

Advanced Data Storage

▼ Indexing and query execution and optimization in relational and no-SQL databases. Explain the differences between index types (B+tree, hash, bitmap, map/reduce) and the process of query optimization.

Indexing in databases refers to the process of creating a data structure that improves the speed of data retrieval operations on a table at the cost of additional space and slower writes. When a query is executed on a database, the query optimizer determines the most efficient way to access the data based on the available indexes and other factors such as table statistics and query complexity.

In relational databases, indexing typically involves using B+ trees, which are balanced tree data structures optimized for disk storage. B+ trees provide efficient range queries and support sorted access to data. Hash indexes, on the other hand, use a hash function to map keys to locations in a hash table, providing constant-time lookup but less efficient for range queries. Bitmap indexes store a bitmap for each distinct value in a column, allowing for fast bitwise operations to filter rows based on multiple criteria. Map/reduce indexes are used in NoSQL databases like MongoDB and involve splitting queries into map and reduce phases to parallelize and distribute processing across a cluster of servers.

Query optimization is the process of determining the most efficient execution plan for a given query. This involves analyzing various factors such as available indexes, table statistics, query predicates, and join conditions to choose the optimal access paths and join methods. The query optimizer may consider different access methods like index scans, table scans, and nested loop joins, as well as join algorithms like hash join and merge join. The goal of query optimization is to minimize the query's execution time and resource consumption.

Let's explore the differences between index types and the process of query optimization in both relational and NoSQL databases.

Relational Databases:

Index Types:

1. B+Tree Index:

- Organizes data in a balanced tree structure.
- Suitable for range queries and sorting operations.
- Commonly used in relational databases like MySQL, PostgreSQL, etc.

2. Hash Index:

- Maps keys to their corresponding values using a hash function.
- Efficient for exact match queries but not suitable for range queries.
- Often used in in-memory databases and some relational databases for hash joins and lookups.

Query Optimization Process:

1. Parsing and Parsing Tree Generation:

- Convert the query into a parse tree representing its syntactic structure.

2. Semantic Analysis:

- Check the parse tree for semantic correctness, such as table and column existence.

3. Query Optimization:

- Rewrite the query into an equivalent but more efficient form.
- Generate multiple execution plans.

4. Cost Estimation:

- Assign costs to different execution plans based on factors like access methods, join algorithms, etc.

5. Plan Selection:

- Choose the execution plan with the lowest estimated cost.

6. Plan Execution:

- Execute the selected plan and fetch results.

NoSQL Databases:

Index Types:

1. B+Tree Index:

- Used in some NoSQL databases like MongoDB for indexing fields.
- Supports range queries and sorting.

2. Hash Index:

- Similar to relational databases, used for exact match lookups.
- Often used in key-value stores like Redis.

3. Bitmap Index:

- Efficient for low cardinality columns with a limited set of distinct values.
- Suitable for analytics workloads.

4. Map/Reduce:

- Not exactly an index type but a programming model for processing large data sets.
- Used in NoSQL databases like Hadoop/Hive for distributed processing and aggregation.

Query Optimization Process:

NoSQL databases often have different optimization strategies due to their varied data models and use cases. However, some common steps include:

1. Query Parsing and Analysis:

- Analyze the query to understand its intent and structure.

2. Access Method Selection:

- Choose the appropriate access method based on the query and data model (e.g., key-value lookup, range scan).

3. Data Distribution and Sharding:

- Optimize data distribution across nodes in distributed NoSQL databases to minimize data movement during query execution.

4. Parallelism and Distributed Processing:

- Exploit parallelism and distributed processing capabilities to execute queries efficiently across nodes.

5. Index Selection and Usage:

- Select suitable indexes based on query patterns and workload characteristics to speed up data retrieval.

In summary, while both relational and NoSQL databases employ indexing and query optimization techniques, their specific implementations vary based on the underlying data model, query language, and system architecture.

▼ Database scaling and big data processing. Describe different approaches to database scaling (Vertical and Horizontal scaling, application level, mirroring, partitioning, sharding) and non-database approaches to big data storing and processing (DW, DataLakes, ETL/ELT).

Database Scaling:

Database scaling refers to the process of increasing a database system's capacity to handle growing amounts of data or increasing workload demands. There are two main approaches to database scaling:

▼ Vertical Scaling (Scaling Up):

- Involves adding more resources (CPU, RAM, storage) to a single machine.
- This approach has limitations because it can only scale up to the capacity of a single server, and there's a practical limit to how much resources can be added.
- Suitable for small to medium-sized workloads or when immediate scalability is needed without significant architectural changes.

▼ Horizontal Scaling (Scaling Out):

- Involves distributing the database workload across multiple machines or nodes.

- Provides better scalability as it allows adding more servers to handle increased load.
- Common techniques for horizontal scaling include:
 - **Application-level Sharding:** The application is responsible for partitioning data and routing queries to the appropriate shard.
 - **Mirroring and Replication:** Copies of the database are maintained on multiple servers, and load balancing distributes queries across these servers.
 - **Partitioning:** Data is divided into partitions based on certain criteria (e.g., range partitioning, hash partitioning).
 - **Sharding:** Divides the dataset into smaller, manageable parts called shards, each stored on a separate server.

▼ **Application-Level Sharding:**

- In application-level sharding, the responsibility for partitioning data and routing queries to the appropriate shard is handled within the application logic rather than being managed by the database itself.
- Applications typically employ techniques like consistent hashing or range-based partitioning to determine which shard a particular piece of data belongs to.
- This approach offers more flexibility and control over how data is distributed and accessed but requires careful design and maintenance of the sharding logic within the application code.
- Application-level sharding is commonly used in scenarios where the database system lacks built-in sharding capabilities or when fine-grained control over data distribution is necessary.

▼ **Mirroring and Replication:**

- Mirroring and replication involve maintaining copies of the database on multiple servers, often in different geographic locations.
- The primary purpose of mirroring and replication is to improve fault tolerance, high availability, and disaster recovery.

- In mirroring, changes made to the primary database are asynchronously or synchronously replicated to one or more secondary replicas.
- Replicas can be used for read-only queries, offloading read workloads from the primary database, or serving as failover targets in case of primary database failure.
- While mirroring and replication enhance resilience and scalability, they may introduce complexity in managing data consistency and synchronization between replicas.

▼ Partitioning:

- Partitioning involves dividing a large dataset into smaller, more manageable segments called partitions based on certain criteria such as range partitioning or hash partitioning.
- Range partitioning distributes data based on a specific range of values (e.g., partitioning by date ranges).
- Hash partitioning uses a hash function to distribute data evenly across partitions, typically based on a hash of a specific column's value.
- Partitioning can improve query performance by reducing the amount of data that needs to be scanned for a given query, especially in scenarios where queries are selective and can be directed to specific partitions.
- Database systems often support automatic partitioning based on predefined rules or manual partitioning defined by the user.

▼ Sharding:

- Sharding is a horizontal partitioning technique that involves dividing a database into smaller, independent units called shards, each residing on a separate server or node.
- Each shard typically contains a subset of the data and is responsible for handling a portion of the overall workload.
- Sharding is particularly effective for distributing data and queries across multiple servers, enabling linear scalability as more shards are

added to the system.

- However, sharding introduces complexities in managing data distribution, shard rebalancing, and ensuring data consistency and integrity across shards.
- Sharding is commonly used in large-scale distributed database systems, especially in NoSQL databases like MongoDB, Cassandra, and sharded MySQL setups.

Non-Database Approaches to Big Data Storing and Processing:

1. Data Warehouses (DW):

- A centralized repository for storing structured data from various sources.
- Optimized for analytics and reporting.
- Uses techniques like star or snowflake schemas for organizing data and typically supports SQL-based querying.

2. Data Lakes:

- A large repository that stores structured, semi-structured, and unstructured data in its raw format.
- Offers flexibility in storing diverse data types and enables data exploration and analytics.
- Data lakes can be implemented using distributed file systems like Hadoop HDFS or cloud-based storage solutions.

3. ETL/ELT (Extract, Transform, Load / Extract, Load, Transform):

- Both are data integration processes.
- ETL involves extracting data from various sources, transforming it into a consistent format, and loading it into a target database or data warehouse.
- ELT is a variation where data is first loaded into the target system and then transformed as needed.

- ETL/ELT processes are crucial for data integration, data quality improvement, and preparing data for analysis.

▼ CAP theorem. Explain basic paradigms of SQL and no-SQL databases (ACID/BASE, Brewer's diagram) and illustrate how the database type corresponds with the data structure and purpose in specific cases.

<https://chat.openai.com/share/eca41628-0739-4e05-8861-d6a073da1821>

The CAP theorem, proposed by computer scientist Eric Brewer, states that in a distributed data store, it is impossible to simultaneously guarantee all of the following three properties:

1. **Consistency (C)**: All nodes in the system have the same data at the same time. When a write operation is performed, all subsequent read operations will reflect that write.
2. **Availability (A)**: Every request made to the system receives a response, regardless of the state of the system. Every non-failing node in the system returns a valid response, either with data or an error message.
3. **Partition tolerance (P)**: The system continues to operate despite network partitions (communication failures) between nodes.

According to the CAP theorem, a distributed system can only ensure two out of the three properties at any given time. Different database systems prioritize different aspects of the CAP theorem, leading to the classification of databases into two main paradigms: SQL (relational) and NoSQL (non-relational) databases.

SQL (Relational) Databases:

Basic Paradigm - ACID:

- **ACID** stands for Atomicity, Consistency, Isolation, and Durability.
- **Atomicity**: Transactions are atomic units of work, meaning they either succeed completely or fail, leaving the database in its initial state.
- **Consistency**: Transactions bring the database from one valid state to another, maintaining consistency at all times.

- **Isolation:** Transactions are isolated from each other, ensuring that concurrent transactions do not interfere with each other's operations.
- **Durability:** Once a transaction is committed, its changes are permanently stored in the database, even in the event of system failures.

Correspondence:

- SQL databases like PostgreSQL, MySQL, and Oracle typically prioritize consistency and isolation over availability and partition tolerance.
- They ensure strong data consistency and integrity through features like transactions, foreign key constraints, and ACID compliance.
- SQL databases are well-suited for applications requiring complex queries, transactions, and data integrity guarantees, such as financial systems, e-commerce platforms, and traditional relational data models.

NoSQL (Non-Relational) Databases:

Basic Paradigm - BASE:

- **BASE** stands for Basically Available, Soft state, Eventually consistent.
- **Basically Available:** The system guarantees availability, meaning every request receives a response, but the response may not always be up-to-date or consistent.
- **Soft state:** The state of the system may change over time due to eventual consistency, where updates are propagated asynchronously.
- **Eventually consistent:** The system will eventually converge to a consistent state, given enough time and absence of further updates.

Correspondence:

- NoSQL databases like MongoDB, Cassandra, and DynamoDB prioritize availability and partition tolerance over strong consistency.
- They are designed to handle large-scale distributed systems, offering high availability, fault tolerance, and scalability.
- NoSQL databases are suitable for use cases with massive volumes of data, real-time analytics, and applications where scalability and

performance are critical, such as social networks, IoT platforms, and content management systems.

In summary, the choice between SQL and NoSQL databases depends on the specific requirements of the application, including factors such as data consistency, availability, scalability, and performance. SQL databases prioritize strong consistency and ACID compliance, making them suitable for transactional workloads. In contrast, NoSQL databases prioritize availability and scalability, offering eventual consistency and BASE properties, making them suitable for distributed and high-throughput applications.

▼ Data analyses and processing. Explain different data storage and processing approaches (SQL, no-SQL, DW/DataLake, RDF and their influence on data analyses. What type of storage is suitable for what type of tasks? Describe basic data-related tasks (visualization, deep data analyses, reporting, data processing) and discuss the relationship between data storage type and the task.

Different data storage and processing approaches have varying characteristics and are suitable for different types of tasks. Let's explore each approach and its influence on data analysis:

1. **SQL Databases:**

- SQL databases store data in structured formats using tables with rows and columns.
- Suitable for transactional workloads and applications requiring strong consistency and ACID compliance.
- Well-suited for complex queries, ad-hoc analysis, and reporting due to their support for SQL-based querying and relational data models.
- Commonly used for OLTP (Online Transaction Processing) applications, such as e-commerce systems, banking, and ERP (Enterprise Resource Planning) systems.

2. **NoSQL Databases:**

- NoSQL databases store data in flexible, schema-less formats and are optimized for horizontal scalability and high availability.

- Suitable for handling large-scale distributed systems, real-time analytics, and unstructured or semi-structured data.
- NoSQL databases are commonly used in applications like social media platforms, IoT (Internet of Things) systems, and content management systems where scalability and performance are critical.

3. **Data Warehouses (DW):**

- Data warehouses are centralized repositories that store structured data from various sources for analysis and reporting.
- They are optimized for OLAP (Online Analytical Processing) workloads and support complex queries, aggregations, and historical analysis.
- Data warehouses typically use star or snowflake schemas for organizing data, enabling efficient querying and analysis.
- Suitable for business intelligence (BI) applications, decision support systems, and historical trend analysis.

4. **Data Lakes:**

- Data lakes are large repositories that store structured, semi-structured, and unstructured data in its raw format.
- They offer flexibility in storing diverse data types and support exploratory analysis, data discovery, and machine learning.
- Data lakes are commonly used in big data analytics, data science projects, and advanced analytics applications.

5. **RDF (Resource Description Framework):**

- RDF is a data model for representing information in the form of subject-predicate-object triples, forming a graph structure.
- RDF databases enable semantic querying and reasoning over interconnected data.
- Suitable for applications requiring semantic web technologies, linked data, and knowledge graphs, such as metadata management, ontology-driven applications, and semantic search.

Relationship Between Storage Type and Data-Related Tasks:

1. Visualization:

- SQL databases and data warehouses are suitable for visualization tasks requiring structured data and support for SQL-based querying.
- NoSQL databases and data lakes are suitable for visualization of unstructured or semi-structured data, offering flexibility in data storage and exploration.

2. Deep Data Analysis:

- Data warehouses are ideal for deep data analysis tasks involving complex queries, aggregations, and historical trend analysis.
- NoSQL databases and data lakes are suitable for deep data analysis requiring large-scale distributed processing, real-time analytics, and machine learning.

3. Reporting:

- SQL databases and data warehouses are commonly used for reporting tasks due to their support for structured data and SQL-based reporting tools.
- NoSQL databases and data lakes can also support reporting, especially for real-time or ad-hoc reporting needs, but may require additional preprocessing steps.

4. Data Processing:

- SQL databases are suitable for transactional data processing tasks requiring ACID compliance and strong consistency guarantees.
- NoSQL databases and data lakes are suitable for batch processing, real-time stream processing, and big data analytics requiring horizontal scalability and flexibility in data storage.

In summary, the choice of data storage type depends on the specific requirements of the data-related tasks, including the nature of the data, scalability needs, performance requirements, and the complexity of the

analysis. Each storage approach offers unique features and trade-offs, allowing organizations to select the most suitable solution for their use cases.

▼ **Methods of data analysis.** Describe the basic data analysis methods (histogram, boxplot, correlation, pairwise covariance, ...) and interpret their results on a simple example. Describe principles of advanced data analysis methods (clustering, regression, classification) and describe the interpretation of their possible outputs.

Basic Data Analysis Methods:

1. Histogram:

- A histogram is a graphical representation of the distribution of numerical data.
- It consists of a series of bars, where each bar represents the frequency of data falling within a specific range or bin.
- Histograms are useful for visualizing the central tendency, dispersion, and skewness of a dataset.

2. Boxplot:

- A boxplot, also known as a box-and-whisker plot, is a graphical summary of the distribution of numerical data.
- It displays the median, quartiles, and potential outliers of the dataset.
- Boxplots are useful for identifying outliers, comparing distributions, and detecting skewness or symmetry in the data.

3. Correlation:

- Correlation measures the strength and direction of the linear relationship between two variables.
- It is represented by the correlation coefficient, which ranges from -1 to 1.
- A positive correlation indicates that the variables move in the same direction, while a negative correlation indicates opposite movement.
- Correlation is useful for identifying associations between variables and understanding their relationships.

4. **Pairwise Covariance:**

- Covariance measures the degree to which two variables change together.
- It indicates the direction of the linear relationship between variables but is not standardized like correlation.
- Positive covariance indicates that the variables tend to increase or decrease together, while negative covariance indicates opposite movement.
- Pairwise covariance is useful for understanding the joint variability of multiple variables in a dataset.

Interpretation Example:

Let's consider a simple example where we have a dataset of students' exam scores in two subjects: Math and Science. We want to analyze the relationship between these scores using the methods described above.

- **Histogram:** We can create histograms for Math and Science scores to visualize their distributions and identify any patterns or outliers.
- **Boxplot:** A boxplot can show us the median, quartiles, and any potential outliers in both Math and Science scores, allowing us to compare their distributions.
- **Correlation:** We can calculate the correlation coefficient between Math and Science scores to determine the strength and direction of their linear relationship. A high positive correlation would indicate that students who perform well in Math also tend to perform well in Science.
- **Pairwise Covariance:** Pairwise covariance can provide insights into how Math and Science scores vary together. A positive covariance would suggest that students with high Math scores also tend to have high Science scores.

Advanced Data Analysis Methods:

1. Clustering:

- Clustering is a technique used to group similar objects or data points together based on their characteristics.
- It aims to identify inherent patterns or structures within the data.
- Common clustering algorithms include K-means, hierarchical clustering, and DBSCAN.

2. **Regression:**

- Regression analysis is used to model the relationship between a dependent variable and one or more independent variables.
- It helps in understanding how changes in the independent variables affect the dependent variable.
- Linear regression, logistic regression, and polynomial regression are common regression techniques.

3. **Classification:**

- Classification assigns predefined categories or labels to data points based on their features.
- It is used for predicting categorical outcomes or classifying data into predefined classes.
- Common classification algorithms include decision trees, support vector machines (SVM), and random forests.

Interpretation of Advanced Methods:

1. **Clustering:** After performing clustering analysis, the output would include clusters of similar data points. Interpretation involves understanding the characteristics of each cluster and identifying meaningful patterns or segments within the data.
2. **Regression:** Regression analysis provides insights into the relationship between the dependent and independent variables. Interpretation involves examining the coefficients of the regression model to understand the magnitude and direction of the effect of each independent variable on the dependent variable.

3. **Classification:** The output of classification algorithms includes a model that can predict the class or category of new data points. Interpretation involves evaluating the model's performance metrics such as accuracy, precision, recall, and F1-score, and understanding the significance of the features in predicting the outcome.

In summary, basic data analysis methods such as histograms, boxplots, correlation, and covariance provide insights into the distribution and relationships within a dataset, while advanced methods like clustering, regression, and classification enable deeper analysis and predictive modeling. Interpretation of the results involves understanding the underlying patterns, relationships, and predictive power of the models.