# Graph Databases

## Using Neo4j and Amazon SNAP Dataset

Database Management Systems Laboratory (CS39202)

## Team CheeseCake

Shivam Choudhury

Ayush Mundada

Tanishka Rahate

Omkar Bhandare

Jadhav Udaykiran

# Problem Statement

The objective of this project is to process large graphs using a graph database. The core tasks include:

i. Install any graph processing system, such as ApacheGraph, Pregel (GoldenOrb), Giraph, Neo4j, or Stanford GPS.

ii. Load a large graph from the Stanford SNAP large graph repository.

iii. Provide an interface to run simple graph queries. Bonus marks may be awarded for implementing PageRank.

iv. Profile and evaluate the performance of the system.

# Background

Graph databases represent a powerful paradigm in database technology that focuses on relationships between data entities. Graph databases are designed to efficiently store, process, and query highly connected data using graph structures with nodes, edges, and properties.

They consist of three fundamental elements:

- **Nodes**: Represent entities or instances such as products, categories, or other items. They function similar to records, like in relational databases.

- **Edges**: Connect nodes, representing the relationships between entities. They can either be directed or undirected, and they often carry properties that describe the nature of the relationship.

- **Properties**: Information associated with nodes and edges that describe their characteristics and attributes.

At its core, a graph database portrays data as it exists conceptually in the real world. For example, in a social network like Twitter, each user would be represented as a node, with "follows" relationships connecting them. This intuitive representation makes graph databases particularly effective for modelling complex, interconnected data.

They are built on graph theory principles but do not require users to understand the mathematical complexities. Instead, they offer an intuitive way to model data that closely resembles how we think about relationships.

They typically use specialised query languages designed for traversing relationships.

Cypher, developed for Neo4j, is one of the most popular graph query languages, allowing users to express complex relationship patterns in a relatively simple syntax.

Graph databases are well-suited for e-commerce applications like the Amazon product recommendation system, which is the one attempted here. Because of their natural affinity for relationship-based data, graph databases are widely adopted across various industries. In modern applications, they play a pivotal role where understanding context and connection is essential - such as recommendation engines, fraud detection systems, and social media analytics. They enable systems to go beyond flat data queries and derive insights through patterns, proximity, and relevance. Their schema flexibility also allows easy evolution with dynamic datasets, a necessity in fast-paced domains like e-commerce and finance. As a result, graph databases are increasingly becoming a foundational choice for building intelligent, adaptive, and scalable data-driven systems.

# Objective

The primary objective of this project is to design and implement a robust graph-based product recommendation system leveraging the Amazon SNAP dataset and Neo4j as the core database technology. The system aims to enhance user experience through intuitive navigation, advanced search capabilities, and intelligent, data-driven recommendations by capitalising on the inherent interconnectedness within product data. This project positions graph databases as a powerful alternative to traditional recommendation systems by exploiting relationships, patterns, and influence within product networks.

To achieve this goal, the project focuses on the following key objectives:

1. **Efficient Product Discovery via Hierarchical Structures**: Implement mechanisms for seamless product discovery by utilising categorical and group-based relationships inherent in the dataset.

2. **Customised PageRank Algorithm for Product Relevance**: Develop and integrate a tailored PageRank-inspired algorithm to rank products based on relevance, popularity, and contextual importance within the graph.

3. **Advanced Multi-Category Search Functionality**: Enable flexible and precise search across multiple product categories simultaneously, allowing for complex and user-specific query resolution.

4. **Performance Monitoring Framework**: Establish a comprehensive performance evaluation and monitoring framework to ensure system scalability, reliability, and responsiveness under varying loads.

5. **User Influence Scoring Mechanism**: Introduce a user influence scoring model that quantifies the impact of user behaviour on product visibility and recommendation outcomes.

6. **Temporal Analysis of Product Trends**: Analyse temporal dynamics in product popularity and user interaction to uncover evolving trends and seasonality within the dataset.

7. **Visualisation of Product Networks**: Design interactive visualisations to represent the structure and relationships among products, categories, and user interactions within the graph.

8. **Integration of a RAG-Based Natural Language Query System**: Incorporate a Retrieval-Augmented Generation (RAG) based natural language processing interface to facilitate intuitive and conversational product discovery.

# Methodology

The development of the graph-based product recommendation system follows a structured and multi-faceted approach, encompassing key stages such as data modelling, algorithm design, system implementation, and user interface development. Each phase is carefully designed to align with the overall objective of leveraging graph structures for enhanced product discovery and recommendation. The methodology ensures a coherent integration of data engineering, graph analytics, and user-centric design to deliver a scalable and effective recommendation platform.

# Data Schema Design

The underlying data schema is designed to effectively represent the complex relationships among products, consumers, categories, and groups within the Amazon SNAP dataset. This schema serves as the structural backbone of the recommendation system, enabling efficient querying, relationship traversal, and algorithmic computations.

**Node Types (along with the attributes)**

- **Product Node**: Id, ASIN, title, salesrank, avg_rating, intrinsic_score, total_score

- **Group Node**: name

- **Category Node**: name

- **Consumer Node**: Id, Customer, name, score

**Edge Types**

- **BELONGS_TO:** Product Node → Category Node

- **PART_OF:** Product Node → Group Node

- **REVIEWED:** Consumer Node → Product Node
  Attributes: Date, Rating, Votes, Helpful

- **SIMILAR:** Product Node → Product Node

- **COPURCHASED_WITH:** Product Node → Product Node
  Attributes: Frequency

## Indexing Strategy

Given the dataset's variety and volume of queries, implementing appropriate indexes is critical for ensuring efficient data access and optimal performance. The following indexes were established, each tailored to anticipated query patterns:

| Name | Type | Entity Type | Labels or Types | Properties |
|------|------|-------------|-----------------|------------|
| `category_name_index` | RANGE | NODE | {Category} | {name} |
| `consumer_customer_index` | RANGE | NODE | {Consumer} | {Customer} |
| `group_name_index` | RANGE | NODE | {Group} | {name} |
| `product_asin_index` | RANGE | NODE | {Product} | {ASIN} |
| `product_id_index` | RANGE | NODE | {Product} | {Id} |

Table 1: Indexing Strategy

# Custom PageRank Algorithm

To accurately reflect the dynamics of product relevance and interconnectedness within the Amazon SNAP dataset, we designed and implemented a novel two-phase PageRank algorithm tailored specifically for graph-based product recommendation. Unlike traditional PageRank—which primarily relies on nodes' structural connectivity—our approach integrates product-specific attributes and co-purchase behaviour to derive a more context-aware and impactful relevance score.

## Phase 1: Computing the Intrinsic Score

$$\text{intrinsic\_score} = (2^{\text{avg\_rating}}) \cdot \left( \frac{\log_{10}(1 + \text{sum of all helpful})}{\log_{10}(2)} \right) \tag{1}$$

This formulation provides a meaningful balance: exponential scaling of the rating gives prominence to top-rated products. In contrast, logarithmic scaling of helpful votes accounts for user trust without allowing it to dominate the score.

## Phase 2: Incorporating Network Influence via Co-Purchase Relationships

$$\text{total\_score} = \text{intrinsic\_score} + 0.5 \cdot \sum \left( \frac{\text{frequency}_j}{\text{total\_frequency}} \cdot \text{total\_score}_j \right) \tag{2}$$

*Note: Only co-purchased products with frequency $\geq 2$ are considered in the summation.*

## Cypher Implementation

### Phase 1: Intrinsic Score Calculation

```
MATCH (p:Product)
OPTIONAL MATCH (u:Consumer)-[r:REVIEWED]->(p)
WITH p, AVG(r.Rating) AS avg_rating, SUM(r.Helpful) AS total_helpful
SET p.intrinsic_score = round(
    CASE
        WHEN avg_rating IS NULL THEN 0.0
        ELSE (2.0 ^ avg_rating) * (LOG(1.0 + COALESCE(total_helpful, 0)) /
    LOG(2))
    END
, 2)
RETURN count(p) AS updated_count
```

### Phase 2: Incorporating Co-Purchase Influence

```
MATCH (p:Product)
WITH p, p.intrinsic_score AS intrinsic_score
OPTIONAL MATCH (p)-[co:COPURCHASED_WITH]->(other:Product)
 WHERE co.Frequency >= 2
WITH p, intrinsic_score,
    collect({
        product: other,
        frequency: co.Frequency,
        score: other.total_score
    }) AS co_purchases
WITH p, intrinsic_score, co_purchases,
```

```
    REDUCE(total = 0, item IN co_purchases | total + item.frequency) AS
    total_frequency
WITH p, intrinsic_score,
    CASE WHEN total_frequency > 0 THEN
        REDUCE(sum = 0, item IN co_purchases |
            sum + (item.frequency / toFloat(total_frequency)) * item.score
        )
    ELSE 0 END AS co_purchase_score
SET p.total_score = round(intrinsic_score + co_purchase_score * 0.5, 2)
RETURN count(p) AS updated_count
```

Only co-purchase edges with frequency $\geq 2$ are considered, effectively filtering out noise and incidental relationships. The final `total_score` represents a balanced combination of intrinsic product quality and network influence, with a damping factor (0.5) applied to avoid disproportionate influence from the co-purchase network.

# User Influence Scoring

To assess the influence of individual reviewers within the system, we developed a custom algorithm that calculates a user-specific influence score based on the helpfulness of their reviews. This approach prioritises users who consistently contribute meaningful and helpful reviews across various products while mitigating the potential impact of outliers.

## Formula

The user influence score is calculated as follows:

$$\text{score} = \sum \left( \min \left( \frac{\log_{10}(1 + \text{helpful})}{\log_{10}(2)}, \ 5.0 \right) \right) \tag{3}$$

## Cypher Implementation

```
MATCH (c:Consumer)-[r:REVIEWED]->()
WHERE c.Customer IS NOT NULL AND r.Helpful IS NOT NULL
RETURN DISTINCT c
MATCH (c)-[r:REVIEWED]->()
WHERE r.Helpful IS NOT NULL
WITH c, collect(r.Helpful) AS helpful_values
WITH c, [val IN helpful_values |
    CASE
```

```
    WHEN val <= 0 THEN 0.0
    ELSE
        CASE
        WHEN log10(1 + val) / log10(2) > 5.0 THEN 5.0
        ELSE toFloat(log10(1 + val) / log10(2))
        END
    END
] AS logs
WITH c, reduce(score = 0.0, x IN logs | score + x) AS final_score
SET c.score = final_score
```

For each consumer, the algorithm computes a base score for each review using logarithmic scaling, caps this score at 5.0 to prevent exceptionally high helpfulness values from disproportionately affecting the overall influence score, and sums up these capped scores across all reviews to derive the final influence score. This method ensures that users who consistently contribute helpful reviews across different products are rewarded, while logarithmic scaling prevents a few highly helpful reviews from dominating the score.

# Trending Product Analysis

To track and identify trending products, we developed a functionality that focuses on identifying products that have seen a notable surge in popularity within a specified time frame. The algorithm performs a time-based analysis of review activity for each product. This approach considers the number of reviews as the primary ranking attribute. It combines it with the product's total score, derived from the custom PageRank algorithm, as the secondary ranking attribute. This dual-ranking system prioritises products with a high volume of recent reviews and strong recommendation scores. By incorporating the quantity and quality of reviews, this analysis provides a robust measure of product popularity, making it highly effective for identifying products trending due to recent customer engagement.

### Formula

$$\text{TrendingScore}(p) = \text{sort}\left(\text{review\_count}_{[t_1,t_2]}(p),\ \text{total\_score}(p)\right) \tag{4}$$

*where $[t_1, t_2]$ represents the time window within which reviews are counted. Products are sorted first by review count (descending), and then by total score (descending) in case of ties.*

## Cypher Implementation

```
MATCH (c:Consumer)-[r:REVIEWED]->(p:Product)
WHERE date(r.Date) >= date($start_date) AND date(r.Date) <= date($end_date)
WITH p, count(r) AS review_count, p.total_score AS total_score
ORDER BY review_count DESC, total_score DESC
LIMIT 10
RETURN p.title AS product_title, p.ASIN AS product_asin, review_count,
    total_score
```

# Multi-Category Search

To empower users with more refined product discovery, we implemented a Multi-Category Search feature that allows seamless filtering across multiple product categories. This functionality is handy when users are interested in items falling under overlapping categories. To make category selection intuitive, especially when users may not know a category's exact spelling or phrasing, we integrated fuzzy matching using the TheFuzz library. Users can type approximate or partial terms and retrieve relevant category names.

## Cypher Implementation

```
MATCH (p:Product)
WHERE EXISTS {
  MATCH (p)-[:BELONGS_TO]->(:Category {name: $categories[0]})
}
AND ALL(category IN $categories[1..] WHERE
    EXISTS {
        MATCH (p)-[:BELONGS_TO]->(:Category {name: category})
    })
AND p.avg_rating >= $min_avg_rating
RETURN p.title AS title, p.ASIN AS asin, p.avg_rating AS rating,
      p.total_score AS total_score
ORDER BY p.total_score DESC
SKIP $skip
LIMIT $limit
```

# RAG-Based Chatbot Implementation

We developed a Retrieve and Generate (RAG)-based chatbot to enhance natural language product discovery, leveraging advanced embedding techniques and efficient similarity search mechanisms. The core of this implementation utilises the 'BAAI/beg-small-en-v1.5' model to generate embeddings for product titles, groups, and categories. These embeddings are then indexed using FAISS (Facebook AI Similarity Search), accelerating the similarity search process and ensuring efficient retrieval of relevant products in response to user queries. This approach effectively translates user queries into embeddings, performs a similarity search across a range of products, categories, and groups, and returns the top matches ranked by similarity. A similarity threshold of 0.3 ensures that only highly relevant results are returned, providing users with accurate and meaningful recommendations.

The RAG engine integrates multiple components to process user queries efficiently:

- **Text Cleaning and Preprocessing:** User input is first cleaned using the `cleantext` library, removing noise such as URLs, emails, and unnecessary punctuation. A stopword filtering process is also applied to ensure that only the most relevant terms remain in the query.

- **Embedding-Based Search:** After cleaning the query, it is passed through an embedding engine to obtain a vector representation. The system then searches for similar products, groups, and categories and ranks them based on their cosine similarity to the query vector.

- **Results Filtering and Response Generation:** The top matching products, categories, and groups are retrieved and enriched with additional product details (e.g., ASIN, title, sales rank). A well-structured response is generated to ensure an engaging user experience, highlighting the most relevant results.

# Performance Monitoring Framework

We have integrated a comprehensive performance monitoring framework to ensure optimal system efficiency and transparency across our application. This framework provides consistent and real-time tracking of key system metrics for every route and query execution.

**Key Metrics Tracked**

- **Query Execution Time (in seconds):** Measures the duration to execute individual route handlers and database interactions.

- **CPU Usage (percentage):** Captures the CPU consumption during processing.

- **Memory Consumption (in megabytes):** Tracks the memory footprint before and after route execution to detect potential memory leaks or spikes.

This monitoring system was implemented using a custom decorator, `@log_performance`, which wraps all relevant route handler functions. This approach ensures minimal intrusion into the application logic while enabling uniform performance tracking across the board.

Each time a route is executed, the decorator logs the stats:



Figure 1: Performance Log Entries

# User Interface Components

The user interface has been thoughtfully designed to provide an intuitive, flexible, and engaging browsing experience across multiple dataset dimensions. Below is a detailed overview of the key components:

- **Home Page:** The landing page presents all available Groups in the dataset. As groups encapsulate broader themes or collections, this is a natural entry point for

product exploration.

- **Products by Group:** Displays all products under a selected Group, sorted by default using the `total_score` in descending order to first surface the most relevant or high-performing products. Users can dynamically sort the product list by either `total_score`, `title`, `avg_rating`, or `salesrank`. Additional filters, such as minimum average rating, are also available to refine results.

- **Products by Categories:** This view displays all products belonging to a selected Category, again sorted by default using `total_score` in descending order. It shares the same interactive capabilities as above.

- **Products by Multiple Categories:** This page enables the discovery of products that belong simultaneously to multiple categories. It leverages a fuzzy-matching powered search bar and returns products that match all selected categories. Sorting and filtering options identical to other product views.

- **Product Detail Page:** Provides a comprehensive view of a product, including core product details, associated categories, co-purchased and similar products, and a complete list of user reviews with metadata.

- **Top Reviewers Page:** Ranks the top reviewers based on a custom reviewer scoring algorithm. This helps surface influential voices within the dataset and provides users with a trusted source for quality opinions.

- **Trending Products Page:** Displays products identified as trending using our internal scoring and trend-detection logic. These may be determined through review activity, popularity spikes, or recent high ratings.

- **User Reviews Page:** Showcases the most helpful reviews written by a selected top reviewer, ranked by the helpful count metric. This promotes valuable user-generated insights and supports more informed decision-making.

## Query Complexity Analysis

Let the number of nodes and edges of defined types be as follows:

| | |
|---|---|
| **Node Types:** | Category = $N$, Consumer = $M$, Group = $G$, Product = $P$ |
| **Edge Types:** | BELONGS_TO = $B$, COPURCHASED_WITH = $C$, PART_OF = $X$, REVIEWED = $R$, SIMILAR = $S$ |

**Finding all Groups (sorted):** $O(G)$

**Products in the same Group (sorted):** $O(1 + X + P \cdot \log P)$

**Products in the same Category (sorted):** $O(1 + B + P \cdot \log P)$

**Products in multiple selected Categories:** $O(N \cdot P + P \cdot \log P)$

**Trending Products:** $O(R + P \cdot \log P)$

**Top Reviewers (with score calculation):** $O(M \cdot \log M + R + M)$

**PageRank score calculation:** $O(P + R + P + C) = O(P + R + C)$

**Reviews of a selected top Reviewer:** $O(R + R \cdot \log R)$

**Selected Product data (all the linked info):** $O(1 + N + S + S \cdot \log S + C + C \cdot \log C + R + R \cdot \log R)$

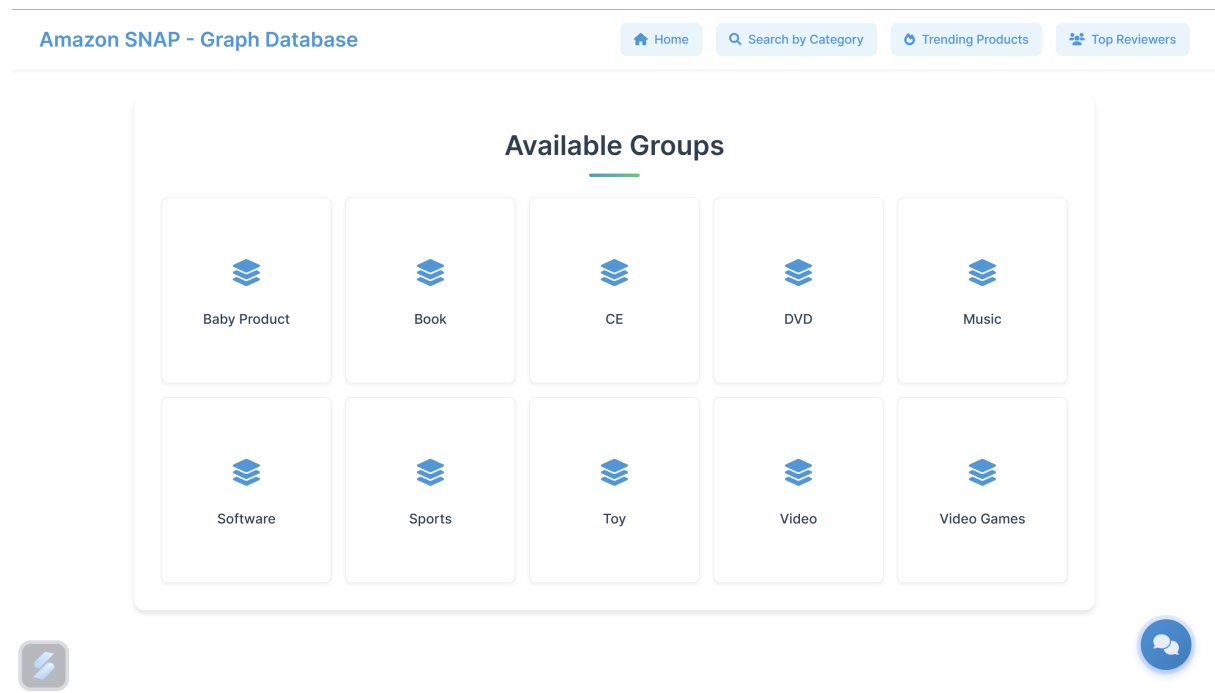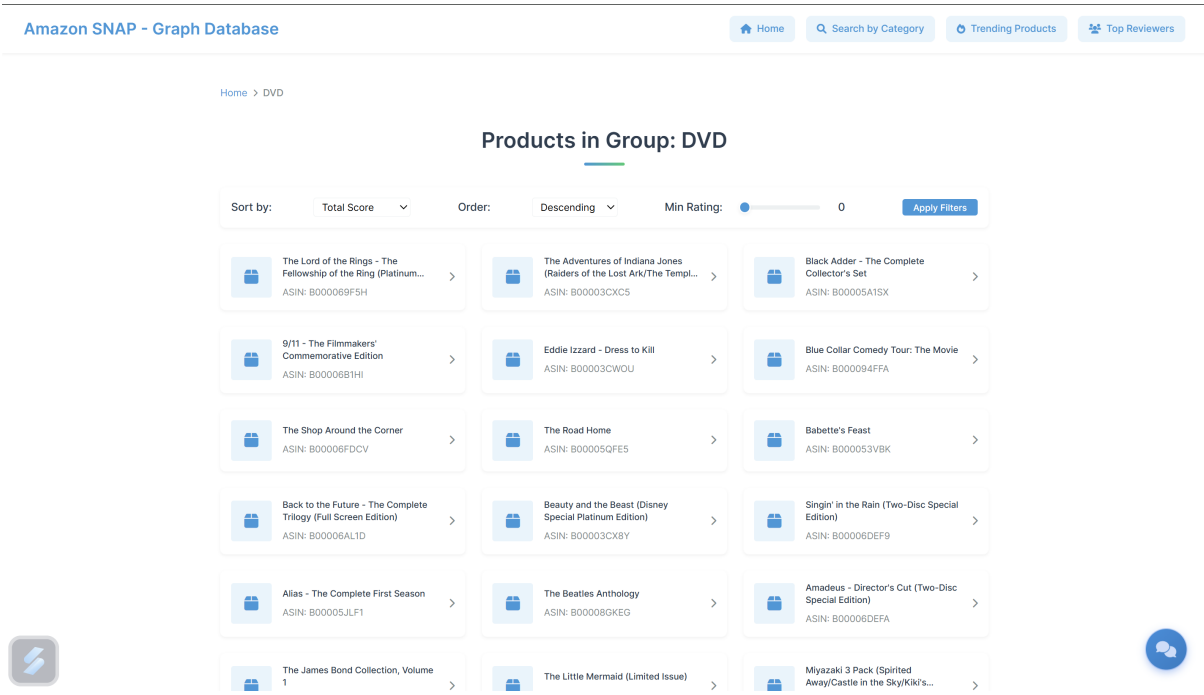# Screenshots of the Demonstration



Figure 2: Home Page

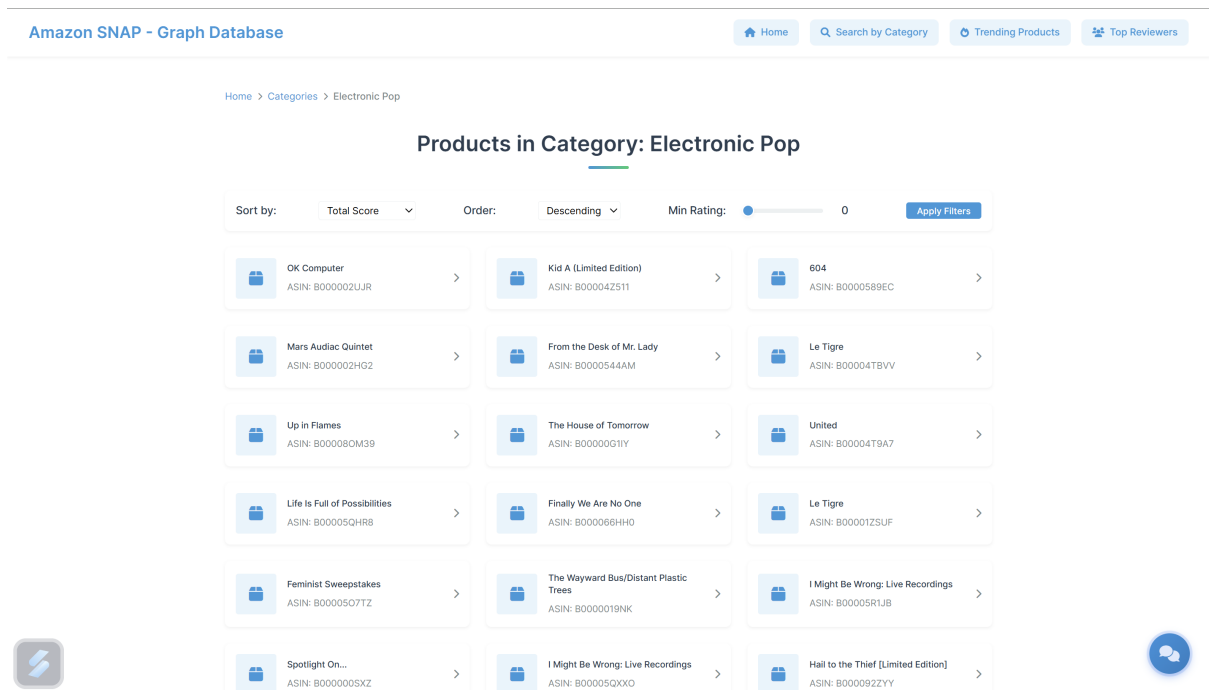Figure 3: Products in selected Group
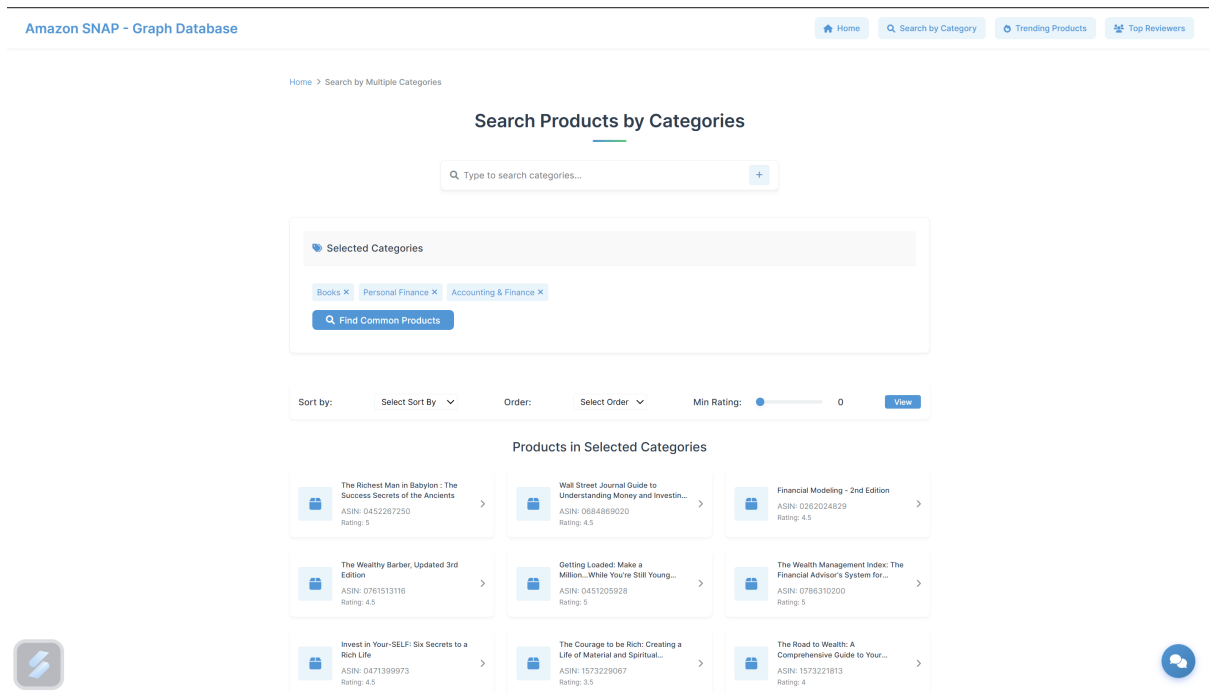


Figure 4: Products in selected Category

Figure 5: Products in selected Categories



Figure 6: Product Details
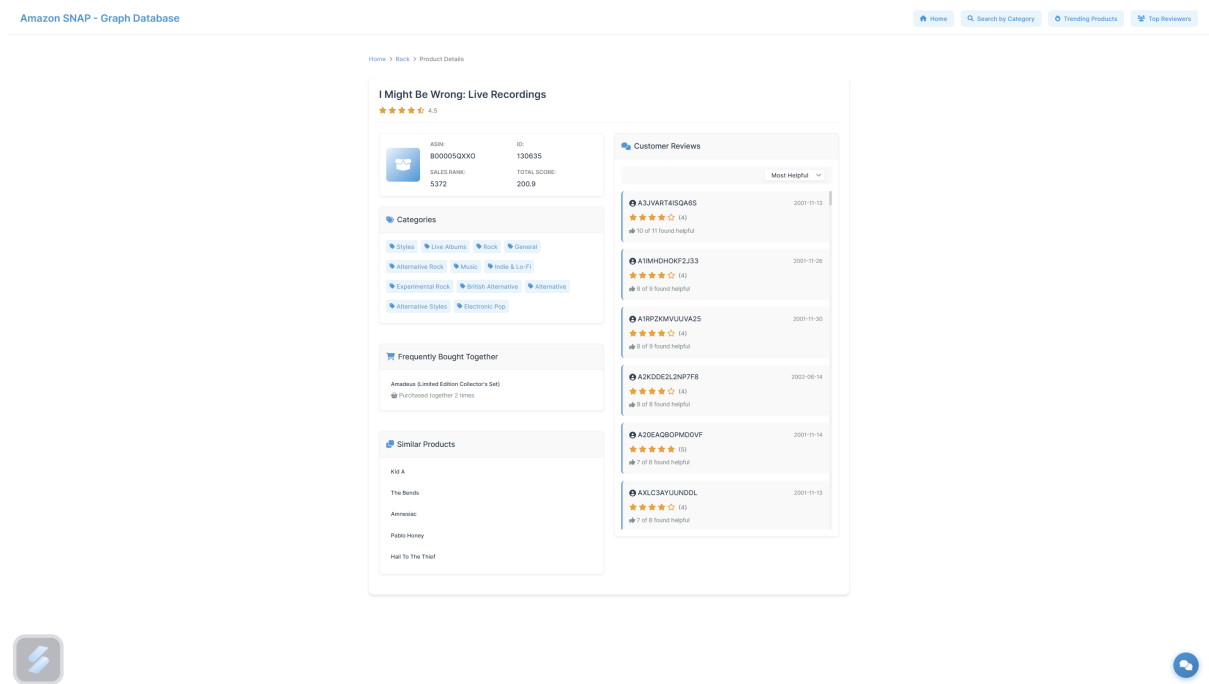
Home > Trending Products

## Trending Products

| | | |
|---|---|---|
| 🔥 Angels & Demons<br>ASIN: 0671027352 — #1 | 🔥 Dance With My Father<br>ASIN: B000099J41 — #2 | 🔥 St. Anger (with Bonus DVD)<br>ASIN: B00008OWZG — #3 |
| 🔥 Angels & Demons<br>ASIN: 0671027360 — #4 | 🔥 Regarding the Fountain (An Avon Camelot Book)<br>ASIN: 0380793474 — #5 | 🔥 Ten<br>ASIN: B0000027RL — #6 |
| 🔥 Star Wars - Episode II, Attack of the Clones (Full Screen Edition)<br>ASIN: B00006HBUI — #7 | 🔥 Star Wars - Episode II, Attack of the Clones (Widescreen Edition)<br>ASIN: B00006HBUJ — #8 | 🔥 To Kill a Mockingbird<br>ASIN: 0446310786 — #9 |
| 🔥 Eragon (Inheritance, Book 1)<br>ASIN: 0966621336 — #10 | | |

Figure 7: Trending Products

Home > Top Reviewers

## Top Reviewers

1. ATVPDKIKX0DER
2. A3UN6WX5RRO2AG
3. A14OJS0VWMOSWO
4. A2NJO6YE954DBH
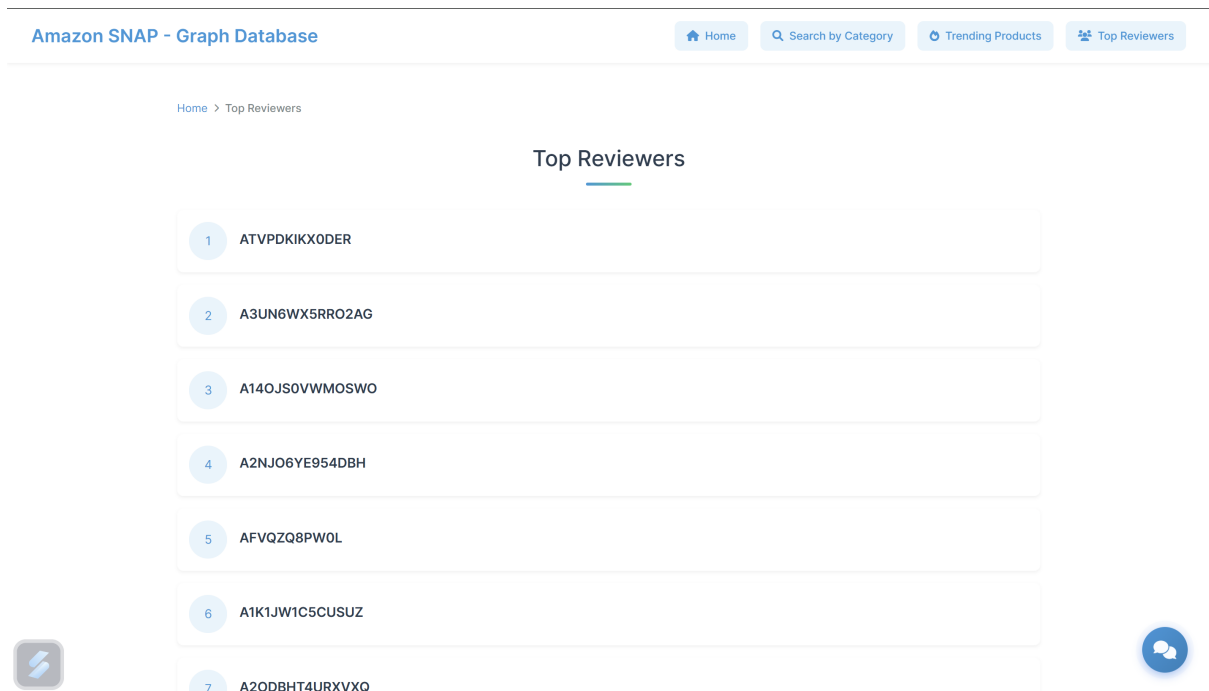5. AFVQZQ8PW0L
6. A1K1JW1C5CUSUZ
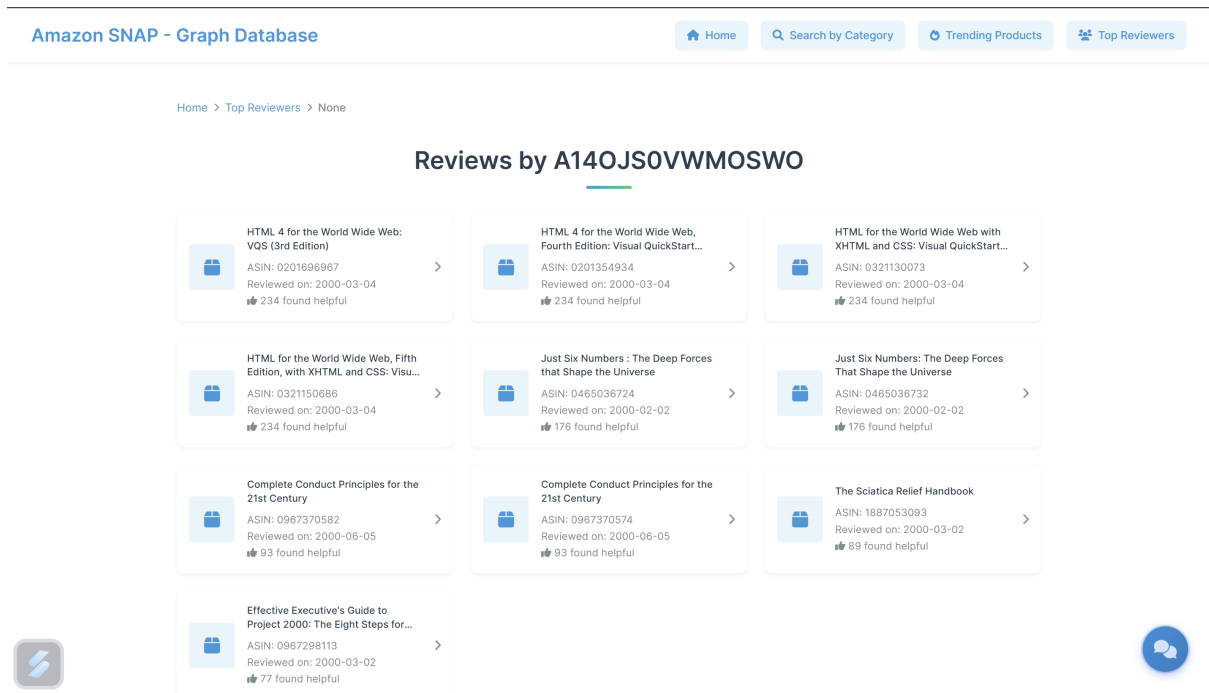7. A2ODBHT4URXVXQ

Figure 8: Top Reviewers

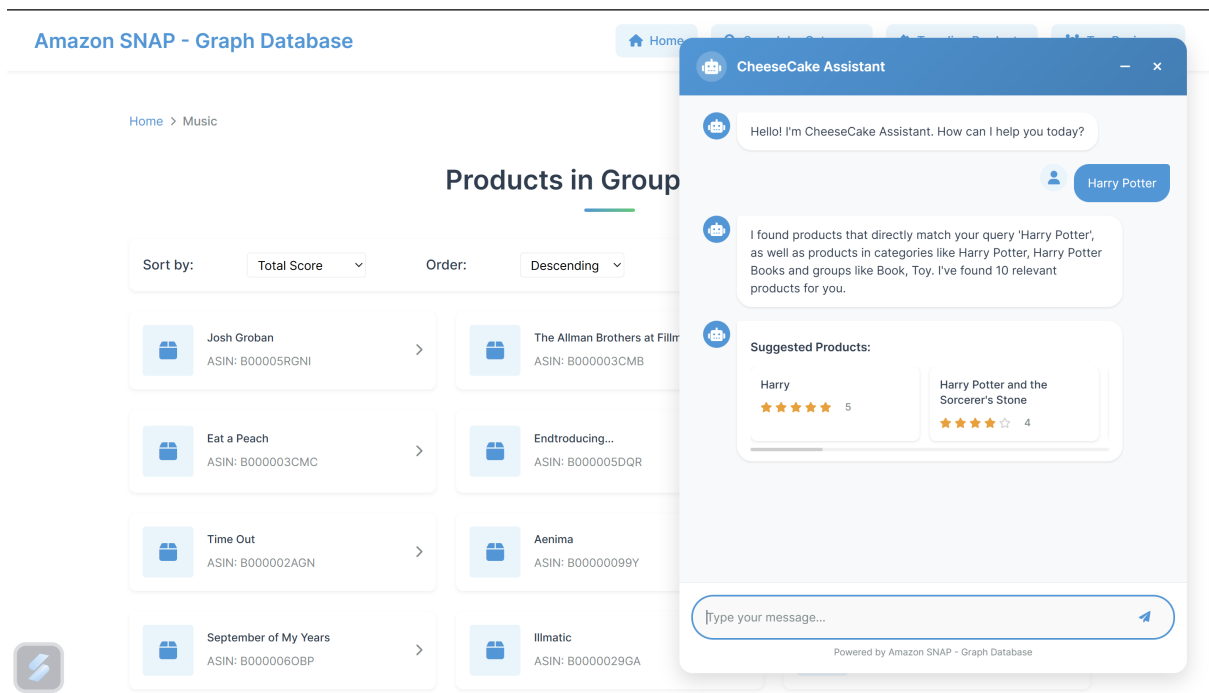Figure 9: Best Reviewed Products of top Reviewers



Figure 10: Integrated RAG-based ChatBox

# Future Scope

## 1. Advanced Hybrid Recommendation Techniques

- Add knowledge graph reasoning for explainable recommendations.

- Integrate real-time session-based recommendations using streaming graph updates.

2. **Enhanced Graph Learning**

- Implement Graph Neural Networks (GNNs) for dynamic relationship modelling.

- Add temporal graph networks for dynamic tracking of evolving product popularity.

- Explore heterogeneous graph transformers for multi-modal data.

3. **Scalability Improvements**

- Implement Graph Partitioning for distributed Neo4j clusters.

- Add incremental graph updates for real-time score recalculations.

- Explore vector-graph hybrid search combining FAISS with Cypher.

4. **RAG System Enhancements**

- Integrate multi-modal embeddings.

- Add conversational memory for chat-based recommendations.

- Implement query intent analysis using graph-aware LLMs.

5. **Social Network Integration**

- Add social graph features for friend-based recommendations.

- Track cross-product network effects through user communities.

# References

- `https://en.wikipedia.org/wiki/Graph_database`

- `https://aws.amazon.com/nosql/graph/`

- `https://youtu.be/y7sXDpffzQQ?si=m-o_b-sak4HKdg29`

- `https://neo4j.com/docs/getting-started/`

- `https://neo4j.com/docs/getting-started/appendix/tutorials/guide-cypher-basics/`

- `https://youtu.be/8jNPelugC2s?si=9TsH3yZc_4r3oeld`

- `https://youtu.be/REVkXVxvMQE?si=SKpIsYs08T-kQVqO`

- `https://youtu.be/oRtVdXvtD3o?si=0CooMOpvgB2FxWyP`