# Final Project Report

Benson Ning qn2

## Notes to Instructors

Due to the nature of this project that it involves connections with databases, it is impossible to clone the project, hit run, and expect everything to work smoothly. Please follow the instructions if you do not have databases installed on your laptop.

## Environment Preparation

### Clone Project

- https://github.com/Aomame-Masami/finalproject480.git

### Extract files needed

- unzip bookcsv.zip into bookcsv directory.

### Database Requirement

- The project requires MySQL 8.0 and Redis 3.0+

### MySQL Installation

- For All Systems: https://dev.mysql.com/downloads/mysql/

### Redis Installation

- For windows: https://github.com/microsoftarchive/redis/releases
- For mac: `brew install redis`

## Data Preparation

### These data should be unzipped to /bookcsv

- Data to be hashed and queried:
  - Locates in bookcsv directory
  - reference: https://www.kaggle.com/datasets/saurabhbagchi/books-dataset/
- Trained vector set for similarity comparison is serialized and store in vector_column.json

## Database Config

- **GO TO *application.properties***
- change the url:

  `spring.datasource.url=jdbc:mysql://localhost:3306/test/useSSL=false&useUnicode=`
  `true&characterEncoding=utf-8&serverTimezone=GMT-5`

  replace "test" with the database you use.
- Change the username and password:

  spring.datasource.username=root -> to your username
  spring.datasource.password=Benson -> to your password
- Change the Redis host and port based on your local settings.

## MySQL data prep

- I recommend using DBeaver for mysql management (it's free)

  https://dbeaver.io/
- Create a table with the following sql cmd:

```
use <Your Database>;

DROP TABLE IF EXISTS book;
CREATE TABLE book (
    ISBN VARCHAR(20) PRIMARY KEY,
    Book_Title VARCHAR(255),
    Book_Author VARCHAR(255),
    Year_Of_Publication INT,
    Publisher VARCHAR(255),
    Image_URL_S VARCHAR(255),
    Image_URL_M VARCHAR(255),
    Image_URL_L VARCHAR(255)
);
```

- Import the csv into the table just created. (You can do it with one-click with DBeaver)
- MySQL would evetually look like this:

| ISBN | Book_Title | Book_Author | Year_Of_Publication | Publisher | Image_URL_S | ABC |
|---|---|---|---|---|---|---|
| 1 | 0000913154 | The Way Things Work | C. van Amerongen (tran | 1,967 | Simon &amp; Schuster | http://images.amazon.com/images/P/000091315 |
| 2 | 0001010565 | Mog's Christmas | Judith Kerr | 1,992 | Collins | http://images.amazon.com/images/P/000101056 |
| 3 | 0001046713 | Twopence to Cross the | Helen Forrester | 1,992 | HarperCollins Publishers | http://images.amazon.com/images/P/000104671 |
| 4 | 0001046934 | The Prime of Miss Jean | Muriel Spark | 1,999 | Trafalgar Square Publishing | http://images.amazon.com/images/P/000104693 |
| 5 | 0001047213 | The Fighting Man | Gerald Seymour | 1,993 | HarperCollins Publishers | http://images.amazon.com/images/P/000104721 |
| 6 | 0001047647 | First Among Equals | Jeffrey Archer | 2,000 | Trafalgar Square Books | http://images.amazon.com/images/P/000104764 |
| 7 | 0001048082 | Made in America | Bill Bryson | 1,995 | HarperCollins Publishers | http://images.amazon.com/images/P/000104808 |
| 8 | 0001048473 | Nothing Can Be Better | Barns &amp; Budd | 1,996 | Atlantic | http://images.amazon.com/images/P/000104847 |
| 9 | 0001052039 | Dark Spectre | Michael Dibdin | 1,997 | HarperCollins Publishers | http://images.amazon.com/images/P/000105203 |
| 10 | 0001056107 | Farmer Giles of Ham: A | J. R. R. Tolkien | 1,999 | Trafalgar Square | http://images.amazon.com/images/P/000105610 |
| 11 | 0001061127 | CHESS FOR YOUNG BE | William T. McLeod | 1,975 | HarperCollins Publishers | http://images.amazon.com/images/P/000106112 |
| 12 | 000123207X | Paddington's Birthday | Michael Bond | 1,942 | HarperCollins Publishers | http://images.amazon.com/images/P/000123207 |
| 13 | 0001232088 | Paddington in the Kitc | Michael Bond | 1,942 | HarperCollins Publishers | http://images.amazon.com/images/P/000123208 |
| 14 | 0001372564 | Which Colour? | Sue Dreamer | 1,989 | HarperCollins Publishers | http://images.amazon.com/images/P/000137256 |
| 15 | 0001382381 | Huck Scarry's Steam Tr | Huck Scarry | 1,979 | HarperCollins Publishers | http://images.amazon.com/images/P/000138238 |
| 16 | 0001711105 | The Book of Riddles (B | Bennett Cerf | 1,983 | HarperCollins Publishers | http://images.amazon.com/images/P/000171110 |
| 17 | 000171421X | It's Not Easy Being a B | Marilyn Sadler | 1,984 | HarperCollins Publishers | http://images.amazon.com/images/P/000171421 |
| 18 | 0001714236 | Spooky Riddles | Mark Brown | 1,984 | HarperCollins Publishers | http://images.amazon.com/images/P/000171423 |
| 19 | 0001811150 | Day Time (Step by Step | Diane Wilmer | 1,987 | HarperCollins Publishers | http://images.amazon.com/images/P/000181115 |
| 20 | 0001811819 | Paddington and the M | Michael Bond | 1,987 | HarperCollins Publishers | http://images.amazon.com/images/P/000181181 |
| 21 | 0001821326 | Paddington at the Tow | Michael Bond | 1,976 | Collin | http://images.amazon.com/images/P/000182132 |
| 22 | 0001837397 | Autumn Story Brambly | Jill Barklem | 0 | William Collins Sons Co Ltd | http://images.amazon.com/images/P/000183739 |
| 23 | 0001840517 | Bess | Robert Leeson | 1,975 | Collins | http://images.amazon.com/images/P/000184051 |
| 24 | 0001841572 | Red Shift | Alan Garner | 1,973 | Collins | http://images.amazon.com/images/P/000184157 |
| 25 | 000184251X | February's Road | John Verney | 1,987 | HarperCollins Publishers | http://images.amazon.com/images/P/000184251 |
| 26 | 0001846086 | There Will Be a Next T | Tony Drake | 1,984 | HarperCollins Publishers | http://images.amazon.com/images/P/000184608 |
| 27 | 0001900277 | Glue (First Facts - First | Harriet Hains | 1,989 | HarperCollins Publishers | http://images.amazon.com/images/P/000190027 |

# Abstract

This project mainly utilizes bloom filter and locality-sensitive hashing to improve database query and recommendations.

- In a scenario where 10,000 clients are simultaneously shopping for books, a primary concern is efficiently determining stock availability. Continuously querying the database for each client request to check stock status can overly burden the database. To mitigate this, we've implemented a layer of Bloom filters for each request. Bloom filters, by their nature, have a zero false-negative rate. This means if the Bloom filter indicates that a book is not in the filter (returns false), we can be certain the book is either out of stock or does not exist. Only for books identified by the Bloom filter as potentially in stock do we then query the database. This approach efficiently screens out unnecessary database queries. Additionally, Bloom filters offer significantly faster query speeds (up to 1300 times faster), providing a substantial performance benefit on the service provider side. This method streamlines the process, reducing database load while maintaining accurate stock information.

- In our system, when customers query for specific books, we aim to offer them additional recommendations. While MySQL's `LIKE` operator allows for approximate queries, it primarily matches books with similar titles. Our objective is broader: to recommend books based on a combination of factors including "Book_Title," "Book_Author," "Year_Of_Publication," and "Publisher." To achieve this, we utilize Locality-Sensitive Hashing (LSH).

  LSH allows us to identify recommendations based on the similarities between dense vectors representing these book attributes. This method is highly efficient, enabling us to generate recommendations rapidly. The process includes obtaining the index of recommended books, retrieving the top-5 matches from the database, and then returning a `List<String>` `bookTitles`. The entire recommendation process is impressively swift, taking only about 0.01 seconds per query. This approach ensures that customers receive relevant, diverse, and quick book suggestions, enhancing their shopping experience.

# Bloom Filter

**The goal for the bloom filter is to return the existence of a Book in our hash set in a much faster speed than querying MySQL database.**

## MySQL Connection:

- Build Book entity

```java
@RedisHash("Book")
public class Book {
    @Id
    private String isbn;
    private String book_Title;
    private String book_Author;
    private int year_Of_Publication;
    private String publisher;
    private String image_UrlS;
    private String image_UrlM;
    private String image_UrlL;


}
```

- Create mapper for query sentences

```java
@Mapper
public interface BookMapper {
    @Select("SELECT * FROM book")
    List<Book> selectAll();

    @Select("SELECT * FROM book WHERE ISBN = #{isbn}")
    Book selectByISBN(@Param("isbn") String isbn);

    @Select("SELECT ISBN FROM book LIMIT 1000")
    List<String> selectFirst1000ISBNs();

    @Select("SELECT ISBN FROM book LIMIT 100")
    List<String> selectFirst100ISBNs();

    @Select("SELECT ISBN FROM book ")
    List<String> selectAllISBNs();
    @Select("SELECT ISBN FROM book")
    List<String> getAllISBNs();


}
```

- In BookService.java, create methods for building a bloom filter for the books and make the query to bloom filter

```java
    public boolean mightContainIsbn(BloomFilter<String> bloomFilter, String isbn)
{
        return bloomFilter.mightContain(isbn);
    }
```

```
    public BloomFilter<String> createBloomFilter(String csvFilePath, int
expectedInsertions) {
        BloomFilter<String> bloomFilter = BloomFilter.create(
                Funnels.unencodedCharsFunnel(),
                expectedInsertions
        );

        List<String> isbnList = bookMapper.selectAllISBNs();
        for (String isbn : isbnList) {
            bloomFilter.put(isbn);
        }
        return bloomFilter;
    }
```
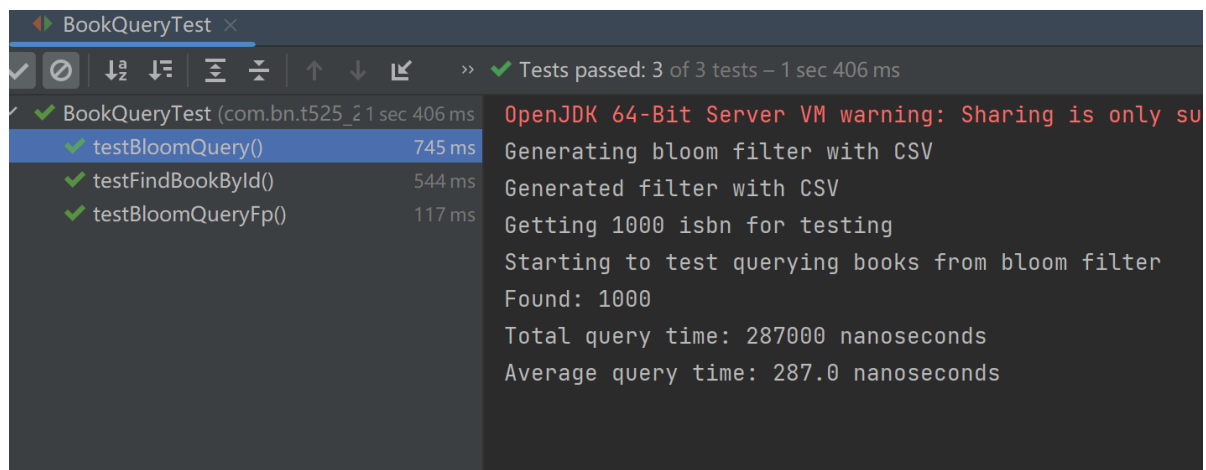
## BookQueryTest

- `testFindBookById()` tests the performance of MySQL query. For 1000 sample isbns, we verify that all of them are in the database, return the times it requires in total, and then calculate the average time for it.

- `public void testBloomQuery()` test the performance of finding 1000 isbns that should exist in the database. Return the total and the average times it takes.

- `public void testBloomQueryFp()` test the performance of finding 1000 isbns that DNE, and return the false positive rate, total time, and average time per query.

## MySQL query

## Bloom Filter Query



```
BookQueryTest ×

Tests passed: 3 of 3 tests – 1 sec 406 ms

BookQueryTest (com.bn.t525_2  1 sec 406 ms    OpenJDK 64-Bit Server VM warning: Sharing is only su
   testBloomQuery()              745 ms       Generating bloom filter with CSV
   testFindBookById()            544 ms       Generated filter with CSV
   testBloomQueryFp()            117 ms       Getting 1000 isbn for testing
                                              Starting to test querying books from bloom filter
                                              Found: 1000
                                              Total query time: 287000 nanoseconds
                                              Average query time: 287.0 nanoseconds
```

## Bloom Filter False Positive Rate



```
Run:    BookQueryTest ×

Tests passed: 3 of 3 tests – 1 sec 406 ms

BookQueryTest (com.bn.t525_2  1 sec 406 ms    Generating bloom filter with CSV
   testBloomQuery()              745 ms       Generated filter with CSV
   testFindBookById()            544 ms       Getting 1000 isbn for testing
   testBloomQueryFp()            117 ms       Starting to test false positive rate of the bloomfilter
                                              Fpr: 0.002
                                              Total query time: 201200 nanoseconds
                                              Average query time: 201.2 nanoseconds
```

## Results

The results were strikingly in favor of the Bloom filter approach, showing a remarkable 1300 times speed increase over the MySQL queries. Moreover, the Bloom filter demonstrated a very low false positive rate of only 0.002, indicating a high level of accuracy in its query responses. These findings underscore the significant efficiency and reliability advantages of using Bloom filters for database querying in high-volume contexts.

# LSH

**The goal for this section is to create a book recommendation system using the idea of lsh**

The development is divided into two steps:

- Train a model in python, making it possible to make recommendations based on similarities between book features: title, author, and publisher.

- Serialize the model and integrate the model to our application

## Training

The model is trained in python, with tools Word2Vec and Doc2Vec.

- The same csv is imported with pandas.
- Tokenize ['Book_Title'], ['Book_Author'] and ['Publisher']
- Convert each book's entities into vectors
- Add the new row df['vector']

```python
def vectorize_book(row):
    print("vectorize_book")
    # Averaging Word2Vec vectors
    word_vecs = [word2vec_model.wv[word] for word in
gensim.utils.simple_preprocess(row['Book_Title'])]
    word_vecs += [word2vec_model.wv[word] for word in
gensim.utils.simple_preprocess(row['Book_Author'])]
    word_vecs += [word2vec_model.wv[word] for word in
gensim.utils.simple_preprocess(row['Publisher'])]
    avg_word_vec = np.mean(word_vecs, axis=0)

    # Doc2Vec vector
    doc_vec = doc2vec_model.infer_vector(
        row['Book_Title'].split() + row['Book_Author'].split() +
row['Publisher'].split())

    # Combine Word2Vec and Doc2Vec vectors
    combined_vec = np.hstack((avg_word_vec, doc_vec))
    print("finished vectorize book")
    return combined_vec
```

## Quick Test With FASTAPI

- Before utilize the model to our application, test with python.
- Start a local server at port 8000

```python
@app.post("/suggest_books")
def suggest_books_api(book_request: BookRequest):
    book_titles = suggest_books(book_index=book_request.book_index,
top_n=book_request.top_n)
    return book_titles

def suggest_books(book_index, top_n=5):
    book_vec = df.iloc[book_index]['vector']
    similarities = cosine_similarity([book_vec], df['vector'].tolist())[0]
    similar_indices = np.argsort(similarities)[::-1][1:top_n + 1]
    return df.iloc[similar_indices]['Book_Title']
```

- Test Result

  Works properly!

- Test with Spring Application

  The controller visits the FastAPI portal

  Access from localhost:8090 will invoke a call in the spring application that send a request to FastAPI (localhost:8000).

```java
@RestController
@RequestMapping("/books")
public class BookController {

    private final BookService bookService;

    @Autowired
    public BookController(BookService bookService) {
        this.bookService = bookService;
    }

    @PostMapping("/suggest")
    public List<String> suggestBooks(@RequestBody SuggestBooksRequest
request) {
        return bookService.suggestBooks(request.getBook_index(),
request.getTop_n());
    }
}
```

- Test result

  Work properly!

**http://localhost:8090/books/suggest**

GET ⌄    http://localhost:8090/books/suggest          Send ⌄

Params    Authorization    Headers (9)    Body ●    Pre-request Script    Tests    Settings          Cookies

○ none    ○ form-data    ○ x-www-form-urlencoded    ● raw    ○ binary    ○ GraphQL    JSON ⌄          Beautify

```
1  {
2      "book_index": 513,
3      "top_n": 10
4  }
```

Body    Cookies    Headers (6)    Test Results          🌐 405 Method Not Allowed   115 ms   304 B   💾 Save as example  ○

Pretty    Raw    Preview    Visualize    JSON ⌄    ⇥

```
1  {
```

Console          All Logs ⌄    Clear    ⧉    ○○○

an-Token:   U910a100-55C0-4104-D510-91D10D/e9//1
"localhost:8000"
t-Encoding: "gzip, deflate, br"
ction: "keep-alive"
nt-Length: "41"
t Body ↗
se Headers
se Body ↗

764":"whats your chinese love sign","3340":"you cant go home again","59397":"you are a millionaire choose
own adventure no 98","1591":"vitamins and your health","59180":"you are a monster choose your own adventu
o 84","338":"first six months getting together with your baby","55314":"ireland your only place","5115":"3
ays to cook pasta","95536":"keeping your cool while sharing your faith","75085":"how you were born"}

## Integrate LSH to Spring Application

Now we have managed to have a recommendation model. We just need to integrate it with our application.

## How Things Work

- The model is serialized, and imported into our application
- Run the cosine similarities locally, get the indexes of books to be recommended.
- Query Database and return the book titles.

## Model Importation

- Ideally, the model should be uploaded to Redis, and our application simply import the model on the server. However, for testing purposes and conveniences, we just serialize into json and store the .json file locally (which is equivalent upload and retrieve from Redis as Redis store data in .json format as well).

- Load the .json model from `/bookcsv` directory.

- In Vector Service, load data:

```java
public void loadVectorData() {
    try {
        ObjectMapper mapper = new ObjectMapper();
        String jsonFilePath = "bookcsv/vector_column.json";
```

```
        // Read JSON file and convert to List of Lists
        bookVectors = mapper.readValue(
                new File(jsonFilePath),
                new TypeReference<List<List<Double>>>() {
                }
        );



    } catch (Exception e) {
        e.printStackTrace();
    }


}
```

## Make Recommendations

- We use cosine similarity to make comparisons:

```
    public static double cosineSimilarity(double[] vectorA, double[] vectorB)
{

        double dotProduct = 0.0;
        double normA = 0.0;
        double normB = 0.0;
        for (int i = 0; i < vectorA.length; i++) {
            dotProduct += vectorA[i] * vectorB[i];
            normA += Math.pow(vectorA[i], 2);
            normB += Math.pow(vectorB[i], 2);
        }
        return dotProduct / (Math.sqrt(normA) * Math.sqrt(normB));
    }
```

- Return a list of book indices given a index

```
    public List<Integer> compareSimilarity(int index) {
        Map<Integer, Double> similarityScores = new HashMap<>();

        double[] targetVector = bookVectors.get(index).stream().mapToDouble(d
-> d).toArray();

        for (int i = 0; i < bookVectors.size(); i++) {
            if (i != index) {
                double[] compareVector =
bookVectors.get(i).stream().mapToDouble(d -> d).toArray();
                double similarity = cosineSimilarity(targetVector,
compareVector);
                similarityScores.put(i, similarity);
            }
        }

        List<Integer> sortedIndices = new ArrayList<>
(similarityScores.keySet());
        sortedIndices.sort((i1, i2) ->
similarityScores.get(i2).compareTo(similarityScores.get(i1)));
```

```
        return sortedIndices;
    }
```

- Make suggestions

```
public List<Book> bookSuggestion(String ISBN, int numOfBooks){
    List<Book> bookList = new ArrayList<>();
    int bookIndex = allISBN.indexOf(ISBN);
    List<Integer> sim = compareSimilarity(bookIndex);
    for (int i = 0; i< numOfBooks; i++){
        Book rBook = bookMapper.selectByISBN(allISBN.get(sim.get(i)));
        bookList.add(rBook);

    }

    return bookList;
}
```

## Test

The average query time for

- Comparing similarities and return book indexes

- Look up for the books

add up to be 0.0189 seconds



## Results

The development process of our book recommendation system was effectively carried out. The performance of the system is noteworthy, with the combined average query time for both comparing similarities (to return book indexes) and subsequently looking up the recommended books clocking in at just 0.0189 seconds. This efficient performance highlights not only the technical proficiency of the model but also its practical effectiveness in enhancing the user experience by providing rapid and relevant book recommendations.

# Conclusion

In addressing the challenge of managing high-volume database queries and book recommendations in an environment with 10,000 concurrent clients, the integration of Bloom filters and Locality-Sensitive Hashing (LSH) has proven to be exceptionally effective. The Bloom filter implementation drastically improved query efficiency, demonstrating a remarkable 1300-fold increase in speed compared to traditional MySQL queries, along with a remarkably low false positive rate. This enhancement significantly reduces the database load and ensures swift and accurate stock availability checks. Concurrently, the use of LSH for book recommendations allowed for the rapid generation of pertinent suggestions, with an impressive average query time of just 0.0189 seconds. This dual-faceted approach not only addresses the initial concerns of database strain and user experience but does so in a manner that is both technically proficient and practically impactful, ensuring a seamless and responsive interaction for the users. These results validate our strategy as a highly effective solution to the challenges posed by high client traffic and complex query requirements in a dynamic book retail context.