# Building a Knowledge Graph Enriched With Large Language Models
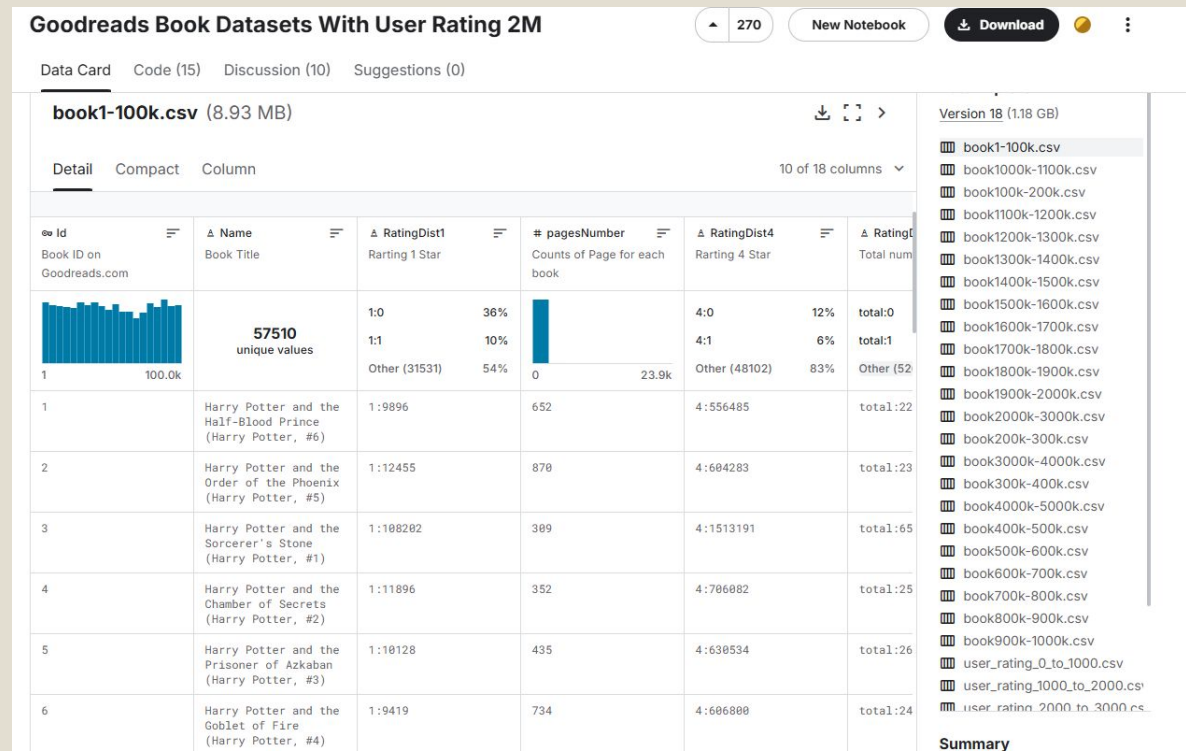
Oğuzhan Güngör 24411006

# Goal of the Presentation

To explain the current state of the project and present the deliverables.

The project is creating a knowledge graph and enhance the Large Language Model (LLM).

Using the knowledge graph to enrich the LLM and link it to DBpedia for enhanced semantic connections and information retrieval.

# Data Preparation

In this project, Goodreads Book Dataset has been used for generating the initial database.

# Data Pipeline Methods

There were **3** approaches when storing data from csv file to Neo4j graph database:

1)   Utilizing Concurrency in Python

2)   Utilizing Multiprocessing in Python

3)   Utilizing Concurrency at Database Level

# Data Pipeline Methods Pros and Cons

Concurrency in Python: It is fast. Comes with race conditions, managing connection pools is hard. Lacks data integrity.

Multiprocessing in Python: Easier solution. It is slow. Don't need to manage connection pools. Data integrity is full.

Concurrency in Neo4j using APOC: It is fastest. Comes with race conditions therefore comes with failed insertions. Lacks data integrity. This method was picked due to its accuracy.

# Data Pipeline Methods- APOC

```
1  CALL apoc.periodic.iterate(
2    'LOAD CSV WITH HEADERS FROM "file:///18/book100k-200k.csv" AS row RETURN row',
3    '
4    CALL apoc.do.when(
5      row.Id IS NULL OR row.Authors IS NULL,
6      "RETURN null",
7      "
8      MERGE (b:Book {id: toInteger(row.Id)})
9      SET b.name = row.Name, b.language = row.Language, b.publisher = row.Publisher, b.publishMonth = row.PublishMonth, b.rating = row.Rating, b.ISBN = row.ISBN
10     MERGE (a:Author {name: row.Authors})
11     MERGE (b)-[:WRITTEN_BY]→(a)
12     ", {row: row})
```

| batches | failedBatches | total | timeTaken | committedOperations |
|---------|---------------|-------|-----------|---------------------|
| 115 | 23 | 57046 | 420 | 45546 |

Started streaming 1 records after 1 ms and completed after 420491 ms.

# Environment Generation

For the development environment, Docker container has been utilized for easier deployment of Neo4j. For Python scripts, Python 3.12 has been picked due to speed benefits of Python 3.12.

For version control, Git and Github has been utilized. Since the project is a solo project with time constraints, no branch controls have been implemented.

# Graph Version 1:
# Generated With Only The Dataset

This version made with only using the contents of the data. In this version, graph database have meaningful connections but lacks the depth.

Nodes:
Books, Authors, Users

Edges:
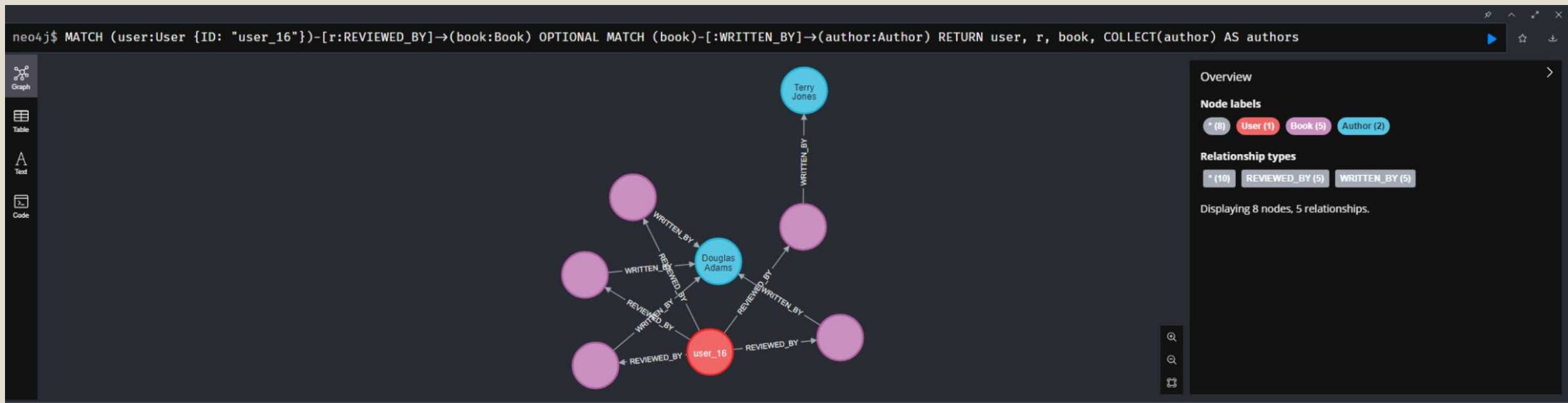WRITTEN_BY: Book -> Author
REVIEWED_BY: Book -> User

# Graph Version 1:
## Generated With Only The Dataset

# Graph Version 1:
## Generated With Only The Dataset

# Graph Version 1: Generated With Only The Dataset

# DBpedia Integration

Since the data for books are lacking Genre of the books, the awards that author or the book got using DBpedia to obtain genres of the books was the solution for improving the graph database.

# Graph Version 11:
# Improved Version With DBpedia

This version made with addition to DBpedia data. In this version, graph database have meaningful connections but still lacks the depth.

Improvements:
Added Genre, Award, Descriptions

Nodes:
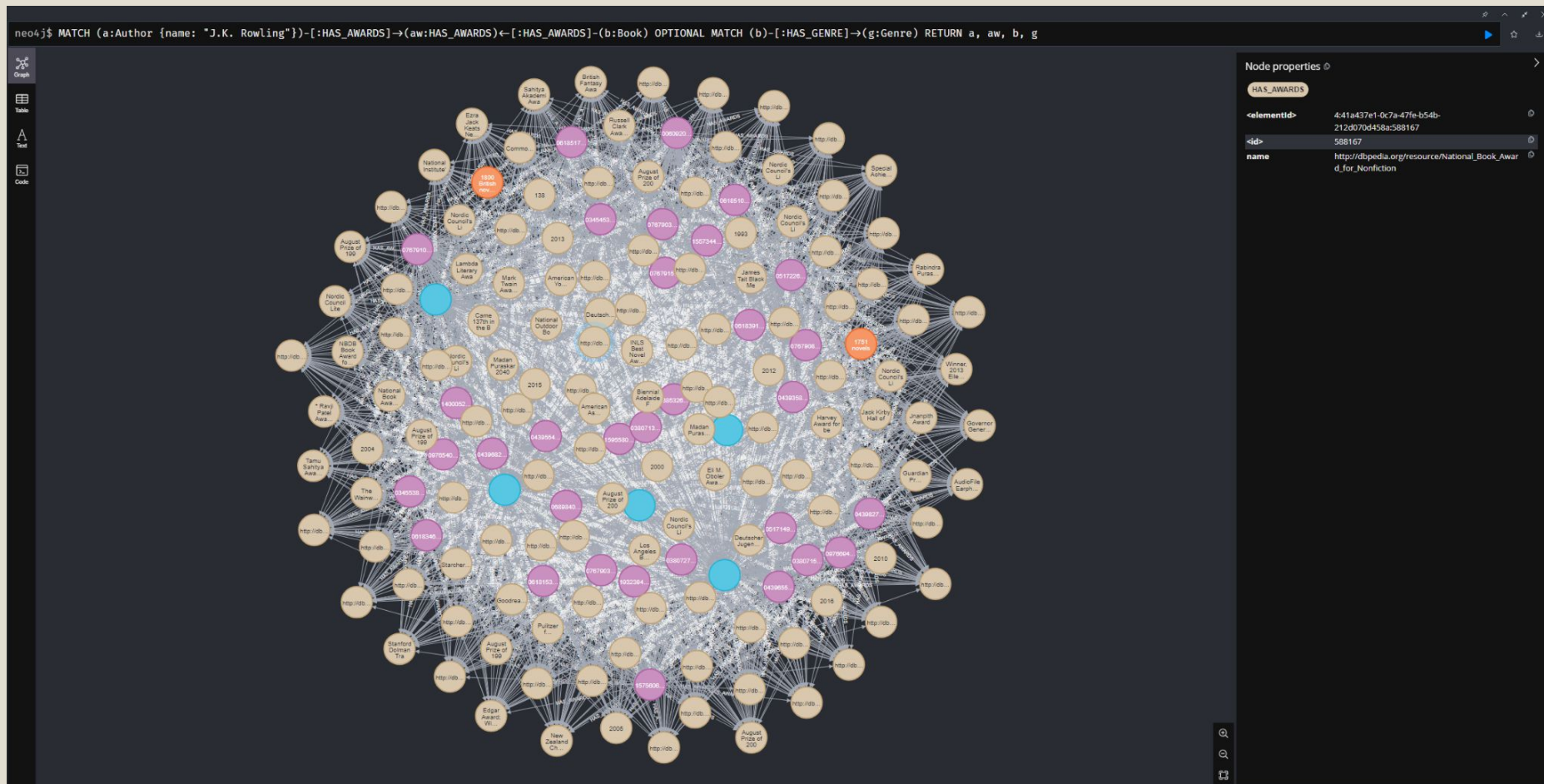Books, Authors, Users, Genre, Award

Edges:
WRITTEN_BY: Book -> Author
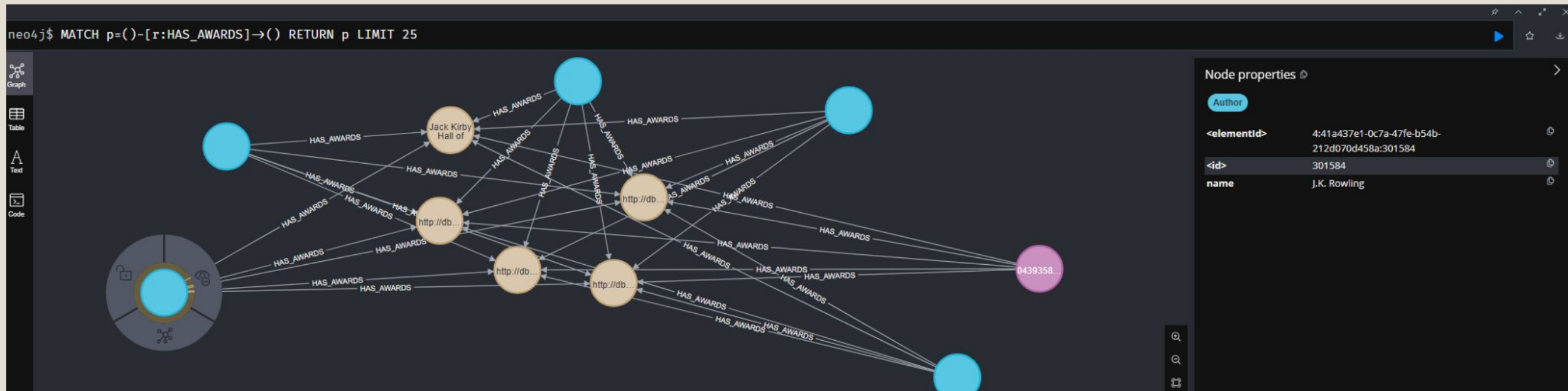REVIEWED_BY: Book -> User
HAS_GENRE: Book -> Genre
HAS_AWARDS: Book -> Award

# Graph Version 11:
# Improved Version With DBpedia

# Graph Version 11:
# Improved Version With DBpedia

# Graph Version 11:
## Improved Version With DBpedia
## (Bugged Version)



```
neo4j$ MATCH (b:Book) WHERE b.description IS NOT NULL RETURN b.id, b.description LIMIT 5
```

| b.id | b.description |
|------|---------------|
| "2" | "Cabal is a 1988 horror novel by the British author Clive Barker. It was originally published in the United States as part of a collection comprising a novel and several short stories from Barker's sixth and final volume of the Books of Blood. The book was adapted into the film Nightbreed in 1990, written and directed by |
| "4" | "Cabal is a 1988 horror novel by the British author Clive Barker. It was originally published in the United States as part of a collection comprising a novel and several short stories from Barker's sixth and final volume of the Books of Blood. The book was adapted into the film Nightbreed in 1990, written and directed by |
| "5" | "Cabal is a 1988 horror novel by the British author Clive Barker. It was originally published in the United States as part of a collection comprising a novel and several short stories from Barker's sixth and final volume of the Books of Blood. The book was adapted into the film Nightbreed in 1990, written and directed by |
| "8" | "Cabal is a 1988 horror novel by the British author Clive Barker. It was originally published in the United States as part of a collection comprising a novel and several short stories from Barker's sixth and final volume of the Books of Blood. The book was adapted into the film Nightbreed in 1990, written and directed by |
| "9" | "Cabal is a 1988 horror novel by the British author Clive Barker. It was originally published in the United States as part of a collection comprising a novel and several short stories from Barker's sixth and final volume of the Books of Blood. The book was adapted into the film Nightbreed in 1990, written and directed by |

Started streaming 5 records after 10 ms and completed after 58 ms.

# Graph Version 11:
## Improved Version With DBpedia
## (Fixed Version)

```
1  MATCH (n) WHERE (n.description) IS NOT NULL
2  RETURN DISTINCT "node" as entity, n.description AS description LIMIT 25
3  UNION ALL
4  MATCH ()-[r]-() WHERE (r.description) IS NOT NULL
5  RETURN DISTINCT "relationship" AS entity, r.description AS description LIMIT 25
```

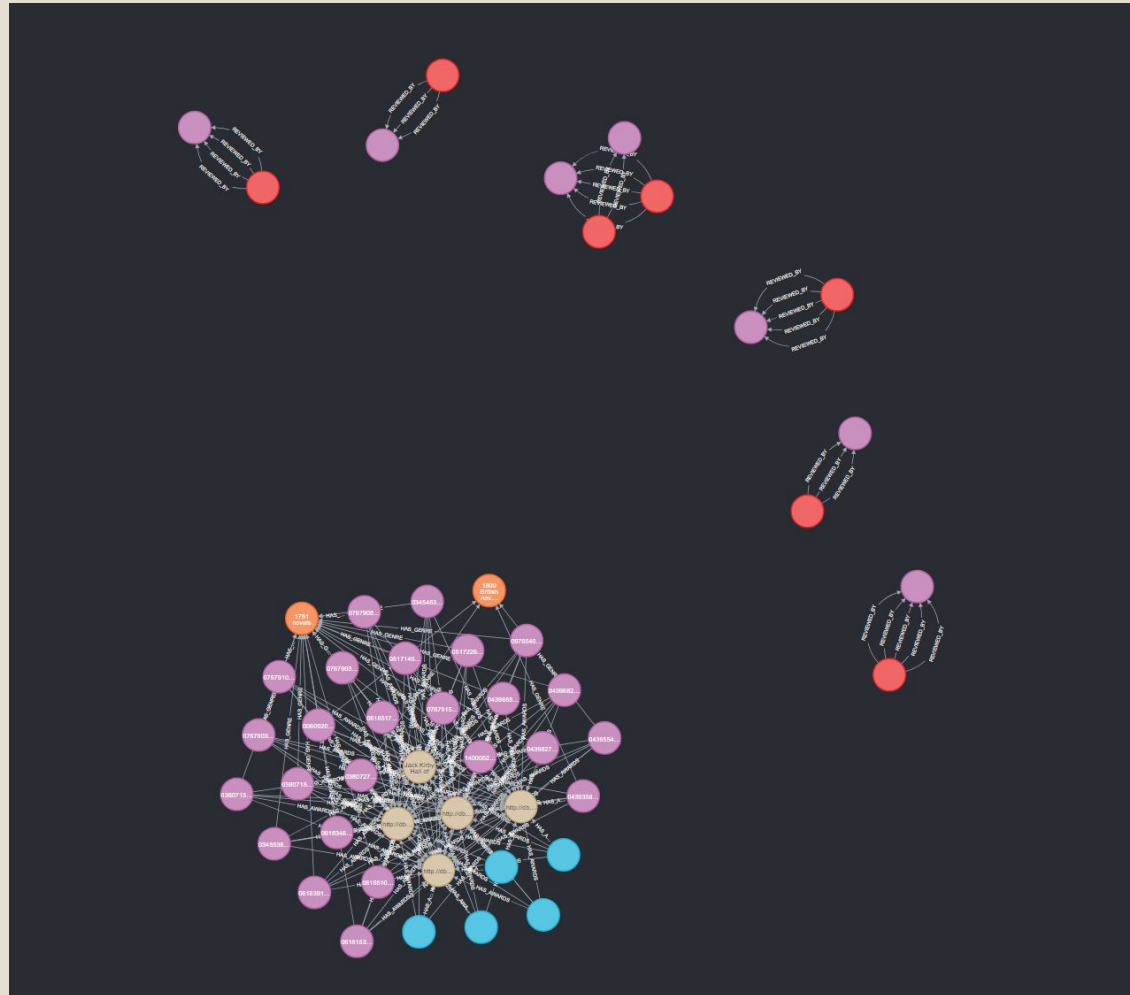| entity | description |
|---|---|
| "node" | "Cabal is a 1988 horror novel by the British author Clive Barker. It was originally published in the United States as part of a collection comprising a novel and several short stories from Barker's sixth and final volume of the Books of Blood. The book was adapted into the film Nightbreed in 1990, written and directed by Barker himself, starring Craig Sheffer and David Cronenberg |
| "node" | "A Short History of Nearly Everything by American-British author Bill Bryson is a popular science book that explains some areas of science, using easily accessible language that appeals more to the general public than many other books dedicated to the subject. It was one of the bestselling popular science books of 2005 in the United Kingdom, selling over 300,000 copies. A S |
| "node" | "Bill Bryson's African Diary is a 2002 book by bestselling travel writer Bill Bryson. The book details a trip Bryson took to Kenya in 2002. Bryson describes his experiences there and observations about Kenyan culture, geography, and politics, as well as his visits to poverty-fighting projects run by CARE International, to which he donated all royalties for the book." |
| "node" | "Neither Here nor There: Travels in Europe is a 1991 humorous travelogue by American writer Bill Bryson. It documents the author's tour of Europe in 1990, with flashbacks to two summer tours he made in 1972 and 1973 in his college days. On his 1973 tour, he travelled with his friend Matt Angerer, pseudonymised in the book as Stephen Katz, who also appeared more prom |
| "node" | "Notes from a Small Island is a humorous travel book on Great Britain by American author Bill Bryson, first published in 1995." |
| "node" | "The Changeling Sea is a fantasy novel for juvenile readers by Patricia A. McKillip. It was first published in hardcover by Atheneum/Macmillan in October 1988, with a paperback edition issued by Del Rey/Ballantine in December 1989. It was subsequently reissued in paperback and ebook by Firebird/Penguin in April 2003. The first British edition was published in hardcover by |
| "node" | "The Known World is a 2003 historical novel by Edward P. Jones. Set in Virginia during the antebellum era, it examines the issues regarding the ownership of Black slaves by both white and Black Americans. The book was published to acclaim, which praised its story and Jones's prose. In particular, his ability to intertwine stories within stories received great praise from The N |

# Graph Version 11:
# Improved Version With DBpedia

# LLM Integration

Since the data for books are lacking semantic relationships with comments and other books, utilizing LLM to improve the graph database was crucial.

This step is work in progress at the moment. This step was missing the crucial connection between similar books.

In this step, sentiment analysis was performed using multiple LLMs, mainly BART and OpenAI. In the sentiment analysis, partial matching is also supported.

# Graph Version III:
## Improved Version With LLM Integration
## BART Only

```
You: What's the sentiment for "Gatsby"
INFO:__main__:Running query:

    MATCH (u:User)-[r:REVIEWED_BY]->(b:Book)
    WHERE toLower(b.name) CONTAINS toLower($name)
    RETURN u.name AS userName, r.review AS review

PARAMS: { name: 'Gatsby' }
Chatbot:
Reviews for any book name containing 'Gatsby':
  - User: user_128
    Review: "liked it"
    BART Sentiment => POSITIVE
  - User: user_128
    Review: "it was amazing"
    BART Sentiment => POSITIVE
  - User: user_128
    Review: "really liked it"
    BART Sentiment => POSITIVE
```

# Graph Version III:
# Improved Version With LLM Integration
# BART and OpenAI

```
Chatbot:
Reviews for any book name containing 'Gatsby':
  - User ID 4:41a437e1-0c7a-47fe-b54b-212d070d458a:585786:
    Review: "liked it"
    [OpenAI] => POSITIVE (4 stars)
    [BART]   => POSITIVE (4 stars)
  - User ID 4:41a437e1-0c7a-47fe-b54b-212d070d458a:585786:
    Review: "it was amazing"
    [OpenAI] => POSITIVE (5 stars)
    [BART]   => POSITIVE (4 stars)
  - User ID 4:41a437e1-0c7a-47fe-b54b-212d070d458a:585786:
    Review: "really liked it"
    [OpenAI] => POSITIVE (4 stars)
    [BART]   => POSITIVE (2 stars)
```

# Storage



Devices and drives
OS-Boot (C:)
4,57 GB free of 237 GB

# References

1) [Goodreads Book Datasets With User Rating 2M](#)
2)

Thanks for listening