

# Distributed Database Management System

A **Distributed Database Management System (DDBMS)** is a software system that manages a distributed database, ensuring that the data is stored across multiple locations in a manner that appears unified to the user.

## Definition -

A **Distributed Database Management System (DDBMS)** is a type of database management system that enables the management, storage, and retrieval of data from a database that is distributed across multiple locations, which can be within the same physical location or across different geographical areas.

## Key Characteristics -

1. **Data Distribution:** Data is stored across multiple sites, which can be geographically dispersed.
2. **Transparency:** Users should interact with the system as if it is a single, unified database, irrespective of the underlying data distribution.
3. **Replication:** Data can be replicated across multiple sites to enhance reliability and availability.
4. **Scalability:** Easily scalable by adding more nodes or locations without significant impact on performance.
5. **Fault Tolerance:** Provides high fault tolerance through data redundancy and replication.
6. **Consistency and Synchronization:** Ensures data consistency and synchronization across all distributed sites.
7. **Concurrency Control:** Manages simultaneous data access and modifications across distributed sites to maintain data integrity.

## Advantages -

- **Improved Reliability:** With data distributed across multiple sites, the failure of one site does not render the entire system inoperative.
- **Enhanced Performance:** Data can be processed locally, reducing access time and network load.
- **Scalability:** Easy to add more locations or nodes to handle growing amounts of data and users.
- **Location Transparency:** Users do not need to know the physical location of data.

## Disadvantages -

- **Complexity:** Managing and maintaining a distributed database system is more complex compared to a centralized system.
- **Security:** Ensuring data security across multiple sites requires robust measures.
- **Cost:** Setting up and maintaining a distributed database can be more expensive due to the need for additional hardware and network infrastructure.

## Uses for distributed databases -

- The corporate management information system makes use of it.
- Multimedia apps utilize it.
- Used in hotel chains, military command systems, etc.
- The production control system also makes use of it

## Examples -

- **Google Spanner:** A globally distributed database service.
- **Amazon Aurora:** A distributed, highly scalable relational database service.
- **Microsoft Azure SQL Database:** Provides cloud-based distributed database services.

**NOTE: *A Distributed Database Management System plays a crucial role in modern data management by ensuring high availability, reliability, and scalability across geographically dispersed locations.***

**Distributed Database Design** – Distributed Database Design is a critical aspect of creating a distributed database system. It involves the layout, architecture, and planning of a database that is spread across multiple locations, which can be within the same site or geographically dispersed.

#### **Objectives -**

1. **Data Availability:** Ensure that data is accessible whenever and wherever it is needed.
2. **Performance:** Optimize data retrieval and update times by distributing the load across multiple sites.
3. **Scalability:** Easily add one or more nodes or databases without significant impact on performance.
4. **Fault Tolerance:** Ensure the system can continue operating even if some nodes fail.

#### **Advantages -**

1. **Data Availability and Reliability:** A Distributed DBMS ensures high availability and reliability because data is replicated across multiple sites. If one site fails, the system can still function using the data from other sites.
2. **Scalability:** Distributed databases are easily scalable. Additional data and users can be accommodated by adding more nodes to the network without significantly impacting performance.
3. **Performance Improvement:** By distributing data closer to where it is needed, Distributed DBMS can reduce data access time and network traffic. This localized data processing leads to faster query responses and overall improved performance.
4. **Resource Sharing:** Different sites can share the computational resources and data, leading to better utilization of hardware and software resources.
5. **Data Integrity and Consistency:** Distributed DBMS enforces data integrity and consistency through distributed transactions and concurrency control mechanisms, ensuring that all replicas of data are synchronized.
6. **Cost Efficiency:** It can be more cost-effective than a large centralized system. Distributed systems allow the use of less expensive, smaller servers rather than one large, costly server.
7. **Flexibility:** It allows for flexible and modular growth, as new locations can be added without requiring significant changes to the existing setup.
8. **Geographic Distribution:** A Distributed DBMS can store data in locations that are geographically dispersed, providing advantages for multinational organizations that need local data accessibility.
9. **Fault Tolerance:** The system can continue to operate even in the event of partial system failures. Data replication and distributed control provide redundancy that enhances fault tolerance.
10. **Data Autonomy:** Local sites can have control over their own data, which can lead to more efficient data management and security at the local level.

#### **Disadvantages -**

1. **Complexity:** Designing, implementing, and maintaining a DDBMS is more complex compared to a centralized database. The complexity arises from managing data distribution, replication, and synchronization across multiple sites.
2. **Cost:** The initial setup cost for a DDBMS can be high. This includes the cost of additional hardware, software, and networking infrastructure. Additionally, ongoing maintenance and management costs can be significant.
3. **Security:** Ensuring data security across multiple locations is challenging. Distributed databases are more vulnerable to security breaches, unauthorized access, and data interception during transmission.

4. **Performance:** The performance of a DDBMS can be affected by network issues such as latency, bandwidth limitations, and network failures. Accessing and updating data across multiple sites can introduce delays and reduce system efficiency.
5. **Consistency and Integrity:** Maintaining data consistency and integrity across distributed sites is complex. Synchronizing replicated data and ensuring that all copies of the data are updated simultaneously can be challenging and resource-intensive.
6. **Reliability:** While DDBMS can be more fault-tolerant, they also introduce additional points of failure. Network outages, server crashes, or hardware failures at any site can disrupt the entire system.
7. **Coordination and Management:** Coordinating and managing a distributed database involves significant administrative overhead. Tasks such as load balancing, query optimization, and transaction management require careful planning and execution.
8. **Backup and Recovery:** Implementing efficient backup and recovery procedures is more complex in a distributed environment. Ensuring that backups are consistent and synchronized across all sites can be challenging.
9. **Data Redundancy:** While replication improves availability, it also introduces data redundancy, leading to increased storage requirements and potential data inconsistency issues.
10. **Data Distribution Decisions:** Deciding how to distribute data effectively based on access patterns, network considerations, and workload can be difficult and may require regular adjustments as conditions change.

### Architecture of DDBMS -

The architecture of a Distributed Database Management System (DDBMS) is designed to manage data that is distributed across multiple locations. This architecture ensures efficient data management, synchronization, and communication between the distributed sites.

- **Database Fragmentation -**
  1. **Horizontal Fragmentation:** Dividing a table into rows and distributing them across different sites.
  2. **Vertical Fragmentation:** Dividing a table into columns and distributing them.
  3. **Hybrid Fragmentation:** Combination of horizontal and vertical fragmentation.
- **Replication and Distribution -**
  1. **Full Replication:** Entire database is replicated at all sites.
  2. **Partial Replication:** Only necessary parts of the database are replicated at certain sites.
  3. **No Replication:** Data is distributed without any duplication.
- **Components of DDBMS Architecture -**
  1. **Local Database System (LDS):**
    - Each site has its own local database system that manages the data stored locally.
    - Responsible for local query processing and transaction management.
  2. **Global Database System (GDS):**
    - Manages global transactions and ensures consistency across all sites.
    - Coordinates distributed queries and transactions.
  3. **Communication Network:**
    - Facilitates communication between different sites.
    - Can include local area networks (LANs), wide area networks (WANs), and the internet.
- **Layers of DDBMS Architecture -**
  1. **Presentation Layer:**
    - Interface for users to interact with the database.
    - Provides query tools, forms, and reports.

2. **Application Layer:**
    - Manages application logic and processing.
    - Handles user requests and communicates with the database layer.
  3. **Database Layer:**
    - Responsible for managing data storage, retrieval, and update operations.
    - Includes the local and global database systems.
- **Distributed Query Processing -**
    1. **Query Decomposition:** Breaking down a query into sub-queries that can be processed at different sites.
    2. **Query Optimization:** Choosing the most efficient way to execute a distributed query.
    3. **Join Processing:** Efficiently joining data from different sites.
  - **Distributed Transaction Management -**
    1. **Concurrency Control:** Ensuring that concurrent transactions do not interfere with each other.
    2. **Two-Phase Commit Protocol:** Ensuring atomicity of distributed transactions.
      - **Prepare Phase:** All sites prepare to commit the transaction.
      - **Commit Phase:** All sites commit the transaction if all are prepared; otherwise, they abort.
  - **Data Consistency and Synchronization -**
    - **Consistency Models:** Ensuring that data remains consistent across distributed sites.
      - **Strong Consistency:** All replicas are updated before a transaction is considered complete.
      - **Eventual Consistency:** Updates propagate to all replicas eventually, but not immediately.
  - **Fault Tolerance and Recovery -**
    1. **Redundancy:** Replicating data to ensure availability in case of site failures.
    2. **Logging and Checkpointing:** Keeping logs of transactions and periodic snapshots of data for recovery.

### Types of DDBMS -

Distributed Database Management Systems (DDBMS) can be classified based on various criteria, such as the degree of distribution, the level of autonomy, and the type of homogeneity. Here are the main types:

- **Based on Degree of Distribution -**
  1. **Fully Distributed DDBMS** - All the data is distributed across multiple locations. Each site has a subset of the database, and no site contains the entire database.
  2. **Partially Distributed DDBMS** - Some data is distributed across multiple locations, but there may also be some centralization of data at specific sites.
- **Based on Autonomy -**
  1. **Homogeneous DDBMS -**
    - **Characteristics:** All sites use the same DBMS software.
    - **Advantage:** Easier to manage and ensure consistency because of uniform software.
    - **Example:** Multiple sites using the same version of Oracle Database.

## 2. Heterogeneous DDBMS -

- **Characteristics:** Different sites may use different DBMS software, potentially even different schemas.
- **Advantage:** Flexibility to integrate diverse systems and software.
- **Example:** Sites using different DBMS software like Oracle, SQL Server, and MySQL.

### • Based on Degree of Homogeneity -

#### 1. Federated DDBMS -

- **Characteristics:** Multiple autonomous DBMSs are integrated, but each retains its autonomy.
- **Advantage:** Allows different databases to work together while maintaining their individual autonomy.
- **Example:** A combination of different databases that cooperate under a unified framework, such as integrating different department databases within a university.

#### 2. Multi-database System (MDBS) -

- **Characteristics:** A more loosely integrated system where multiple autonomous databases work together without a common schema.
- **Advantage:** Greater independence and flexibility.
- **Example:** Different organizational units with their own databases sharing data occasionally.

### • Based on Level of Transparency -

#### 1. Transparent DDBMS -

- a. **Characteristics:** The distributed nature of the database is hidden from users. They interact as if it were a single database.
- b. **Advantage:** Simplifies user interaction and data access.
- c. **Example:** Users can perform queries without knowing where the data is physically located.

#### 2. Non-Transparent DDBMS -

- a. **Characteristics:** Users are aware of the distribution and must specify where the data is located.
- b. **Advantage:** Greater control and efficiency for users who understand the data distribution.
- c. **Example:** Users need to write queries that specify the location of data.

■ **Data Fragmentation** - Data Fragmentation is a process in database management, particularly in distributed database systems, where a database is divided into smaller, more manageable pieces called **fragments**. These fragments can then be stored across multiple locations to enhance performance, availability, and scalability. Data fragmentation is crucial for efficient data distribution and management in distributed databases.

*Data fragmentation is a crucial technique in distributed database management that involves dividing a database into smaller, more manageable fragments. While it offers significant advantages in terms of performance, availability, and scalability, it also introduces complexity in data management and consistency. Effective data fragmentation strategies are essential for optimizing distributed database systems.*

---

## ■ Types of Data Fragmentation -

1. **Horizontal Fragmentation:** Dividing a table into subsets of rows (tuples) and distributing these subsets across different sites. Each fragment is a subset of the rows of the original table.

*Example: If a table contains customer data, horizontal fragmentation might involve storing rows for customers in different regions at different sites.*

2. **Vertical Fragmentation:** Dividing a table into subsets of columns (attributes) and distributing these subsets across different sites. Each fragment is a subset of the columns of the original table.

*Example: If a table contains employee data, vertical fragmentation might involve storing personal information at one site and job-related information at another.*

3. **Hybrid Fragmentation:** A combination of horizontal and vertical fragmentation. A table is first divided into horizontal fragments, and then each horizontal fragment is further divided into vertical fragments (or vice versa).

*Example: First dividing customer data by region (horizontal fragmentation) and then splitting each region's data into personal and order-related columns (vertical fragmentation).*

## ■ Advantages of Data Fragmentation -

- **Improved Performance:** Localized data access reduces the need for extensive data transfer over the network, resulting in faster query processing and transaction times.
- **Increased Availability:** Data fragmentation can enhance the availability of the database by distributing fragments across multiple sites. If one site fails, other sites can still function using their fragments.
- **Scalability:** Data fragmentation allows the database to grow and scale by adding new fragments and distributing them across additional sites.
- **Reduced Storage Costs:** By storing only relevant fragments at specific sites, data storage costs can be optimized.

## ■ Disadvantages of Data Fragmentation -

- 1) **Complexity:** Managing and maintaining fragmented data can be complex, requiring sophisticated algorithms for data distribution, synchronization, and querying.
- 2) **Data Consistency:** Ensuring data consistency across multiple fragments and sites can be challenging, especially in systems with frequent updates.
- 3) **Query Processing Overhead:** Fragmented data may lead to increased complexity in query processing, as queries may need to be reconstructed from multiple fragments.

Feature	RDBMS	DDBMS	OODBMS
<b>Definition</b>	Centralized database system	Database spread across multiple locations	Database system using object-oriented principles
<b>Data Storage</b>	Tables (relations)	Tables (relations) spread across sites	Objects and classes
<b>Architecture</b>	Simple, centralized	Complex, distributed	Can be centralized or distributed
<b>Data Integrity and Consistency</b>	High, centrally managed	More challenging, requires synchronization	High, through encapsulation and inheritance
<b>Data Model</b>	Relational (tables)	Relational (tables)	Object-oriented (objects/classes)
<b>Examples</b>	MySQL, PostgreSQL, Oracle	Google Spanner, Amazon Aurora, Cassandra	db4o, ObjectDB, Versant
<b>Performance</b>	May decline with scale	Generally better for large scale, but network latency can	Efficient for complex data and relationships
<b>Scalability</b>	Limited	High	High, especially for complex data structures
<b>Fault Tolerance</b>	Lower	Higher	Can be high with distributed architecture
<b>Maintenance</b>	Easier	More complex	Can be complex due to object-oriented principles
<b>Query Language</b>	SQL	SQL with distributed query support	OQL (Object Query Language) or extended SQL
<b>Suitable For</b>	Traditional business applications	Large-scale, geographically dispersed applications	Applications with complex data and relationships (e.g., CAD/CAM, multimedia)
<b>Data Relationships</b>	Managed through foreign keys	Managed through foreign keys and replication	Managed through object relationships and inheritance



#### NOTE:

- **RDBMS:** Best suited for applications with straightforward, well-defined relationships and centralized data needs. It is easier to maintain but can face performance and scalability issues as the system grows.
- **DDBMS:** Ideal for applications requiring high availability, fault tolerance, and scalability, especially those spread across multiple geographic locations. It has complex management requirements.
- **OODBMS:** Suitable for applications with complex data models and relationships, such as CAD/CAM, multimedia, and scientific applications. It leverages object-oriented principles to handle complex data but can be challenging to implement and manage.

■ **Data Replication** - Data Replication is a process in database management where copies of data are maintained at multiple locations to ensure data availability, reliability, and fault tolerance. This technique is essential in distributed systems where data consistency and accessibility are critical.

*Data replication is a fundamental technique in database management, particularly in distributed systems. It ensures high availability, fault tolerance, and improved performance. However, it also introduces complexity in terms of data management and consistency. Balancing these factors is crucial for effective replication strategies.*

---

#### ■ Key Concepts –

- **Replication Types:**
  - **Full Replication:** The entire database is replicated at each site. This ensures high availability but can be resource-intensive.
  - **Partial Replication:** Only a subset of the database is replicated at various sites. This approach balances the benefits of replication with resource considerations.
- **Replication Strategies:**
  - **Synchronous Replication:** Changes are replicated to all copies of the data simultaneously. This ensures consistency but can impact performance due to the overhead of maintaining synchronization.
  - **Asynchronous Replication:** Changes are propagated to replicas after the primary operation is completed. This improves performance but may lead to temporary inconsistencies.
- **Replication Schemes:**
  - **Master-Slave Replication:** One node acts as the master, handling all write operations, while other nodes (slaves) handle read operations and replicate the master's data.
  - **Multi-Master Replication:** Multiple nodes can handle write operations, with data being replicated among all masters. This provides higher availability and write throughput but requires more complex conflict resolution mechanisms.



## ■ Advantages of Data Replication –

- a) **Increased Availability:** Ensures that data is available even if one or more sites fail, enhancing system reliability.
- b) **Improved Performance:** By placing copies of data closer to users, read and write operations can be performed more quickly.
- c) **Fault Tolerance:** Replicated data ensures that the system can continue functioning in the event of a hardware or software failure at any site.
- d) **Load Balancing:** Distributes the load across multiple sites, preventing any single node from becoming a bottleneck.

## ■ Disadvantages of Data Replication -

- a) **Increased Storage Costs:** Maintaining multiple copies of data requires additional storage resources.
- b) **Complexity in Data Management:** Ensuring consistency and synchronization across all replicas can be complex and resource-intensive.
- c) **Potential for Data Inconsistency:** Especially in asynchronous replication, there can be a lag between updates, leading to temporary inconsistencies.
- d) **Higher Maintenance Overhead:** Regularly updating and synchronizing replicas requires ongoing effort and resources.

## ■ Use Cases -

- a) **Content Delivery Networks (CDNs):** Replicating web content across multiple servers worldwide to ensure fast access for users.
- b) **Distributed Databases:** Ensuring data availability and fault tolerance in distributed database systems.
- c) **Disaster Recovery:** Replicating data to geographically distant locations to safeguard against data loss due to natural disasters or other catastrophic events.

■ **Data Allocation** - Data Allocation is a crucial aspect of database management, especially in distributed database systems. It refers to the process of deciding where to store data fragments across multiple sites to optimize performance, availability, and resource utilization. Effective data allocation strategies ensure that the distributed database system operates efficiently and meets the desired requirements.

*Data allocation is a vital process in distributed database management, involving the strategic placement of data across multiple sites to optimize performance, availability, and resource utilization. Effective data allocation strategies are essential for ensuring the efficient operation of distributed database systems while balancing factors like access patterns, network latency, and storage costs.*

---

## ■ Key Concepts -

- a) **Objective:** To distribute data in a manner that minimizes access time, balances the load, ensures availability, and reduces communication costs between sites.
- b) **Factors Influencing Data Allocation:**
  - **Access Patterns:** Understanding how frequently data is accessed and by which sites helps in placing data where it is most needed.
  - **Network Latency:** Reducing the time data takes to travel between sites by placing data closer to where it is frequently accessed.
  - **Storage Costs:** Balancing the costs associated with storing data at different sites.
  - **Processing Power:** Distributing data to leverage the processing capabilities of different sites.
  - **Redundancy Requirements:** Deciding how much replication is needed for fault tolerance and high availability.
- c) **Data Allocation Strategies:**
  - **Centralized Allocation:** All data is stored in a central location. This simplifies management but can lead to bottlenecks and high access times for remote sites.
  - **Partitioned Allocation:** Data is divided into fragments, and each fragment is stored at different sites based on access patterns and other factors.
  - **Replicated Allocation:** Copies of data are stored at multiple sites to ensure high availability and fault tolerance. This can lead to increased storage costs and complexity in maintaining consistency.
  - **Hybrid Allocation:** A combination of partitioning and replication to balance performance, availability, and cost considerations.

## ■ Advantages -

- **Improved Performance:** Reduces data access times by placing data closer to where it is needed, enhancing the overall system performance.
- **Load Balancing:** Distributes the workload evenly across multiple sites, preventing any single site from becoming a bottleneck.
- **Enhanced Availability:** Ensures that data is available even if some sites fail, through replication and redundancy.

- **Reduced Communication Costs:** Minimizes the need for data transfer between sites, reducing network traffic and associated costs.

#### ■ Disadvantages -

- **Complexity:** Deciding the optimal allocation strategy can be complex, requiring sophisticated algorithms and a deep understanding of access patterns and network topology.
- **Consistency:** Ensuring data consistency across multiple sites, especially in the presence of replicated data, can be challenging.
- **Dynamic Changes:** Adjusting the allocation strategy in response to changing access patterns and workloads can be resource-intensive.

#### ■ Query Processing in Distributed Databases - Query Processing in distributed databases involves managing and executing database queries across multiple distributed sites efficiently and effectively. The goal is to retrieve and manipulate data as if it were located in a single, centralized database, despite being physically distributed. Here are the key aspects and challenges involved:

#### ■ Key Concepts:

- **Query Decomposition:** The process of breaking down a high-level query into simpler sub-queries that can be executed at individual sites. Ensures that each sub-query is optimized for the specific data fragment it will access.
- **Query Optimization:** The process of finding the most efficient way to execute a query by considering multiple possible execution plans and choosing the one with the lowest cost in terms of resources and time. In a distributed context, this includes minimizing data transfer between sites, which can be a significant overhead.
- **Data Localization:** Translating a global query into local queries that operate on specific fragments of data stored at different sites. Ensures that the query processing is efficient and leverages local data access.
- **Join Processing:** Handling joins between tables that may be located at different sites. Techniques include semijoin, bloom join, and shipping whole tables to a single site to perform the join.
- **Query Coordination:** Managing the execution of the sub-queries and combining their results to produce the final output. Coordination is crucial for ensuring the correct sequence of operations and aggregation of results.

#### ■ Steps in Distributed Query Processing -

- **Query Parsing:** The high-level query is parsed to check its syntax and convert it into an internal representation, often a query tree or query graph.
- **Query Decomposition:** The parsed query is decomposed into smaller sub-queries that correspond to the fragments of data distributed across various sites. These sub-queries are optimized locally at each site.

- **Data Localization:** The sub-queries are translated into queries that can be executed on the specific data fragments located at different sites. Local optimization techniques are applied to minimize resource usage and response time.
- **Global Optimization:** The execution plans for the sub-queries are optimized globally by considering the cost of data transfer between sites.
- **Query Execution:** The optimized sub-queries are executed at the respective sites. The results are transferred back to the coordinating site, where they are combined to form the final result.
- **Result Assembly:** The partial results from different sites are assembled into a coherent final result and returned to the user.

#### ■ Challenges in Distributed Query Processing -

- **Network Latency and Bandwidth:** High network latency and limited bandwidth can significantly impact query performance.
- **Data Consistency:** Ensuring that data remains consistent across distributed sites during query execution, especially in the presence of concurrent updates.
- **Fault Tolerance:** Handling site failures and ensuring the system can recover without data loss or corruption.
- **Complexity:** The added complexity of coordinating and optimizing queries across multiple sites.

***NOTE: Query processing in distributed databases involves decomposing, optimizing, and coordinating queries to efficiently retrieve and manipulate data spread across multiple sites. This process aims to provide users with seamless access to distributed data while minimizing resource usage and response time. Effective query processing is essential for the performance and reliability of distributed database systems.***

=====