# Big Data Modeling (UCAML403 (F))

## UNIT I: INTRODUCTION

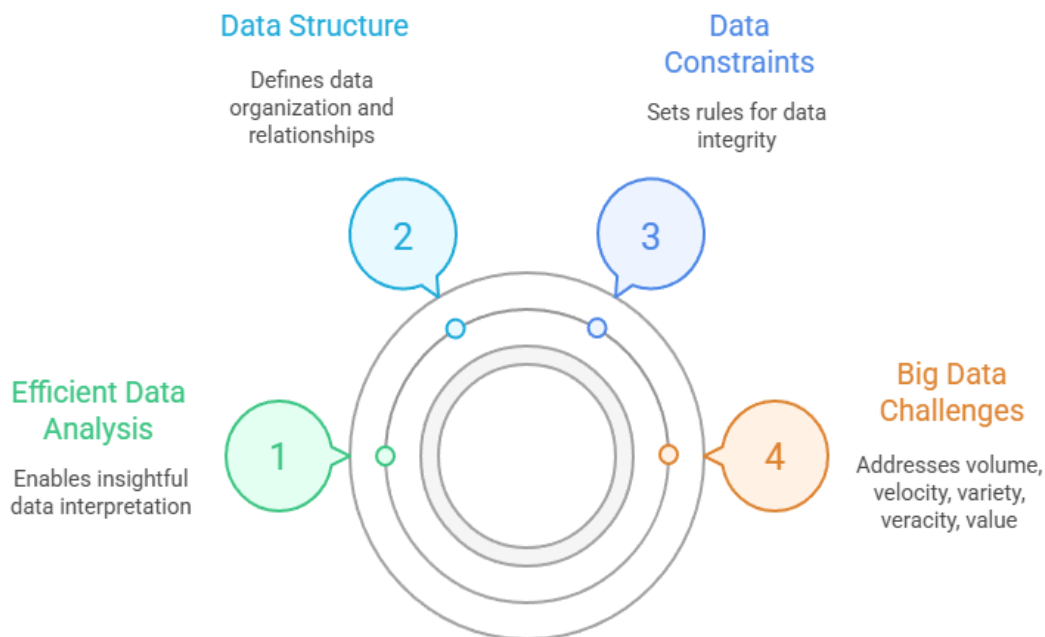| UNIT I: INTRODUCTION |
| --- |
| Understanding Big Data: Concepts and terminology, Big Data Characteristics, Different types of Data, Identifying Data Characteristics - Big Data Architecture - Big Data Storage: File system and Distributed File System, Sharding, Replication, Sharding and Replication, ACID and BASE Properties. |

Big data modeling is the process of designing a logical representation of data that is too large or complex for traditional data processing application software.

It involves defining the structure, relationships, and constraints of data to facilitate efficient storage, retrieval, and analysis.

Unlike traditional data modeling, big data modeling must address the challenges of volume, velocity, variety, veracity, and value (the 5 Vs of big data).



**Big Data Modeling Process**

Data Structure
Defines data organization and relationships
2

Data Constraints
Sets rules for data integrity
3

Efficient Data Analysis
Enables insightful data interpretation
1

Big Data Challenges
Addresses volume, velocity, variety, veracity, value
4

Prof. Shweta Lilhare

## Importance of Big Data Modeling

Effective big data modeling is crucial for several reasons:

- **Improved Data Understanding:** A well-defined data model provides a clear and concise representation of the data, making it easier for users to understand its structure and meaning.

- **Enhanced Data Quality:** By defining data types, constraints, and relationships, data modeling helps to ensure data consistency and accuracy.

- **Optimized Query Performance:** A properly designed data model can significantly improve query performance by enabling efficient data access and retrieval.

- **Simplified Data Integration:** Data modeling facilitates the integration of data from multiple sources by providing a common framework for understanding and transforming data.

- **Better Decision Making:** By providing a clear and accurate view of the data, data modeling supports better decision-making based on data-driven insights.

## Challenges of Big Data Modeling

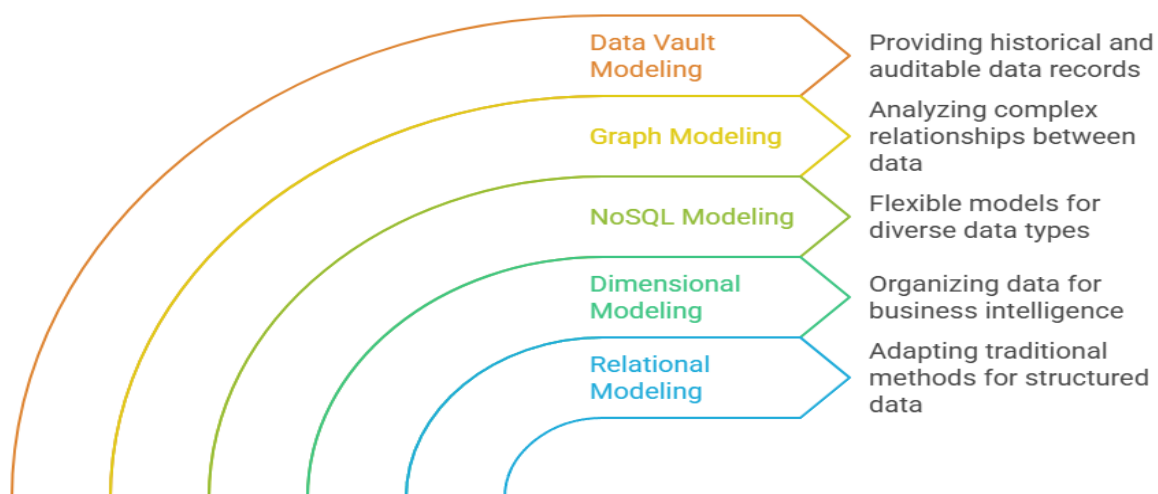Big data modeling presents several unique challenges:

- **Volume:** The sheer volume of data can make it difficult to design a data model that can efficiently store and process the data.

- **Velocity:** The high velocity of data streams requires data models that can handle real-time or near real-time data ingestion and processing.

- **Variety:** The variety of data types and formats requires data models that can accommodate structured, semi-structured, and unstructured data.

- **Veracity:** The uncertainty and inconsistency of data quality requires data models that can handle missing, incomplete, or inaccurate data.

- **Scalability:** The data model must be able to scale to accommodate future growth in data volume and complexity.

- **Evolving Data:** Big data is often dynamic and evolving, requiring flexible data models that can adapt to changing data structures and requirements.

- **Complexity:** The complexity of relationships between data elements can make it difficult to design a data model that accurately reflects the real-world phenomena being modeled.

Prof. Shweta Lilhare

## Methodologies for Big Data Modeling

Several methodologies can be used for big data modeling, each with its own strengths and weaknesses:

- **Relational Modeling:** Traditional relational modeling techniques can be adapted for big data, but they may not be suitable for all types of big data. Relational models are well-suited for structured data with well-defined relationships, but they can struggle with unstructured or semi-structured data.

- **Dimensional Modeling:** Dimensional modeling is a technique that is commonly used for data warehousing and business intelligence. It involves organizing data into facts and dimensions, which can be used to analyze data from different perspectives.

- **NoSQL Modeling:** NoSQL databases offer a variety of data models, including document, key-value, column-family, and graph models. These models are often more flexible and scalable than relational models, making them well-suited for big data applications.

- **Graph Modeling:** Graph databases are designed to store and analyze data as a network of nodes and edges. They are particularly well-suited for applications that involve complex relationships between data elements, such as social networks, knowledge graphs, and recommendation systems.

- **Data Vault Modeling:** Data Vault is a data modeling technique designed to provide a historical and auditable record of data. It is often used in data warehousing environments to track changes to data over time.

### Overview of Big Data Modeling Techniques

| | |
|---|---|
| Data Vault Modeling | Providing historical and auditable data records |
| Graph Modeling | Analyzing complex relationships between data |
| NoSQL Modeling | Flexible models for diverse data types |
| Dimensional Modeling | Organizing data for business intelligence |
| Relational Modeling | Adapting traditional methods for structured data |

- **Big Data:**

Refers to datasets that are too large, fast-moving, or varied for traditional data processing systems to handle effectively.

- **Data Modeling:**

The process of creating a visual representation of data structures, relationships, and rules. It provides a blueprint for how data is stored, accessed, and used.

- **Conceptual Data Model:**

The initial, high-level view of the data, focusing on entities, their attributes, and relationships, often represented using Entity-Relationship diagrams (ERDs).

- **Logical Data Model:**

A more detailed model that defines data types, relationships, and constraints, but still independent of any specific database technology.

- **Physical Data Model:**

The specific implementation of the data model within a particular database system, including tables, columns, and indexes.

**What are the Five "Vs" of Big Data?**

**Volume**

Volume refers to the sheer amount of data being generated and collected. This is perhaps the most commonly understood aspect of Big Data. We're talking about data volumes that are far beyond the capacity of traditional database systems to handle efficiently.

- **Scale:** Big Data deals with massive datasets, often measured in terabytes (TB), petabytes (PB), and even exabytes (EB).

- **Sources:** This data comes from a multitude of sources, including social media, sensors, transaction records, web logs, and many more.

- **Challenges:** Storing, processing, and analyzing such large volumes of data present significant challenges in terms of infrastructure, storage solutions, and processing power.

- **Examples:** Consider the data generated by Facebook users daily (posts, likes, shares), or the data collected by sensors in a smart city.

**Velocity**

Velocity refers to the speed at which data is generated and processed. It's not just about the amount of data, but also how quickly it's arriving and how quickly we need to process it.

- **Speed of Generation:** Data is being generated at an unprecedented rate, often in real-time or near real-time.

- **Processing Requirements:** Many applications require immediate analysis and action based on incoming data streams.

- **Streaming Data:** This often involves dealing with streaming data, which requires specialized processing techniques.

- **Examples:** Think about stock market data, where decisions need to be made in milliseconds, or real-time traffic monitoring systems.

**Variety**

Variety refers to the different types and formats of data. Big Data is not just structured data (like data in a relational database), but also includes unstructured and semi-structured data.

- **Structured Data:** This is data that fits neatly into a predefined format, like a database table.

- **Unstructured Data:** This includes data like text documents, images, audio, and video, which don't have a predefined format.

- **Semi-structured Data:** This includes data like XML, JSON, and log files, which have some organizational properties but are not as rigid as structured data.

- **Challenges:** Dealing with variety requires tools and techniques that can handle different data formats and integrate them for analysis.

- **Examples:** A company might need to analyze customer reviews (unstructured text), sales data (structured data), and website clickstream data (semi-structured data) to get a complete picture of customer behavior.

**Veracity**

Veracity refers to the accuracy and reliability of the data. With so much data coming from so many sources, it's important to ensure that the data is trustworthy.

- **Data Quality:** Big Data often includes noisy, incomplete, and inconsistent data.

- **Data Cleansing:** Data cleansing and validation are crucial steps in ensuring data quality.

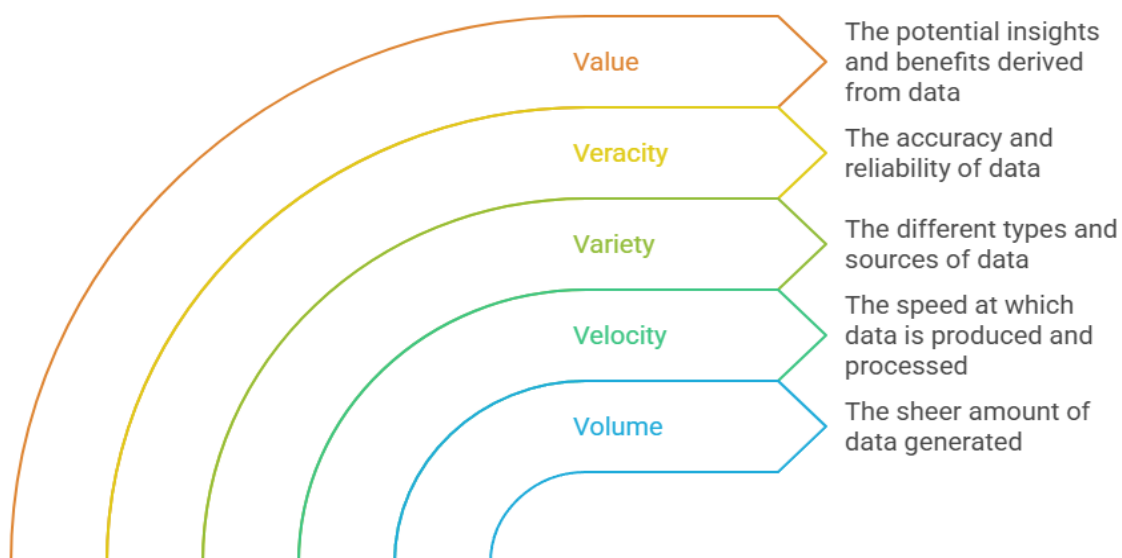- **Source Reliability:** The reliability of the data source also needs to be considered.

- **Impact on Analysis:** Inaccurate data can lead to incorrect conclusions and poor decision-making.

- **Examples:** Social media data can be full of misinformation or biased opinions. Sensor data can be affected by environmental factors or equipment malfunctions.

## Value

Value refers to the ability to turn Big Data into actionable insights and business value. This is the ultimate goal of Big Data initiatives.

- **Business Insights:** The real value of Big Data lies in its ability to uncover hidden patterns, trends, and relationships that can inform business decisions.

- **Decision Making:** These insights can be used to improve operational efficiency, enhance customer experience, develop new products and services, and gain a competitive advantage.

- **Return on Investment:** Organizations need to focus on extracting value from their Big Data investments.

- **Examples:** A retailer might use Big Data to personalize product recommendations, optimize pricing, or predict future demand. A healthcare provider might use Big Data to improve patient outcomes or reduce costs.

## Understanding the Five "Vs" of Big Data

Value — The potential insights and benefits derived from data

Veracity — The accuracy and reliability of data

Variety — The different types and sources of data

Velocity — The speed at which data is produced and processed

Volume — The sheer amount of data generated

- **Conceptual Data Model:**

    - This model provides a high-level, business-oriented view of the data.

    - It focuses on identifying key entities, their attributes, and the relationships between them, without delving into technical implementation details.

    - Its purpose is to establish a shared understanding of the data requirements among business stakeholders and technical teams.

    - It is technology-independent.

- **Logical Data Model:**

    - This model builds upon the conceptual model by adding more detail, defining specific attributes for each entity, and establishing primary and foreign keys to represent relationships.

    - It describes what data is needed and how it is organized, but still remains independent of a specific database management system (DBMS).

    - It provides a blueprint for database design and helps in understanding data relationships and constraints.

- **Physical Data Model:**

    - This model represents the actual implementation of the data in a specific DBMS.

    - It includes technical details such as table names, column names, data types, constraints (e.g., primary keys, foreign keys, unique constraints), indexes, and storage considerations.

## Key data characteristics identified in a data model include:

- **Entities:**

The main objects or concepts about which data is stored (e.g., Customer, Product, Order).

- **Attributes:**

The properties or characteristics that describe each entity (e.g., for Customer: CustomerID, FirstName, LastName, Address).

- **Data Types:**

The specific type of data an attribute can hold, dictating its format and permissible values (e.g., text, integer, decimal, date, boolean).

- **Relationships:**

The associations or connections between different entities, indicating how they interact (e.g., a Customer places an Order).

- **Cardinality:** Defines the number of instances of one entity related to instances of another (e.g., one-to-one, one-to-many, many-to-many).

- **Nullability:** Specifies whether an attribute can contain a null value (i.e., no value).

- **Keys:**

Attributes or sets of attributes that uniquely identify records within an entity and establish relationships:

- **Primary Key (PK):** Uniquely identifies each record within a table.

- **Foreign Key (FK):** Establishes a link between two tables by referencing the primary key of another table.

- **Constraints:**

Rules or limitations imposed on the data to maintain its integrity and consistency (e.g., uniqueness constraints, check constraints, default values).

- **Domains:**

The set of all possible values that an attribute can take, often defining a specific range or enumeration.

- **Metadata:**

Data about the data, including definitions, sources, ownership, and quality metrics, which provides context and aids understanding.

## What is Big Data Architecture?

Big data architecture is specifically designed to manage data ingestion, data processing, and analysis of data that is too large or complex. A big size data cannot be store, process and manage by conventional relational databases. The solution is to organize technology into a structure of big data architecture. Big data architecture is able to manage and process data.
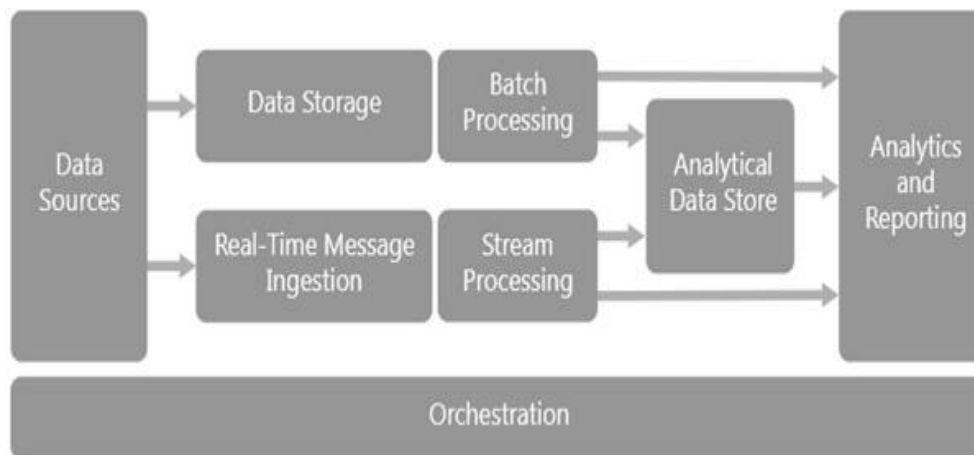
Key Aspects of Big Data Architecture

The following are some key aspects of big data architecture −

- To store and process large size data like 100 GB in size.

- To aggregates and transform of a wide variety of unstructured data for analysis and reporting.

- Access, processing and analysis of streamed data in real time.



- **Data sources:** These are the various places big data sets pull from and then ingest into the big data architecture framework. Data sources include social media, websites, and the internet.

- **Data storage:** Data storage refers to the software framework capable of holding massive amounts of either structured or unstructured data. Storing large volumes of structured or unstructured data is frequently labeled as a data lake.

- **Batch processing:** This is the software used to sort and convert the big data sets into usable files that are ready for analysis.

- **Real-time message ingestion:** This part of big data architecture places the data into groups, allowing an easier shift into the subsequent stages of the storage process. Essentially, this procedure permits the ingestion and storage of streaming data, which needs real-time processing.

- **Stream processing:** This process takes the real-time messages, filters them, and prepares them for analysis by adding the processed messages to an output sink.

- **Analytical datastore:** This is where the processed and cleaned data becomes accessible for use by different analytical techniques and tools.

## File Systems

A file system is a method of organizing and storing data on a storage device, such as a hard drive or solid-state drive. It provides a hierarchical structure of directories and files, allowing users and applications to easily access and manage data. Traditional file systems, like NTFS, ext4, and HFS+, are designed for single-machine environments and are optimized for performance and reliability within those constraints.

**Key Characteristics of Traditional File Systems:**

- **Hierarchical Structure:** Data is organized into directories and subdirectories, creating a tree-like structure.

- **Metadata Management:** File systems store metadata about files, such as name, size, creation date, and permissions.

- **Access Control:** File systems provide mechanisms for controlling access to files and directories, ensuring data security.

- **Data Integrity:** File systems implement techniques to ensure data integrity, such as checksums and journaling.

**Limitations for Big Data:**

Traditional file systems are not well-suited for big data storage due to their limitations in scalability, fault tolerance, and parallel processing..

## Distributed File Systems (DFS)

A distributed file system (DFS) is a file system that spans multiple machines, allowing data to be stored and accessed across a network. DFSs are designed to address the limitations of traditional file systems in the context of big data by providing scalability, fault tolerance, and parallel processing capabilities.

**Key Characteristics of Distributed File Systems:**

- **Scalability:** DFSs can scale to store petabytes or even exabytes of data by adding more machines to the cluster.

- **Fault Tolerance:** DFSs replicate data across multiple machines, ensuring that data is available even if some machines fail.

- **Parallel Processing:** DFSs allow data to be processed in parallel across multiple machines, significantly improving performance.

- **Data Locality:** DFSs strive to store data close to the processing nodes, minimizing network traffic and improving performance.

- **Data Replication:** DFSs replicate data across multiple nodes to ensure fault tolerance and high availability.

- **Data Partitioning:** DFSs divide data into smaller blocks and distribute them across the cluster.
- **Metadata Management:** DFSs maintain metadata about files and blocks, such as location and ownership.

## Common Distributed File Systems:

- **Hadoop Distributed File System (HDFS):** HDFS is a widely used DFS designed for storing and processing large datasets in Hadoop clusters. It is known for its scalability, fault tolerance, and support for batch processing.
- **GlusterFS:** GlusterFS is an open-source DFS that provides a scalable and distributed storage solution. It supports various storage configurations and is suitable for a wide range of applications.
- **Ceph:** Ceph is a distributed object storage system that provides a unified storage platform for block storage, object storage, and file storage. It is known for its scalability, reliability, and performance.
- **Lustre:** Lustre is a high-performance DFS designed for large-scale scientific computing and high-performance computing (HPC) environments.

## Hadoop Distributed File System (HDFS) in Detail:

HDFS is a core component of the Hadoop ecosystem and is designed to store and manage large datasets reliably and efficiently. It follows a master-slave architecture, with a NameNode managing the file system metadata and DataNodes storing the actual data blocks.

- **Architecture:** HDFS consists of a NameNode and multiple DataNodes. The NameNode stores the file system metadata, such as the directory structure, file names, and block locations. The DataNodes store the actual data blocks.
- **Data Replication:** HDFS replicates data blocks across multiple DataNodes to ensure fault tolerance. The default replication factor is 3, meaning that each data block is stored on three different DataNodes.
- **Data Locality:** HDFS strives to store data blocks close to the processing nodes, minimizing network traffic and improving performance. When a MapReduce job is submitted, the JobTracker attempts to schedule tasks on the DataNodes that contain the input data.
- **Fault Tolerance:** HDFS is designed to be fault-tolerant. If a DataNode fails, the NameNode detects the failure and instructs other DataNodes to replicate the missing data blocks.
- **Scalability:** HDFS can scale to store petabytes or even exabytes of data by adding more DataNodes to the cluster. The NameNode can also be scaled horizontally using techniques such as federation.

## Sharding

Sharding, also known as horizontal partitioning, involves dividing a large database into smaller, more manageable pieces called shards.

Each shard contains a subset of the overall data and resides on a separate database server. This distribution allows for parallel processing and reduces the load on any single server.

**Benefits**

- **Improved Performance:** By distributing data across multiple servers, queries can be executed in parallel, leading to faster response times and increased throughput.

- **Enhanced Scalability:** Sharding allows you to scale your database horizontally by adding more shards as your data grows. This avoids the limitations of vertical scaling (increasing the resources of a single server).

- **Reduced Downtime:** If one shard goes down, only a subset of the data is affected, and the rest of the database remains operational.

- **Geographic Distribution:** Shards can be located in different geographic regions to improve performance for users in those regions and comply with data sovereignty regulations.

**Drawbacks**

- **Increased Complexity:** Implementing and managing a sharded database is more complex than managing a single database.

- **Data Distribution Challenges:** Choosing the right sharding key (the attribute used to determine which shard a piece of data belongs to) is crucial for even data distribution and optimal performance. Poorly chosen sharding keys can lead to hotspots, where some shards are overloaded while others are underutilized.

- **Cross-Shard Queries:** Queries that require data from multiple shards can be more complex and slower than queries that can be satisfied from a single shard.

- **Data Consistency:** Maintaining data consistency across shards can be challenging, especially in distributed environments.

**Sharding Strategies**

Several sharding strategies can be employed, each with its own trade-offs:

- **Range-Based Sharding:** Data is divided into shards based on a range of values for the sharding key. For example, users with IDs between 1 and 1000 might be stored in shard 1, users with IDs between 1001 and 2000 in shard 2, and so on. This strategy is simple to implement but can lead to hotspots if data is not evenly distributed across the ranges.

- **Hash-Based Sharding:** A hash function is used to map the sharding key to a shard. This strategy generally provides better data distribution than range-based sharding but can make range queries more difficult.

- **Directory-Based Sharding:** A lookup table or directory is used to map sharding keys to shards. This strategy provides flexibility but adds an extra layer of indirection.

- **Dynamic Sharding:** Shards are dynamically created and rebalanced as data grows or shrinks. This strategy requires more sophisticated management but can provide optimal performance and scalability.

  **Use Cases**

- **E-commerce platforms:** Sharding can be used to distribute product catalogs, user accounts, and order data across multiple servers.

- **Social media networks:** Sharding can be used to distribute user profiles, posts, and connections across multiple servers.

- **Gaming platforms:** Sharding can be used to distribute game data, player profiles, and leaderboard information across multiple servers.

- **Financial institutions:** Sharding can be used to distribute transaction data, account information, and customer data across multiple servers.


## Replication

Replication involves creating multiple copies of a database and storing them on different servers. These copies, called replicas, are kept synchronized with each other. Replication provides redundancy and improves availability and read performance.

**Benefits**

- **High Availability:** If one replica fails, the other replicas can continue to serve requests, ensuring minimal downtime.

- **Improved Read Performance:** Read requests can be distributed across multiple replicas, reducing the load on any single server and improving response times.

- **Disaster Recovery:** Replicas can be located in different geographic regions to provide disaster recovery capabilities.

## Drawbacks

- **Increased Storage Costs:** Replication requires storing multiple copies of the data, which increases storage costs.

- **Data Consistency Challenges:** Maintaining data consistency across replicas can be challenging, especially in asynchronous replication scenarios.

- **Write Performance Overhead:** Write operations must be propagated to all replicas, which can increase write latency.

**Replication Types**

Several replication types exist, each with its own trade-offs:

- **Master-Slave Replication:** One replica is designated as the master, and all write operations are directed to the master. The master then replicates the changes to the other replicas, which are called slaves. This is a simple and widely used replication type, but it can suffer from write contention on the master.

- **Master-Master Replication:** All replicas can accept write operations. Changes made to one replica are propagated to all other replicas. This replication type provides higher availability and write performance but is more complex to manage and can lead to data conflicts.

- **Multi-Master Replication:** A variation of master-master replication where multiple masters exist, each responsible for a subset of the data. This can improve write performance and reduce the risk of data conflicts.

- **Snapshot Replication:** A snapshot of the data is taken at a specific point in time and replicated to other servers. This replication type is suitable for read-only data or data that changes infrequently.

- **Transactional Replication:** Changes are replicated as transactions, ensuring data consistency across replicas. This replication type is more complex but provides the highest level of data consistency.

## ACID

ACID is a set of properties that guarantee database transactions are processed reliably.

It's a cornerstone of traditional relational database management systems (RDBMS).

- **Atomicity:** A transaction is treated as a single, indivisible unit of work. Either all changes within the transaction are applied successfully, or none are. If any part of the transaction fails, the entire transaction is rolled back to its initial state, ensuring data integrity. Think of it as an "all or nothing" principle.

- **Consistency:** A transaction must maintain the integrity of the database. It ensures that the database transitions from one valid state to another. This means that the transaction must adhere to all defined rules, constraints, and integrity checks. If a transaction violates any of these rules, it's rolled back, preventing the database from entering an invalid state.

- **Isolation:** Transactions are isolated from each other, meaning that concurrent transactions do not interfere with each other's execution. Each transaction operates as if it were the only transaction running on the system. This prevents issues like dirty reads, non-repeatable

reads, and phantom reads, ensuring data accuracy and predictability. Isolation levels define the degree to which transactions are isolated from each other.

- **Durability:** Once a transaction is committed, the changes are permanent and will survive even system failures such as power outages or crashes. The database ensures that the committed data is stored persistently and can be recovered in case of a failure. This is typically achieved through techniques like write-ahead logging and data replication.
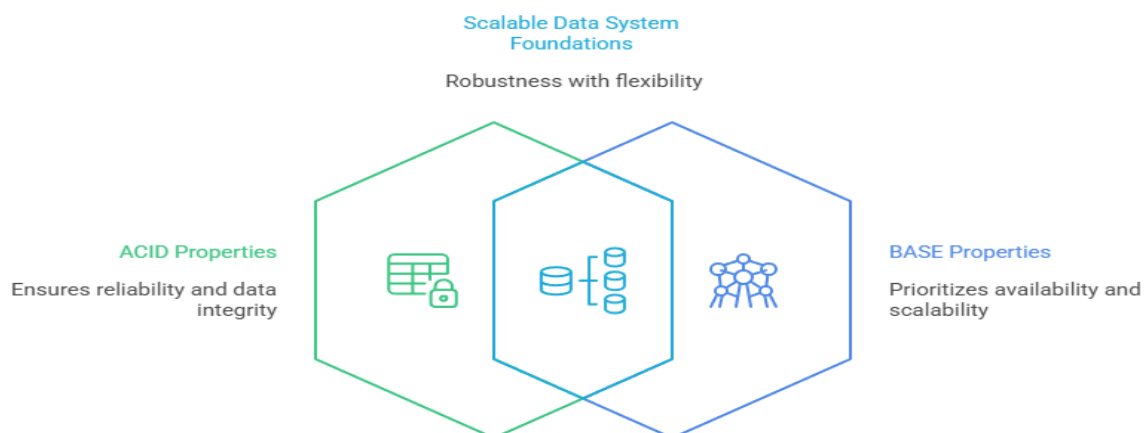
**Advantages of ACID:**

- **Data Integrity:** ACID properties guarantee a high level of data integrity, making them suitable for applications where data accuracy is paramount, such as financial systems and healthcare records.

- **Reliability:** The atomicity and durability properties ensure that transactions are processed reliably, even in the face of system failures.

- **Predictability:** The isolation property ensures that concurrent transactions do not interfere with each other, making the behavior of the system predictable.

**Disadvantages of ACID:**

- **Scalability:** Achieving full ACID compliance can be challenging in distributed systems, as it often requires complex coordination and synchronization mechanisms. This can limit the scalability and performance of the system.

- **Performance:** The strict consistency and isolation requirements of ACID can lead to performance overhead, especially in high-concurrency environments.

- **Complexity:** Implementing and managing ACID-compliant systems can be complex, requiring specialized expertise and tools.



**Comparing ACID and BASE Data Model Properties**

Scalable Data System Foundations

Robustness with flexibility

ACID Properties

Ensures reliability and data integrity

BASE Properties

Prioritizes availability and scalability

BASE is an alternative set of properties that prioritize availability and scalability over strict consistency. It's often used in NoSQL databases and distributed systems where high availability and performance are more important than immediate consistency. Let's examine each property:

- **Basically Available:** The system should be available most of the time. This means that the system should respond to requests even if some nodes are unavailable or experiencing network issues. Availability is prioritized over consistency.

- **Soft State:** The state of the system can change over time, even without any input. This means that the data may not be immediately consistent across all nodes. The system relies on eventual consistency to propagate updates and resolve conflicts.

- **Eventually Consistent:** The system will eventually become consistent, given enough time and no further updates. This means that there may be a period of inconsistency after an update, but the system will eventually converge to a consistent state. The time it takes for the system to become consistent is known as the consistency window.

**Advantages of BASE:**

- **Scalability:** BASE properties allow for greater scalability in distributed systems, as they relax the strict consistency requirements of ACID. This allows the system to handle a larger volume of data and traffic.

- **Availability:** BASE prioritizes availability, ensuring that the system remains responsive even in the face of failures.

- **Performance:** The relaxed consistency requirements of BASE can lead to improved performance, especially in high-concurrency environments.

**Disadvantages of BASE:**

- **Data Inconsistency:** The eventual consistency model can lead to data inconsistency, which can be problematic for applications that require strict data accuracy.

- **Complexity:** Dealing with eventual consistency requires careful consideration of conflict resolution strategies and data reconciliation mechanisms.

- **Unpredictability:** The eventual consistency model can make the behavior of the system less predictable, as the state of the data may not be immediately apparent.