# ADVANCED DATABASE MANAGEMENT SYSTEM
# ASSIGNMENT II

**ABHINAND K S**
**ROLL NO : 03**
**RMCA-A**

# NEXT GENERATION DATABASE

The phrase "next-generation database" often refers to a new wave of database technologies and techniques with the goal of addressing the shortcomings of conventional relational databases and offering better capabilities to manage contemporary data requirements. These next-generation databases frequently place a strong emphasis on scalability, performance, flexibility, and the capacity to manage a variety of data types, including structured, semi-structured, and unstructured data.

Next-generation databases aim to address the challenges posed by big data, real-time data processing, high concurrency, scalability, and other modern data requirements. They offer alternatives and extensions to traditional relational databases, providing organizations with more options to meet their specific data management needs.

Some key characteristics and technologies associated with next-generation databases include:

1. NoSQL Databases: NoSQL (Not Only SQL) databases emerged as an alternative to traditional relational databases. They are designed to handle large volumes of data and offer flexible schema designs, horizontal scalability, and high availability. NoSQL databases include various types such as key-value stores, document stores, column-family stores, and graph databases.

2. NewSQL Databases: NewSQL databases combine the benefits of traditional relational databases with the scalability and performance advantages of NoSQL databases. They aim to provide high throughput, low latency, and strong consistency, while still supporting standard SQL and ACID (Atomicity, Consistency, Isolation, Durability) transactions.

3. Distributed Databases: Next-generation databases often employ distributed architectures that allow data to be spread across multiple nodes or clusters, enabling horizontal scalability, fault tolerance, and high availability. Distributed databases employ techniques like sharding, replication, and consistent hashing to partition and replicate data across nodes.

4. In-Memory Databases: In-memory databases store data in the main memory (RAM) rather than on disk, allowing for extremely fast read and write operations. They are suitable for applications that require low-latency access to data, such as real-time analytics, caching, and high-speed transactions.

5. Cloud Databases: Cloud databases are designed to operate in cloud computing environments and take advantage of the scalability, elasticity, and pay-as-you-go models offered by cloud providers. These databases can be easily provisioned, scaled, and managed, and oftEn integrate with other cloud services.

6. Hybrid Databases: Hybrid databases combine different database models or approaches to provide a unified platform for managing various types of data. For example, a hybrid database might combine a document store and a graph database to handle both document-oriented and graph-related use cases.

# DISTRIBUTED RELATIONAL DATABASE

A distributed relational database is a type of database system that spans across multiple nodes or servers, allowing data to be partitioned and replicated for improved scalability, fault tolerance, and performance.

In a distributed relational database, data is distributed across multiple nodes, with each node responsible for storing and processing a subset of the data. This distribution enables horizontal scalability, as more nodes can be added to accommodate growing data volumes and user loads.

To ensure data availability and fault tolerance, the database system replicates data across multiple nodes. This replication provides redundancy, allowing the system to continue functioning even if some nodes fail.

Distributed relational databases employ various techniques to manage data consistency and ensure that updates to the database are properly synchronized across nodes. These techniques include distributed transaction processing, distributed locking mechanisms, and consensus protocols.

The benefits of distributed relational databases include improved scalability to handle large datasets and high user loads, increased fault tolerance and availability through data replication, and the ability to provide low-latency access to data by distributing it geographically closer to users.

However, distributed relational databases also introduce additional complexity in terms of data partitioning, replication management, and distributed query processing. It requires careful planning and design to ensure proper data distribution, consistency, and efficient data retrieval.

# NON- RELATIONAL DISTRIBUTED DATABASE

A non-relational distributed database, also known as a NoSQL database, is a type of database system that provides a flexible and scalable approach to storing and managing data. Unlike traditional relational databases, NoSQL databases do not rely on a fixed schema or use SQL for data manipulation.

In a non-relational distributed database, data is typically stored in a denormalized manner, allowing for greater flexibility and performance in handling large volumes of structured, semi-structured, and unstructured data. These databases are designed to handle massive amounts of data and provide horizontal scalability by distributing data across multiple nodes or servers.

NoSQL databases employ various data models, including key-value stores, document stores, column-family stores, and graph databases, to cater to different types of data and use cases. Each data model offers unique advantages and trade-offs in terms of data access patterns, data organization, and query capabilities.

Distributed NoSQL databases utilize a distributed architecture, where data is partitioned and replicated across nodes to ensure fault tolerance, high availability, and scalability. They often employ techniques like sharding, consistent hashing, and replication mechanisms to distribute data and handle data consistency.

The benefits of non-relational distributed databases include the ability to handle large-scale and high-velocity data, support for flexible schema designs, scalability to accommodate growing data volumes and user loads, and the ability to operate in distributed and cloud environments.

However, non-relational distributed databases may have trade-offs in terms of data consistency models, limited query capabilities compared to SQL, and increased complexity in data management. Each NoSQL database model has its

own strengths and weaknesses, and the choice of a specific NoSQL database depends on the requirements of the application.

## MONGODB SHARDING AND REPLICATION

MongoDB sharding and replication are two important features that enhance the scalability, availability, and fault tolerance of MongoDB, a popular NoSQL database.

## 1. Sharding:

Sharding is the process of horizontally partitioning data across multiple MongoDB instances called shards. Each shard contains a subset of the data, allowing the database to handle large datasets and high workloads. Sharding enables MongoDB to distribute data across multiple servers, providing improved scalability and performance.

Key points about MongoDB sharding:

- Sharding is achieved by defining a shard key, which determines how data is partitioned across shards based on a chosen field or set of fields.

- The sharded cluster consists of three main components: shards (servers that store data subsets), mongos (routers that direct client requests to the appropriate shard), and config servers (metadata storage for the sharded cluster).

- MongoDB's balancer automatically moves data between shards to maintain an even data distribution.

- Sharding enables horizontal scalability as additional shards can be added to accommodate data growth and increased user loads.

- It allows for efficient parallel execution of queries across shards, improving query performance.

- Sharding does not provide data redundancy or high availability on its own. Replication is used in conjunction with sharding to ensure fault tolerance and data durability.

## 2. Replication:

Replication in MongoDB involves creating multiple copies of data, called replicas, to ensure data availability, durability, and fault tolerance. Replication provides redundancy and allows for automatic failover in case of primary replica failure.

Key points about MongoDB replication:

- Replication is achieved through a primary-secondary replication model. Each replica set consists of a primary node that accepts write operations and multiple secondary nodes that replicate data from the primary.

- When the primary node fails, an election process selects a new primary from the available secondary nodes.

- Read operations can be directed to secondary nodes to distribute the read workload and improve read scalability. However, secondary nodes may have some degree of replication lag, so read preference settings need to be configured accordingly.

- Replication provides data durability, as write operations are recorded in the primary node's oplog (operation log) and replicated to secondary nodes. In the event of a primary failure, a secondary can be promoted to primary without data loss.

- Replication can be used independently of sharding and is often combined with sharding to achieve both scalability and fault tolerance.

By combining sharding and replication, MongoDB can scale horizontally to handle large amounts of data and high workloads while ensuring data availability, durability, and fault tolerance. This makes MongoDB a powerful choice for data-intensive applications that require both scalability and reliability.


# HBASE

HBase is an open-source, distributed, columnar NoSQL database built on top of the Apache Hadoop ecosystem. It is designed to provide real-time read and write access to large-scale, structured and semi-structured data.


Key characteristics of HBase include:


1. Distributed Architecture: HBase is designed to operate on a distributed cluster of commodity hardware. It leverages the Hadoop Distributed File System (HDFS) for storage, allowing it to handle large datasets that are partitioned and replicated across multiple nodes.


2. Columnar Storage: HBase organizes data in a columnar manner, where data is stored by column rather than by row. This enables efficient read and write operations on specific columns, making it well-suited for applications that require random access to specific attributes or fields of a record.


3. Scalability and High Availability: HBase provides automatic sharding and load balancing capabilities, allowing data to be distributed across multiple nodes. It can scale horizontally by adding more nodes to the cluster, accommodating growing data volumes and increasing workloads. HBase also offers fault tolerance and high availability by replicating data across multiple nodes.


4. Strong Consistency: HBase guarantees strong consistency within a row, meaning that concurrent read and write operations on the same row will always see the latest committed version of the data.


5. Flexible Schema: HBase is schema-less, meaning that it does not enforce a fixed schema on the data. This provides flexibility in terms of data structure and allows for the addition or modification of columns dynamically.

6. Integrated with Hadoop Ecosystem: HBase integrates well with other components of the Hadoop ecosystem, such as Apache Hive for data warehousing, Apache Spark for data processing, and Apache Kafka for real-time streaming.

HBase is commonly used for applications that require random real-time access to large datasets, such as social media analytics, time-series data analysis, sensor data storage, and log data processing. Its scalability, fault tolerance, and efficient columnar storage make it a suitable choice for handling big data workloads in a distributed environment.

# CASSANDRA

Cassandra is an open-source distributed NoSQL database designed to handle large amounts of data across multiple commodity servers with high availability and scalability. It was developed at Facebook and is now maintained by the Apache Software Foundation.

Cassandra's distributed, fault-tolerant, and highly scalable nature makes it an ideal choice for applications that require continuous availability, massive scalability, and high-performance data operations. It has been widely adopted by organizations dealing with big data and mission-critical applications where data consistency, reliability, and scalability are crucial.

Key features and characteristics of Cassandra include:

1. Distributed Architecture: Cassandra operates on a distributed cluster of nodes, with data distributed across multiple servers. This architecture allows for horizontal scalability, enabling seamless scaling by adding or removing nodes to accommodate growing data volumes and user loads.

2. No Single Point of Failure: Cassandra is designed to be highly available and fault-tolerant. It achieves this by employing a peer-to-peer distributed model, where all nodes are equal and there is no single point of failure. Data is automatically replicated across multiple nodes to ensure durability and availability.

3. High Performance: Cassandra is optimized for fast read and write operations, making it suitable for use cases that require high throughput and low-latency access to data. It achieves this through features such as log-structured storage, asynchronous replication, and distributed query processing.

4. Linear Scalability: Cassandra's architecture allows it to scale linearly as additional nodes are added to the cluster. It leverages consistent hashing to distribute data evenly across nodes, ensuring load balancing and efficient data distribution.

5. Tunable Consistency: Cassandra offers tunable consistency, allowing developers to define the level of consistency required for read and write operations. It provides various consistency levels, such as strong consistency, eventual consistency, and tunable consistency in between, enabling flexibility based on application requirements

6. Data Model: Cassandra follows a column-family data model, similar to other NoSQL databases. It allows for flexible schemas, enabling dynamic addition or modification of columns without disrupting existing data.

7. Wide Range of Use Cases: Cassandra is well-suited for use cases that involve high data volumes, high write throughput, and low-latency data access, such as real-time analytics, time-series data, recommendation systems, messaging platforms, and IoT applications.

# CAP THEORAM

The CAP theorem, also known as Brewer's theorem, is a fundamental concept in distributed computing that states that it is impossible for a distributed system to simultaneously guarantee all of the following three properties:

1. Consistency (C): Every read operation in the system will return the most recent write or an error. Consistency ensures that all nodes in a distributed system see the same data at the same time.

2. Availability (A): Every request made to the system will receive a response, even in the presence of failures. Availability implies that the system remains operational and accessible despite individual node failures.

3. Partition tolerance (P): The system can continue to operate and provide meaningful responses even when network partitions occur, causing communication failures between nodes.

According to the CAP theorem, a distributed system can only prioritize two out of the three properties: consistency, availability, or partition tolerance. This means that in the event of a network partition, a system must choose between maintaining consistency and sacrificing availability (CP systems) or sacrificing consistency to ensure availability (AP systems).

It's important to note that the CAP theorem assumes that network partitions are inevitable in a distributed system. In real-world scenarios, network partitions can occur due to network failures, latency, or other issues.

Many distributed databases and systems, such as relational databases, adhere to the ACID (Atomicity, Consistency, Isolation, Durability) properties, prioritizing consistency over availability in the presence of partitions. On the other hand, NoSQL databases and systems often prioritize availability and partition tolerance over strict consistency, offering eventual consistency or relaxed consistency models.

It's worth mentioning that while the CAP theorem provides a theoretical framework for understanding the trade-offs in distributed systems, it does not dictate the specific design decisions or guarantees of a particular system. Different systems make different choices based on their specific requirements and use cases.