  },
  {
    title: "Harry Potter Collection (Harry Potter, #1-6)",
    average_rating: 4.73
  },
  {
    title: "Harry Potter Collection (Harry Potter, #1-6)",
    average_rating: 4.73
  },
  {
    title: "Harry Potter and the Deathly Hallows",
    average_rating: 4.61
  },
  {
    title: "Harry Potter and the Deathly Hallows",
    average_rating: 4.61
  },
  {
    title: "The Hobbit and The Lord of the Rings",
    average_rating: 4.59
  },
  {
    title: "The Hobbit and The Lord of the Rings",
    average_rating: 4.59
  },
  {
    title: "Cuentos completos", average_rating: 4.58 },
  {
    title: "Cuentos completos", average_rating: 4.58 },
  {
    title: "A Storm of Swords", average_rating: 4.54 },
  {
    title: "A Storm of Swords", average_rating: 4.54 },
  ]
```

5.4.2 2

Description

Delete all users who have never rated a book and never marked a book as `to_read`. Since this database has only a subset of the full dataset it seems that no users were found to both not have rated or marked a book as to read, so no data was deleted.

Query

```
db.users.deleteMany({
  $and: [{to_read: { $exists: false }},
    {ratings: { $exists: false }}]
})
```

Output

```
bookstore> db.users.deleteMany({$and: [{ to_read: { $exists: false } }, { ratings: { $exists: false } }]}
{ acknowledged: true, deletedCount: 0 }
bookstore> |
```

5.4.3 3

Description

Delete all books with an average rating less than 2 and more than 300 ratings. Similar to the previous query, since this is only a select subset of the dataset no books were found that matched the query. So no books were deleted in this case.

Query

```
db.books.deleteMany({
  $and: [ {
    "average_rating": { $lt: 2 }
  }, {
    "total_ratings": { $gt: 300 }
  }]
})
```

Output

```
bookstore> db.books.deleteMany({ $and: [ { "average_rating": { $lt: 2 } }, { "total_ratings": { $gt: 300 } } ] })
{ acknowledged: true, deletedCount: 0 }
bookstore> |
```

5.4.4 4

Description

Find the name of the user who has rated the most books.

Query

```
db.users.aggregate({
  $project: { user_id: 1, user_name: 1, num_ratings: { $size: "$ratings" } }
},{
  $sort: { num_ratings: -1 }
},{
  $limit: 1
}, {
  $project: {
    _id: 0, user_name: 1
  }
})
```

Output

```
bookstore> db.users.aggregate({ $project: { user_id: 1, user_name: 1, num_ratings: { $size: "$ratings" } } }, { $sort: {
num_ratings: -1 } }, { $limit: 1 }, { $project: { _id: 0, user_name: 1 } })
[ { user_name: 'Catherine Zimmerman' } ]
bookstore> |
```

5.5 KFWJOR001

5.5.1 1

Description

Finding users who have rated the most diverse genres

Query

```
db.users.aggregate([
  { $unwind: "$ratings" },
  { $lookup: {
    from: "books",
    localField: "ratings.book.title", foreignField: "title",
    as: "rated_books" } },
  { $unwind: "$rated_books" },
  { $unwind: "$rated_books.tags" },
  { $group: { _id: "$user_id", user_name: { $first: "$user_name" },
    genres: { $addToSet: "$rated_books.tags.tag_name" } } },
  { $addFields: { num_genres: { $size: "$genres" } } },
  { $sort: { num_genres: -1 } },
  { $limit: 10 }
])
```

Output

```

[
  {
    _id: 89,
    user_name: 'Justin Hall',
    genres: [
      'kids-books',      'currently-reading', 'childhood',
      'literature',      'mythology',         'thriller',
      'harry-potter',    'ثاني اور',          'books-i-own',
      'comics',          'biography',         'young-adult',
      'psychology',      'chicklit',          'nonfiction',
      'ya',              'chick-lit',         'india',
      'magical-realism', 'humor',              'childrens',
      'food',            'history',            'classics',
      'finance',         'mystery',            'spirituality',
      'fantasy',         'japan',              'children',
      'memoir',          'dystopian',          'christian',
      'china',           'romance',            'historical',
      'roald-dahl',      'owned',              'default',
      'poetry',          'french',              'feminism',
      'war',             'dystopia',           'self-help',
      'science-fiction', 'fiction',             'school',
      'business',        'spanish',             'classic',
      'shakespeare',    'children-s',          'to-read',
      'gothic',          'writing',             'novels',
      'sci-fi',          'plays',               'favorites',
      'classics',        'time-travel',         'book-club',
      'kids',            'historical-fiction',  'childhood-favorites',
      'non-fiction',     'childrens-books',     'bbc-big-read',
      'greece'
    ],
    num_genres: 70
  },
  {
    id: 247,

```


5.5.2 2

Description

Find highly rated books and only show their title, authors and average rating

Query

```
db.books.find({ "average_rating": { $gte: 4.5 } }, {title: 1, authors:1, average_rating: 1})
```

Output

```
goodbooks> db.books.find({ "average_rating": { $gte: 4.5 } }, {title: 1, authors:1, average_rating: 1})
[
  {
    _id: ObjectId('65ddcaa866c3549e7ba220ec'),
    authors: [
      'Neil Gaiman',
      'Marc Hempel',
      'Richard Case',
      'DIsraeli',
      'Teddy Kristiansen',
      'Glyn Dillon',
      'Charles Vess',
      'Dean Ormston',
      'Kevin Nowlan',
      'Todd Klein',
      'Frank McConnell'
    ],
    title: 'The Sandman: The Kindly Ones',
    average_rating: 4.59
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba22137'),
    authors: [ 'Bill Watterson' ],
    title: 'The Complete Calvin and Hobbes',
    average_rating: 4.82
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba22138'),
    authors: [ 'Bill Watterson' ],
    title: "It's a Magical World: A Calvin and Hobbes Collection",
    average_rating: 4.75
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba22139'),
```

5.5.3 3

Description

Give all books published before 2010 the “classic” tag

Query

```
db.books.updateMany( {  
  "original_publication_year": { $lt: 2010 } }, {  
    addToSet: { "tags": {tag_id: 12345, tag_name: "classic"} } }
```

Output

- 3 - Query

```
goodbooks> db.books.updateMany( { "original_publication_year": { $lt: 2010 } }, { $addToSet: { "tags": {tag_id: 12345, tag_name: "classic"} } } )  
{  
  acknowledged: true,  
  insertedId: null,  
  matchedCount: 2184,  
  modifiedCount: 2184,  
  upsertedCount: 0  
}  
goodbooks> █
```

- 3 - Before

```
goodbooks> db.books.find({"original_publication_year": { $lt: 2010 }}, {title: 1, tags: 1, original_publication_year: 1})  
[  
  {  
    _id: ObjectId('65ddcaa866c3549e7ba220ce'),  
    original_publication_year: 2004,  
    title: 'Bergdorf Blondes',  
    tags: [  
      { tag_id: 11743, tag_name: 'fiction' },  
      { tag_id: 8717, tag_name: 'currently-reading' },  
      { tag_id: 30574, tag_name: 'to-read' },  
      { tag_id: 11557, tag_name: 'favorites' },  
      { tag_id: 8055, tag_name: 'contemporary' }  
    ]  
  },  
  {  
    _id: ObjectId('65ddcaa866c3549e7ba220cf'),  
    original_publication_year: 1987,  
    title: 'The Black Dahlia',  
    tags: [  
      { tag_id: 30574, tag_name: 'to-read' },  
      { tag_id: 20939, tag_name: 'mystery' },  
      { tag_id: 8517, tag_name: 'crime' },  
      { tag_id: 11743, tag_name: 'fiction' },  
      { tag_id: 21672, tag_name: 'noir' }  
    ]  
  },  
  {  
    _id: ObjectId('65ddcaa866c3549e7ba220d0'),  
    original_publication_year: 1996,  
    title: 'A Supposedly Fun Thing I'll Never Do Again: Essays and Arguments',  
    tags: [  
      { tag_id: 30574, tag_name: 'to-read' },  
    ]  
  }  
]
```

- 3 - After

```
goodbooks> db.books.find({"original_publication_year": { $lt: 2010 }}, {title: 1, tags: 1, original_publication_year: 1})
[
  {
    _id: ObjectId('65ddcaa866c3549e7ba220ce'),
    original_publication_year: 2004,
    title: 'Bergdorf Blondes',
    tags: [
      { tag_id: 11743, tag_name: 'fiction' },
      { tag_id: 8717, tag_name: 'currently-reading' },
      { tag_id: 30574, tag_name: 'to-read' },
      { tag_id: 11557, tag_name: 'favorites' },
      { tag_id: 8055, tag_name: 'contemporary' },
      { tag_id: 12345, tag_name: 'classic' }
    ]
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba220cf'),
    original_publication_year: 1987,
    title: 'The Black Dahlia',
    tags: [
      { tag_id: 30574, tag_name: 'to-read' },
      { tag_id: 20939, tag_name: 'mystery' },
      { tag_id: 8517, tag_name: 'crime' },
      { tag_id: 11743, tag_name: 'fiction' },
      { tag_id: 21672, tag_name: 'noir' },
      { tag_id: 12345, tag_name: 'classic' }
    ]
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba220d0'),
    original_publication_year: 1996,
    title: 'A Supposedly Fun Thing I'll Never Do Again: Essays and Arguments',
    tags: [
      { tag_id: 30574, tag_name: 'to-read' }
    ]
  }
]
```

5.5.4 4

Description

Updates all books authored by J.K. Rowling (all Harry Potter) by removing the incorrect author “Mary Grandpré” while also removing any ratings below 5 stars for those books

Query

```
db.books.updateMany(  
  { "authors": "JK Rowling" },  
  {  
    $set: { "authors": ["JK Rowling"] },  
    $pull: { "ratings": { "rating": { $lt: 5 } } }  
  }  
);
```

Output

```
goodbooks> db.books.updateMany(  
...   { "authors": "JK Rowling" },  
...   {  
...     $set: { "authors": ["JK Rowling"] },  
...     $pull: { "ratings": { "rating": { $lt: 5 } } }  
...   }  
... );  
{  
  acknowledged: true,  
  insertedId: null,  
  matchedCount: 10,  
  modifiedCount: 9,  
  upsertedCount: 0  
}
```

The Query

```

goodbooks> db.books.find({ "authors": "JK Rowling", "ratings.rating": { $lt: 5 } })
[
  {
    _id: ObjectId('65ddcaa866c3549e7ba226ef'),
    book_id: '2402163',
    isbn: '043965548X',
    isbn13: '9780439655480.0',
    authors: [ 'JK Rowling', 'Mary GrandPr', 'Rufus Beck' ],
    original_publication_year: 1999,
    title: 'Harry Potter and the Prisoner of Azkaban',
    language_code: 'eng',
    average_rating: 4.53,
    total_ratings: 1832823,
    rating_counts: {
      ratings_1: 6716,
      ratings_2: 20413,
      ratings_3: 166129,
      ratings_4: 509447,
      ratings_5: 1266670
    },
    image_url: 'https://images.gr-assets.com/books/1499277281m/5.jpg',
    tags: [
      { tag_id: 11557, tag_name: 'favorites' },
      { tag_id: 11305, tag_name: 'fantasy' },
      { tag_id: 8717, tag_name: 'currently-reading' },
      { tag_id: 33114, tag_name: 'young-adult' },
      { tag_id: 30574, tag_name: 'to-read' },
      { tag_id: 12345, tag_name: 'classic' }
    ],
    ratings: [
      { user: { user_id: 4, user_name: 'Tracy Howard' }, rating: 5 },
      { user: { user_id: 15, user_name: 'Matthew Gomez' }, rating: 4 },
      {

```

Before ratings

```

27
  {
    _id: ObjectId('65ddcaa866c3549e7ba22950'),
    authors: [ 'JK Rowling' ],
    title: 'Harry Potter and the Half-Blood Prince'
  }
]
goodbooks> db.books.find({ "authors": "JK Rowling", "ratings.rating": { $lt: 5 } })

goodbooks> █

```

After ratings

```

goodbooks> db.books.find({ "authors": "JK Rowling" }, {title: 1, authors: 1})
[
  {
    _id: ObjectId('65ddcaa866c3549e7ba226ef'),
    authors: [ 'JK Rowling', 'Mary GrandPr', 'Rufus Beck' ],
    title: 'Harry Potter and the Prisoner of Azkaban'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba2277c'),
    authors: [ 'JK Rowling', 'Mary GrandPr' ],
    title: 'Harry Potter and the Order of the Phoenix'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba227c0'),
    authors: [ 'JK Rowling', 'Mary GrandPr' ],
    title: 'Harry Potter and the Deathly Hallows'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba227ea'),
    authors: [ 'JK Rowling', 'Mary GrandPr' ],
    title: 'Harry Potter and the Goblet of Fire'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba228d9'),
    authors: [ 'Newt Scamander', 'JK Rowling', 'Albus Dumbledore' ],
    title: 'Fantastic Beasts and Where to Find Them'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba228ec'),
    authors: [ 'JK Rowling', 'Mary GrandPr' ],
    title: "Harry Potter and the Philosopher's Stone"
  },
]

```

Before books

```
goodbooks> db.books.find({ "authors": "JK Rowling" }, {title: 1, authors: 1})
[
  {
    _id: ObjectId('65ddcaa866c3549e7ba226ef'),
    authors: [ 'JK Rowling' ],
    title: 'Harry Potter and the Prisoner of Azkaban'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba2277c'),
    authors: [ 'JK Rowling' ],
    title: 'Harry Potter and the Order of the Phoenix'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba227c0'),
    authors: [ 'JK Rowling' ],
    title: 'Harry Potter and the Deathly Hallows'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba227ea'),
    authors: [ 'JK Rowling' ],
    title: 'Harry Potter and the Goblet of Fire'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba228d9'),
    authors: [ 'JK Rowling' ],
    title: 'Fantastic Beasts and Where to Find Them'
  },
  {
    _id: ObjectId('65ddcaa866c3549e7ba228ec'),
    authors: [ 'JK Rowling' ],
    title: "Harry Potter and the Philosopher's Stone"
  },
  {

```

After books

6 Link the Database to a Program

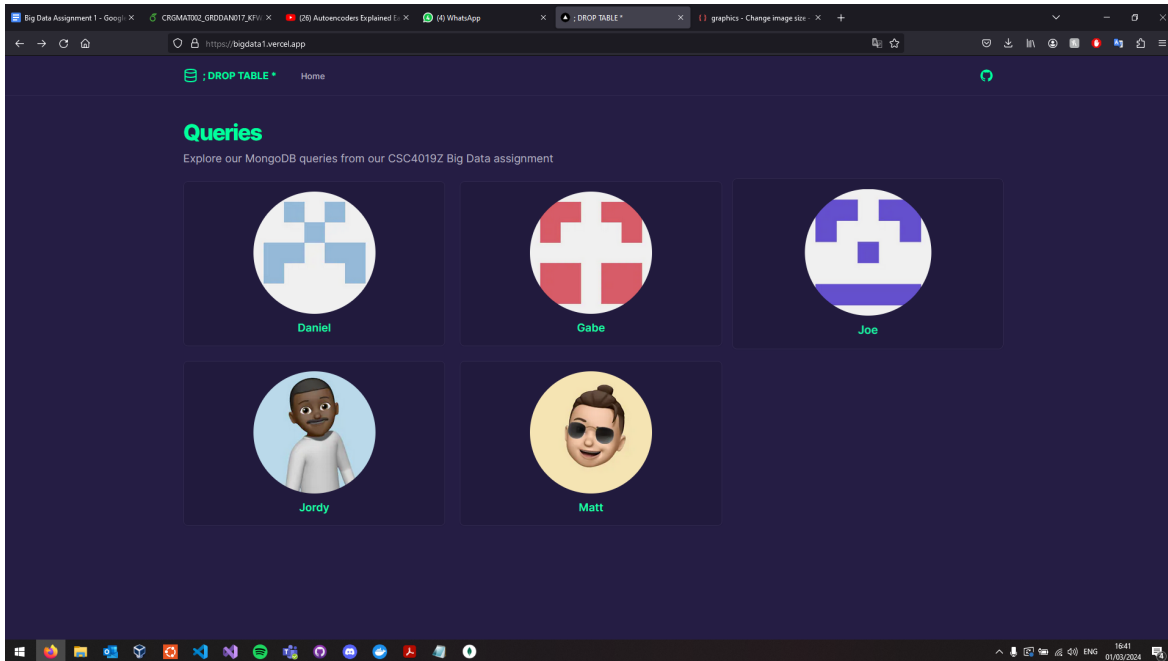
We created a web application which can run on your local machine but can also be extended to run on a live server. Instructions for running the application locally can be found in the project readme. The app is also deployed and available at <https://bigdata1.vercel.app>. The deployed version of the application uses MongoDB Atlas (a fully-managed cloud version of MongoDB).

The stack used in the creation of the program are: TypeScript, Next.js Tailwind CSS was used for styling.

The frontend was implemented in Next.js and TypeScript and the backend uses Next.js server actions, together with Mongoose (Object Data Modelling library for MongoDB) for connecting to and querying the database. Server actions are simply asynchronous functions that run on the server side of a Next.js app. The implementation of the server actions can be seen in the server directory of the frontend folder. Additionally, the calling of these server actions can be seen in the page.tsx files in the Person directory inside frontend/app. The actual connection code can be found in app/lib in db-connect.ts.

Instructions for running the program can also be found in **README.md**.

Screenshots:



Screenshot of the home page. Click on a student and wait to view query and output.

Jordy's Query

Find the top 10 users who have rated the most diverse range of genres

```
db.users.aggregate([
  { $wind: "$ratings" }, { $lookup: { from: "books", localField: "ratings.book.title", foreignField: "title", as: "rated_books" } }, { $unwind: "$rated_books" }, { $unwind: "$rated_books.tags" }, { $group: { _id: "$user_id", user_name: { $first: "$user_name" }, genres: { $addToSet: "$rated_books.tags.tag_name" } } }, { $addFields: { num_genres: { $size: "$genres" } } }, { $sort: { num_genres: -1 } }, { $limit: 10 } ])
```

Justin Hall
70 genres

biography war harry-potter historical to-read gothic writing kids psychology fantasy classics young-adult
 childhood india childrens-books comics school spanish poetry feminism china 中国 book-club science-fiction
 mystery dystopian non-fiction plays shakespeare bbc-big-read christian books-i-own chicklit sci-fi childrens
 magical-realism owned thriller roald-dahl classic time-travel food kids-books childhood-favorites mythology literature
 novels favorites classics history business self-help humor dystopia memoir children-s romance japan ya
 currently-reading historical-fiction nonfiction finance greece default children french spirituality chick-lit fiction

Jordy's Query with output.

Daniel's Query

Find users who have rated a book but have an empty read list

```
db.users.find({ ratings: { $exists: true, $not: { $size: 0 } }, to_read: { $size: 0 }, { _id: false, ratings: false } })
```

Name
Tracy Howard
Dale Peterson
Alan Dalton
Cheryl Ross
Felicia Turner
Courtney Jones
Robert Sullivan
Hannah Hardy
Xavier Mora

Daniel's Query with output.

The screenshot shows a web application titled "Gabe's Query" with a red cross icon. The query displayed is:

```
db.users.find({"to_read.book.title": "The Book Thief", "_id": 0, "user_name": 1})
```

The output is a table with two rows:

Name
Mrs. Cynthia Martinez DVM
Stephanie Pham

Gabe's Query with output.

The screenshot shows a web application titled "Joe's Query" with a blue cross icon. The query displayed is:

```
db.books.aggregate([{$match: { "total_ratings": { $gte: 1000 }, "tags.tag_name": "fiction" }}, {$sort: { "average_rating": -1 }}, {$limit: 50}, {$project: { _id: 0, title: 1, average_rating: 1 } }])
```

The output is a table with two columns: Title and Average Rating.


	Title	Average Rating
1	It's a Magical World: A Calvin and Hobbes Collection	★ 4.75
2	There's Treasure Everywhere: A Calvin and Hobbes Collection	★ 4.74
3	Harry Potter Collection (Harry Potter, #1-6)	★ 4.73
4	Harry Potter and the Deathly Hallows	★ 4.61
5	The Hobbit and The Lord of the Rings	★ 4.59
6	Cuentos completos	★ 4.58

Joe's Query with output.

https://bigdata1.vercel.app/person/matt

DROP TABLE *

Home



Matt's Query

Aggregate user IDs alongside the total number of ratings that user has submitted. The result is sorted by the number of ratings they have submitted and the top 10 are shown

```
db.users.aggregate([{$project: {_id: 0, user_id: 1, numRatings: { $size: "$ratings" } }}, { $sort: { numRatings: -1 }}, { $limit: 10 }])
```

User ID	Number of Ratings
247	100
395	100
324	100
207	100
116	100
370	100
446	100
158	100

Matt's Query with output.

7 Contribution Statement

Task	Who did it
Design Schema/JSON	Matthew & Gabe
Setup Devops (Github, Docker)	Jordy
Preprocessing	Joe & Daniel
Create the App	Jordy
Query	Everyone
Load into mongodb	Mathew
Write Up	Gabe & Matthew & Daniel