

Introduction to Relational Databases in Cloud Computing

Recap of Previous Lecture

- Discussed **Relational Databases (RDBMS)**.
 - Defined how RDBMS works in traditional systems.
-

Relational Databases in Cloud

- Many databases are now hosted on the cloud:
 - Examples: Google Cloud SQL, AWS RDS, Azure SQL Database, IBM Cloud DB, etc.
 - Before diving into cloud databases, it's important to understand how **traditional RDBMS** functioned.
-

Traditional RDBMS Overview

- A **Relational Database Management System (RDBMS)** is software used to **store, organize, and manage data**.
- Stores data in **tabular form** using **rows and columns**.
- Interactions happen through **SQL (Structured Query Language)**.
- Components:
 - Application/program interacts with RDBMS.
 - Disk space is managed using **disk blocks** and **memory buffers**.

- As data grows (bulk data), accessibility can slow down.
-

Functioning of RDBMS

- Data stored in a **structured manner** using **disk pages**.
 - Uses **indexing** for efficient access.
 - Index helps locate where a record is stored.
 - Tables relate to each other through **predefined relationships**.
 - **Operations like Create, Select, Insert, Delete** are performed via SQL.
-

Memory Management and Data Storage

- Cloud RDBMS platforms use their own disk management systems.
 - Data stored in:
 - **Disk pages** (continuous pages)
 - **Buffer memory** fetches data via **page requests**.
 - **Page Replacement Policy** (example: in Google BigTable) ensures optimal data retrieval.
-

Indexing in Databases

- Index = **Data Pointer** (सूचक).
- Helps in efficient **search and access** of records.

- Types of indexing:
 - **Primary Index**
 - **B+ Tree Index**
 - **Bitmap Index**
 - **Hash Index**
-

Row vs Column Orientation

Aspect	Row-Oriented	Column-Oriented
Structure	Rows (Tuples)	Columns (Attributes)
Storage	Organized by rows	Organized by columns
Index	Works with primary keys and unique values	Better for analytical queries
Usage	Good for OLTP (Transactional systems)	Good for OLAP (Analytical systems)

SQL (Structured Query Language)

- Language used to **interact with RDBMS**.
- Helps:
 - Collect data
 - Perform queries
 - Create new records
 - Delete or update existing records
- Example operations:
 - `SELECT * FROM users;`

- `INSERT INTO products VALUES (...);`
 - `DELETE FROM orders WHERE id=...;`
-

Summary

- RDBMS are essential to both traditional and cloud-based storage.
- Efficient data management relies on:
 - **Disk storage**
 - **Indexing**
 - **Memory buffers**
 - **SQL queries**
- Cloud platforms manage their own database systems but follow these fundamental principles.

Introduction to SQL

- SQL stands for **Structured Query Language**.
 - It is used for managing and accessing **data in a database**.
 - Developed by **IBM in the 1970s** for data collection and maintenance.
 - Acts as a **primary interface** for working with **Relational Databases (RDBMS)**.
-

Purpose and Functionality of SQL

- SQL is used to:

- **Select** data
 - **Update** data
 - **Create** tables
 - **Delete** records
 - It supports **data retrieval** and **analytical operations**.
 - SQL organizes data in **rows and columns** within **tables**.
-

Understanding Tables

- Tables consist of:
 - **Rows** (each representing a record)
 - **Columns** (each representing an attribute/field)
 - Example:
 - A **Customer** table with columns like **Customer_ID**, **Customer_Name**, **Product_Purchased**.
 - Each record (row) has a **unique identifier**, often defined using a **Primary Key**.
-

Keys and Constraints

- **Primary Key**:
 - Ensures each row in a table is unique.
- **Constraints** in SQL ensure data integrity:
 - **NOT NULL**

- **DEFAULT**
 - **CHECK**
 - **UNIQUE**
 - **FOREIGN KEY**
- Constraints ensure that the data is **accurate, consistent, and reliable**.
-

SQL Queries

- Common SQL operations include:
 - **SELECT**: to retrieve data
 - **INSERT**: to add new data
 - **UPDATE**: to modify existing data
 - **DELETE**: to remove data
 - **CREATE**: to create tables or databases
 - These help manage and manipulate the data within the database.
-

Data Integrity

- Refers to the **completeness and accuracy** of the data.
 - Ensured through constraints and validation rules applied to tables and columns.
-

Transactions in SQL

- A **Transaction** is a single unit of work that must be executed completely or not at all.
 - Examples: **COMMIT**, **ROLLBACK**
 - Transactions ensure **data consistency and reliability**.
-

ACID Properties

- SQL supports **ACID** compliance for transactions:
 - **Atomicity** – All steps succeed or fail together.
 - **Consistency** – Database remains in a valid state.
 - **Isolation** – Transactions don't interfere with each other.
 - **Durability** – Once committed, data is saved permanently.
-

Why SQL is Important

- Essential for understanding **Data Analytics**, **AI**, and **Machine Learning**.
- SQL is often required to:
 - Query backend databases.
 - Understand data patterns.
 - Build applications involving data.

detailed comparison between **GFS (Google File System)** and **HDFS (Hadoop Distributed File System)** in **column format** with deeper explanations:

Property	GFS (Google File System)	HDFS (Hadoop Distributed File System)
----------	--------------------------	---------------------------------------

Design Goals	<ul style="list-style-type: none"> - Designed to support very large files - Optimized for batch processing rather than user interaction - Focus on high throughput and fault tolerance - Built for Google's internal use 	<ul style="list-style-type: none"> - Also supports large files and batch processing - Optimized for data-intensive computing tasks like analytics - Designed for reliable storage, even with frequent hardware failures
Failure Tolerance	<ul style="list-style-type: none"> - Fault tolerance through replication across chunk servers - Can handle failures in master, chunk servers, or network 	<ul style="list-style-type: none"> - Uses block replication (default 3 copies) - Can handle failures in namenode, datanode, or network
File Management	<ul style="list-style-type: none"> - File system structured with flat namespace - Files organized using path names - Primarily for Google-only use 	<ul style="list-style-type: none"> - Follows traditional hierarchical structure (folders, directories) - Designed to work with third-party systems (e.g., cloud tools)
Architecture	<ul style="list-style-type: none"> - Components: Master node + Chunk Servers - Master holds metadata; chunk servers store file chunks 	<ul style="list-style-type: none"> - Components: Namenode (metadata) + Datanodes (actual data blocks) - Namenode is central control unit
Security	<ul style="list-style-type: none"> - GFS operates with restricted access - Allows access only to authorized employees and vendors - Data centers are undisclosed and protected for redundancy 	<ul style="list-style-type: none"> - Based on POSIX-style permissions - Allows simple file access controls using users and groups - Less strict than GFS
Database File Support	<ul style="list-style-type: none"> - Integrated with Bigtable, a distributed storage system for structured data 	<ul style="list-style-type: none"> - Works with HBase, an open-source implementation of Bigtable over HDFS
File Chunking	<ul style="list-style-type: none"> - Files split into fixed-size chunks of 64 MB - Chunks are compressed and distributed 	<ul style="list-style-type: none"> - Files split into large blocks (default 64 MB or 128 MB), configurable - Stored across datanodes for efficiency

Cache Management	<ul style="list-style-type: none"> - Clients cache metadata (e.g., file locations) locally - No client or server-level file data caching 	<ul style="list-style-type: none"> - Uses a Distributed Cache system - Useful for MapReduce jobs with large read-only files like text, JARs
File Storage	<ul style="list-style-type: none"> - File chunks stored as local files in the Linux file system 	<ul style="list-style-type: none"> - Distributed cache supports private and public modes for shared or isolated access
Data Access Model	<ul style="list-style-type: none"> - Read/Write via direct communication with chunk servers 	<ul style="list-style-type: none"> - Clients contact namenode for metadata, then datanodes for actual data
Communication Protocol	<ul style="list-style-type: none"> - Uses TCP connections for all communication, including file transfers 	<ul style="list-style-type: none"> - Uses RPC (Remote Procedure Call) protocol over TCP/IP for communication between nodes
Integration with Tools	<ul style="list-style-type: none"> - Tightly coupled with Google's internal systems like MapReduce and Bigtable 	<ul style="list-style-type: none"> - Designed for integration with Hadoop ecosystem tools like MapReduce, Hive, Pig, Spark
Flexibility	<ul style="list-style-type: none"> - Not open source - Internal to Google systems only 	<ul style="list-style-type: none"> - Open source (Apache Hadoop) - Widely adopted in industry and research

Bigtable Overview

What is Bigtable?

- Bigtable is a **distributed structured storage system** developed by **Google**.
 - It is built on top of **Google File System (GFS)**.
 - Bigtable is designed to handle **sparse, distributed, persistent, and multi-dimensional sorted maps**.
-

Structure and Storage

- **Data Structure:**
 - Bigtable stores data in **multi-dimensional sorted maps**.
 - It maps data using:
 - **Row key**
 - **Column key**
 - **Timestamp**
 - **Rows & Columns:**
 - A **row** is uniquely identified by a **row key**.
 - Each row can have **multiple column families**.
 - Each **column family** contains **column qualifiers** and **values** (name-value pairs).
 - Columns in a column family are typically stored together on disk.
 - **Timestamps:**
 - Each cell (intersection of a row and column) can store **multiple versions of data**, identified by **timestamps**.
 - These versions are stored in **decreasing order** of time.
-

Integration and Ecosystem

- **Client Libraries:**
 - Bigtable can be accessed using various **client libraries**, especially in **Java** using **Apache HBase libraries**.
- **Compatibility:**

- Bigtable is **compatible with the open-source Apache ecosystem**.
 - Works smoothly with **Apache HBase API**.
-

Components of Bigtable

- **Row Key:**
 - Uniquely identifies a row.
 - Used for locating and organizing data.
 - **Column Family:**
 - A group of related columns.
 - All columns within a family are stored together.
 - Example: Column family = “Product”, columns = “Detergent”, “Soap”, etc.
 - **Cell:**
 - Intersection of row key and column.
 - Stores the actual **value**, possibly with **multiple versions**.
 - **Data Storage:**
 - Stored using a **column-oriented database model**.
 - Supports **versioned data storage** using timestamps.
-

Usage and Purpose

- **Efficient storage** of large-scale structured data.
- Can **scale horizontally** across many machines.
- Used in systems where:

- Massive data needs to be stored and quickly accessed.
 - Time-based versions of data are important.
 - Data is **sparse** and **heterogeneous**.
-

Key Points

- Designed for **high scalability, performance, and flexibility**.
- Rows are kept in **lexicographic order** by row key.
- Ideal for **real-time data processing, analytics, and data indexing**.
- **Column families** must be declared at table creation and are **schema-defined**.
- Each cell can hold **multiple versions** based on time, useful for tracking changes.

Introduction to HBase

- **HBase** is an **open-source, NoSQL, distributed big data** storage system.
- It is designed for **random, real-time read/write access** to data at the **petabyte scale**.
- It is highly efficient in handling **large, sparse datasets**.
- HBase is:
 - **Column-oriented**
 - **Horizontally scalable**
 - Operates as a **set of tables** storing data in **key-value format**.

Features of HBase

- Stores data in a structured way with **similar types grouped together**.
 - Provides **APIs** for application development in multiple **programming languages**.
 - Part of the **Hadoop ecosystem**.
 - Offers **random, real-time read/write access** to data stored in the **HDFS** (Hadoop Distributed File System).
 - **Traditional RDBMS** systems become **slower** with large datasets, while HBase handles them efficiently.
 - HBase can manage **null values** effectively, reducing **overhead** in sparse datasets.
 - Does not require **schema change downtime** like RDBMS.
-

HBase vs RDBMS

- RDBMS:
 - Suited for **highly structured data** with **well-defined schemas**.
 - Slows down with **very large datasets**.
 - Schema changes may require **downtime**.
 - HBase:
 - Optimized for **sparse, distributed data**.
 - Handles **schema flexibility** better.
 - Performs well with **real-time access** needs.
-

HBase Architecture

HBase architecture is composed of **three main components**:

1. **HMaster**
 2. **Region Servers**
 3. **ZooKeeper**
-

1. HMaster

- Acts as the **main server** in HBase.
 - Interacts with multiple **Region Servers**.
 - Responsible for:
 - Assigning regions to region servers.
 - Managing **DDL operations** (e.g., CREATE, DELETE TABLE).
 - **Monitoring** all region servers.
 - Running **background tasks** in a distributed environment.
 - Handling features like:
 - **Controlling**
 - **Load balancing**
 - **Failure recovery**
-

2. Region Server

- Stores **data in rows**, which are horizontally divided into **regions**.
- Regions are the **building blocks** of HBase clusters.
- Each region:

- Contains **column families** with distributed tables.
 - Region servers **run on HDFS DataNodes** in a Hadoop cluster.
 - Responsible for:
 - Handling **read/write requests**
 - **Managing and executing** data operations
 - Default block size: **256 MB**
-

3. ZooKeeper

- Acts as the **coordinator** in the HBase ecosystem.
- Provides services such as:
 - Maintaining **configuration**
 - **Synchronizing** information
 - Coordinating the **interaction** between client and servers

Advantages of HBase

- Handles **large and sparse datasets** very efficiently.
 - **Real-time access** to data (read/write).
 - **Scalable horizontally**, i.e., add more machines as data grows.
 - No schema change downtime.
 - Well-suited for **unstructured** and **semi-structured** data.
-

Disadvantages of HBase

- Not suitable for **small datasets** or traditional relational data.
- **Complex architecture** and configuration.
- **No built-in support** for SQL queries (unlike RDBMS).
- Slower for operations involving **small data size**.
- Requires **Zookeeper** for coordination.



HBase Components – Notes in English



1. Master Server

- **Definition:** A Master Server is a central component responsible for managing the HBase cluster.
- **Role:**
 - Controls the **installation of software** related to **authentication** and **authorization**.
 - Maintains **centralized policy management**.
 - **Assigns tasks** to Region Servers.
 - Decides how Region Servers will work and distributes the load.
 - Handles **load balancing** and **region access** across the cluster.
 - Works **in coordination with ZooKeeper** to manage tasks and cluster communication.
- **Key Responsibilities:**
 - **Load Balancing:** Redistributes work from overloaded Region Servers to underloaded ones.

- Maintains the **state of the cluster** – monitors health and performance.
 - Manages **schema changes**, such as:
 - Metadata operations
 - Table creation
 - Row creation
 - Column family creation
-

2. Region Server

- **Definition:** Handles client requests and performs operations related to tables and data.
- **Primary Function:** Communication with clients to **read/write data** in HBase.
- **Structure:**
 - Tables are split into smaller **regions**, and these are handled by Region Servers.
 - Region Servers manage how these regions are **distributed across the file system or cloud**.
- **Responsibilities:**
 - Handles **data-related operations**.
 - Manages **read and write requests** from the file system.
 - Determines the **size of each region** based on defined threshold limits.
 - Controls how **HFiles** (HDFS files) are stored and managed in memory.
 - Ensures efficient **data handling, splitting, and storage** of tables.

What is ZooKeeper?

- **ZooKeeper** is an **open-source project**.
 - It provides services for:
 - **Maintaining configuration**
 - **Managing distributed synchronization**
 - **Storing information**
 - **Naming services**
-

Purpose of ZooKeeper in HBase

- Helps maintain the **state and coordination** between distributed components in the HBase system.
 - Ensures **synchronization** between Master Server and Region Servers.
 - Used for **server availability tracking** and **failure detection**.
 - Allows tracking of:
 - **Server failures**
 - **Network partitioning**
-

ZooKeeper and Nodes

- ZooKeeper uses **nodes** (called **znodes**) to represent and monitor different components (like Region Servers).
- These nodes help the **Master Server** to:
 - Discover available servers

- Handle failover and recovery
 - **Master Server** interacts with ZooKeeper to manage **cluster state** dynamically.
-

Communication Flow

- Region Servers **communicate through ZooKeeper**.
 - Clients also communicate with **HBase via ZooKeeper**.
 - Acts as a **coordinator** between all components.
-

ZooKeeper in HBase Setup

- HBase uses ZooKeeper **by default** as a **pseudo-standalone model**.
- Even in simple setups, **HBase internally depends on ZooKeeper** to manage its components efficiently.

MapReduce and Its Extension – Notes

What is MapReduce?

- **MapReduce** is a **programming model** used for **processing and generating large data sets**.
- It was originally developed by **Google**.
- Designed to handle **large-scale web-based data** like:
 - **Text files**

- **Search data**
 - It works with technologies like:
 - **BigTable**
 - **Google Distributed File System (GFS)**
-

Purpose of MapReduce

- Enables **parallel processing** of massive data collections.
 - Designed to process **large volumes of web data** efficiently.
 - Solves challenges related to:
 - **Data storage**
 - **Data retrieval**
 - **Data processing at scale**
-

Core Concepts

- **Map Phase:** Breaks down tasks and distributes them across nodes.
 - **Reduce Phase:** Aggregates the results from the Map phase.
 - Supports **parallel computing** using **thousands of processors**.
 - Ideal for dividing a huge task into smaller subtasks and executing them simultaneously.
-

Extension and Infrastructure

- The **MapReduce extension** refers to its application in **parallel computing** and **distributed environments**.
 - Designed to handle:
 - **Large-scale computation**
 - **Network failures**
 - **Fault tolerance** during data processing
-

Infrastructure Support

- Supports fault-tolerant infrastructure.
 - If network issues or failures occur during computation, the system ensures continuity and reliability.
 - The infrastructure is capable of **handling big data processing** efficiently.
-

Open Source Implementation

- **Hadoop** is the popular **open-source implementation** of MapReduce.
 - Developed to replicate Google's internal MapReduce model.
 - Initially used to understand and manage how **files are stored and processed** in large-scale systems.
-

Summary

- **MapReduce** was developed by **Google** to process web-scale data using **BigTable** and **GFS**.
- It enables **parallel processing** for handling **massive datasets**.



Parallel Computing – Notes



Introduction

- This topic is a continuation of the previous topic: **MapReduce Extension**.
 - **MapReduce** was introduced specifically to support **Parallel Computing**.
 - It is used to handle complex computations across multiple systems simultaneously.
-



What is Parallel Computing?

- A method of computation where **multiple processors** work **simultaneously** to solve a problem.
 - Introduced to solve complex problems using **scientific techniques**.
 - Has historical roots dating back to the **1960s and 1970s**.
-



Why Parallel Computing?

- Designed for systems that use **multiple processors**.
 - Helps solve problems that are too large or complex for a single processor to handle efficiently.
 - Especially useful for **large-scale server environments** and **web-based data**.
-



How It Works

- Large problems are **broken down into smaller sub-problems**.

- These sub-problems are distributed among **multiple processors** to be solved in parallel.
 - These smaller parts are managed across **servers and file systems** using:
 - **BigTable**
 - **Google Distributed File System (GFS)**
-

Application in File Systems

- Involves breaking data into **sections** or **chunks**.
 - Each server or processor handles a **small section** of the larger problem.
 - This division and distribution of tasks is what constitutes **parallel computation**.
-

Resource Utilization

- Resources are **allocated dynamically** based on the **need of the task at hand**.
 - The system determines:
 - **Which task needs which resource**
 - **When and how to allocate it**
 - This ensures **efficiency** and **optimized performance**.
-

Summary

- **Parallel Computing** is essential for handling **massive data** and complex tasks.
- Works by **dividing tasks** and distributing them among **multiple processors**.

- **MapReduce** plays a key role in implementing parallel computing in modern systems.
- Supports better **scalability**, **efficiency**, and **fault tolerance** in server-based architectures.
-
- Its **extension** supports **fault tolerance**, **network failure handling**, and **distributed processing**.
- **Hadoop** is the open-source version widely used today.

Efficiency of MapReduce – Notes

Background

- Previously, we studied **MapReduce** and **Parallel Computing**.
 - Now, the focus is on **how efficient MapReduce is** in practical applications.
-

What Does Efficiency Mean in MapReduce?

- **Efficiency** here refers to how well MapReduce handles:
 - **Communication costs**
 - **Computational costs**
 - **Synchronization overhead**
-

MapReduce in Cloud Computing

- MapReduce is typically executed on a **Cloud Computing Framework**.
- It was originally designed for:

- **Large-scale data analysis**
 - **Application-level data processing**
-

Performance Challenges

- In real-world scenarios, MapReduce can suffer from **performance degradation**.
 - Causes of inefficiency include:
 - **High adoption and frequent sequencing** of processes.
 - **Lack of optimization or modification** in workflow.
 - **Underutilization of cloud resources** during execution.
-

Key Factors Impacting Efficiency

- **Communication Cost**: Data transfer between nodes can slow down performance.
 - **Synchronization Cost**: Time spent ensuring all tasks/processes are aligned.
 - These hidden (implicit) costs affect the **overall efficiency** of MapReduce operations.
-

Conclusion

- The **efficiency** of MapReduce is closely tied to how well it manages **resource usage**, **communication**, and **process synchronization**.
- Without careful handling, it may **underutilize cloud resources** and experience **performance drops**.

MapReduce Model – Key Concepts

1. Introduction

- MapReduce is a **programming model** used for **parallel computing**.
- It is mainly used to **process large-scale data** by dividing it into smaller tasks.

2. Parallel Computing

- In the previous lecture, parallel computing was discussed.
- Parallel computing was originally developed for systems with a **small number of processors**.
- It ensures proper **communication and synchronization** between processors.

3. MapReduce – Basics

- MapReduce follows a **technical and programmatic approach**.
- It divides tasks into **small units**, enabling efficient processing on multiple processors.
- It is **Java-based** and contains two primary tasks:
 - **Map Task (Mapper)**
 - **Reduce Task (Reducer)**

4. Purpose of MapReduce

- Designed to work with a **large number of small processors**.
- Helps break down large data into small, manageable chunks.
- The goal is to **reduce** the data and **aggregate** results efficiently.

5. How MapReduce Works

- The **Map function**:

- Takes input data (e.g., documents).
- Breaks them into **key-value pairs**.
- Each document or element is split into small tuples (like word and count).
- The **Reduce function**:
 - Collects all key-value pairs from the map phase.
 - Combines and **reduces** them by counting occurrences of each word (for example).

6. Example – Word Count

- Multiple documents are stored in **HDFS (Hadoop Distributed File System)**.
- These documents are split into smaller chunks:
 - For example, 3 input files are divided into 3 mappers.
- Each mapper processes part of the data and emits key-value pairs (word, count).
- These outputs are then passed to **reducers**, which:
 - Aggregate the counts for each word.
 - Combine the data from mappers.
- The result is a **summary** showing how many times each word appeared.

7. Final Output

- The final output contains **unique key-value combinations** (no repetition).
- Each unique word is mapped with the total count of its occurrence.
- This is a simplified form of the data, suitable for further analysis.

Summary

- **MapReduce** = Map + Reduce (2 main tasks).
- Useful for **large-scale data processing**.
- Performs:
 - **Mapping**: Breaking data into key-value pairs.
 - **Reducing**: Aggregating those pairs.
- Mainly used in **Hadoop ecosystem** for tasks like **word counting**, **log analysis**, etc.
- Helps break down a complex job into smaller tasks for **parallel execution**.

Topic: Applications of MapReduce

MapReduce is widely used in processing large datasets. Below are the key areas where it is applied:

1. Social Networking

- Used to analyze user behavior.
- Helps in understanding user interests, preferences, and interactions.
- Example: Suggesting friends or content based on user activity.

2. Entertainment Industry

- Applied in platforms like streaming services to recommend content.
 - Works based on users' viewing history and preferences.
-

3. Electronic Communication / E-Commerce

- Used in platforms like Amazon, Flipkart, eBay.
 - Analyzes customer behavior: favorite products, preferences, and interests.
 - Helps understand user needs and optimize user experience.
 - Supports product recommendations and marketing strategies.
-

4. Fraud Detection

- Widely used in the **finance sector**, including:
 - Banks
 - Insurance companies
 - Payment services
 - Helps detect and reduce fraudulent transactions.
 - Analyzes transaction patterns and identifies anomalies.
-

5. Search and Advertisement Mechanism

- Helps in understanding:
 - Popular search terms
 - Most viewed or clicked ads
 - Regional user engagement
 - Performance of various search keywords
 - Used by search engines like Google to improve search accuracy and ad targeting.
-

Conclusion:

MapReduce is extensively used in our daily digital life in:

- Social Networking
- Entertainment
- E-Commerce
- Fraud Detection in Finance
- Search and Advertisement Mechanisms

Topic: Relational Operators

To understand **Relational Operators**, it's important to first know what **Relational Algebra** is and how **relations (tables)** work in a database.

- **Relational Algebra** works on **sets of data** (relations) and ensures that **duplicate rows are not allowed**.
 - These operations are used in SQL as well to manipulate and retrieve data from tables.
-

Types of Relational Operators:

1. Selection (σ)

- Used to select a **subset of tuples (rows)** from a relation.
 - It applies a condition to filter rows.
 - **Example:** Selecting rows where **age > 25**.
-

2. Projection (π)

- Used to select **specific columns** from a relation.
 - Removes duplicate values and shows only selected attributes.
 - Similar to the `SELECT column_name` command in SQL.
-

3. Union (\cup)

- Combines two relations (tables) **vertically**.
 - Only works when both tables have **same number of attributes and same data types**.
 - **Removes duplicate rows** in the result.
-

4. Intersection (\cap)

- Retrieves **common rows** between two relations.
 - Similar to SQL's `INTERSECT` operation.
 - Output includes only the data **present in both tables**.
-

5. Difference / Minus ($-$)

- Returns rows that are **present in the first table but not in the second**.
 - Helps in finding the difference between two datasets.
-

6. Natural Join (\bowtie)

- Combines two relations **based on common attributes**.
- It automatically merges tables on columns with the **same name**.

- Similar to SQL's **INNER JOIN**.
-

7. Grouping and Aggregation

- Groups rows based on certain attributes.
 - Applies **aggregate functions** like:
 - **SUM**
 - **COUNT**
 - **MIN**
 - **MAX**
 - Similar to SQL's **GROUP BY** clause.
-

Summary of Operators Covered:

- Selection
- Projection
- Union
- Intersection
- Difference
- Natural Join
- Grouping and Aggregation

These are key relational operators used in **Relational Algebra** to perform various data retrieval and transformation operations in databases.

Topic: Enterprise Batch Processing & Its Applications

♦ What is Enterprise Batch Processing?

- Enterprise Batch Processing refers to the **execution of tasks** in **batches** without any **user interaction**.
 - It runs **automatically** without the need for a person to trigger it.
-

♦ Key Characteristics:

- Batch processing **automates** tasks that are **repeated periodically**.
 - These tasks are designed to **run in the background**.
 - They are executed when:
 - **User interaction is minimal or not required.**
 - **Large amounts of data or resources** need to be processed.
-

♦ Applications of Batch Processing:

Used when the system needs to handle:

- **Large files or datasets**
- **Database creation or management**
- **Image processing**
- **Report generation** (e.g., monthly building reports)
- **Data format conversion**
- **Mass data uploads or transformations**

◆ How It Works:

- It involves **running batch jobs** in a **common system environment**.
- These batch jobs operate on a **shared infrastructure**, typically using:
 - **Java platforms**
 - **Business logic systems**

◆ Infrastructure & Platform:

- Batch processing applications run on an **enterprise-grade platform**.
- These platforms are:
 - **Reusable across multiple applications**
 - Built to **support heavy business logic**
 - Designed for **efficient task automation**

✓ Summary:

- Enterprise Batch Processing is crucial for **automating routine, large-scale tasks** without user intervention.
- Commonly used in business environments for **report generation, image processing, database handling, and data transformation**.
- Runs on a **common, Java-based infrastructure** to support enterprise needs.