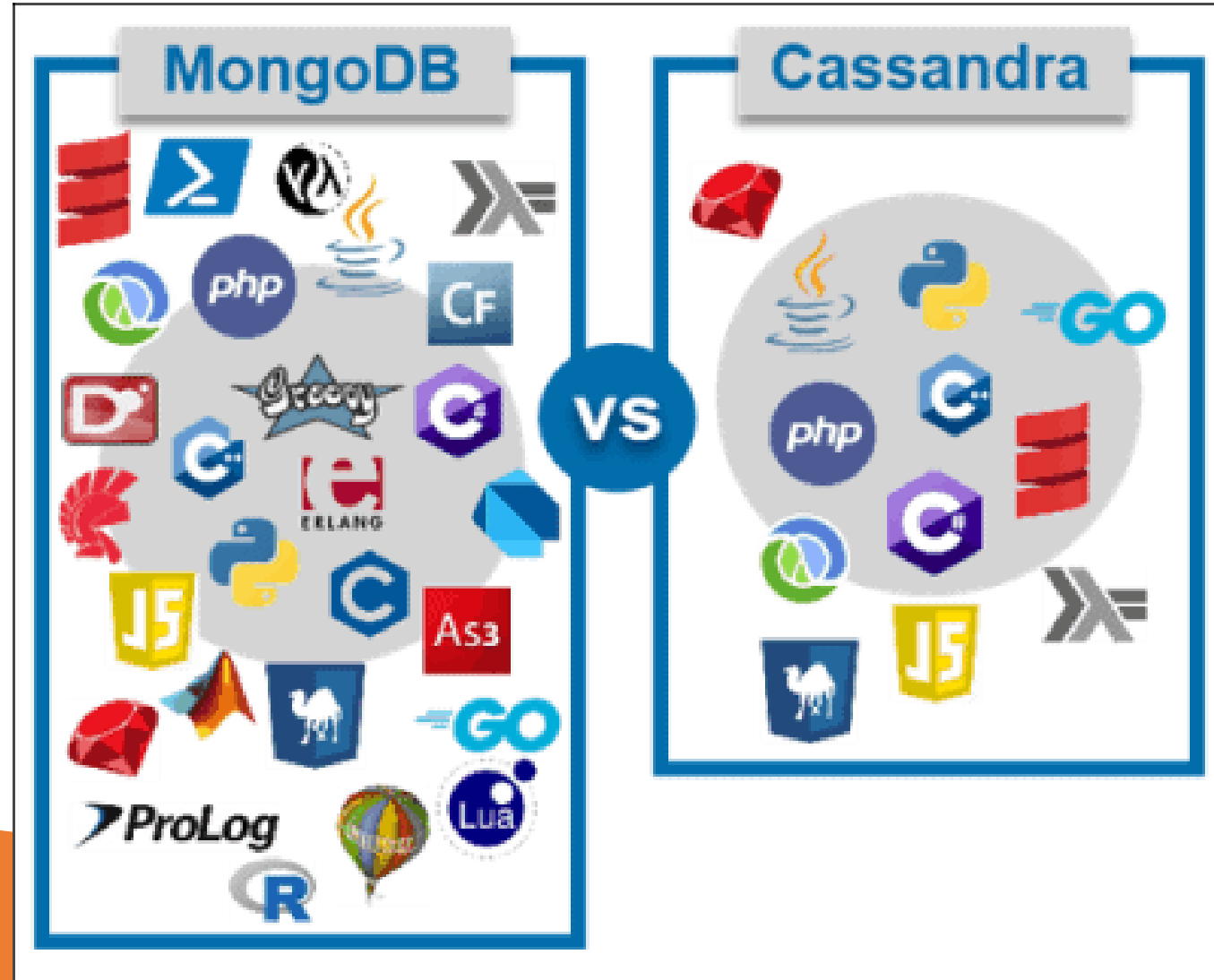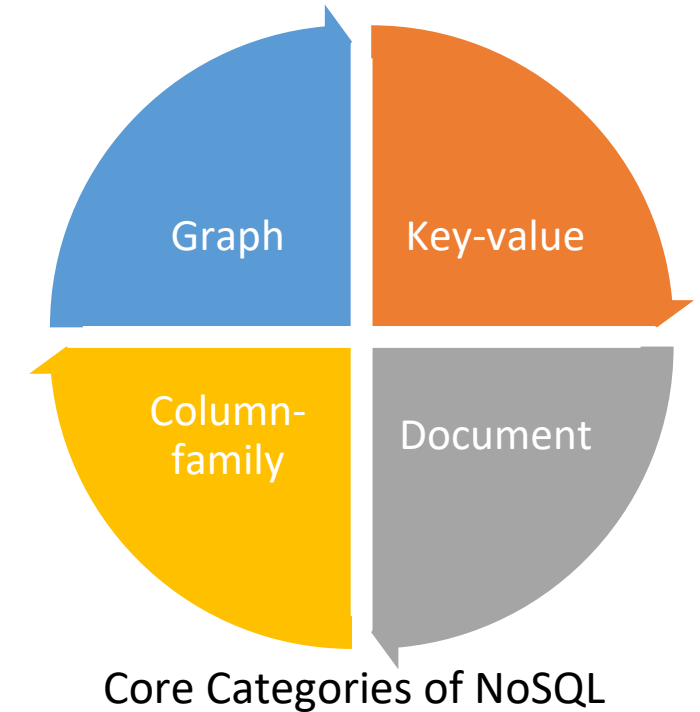# Comparative Performance Analysis between Cassandra and MongoDB in tracking efficiency of NBA Players
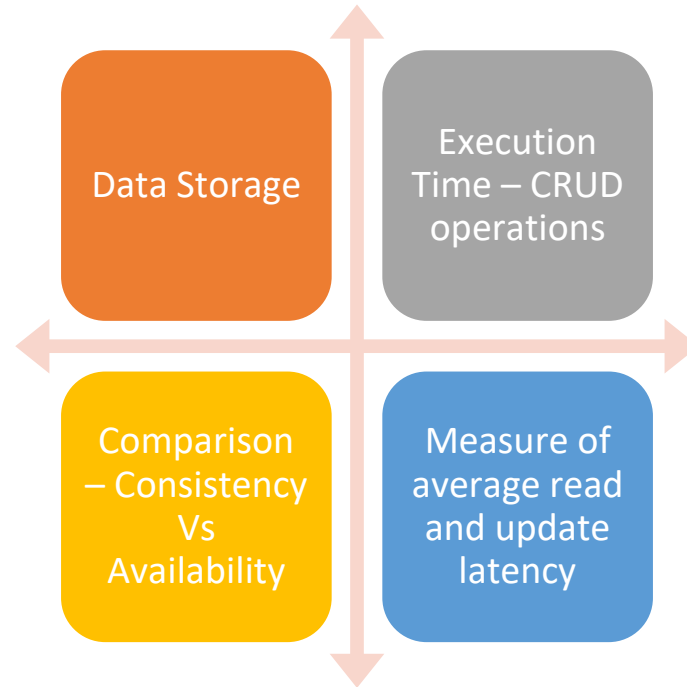
# Introduction: What is the problem?

- Finding a right database for a specific workload or an application can be very challenging due to the availability of a huge number of NoSQL databases.

- **Hence, evaluating the performance of NOSQL in-order to select a right database as per the requirement of each application is the need of the hour.**

- Each application communicates with the database using the CRUD operations

- The amount of data, the type of the data to be handled, the run-time for CRUD operations can affect the performance of the NOSQL.

- MongoDB and Cassandra is selected for evaluating their performance in terms of data storage, query handling and run-time for each operation.

Graph

Key-value

Column-family

Document

Core Categories of NoSQL

# Introduction:  Scope of Investigation

❏ The scope of this project is to do an in-depth analysis on the performance of Cassandra and MongoDB as these two databases are widely used in many industrial applications.

❏ The dataset consists of performance parameters of NBA players and other game measures.

❏ The comparison will be done based on the following parameters:

Data Storage

Execution Time – CRUD operations

Comparison – Consistency Vs Availability

Measure of average read and update latency

# Theoretical bases & literature review: Our Solution to solve this problem

- The research papers studied showed performance analysis between different databases for synthetic data that was generated using YCSB benchmark.

- Hence, a real-life scenario dataset was selected to measure the efficiency of NBA players to make a performance comparison between MongoDB and Cassandra.

- This will give a better and meaningful insights about the effectiveness of both databases.

# Theoretical bases & literature review: Where our solution different from others?

## Data Storage

- Analyze which data storage model makes querying easier

## CAP

- Understand which database is better option when focusing on consistency over availability or vice versa.

## Real-life Scenario dataset

- Simulating a real-scenario dataset

# Theoretical bases & literature review: Why our solution is better?

Comparing the time taken to import data in both the databases

Analysis on CAP theorem

Comparing the execution time of READ, UPDATE, DELETE operation on both the databases

Exploring the average read and update latency of both the databases along with throughput rate.

# Hypothesis/Goals

To test and verify if MongoDB performs better in read operations and Cassandra performs better in write operations.

- From the research papers studied, it was observed that MongoDB performs better for read and delete operations whereas Cassandra performs better for write operations. We would like to verify if the same results are obtained for our analysis.

To test and verify if Cassandra has low update and read latency as compared to MongoDB, when the load increases.

- As per the research paper, Cassandra has low average read/update latency and high throughput rate when the load increases. This is due to the fact that they have peer-peer replication. On the other hand, MongoDB latency increases exponentially when the load is increased.

# Methodology: How to generate/collect input data

## Downloading Data

| Dataset downloaded from Kaggle website | Upload the data in Pandas for data wrangling, if required | Check the shape and size of the dataset for better understanding |
|---|---|---|

## Data Loading in MongoDB

Installing the MongoDB Module and MongoDB compass → Importing required modules in notebook and connecting to the mongo client i.e. local host and port → Converting the uploaded CSV files to Json format with the help of Python-Pandas

Imported data can be viewed using MongoDB Compass and execution time in data loading shall be noted ← Inserting data in the database and collection in MongoDB compass by using command **insert_many()** ← Creating database and collection in the Mongo cluster

# Data Loading in MongoDB

# Methodology: How to generate/collect input data

**Data Loading in Cassandra**



```
Installation of Cassandra 4.0.3 on docker  →  Creating Keyspace and column families  →  Importing the csv file to the column family using the COPY command
                                                                                                      ↓
Time taken to load data can be seen in the command prompt  ←  Data can be seen from the keyspace and specific column family by using select *  command
```

KEYSPACE

TABLE 1
TABLE 2
TABLE 3
TABLE 4
TABLE 5
..............

NODE

10.0.0.7

# Data Loading in Cassandra

# Methodology: Algorithm Design

Data Loading → Setting the replica factor → Recording the loading time

Performing the READ,UPDATE,DELETE operation → Recording the execution time for queries in atomic operation → Measuring the importance of Consistency Vs Availability

Measuring the average read/update latency by varying the number of records → Performance Comparison between MongoDB and Cassandra based on these results → Tracking the efficiency of NBA players and comparing the effectiveness of both the databases

# Methodology: Language Used

| Language | Purpose |
|---|---|
| **CQL** | For querying in Cassandra |
| **Python - Pandas** | For ETL and data wrangling |
| **PyMongo** | For querying and inserting data in MongoDB |

PyMongo is a Python distribution containing tools for working with MongoDB, and is the recommended way to work with MongoDB from Python

# Methodology: Tools Used

| Tools | Purpose |
|---|---|
| **Jupyter notebook** | For converting csv file to json format, data visualization and data cleaning |
| **MongoDB Compass** | For managing collection and document. Also, for querying, aggregating, and analyzing MongoDB data in a visual environment |
| **MongoDB Atlas** | To access the data virtually in the cloud, we will establish a connection to the cloud using MongoDB Atlas |
| **Docker** | Docker Image is one of the methods of installing Cassandra |
| **Cassandra** | For data storing and data analysis |

# Output Generation



**SQLite DB browser**

**Cassandra**

**MongoDB**

# Output Analysis & Discussion

## Comparison of data loading time

- The execution time required for loading data for both the databases

| Data Loading | Number of Records | Cassandra (ms) | MongoDB (ms) |
|---|---|---|---|
| Team_analysis | 29 | 398 | 1 |
| Player_analysis | 394 | 659 | 0 |
| Game_analysis | 2720 | 4275 | 4 |
| Basketball_analysis | 34567 | 103549 | 26 |

### Time for Data Loading

# Output Analysis & Discussion

## Comparison of execution time for various CRUD operations

- The time taken for various CRUD operations for both the databases

| CRUD Operations | Cassandra (ms) | MongoDB (ms) |
|---|---|---|
| Insert | 41.616 | 42 |
| Read | 15.403 | 7 |
| Update | 54.020 | 63 |
| Delete | 33.122 | 5 |

**Execution time for CRUD operations**

# Output Analysis & Discussion

## Comparison of execution time for different consistency levels

| CRUD | Consistency levels | Cassandra (ms) | MongoDB (ms) |
|------|--------------------|----------------|--------------|
| Read | One | 225.147 | 1 |
| | Quorum | 228.763 | 5 |
| Update | One | 108.951 | 70 |
| | Quorum | 162.396 | 556 |
| Insert | One | 85.303 | 2 |
| | Quorum | 111.234 | 36 |

**Consistency levels for Cassandra**

**Consistency levels for MongoDB**

# Output Analysis

## Cassandra: Comparison of Read and Write latency

| Read Operation | | Write Operation | |
|---|---|---|---|
| Count | Time Taken (ms) | Count | Time Taken (ms) |
| 9 | 11.994 | 2 | .0118 |
| 5 | 10.532 | 3 | 0.239 |
| 14 | 10.311 | 10 | 0.8415 |

### Read Latency for Cassandra

### Write Latency for Cassandra

# Output Analysis

## MongoDB: Comparison of Read and Write latency

| Read Operation | | Write Operation | |
|---|---|---|---|
| Count | Time Taken (ms) | Count | Time Taken (ms) |
| 2 | 0.64 | 2 | 2.048 |
| 4 | 1.28 | 4 | 2.56 |
| 6 | 1.28 | 6 | 49.152 |
| 8 | 1.024 | 10 | 98.304 |
| 10 | 5.12 | 12 | 166.644 |
| 20 | 5.434 | 15 | 148.936 |



Read and Write latency for MongoDB

# Compare output against hypothesis

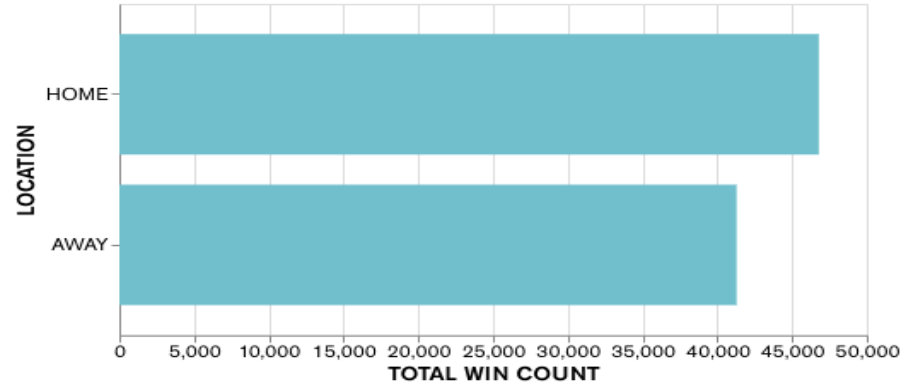**Hypothesis 1 :** "To test and verify if MongoDB performs better in CRUD operations and Cassandra performs better in write operations"

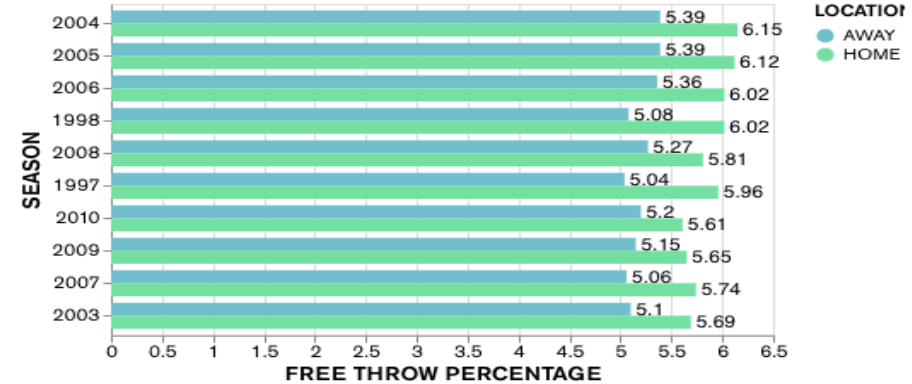# Performance evaluation for NBA Players

# Summary & Conclusions

❖ **CRUD Operation**

- MongoDB loads the data faster in comparison to Cassandra
- Except for UPDATE operation, MongoDB is more efficient than Cassandra for READ and DELETE operations.

❖ **Latency**

- Cassandra performed much better in terms of latency (lower latency) as the workload increased
- Write latency for MongoDB increased significantly with increase in workload

❖ **Consistency**

- Cassandra and MongoDB: At consistency level one, the time required for each CRUD operation to be executed is low as compared to the time required for Quorum consistency

❖ **Data Visualization**

- Evaluated the performance of the NBA players using the MONGO charts in Atlas
- Data analysis using mongo charts is more efficient as data is updated at regular intervals

# Recommendations

- More number of read and write operations can be performed to better understand the latency patterns as well as the CRUD operations of Cassandra and MongoDB

- Analysis can be performed on larger datasets which shall help to get better understanding about performance of both the databases

- Comparison between Cassandra and MongoDB can be done on real time data such as for on-field game analytics to get better insights of their performance metrics in real time

- On-field game analytics can make use of Mongo charts for data analysis as the data is updated at regular interval

- Horizontal scalability can also be measured by increasing the number of nodes

Thank you