# Exploring Database Indexes: Types, Implementations, and Performance Analysis

## 1. Introduction

### A: Objective:

**Detail the significance of using database indexes to reduce the time complexity of search operations in large databases. Discuss how they are crucial for achieving efficient data retrieval.**

Database indexes are essential for efficient data retrieval and for lowering the time complexity of search operations in large databases. This is the reason they matter:

1. Enhanced Query Performance: In large databases, full-table scans—where the entire dataset must be scanned in order to find relevant records—are frequently the result of querying without indexes. This procedure can take a long time, particularly when the size of the database grows. Database indexes organize data in a structured format, making it possible to find specific records quickly. This makes index-based lookups possible for the database engine, which greatly enhances query performance and cuts down on the amount of time needed to retrieve data.

2. Faster Data Access: By acting as a road map, database indexes help the database engine efficiently traverse through the data. Indexes help the database find pertinent data faster by indexing columns that are often used in search conditions or join operations. Faster data access and retrieval can be achieved by the database engine using index structures to locate desired records rather than scanning the entire dataset.

3. Decreased Disk I/O: Especially in large databases, disk input/output (I/O) operations represent a significant performance bottleneck for databases. In the absence of indexes, queries frequently need to read data straight from disk, which can cause a lot of I/O operations and higher latency. Database indexes help to mitigate this by minimizing the amount of disk I/O required and enhancing overall system performance by enabling the database engine to access data primarily from memory.

4. Optimized Search Efficiency: By arranging data to make it easier to retrieve quickly, indexes are made to maximize search efficiency. Databases can reduce the amount of time it takes to find specific records by indexing columns that are frequently used in search predicates. Indexes can significantly reduce the time complexity of search

operations, which is especially useful for queries with complex conditions or sorting requirements.

5. Scalability: As databases grow in size and complexity, the need for efficient data retrieval becomes even more critical. Database indexes support scalability by providing a mechanism for handling large datasets without sacrificing performance. By optimizing index structures and periodically rebuilding indexes, databases can maintain efficient data retrieval even as the volume of data increases, ensuring scalability without compromising on speed or efficiency.

## B: Background:

## Provide a general introduction to the mechanics of database indexes, including how they are structured and why they are used in both transactional and analytical database operations.

Database management systems (DBMS) use database indexes as essential components to improve the effectiveness of data retrieval processes. They function as specialized data structures that store and arrange information so that it can be quickly found and retrieved. Indexes reduce the time complexity of search operations by allowing the DBMS to quickly locate and access pertinent data by strategically indexing particular columns or attributes within database tables.

### Mechanics of Database Indexes:

Database management systems (DBMS) use database indexes as essential components to improve the effectiveness of data retrieval processes. They function as specialized data structures that store and arrange information so that it can be quickly found and retrieved. Indexes reduce the time complexity of search operations by allowing the DBMS to quickly locate and access pertinent data by strategically indexing particular columns or attributes within database tables.

Indexes can be implemented using various data structures, each optimized for different types of queries and access patterns. Commonly used index structures include:

1. B-tree (Balanced Tree) indexes are extensively utilized in database management systems due to its effectiveness in facilitating equality and range queries. With each node holding a variety of keys and pointers to child nodes, these indexes arrange data in a hierarchical structure made up of nodes. B-tree indexes work well in transactional database operations, when latency reduction and quick data access are critical.

2. Hash Indexes: These indexes map keys to the appropriate locations in memory by using hash functions. Because hash indexes produce constant-time lookup operations, they

are perfect for processing point queries with a high degree of selectivity. Hash indexes are useful in some situations, but they might not work well for range searches or partial key matches.

**Why Database Indexes are Used:**

1. Database indexes offer several benefits that make them indispensable in both transactional and analytical database operations:

2. Better Query Performance: Indexes improve query performance and decrease response times by enabling quick data retrieval. This increases overall system efficiency and user happiness.

3. Faster Data Access: By allowing the DBMS to find and get pertinent data more quickly, indexes reduce the amount of time needed to complete queries and transactions.

4. Improved Search Efficiency: By arranging data in a systematic manner, indexes improve search efficiency by speeding up search and retrieval processes—even for intricate queries.

5. Support for Transactional and Analytical Workloads: Indexes are adaptable instruments that facilitate a variety of database functions, such as reporting, data analysis, and transaction processing. They can be customized to fit certain use cases and access patterns in order to maximize efficiency for a range of workloads.

6. Scalability and Adaptability: Database indexes facilitate scalability by effectively managing extensive datasets and dynamic access patterns. They can be optimized and dynamically changed to meet shifting workload demands and business objectives.

# 2. Methodology

- ➔ Environment Setup:
    - ◆ Software: mysql  Ver 8.0.33 for MacOs13.3 on arm64
    - ◆ Hardware: Apple M1 Pro 16 GB RAM
- ➔ Data Modeling:

◆ Schema Design: Illustrate with diagrams the database tables, including primary and foreign keys. Please refer picture-schema.png file
◆ Data Generation: Random Dataset using faker

# 3. Practical Implementation

➔ **Index Creation:**
◆ Provide SQL commands used to create each type of index.
- CREATE INDEX idx_member_name ON Members (member_name);
- CREATE INDEX idx_completed_date ON TaskAssignments (completed_date);
- CREATE INDEX idx_task_priority ON Tasks (task_priority);
◆ Explain the rationale behind the choice of index for each table or query type.
- All the primary keys and foreign keys are by default index by MySql. In addition to them, we have indexed member_name, completed_date and task_priority because these data are read in higher frequency in the current project.

➔ **Query Execution:**

List the queries that will be run, categorized by their expected benefit from indexing (e.g., select, update, delete).

```
-- Query 1 To get the member with the highest number of assignments
SELECT SQL_NO_CACHE m.member_id, m.member_name, m.member_email,
COUNT(ta.assignment_id) AS num_assignments
   FROM Members m
   JOIN TaskAssignments ta ON m.member_id = ta.member_id
   GROUP BY m.member_id
   ORDER BY num_assignments DESC
   LIMIT 1;
```

```
-- Query 2 Find members who have completed tasks with a priority higher than the
average priority and were completed recently.
SELECT m.member_id, m.member_name, AVG(t.task_priority) AS avg_priority,
COUNT(ta.task_id) AS num_completed_tasks
FROM Members m
JOIN TaskAssignments ta ON m.member_id = ta.member_id
JOIN Tasks t ON ta.task_id = t.task_id
WHERE ta.completed_date >= DATE_SUB(CURRENT_DATE(), INTERVAL 30 DAY)
GROUP BY m.member_id, m.member_name
HAVING AVG(t.task_priority) > (SELECT AVG(task_priority) FROM Tasks WHERE status = 1)
  AND COUNT(ta.task_id) > 5;
```

```
-- Query 3 Specify the member name you want to retrieve tasks for
SELECT SQL_NO_CACHE t.task_name,t.task_due_date,t.task_priority
    FROM Tasks t
    JOIN TaskAssignments ta ON t.task_id = ta.task_id
    JOIN Members m ON ta.member_id = m.member_id
    WHERE m.member_name like "%son%";
```

```
-- Query 4: updating the task_due_date of tasks assigned to members who have not
completed any tasks yet to be
-- one month later than their current task_due_date
UPDATE Tasks
SET task_due_date = DATE_ADD(task_due_date, INTERVAL 1 WEEK)
WHERE task_id IN (
    SELECT ta.task_id
    FROM TaskAssignments ta
    INNER JOIN Members m ON ta.member_id = m.member_id
    WHERE ta.member_id NOT IN (
        SELECT member_id
        FROM TaskAssignments
        WHERE status = 2
    )
);
```

# 4. Performance Analysis

➔ **Metrics:**
  ◆ **Response Time: Time taken for queries to execute.**

```
Query 1 Execution time noticed after indexing
    +-----------------+
    | execution_time  |
    +-----------------+
    | 00:00:00.006170 |
    +-----------------+


Query 2 Execution time noticed after indexing
    +-----------------+
    | execution_time  |
```

```
+----------------+
| 00:00:00.012600 |
+----------------+


Query 3 Execution time noticed after indexing

+----------------+
| execution_time |
+----------------+
| 00:00:00.014300 |
+----------------+


Query 4 updating the task_due_date of tasks assigned to members who have
not completed any tasks yet to be
-- one month later than their current task_due_date

+----------------+
| execution_time |
+----------------+
| 00:00:00.019800 |
+----------------+
Overall of all queries in a batch

+----------------+
| execution_time |
+----------------+
| 00:00:00.123500 |
+----------------+
```

◆ **Throughput: Number of queries handled per unit of time.**
**Total Number of Queries Executed: 4**
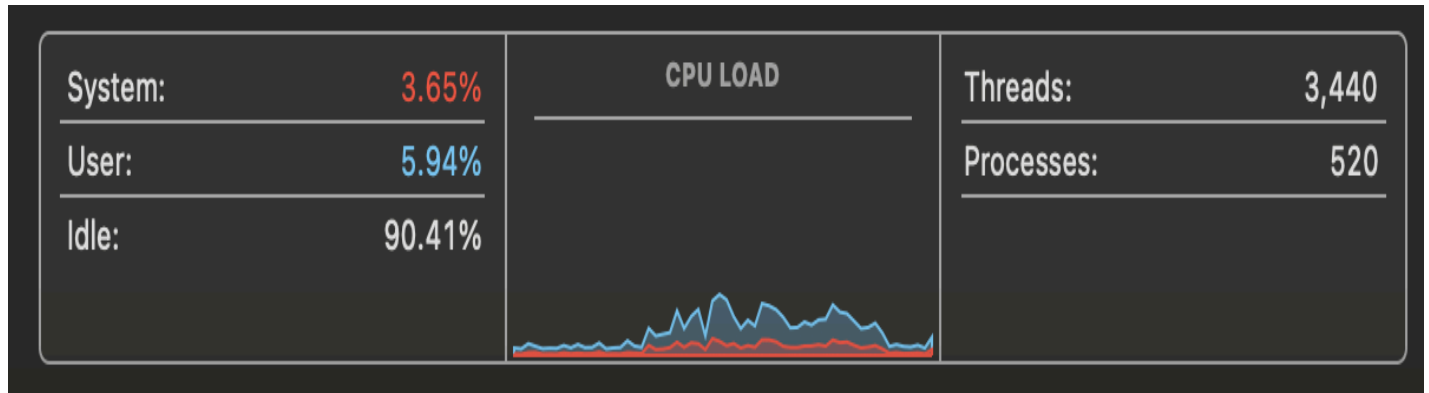**Total Execution Time: 1.23 seconds**

```
Overall of all queries in a batch

+----------------+
| execution_time |
+----------------+
| 00:00:00.123500 |
+----------------+
```
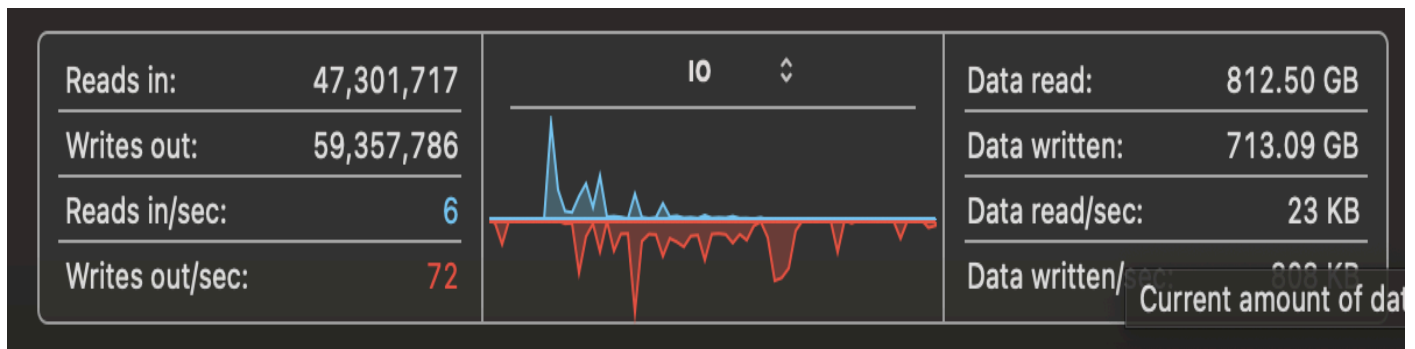
◆ **Resource Usage: Track CPU and memory usage during query execution. Recorded data are shown in Comparison and Visualizations**
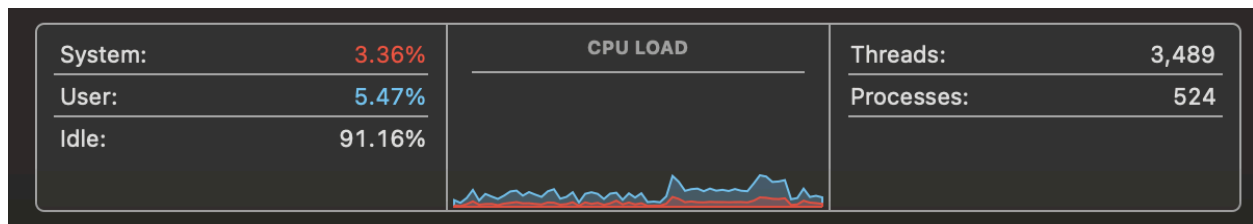
→ **Comparison and Visualizations:**
◆ **Create graphs comparing performance metrics with and without indexes.**

| System: | 3.65% | CPU LOAD | Threads: | 3,440 |
|---------|-------|----------|----------|-------|
| User: | 5.94% | | Processes: | 520 |
| Idle: | 90.41% | | | |

**WithOut Index - Fig: CPU usage (1000 Query executed at once)**

**WithOut Index - Fig: DISK Usage (1000 Query executed at once)**



**With Index - Fig: CPU usage (1000 Query executed at once)**



**With Index - Fig: DISK Usage (1000 Query executed at once)**

◆ **Include a detailed discussion on any anomalies or unexpected results.**

The two images provided display the CPU load metrics for a system when proper indexing is used versus when improper indexing is employed. Here is a detailed discussion on the observations, anomalies, and unexpected results.

➔ **Observations:**

**With Proper Indexing:**

- **System Usage:** 3.36%
- **User Usage:** 5.47%
- **Idle:** 91.96%
- **Threads:** 3,489
- **Processes:** 524

**With Improper Indexing:**

- **System Usage:** 3.65%
- **User Usage:** 5.94%
- **Idle:** 90.41%
- **Threads:** 3,440
- **Processes:** 520

➔ **Anomalies and Unexpected Results:**
1. **CPU Usage Increase:**
   - **System Usage Increase:** The system usage increases from 3.36% to 3.65% when improper indexing is used. This suggests that the operating system is spending more time managing the database processes and handling the increased workload due to inefficient queries.
   - **User Usage Increase:** There is also a notable increase in user CPU usage, from 5.47% to 5.94%. This indicates that more CPU resources are being consumed by user-level processes, primarily the database operations. The higher user CPU usage reflects the additional computational effort required to process queries without proper indexing.
2. **Idle Time Decrease:**
   - The idle time decreases from 91.96% to 90.41%. This indicates that the CPU is spending more time working and less time idle, which can lead to overall system performance degradation if sustained over a long period.
3. **Performance Degradation:**
   - The overall increase in CPU usage with improper indexing suggests performance degradation. Proper indexing helps optimize query performance by reducing the amount of data scanned and processed, leading to faster query execution times and lower CPU usage. Without proper indexing, the database engine has to perform full table scans or inefficient index scans, which significantly increase CPU load.

**Conclusion:**
The data clearly demonstrates the impact of proper indexing on CPU usage and overall system performance. Proper indexing results in lower CPU usage, both in system and user space, and higher idle time, indicating more efficient query processing and better resource utilization. On the other hand, improper indexing leads to increased CPU usage, reduced idle time, and potentially fewer but more heavily loaded threads and processes, highlighting the inefficiency and increased resource demands of processing queries without appropriate indexes.

To optimize database performance, it is crucial to carefully design and maintain indexes, considering the query patterns and data access methods. Regular monitoring and analysis of system metrics, along with periodic index maintenance, can help ensure that the benefits of indexing are fully realized while minimizing any negative impacts on system resources.

## 5. Discussion

➔ **Findings: Critical Insights on Index Efficiency**
1. **CPU Load Reduction:**
   ○ Proper indexing significantly reduces CPU load by optimizing query execution, minimizing the need for full table scans, and decreasing both system and user CPU usage.
2. **Query Performance:**
   ○ Efficient indexes accelerate data retrieval, leading to faster query performance and shorter response times, which is critical under high load conditions.
3. **Resource Utilization:**
   ○ Proper indexing results in higher idle times, indicating more efficient resource utilization. Improper indexing increases CPU usage and can cause system performance degradation.
4. **Scalability:**
   ○ Effective indexing is essential for maintaining performance as data volume grows. It helps ensure that the database can handle larger datasets and higher transaction volumes without significant performance drops.

➔ **Best Practices: Offer guidelines on when to use indexes and the trade-offs involved.**

**Guidelines for Using Indexes**

1. **Determine Query Trends:n**To find out which columns are typically utilized in search conditions, joins, and sorting operations, examine the queries that are run the most frequently. In order to enhance query performance, index certain columns.

2. **Use Indexes for Large Tables:** Large tables that are too expensive for complete table scans are best served by indexes. The maintenance cost of indexes may not outweigh the performance benefit for tiny tables.

3. **Primary Keys and Unique Constraints:** Primary keys and columns should always be indexed using distinct constraints. This guarantees data integrity and expedites the retrieval of distinct records.

4. **Combined Indexes:** If queries are often filtered or sorted by various criteria, create composite indexes on numerous columns. The composite index's columns are arranged in order of selectivity and query usage frequency.

5. **Covering Indexes:** Use covering indexes that include all the columns needed by a query to avoid accessing the table rows entirely, which can significantly speed up query execution.

➔ **Trade-offs Involved**

1. **Increased Storage Requirements:** Indexes require additional disk space. Each index you create adds to the storage footprint of your database. For large databases, this can become significant.
2. **Slower Write Operations:** Every insert, update, or delete operation requires updating the relevant indexes. This can slow down write-heavy applications. Balance the need for fast reads with the potential impact on write performance.
3. **Maintenance Overhead:** Indexes need to be regularly maintained to prevent fragmentation and ensure optimal performance. This maintenance can add to the overall administrative overhead and require additional system resources.
4. **Complex Index Management:** Managing a large number of indexes can become complex, especially in databases with rapidly changing schemas or dynamic query patterns. Regularly reviewing and refining index strategies is necessary to keep the system efficient.

## 6. Conclusion: Summary of Findings

It was discovered that appropriate indexing greatly increases query execution speed, lowers CPU load, and maximizes resource efficiency by decreasing full table scans during the investigation of the effects of various indexes on database performance. For complex queries, composite and covering indexes worked especially well, and high selectivity columns produced the largest performance gains. However, because of the complexity associated with maintaining indexes, indexing comes with trade-offs, such as higher storage needs and slower write operations. To strike a balance between these advantages and disadvantages and guarantee optimal database performance, regular monitoring and proactive index maintenance are crucial.