# NETFLIX TITLES DATASET

BY:

RAJAT TIWARI
KALYAN VANNE
JYOTHSNA LENKA
SAI TEJA MALLADI

# PROBLEM STATEMENT

*The problem is to analyze the Netflix titles dataset to gain insights into the content available on the platform. This includes understanding the distribution of titles by various attributes such as type (movie or TV show), genre, country, release year, and more.*

# WHY GRAPH DATABASE?

*A graph database like Neo4j is well-suited for this problem because it allows us to model complex relationships between entities easily.*

*The Netflix titles dataset contains various entities such as titles, directors, actors, etc., and their relationships with each other (e.g., a title directed by a director, a title belonging to a genre, etc.). Graph databases excel at representing such interconnected data and querying relationships efficiently.*

# MODELING THE GRAPH DATABASE

- **Title**: *Represents a movie or TV show.*
- **Director**: *Represents the director of a title.*
- **Actor**: *Represents an actor/actress appearing in a title.*
- **Country**: *Represents the country associated with a title.*

# NODES AND PROPERTIES

- **Title**: Properties may include `show_id`, `title`, `type`, `date_added`, `release_year`, `rating`, `duration`, `description`.
- **Director**: Property include `name`.
- **Actor**: Property include `name`.
- **Genre**: Property include `name`.
- **Country**: Property include `name`.

# RELATIONSHIPS

- **_DIRECTED_BY_**_: Relationship between Title and Director._
- **_ACTED_BY_**_: Relationship between Title and Actor._
- **_AVAILABLE_IN_**_: Relationship between Title and Country._
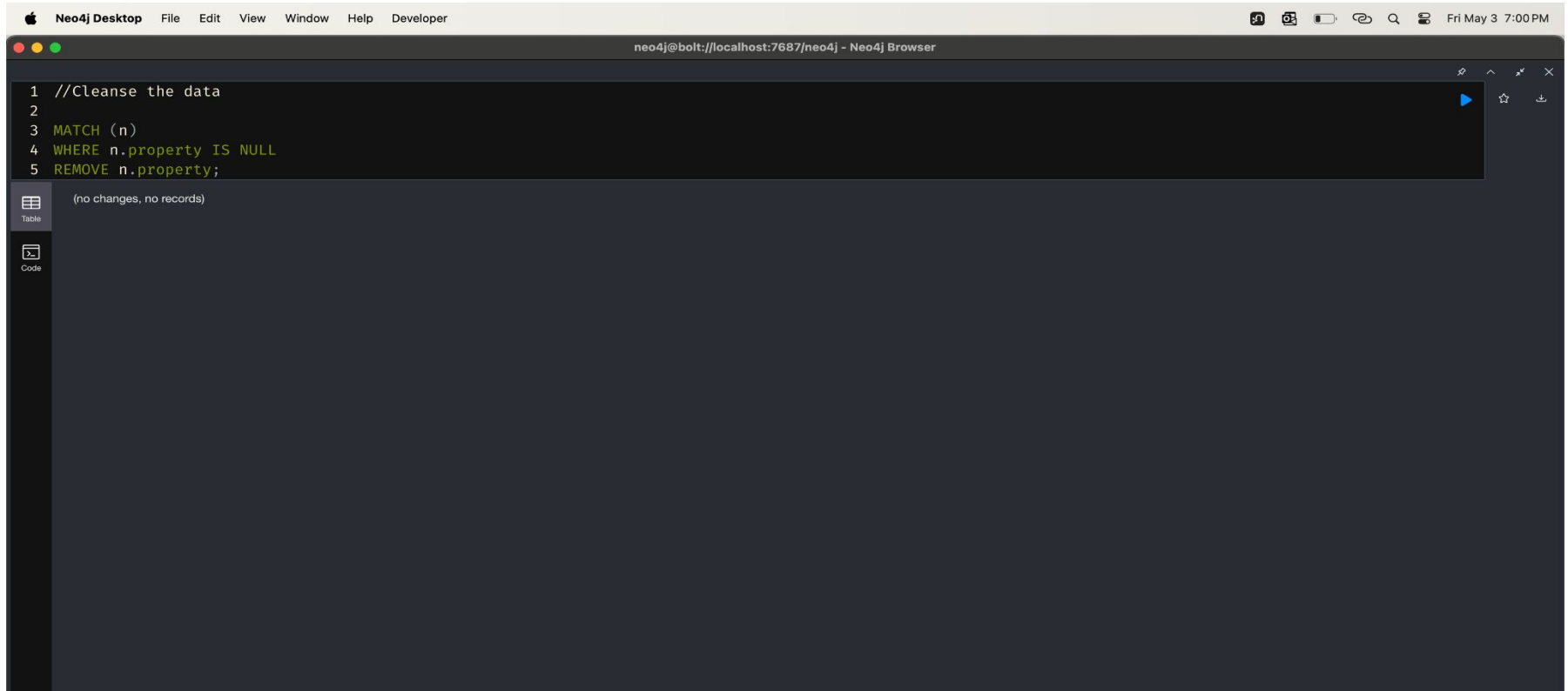
# CODING & IMPLEMENTATION

# Creating and labeling Nodes and define Relationships

```
 1  //Creating Nodes
 2  LOAD CSV WITH HEADERS FROM "file:///netflix_titles.csv" AS row
 3
 4  // Create Title nodes
 5  CREATE (:Title {
 6  show_id: row.show_id,
 7  type: row.type,
 8  title: row.title,
 9  director: row.director,
10  cast: row.cast,
11  country: row.country,
12  date_added: row.date_added,
```

```
13  release_year: row.release_year,
14  rating: row.rating,
15  duration: row.duration,
16  listed_in: row.listed_in,
17  description: row.description
18  });
19
20  // Create Director nodes and relationships
21  MATCH (t:Title)
22  WITH t, split(t.director, ', ') AS directors
23  UNWIND directors AS director
24  MERGE (d:Director {name: director})
```

```
25  MERGE (t)-[:DIRECTED_BY]→(d);
26
27  // Create Actor nodes and relationships
28  MATCH (t:Title)
29  WITH t, split(t.cast, ', ') AS actors
30  UNWIND actors AS actor
31  MERGE (a:Actor {name: actor})
32  MERGE (t)-[:ACTED_BY]→(a);
33
34  // Create Country nodes and relationships
35  MATCH (t:Title)
36  WITH t, split(t.country, ', ') AS countries
```

```
28  MATCH (t:Title)
29  WITH t, split(t.cast, ', ') AS actors
30  UNWIND actors AS actor
31  MERGE (a:Actor {name: actor})
32  MERGE (t)-[:ACTED_BY]→(a);
33
34  // Create Country nodes and relationships
35  MATCH (t:Title)
36  WITH t, split(t.country, ', ') AS countries
37  UNWIND countries AS country
38  MERGE (c:Country {name: country})
39  MERGE (t)-[:AVAILABLE_IN]→(c);
```

# Cleanse data if there are any null values:

```
1  //Cleanse the data
2
3  MATCH (n)
4  WHERE n.property IS NULL
5  REMOVE n.property;
```

(no changes, no records)

# Transform data

```
1  //Transform Data
2
3  // Split cast, listed_in, and country into arrays
4
5  MATCH (t:Title)
6  WITH t, split(t.cast, ', ') AS actors, split(t.listed_in, ', ') AS genres, split(t.country, ', ') AS countries
7  SET t.actors = actors, t.genres = genres, t.countries = countries;
8
9  // Convert release_year to integer
10
11 MATCH (t:Title)
12 SET t.release_year = toInteger(t.release_year);
13
14 // Convert titles to lowercase
15
16 MATCH (t:Title)
17 SET t.title = toLower(t.title);
```

```
neo4j$ MATCH (t:Title) WITH t, split(t.cast, ', ') AS actors, split(t.listed_in, ', ') AS genres, split(t.country, ', ') AS countries SET t.act...
neo4j$ MATCH (t:Title) SET t.release_year = toInteger(t.release_year)
neo4j$ MATCH (t:Title) SET t.title = toLower(t.title)
```

# Perform Aggregation Operations:

```
1  //Aggregation Operation 1
2
3  //Count the number of titles based on their type (movie or TV show).
4
5  MATCH (t:Title)
6  RETURN t.type AS type, count(*) AS num_titles
7
```

| type | num_titles |
|------|-----------|
| "Movie" | 50 |
| "TV Show" | 9 |

```
1  //Aggregation Operation 3
2
3  //Average number of actors per title
4
5  MATCH (t:Title)
6  RETURN avg(size(t.actors)) AS avg_actors_per_title
```

| avg_actors_per_title |
|---------------------|
| 6.8181818181818175 |

```
1  //Aggregation Operation 2
2
3  //Distribution of titles by release year
4
5  MATCH (t:Title)
6  RETURN t.release_year AS release_year, count(*) AS num_titles
7  ORDER BY t.release_year
```

| release_year | num_titles |
|-------------|-----------|
| 2004 | 1 |
| 2009 | 3 |
| 2010 | 1 |
| 2011 | 2 |
| 2012 | 2 |

# Query your database:

```
1  //Query the Database
2
3  //Find all titles directed by a specific director
4
5  MATCH (t:Title)-[:DIRECTED_BY]→(d:Director {name: 'Fernando Lebrija'})
6  RETURN t.title AS title;
```

| | title |
|---|---|
| Table | |
| | ¹ "#realityhigh" |
| A | |

```
1  //Query the database 2
2
3  // Count the number of titles available in a particular country
4
5  MATCH (t:Title)-[:AVAILABLE_IN]→(c:Country {name: 'United States'})
6  RETURN count(t) AS num_titles;
7
```

| | num_titles |
|---|---|
| Table | |
| | ¹ 22 |
| A | |

# Query your database:

```
1  //Query the database 3
2
3  //Find all the titles in a specific year
4
5  MATCH (t:Title)
6  WHERE t.release_year = 2009
7  RETURN t.title AS title, t.type AS type, t.release_year AS release_year;
```

| | title | type | release_year |
|---|---|---|---|
| 1 | "krish trish and baltiboy" | "Movie" | 2009 |
| 2 | "bangkok traffic (love) story" | "Movie" | 2009 |
| 3 | "phobia 2" | "Movie" | 2009 |

```
1  //Query the database 4
2
3  //Show all titles where a specific actor/actress appears
4
5  MATCH (t:Title)-[:ACTED_BY]→(a:Actor {name: 'Nicolas Cage'})
6  RETURN t.title AS title, t.type AS type, t.release_year AS release_year;
7
```
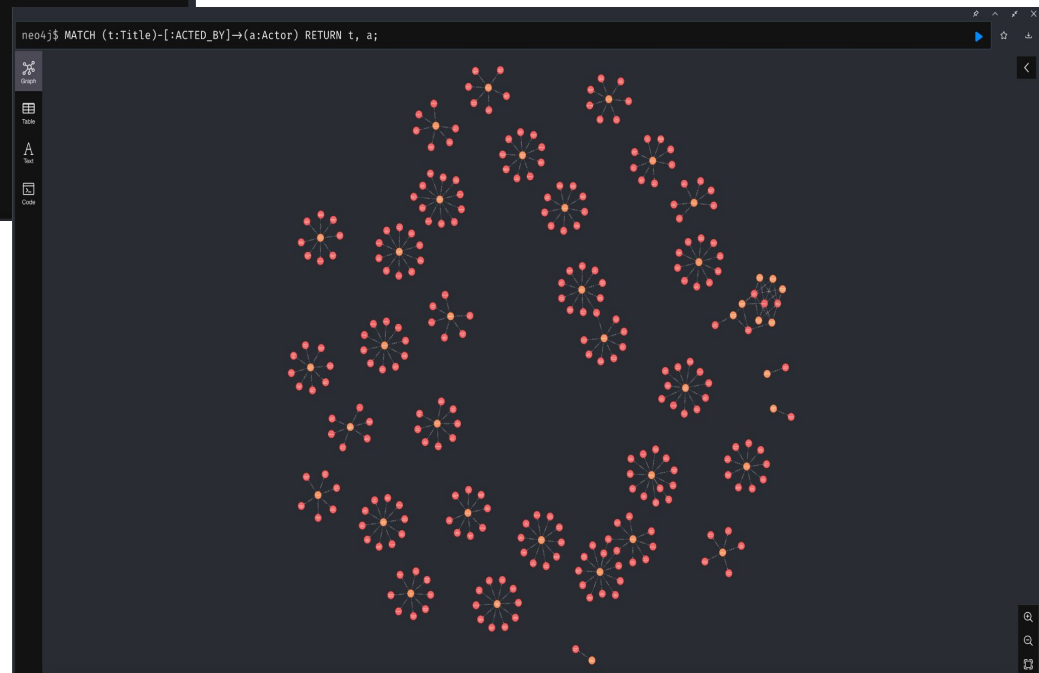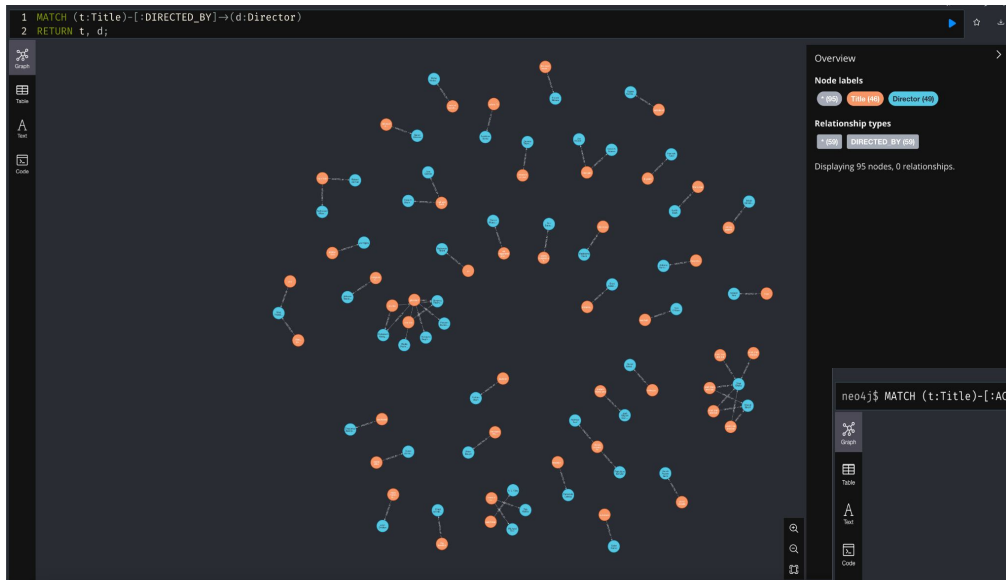
| | title | type | release_year |
|---|---|---|---|
| 1 | "the runner" | "Movie" | 2015 |

```
1  MATCH (t:Title)-[:DIRECTED_BY]→(d:Director)
2  RETURN t, d;
```
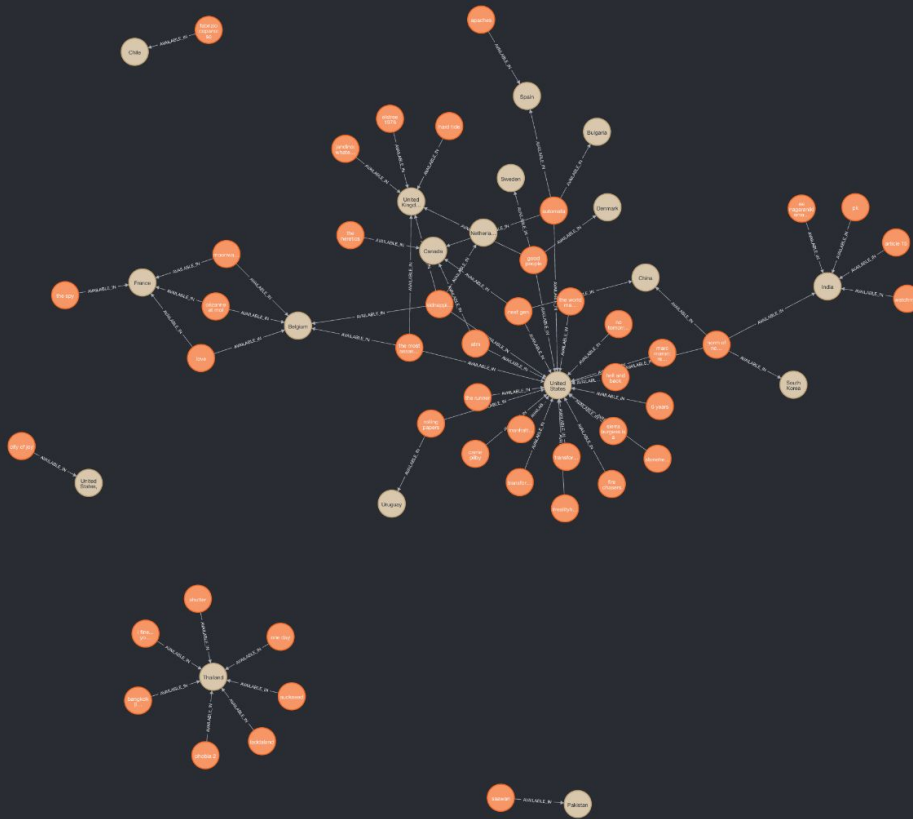
Overview

**Node labels**

* (95)    Title (46)    Director (49)

**Relationship types**

* (50)    DIRECTED_BY (50)

Displaying 95 nodes, 0 relationships.

```
neo4j$ MATCH (t:Title)-[:ACTED_BY]→(a:Actor) RETURN t, a;
```

# CONCLUSION

*In this project, we explored the Netflix titles dataset using Neo4j, a graph database. We leveraged the graph database model to represent the relationships between different entities such as titles, actors, directors, genres, and countries. Through a series of steps, we performed data modeling, data cleansing, transformation, aggregation operations, and querying to gain insights into the dataset.*

*The use of Neo4j proved to be advantageous for this project due to its native support for graph data modeling and querying. It allowed us to represent complex relationships between entities in a natural and intuitive way, making it easier to analyze and query the dataset.*

*Throughout the project, we utilized Cypher queries to interact with the graph database, performing various tasks such as retrieving nodes and relationships, cleansing data, transforming data, and performing aggregation operations. We visualized the data using the Neo4j Browser interface, which provided an interactive way to explore the graph and gain insights into the relationships between different entities.*