# MariaDB vs Neo4j
## Comparative Performance Analysis of Relational and NoSQL Database Systems

Cristian Santaroni 1800659

Francesco Fortunato 1848527

Data Management Course

Master in Engineering in Computer Science

A.y. 2022/2023

# OVERVIEW

## Topic

Pick up a relational database system, a noSQL database system and a dataset, to compare their performances w.r.t. the execution of some relevant queries.

## Goal

Find queries where relational database systems perform better and queries where noSQL database system perform better.

# TOOLS

- Relational Database System: MariaDB
- NoSQL Database System: Neo4j
- Python
- Dbeaver

# DATASET DESCRIPTION

**Dataset Details:**

- Clean, Structured, and Updated Football Data
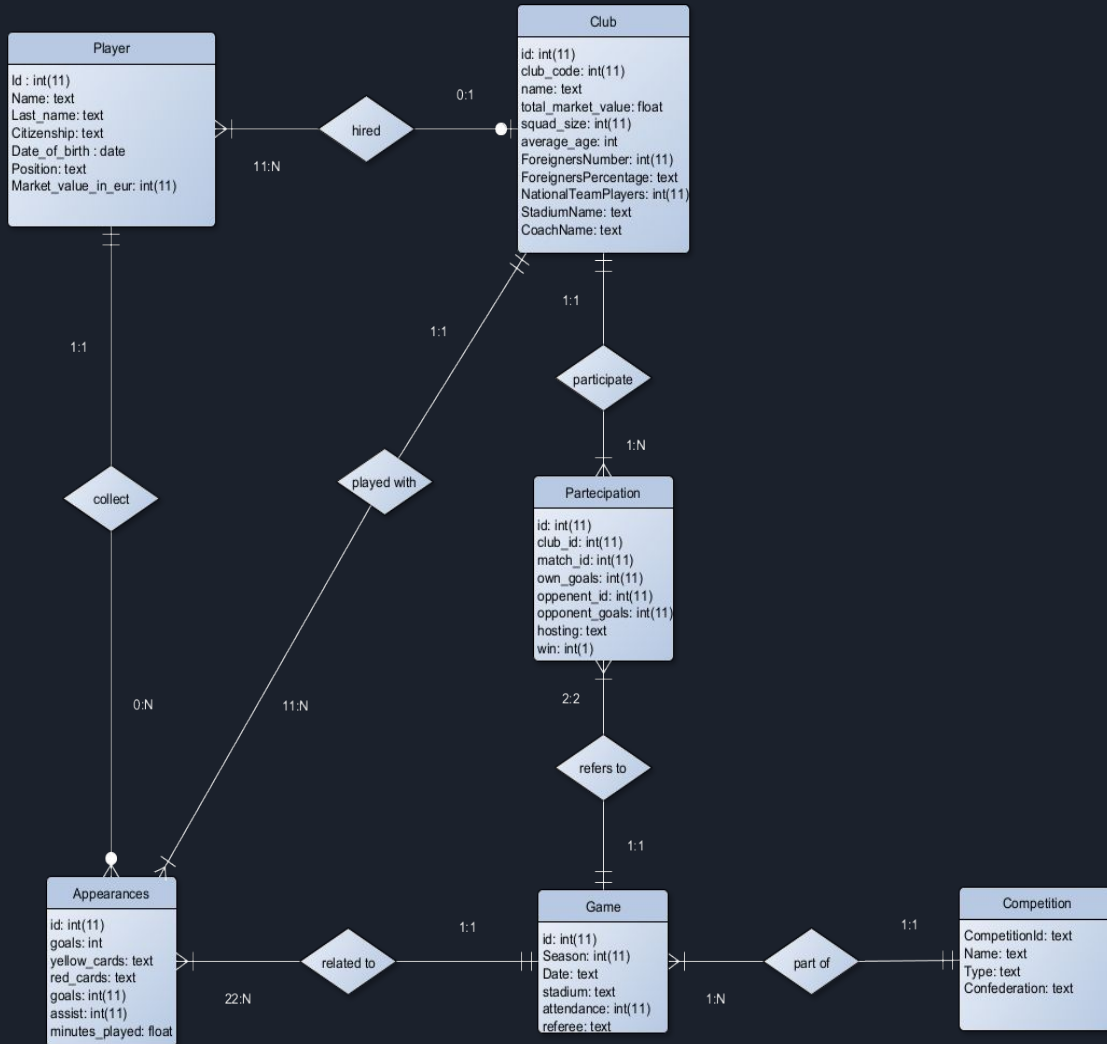- Sourced from Transfermarkt ([Kaggle link](#))

**Dataset Contents:**

- 63.281 Games Spanning Multiple Seasons
- 106.074 Participations Related to Games
- Representing All 43 Major Competitions
- Data on 426 Football Clubs
- Information on 29.455 Football Players
- Extensive Data on 1.179.060 Player Appearances
- Total numbers of records: 1.378.339
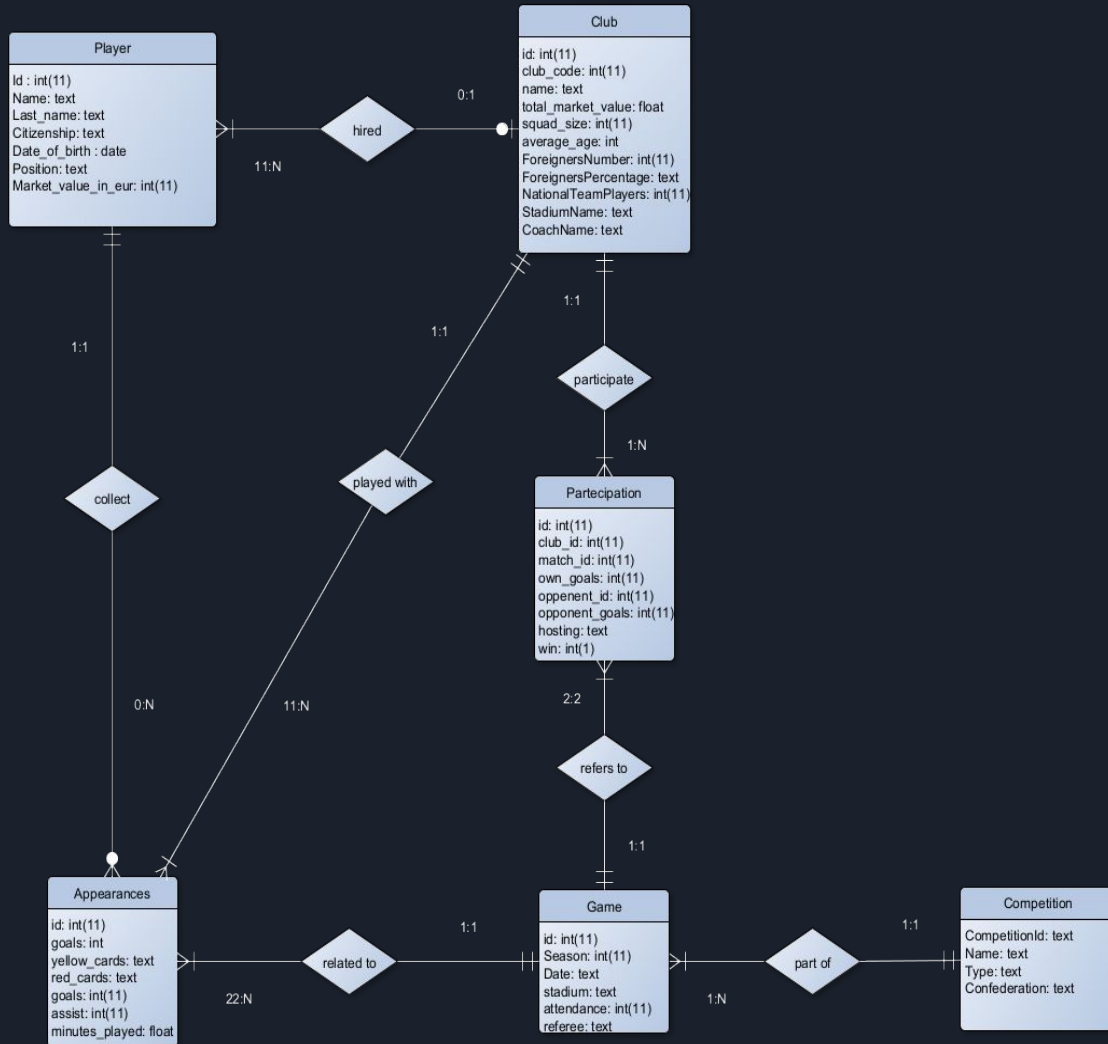
# ER SCHEMA

## Tables:

- **Player**: This table stores information about football players.
- **Appearances**: This table keeps records of player appearances in various games.
- **Game**: This table contains details about football games.
- **Club**: This table holds information about football clubs.
- **Participation**: This table links players and clubs to the games in which they participated.
- **Competition**: This table contains data about football competitions.



**Player**
Id : int(11)
Name: text
Last_name: text
Citizenship: text
Date_of_birth : date
Position: text
Market_value_in_eur: int(11)

**Club**
id: int(11)
club_code: int(11)
name: text
total_market_value: float
squad_size: int(11)
average_age: int
ForeignersNumber: int(11)
ForeignersPercentage: text
NationalTeamPlayers: int(11)
StadiumName: text
CoachName: text

**Partecipation**
id: int(11)
club_id: int(11)
match_id: int(11)
own_goals: int(11)
oppoent_id: int(11)
opponent_goals: int(11)
hosting: text
win: int(1)

**Appearances**
id: int(11)
goals: int
yellow_cards: text
red_cards: text
goals: int(11)
assist: int(11)
minutes_played: float

**Game**
id: int(11)
Season: int(11)
Date: text
stadium: text
attendance: int(11)
referee: text

**Competition**
CompetitionId: text
Name: text
Type: text
Confederation: text

hired — 0:1 — 11:N
collect — 1:1 — 0:N
played with — 1:1 — 11:N
participate — 1:1 — 1:N
refers to — 2:2 — 1:1
related to — 1:1 — 22:N
part of — 1:1 — 1:N

# ER SCHEMA

## Relations:

- **Collect**: player collects an appearance
- **Hired**: a player can be hired by a club
- **Refers to**: a game must have 2 clubs
- **Part of**: a game belongs to a competition
- **Played with**: an appearance is always played with a club
- **Participate**: a participation is always related to a club



**Player**
Id : int(11)
Name: text
Last_name: text
Citizenship: text
Date_of_birth : date
Position: text
Market_value_in_eur: int(11)

**Club**
id: int(11)
club_code: int(11)
name: text
total_market_value: float
squad_size: int(11)
average_age: int
ForeignersNumber: int(11)
ForeignersPercentage: text
NationalTeamPlayers: int(11)
StadiumName: text
CoachName: text

**Partecipation**
id: int(11)
club_id: int(11)
match_id: int(11)
own_goals: int(11)
oppoenent_id: int(11)
opponent_goals: int(11)
hosting: text
win: int(1)

**Appearances**
id: int(11)
goals: int
yellow_cards: text
red_cards: text
goals: int(11)
assist: int(11)
minutes_played: float

**Game**
id: int(11)
Season: int(11)
Date: text
stadium: text
attendance: int(11)
referee: text

**Competition**
CompetitionId: text
Name: text
Type: text
Confederation: text

hired — 0:1 / 11:N
collect — 1:1 / 0:N
played with — 1:1 / 11:N
participate — 1:1 / 1:N
refers to — 2:2 / 1:1
related to — 1:1 / 22:N
part of — 1:1 / 1:N

# DBeaver: EXCERPT

# DBeaver: FOREIGN KEYS

# Python scripts: importing data

In order to import all the data in our graph db, to create nodes and edges, we designed a python script. To improve the creation of the nodes and the respective edges we used a "**Uniqueness Constraint** " which ensure the uniqueness of the nodes based on the ID. This is used by Neo4j which automatically define an index to speed up the loading process and queries.

The python code can be divided in 3 parts:

- **Data Extraction** : opening all the csv files taken from kaggle and organizing them into dictionaries.
- **Node Creation**: using cypher queries, nodes are created for each data type.
- **Edges Creation**: using cypher, the relationships are created between the nodes

After this phase, data consistency was achieved with a total number of nodes of 1.378.339 and 3.842.064 relationships
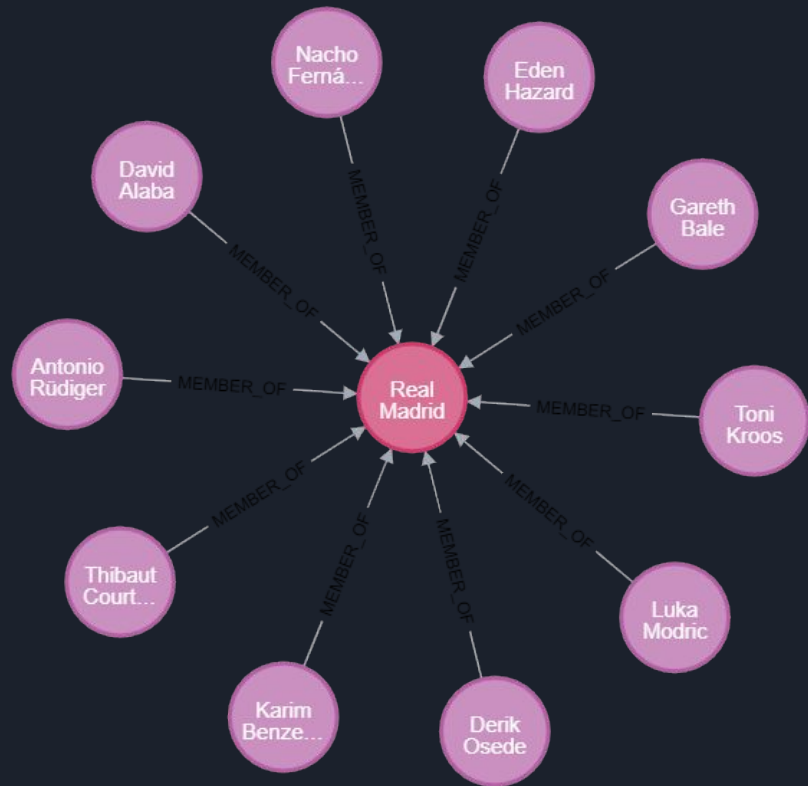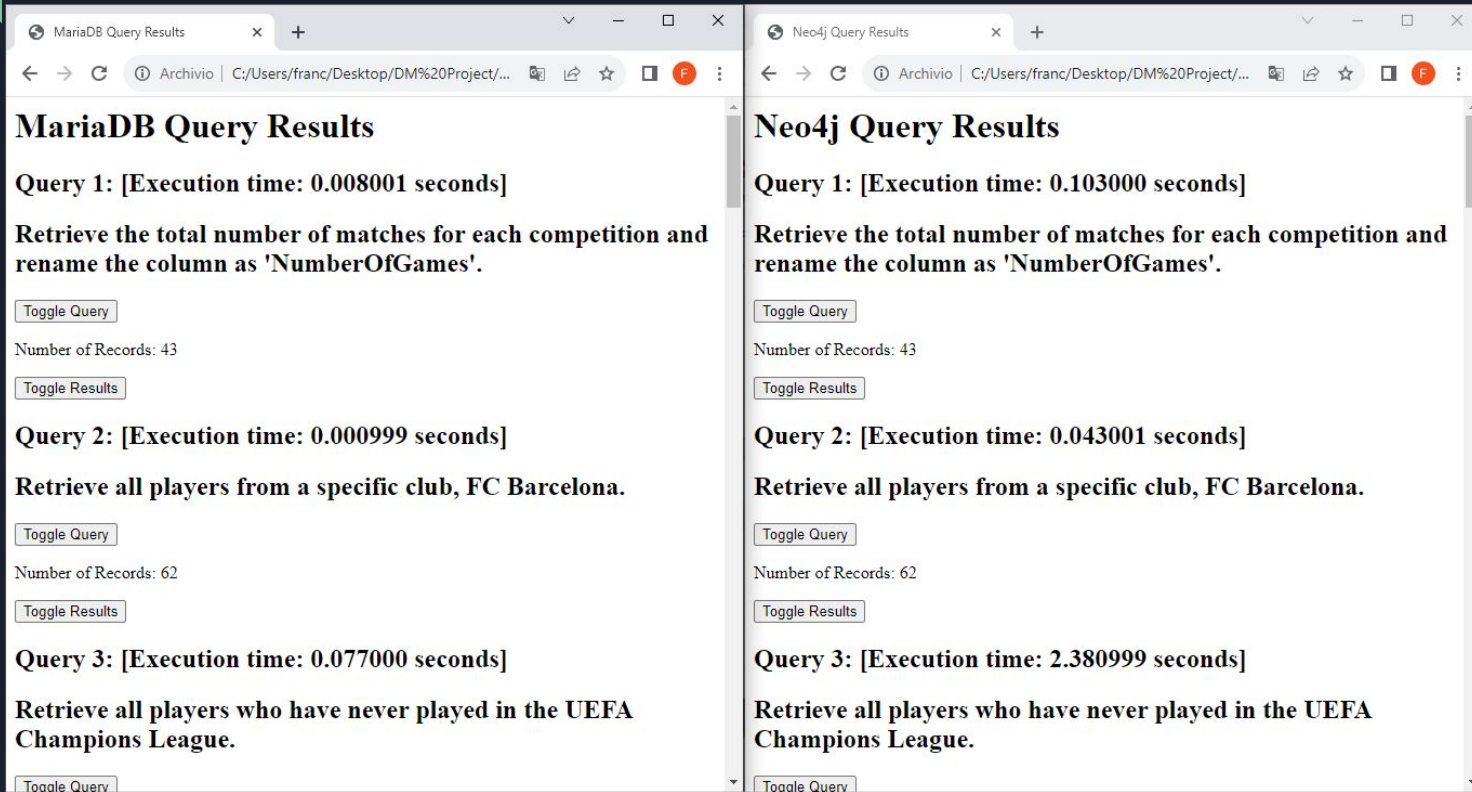
# Python scripts: average execution time and demo

- For the live demo, we designed a Python script which runs 20 queries and stores the results inside two different html files, one for MariaDB and one for Neo4j. The code can be reachable [here](#).

- Regarding the performance evaluation, to measure and compare the average execution time between Mariadb and Neo4j across various queries, we've developed a dedicated Python script. The script runs each query 1000 times, and then computes the average execution time of each query. In the subsequent slides, we will present the results of this performance evaluation, providing  the average execution times for each query comparing MariaDB and Neo4j.

# Python scripts: Live Demo

# Python scripts: Average execution time

## MariaDB

```python
counter = 0

for query in arr_string:
    total_time =0
    for i in range(1000):
        start = time.time()
        cur.execute(query)
        conn.commit()
        result = list(cur)
        total_time += time.time() - start
        num_records = len(result)   # Count the number of records
    counter += 1
    avg_time= total_time/1000
    print(f"Query {counter} completed. Average execution time: {avg_time:.6f} seconds")
conn.close()
```

## Neo4j

```python
# Function to execute query and measure execution time
def execute_query(query):
    with driver.session() as session:
        total_time = 0
        for i in range(1000):
            with session.begin_transaction() as tx:
                start_time = time.time()
                result = tx.run(query)
                # Fetch and store the results in a list
                results_list = list(result)
                tx.commit()
                total_time += time.time() - start_time
            num_records = len(results_list)   # Count the number of records

        return total_time, results_list, num_records

for i, query in enumerate(queries):
    execution_time, result, num_records = execute_query(query)
    avg_time = execution_time/1000
    print(f"Query {i+1} completed. average execution time: {execution_time:.6f} seconds")
```

# Average execution time (1/2)

Specs:
- AMD Ryzen 5800x3d
- Ram 16 GB
- GTX 1080

| # Query | MariaDB | Neo4j |
| --- | --- | --- |
| Query 1 | 8 ms | 14.8 ms |
| Query 2 | 0.3 ms | 18 ms |
| Query 3 | 80.9 ms | 1.912253 s |
| Query 4 | 184.1 ms | 637.721 ms |
| Query 5 | 1.604314 s | 1.101904 s |
| Query 6 | 13.5 ms | 299.801 ms |
| Query 7 | 1.365251 s | 362.8 ms |
| Query 8 | 5.790 ms | 9.4 ms |
| Query 9 | 49.865 ms | 18.3 ms |
| Query 10 | 2.340601 s | 116.9 ms |

# Average execution time (2/2)

Specs:
- AMD Ryzen 5800x3d
- Ram 16 GB
- GTX 1080

| # Query | MariaDB | Neo4j |
|---------|---------|-------|
| Query 11 | 1.754104 s | 2.798699 s |
| Query 12 | 23.5 ms | 287.551 ms |
| Query 13 | 998.044 ms | 313.567 ms |
| Query 14 | 2.5874 s | 1.466901 s |
| Query 15 | 69.4 ms | 67.4 ms |
| Query 16 | 3.0753 s | 2.754757 s |
| Query 17 | 172.201 ms | 96.407 ms |
| Query 18 | 61.505777 s | 8.48902 s |
| Query 19 | 48.671087 s | 12.008614 s |
| Query 20 | 9.5 ms | 28.7 ms |

# Conclusions: Neo4j wins

## Mariadb

**Pros**:

- Fast Data Loading
- Excelled in simpler, structured queries
- Effective filtering

**Cons**:

- Sometimes queries are more complex to write
- Performance tapered off with more complex query
- Query Optimization Challenges

## Neo4j

**Pros**:

- User-Friendly Queries
- Competitive Performance for Complex Graph-related Queries
- Most of the time is faster than RDBMS
- Consistent Performance

**Cons**:

- Filtering Challenges (Not great at filtering out useless nodes)
- Initial Loading Complexity.

Thank you