

# Database Schema Manual

324827328, 213713522

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>The Core Entities</b>	<b>2</b>
2.1	Movies . . . . .	2
2.2	People . . . . .	3
<b>3</b>	<b>Relationships</b>	<b>3</b>
3.1	The Cast & Crew . . . . .	3
3.2	Categorization (Genres, Keywords) . . . . .	3
<b>4</b>	<b>Production Details</b>	<b>4</b>
<b>5</b>	<b>User Interaction</b>	<b>4</b>
5.1	Ratings . . . . .	4
<b>6</b>	<b>Binary Storage</b>	<b>4</b>
6.1	Storing Images in the DB . . . . .	4
<b>7</b>	<b>Application Logic Queries</b>	<b>5</b>
7.1	Movie Retrieval . . . . .	5
7.2	Movie Details Enrichment . . . . .	5
7.3	Search . . . . .	5
7.4	Rating System . . . . .	5
7.5	Assets . . . . .	6
<b>8</b>	<b>Analytics &amp; Reporting</b>	<b>6</b>
8.1	Genre Popularity . . . . .	6
8.2	Actor Reuse (Prolific Actors) . . . . .	6
8.3	Runtime Analysis . . . . .	6
8.4	Financial Underperformance ("Flops") . . . . .	6
8.5	Director-Actors . . . . .	7
8.6	Global Reach . . . . .	7
8.7	Content Richness . . . . .	7
8.8	Genre Complexity . . . . .	7
8.9	Golden Era . . . . .	7
8.10	Multiskilled Crew . . . . .	8

# 1 Introduction

The database design was a critical component of the project, establishing a foundation for API performance and application integrity. We aimed for a normalized schema to minimize redundancy while maintaining the flexibility required for analytical queries.

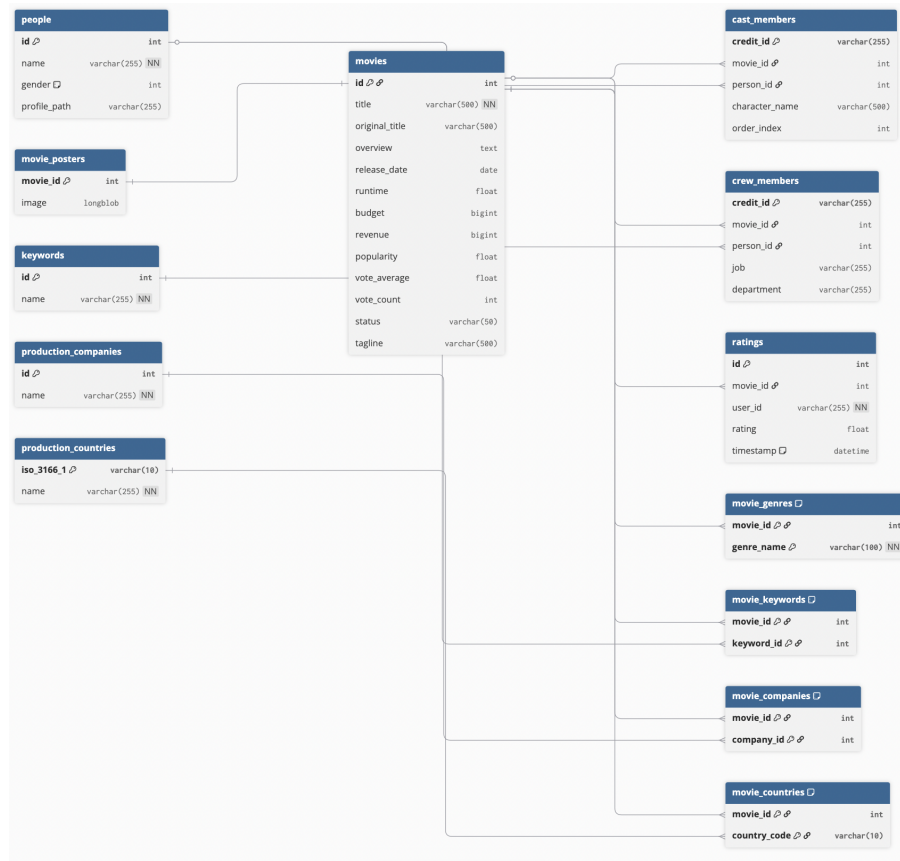


Figure 1: Full Database Schema Diagram

## 2 The Core Entities

### 2.1 Movies

The **movies** table is the central entity of the database. We explicitly included **vote\_average** and **vote\_count** fields in this table.

This denormalization strategy was chosen to optimize read performance. Instead of calculating averages from the **ratings** table for every query, we maintain these aggregated values transactionally.

```
1 CREATE TABLE IF NOT EXISTS movies (  
2     id INT PRIMARY KEY,  
3     title VARCHAR(500) NOT NULL,  
4     -- We keep the aggregated stats right here for speed  
5     vote_average FLOAT,  
6     vote_count INT,  
7     -- ... other metadata  
8     budget BIGINT,  
9     revenue BIGINT,  
10    status VARCHAR(50)  
11 );
```

Listing 1: The heart of the system

## 2.2 People

We unified actors, directors, and writers into a single `people` table. This approach avoids data duplication for individuals who fulfill multiple roles (e.g., both acting and directing).

```
1 CREATE TABLE IF NOT EXISTS people (  
2     id INT PRIMARY KEY,  
3     name VARCHAR(255) NOT NULL,  
4     gender INT,  
5     profile_path VARCHAR(255)  
6 );
```

Listing 2: A unified People table

## 3 Relationships

The database relational structure defines the connections between entities.

### 3.1 The Cast & Crew

Connecting movies to people required more than a simple link. We needed to know *what role* they played. For the **Cast**, we stored the `character_name` and an `order_index`. The index is crucial because it allows the frontend to validly display the "Top Billed" actors first, rather than a random alphabetical list. For the **Crew**, it's about the job and `department`. This granular detail allows us to answer interesting questions like "Which directors also acted in their own movies?"

```
1 CREATE TABLE IF NOT EXISTS cast_members (  
2     credit_id VARCHAR(255) PRIMARY KEY,  
3     movie_id INT,  
4     person_id INT,  
5     character_name VARCHAR(500),  
6     order_index INT,  
7     FOREIGN KEY (movie_id) REFERENCES movies(id),  
8     FOREIGN KEY (person_id) REFERENCES people(id)  
9 );
```

Listing 3: Linking People to Movies

### 3.2 Categorization (Genres, Keywords)

We handled categories like Genres and Keywords as Many-to-Many relationships. This means a movie can be both "Action" and "Comedy". We separated the names into their own lookup tables (`genres`, `keywords`) and used junction tables (`movie_genres`) to link them. This makes searching for "All Action movies" incredibly fast because it's just a simple index lookup.

```
1 CREATE TABLE IF NOT EXISTS movie_genres (  
2     movie_id INT,  
3     genre_name VARCHAR(100) NOT NULL,  
4     PRIMARY KEY (movie_id, genre_name),  
5     FOREIGN KEY (movie_id) REFERENCES movies(id)  
6 );
```

Listing 4: Genre Association

```
1 CREATE TABLE IF NOT EXISTS keywords (  
2     id INT PRIMARY KEY,  
3     name VARCHAR(255) NOT NULL  
4 );
```

Listing 5: Keywords Definition

```
1 CREATE TABLE IF NOT EXISTS movie_keywords (  
2     movie_id INT,  
3     keyword_id INT,  
4     PRIMARY KEY (movie_id, keyword_id),  
5     FOREIGN KEY (movie_id) REFERENCES movies(id),  
6     FOREIGN KEY (keyword_id) REFERENCES keywords(id)  
7 );
```

Listing 6: Keyword Association

## 4 Production Details

To provide a comprehensive view of the film's origin, we track both the production companies and the countries involved. These are modeled as many-to-many relationships, acknowledging that modern films are often co-productions.

```
1 CREATE TABLE IF NOT EXISTS production_companies (  
2     id INT PRIMARY KEY,  
3     name VARCHAR(255) NOT NULL  
4 );  
5  
6 CREATE TABLE IF NOT EXISTS movie_companies (  
7     movie_id INT,  
8     company_id INT,  
9     PRIMARY KEY (movie_id, company_id),  
10    FOREIGN KEY (movie_id) REFERENCES movies(id),  
11    FOREIGN KEY (company_id) REFERENCES production_companies(id)  
12 );
```

Listing 7: Production Companies

```
1 CREATE TABLE IF NOT EXISTS production_countries (  
2     iso_3166_1 VARCHAR(10) PRIMARY KEY,  
3     name VARCHAR(255) NOT NULL  
4 );  
5  
6 CREATE TABLE IF NOT EXISTS movie_countries (  
7     movie_id INT,  
8     country_code VARCHAR(10),  
9     PRIMARY KEY (movie_id, country_code),  
10    FOREIGN KEY (movie_id) REFERENCES movies(id),  
11    FOREIGN KEY (country_code) REFERENCES production_countries(iso_3166_1)  
12 );
```

Listing 8: Production Countries

## 5 User Interaction

### 5.1 Ratings

We built the `ratings` table. We decided to use a string-based `user_id` to align with our "Soft Session" strategy on the client. This table grows the fastest, so we kept it lightweight: just who, what, what score, and when.

```
1 CREATE TABLE IF NOT EXISTS ratings (  
2     id INT AUTO_INCREMENT PRIMARY KEY,  
3     movie_id INT,  
4     user_id VARCHAR(255) NOT NULL,  
5     rating FLOAT,  
6     timestamp DATETIME DEFAULT CURRENT_TIMESTAMP,  
7     FOREIGN KEY (movie_id) REFERENCES movies(id)  
8 );
```

Listing 9: Capturing user votes

## 6 Binary Storage

### 6.1 Storing Images in the DB

We opted to store movie posters directly in the database using the `LONGBLOB` data type. While external storage (like S3) is a common pattern, storing images within the database ensures the dataset is self-contained. This simplifies backup and deployment procedures by keeping the entire application state in a single location.

```
1 CREATE TABLE IF NOT EXISTS movie_posters (  
2     movie_id INT PRIMARY KEY,  
3     image LONGBLOB,
```

```

4 FOREIGN KEY (movie_id) REFERENCES movies(id)
5 );

```

Listing 10: Self-contained image storage

## 7 Application Logic Queries

This section documents the SQL queries used within the application's backend logic.

### 7.1 Movie Retrieval

For the main listing, we fetch movies with pagination logic.

```

1 SELECT * FROM movies LIMIT %s OFFSET %s

```

### 7.2 Movie Details Enrichment

To efficiently load related data for a list of movies (or a single movie), we use IN clauses or direct equality checks to fetch Genres, Cast, and Crew.

```

1 SELECT mg.movie_id, mg.genre_name
2 FROM movie_genres mg
3 WHERE mg.movie_id IN ({placeholders})

```

Listing 11: Fetching Genres

```

1 SELECT c.movie_id, p.id, p.name, p.gender, p.profile_path,
2        c.character_name, c.order_index
3 FROM people p
4 JOIN cast_members c ON p.id = c.person_id
5 WHERE c.movie_id IN ({placeholders})
6 ORDER BY c.movie_id, c.order_index

```

Listing 12: Fetching Cast

```

1 SELECT c.movie_id, p.id, p.name, p.gender, p.profile_path,
2        c.job, c.department
3 FROM people p
4 JOIN crew_members c ON p.id = c.person_id
5 WHERE c.movie_id IN ({placeholders})

```

Listing 13: Fetching Crew

### 7.3 Search

We implement a simple but effective text search using LIKE on the movie title.

```

1 SELECT * FROM movies WHERE title LIKE %s

```

### 7.4 Rating System

The rating system involves a transactional flow to ensure data consistency.

First, we check if the user has already rated the movie:

```

1 SELECT rating FROM ratings
2 WHERE movie_id = %s AND user_id = %s

```

Then, we perform an "Upsert" operation. If the record exists, we update the timestamp and rating; otherwise, we insert a new record.

```

1 INSERT INTO ratings (movie_id, user_id, rating)
2 VALUES (%s, %s, %s)
3 ON DUPLICATE KEY UPDATE
4     rating = VALUES(rating),
5     timestamp = CURRENT_TIMESTAMP

```

Finally, to keep the read-heavy `movies` table fast, we recalculate and update the aggregated statistics on the movie record itself.

```
1 UPDATE movies
2 SET vote_average = %s, vote_count = %s
3 WHERE id = %s
```

## 7.5 Assets

Retrieving the binary image data for a movie poster.

```
1 SELECT image FROM movie_posters WHERE movie_id = %s
```

# 8 Analytics & Reporting

The system enables complex analytical queries to derive insights from the movie database. We create 10 different queries to analyze the data. This section details the SQL formulation for these reports.

## 8.1 Genre Popularity

To determine which genres are most prevalent in the database, we aggregate the count of movies for each genre and order them descending.

```
1 SELECT genre_name AS genre, COUNT(movie_id) AS movie_count
2 FROM movie_genres
3 GROUP BY genre_name
4 ORDER BY movie_count DESC
5 LIMIT 5;
```

## 8.2 Actor Reuse (Prolific Actors)

This query identifies actors who have appeared in more than one movie. It uses a `HAVING` clause to filter groups (actors) after aggregation.

```
1 SELECT p.name as actor_name, COUNT(c.movie_id) as movie_count
2 FROM people p
3 JOIN cast_members c ON p.id = c.person_id
4 GROUP BY p.id, p.name
5 HAVING movie_count > 1
6 ORDER BY movie_count DESC
7 LIMIT 10
```

## 8.3 Runtime Analysis

We calculate the average runtime for each genre. Since a movie can have multiple genres, it contributes to the average of each genre it belongs to.

```
1 SELECT mg.genre_name AS genre, AVG(m.runtime) AS avg_runtime
2 FROM movie_genres mg
3 JOIN movies m ON mg.movie_id = m.id
4 WHERE m.runtime IS NOT NULL
5 GROUP BY mg.genre_name
6 ORDER BY avg_runtime DESC
7 LIMIT 10
```

## 8.4 Financial Underperformance ("Flops")

We identify movies with a significant budget (over \$1M) that failed to recover their costs. The loss is calculated directly in the projection as `budget - revenue`.

```
1 SELECT title, budget, revenue, (budget - revenue) as loss
2 FROM movies
3 WHERE budget > 1000000 AND revenue < budget AND revenue > 0
4 ORDER BY loss DESC
5 LIMIT 10
```

## 8.5 Director-Actors

This complex join finds individuals who have both directed and acted in the *same* movie. It requires joining people to cast\_members and crew\_members, and ensuring the movie\_id matches across both roles.

```
1 SELECT p.name, m.title
2 FROM people p
3 JOIN cast_members c ON p.id = c.person_id
4 JOIN crew_members cr ON p.id = cr.person_id AND c.movie_id = cr.movie_id
5 JOIN movies m ON c.movie_id = m.id
6 WHERE cr.job = 'Director'
```

## 8.6 Global Reach

We aggregate production countries to see which nations produce the most content in our database.

```
1 SELECT pc.name, COUNT(mc.movie_id) as movie_count
2 FROM production_countries pc
3 JOIN movie_countries mc ON pc.iso_3166_1 = mc.country_code
4 GROUP BY pc.name
5 ORDER BY movie_count DESC
6 LIMIT 10
```

## 8.7 Content Richness

We rank movies by the number of associated keywords, using this as a proxy for the depth of metadata available.

```
1 SELECT m.title, COUNT(mk.keyword_id) as keyword_count
2 FROM movies m
3 JOIN movie_keywords mk ON m.id = mk.movie_id
4 GROUP BY m.id, m.title
5 ORDER BY keyword_count DESC
6 LIMIT 10
```

## 8.8 Genre Complexity

We identify movies that are categorized into more than 3 genres. This helps find films with complex or cross-cutting themes.

```
1 SELECT m.title, COUNT(mg.genre_name) as genre_count
2 FROM movies m
3 JOIN movie_genres mg ON m.id = mg.movie_id
4 GROUP BY m.id, m.title
5 HAVING genre_count > 3
6 ORDER BY genre_count DESC
```

## 8.9 Golden Era

We determine the "Best Year" for cinema based on the average user rating. We filter for years with at least 3 movies to avoid statistical outliers from years with very few releases.

```
1 SELECT YEAR(release_date) as year, AVG(vote_average) as avg_vote, COUNT(*) as
   movie_count
2 FROM movies
3 WHERE release_date IS NOT NULL
4 GROUP BY year
5 HAVING movie_count >= 3
6 ORDER BY avg_vote DESC
7 LIMIT 1
```

## 8.10 Multiskilled Crew

We find crew members who are talented enough to work in multiple different departments (e.g., Sound and Editing). We use `COUNT(DISTINCT cr.department)` to ensure we are counting different functional areas, not just multiple credits in the same department.

```
1 SELECT p.name, COUNT(DISTINCT cr.department) as dept_count
2 FROM people p
3 JOIN crew_members cr ON p.id = cr.person_id
4 GROUP BY p.id, p.name
5 HAVING dept_count > 1
6 ORDER BY dept_count DESC
7 LIMIT 10
```