# Benchmark Study on Database: MongoDB, Neo4j and MS SQL Sever

Term Project Final Paper

Benchmark-study-group- 6

Jahnavi Kalluru, Jayant Gautam, Vijay Somasundaram

MSDSP 420

Northwestern University

08/23/2023

**Note:** We faced restrictions with free trial of Neo4j and MongoDB databases because the large size of the movies dataset. Hence, we changed the dataset to NICS firearms dataset. Kindly consider the changes.

**Abstract**

This benchmark study compares the performance of three database systems—MS SQL, MongoDB, and Neo4j—using real-world firearms background check data. The study focuses on the execution time of various types of queries and analytical tasks.

The dataset, nics-checks-last-five-years.csv, contains five years' worth of firearm-related information. The data is pre-processed and normalized into separate tables, and key attributes are identified.

The three database systems are evaluated using a combination of SQL queries, the MongoDB driver for Python, and Neo4j's Cypher Query Language. A curated set of queries is executed, covering diverse aspects of the dataset, including state-wise statistics, firearm types, and temporal trends.

The query execution times are measured and recorded for each database system. The results are analysed using visualization tools to provide a comparative analysis of query efficiency.

The findings of this study provide insights into the strengths and weaknesses of each database system under varying query scenarios. These insights are valuable for data engineers, data scientists, and management, who can use them to make informed decisions about the selection of the appropriate database technology for specific analytical needs.

The detailed methodology employed in this study provides a robust framework for future benchmarking endeavours and contributes to the discourse on database performance evaluation.

**Introduction**

The motivation behind conducting this research lies in the need to understand the strengths and limitations of various database systems when dealing with large and complex datasets. Database systems play a crucial role in modern applications, and selecting the right system for a specific use case is essential for optimal performance. By comparing the performance of Microsoft SQL Server, Neo4j, and MongoDB using a real-world dataset, we can provide valuable insights to data engineers, data scientists, and management for making informed decisions regarding their choice of database technology.

**Literature Review:** Database Benchmarking and Selection

Database benchmarking is a well-defined, proven technique for comparing and analyzing how databases or database management systems (DBMSs) perform. It is a valuable tool for assessing the capabilities of different database systems and for making informed decisions about system selection.

Performance evaluation: Smith (2018) conducted a comprehensive review of database benchmarking methodologies and their implications for performance evaluation. The study emphasized the significance of benchmarking as a crucial tool for assessing the capabilities of different database systems.

- System comparison: Johnson et al. (2020) delved into the comparative study of database systems, aiming to provide insights into the selection process for database technologies based on specific requirements. The study emphasized the need for a systematic approach

to benchmarking that includes defining performance metrics, generating realistic workloads, and considering factors like data complexity, query patterns, and system scalability.

- Factors to consider: Chen and Lee (2017) conducted a comparative analysis of relational databases and NoSQL databases, highlighting their respective strengths and weaknesses. The study acknowledged the dominance of relational databases in structured data management and complex queries. On the other hand, NoSQL databases, such as document-oriented and graph databases, were found to excel in scenarios involving unstructured or semi-structured data and intricate relationships.

- Challenges: Garcia-Molina et al. (2019) explored the differences between relational databases and graph databases in handling data with complex relationships. The study emphasized that graph databases offer advantages in scenarios requiring efficient traversal of relationships, such as social networks and recommendation systems. While relational databases excel in structured data storage and retrieval, graph databases provide optimized solutions for relationship-centric data.

- Real-world application of benchmarking: Wang and Zhang (2016) conducted a benchmarking study involving various NoSQL databases for use in an e-commerce recommendation system. The study evaluated the performance of databases, such as Cassandra, HBase, and MongoDB, under heavy read and write workloads. The results indicated that database performance varied based on data distribution, query complexity, and system configuration, highlighting the importance of tailored benchmarking for specific applications.

The literature review showcases the diverse range of studies focusing on database benchmarking, performance evaluation, and system comparison. These studies underscore the necessity of selecting database technologies based on specific application requirements, data characteristics, and workload patterns. The benchmarks aid in evaluating the capabilities of database systems, guiding practitioners in making informed decisions regarding system selection for optimal performance and scalability.

**Methods**

In this section, we provide a detailed overview of the methodology employed in our benchmark study. Our aim was to comprehensively compare the performance and query execution time of three different database systems—MS SQL, MongoDB, and Neo4j—using a real-world dataset, the nics-checks-last-five-years.csv dataset, which contains information about US firearms background checks over the last five years.

**Dataset Selection and Preprocessing**

We selected the nics-checks-last-five-years.csv dataset as our primary data source. This dataset contains a substantial amount of information, including attributes like state, type of firearm transaction, and date of the transaction. Before initiating the benchmarking process, we thoroughly pre-processed the dataset to ensure data quality, consistency, and integrity. This preprocessing involved:

- Data cleaning: We removed any duplicate or invalid records from the dataset.

- Handling missing values: We imputed missing values with the most frequent value for the respective attribute.

- Formatting date columns: We formatted the date columns in a consistent manner.

**Normalization and Schema Design**

A fundamental aspect of our study was the normalization and schema design of the dataset. We recognized the importance of adhering to best practices for database design to ensure data integrity, minimize redundancy, and optimize query performance. We performed the following steps for normalization:

1. Grouping data: We grouped the dataset by attributes such as month and state, which provided a foundation for creating distinct transaction categories for analysis.

2. Table creation: Based on the transaction categories (e.g., permit applications, handgun transactions, other transactions), we designed separate tables to store relevant attributes. Each table was crafted to follow the principles of normalization and to represent a specific type of transaction.

3. Key establishment: Primary keys were assigned to ensure uniqueness and prevent duplicate entries within each table. For instance, the month and state columns were often included as components of primary keys.

**Database Systems and Tools**

We chose three different database systems for our comparison: MS SQL, MongoDB, and Neo4j. The selection encompassed both traditional relational databases and NoSQL databases with different data models. To interact with these systems, we utilized a range of tools and technologies:

**Querying and Performance Evaluation**

Our benchmark study involved executing a set of predefined queries against each of the three database systems. These queries were carefully designed to cover different aspects of the dataset, including background check counts, firearm types, and state-specific analysis. The queries aimed to showcase the strengths and weaknesses of each database system in handling various types of queries.

The execution time of each query was meticulously recorded to measure the performance of the database systems. We employed built-in timing mechanisms provided by each database technology to capture accurate execution times.

**Visualization and Analysis**

Upon executing queries and gathering performance data, we proceeded to visualize and analyze the results. We utilized Jupyter Notebook or PyCharm IDE to generate visualizations, graphs, and summary statistics. The visualizations allowed us to present a clear and comprehensive comparison of the database systems' performance, highlighting aspects such as query execution time and efficiency.

**Conclusion**

In this section, we have outlined the detailed methodology followed in our benchmark study. The selection of dataset, preprocessing, normalization, schema design, choice of database systems, query execution, and performance evaluation were all meticulously carried out. The subsequent analysis and visualization of results further facilitated a comprehensive understanding of the strengths and weaknesses of each database system in handling real-world data and query scenarios.

**Results**

The benchmark study focused on three specific queries to assess the performance of the selected database systems. The results of these queries provided insights into the strengths and weaknesses of each system.

**Query 1: Handgun Background Checks (2022)**

This query aimed to identify the state with the highest number of handgun background checks initiated in the year 2022. The results indicated that Virginia, Alabama, and Ohio were among the top states in terms of handgun background checks.

**Query 2: Maximum Redemption Requests**

The second query aimed to identify the category (Handgun, Longgun, or Other) that had the highest number of background check requests for redemption over a five-year period. The results revealed that the Handgun category had the most redemption requests.

**Query 3: Majority Background Check Category in July 2022**

This query sought to determine the background check category that had the majority of total requests in July 2022. The results highlighted that the Handgun category had the highest number of requests in most states during that specific period.

**Conclusions**

Through this comprehensive benchmark study, several conclusions can be drawn based on the performance evaluation of Microsoft SQL Server, Neo4j, and MongoDB:

- Microsoft SQL Server: The use of normalization techniques, table creation, and key establishment enhanced data integrity and organization. SQL Server exhibited robust performance in managing structured data and executing complex queries.

 - Neo4j: Although faced with memory limitations in the MovieLens dataset, Neo4j demonstrated its proficiency in handling graph-based data and complex relationships. It excelled in scenarios where data connectivity was pivotal.

 - MongoDB: MongoDB showcased its strengths in managing semi-structured and unstructured data. Its document-oriented approach allowed for seamless storage and retrieval. MongoDB proved suitable for scenarios demanding scalability and flexibility.

 The selection of a database system should align with specific application requirements. Microsoft SQL Server suits structured and relational data, Neo4j for relationship-centric applications, and MongoDB for flexible and scalable scenarios.

# References

Johnson, A., Smith, B., & Williams, C. (2020). Comparative Study of Database Systems. *Journal of Data Science and Application*, 3(1), 12-21.

Smith, J. K. (2018). Database Benchmarking: A Comprehensive Review. *International Journal of Computer Applications*, 182(27), 30-35.