# III/I (R16)

# DATABASE MANAGEMENT SYSTEMS

# UNIT -1

# SYLLABUS

**UNIT-I: An Overview of Database Management,**
- ➢ Introduction
- ➢ What is Database System
- ➢ What is Database
- ➢ Why Database
- ➢ Data Independence
- ➢ Relation Systems and Others Summary

**Database system architecture, Introduction**

- ➢ The Three Levels of Architecture
- ➢ The External Level
- ➢ The Conceptual Level
- ➢ The Internal Level
- ➢ Mapping
- ➢ The Database Administrator
- ➢ The Database Management Systems
- ➢ Client/Server Architecture.

**INTRODUCTION:-**

**Evolution of databases**

Databases have evolved since their inception in the 1960s, beginning with hierarchical and network databases, through the 1980s with object-oriented databases, and today with SQL and NoSQL databases and cloud databases.

In one view, databases can be classified according to content type: bibliographic, full text, numeric and images. In computing, databases are sometimes classified according to their organizational approach. There are many different kinds of databases, ranging from the most prevalent approach, the relational database, to a distributed database, cloud database or NoSQL database

A database is a collection of information that is organized so that it can be easily accessed, managed and updated.

Data is organized into rows, columns and tables, and it is indexed to make it easier to find relevant information. Data gets updated, expanded and deleted as new information is added. Databases process workloads to create and update themselves, querying the data they contain and running applications against it.

### What is Database System?

Database System in short refers to the technology of storing and retrieving users' data with utmost efficiency along with appropriate security measures. l

### What is Data?

In simple words data can be facts related to any object in consideration.
For example your name, age, height, weight, etc are some data related to you.
A picture , image , file , pdf etc can also be considered data.

# What is a Database?

Database is a systematic collection of data. Databases support storage and manipulation of data. Databases make data management easy. Let's discuss few examples.

An online telephone directory would definitely use database to store data pertaining to people, phone numbers, other contact details, etc.

Your electricity service provider is obviously using a database to manage billing , client related issues, to handle fault data, etc.

Let's also consider the facebook. It needs to store, manipulate and present data related to members, their friends, member activities, messages, advertisements and lot more.
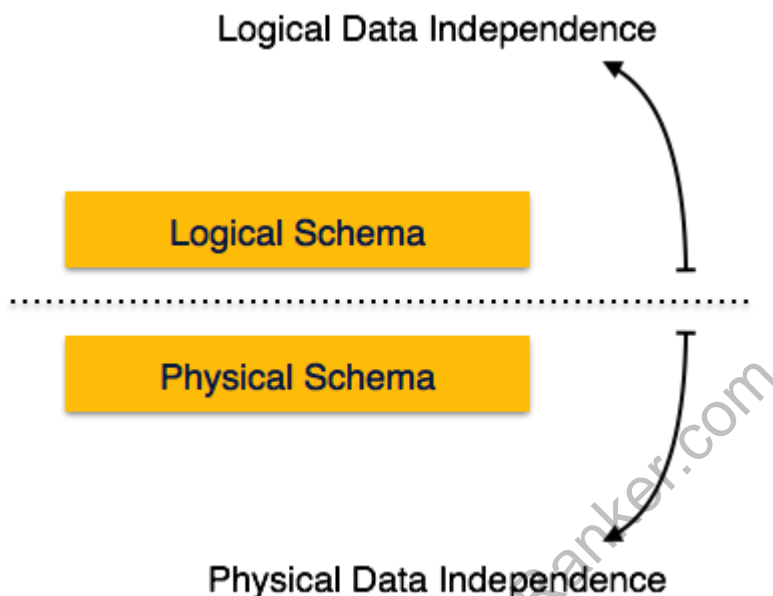
**Why Database?**

A database is an organized collection of various forms of data. It is also known as a structured set of data that is accessible in many ways through your computer.

The various reasons for which we require databases are:

- To manage large chunks of data: Yes, you can store data into a spreadsheet, but if you add large chunks of data into the sheet, it will simply not work. For instance: if your size of data increases into thousands of records, it will simply create a problem of speed.
- Accuracy: When doing data entry files in a spreadsheet, it becomes difficult to manage the accuracy as there are no validations present in it.
- Ease of updating data: With the database, you can flexibly update the data according to your convenience. Moreover, multiple people can also edit data at same time.
- Security of data: There is no denying the fact that your data is less secure in spreadsheets. Anyone can easily get access to file and can make changes to it. With databases you have security groups and privileges you set to restrict access.
- Data integrity: Data integrity also becomes a question when storing data in spreadsheets. In databases, you can be assured of accuracy and consistency of data due to the built in integrity checks and access controls.

# Data Independence

A database system normally contains a lot of data in addition to users' data. For example, it stores data about data, known as metadata, to locate and retrieve data easily. It is rather difficult to modify or update a set of metadata once it is stored in the database. But as a DBMS expands, it needs to change over time to satisfy the requirements of the users. If the entire data is dependent, it would become a tedious and highly complex job.

Metadata itself follows a layered architecture, so that when we change data at one layer, it does not affect the data at another level. This data is independent but mapped to each other.

## Logical Data Independence

Logical data is data about database, that is, it stores information about how data is managed inside. For example, a table (relation) stored in the database and all its constraints, applied on that relation.

Logical data independence is a kind of mechanism, which liberalizes itself from actual data stored on the disk. If we do some changes on table format, it should not change the data residing on the disk.

## Physical Data Independence

All the schemas are logical, and the actual data is stored in bit format on the disk. Physical data independence is the power to change the physical data without impacting the schema or logical data.

For example, in case we want to change or upgrade the storage system itself − suppose we want to replace hard-disks with SSD − it should not have any impact on the logical data or schemas.

**REALTIONAL SYSTEMS AND OTHERS:**

Relational data model is the primary data model in the Relational Systems, which is used widely around the world for data storage and processing. This model is simple and it has all the properties and capabilities required to process data with storage efficiency.

# Concepts

**Tables** − In relational data model, relations are saved in the format of Tables. This format stores the relation among entities. A table has rows and columns, where rows represents records and columns represent the attributes.

**Tuple** − A single row of a table, which contains a single record for that relation is called a tuple.

**Relation instance** − A finite set of tuples in the relational database system represents relation instance. Relation instances do not have duplicate tuples.

**Relation schema** − A relation schema describes the relation name (table name), attributes, and their names.

**Relation key** − Each row has one or more attributes, known as relation key, which can identify the row in the relation (table) uniquely.

**Attribute domain** − Every attribute has some pre-defined value scope, known as attribute domain.
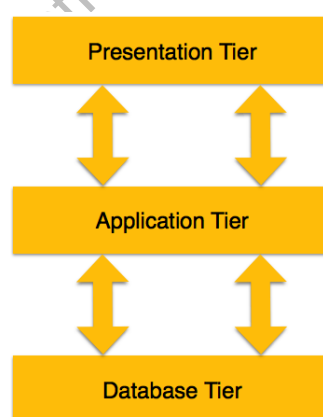
## DATABASE SYSTEM ARCHITECTURE:-

The design of a DBMS depends on its architecture. It can be centralized or decentralized or hierarchical. The architecture of a DBMS can be seen as either single tier or multi-tier. An n-tier architecture divides the whole system into related but independent **n** modules, which can be independently modified, altered, changed, or replaced.

In 1-tier architecture, the DBMS is the only entity where the user directly sits on the DBMS and uses it. Any changes done here will directly be done on the DBMS itself. It does not provide handy tools for end-users. Database designers and programmers normally prefer to use single-tier architecture.

If the architecture of DBMS is 2-tier, then it must have an application through which the DBMS can be accessed. Programmers use 2-tier architecture where they access the DBMS by means of an application. Here the application tier is entirely independent of the database in terms of operation, design, and programming.

## 3-tier Architecture

3-tier architecture separates its tiers from each other based on the complexity of the users and how they use the data present in the database. It is the most widely used architecture to design a DBMS.



- **Database (Data) Tier** − At this tier, the database resides along with its query processing languages. We also have the relations that define the data and their constraints at this level.

- **Application (Middle) Tier** − At this tier reside the application server and the programs that access the database. For a user, this application tier presents an abstracted view of the database. End-users are unaware of any existence of the database beyond the application. At the other end, the database tier is not aware of any other user beyond the application tier. Hence, the application layer

sits in the middle and acts as a mediator between the end-user and the database.

- **User (Presentation) Tier** – End-users operate on this tier and they know nothing about any existence of the database beyond this layer. At this layer, multiple views of the database can be provided by the application. All views are generated by applications that reside in the application tier.

**The Three Levels of Architecture:-**

**Following are the three levels of database architecture,**

1. Physical Level
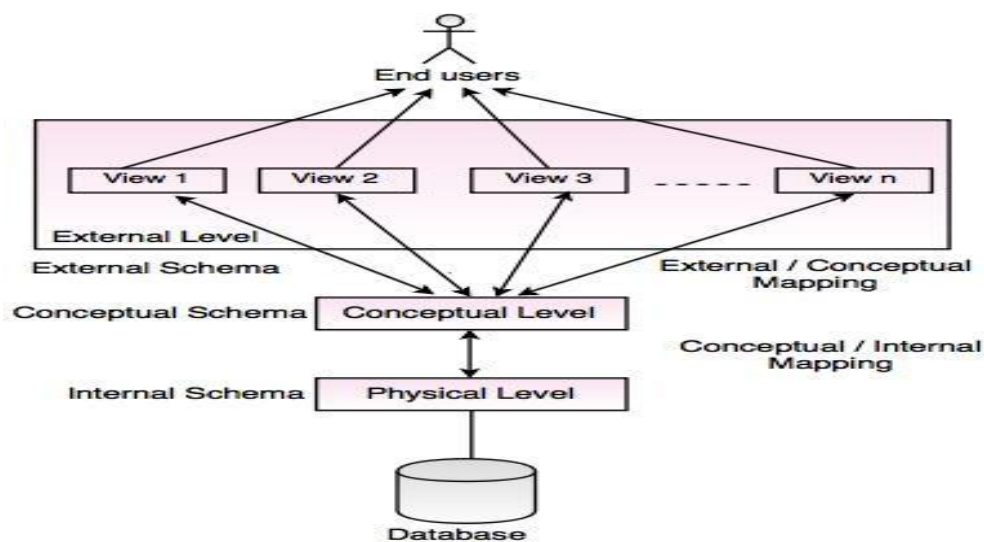
2. Conceptual Level

3. External Level



Fig. Three Level Architechture of DBMS

**In the above diagram,**
It shows the architecture of DBMS.
Mapping is the process of transforming request response between various database levels of architecture.
Mapping is not good for small database, because it takes more time.
In External / Conceptual mapping, DBMS transforms a request on an external schema against the conceptual schema.
In Conceptual / Internal mapping, it is necessary to transform the request from the conceptual to internal levels.

## 1. Physical Level

Physical level describes the physical storage structure of data in database.

It is also known as Internal Level.

This level is very close to physical storage of data.

At lowest level, it is stored in the form of bits with the physical addresses on the secondary storage device.

At highest level, it can be viewed in the form of files.

The internal schema defines the various stored data types. It uses a physical data model.

## 2. Conceptual Level

Conceptual level describes the structure of the whole database for a group of users.

It is also called as the data model.

Conceptual schema is a representation of the entire content of the database.

These schema contains all the information to build relevant external records.

It hides the internal details of physical storage.

## 3. External Level

External level is related to the data which is viewed by individual end users.

This level includes a no. of user views or external schemas.

This level is closest to the user.

External view describes the segment of the database that is required for a particular user group and hides the rest of the database from that user group.

**DATABASE  ADMINISTRATOR:-**

A database administrator (DBA) directs or performs all activities related to maintaining a successful databaseenvironment. Responsibilities include designing, implementing, and maintaining the database system; establishing policies and procedures pertaining to the management, security, maintenance, and use of thedatabase management system; and training employees in database management and use.

**ROLE OF DBA:-**

- Actively participates in planning the installation of new organization-wide systems and applications. He also assists during the installation as per specific functions.

- Implements the work plan for the Department. He meets the staff to identify any problems, take measures, and resolves it.

- Monitors the efficiency and effectiveness of all database resources and thus, keep the flow of work uninterrupted owing to technology.

- Ensures maximum service through identification of opportunities for improvement and make new recommendations.
- Continuous review and evaluation of the software, hardware, service delivery, and updates as and when required.
- Troubleshoot the problems [if any]. It includes a quick understanding of the problem and its resolution, restoration of the data, rectify the issue and minimize the damage.

## The Database Management Systems:-

A **database management system** stores data in such a way that it becomes easier to retrieve, manipulate, and produce information.

## Characteristics:-

Traditionally, data was organized in file formats. DBMS was a new concept then, and all the research was done to make it overcome the deficiencies in traditional style of data management. A modern DBMS has the following characteristics −

- **Real-world entity** − A modern DBMS is more realistic and uses real-world entities to design its architecture. It uses the behavior and attributes too. For example, a school database may use students as an entity and their age as an attribute.

- **Relation-based tables** − DBMS allows entities and relations among them to form tables. A user can understand the architecture of a database just by looking at the table names.

- **Isolation of data and application** − A database system is entirely different than its data. A database is an active entity, whereas data is said to be passive, on which the database works and organizes. DBMS also stores metadata, which is data about data, to ease its own process.

- **Less redundancy** − DBMS follows the rules of normalization, which splits a relation when any of its attributes is having redundancy in values. Normalization is a mathematically rich and scientific process that reduces data redundancy.

- **Consistency** − Consistency is a state where every relation in a database remains consistent. There exist methods and techniques, which can detect attempt of leaving database in inconsistent state. A DBMS can provide greater consistency as compared to earlier forms of data storing applications like file-processing systems.

- **Query Language** – DBMS is equipped with query language, which makes it more efficient to retrieve and manipulate data. A user can apply as many and as different filtering options as required to retrieve a set of data. Traditionally it was not possible where file-processing system was used.

- **ACID Properties** – DBMS follows the concepts of **A**tomicity, **C**onsistency, **I**solation, and **D**urability (normally shortened as ACID). These concepts are applied on transactions, which manipulate data in a database. ACID properties help the database stay healthy in multi-transactional environments and in case of failure.

- **Multiuser and Concurrent Access** – DBMS supports multi-user environment and allows them to access and manipulate data in parallel. Though there are restrictions on transactions when users attempt to handle the same data item, but users are always unaware of them.

- **Multiple views** – DBMS offers multiple views for different users. A user who is in the Sales department will have a different view of database than a person working in the Production department. This feature enables the users to have a concentrate view of the database according to their requirements.

- **Security** – Features like multiple views offer security to some extent where users are unable to access data of other users and departments. DBMS offers methods to impose constraints while entering data into the database and retrieving the same at a later stage. DBMS offers many different levels of security features, which enables multiple users to have different views with different features. For example, a user in the Sales department cannot see the data that belongs to the Purchase department. Additionally, it can also be managed how much data of the Sales department should be displayed to the user. Since a DBMS is not saved on the disk as traditional file systems, it is very hard for miscreants to break the code.

**Client/Server Architecture:-**

**Client-server architecture** (**client/server**) is a network **architecture** in which each computer or process on the network is either a **client** or a **server**. Servers are powerful computers or processes dedicated to managing disk drives (file servers), printers (print servers), or network traffic (network servers).

In client/server architecture, the database application and the database are separated into two parts: a front-end or client portion, and a back-end or server portion. The client executes the database application that accesses database information and interacts with a user through the keyboard, screen, and pointing device such as a mouse. The server executes the Oracle software and handles the functions required for concurrent, shared data access to an Oracle database.

Although the client application and Oracle can be executed on the same computer, it may be more efficient and effective when the client portion(s) and server portion are executed by different computers connected via a network.

The below figure shows the interaction between Client and |Server:-
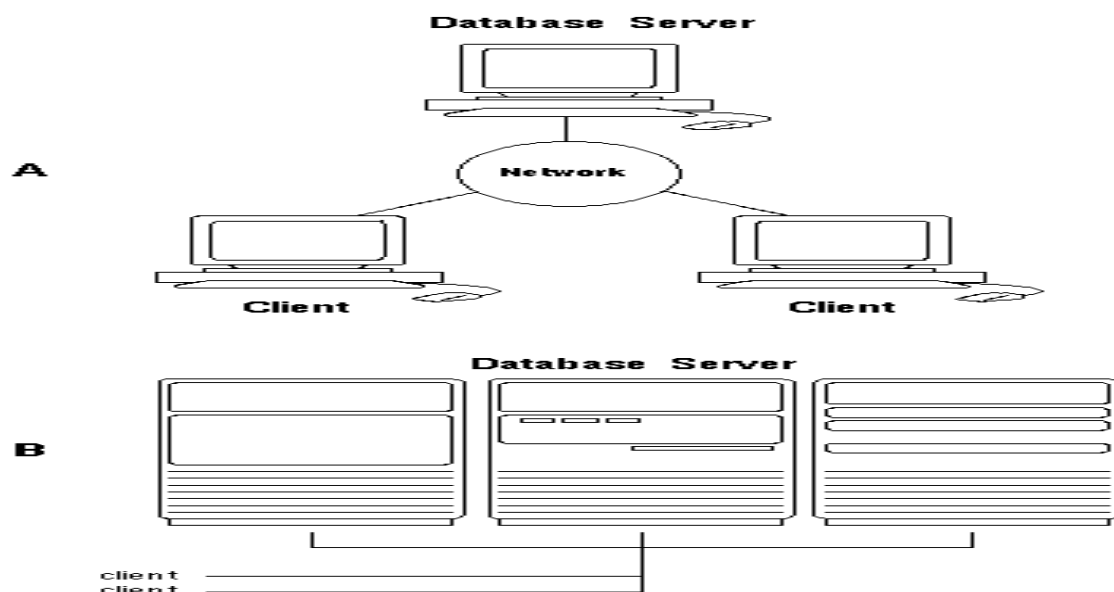


Figure:- **The Client/Server Architecture and Distributed Processing**

**Benefits of the Oracle client/server architecture in a distributed processing environment include the following:**

- Client applications are not responsible for performing any data processing. Client applications can concentrate on requesting input from users, requesting desired data from the server, and then analyzing and presenting this data using the display capabilities of the client workstation or the terminal (for example, using graphics or spreadsheets).

- Client applications can be designed with no dependence on the physical location of the data. If the data is moved or distributed to other database servers, the application continues to function with little or no modification.

- Oracle exploits the multitasking and shared-memory facilities of its underlying operating system. As a result, it delivers the highest possible degree of concurrency, data integrity, and performance to its client applications.

- Client workstations or terminals can be optimized for the presentation of data (for example, by providing graphics and mouse support) and the server can be optimized for the processing and storage of data (for example, by having large amounts of memory and disk space).

- If necessary, Oracle can be *scaled*. As your system grows, you can add multiple servers to distribute the database processing load throughout the network (*horizontally scaled*). Alternatively, you can replace Oracle on a less powerful computer, such as a microcomputer, with Oracle running on a minicomputer or mainframe, to take advantage of a larger system's performance (*vertically scaled*). In either case, all data and applications are maintained with little or no modification, since Oracle is portable between systems.

- In networked environments, shared data is stored on the servers, rather than on all computers in the system. This makes it easier and more efficient to manage concurrent access.

- In networked environments, inexpensive, low-end client workstations can be used to access the remote data of the server effectively.

- In networked environments, client applications submit database requests to the server using SQL statements. Once received, the SQL statement is processed by the server, and the results are returned to the client application. Network traffic is kept to a minimum because only the requests and the results are shipped over the network.

# UNIT -2

# SYLLABUS

**UNIT-II:**
  - The E/R Models
  - The Relational Model
  - Relational Calculus
  - Introduction to Database Design
  - Database Design and E/R Diagrams
  - Entities Attributes and Entity Sets-Relationship and Relationship Sets
  - Conceptual Design with the E/R Models
  - The Relational Model Integrity
  - Constraints over Relations
  - Key Constraints
  - Foreign Key Constraints
  - General Constraints
  - Relational Algebra and Calculus
  - Relational Algebra
  - Selection and Projection
  - Set Operation
  - Renaming
  - Joins
  - Division
  - More Examples of Queries
  - Relational Calculus
  - Tuple Relational Calculus
  - Domain Relational Calculus

### ER Model - Basic Concepts:-

The ER model defines the conceptual view of a database. It works around real-world entities and the associations among them. At view level, the ER model is considered a good option for designing databases.

## Entity

An entity can be a real-world object, either animate or inanimate, that can be easily identifiable. For example, in a school database, students, teachers, classes, and courses offered can be considered as entities. All these entities have some attributes or properties that give them their identity.

An entity set is a collection of similar types of entities. An entity set may contain entities with attribute sharing similar values. For example, a Students set may contain all the students of a school; likewise a Teachers set may contain all the teachers of a school from all faculties. Entity sets need not be disjoint.

## Attributes

Entities are represented by means of their properties, called **attributes**. All attributes have values. For example, a student entity may have name, class, and age as attributes.

There exists a domain or range of values that can be assigned to attributes. For example, a student's name cannot be a numeric value. It has to be alphabetic. A student's age cannot be negative, etc.

## Types of Attributes

- **Simple attribute** – Simple attributes are atomic values, which cannot be divided further. For example, a student's phone number is an atomic value of 10 digits.

- **Composite attribute** – Composite attributes are made of more than one simple attribute. For example, a student's complete name may have first_name and last_name.

- **Derived attribute** – Derived attributes are the attributes that do not exist in the physical database, but their values are derived from other attributes present in the database. For example, average_salary in a department should not be saved directly in the database, instead it can be derived. For another example, age can be derived from data_of_birth.

- **Single-value attribute** – Single-value attributes contain single value. For example – Social_Security_Number.

- **Multi-value attribute** – Multi-value attributes may contain more than one values. For example, a person can have more than one phone number, email_address, etc.

These attribute types can come together in a way like –

- simple single-valued attributes

- simple multi-valued attributes

- composite single-valued attributes

- composite multi-valued attributes

## Entity-Set and Keys

Key is an attribute or collection of attributes that uniquely identifies an entity among entity set.

For example, the roll_number of a student makes him/her identifiable among students.

- **Super Key** – A set of attributes (one or more) that collectively identifies an entity in an entity set.

- **Candidate Key** – A minimal super key is called a candidate key. An entity set may have more than one candidate key.

- **Primary Key** – A primary key is one of the candidate keys chosen by the database designer to uniquely identify the entity set.

# Relationship

The association among entities is called a relationship. For example, an employee **works_at** a department, a student**enrolls** in a course. Here, Works_at and Enrolls are called relationships.

## Relationship Set

A set of relationships of similar type is called a relationship set. Like entities, a relationship too can have attributes. These attributes are called **descriptive attributes**.
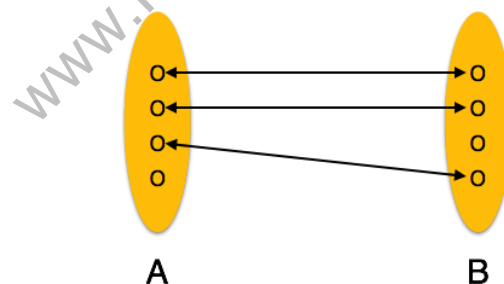
## Degree of Relationship

The number of participating entities in a relationship defines the degree of the relationship.

- Binary = degree 2
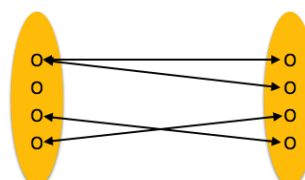- Ternary = degree 3
- n-ary = degree

## Mapping Cardinalities

**Cardinality** defines the number of entities in one entity set, which can be associated with the number of entities of other set via relationship set.
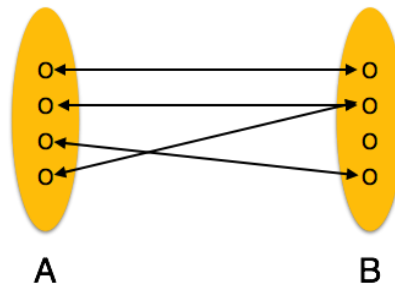
- **One-to-one** – One entity from entity set A can be associated with at most one entity of entity set B and vice versa.
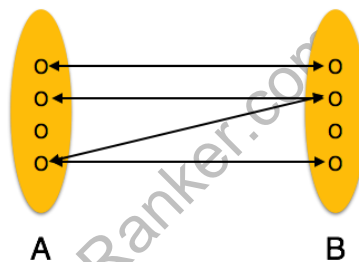


- **One-to-many** – One entity from entity set A can be associated with more than one entities of entity set B however an entity from entity set B, can be associated with at most one entity.

- **Many-to-one** – More than one entities from entity set A can be associated with at most one entity of entity set B, however an entity from entity set B can be associated with more than one entity from entity set A.



- **Many-to-many** – One entity from A can be associated with more than one entity from B and vice versa.



## Relational Model:-

The most popular data model in DBMS is the Relational Model. It is more scientific a model than others. This model is based on first-order predicate logic and defines a table as an **n-ary relation**.

The main highlights of this model are −

- Data is stored in tables called **relations**.
- Relations can be normalized.
- In normalized relations, values saved are atomic values.
- Each row in a relation contains a unique value.
- Each column in a relation contains values from a same domain.

## Relational Calculus

In contrast to Relational Algebra, Relational Calculus is a non-procedural query language, that is, it tells what to do but never explains how to do it.

Relational calculus exists in two forms −

## Tuple Relational Calculus (TRC)

Filtering variable ranges over tuples

**Notation** − {T | Condition}

Returns all tuples T that satisfies a condition.

**For example** −

```
{ T.name |  Author(T) AND T.article = 'database' }
```

**Output** − Returns tuples with 'name' from Author who has written article on 'database'.

TRC can be quantified. We can use Existential (∃) and Universal Quantifiers (∀).

**For example** −

```
{ R| ∃T   ∈ Authors(T.article='database' AND R.name=T.name)}
```

**Output** − The above query will yield the same result as the previous one.

## Domain Relational Calculus (DRC)

In DRC, the filtering variable uses the domain of attributes instead of entire tuple values (as done in TRC, mentioned above).

**Notation** −

{ $a_1$, $a_2$, $a_3$, ..., $a_n$ | P ($a_1$, $a_2$, $a_3$, ... ,$a_n$)}

Where a1, a2 are attributes and **P** stands for formulae built by inner attributes.

**For example** −

```
{< article, page, subject > |  ∈ TutorialsPoint ∧ subject = 'database'}
```

**Output** − Yields Article, Page, and Subject from the relation TutorialsPoint, where subject is database.

Just like TRC, DRC can also be written using existential and universal quantifiers. DRC also involves relational operators.

The expression power of Tuple Relation Calculus and Domain Relation Calculus is equivalent to Relational Algebra.


# Database Design:

### Introduction to Database Design:-

Database design is the design of the database structure that will be used to store and manage data rather than the design of the DBMS software. Once the database design is completed, the DBMS handles all the complicated activities required to translate the designer's view of the structures into structures that are usable to the computer.

A poorly designed database tends to generate errors that are likely to lead to bad decisions. A bad database design eventually can be self correcting: organizations using poorly designed databases often fail because their managers do not have access to timely (or even correct) information, thereby dominating the bad database design.

The availability of a DBMS makes it possible to tackle far more sophisticated uses of the data resources, if the database is designed to make use of that available power. The kinds of data structures created within the database and the extent of the relationships among them play a powerful role in determining how effective the DBMS is. Therefore, database design become a crucial activity in the database environment.

Database design is made much simpler when we use **models.** A Database model is a collection of a logical constructs used to represent the data structure and the data relationships found within the database i.e. simplified abstractions of real-world events or conditions. If the models are not logically sound, the database designs derived from them will not deliver the database system's promise to effective information drawn from an efficient database. "*Good models yield good database design that are the basis for good applications*".

### Goals of Database Design :

Database Design normally involves defining the logical attributes of the database designing the layout of the database file structure.
The main objectives of database design is
1. To satisfy the information content requirement of the specified user and application.
2. To provide a natural and easy way to understand structuring of the information.
3. To support processing requirements and any performance objectives such as

i. Response time
ii. Processing time
iii. Storage space
The main objective of the database design is to ensure that the database meets the reporting and information requirements of the users efficiently. The database should be designed in such a way that :
i. It eliminates or minimizes data redundancy.
ii. Maintains the integrity and independence of the data.

## Phases in Database Design :

First phase: The overall purpose of the database initial study is to

a. analyze the organization/system situation

b. define problem and constraints

c. define objectives

d. define scope and boundaries.

Second phase : The second phase focuses on the design of the database model that will support organization operations and objectives.

In this phase, we can identify six main phases of the database design :

I. Requirements collection and analysis

II. Conceptual database design

III. Choice of DBMS

IV. Data model mapping

V. Physical database design

VI. Database system implement

**E-R Diagrams:-**

## ER Diagram Representation

Let us now learn how the ER Model is represented by means of an ER diagram. Any object, for example, entities, attributes of an entity, relationship sets, and attributes of relationship sets, can be represented with the help of an ER diagram.

## Entity

Entities are represented by means of rectangles. Rectangles are named with the entity set they represent.



## Attributes

Attributes are the properties of entities. Attributes are represented by means of ellipses. Every ellipse represents one attribute and is directly connected to its entity (rectangle).



If the attributes are **composite**, they are further divided in a tree like structure. Every node is then connected to its attribute. That is, composite attributes are represented by ellipses that are connected with an ellipse.

**Multivalued** attributes are depicted by double ellipse.


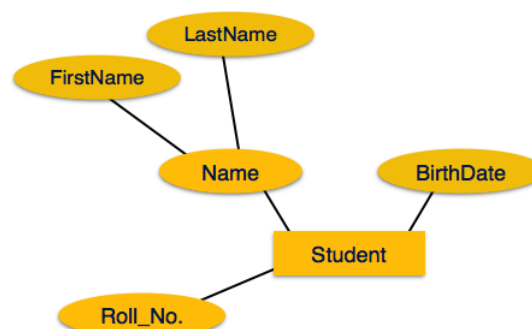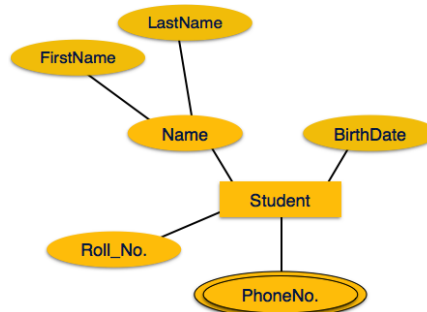
**Derived** attributes are depicted by dashed ellipse.



# Relationship

Relationships are represented by diamond-shaped box. Name of the relationship is written inside the diamond-box. All the entities (rectangles) participating in a relationship, are connected to it by a line.

## Binary Relationship and Cardinality

A relationship where two entities are participating is called a**binary relationship**. Cardinality is the number of instance of an entity from a relation that can be associated with the relation.

- **One-to-one** – When only one instance of an entity is associated with the relationship, it is marked as '1:1'. The following image reflects that only one instance of each entity should be associated with the relationship. It depicts one-to-one relationship.

- **One-to-many** – When more than one instance of an entity is associated with a relationship, it is marked as '1:N'. The following image reflects that only one instance of entity on the left and more than one instance of an entity on the right can be associated with the relationship. It depicts one-to-many relationship.



- **Many-to-one** – When more than one instance of entity is associated with the relationship, it is marked as 'N:1'. The following image reflects that more than one instance of an entity on the left and only one instance of an entity on the right can be associated with the relationship. It depicts many-to-one relationship.



- **Many-to-many** – The following image reflects that more than one instance of an entity on the left and more than one instance of an entity on the right can be associated with the relationship. It depicts many-to-many relationship.



## Participation Constraints

- **Total Participation** – Each entity is involved in the relationship. Total participation is represented by double lines.

- **Partial participation** – Not all entities are involved in the relationship. Partial participation is represented by single lines.

The ER Model has the power of expressing database entities in a conceptual hierarchical manner. As the hierarchy goes up, it generalizes the view of entities, and as we go deep in the hierarchy, it gives us the detail of every entity included.

Going up in this structure is called **generalization**, where entities are clubbed together to represent a more generalized view. For example, a particular student named Mira can be generalized along with all the students. The entity shall be a student, and further, the student is a person. The reverse is called **specialization** where a person is a student, and that student is Mira.

## Generalization

As mentioned above, the process of generalizing entities, where the generalized entities contain the properties of all the generalized entities, is called generalization. In generalization, a number of entities are brought together into one generalized entity based on their similar characteristics. For example, pigeon, house sparrow, crow and dove can all be generalized as Birds.



## Specialization

Specialization is the opposite of generalization. In specialization, a group of entities is divided into sub-groups based on their characteristics. Take a group 'Person' for example. A person has name, date of birth, gender, etc. These properties are common in all persons, human beings. But in a company, persons can be identified as employee, employer, customer, or vendor, based on what role they play in the company.

Similarly, in a school database, persons can be specialized as teacher, student, or a staff, based on what role they play in school as entities.

## Conceptual Design with E-R Model:-

- v   Requirement analysis

    – Data to be stored

    – Applications to be built

    – Operations (most frequent) subject to performance requirement

- v   Conceptual db design

    – Description of the data (including constraints)

    – By high level model such as ER

- v   Logical db design

    – Choose DBMS to implement

    – Convert conceptual db design into database schema

Issues to consider:   *(ER Model is used at this stage.)*

    – What are the *entities* and *relationships* in the enterprise?

    – What information about these entities and relationships should we store in the database (i.e., attributes)?

    – What are the *integrity constraints* or *business rules* that hold?

Solution:

    – A database `schema' in the ER Model can be represented pictorially (*ER diagrams*).

    – Can map an ER diagram into a relational schema.

**The Relational Model Integrity: Constraints Over Relations**

Constraints:-

Every relation has some conditions that must hold for it to be a valid relation. These conditions are called **Relational Integrity Constraints**. There are three main integrity constraints −

- Key constraints
- Domain constraints
- Referential integrity constraints

# Key Constraints:-

There must be at least one minimal subset of attributes in the relation, which can identify a tuple uniquely. This minimal subset of attributes is called **key** for that relation. If there are more than one such minimal subsets, these are called *candidate keys*.

Key constraints force that −

- in a relation with a key attribute, no two tuples can have identical values for key attributes.
- a key attribute can not have NULL values.

Key constraints are also referred to as Entity Constraints.

# Domain Constraints:-

Attributes have specific values in real-world scenario. For example, age can only be a positive integer. The same constraints have been tried to employ on the attributes of a relation. Every attribute is bound to have a specific range of values. For example, age cannot be less than zero and telephone numbers cannot contain a digit outside 0-9.

# Referential integrity Constraints:-

Referential integrity constraints work on the concept of Foreign Keys. A foreign key is a key attribute of a relation that can be referred in other relation.

Referential integrity constraint states that if a relation refers to a key attribute of a different or same relation, then that key element must exist.

## FOREIGN KEY Constraint:-

A FOREIGN KEY is a key used to link two tables together.

A FOREIGN KEY is a field (or collection of fields) in one table that refers to the PRIMARY KEY in another table.

The table containing the foreign key is called the child table, and the table containing the candidate key is called the referenced or parent table.

Look at the following two tables:

Look at the following two tables:

"Persons" table:

| PersonID | LastName | FirstName | Age |
|----------|-----------|-----------|-----|
| 1 | Hansen | Ola | 30 |
| 2 | Svendson | Tove | 23 |
| 3 | Pettersen | Kari | 20 |

"Orders" table:

| OrderID | OrderNumber | PersonID |
|---------|-------------|----------|
| 1 | 77895 | 3 |
| 2 | 44678 | 3 |

| 3 | 22456 | 2 |
| 4 | 24562 | 1 |

Notice that the "PersonID" column in the "Orders" table points to the "PersonID" column in the "Persons" table.

The "PersonID" column in the "Persons" table is the PRIMARY KEY in the "Persons" table.

The "PersonID" column in the "Orders" table is a FOREIGN KEY in the "Orders" table.

The FOREIGN KEY constraint is used to prevent actions that would destroy links between tables.

The FOREIGN KEY constraint also prevents invalid data from being inserted into the foreign key column, because it has to be one of the values contained in the table it points to.

**FOREIGN KEY on CREATE TABLE:-**

The following SQL creates a FOREIGN KEY on the "PersonID" column when the "Orders" table is created:

**MySQL:**

```
CREATE TABLE Orders (
    OrderID int NOT NULL,
    OrderNumber int NOT NULL,
    PersonID int,
    PRIMARY KEY (OrderID),
    FOREIGN KEY (PersonID) REFERENCES Persons(PersonID)
);
```

**Relational Algebra and Calculus**:-

**Introduction:-**

Relational database systems are expected to be equipped with a query language that can assist its users to query the database instances. There are two kinds of query languages − relational algebra and relational calculus.

**Relational Algebra**

Relational algebra is a procedural query language, which takes instances of relations as input and yields instances of relations as output. It uses operators to perform queries. An operator can be either **unary** or **binary**. They accept relations as their input and yield relations as their output. Relational algebra is performed recursively on a relation and intermediate results are also considered relations.

The fundamental operations of relational algebra are as follows −

- Select
- Project
- Union
- Set different
- Cartesian product
- Rename

We will discuss all these operations in the following sections.

# Select Operation (σ)

It selects tuples that satisfy the given predicate from a relation.

**Notation** − $\sigma_p(r)$

Where **σ** stands for selection predicate and **r** stands for relation. $p$ is prepositional logic formula which may use connectors like **and, or,** and **not**. These terms may use relational operators like − $=, \neq, \geq, <, >, \leq$.

**For example** −

$\sigma_{subject = "database"}(Books)$

**Output** − Selects tuples from books where subject is 'database'.

$\sigma_{subject = "database" \text{ and } price = "450"}(Books)$

**Output** − Selects tuples from books where subject is 'database' and 'price' is 450.

$$\sigma_{subject = "database" \text{ and } price = "450" \text{ or } year > "2010"}(Books)$$

**Output** − Selects tuples from books where subject is 'database' and 'price' is 450 or those books published after 2010.

# Project Operation (∏)

It projects column(s) that satisfy a given predicate.

Notation − $\prod_{A_1, A_2, A_n} (r)$

Where $A_1, A_2, A_n$ are attribute names of relation **r**.

Duplicate rows are automatically eliminated, as relation is a set.

**For example** −

$$\prod_{subject, author} (Books)$$

Selects and projects columns named as subject and author from the relation Books.

# Union Operation (∪)

It performs binary union between two given relations and is defined as −

$$r \cup s = \{ t \mid t \in r \text{ or } t \in s\}$$

**Notation** − r U s

Where **r** and **s** are either database relations or relation result set (temporary relation).

For a union operation to be valid, the following conditions must hold −

- **r**, and **s** must have the same number of attributes.
- Attribute domains must be compatible.
- Duplicate tuples are automatically eliminated.

$$\prod_{author} (Books) \cup \prod_{author} (Articles)$$

**Output** − Projects the names of the authors who have either written a book or an article or both.

# Set Difference (−)

The result of set difference operation is tuples, which are present in one relation but are not in the second relation.

**Notation** − **r** − **s**

Finds all the tuples that are present in **r** but not in **s**.

> ∏ ₐᵤₜₕₒᵣ (Books) − ∏ ₐᵤₜₕₒᵣ (Articles)

**Output** − Provides the name of authors who have written books but not articles.

# Cartesian Product (X)

Combines information of two different relations into one.

**Notation** − r X s

Where **r** and **s** are relations and their output will be defined as −

r X s = { q t | q ∈ r and t ∈ s}

> σ₍ₐᵤₜₕₒᵣ = 'tutorialspoint'₎(Books X Articles)

**Output** − Yields a relation, which shows all the books and articles written by tutorialspoint.

# Rename Operation (ρ)

The results of relational algebra are also relations but without any name. The rename operation allows us to rename the output relation. 'rename' operation is denoted with small Greek letter **rho**ρ.

**Notation** − ρ ₓ (E)

Where the result of expression **E** is saved with name of **x**.

Additional operations are −

- Set intersection
- Assignment
- Natural join

### Joins:-

**Join** is a combination of a Cartesian product followed by a selection process. A Join operation pairs two tuples from different relations, if and only if a given join condition is satisfied.

# Theta (θ) Join

Theta join combines tuples from different relations provided they satisfy the theta condition. The join condition is denoted by the symbol **θ**.

## Notation

```
R1 ⋈θ R2
```

R1 and R2 are relations having attributes (A1, A2, .., An) and (B1, B2,.. ,Bn) such that the attributes don't have anything in common, that is R1 ∩ R2 = Φ.

Theta join can use all kinds of comparison operators.

| Student | | |
|---------|---|---|
| **SID** | **Name** | **Std** |
| 101 | Alex | 10 |
| 102 | Maria | 11 |

| Subjects | |
|----------|---|
| **Class** | **Subject** |
| 10 | Math |
| 10 | English |
| 11 | Music |

| 11 | Sports |
|----|--------|

Student_Detail −

STUDENT ⋈Student.Std = Subject.Class SUBJECT

| Student_detail | | | | |
|----|----|----|----|----|
| **SID** | **Name** | **Std** | **Class** | **Subject** |
| 101 | Alex | 10 | 10 | Math |
| 101 | Alex | 10 | 10 | English |
| 102 | Maria | 11 | 11 | Music |
| 102 | Maria | 11 | 11 | Sports |

# Equijoin

When Theta join uses only **equality** comparison operator, it is said to be equijoin. The above example corresponds to equijoin.

# Natural Join (⋈)

Natural join does not use any comparison operator. It does not concatenate the way a Cartesian product does. We can perform a Natural Join only if there is at least one common attribute that exists between two relations. In addition, the attributes must have the same name and domain.

Natural join acts on those matching attributes where the values of attributes in both the relations are same.

| Courses | | |
|----|----|----|
| **CID** | **Course** | **Dept** |

| CS01 | Database | CS |
| ME01 | Mechanics | ME |
| EE01 | Electronics | EE |

| HoD | |
|---|---|
| **Dept** | **Head** |
| CS | Alex |
| ME | Maya |
| EE | Mira |

| Courses ⋈ HoD | | | |
|---|---|---|---|
| **Dept** | **CID** | **Course** | **Head** |
| CS | CS01 | Database | Alex |
| ME | ME01 | Mechanics | Maya |
| EE | EE01 | Electronics | Mira |

# Outer Joins

Theta Join, Equijoin, and Natural Join are called inner joins. An inner join includes only those tuples with matching attributes and the rest are discarded in the resulting relation. Therefore, we need to use outer joins to include all the tuples from the participating relations in the resulting relation. There are three kinds of outer joins − left outer join, right outer join, and full outer join.

# Left Outer Join(R ⋈ S)

All the tuples from the Left relation, R, are included in the resulting relation. If there are tuples in R without any matching tuple in the Right relation S, then the S-attributes of the resulting relation are made NULL.

| Left | |
|---|---|
| **A** | **B** |
| 100 | Database |
| 101 | Mechanics |
| 102 | Electronics |

| Right | |
|---|---|
| **A** | **B** |
| 100 | Alex |
| 102 | Maya |
| 104 | Mira |

| Courses ⋈ HoD | | | |
|---|---|---|---|
| **A** | **B** | **C** | **D** |
| 100 | Database | 100 | Alex |
| 101 | Mechanics | --- | --- |

| 102 | Electronics | 102 | Maya |
|-----|-------------|-----|------|

## Right Outer Join: ( R ⋈ S )

All the tuples from the Right relation, S, are included in the resulting relation. If there are tuples in S without any matching tuple in R, then the R-attributes of resulting relation are made NULL.

| Courses ⋈ HoD | | | |
|---|---|---|---|
| **A** | **B** | **C** | **D** |
| 100 | Database | 100 | Alex |
| 102 | Electronics | 102 | Maya |
| --- | --- | 104 | Mira |

## Full Outer Join: ( R ⋈ S)

All the tuples from both participating relations are included in the resulting relation. If there are no matching tuples for both relations, their respective unmatched attributes are made NULL.

| Courses ⋈ HoD | | | |
|---|---|---|---|
| **A** | **B** | **C** | **D** |
| 100 | Database | 100 | Alex |
| 101 | Mechanics | --- | --- |
| 102 | Electronics | 102 | Maya |
| --- | --- | 104 | Mira |

### DIVISION:-

Division is typically required when you want to find out entities that are interacting with **all entities** of a set of different type entities.

**Some instances where division operator is used are:**
- Which person has account in all the banks of a particular city?
- Which students have taken all the courses required to graduate?
- 

In all these queries, the description after the keyword 'all' defines a set which contains some elements and the final result contains those units who satisfy these requirements.

**Important: Division is not supported by SQL implementations. However, it can be represented using other operations.(like cross join, Except, In )SQL Implementation of Division.**

**Given two relations(tables): R(x,y) , S(y).**
**R and S** : tables
**x and y** : column of R
**y** : column of S

**R(x,y) div S(y)** means gives all distinct values of x from R that are associated with all values of y in S.

**Computation of Division:** R(x,y) div S(y)
**Steps:**
- Find out all possible combinations of S(y) with R(x) by computing R(x) x(cross join) S(y), say r1
- Subtract actual R(x,y) from r1, say r2
- x in r2 are those that are not associated with every value in S(y); therefore R(x)-r2(x) gives us x that are associated with all values in S

### Relational Calculus

In contrast to Relational Algebra, Relational Calculus is a non-procedural query language, that is, it tells what to do but never explains how to do it.

Relational calculus exists in two forms −

## Tuple Relational Calculus (TRC)

Filtering variable ranges over tuples

**Notation** − {T | Condition}

Returns all tuples T that satisfies a condition.

**For example** −

```
{ T.name |  Author(T) AND T.article = 'database' }
```

**Output** − Returns tuples with 'name' from Author who has written article on 'database'.

TRC can be quantified. We can use Existential (∃) and Universal Quantifiers (∀).

**For example** −

```
{ R| ∃T  ∈ Authors(T.article='database' AND R.name=T.name)}
```

**Output** − The above query will yield the same result as the previous one.

# Domain Relational Calculus (DRC)

In DRC, the filtering variable uses the domain of attributes instead of entire tuple values (as done in TRC, mentioned above).

**Notation** −

$\{ a_1, a_2, a_3, ..., a_n \mid P (a_1, a_2, a_3, ... ,a_n)\}$

Where a1, a2 are attributes and **P** stands for formulae built by inner attributes.

**For example** −

```
{< article, page, subject > |  ∈ TutorialsPoint ∧ subject = 'database'}
```

**Output** − Yields Article, Page, and Subject from the relation TutorialsPoint, where subject is database.

Just like TRC, DRC can also be written using existential and universal quantifiers. DRC also involves relational operators.

The expression power of Tuple Relation Calculus and Domain Relation Calculus is equivalent to Relational Algebra.

# UNIT -3

## Queries, Constraints, Triggers

## SYLLABUS

**UNIT-III**

- ➤ The Form of Basic SQL Query
- ➤ Union
- ➤ Intersect and Except
- ➤ Nested Queries
- ➤ Aggregate Operators
- ➤ Null Values
- ➤ Complex Integrity Constraints in SQL
- ➤ Triggers and Active Database.

**The Form of Basic SQL Query:-**

**Introduction to SQL:-**

SQL is a database computer language designed for the retrieval and management of data in a relational database. SQL stands for Structured Query Language. This tutorial will give you a quick start to SQL. It covers most of the topics required for a basic understanding of SQL and to get a feel of how it works.

SQL is a language to operate databases; it includes database creation, deletion, fetching rows, modifying rows, etc. SQL is an**ANSI** (American National Standards Institute) standard language, but there are many different versions of the SQL language.

# What is SQL?

SQL is Structured Query Language, which is a computer language for storing, manipulating and retrieving data stored in a relational database.

SQL is the standard language for Relational Database System. All the Relational Database Management Systems (RDMS) like MySQL, MS Access, Oracle, Sybase, Informix, Postgres and SQL Server use SQL as their standard database language.

Also, they are using different dialects, such as −

- MS SQL Server using T-SQL,
- Oracle using PL/SQL,
- MS Access version of SQL is called JET SQL (native format) etc.

# Why SQL?

SQL is widely popular because it offers the following advantages −

- Allows users to access data in the relational database management systems.

- Allows users to describe the data.

- Allows users to define the data in a database and manipulate that data.

- Allows to embed within other languages using SQL modules, libraries & pre-compilers.

- Allows users to create and drop databases and tables.

- Allows users to create view, stored procedure, functions in a database.

- Allows users to set permissions on tables, procedures and views.

A Brief History of SQL

- **1970** – Dr. Edgar F. "Ted" Codd of IBM is known as the father of relational databases. He described a relational model for databases.

- **1974** – Structured Query Language appeared.

- **1978** – IBM worked to develop Codd's ideas and released a product named System/R.

- **1986** – IBM developed the first prototype of relational database and standardized by ANSI. The first relational database was released by Relational Software which later came to be known as Oracle.

# SQL Process

When you are executing an SQL command for any RDBMS, the system determines the best way to carry out your request and SQL engine figures out how to interpret the task.

There are various components included in this process.

These components are –

- Query Dispatcher
- Optimization Engines
- Classic Query Engine
- SQL Query Engine, etc.

A classic query engine handles all the non-SQL queries, but a SQL query engine won't handle logical files.

Following is a simple diagram showing the SQL Architecture –

# SQL Commands

The standard SQL commands to interact with relational databases are CREATE, SELECT, INSERT, UPDATE, DELETE and DROP. These commands can be classified into the following groups based on their nature −

## DDL - Data Definition Language

| Sr.No. | Command & Description |
|---|---|
| 1 | **CREATE**<br><br>Creates a new table, a view of a table, or other object in the database. |
| 2 | **ALTER**<br><br>Modifies an existing database object, such as a table. |
| 3 | **DROP**<br><br>Deletes an entire table, a view of a table or other objects in the database. |

## DML - Data Manipulation Language

| Sr.No. | Command & Description |
|--------|----------------------|
| 1 | **SELECT**<br><br>Retrieves certain records from one or more tables. |
| 2 | **INSERT**<br><br>Creates a record. |
| 3 | **UPDATE**<br><br>Modifies records. |
| 4 | **DELETE**<br><br>Deletes records. |

## DCL - Data Control Language

| Sr.No. | Command & Description |
|--------|----------------------|
| 1 | **GRANT**<br><br>Gives a privilege to user. |
| 2 | **REVOKE**<br><br>Takes back privileges granted from user. |

The SQL UNION clause/operator is used to combine the results of two or more SELECT statements without returning any duplicate rows.

To use this UNION clause, each SELECT statement must have

- The same number of columns selected
- The same number of column expressions
- The same data type and
- Have them in the same order

But they need not have to be in the same length.

# Syntax

The basic syntax of a **UNION** clause is as follows −

```
SELECT column1 [, column2 ]
FROM table1 [, table2 ]
[WHERE condition]


UNION


SELECT column1 [, column2 ]
FROM table1 [, table2 ]
[WHERE condition]
```

Here, the given condition could be any given expression based on your requirement.

# Example

Consider the following two tables.

**Table 1** − CUSTOMERS Table is as follows.

```
+----+----------+-----+-----------+----------+
| ID | NAME     | AGE | ADDRESS   | SALARY   |
+----+----------+-----+-----------+----------+
|  1 | Ramesh   |  32 | Ahmedabad |  2000.00 |
|  2 | Khilan   |  25 | Delhi     |  1500.00 |
|  3 | kaushik  |  23 | Kota      |  2000.00 |
|  4 | Chaitali |  25 | Mumbai    |  6500.00 |
|  5 | Hardik   |  27 | Bhopal    |  8500.00 |
```

```
|  6 | Komal    |  22 | MP        |  4500.00 |

|  7 | Muffy    |  24 | Indore    | 10000.00 |

+----+----------+-----+-----------+----------+
```

**Table 2** – ORDERS Table is as follows.

```
+-----+---------------------+-------------+--------+

|OID  | DATE                | CUSTOMER_ID | AMOUNT |

+-----+---------------------+-------------+--------+

| 102 | 2009-10-08 00:00:00 |           3 |   3000 |

| 100 | 2009-10-08 00:00:00 |           3 |   1500 |

| 101 | 2009-11-20 00:00:00 |           2 |   1560 |

| 103 | 2008-05-20 00:00:00 |           4 |   2060 |

+-----+---------------------+-------------+--------+
```

Now, let us join these two tables in our SELECT statement as follows –

```
SQL> SELECT  ID, NAME, AMOUNT, DATE

   FROM CUSTOMERS

   LEFT JOIN ORDERS

   ON CUSTOMERS.ID = ORDERS.CUSTOMER_ID

UNION

   SELECT  ID, NAME, AMOUNT, DATE

   FROM CUSTOMERS

   RIGHT JOIN ORDERS

   ON CUSTOMERS.ID = ORDERS.CUSTOMER_ID;
```

This would produce the following result –

```
+------+----------+--------+---------------------+

| ID   | NAME     | AMOUNT | DATE                |

+------+----------+--------+---------------------+

|    1 | Ramesh   |   NULL | NULL                |

|    2 | Khilan   |   1560 | 2009-11-20 00:00:00 |

|    3 | kaushik  |   3000 | 2009-10-08 00:00:00 |

|    3 | kaushik  |   1500 | 2009-10-08 00:00:00 |

|    4 | Chaitali |   2060 | 2008-05-20 00:00:00 |
```

```
|    5 | Hardik   |   NULL | NULL             |
|    6 | Komal    |   NULL | NULL             |
|    7 | Muffy    |   NULL | NULL             |
+------+----------+--------+------------------+
```

## The UNION ALL Clause

The UNION ALL operator is used to combine the results of two SELECT statements including duplicate rows.

The same rules that apply to the UNION clause will apply to the UNION ALL operator.

# Syntax

The basic syntax of the **UNION ALL** is as follows.

```
SELECT column1 [, column2 ]
FROM table1 [, table2 ]
[WHERE condition]


UNION ALL


SELECT column1 [, column2 ]
FROM table1 [, table2 ]
[WHERE condition]
```

Here, the given condition could be any given expression based on your requirement.

# Example

Consider the following two tables,

**Table 1** – CUSTOMERS Table is as follows.

```
+----+----------+-----+-----------+----------+
| ID | NAME     | AGE | ADDRESS   | SALARY   |
+----+----------+-----+-----------+----------+
|  1 | Ramesh   |  32 | Ahmedabad |  2000.00 |
|  2 | Khilan   |  25 | Delhi     |  1500.00 |
|  3 | kaushik  |  23 | Kota      |  2000.00 |
|  4 | Chaitali |  25 | Mumbai    |  6500.00 |
|  5 | Hardik   |  27 | Bhopal    |  8500.00 |
|  6 | Komal    |  22 | MP        |  4500.00 |
```

```
|  7 | Muffy    |  24 | Indore   | 10000.00 |

+----+----------+-----+----------+----------+
```

**Table 2** − ORDERS table is as follows.

```
+-----+---------------------+-------------+--------+

|OID  | DATE                | CUSTOMER_ID | AMOUNT |

+-----+---------------------+-------------+--------+

| 102 | 2009-10-08 00:00:00 |           3 |   3000 |

| 100 | 2009-10-08 00:00:00 |           3 |   1500 |

| 101 | 2009-11-20 00:00:00 |           2 |   1560 |

| 103 | 2008-05-20 00:00:00 |           4 |   2060 |

+-----+---------------------+-------------+--------+
```

Now, let us join these two tables in our SELECT statement as follows −

```
SQL> SELECT ID, NAME, AMOUNT, DATE

   FROM CUSTOMERS

   LEFT JOIN ORDERS

   ON CUSTOMERS.ID = ORDERS.CUSTOMER_ID

UNION ALL

   SELECT ID, NAME, AMOUNT, DATE

   FROM CUSTOMERS

   RIGHT JOIN ORDERS

   ON CUSTOMERS.ID = ORDERS.CUSTOMER_ID;
```

This would produce the following result −

```
+------+----------+--------+---------------------+
| ID   | NAME     | AMOUNT | DATE                |
+------+----------+--------+---------------------+
|    1 | Ramesh   |   NULL | NULL                |
|    2 | Khilan   |   1560 | 2009-11-20 00:00:00 |
|    3 | kaushik  |   3000 | 2009-10-08 00:00:00 |
|    3 | kaushik  |   1500 | 2009-10-08 00:00:00 |
|    4 | Chaitali |   2060 | 2008-05-20 00:00:00 |
|    5 | Hardik   |   NULL | NULL
```

```
|   6 | Komal    |  NULL | NULL                |
|   7 | Muffy    |  NULL | NULL                |
|   3 | kaushik  |  3000 | 2009-10-08 00:00:00 |
|   3 | kaushik  |  1500 | 2009-10-08 00:00:00 |
|   2 | Khilan   |  1560 | 2009-11-20 00:00:00 |
|   4 | Chaitali |  2060 | 2008-05-20 00:00:00 |
+------+----------+--------+---------------------+
```

There are two other clauses (i.e., operators), which are like the UNION clause.

- SQL INTERSECT Clause – This is used to combine two SELECT statements, but returns rows only from the first SELECT statement that are identical to a row in the second SELECT statement.

- SQL EXCEPT Clause – This combines two SELECT statements and returns rows from the first SELECT statement that are not returned by the second SELECT statement.

**INTERSECT OPERATION:-**

## Description

The SQL INTERSECT operator is used to return the results of 2 or more SELECT statements. However, it only returns the rows selected by all queries or data sets. If a record exists in one query and not in the other, it will be omitted from the INTERSECT results.

**Explanation:** The INTERSECT query will return the records in the blue shaded area. These are the records that exist in both Dataset1 and Dataset2.

Each SQL statement within the SQL INTERSECT must have the same number of fields in the result sets with similar data types.

## Syntax

The syntax for the INTERSECT operator in SQL is:

```
SELECT expression1, expression2, ... expression_n

FROM tables

[WHERE conditions]

INTERSECT

SELECT expression1, expression2, ... expression_n

FROM tables

[WHERE conditions];
```

Parameters or Arguments

**expression1, expression2, expression_n**

The columns or calculations that you wish to retrieve.

**tables**

The tables that you wish to retrieve records from. There must be at least one table listed in the FROM clause.

**WHERE conditions**

Optional. These are conditions that must be met for the records to be selected.

## Note

- There must be same number of expressions in both SELECT statements.
- The corresponding expressions must have the same data type in the SELECT statements. For example: *expression1* must be the same data type in both the first and second SELECT statement.

**Except Operator:-**

## Description

The SQL Server (Transact-SQL) EXCEPT operator is used to return all rows in the first SELECT statement that are not returned by the second SELECT statement. Each SELECT statement will define a dataset. The EXCEPT operator will retrieve all records from the first dataset and then remove from the results all records from the second dataset.

**Except Query:-**



TechOnTheNet.com

**Explanation:** The EXCEPT query will return the records in the blue shaded area. These are the records that exist in Dataset1 and not in Dataset2.

Each SELECT statement within the EXCEPT query must have the same number of fields in the result sets with similar data types.

## Syntax

The syntax for the EXCEPT operator in SQL Server (Transact-SQL) is:

```
SELECT expression1, expression2, ... expression_n

FROM tables

[WHERE conditions]

EXCEPT

SELECT expression1, expression2, ... expression_n

FROM tables

[WHERE conditions];
```

Parameters or Arguments

**expressions**

The columns or calculations that you wish to compare between the two SELECT statements. They do not have to be the same fields in each of the SELECT statements, but the corresponding columns must be similar data types.

**tables**

The tables that you wish to retrieve records from. There must be at least one table listed in the FROM clause.

**WHERE conditions**
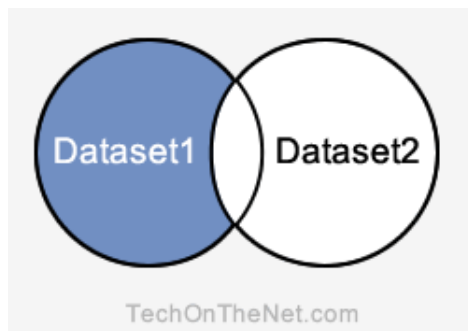
Optional. The conditions that must be met for the records to be selected.

# Note

- There must be same number of expressions in both SELECT statements.
- The corresponding columns in each of the SELECT statements must have similar data types.
- The EXCEPT operator returns all records from the first SELECT statement that are not in the second SELECT statement.
- The EXCEPT operator in SQL Server is equivalent to the MINUS operator in Oracle.

**Nested Queries:-**

A Subquery or Inner query or a Nested query is a query within another SQL query and embedded within the WHERE clause.

A subquery is used to return data that will be used in the main query as a condition to further restrict the data to be retrieved.

Subqueries can be used with the SELECT, INSERT, UPDATE, and DELETE statements along with the operators like =, <, >, >=, <=, IN, BETWEEN, etc.

There are a few rules that subqueries must follow −

- Subqueries must be enclosed within parentheses.

- A subquery can have only one column in the SELECT clause, unless multiple columns are in the main query for the subquery to compare its selected columns.

- An ORDER BY command cannot be used in a subquery, although the main query can use an ORDER BY. The GROUP BY command can be used to perform the same function as the ORDER BY in a subquery.

- Subqueries that return more than one row can only be used with multiple value operators such as the IN operator.

- The SELECT list cannot include any references to values that evaluate to a BLOB, ARRAY, CLOB, or NCLOB.

- A subquery cannot be immediately enclosed in a set function.

- The BETWEEN operator cannot be used with a subquery. However, the BETWEEN operator can be used within the subquery.

## Subqueries with the SELECT Statement

Subqueries are most frequently used with the SELECT statement. The basic syntax is as follows —

```
SELECT column_name [, column_name ]
FROM   table1 [, table2 ]
WHERE  column_name OPERATOR
   (SELECT column_name [, column_name ]
   FROM table1 [, table2 ]
   [WHERE])
```

# Example

Consider the CUSTOMERS table having the following records —

```
+----+----------+-----+-----------+----------+
| ID | NAME     | AGE | ADDRESS   | SALARY   |
+----+----------+-----+-----------+----------+
|  1 | Ramesh   |  35 | Ahmedabad |  2000.00 |
|  2 | Khilan   |  25 | Delhi     |  1500.00 |
|  3 | kaushik  |  23 | Kota      |  2000.00 |
|  4 | Chaitali |  25 | Mumbai    |  6500.00 |
|  5 | Hardik   |  27 | Bhopal    |  8500.00 |
|  6 | Komal    |  22 | MP        |  4500.00 |
|  7 | Muffy    |  24 | Indore    | 10000.00 |
+----+----------+-----+-----------+----------+
```

Now, let us check the following subquery with a SELECT statement.

```
SQL> SELECT *
   FROM CUSTOMERS
   WHERE ID IN (SELECT ID
        FROM CUSTOMERS
        WHERE SALARY > 4500) ;
```

This would produce the following result.

```
+----+----------+-----+---------+----------+
| ID | NAME     | AGE | ADDRESS | SALARY   |
+----+----------+-----+---------+----------+
|  4 | Chaitali |  25 | Mumbai  |  6500.00 |
|  5 | Hardik   |  27 | Bhopal  |  8500.00 |
|  7 | Muffy    |  24 | Indore  | 10000.00 |
+----+----------+-----+---------+----------+
```

# Subqueries with the INSERT Statement

Subqueries also can be used with INSERT statements. The INSERT statement uses the data returned from the subquery to insert into another table. The selected data in the subquery can be modified with any of the character, date or number functions.

The basic syntax is as follows.

```
INSERT INTO table_name [ (column1 [, column2 ]) ]
   SELECT [ *|column1 [, column2 ]
   FROM table1 [, table2 ]
   [ WHERE VALUE OPERATOR ]
```

## Example

Consider a table CUSTOMERS_BKP with similar structure as CUSTOMERS table. Now to copy the complete CUSTOMERS table into the CUSTOMERS_BKP table, you can use the following syntax.

```
SQL> INSERT INTO CUSTOMERS_BKP
   SELECT * FROM CUSTOMERS
   WHERE ID IN (SELECT ID
   FROM CUSTOMERS) ;
```

**Subqueries with the UPDATE Statement:-**

The subquery can be used in conjunction with the UPDATE statement. Either single or multiple columns in a table can be updated when using a subquery with the UPDATE statement.

The basic syntax is as follows.

```
UPDATE table
SET column_name = new_value
[ WHERE OPERATOR [ VALUE ]
   (SELECT COLUMN_NAME
   FROM TABLE_NAME)
   [ WHERE) ]
```

# Example

Assuming, we have CUSTOMERS_BKP table available which is backup of CUSTOMERS table. The following example updates SALARY by 0.25 times in the CUSTOMERS table for all the customers whose AGE is greater than or equal to 27.

```
SQL> UPDATE CUSTOMERS
   SET SALARY = SALARY * 0.25
   WHERE AGE IN (SELECT AGE FROM CUSTOMERS_BKP
     WHERE AGE >= 27 );
```

This would impact two rows and finally CUSTOMERS table would have the following records.

```
+----+----------+-----+-----------+----------+
| ID | NAME     | AGE | ADDRESS   | SALARY   |
+----+----------+-----+-----------+----------+
|  1 | Ramesh   |  35 | Ahmedabad |   125.00 |
|  2 | Khilan   |  25 | Delhi     |  1500.00 |
|  3 | kaushik  |  23 | Kota      |  2000.00 |
|  4 | Chaitali |  25 | Mumbai    |  6500.00 |
|  5 | Hardik   |  27 | Bhopal    |  2125.00 |
|  6 | Komal    |  22 | MP        |  4500.00 |
|  7 | Muffy    |  24 | Indore    | 10000.00 |
+----+----------+-----+-----------+----------+
```

**Subqueries with the DELETE Statement:-**

The subquery can be used in conjunction with the DELETE statement like with any other statements mentioned above.

The basic syntax is as follows.

```
DELETE FROM TABLE_NAME
[ WHERE OPERATOR [ VALUE ]
   (SELECT COLUMN_NAME
   FROM TABLE_NAME)
   [ WHERE) ]
```

# Example

Assuming, we have a CUSTOMERS_BKP table available which is a backup of the CUSTOMERS table. The following example deletes the records from the CUSTOMERS table for all the customers whose AGE is greater than or equal to 27.

```
SQL> DELETE FROM CUSTOMERS
    WHERE AGE IN (SELECT AGE FROM CUSTOMERS_BKP
      WHERE AGE >= 27 );
```

This would impact two rows and finally the CUSTOMERS table would have the following records.

```
+----+----------+-----+---------+----------+
| ID | NAME     | AGE | ADDRESS | SALARY   |
+----+----------+-----+---------+----------+
|  2 | Khilan   |  25 | Delhi   |  1500.00 |
|  3 | kaushik  |  23 | Kota    |  2000.00 |
|  4 | Chaitali |  25 | Mumbai  |  6500.00 |
|  6 | Komal    |  22 | MP      |  4500.00 |
|  7 | Muffy    |  24 | Indore  | 10000.00 |
+----+----------+-----+---------+----------+
```

### Aggregate Operators:-

An aggregate function allows you to perform a calculation on a set of values to return a single scalar value. We often use aggregate functions with the GROUP BY and HAVING clauses of the SELECT statement.

The following are the most commonly used SQL aggregate functions:

- AVG – calculates the average of a set of values.
- COUNT – counts rows in a specified table or view.
- MIN – gets the minimum value in a set of values.
- MAX – gets the maximum value in a set of values.
- SUM – calculates the sum of values.

Notice that all aggregate functions above ignore NULL values except for the COUNT function.

## SQL aggregate functions syntax

To call an aggregate function, you use the following syntax:

```
1   aggregate_function (DISTINCT | ALL expression)
```

Let's examine the syntax above in greater detail:

- First, specify an aggregate function that you want to use e.g., MIN, MAX, AVG, SUM or COUNT.
- Second, put DISTINCT or ALL modifier followed by an expression inside parentheses. If you explicitly use DISTINCT modifier, the aggregate function ignores duplicate values and only consider the unique values. If you use the ALL modifier, the aggregate function uses all values for calculation or evaluation. The ALL modifier is used by default if you do not specify any modifier explicitly.

## SQL aggregate function examples

Let's take a look some examples of using SQL aggregate functions.

## COUNT function example

To get the number of the products in the products table, you use the COUNT function as follows:

```
1 SELECT

2   COUNT(*)

3 FROM

4   products;
```

| COUNT(*) |
|----------|
| 77 |

More information on the COUNT function.

## AVG function example

To calculate the average units in stock of the products, you use the AVG function as follows:

1 SELECT

2    AVG(unitsinstock)

3 FROM

4    products;

| AVG(unitsinstock) |
|-------------------|
| 40.5065 |

To calculate units in stock by product category, you use the AVG function with the GROUP BYclause as follows:

1 SELECT

2    categoryid, AVG(unitsinstock)

3 FROM

4    products

5 GROUP BY categoryid;

| categoryid | AVG(unitsinstock) |
|------------|-------------------|
| 1 | 46.5833 |
| 2 | 42.2500 |
| 3 | 29.6923 |
| 4 | 39.3000 |
| 5 | 44.0000 |
| 6 | 27.5000 |
| 7 | 20.0000 |
| 8 | 58.4167 |

More information on AVG function.

## SUM function example

To calculate the sum of units in stock by product category, you use the SUM function with the GROUP BY clause as the following query:

1 SELECT

2    categoryid, SUM(unitsinstock)

3 FROM

4    products

5 GROUP BY categoryid;

| categoryid | sum(unitsinstock) |
|---|---|
| 1 | 559 |
| 2 | 507 |
| 3 | 386 |
| 4 | 393 |
| 5 | 308 |
| 6 | 165 |
| 7 | 100 |
| 8 | 701 |

Check it out the SUM function tutorial for more information on how to use the SUM function.

## MIN function example

To get the minimum units in stock of products in the products table, you use the MIN function as follows:

1 SELECT

2    MIN(unitsinstock)

3 FROM

4    products;

| MIN(unitsinstock) |
|---|
| 0 |

More information on the MIN function.

**MAX function example**

To get the maximum units in stock of products in the products table, you use the MAX function as shown in the following query:

```
1 SELECT
2    MAX(unitsinstock)
3 FROM
4    products;
```

| | Max(unitsinstock) |
|---|---|
| ▶ | 125 |

Check it out the MAX function tutorial for more information.

In this tutorial, we have introduced you to the SQL aggregate functions including the most commonly used functions: AVG, COUNT, MIN, MAX, and SUM.

**NULL VALUES:-**

The SQL **NULL** is the term used to represent a missing value. A NULL value in a table is a value in a field that appears to be blank.

A field with a NULL value is a field with no value. It is very important to understand that a NULL value is different than a zero value or a field that contains spaces.

# Syntax

The basic syntax of **NULL** while creating a table.

```
SQL> CREATE TABLE CUSTOMERS(
   ID   INT            NOT NULL,
   NAME VARCHAR (20)    NOT NULL,
   AGE  INT            NOT NULL,
   ADDRESS  CHAR (25) ,
   SALARY   DECIMAL (18, 2),
   PRIMARY KEY (ID)
);
```

Here, **NOT NULL** signifies that column should always accept an explicit value of the given data type. There are two columns where we did not use NOT NULL, which means these columns could be NULL.

A field with a NULL value is the one that has been left blank during the record creation.

# Example

The NULL value can cause problems when selecting data. However, because when comparing an unknown value to any other value, the result is always unknown and not included in the results. You must use the **IS NULL** or **IS NOT NULL** operators to check for a NULL value.

Consider the following CUSTOMERS table having the records as shown below.

```
+----+----------+-----+-----------+----------+
| ID | NAME     | AGE | ADDRESS   | SALARY   |
+----+----------+-----+-----------+----------+
|  1 | Ramesh   |  32 | Ahmedabad |  2000.00 |
|  2 | Khilan   |  25 | Delhi     |  1500.00 |
|  3 | kaushik  |  23 | Kota      |  2000.00 |
|  4 | Chaitali |  25 | Mumbai    |  6500.00 |
|  5 | Hardik   |  27 | Bhopal    |  8500.00 |
|  6 | Komal    |  22 | MP        |          |
|  7 | Muffy    |  24 | Indore    |          |
+----+----------+-----+-----------+----------+
```

Now, following is the usage of the **IS NOT NULL** operator.

```
SQL> SELECT  ID, NAME, AGE, ADDRESS, SALARY

   FROM CUSTOMERS

   WHERE SALARY IS NOT NULL;
```

This would produce the following result −

```
+----+----------+-----+-----------+----------+
| ID | NAME     | AGE | ADDRESS   | SALARY   |
+----+----------+-----+-----------+----------+
|  1 | Ramesh   |  32 | Ahmedabad |  2000.00 |
|  2 | Khilan   |  25 | Delhi     |  1500.00 |
|  3 | kaushik  |  23 | Kota      |  2000.00 |
|  4 | Chaitali |  25 | Mumbai    |  6500.00 |
|  5 | Hardik   |  27 | Bhopal    |  8500.00 |
+----+----------+-----+-----------+----------+
```

Now, following is the usage of the **IS NULL** operator.

```
SQL> SELECT  ID, NAME, AGE, ADDRESS, SALARY

   FROM CUSTOMERS

   WHERE SALARY IS NULL;
```

This would produce the following result −

```
+----+----------+-----+----------+----------+
| ID | NAME     | AGE | ADDRESS  | SALARY   |
+----+----------+-----+----------+----------+
|  6 | Komal    |  22 | MP       |          |
|  7 | Muffy    |  24 | Indore   |          |
+----+----------+-----+----------+----------+
```

**Complex Integrity Constraints in SQL:-**

Integrity Constraints are used to apply business rules for the database tables.

The constraints available in SQL are **Foreign Key**,**Not Null**, **Unique**, **Check**.

Constraints can be defined in two ways

1) The constraints can be specified immediately after the column definition. This is called column-level definition.

2) The constraints can be specified after all the columns are defined. This is called table-level definition.

# 1) SQL Primary key:

This constraint defines a column or combination of columns which uniquely identifies each row in the table.

**Syntax to define a Primary key at column level:**

```
column name datatype [CONSTRAINT constraint_name] PRIMARY KEY
```

**Syntax to define a Primary key at table level:**

```
[CONSTRAINT          constraint_name]          PRIMARY          KEY
(column_name1,column_name2,..)
```

- **column_name1, column_name2** are the names of the columns which define the primary Key.
- The syntax within the bracket i.e. [CONSTRAINT constraint_name] is optional.
  - **For Example:** To create an employee table with Primary Key constraint, the query would be like.
  - **Primary Key at column level:**
  -
    ```
    CREATE                    TABLE                    employee

    (         id         number(5)      PRIMARY          KEY,

    name                                            char(20),

    dept                                            char(10),

    age                                             number(2),

    salary                                          number(10),

    location                                        char(10)

    );
    ```

  - or

    ```
    CREATE                    TABLE                    employee

    (  id   number(5)   CONSTRAINT   emp_id_pk   PRIMARY   KEY,

    name                                            char(20),

    dept                                            char(10),

    age                                             number(2),

    salary                                          number(10),

    location                                        char(10)

    );
    ```

**Primary Key at column level:**

- ```
  CREATE                    TABLE                    employee

  (                    id                    number(5),

  name                                        char(20),

  dept                                        char(10),

  age                                         number(2),

  salary                                      number(10),

  location                                    char(10),

  CONSTRAINT      emp_id_pk      PRIMARY      KEY      (id)

  );
  ```

**Primary Key at table level:**

```
CREATE                    TABLE                    employee

(         id        number(5),        NOT         NULL,

name                                        char(20),

dept                                        char(10),

age                                         number(2),

salary                                      number(10),

location                                    char(10),

ALTER TABLE employee ADD CONSTRAINT PK_EMPLOYEE_ID PRIMARY

KEY                                                (id)

);
```

**Foreign key:-**

A FOREIGN KEY is a key used to link two tables together.

A FOREIGN KEY is a field (or collection of fields) in one table that refers to the PRIMARY KEY in another table.

The table containing the foreign key is called the child table, and the table containing the candidate key is called the referenced or parent table.

Look at the following two tables:

"Persons" table:

| PersonID | LastName | FirstName | Age |
|----------|----------|-----------|-----|
| 1 | Hansen | Ola | 30 |
| 2 | Svendson | Tove | 23 |
| 3 | Pettersen | Kari | 20 |

"Orders" table:

| OrderID | OrderNumber | PersonID |
|---------|-------------|----------|
| 1 | 77895 | 3 |
| 2 | 44678 | 3 |

| 3 | 22456 | 2 |
|---|-------|---|
| 4 | 24562 | 1 |

Notice that the "PersonID" column in the "Orders" table points to the "PersonID" column in the "Persons" table.

The "PersonID" column in the "Persons" table is the PRIMARY KEY in the "Persons" table.

The "PersonID" column in the "Orders" table is a FOREIGN KEY in the "Orders" table.

The FOREIGN KEY constraint is used to prevent actions that would destroy links between tables.

The FOREIGN KEY constraint also prevents invalid data from being inserted into the foreign key column, because it has to be one of the values contained in the table it points to.

Active Database Concepts and Triggers:-

➢ Rules that specify actions that are automatically triggered by certain events have been considered important enhancements to database systems for quite some time. In fact, the concept of **triggers**—a technique for specifying certain types of active rules—has existed in early versions of the SQL specification for relational databases and triggers are now part of the SQL-99 and later standards.

**Generalized Model for Active Databases and Oracle Triggers:-**

The model that has been used to specify active database rules is referred to as the **Event-Condition-Action** (**ECA)** model. A rule in the ECA model has three components:

➢  1. The **event(s)** that triggers the rule: These events are usually database update operations that are explicitly applied to the database. However, in the general model, they could also be temporal events[2] or other kinds of external events.

2. The **condition** that determines whether the rule action should be executed: Once the triggering event has occurred, an *optional* condition may be evaluated. If *no condition* is specified, the action will be executed once the event occurs. If a condition is specified, it is first evaluated, and only *if it evaluates to true* will the rule action be executed.

➢  3. The **action** to be taken: The action is usually a sequence of SQL statements, but it could also be a database transaction or an external program that will be automatically executed.

**SYNTAX:-**

```
-- Trigger on an INSERT, UPDATE, or DELETE statement to a table or view (DML
Trigger)
```

A syntax summary for specifying triggers in the Oracle system (main options only).

```
<trigger>        ::=     CREATE TRIGGER <trigger name>

                 ( AFTER I BEFORE ) <triggering events> ON <table name>

                 [ FOR EACH ROW ]

                 [ WHEN <condition> ]

                 <trigger actions> ;

<triggering events>  ::= <trigger event> {OR <trigger event> }

<trigger event>              ::= INSERT I DELETE I UPDATE [ OF <column name> { ,
<column name> } ]

<trigger action>             ::= <PL/SQL block>
```

## Advantages

Followings are the advantages of using MySQL triggers:

1.  **Integrity of data:** With the help of MySQL trigger we can check the integrity of data in the table. In other words, MySQL triggers are the alternative way to check the integrity of data.
2.  **Useful for catching errors:** MySQL triggers can catch errors in business logic in the database layer.
3.  **Alternative way to run scheduled tasks:** Actually by using MySQL triggers we do not have to wait to run the scheduled tasks because the triggers are invoked automatically 'before' or 'after' a modification is done to the data in the table.
4.  **Auditing:** Actually MySQL triggers are very much useful for the purpose of auditing of the changes made in the table.
5.  **Prevention of invalid transactions:** MySQL triggers are very useful in the prevention of invalid transactions.
6.  **Logging of event:** MySQL triggers can log an event and can also store the information on the access of table.

## Disadvantages

Followings are the disadvantages of using MySQL triggers:

1.  **Cannot replace all validations:** Actually, MySQL triggers cannot replace all the validations and can only provide an extended validation.
2.  **Invisible from client applications:** Basically MySQL triggers are invoked and executed invisible from the client applications hence it is very much difficult to figure out what happens in the database layer.
3.  **Impose load on server:** Triggers can impose a high load on the database server.
4.  **Not recommended for high velocity of data:** Triggers are not beneficial for use with high-velocity data i.e. the data when a number of events per second are high. It is because in case of high-velocity data the triggers get triggered all the time.

## UNIT-IV

**Syllabus:**

**Schema Refinement (Normalization) :** Purpose of Normalization or schema refinement, concept of functional dependency, normal forms based on functional dependency(1NF, 2NF and 3 NF), concept of surrogate key, Boyce-codd normal form(BCNF), Lossless join and dependency preserving decomposition, Fourth normal form(4NF).

### Introduction

When designing a database, you have to make decisions regarding how best to take some system in the real world and model it in a database. This consists of deciding which tables to create, what columns they will contain, as well as the relationships between the tables. While it would be nice if this process was totally intuitive and obvious, or even better automated, this is simply not the case. A well-designed database takes time and effort to conceive, build and refine.

The benefits of a database that has been designed according to the relational model are numerous. Some of them are:

• Data entry, updates and deletions will be efficient.

• Data retrieval, summarization and reporting will also be efficient.

• Since the database follows a well-formulated model, it behaves predictably.

• Since much of the information is stored in the database rather than in the application, the database is somewhat self-documenting.

• Changes to the database schema are easy to make.

### Purpose of Normalization:-

If a database design is not perfect, it may contain anomalies, which are like a bad dream for any database administrator. Managing a database with anomalies is next to impossible.

- **Update anomalies** − If data items are scattered and are not linked to each other properly, then it could lead to strange situations. For example, when we try to update one data item having its copies scattered over several places, a few instances get updated properly while a few others are left with old values. Such instances leave the database in an inconsistent state.

- **Deletion anomalies** − We tried to delete a record, but parts of it was left undeleted because of unawareness, the data is also saved somewhere else.

- **Insert anomalies** − We tried to insert data in a record that does not exist at all.

Normalization is a method to remove all these anomalies and bring the database to a consistent state.

### Concept of Functional Dependency:-

Functional dependency (FD) is a set of constraints between two attributes in a relation. Functional dependency says that if two tuples have same values for attributes A1, A2,..., An, then those two tuples must have to have same values for attributes B1, B2, ..., Bn.

Functional dependency is represented by an arrow sign ($\rightarrow$) that is, X$\rightarrow$Y, where X functionally determines Y. The left-hand side attributes determine the values of attributes on the right-hand side.

**Normal Forms based on Functional Dependencies(1NF,2NF and 3 NF):-**

**First Normal Form**

First Normal Form is defined in the definition of relations (tables) itself. This rule defines that all the attributes in a relation must have atomic domains. The values in an atomic domain are indivisible units.

| Course | Content |
|--------|---------|
| Programming | Java, c++ |
| Web | HTML, PHP, ASP |

We re-arrange the relation (table) as below, to convert it to First Normal Form.

| Course | Content |
|--------|---------|
| Programming | Java |
| Programming | c++ |
| Web | HTML |
| Web | PHP |
| Web | ASP |

Each attribute must contain only a single value from its pre-defined domain.

**Second Normal Form**

Before we learn about the second normal form, we need to understand the following −

- **Prime attribute** − An attribute, which is a part of the candidate-key, is known as a prime attribute.

- **Non-prime attribute** − An attribute, which is not a part of the prime-key, is said to be a non-prime attribute.

If we follow second normal form, then every non-prime attribute should be fully functionally dependent on prime key attribute. That is, if X → A holds, then there should not be any proper subset Y of X, for which Y → A also holds true.

## Student_Project

| Stu_ID | Proj_ID | Stu_Name | Proj_Name |
|--------|---------|----------|-----------|

We see here in Student_Project relation that the prime key attributes are Stu_ID and Proj_ID. According to the rule, non-key attributes, i.e. Stu_Name and Proj_Name must be dependent upon both and not on any of the prime key attribute individually. But we find that Stu_Name can be identified by Stu_ID and Proj_Name can be identified by Proj_ID independently. This is called **partial dependency**, which is not allowed in Second Normal Form.

## Student

| Stu_ID | Stu_Name | Proj_ID |
|--------|----------|---------|

## Project

| Proj_ID | Proj_Name |
|---------|-----------|

We broke the relation in two as depicted in the above picture. So there exists no partial dependency.

### Third Normal Form

For a relation to be in Third Normal Form, it must be in Second Normal form and the following must satisfy –

- No non-prime attribute is transitively dependent on prime key attribute.
- For any non-trivial functional dependency, $X \rightarrow A$, then either –
  - X is a superkey or,
  - A is prime attribute.

## Student_Detail

| Stu_ID | Stu_Name | City | Zip |
|--------|----------|------|-----|

We find that in the above Student_detail relation, Stu_ID is the key and only prime key attribute. We find that City can be identified by Stu_ID as well as Zip itself. Neither Zip is a superkey nor is City a prime attribute. Additionally, Stu_ID → Zip → City, so there exists **transitive dependency**.

To bring this relation into third normal form, we break the relation into two relations as follows −

**Student_Detail**

| Stu_ID | Stu_Name | Zip |
|--------|----------|-----|

**ZipCodes**

| Zip | City |
|-----|------|

### Concept of a surrogate key in DBMS:-

A surrogate key is any column or set of columns that can be declared as the **primary key** instead of a "real" or natural key. Sometimes there can be several natural keys that could be declared as the **primary key**, and these are all called candidate keys. So a surrogate is a **candidate key**.

A **candidate key** is a column, or set of columns, in a table that can uniquely identify any **database** record without referring to any other data. Each table may have one or more **candidate keys**, but one **candidate key** is unique, and it is called the **primary key**.

### Boyce-Codd normal form(BCNF):-

Boyce-Codd Normal Form (BCNF) is an extension of Third Normal Form on strict terms. BCNF states that −

- For any non-trivial functional dependency, X → A, X must be a super-key.

**Student_Detail**

| Stu_ID | Stu_Name | Zip |
|--------|----------|-----|

**ZipCodes**

| Zip | City |
|-----|------|

In the above image, Stu_ID is the super-key in the relation Student_Detail and Zip is the super-key in the relation ZipCodes. So,

Stu_ID → Stu_Name, Zip

and

Zip → City

Which confirms that both the relations are in BCNF.

<u>**Lossless join and dependency preserving decomposition:-**</u>

## What is decomposition?

- Decomposition is the process of breaking down in parts or elements.
- It replaces a relation with a collection of smaller relations.
- It breaks the table into multiple tables in a database.
- It should always be lossless, because it confirms that the information in the original relation can be accurately reconstructed based on the decomposed relations.
- If there is no proper decomposition of the relation, then it may lead to problems like loss of information.

## Properties of Decomposition

**Following are the properties of Decomposition,**
1. Lossless Decomposition
2. Dependency Preservation
3. Lack of Data Redundancy.

### 1. Lossless Decomposition

- Decomposition must be lossless. It means that the information should not get lost from the relation that is decomposed.
- It gives a guarantee that the join will result in the same relation as it was decomposed.

**Example:**
Let's take 'E' is the Relational Schema, With instance 'e'; is decomposed into: E1, E2, E3, . . . . En; With instance: e1, e2, e3, . . . . en, If e1 ⋈ e2 ⋈ e3 . . . . ⋈ en, then it is called as **'Lossless Join Decomposition'.**

- In the above example, it means that, if natural joins of all the decomposition give the original relation, then it is said to be lossless join decomposition.

## Example: <Employee_Department> Table

| Eid | Ename | Age | City | Salary | Deptid | DeptName |
|-----|-------|-----|------|--------|--------|----------|
| E001 | ABC | 29 | Pune | 20000 | D001 | Finance |
| E002 | PQR | 30 | Pune | 30000 | D002 | Production |
| E003 | LMN | 25 | Mumbai | 5000 | D003 | Sales |
| E004 | XYZ | 24 | Mumbai | 4000 | D004 | Marketing |

| E005 | STU | 32 | Bangalore | 25000 | D005 | Human Resource |

- Decompose the above relation into two relations to check whether a decomposition is lossless or lossy.

- Now, we have decomposed the relation that is Employee and Department.

### Relation 1 : <Employee> Table

| Eid | Ename | Age | City | Salary |
|-----|-------|-----|------|--------|
| E001 | ABC | 29 | Pune | 20000 |
| E002 | PQR | 30 | Pune | 30000 |
| E003 | LMN | 25 | Mumbai | 5000 |
| E004 | XYZ | 24 | Mumbai | 4000 |
| E005 | STU | 32 | Bangalore | 25000 |

- Employee Schema contains (Eid, Ename, Age, City, Salary).

### Relation 2 : <Department> Table

| Deptid | Eid | DeptName |
|--------|-----|----------|
| D001 | E001 | Finance |
| D002 | E002 | Production |
| D003 | E003 | Sales |
| D004 | E004 | Marketing |
| D005 | E005 | Human Resource |

- Department Schema contains (Deptid, Eid, DeptName).

- So, the above decomposition is a Lossless Join Decomposition, because the two relations contains one common field that is 'Eid' and therefore join is possible.

- Now apply natural join on the decomposed relations.

### Employee ⋈ Department

| Eid | Ename | Age | City | Salary | Deptid | DeptName |
|-----|-------|-----|------|--------|--------|----------|
| E001 | ABC | 29 | Pune | 20000 | D001 | Finance |
| E002 | PQR | 30 | Pune | 30000 | D002 | Production |
| E003 | LMN | 25 | Mumbai | 5000 | D003 | Sales |
| E004 | XYZ | 24 | Mumbai | 4000 | D004 | Marketing |

| E005 | STU | 32 | Bangalore | 25000 | D005 | Human Resource |

Hence, the decomposition is Lossless Join Decomposition.

- If the <Employee> table contains (Eid, Ename, Age, City, Salary) and <Department> table contains (Deptid and DeptName), then it is not possible to join the two tables or relations, because there is no common column between them. And it becomes **Lossy Join Decomposition.**

### Dependency Preservation

Dependency is an important constraint on the database.

- Every dependency must be satisfied by at least one decomposed table.
- If {A → B} holds, then two sets are functional dependent. And, it becomes more useful for checking the dependency easily if both sets in a same relation.
  - This decomposition property can only be done by maintaining the functional dependency.
  - In this property, it allows to check the updates without computing the natural join of the database structure.

### Fourth Normal Form(4NF):-

Fourth Normal Form comes into picture when **Multi-valued Dependency** occur in any relation. In this tutorial we will learn about Multi-valued Dependency, how to remove it and how to make any table satisfy the fourth normal form.

### Rules for 4th Normal Form

For a table to satisfy the Fourth Normal Form, it should satisfy the following two conditions:

1.  It should be in the **Boyce-Codd Normal Form**.
2.  And, the table should not have any **Multi-valued Dependency**.

### What is Multi-valued Dependency?

A table is said to have multi-valued dependency, if the following conditions are true,

1.  For a dependency A → B, if for a single value of A, multiple value of B exists, then the table may have multi-valued dependency.
2.  Also, a table should have at-least 3 columns for it to have a multi-valued dependency.
3.  And, for a relation R(A,B,C), if there is a multi-valued dependency between, A and B, then B and C should be independent of each other.

If all these conditions are true for any relation(table), it is said to have multi-valued dependency.

### Example

Below we have a college enrolment table with columns `s_id`, `course` and `hobby`.

| s_id | course | hobby |
|------|--------|-------|
| 1 | Science | Cricket |
| 1 | Maths | Hockey |
| 2 | C# | Cricket |
| 2 | Php | Hockey |

As you can see in the table above, student with `s_id` **1** has opted for two courses, **Science** and **Maths**, and has two hobbies, **Cricket** and **Hockey**.

You must be thinking what problem this can lead to, right?

Well the two records for student with `s_id` **1**, will give rise to two more records, as shown below, because for one student, two hobbies exists, hence along with both the courses, these hobbies should be specified.

| s_id | course | hobby |
|------|--------|-------|
| 1 | Science | Cricket |
| 1 | Maths | Hockey |
| 1 | Science | Hockey |
| 1 | Maths | Cricket |

And, in the table above, there is no relationship between the columns `course` and `hobby`. They are independent of each other.

So there is multi-value dependency, which leads to un-necessary repetition of data and other anomalies as well.

**How to satisfy 4th Normal Form?**

To make the above relation satify the 4th normal form, we can decompose the table into 2 tables.

**CourseOpted Table**

| s_id | course |
|------|--------|
| 1 | Science |
| 1 | Maths |
| 2 | C# |
| 2 | Php |

And, **Hobbies Table**,

| s_id | hobby |
|------|-------|
| 1 | Cricket |
| 1 | Hockey |
| 2 | Cricket |
| 2 | Hockey |

Now this relation satisfies the fourth normal form.

A table can also have functional dependency along with multi-valued dependency. In that case, the functionally dependent columns are moved in a separate table and the multi-valued dependent columns are moved to separate tables.

If you design your database carefully, you can easily avoid these issues.

# UNIT – V

**SYLLABUS**

**Transaction Management and Concurrency Control:**
Transaction, properties of transactions, transaction log, and transaction management with SQL using commit rollback and save point. Concurrency control for lost updates, uncommitted data, inconsistent retrievals and the Scheduler. Concurrency control with locking methods : lock granularity, lock types, two phase locking for ensuring serializability, deadlocks, Concurrency control with time stamp ordering : Wait/Die and Wound/Wait Schemes, Database Recovery management : Transaction recovery.

## Transaction:-

A transaction can be defined as a group of tasks. A single task is the minimum processing unit which cannot be divided further.

Let's take an example of a simple transaction. Suppose a bank employee transfers Rs 500 from A's account to B's account. This very simple and small transaction involves several low-level tasks.

## A's Account

```
Open_Account(A)
Old_Balance = A.balance
New_Balance = Old_Balance - 500
A.balance = New_Balance
Close_Account(A)
```

## B's Account

```
Open_Account(B)
Old_Balance = B.balance
New_Balance = Old_Balance + 500
B.balance = New_Balance
Close_Account(B)
```

### Properties of Transactions:-

A transaction is a very small unit of a program and it may contain several lowlevel tasks. A transaction in a database system must maintain **A**tomicity, **C**onsistency, **I**solation, and **D**urability − commonly known as ACID properties − in order to ensure accuracy, completeness, and data integrity.

- **Atomicity** − This property states that a transaction must be treated as an atomic unit, that is, either all of its operations are executed or none. There must be no state in a database where a transaction is left partially completed. States should be defined either before the execution of the transaction or after the execution/abortion/failure of the transaction.

- **Consistency** − The database must remain in a consistent state after any transaction. No transaction should have any adverse effect on the data residing in the database. If the database was in a consistent state before the execution of a transaction, it must remain consistent after the execution of the transaction as well.

- **Durability** − The database should be durable enough to hold all its latest updates even if the system fails or restarts. If a transaction updates a chunk of data in a database and commits, then the database will hold the modified data. If a transaction commits but the system fails before the data could be written on to the disk, then that data will be updated once the system springs back into action.

- **Isolation** − In a database system where more than one transaction are being executed simultaneously and in parallel, the property of isolation states that all the transactions will be carried out and executed as if it is the only transaction in the system. No transaction will affect the existence of any other transaction.

## Transaction Log:-

A **transaction log** (also **transaction journal**, **database log**, **binary log** or **audit trail**) is a history of actions executed by a database management system used to guarantee ACID properties over crashes or hardware failures. Physically, a log is a file listing changes to the database, stored in a stable storage format.

If, after a start, the database is found in an inconsistent state or not been shut down properly, the database management system reviews the database logs for uncommittedtransactions and rolls back the changes made by these transactions. Additionally, all transactions that are already committed but whose changes were not yet materialized in the database are re-applied. Both are done to ensure atomicity and durability of transactions.

This term is not to be confused with other, human-readable logs that a database management system usually provides.

## Transaction Management with SQL Commit, Rollback and Savepoint :-

Transaction Control Language(TCL) commands are used to manage transactions in the database. These are used to manage the changes made to the data in a table by DML statements. It also allows statements to be grouped together into logical transactions.

## COMMIT command:-

`COMMIT` command is used to permanently save any transaction into the database.

When we use any DML command like `INSERT`, `UPDATE` or `DELETE`, the changes made by these commands are not permanent, until the current session is closed, the changes made by these commands can be rolled back.

To avoid that, we use the `COMMIT` command to mark the changes as permanent.

Following is commit command's syntax,

COMMIT;

**ROLLBACK command:-**

This command restores the database to last commited state. It is also used with SAVEPOINT command to jump to a savepoint in an ongoing transaction.

If we have used the UPDATE command to make some changes into the database, and realise that those changes were not required, then we can use the ROLLBACK command to rollback those changes, if they were not commited using the COMMIT command.

Following is rollback command's syntax,

ROLLBACK TO savepoint_name;

**SAVEPOINT command:-**

SAVEPOINT command is used to temporarily save a transaction so that you can rollback to that point whenever required.

Following is savepoint command's syntax,

SAVEPOINT savepoint_name;

**Concurrency Control:-**

In a multiprogramming environment where multiple transactions can be executed simultaneously, it is highly important to control the concurrency of transactions. We have concurrency control protocols to ensure atomicity, isolation, and serializability of concurrent transactions. Concurrency control protocols can be broadly divided into two categories

- Lock based protocols
- Time stamp based protocols

**Lock-based Protocols**

Database systems equipped with lock-based protocols use a mechanism by which any transaction cannot read or write data until it acquires an appropriate lock on it. Locks are of two kinds −

- **Binary Locks** − A lock on a data item can be in two states; it is either locked or unlocked.

- **Shared/exclusive** − This type of locking mechanism differentiates the locks based on their uses. If a lock is acquired on a data item to perform a write operation, it is an exclusive lock. Allowing more than one transaction to write on the same data item would lead the database into an inconsistent state. Read locks are shared because no data value is being changed.

There are four types of lock protocols available −

### Simplistic Lock Protocol

Simplistic lock-based protocols allow transactions to obtain a lock on every object before a 'write' operation is performed. Transactions may unlock the data item after completing the 'write' operation.

### Pre-claiming Lock Protocol

Pre-claiming protocols evaluate their operations and create a list of data items on which they need locks. Before initiating an execution, the transaction requests the system for all the locks it needs beforehand. If all the locks are granted, the transaction executes and releases all the locks when all its operations are over. If all the locks are not granted, the transaction rolls back and waits until all the locks are granted.



### Two-Phase Locking 2PL

This locking protocol divides the execution phase of a transaction into three parts. In the first part, when the transaction starts executing, it seeks permission for the locks it requires. The second part is where the transaction acquires all the locks. As soon as the transaction releases its first lock, the third phase starts. In this phase, the transaction cannot demand any new locks; it only releases the acquired locks.



Two-phase locking has two phases, one is **growing**, where all the locks are being acquired by the transaction; and the second phase is shrinking, where the locks held by the transaction are being released.

To claim an exclusive (write) lock, a transaction must first acquire a shared (read) lock and then upgrade it to an exclusive lock.

### Strict Two-Phase Locking

The first phase of Strict-2PL is same as 2PL. After acquiring all the locks in the first phase, the transaction continues to execute normally. But in contrast to 2PL, Strict-2PL does not release a lock after using it. Strict-2PL holds all the locks until the commit point and releases all the locks at a time.

Strict-2PL does not have cascading abort as 2PL does.

## Timestamp-based Protocols

The most commonly used concurrency protocol is the timestamp based protocol. This protocol uses either system time or logical counter as a timestamp.

Lock-based protocols manage the order between the conflicting pairs among transactions at the time of execution, whereas timestamp-based protocols start working as soon as a transaction is created.

Every transaction has a timestamp associated with it, and the ordering is determined by the age of the transaction. A transaction created at 0002 clock time would be older than all other transactions that come after it. For example, any transaction 'y' entering the system at 0004 is two seconds younger and the priority would be given to the older one.

In addition, every data item is given the latest read and write-timestamp. This lets the system know when the last 'read and write' operation was performed on the data item.

## Timestamp Ordering Protocol

The timestamp-ordering protocol ensures serializability among transactions in their conflicting read and write operations. This is the responsibility of the protocol system that the conflicting pair of tasks should be executed according to the timestamp values of the transactions.

- The timestamp of transaction $T_i$ is denoted as $TS(T_i)$.
- Read time-stamp of data-item X is denoted by R-timestamp(X).
- Write time-stamp of data-item X is denoted by W-timestamp(X).

Timestamp ordering protocol works as follows −

- **If a transaction Ti issues a read(X) operation −**

  o If $TS(Ti) < $ W-timestamp(X)

    ▪ Operation rejected.

  o If $TS(Ti) >= $ W-timestamp(X)

    ▪ Operation executed.

  o All data-item timestamps updated.

- **If a transaction Ti issues a write(X) operation −**

  o If $TS(Ti) < $ R-timestamp(X)

- Operation rejected.
  - If $TS(Ti) < W\text{-}timestamp(X)$
    - Operation rejected and Ti rolled back.
  - Otherwise, operation executed.

## Serializability

When multiple transactions are being executed by the operating system in a multiprogramming environment, there are possibilities that instructions of one transactions are interleaved with some other transaction.

- **Schedule** − A chronological execution sequence of a transaction is called a schedule. A schedule can have many transactions in it, each comprising of a number of instructions/tasks.

- **Serial Schedule** − It is a schedule in which transactions are aligned in such a way that one transaction is executed first. When the first transaction completes its cycle, then the next transaction is executed. Transactions are ordered one after the other. This type of schedule is called a serial schedule, as transactions are executed in a serial manner.

In a multi-transaction environment, serial schedules are considered as a benchmark. The execution sequence of an instruction in a transaction cannot be changed, but two transactions can have their instructions executed in a random fashion. This execution does no harm if two transactions are mutually independent and working on different segments of data; but in case these two transactions are working on the same data, then the results may vary. This ever-varying result may bring the database to an inconsistent state.

To resolve this problem, we allow parallel execution of a transaction schedule, if its transactions are either serializable or have some equivalence relation among them.

## Equivalence Schedules

An equivalence schedule can be of the following types −

## Result Equivalence

If two schedules produce the same result after execution, they are said to be result equivalent. They may yield the same result for some value and different results for another set of values. That's why this equivalence is not generally considered significant.

## View Equivalence

Two schedules would be view equivalence if the transactions in both the schedules perform similar actions in a similar manner.

For example −

- If T reads the initial data in S1, then it also reads the initial data in S2.

- If T reads the value written by J in S1, then it also reads the value written by J in S2.

- If T performs the final write on the data value in S1, then it also performs the final write on the data value in S2.

## Conflict Equivalence

Two schedules would be conflicting if they have the following properties

- Both belong to separate transactions.
- Both accesses the same data item.
- At least one of them is "write" operation.

Two schedules having multiple transactions with conflicting operations are said to be conflict equivalent if and only if −

- Both the schedules contain the same set of Transactions.
- The order of conflicting pairs of operation is maintained in both the schedules.

**Note** − View equivalent schedules are view serializable and conflict equivalent schedules are conflict serializable. All conflict serializable schedules are view serializable too.

## Concurrency control with time stamp ordering : Wait/Die and Wound/Wait Schemes:-

In a multi-process system, deadlock is an unwanted situation that arises in a shared resource environment, where a process indefinitely waits for a resource that is held by another process.

For example, assume a set of transactions $\{T_0, T_1, T_2, ...,T_n\}$. $T_0$ needs a resource X to complete its task. Resource X is held by $T_1$, and $T_1$ is waiting for a resource Y, which is held by $T_2$. $T_2$ is waiting for resource Z, which is held by $T_0$. Thus, all the processes wait for each other to release resources. In this situation, none of the processes can finish their task. This situation is known as a deadlock.

Deadlocks are not healthy for a system. In case a system is stuck in a deadlock, the transactions involved in the deadlock are either rolled back or restarted.

## Deadlock Prevention

To prevent any deadlock situation in the system, the DBMS aggressively inspects all the operations, where transactions are about to execute. The DBMS inspects the operations and analyzes if they can create a deadlock situation. If it finds that a deadlock situation might occur, then that transaction is never allowed to be executed.

There are deadlock prevention schemes that use timestamp ordering mechanism of transactions in order to predetermine a deadlock situation.

## Wait-Die Scheme

In this scheme, if a transaction requests to lock a resource (data item), which is already held with a conflicting lock by another transaction, then one of the two possibilities may occur −

- If $TS(T_i) < TS(T_j)$ − that is $T_i$, which is requesting a conflicting lock, is older than $T_j$ − then $T_i$ is allowed to wait until the data-item is available.

www.FirstRanker.com          www.FirstRanker.com

- If $TS(T_i) > TS(t_j)$ − that is $T_i$ is younger than $T_j$ − then $T_i$ dies. $T_i$ is restarted later with a random delay but with the same timestamp.

This scheme allows the older transaction to wait but kills the younger one.

### Wound-Wait Scheme

In this scheme, if a transaction requests to lock a resource (data item), which is already held with conflicting lock by some another transaction, one of the two possibilities may occur −

- If $TS(T_i) < TS(T_j)$, then $T_i$ forces $T_j$ to be rolled back − that is $T_i$wounds $T_j$. $T_j$ is restarted later with a random delay but with the same timestamp.

- If $TS(T_i) > TS(T_j)$, then $T_i$ is forced to wait until the resource is available.

This scheme, allows the younger transaction to wait; but when an older transaction requests an item held by a younger one, the older transaction forces the younger one to abort and release the item.

In both the cases, the transaction that enters the system at a later stage is aborted.

### Transaction recovery:-

### Crash Recovery

DBMS is a highly complex system with hundreds of transactions being executed every second. The durability and robustness of a DBMS depends on its complex architecture and its underlying hardware and system software. If it fails or crashes amid transactions, it is expected that the system would follow some sort of algorithm or techniques to recover lost data.

### Failure Classification

To see where the problem has occurred, we generalize a failure into various categories, as follows −

### Transaction failure

A transaction has to abort when it fails to execute or when it reaches a point from where it can't go any further. This is called transaction failure where only a few transactions or processes are hurt.

Reasons for a transaction failure could be −

- **Logical errors** − Where a transaction cannot complete because it has some code error or any internal error condition.

- **System errors** − Where the database system itself terminates an active transaction because the DBMS is not able to execute it, or it has to stop because of some system condition. For example, in case of deadlock or resource unavailability, the system aborts an active transaction.

### System Crash

There are problems − external to the system − that may cause the system to stop abruptly and cause the system to crash. For example, interruptions in power supply may cause the failure of underlying hardware or software failure.

Examples may include operating system errors.

### Disk Failure

In early days of technology evolution, it was a common problem where hard-disk drives or storage drives used to fail frequently.

Disk failures include formation of bad sectors, unreachability to the disk, disk head crash or any other failure, which destroys all or a part of disk storage.

## UNIT-VI

**SYLLABUS**

Overview of Storages and Indexing, Data on External Storage- File Organization and Indexing –Clustered Indexing – Primary and Secondary Indexes, Index Data Structures, Hash-Based Indexing – Tree-Based Indexing, Comparison of File Organization.

### Overview of Storages and Indexing:-

Databases are stored in file formats, which contain records. At physical level, the actual data is stored in electromagnetic format on some device. These storage devices can be broadly categorized into three types −



- **Primary Storage** − The memory storage that is directly accessible to the CPU comes under this category. CPU's internal memory (registers), fast memory (cache), and main memory (RAM) are directly accessible to the CPU, as they are all placed on the motherboard or CPU chipset. This storage is typically very small, ultra-fast, and volatile. Primary storage requires continuous power supply in order to maintain its state. In case of a power failure, all its data is lost.

- **Secondary Storage** − Secondary storage devices are used to store data for future use or as backup. Secondary storage includes memory devices that are not a part of the CPU chipset or motherboard, for example, magnetic disks, optical disks (DVD, CD, etc.), hard disks, flash drives, and magnetic tapes.

- **Tertiary Storage** − Tertiary storage is used to store huge volumes of data. Since such storage devices are external to the computer system, they are the slowest in speed. These storage devices are mostly used to take the back up of an entire system. Optical disks and magnetic tapes are widely used as tertiary storage.

### Primary and Secondary Indexes:-

We know that data is stored in the form of records. Every record has a key field, which helps it to be recognized uniquely.

Indexing is a data structure technique to efficiently retrieve records from the database files based on some attributes on which the indexing has been done. Indexing in database systems is similar to what we see in books.

Indexing is defined based on its indexing attributes. Indexing can be of the following types −

- **Primary Index** − Primary index is defined on an ordered data file. The data file is ordered on a **key field**. The key field is generally the primary key of the relation.

- **Secondary Index** − Secondary index may be generated from a field which is a candidate key and has a unique value in every record, or a non-key with duplicate values.

- **Clustering Index** − Clustering index is defined on an ordered data file. The data file is ordered on a non-key field.

Ordered Indexing is of two types −

- Dense Index
- Sparse Index

### Dense Index

In dense index, there is an index record for every search key value in the database. This makes searching faster but requires more space to store index records itself. Index records contain search key value and a pointer to the actual record on the disk.



### Sparse Index

In sparse index, index records are not created for every search key. An index record here contains a search key and an actual pointer to the data on the disk. To search a record, we first proceed by index record and reach at the actual location of the data. If the data we are looking for is not where we directly reach by following the index, then the system starts sequential search until the desired data is found.



### Multilevel Index

Index records comprise search-key values and data pointers. Multilevel index is stored on the disk along with the actual database files. As the size of the database grows, so does the size of the indices. There is an immense need to keep the index records in the main memory so as to speed up the search operations. If single-level index is used, then a large size index cannot be kept in memory which leads to multiple disk accesses.

Multi-level Index helps in breaking down the index into several smaller indices in order to make the outermost level so small that it can be saved in a single disk block, which can easily be accommodated anywhere in the main memory.

## Hash-Based Indexing:-

For a huge database structure, it can be almost next to impossible to search all the index values through all its level and then reach the destination data block to retrieve the desired data. Hashing is an effective technique to calculate the direct location of a data record on the disk without using index structure.

Hashing uses hash functions with search keys as parameters to generate the address of a data record.

## Hash Organization

- **Bucket** − A hash file stores data in bucket format. Bucket is considered a unit of storage. A bucket typically stores one complete disk block, which in turn can store one or more records.

- **Hash Function** − A hash function, **h,** is a mapping function that maps all the set of search-keys **K** to the address where actual records are placed. It is a function from search keys to bucket addresses.

## Static Hashing

In static hashing, when a search-key value is provided, the hash function always computes the same address. For example, if mod-4 hash function is used, then it shall generate only 5 values. The output address shall always be same for that function. The number of buckets provided remains unchanged at all times.

### Operation

- **Insertion** − When a record is required to be entered using static hash, the hash function **h** computes the bucket address for search key **K**, where the record will be stored.

  Bucket address = h(K)

- **Search** − When a record needs to be retrieved, the same hash function can be used to retrieve the address of the bucket where the data is stored.

- **Delete** − This is simply a search followed by a deletion operation.

### Bucket Overflow

The condition of bucket-overflow is known as **collision**. This is a fatal state for any static hash function. In this case, overflow chaining can be used.

- **Overflow Chaining** − When buckets are full, a new bucket is allocated for the same hash result and is linked after the previous one. This mechanism is called **Closed Hashing**.



- **Linear Probing** − When a hash function generates an address at which data is already stored, the next free bucket is allocated to it. This mechanism is called **Open Hashing**.

<u>Dynamic Hashing</u>

The problem with static hashing is that it does not expand or shrink dynamically as the size of the database grows or shrinks. Dynamic hashing provides a mechanism in which data buckets are added and removed dynamically and on-demand. Dynamic hashing is also known as **extended hashing**.

Hash function, in dynamic hashing, is made to produce a large number of values and only a few are used initially.



## Tree-based Indexing:-

### B⁺ Tree

A B⁺ tree is a balanced binary search tree that follows a multi-level index format. The leaf nodes of a B⁺ tree denote actual data pointers. B⁺ tree ensures that all leaf nodes remain at the same height, thus balanced. Additionally, the leaf nodes are linked using a link list; therefore, a B⁺ tree can support random access as well as sequential access.

### Structure of B⁺ Tree

Every leaf node is at equal distance from the root node. A B⁺ tree is of the order **n** where **n** is fixed for every B⁺ tree.



**Internal nodes** −

- Internal (non-leaf) nodes contain at least ⌈n/2⌉ pointers, except the root node.
- At most, an internal node can contain **n** pointers.

**Leaf nodes** −

- Leaf nodes contain at least ⌈n/2⌉ record pointers and ⌈n/2⌉ key values.
- At most, a leaf node can contain **n** record pointers and **n** key values.

- Every leaf node contains one block pointer **P** to point to next leaf node and forms a linked list.

### B$^+$ Tree Insertion

- B$^+$ trees are filled from bottom and each entry is done at the leaf node.

- If a leaf node overflows −

  - Split node into two parts.

  - Partition at **i = ⌊(m+1)$_{/2}$⌋.**

  - First **i** entries are stored in one node.

  - Rest of the entries (i+1 onwards) are moved to a new node.

  - **i$^{th}$** key is duplicated at the parent of the leaf.

- If a non-leaf node overflows −

  - Split node into two parts.

  - Partition the node at **i = ⌈(m+1)$_{/2}$⌉**.

  - Entries up to **i** are kept in one node.

  - Rest of the entries are moved to a new node.

### B$^+$ Tree Deletion

- B$^+$ tree entries are deleted at the leaf nodes.

- The target entry is searched and deleted.

  - If it is an internal node, delete and replace with the entry from the left position.

- After deletion, underflow is tested,

  - If underflow occurs, distribute the entries from the nodes left to it.

- If distribution is not possible from left, then

  - Distribute from the nodes right to it.

- If distribution is not possible from left or from right, then

  - Merge the node with left and right to it.

**Comparison of File Organization:-**

Relative data and information is stored collectively in file formats. A file is a sequence of records stored in binary format. A disk drive is formatted into several blocks that can store records. File records are mapped onto those disk blocks.

File Organization defines how file records are mapped onto disk blocks. We have four types of File Organization to organize file records –



## Heap File Organization

When a file is created using Heap File Organization, the Operating System allocates memory area to that file without any further accounting details. File records can be placed anywhere in that memory area. It is the responsibility of the software to manage the records. Heap File does not support any ordering, sequencing, or indexing on its own.

## Sequential File Organization

Every file record contains a data field (attribute) to uniquely identify that record. In sequential file organization, records are placed in the file in some sequential order based on the unique key field or search key. Practically, it is not possible to store all the records sequentially in physical form.

## Hash File Organization

Hash File Organization uses Hash function computation on some fields of the records. The output of the hash function determines the location of disk block where the records are to be placed.

## Clustered File Organization

Clustered file organization is not considered good for large databases. In this mechanism, related records from one or more relations are kept in the same disk block, that is, the ordering of records is not based on primary key or search key.

## File Operations

Operations on database files can be broadly classified into two categories

- **Update Operations**

- **Retrieval Operations**

Update operations change the data values by insertion, deletion, or update. Retrieval operations, on the other hand, do not alter the data but retrieve them after optional conditional filtering. In both types of operations, selection plays a significant role. Other than creation and deletion of a file, there could be several operations, which can be done on files.

- **Open** − A file can be opened in one of the two modes, **read mode** or **write mode**. In read mode, the operating system does not allow anyone to alter data. In other words, data is read only. Files opened in read mode can be shared among several entities. Write mode allows data modification. Files opened in write mode can be read but cannot be shared.

- **Locate** − Every file has a file pointer, which tells the current position where the data is to be read or written. This pointer can be adjusted accordingly. Using find (seek) operation, it can be moved forward or backward.

- **Read** − By default, when files are opened in read mode, the file pointer points to the beginning of the file. There are options where the user can tell the operating system where to locate the file pointer at the time of opening a file. The very next data to the file pointer is read.

- **Write** − User can select to open a file in write mode, which enables them to edit its contents. It can be deletion, insertion, or modification. The file pointer can be located at the time of opening or can be dynamically changed if the operating system allows to do so.

- **Close** − This is the most important operation from the operating system's point of view. When a request to close a file is generated, the operating system

  o removes all the locks (if in shared mode),

  o saves the data (if altered) to the secondary storage media, and

  o releases all the buffers and file handlers associated with the file.

The organization of data inside a file plays a major role here. The process to locate the file pointer to a desired record inside a file various based on whether the records are arranged sequentially or clustered.

## VSM COLLEGE OF ENGINEERING

## DEPARTMENT OF CSE

## III/IV B.TECH – I SEM (R16)

## DATABASE MANAGEMENT SYSTEMS

## FREQUENTLY ASKED QUESTIONS

# UNIT -I

1. Define DBMS? Explain in detail about DBMS architecture with a neat sketch.

2. Define Data Independence? Explain about two types of data independence with a relevant examples.

3. Define Mapping? Discuss Three Level Schema Archictecture with relevant examples.

4. Explain the significance of DBA in Database Environment?

5. Explain the concept of Client/Server Interaction in Client/Server Architecture.

# UNIT -II

1. Define Relation? Explain in detail about Relational model and its related terms with examples?

2. Explain various steps to design a good database?

3. Draw E-R diagram for an employee database? Explain about the key terms involved in E-R diagram with relevant examples?

4. Define Constraint? Explain various integrity constraints with suitable examples?

5. Differentiate Tuple Relational Calculus with Domain Relational Calculus with examples?

6. Define Relational Algebra? Explain in detail about various operations in relational algebra with suitable examples?

# UNIT -III

1. Define Query? Explain various forms that are present in SQL?

2. Discuss in detail the operators SELECT, PROJECT, UNION with suitable examples.

3. Explain the importance of Null values in Relational Model.

4. Discuss various Aggregate operators with suitable example?

5. Explain the concept of triggers with relevant example?

6. Differentiate between independent and correlated nested queries.

# UNIT –IV

1) **Illustrate** redundancy and the problems that it can cause?

2) **Define** decomposition and how does it address redundancy? Discuss the problem s that may be caused by the use of decompositions?

3) **Define** functional dependencies. How are primary keys related to FD's?

4) **Define** normalization? Explain 1NF, 2NF, 3NF Normal forms?

5) **Compare and contrast** BCNF with 3NF?

6) **Describe properties of decompositions?**

7) **Explain** about Schema refinement in Database design?

8) **Illustrate** Multivalued dependencies and Fourth normal form with example?

# UNIT –V

1) **Explain** ACID properties and Illustrate them through examples?

2) **Discuss** How do you implement Atomicity and Durability?

3) **Illustrate** Concurrent execution of transaction with examples?

4) **Discuss** Serializability in detail?

5) **Discuss** two phase locking protocol and strict two phase locking protocols?

6) **Describe Timestamp based locking protocols?**

7) **Explain** in detail Storage structure?

8) **Discuss** how do you recover from Concurrent transactions?

# UNIT –VI

1) **Write** in detail about Hash based Indexing and Tree based Indexing?

2) Compare **various** File Organizations ?

3) **Explain** B+ trees? Discuss about this Dynamic Index Structure?

4) **Demonstrate** searching a given element in B+ trees? Explain with example?

5) **Illustrate** insertion and deletion of an element in B+ trees with example?

6) **Write** in detail about Static Hashing?

7) **Explain** in detail about Extendible Hashing?

8) **Explain** in detail about Linear Hashing?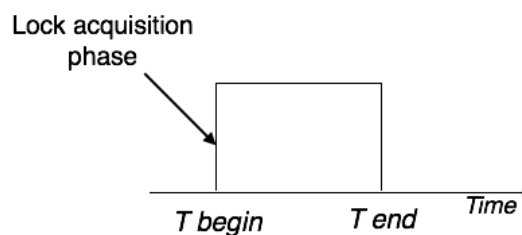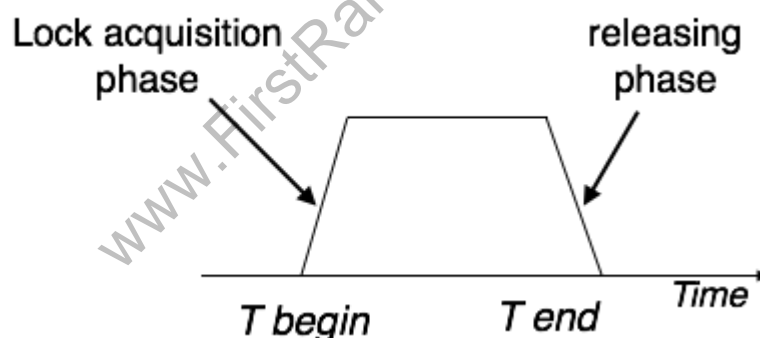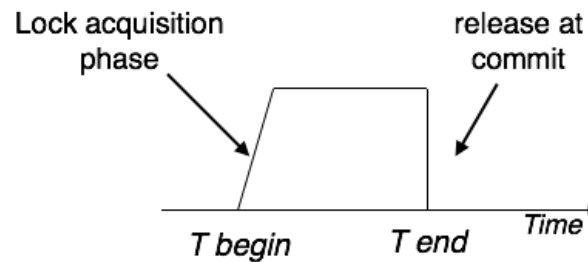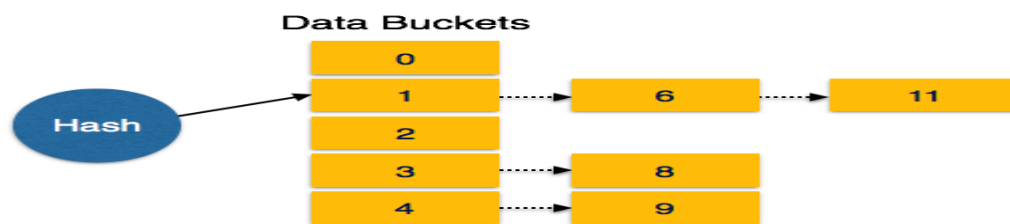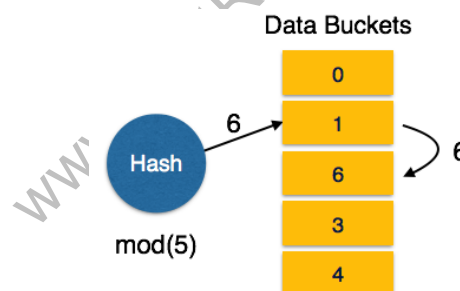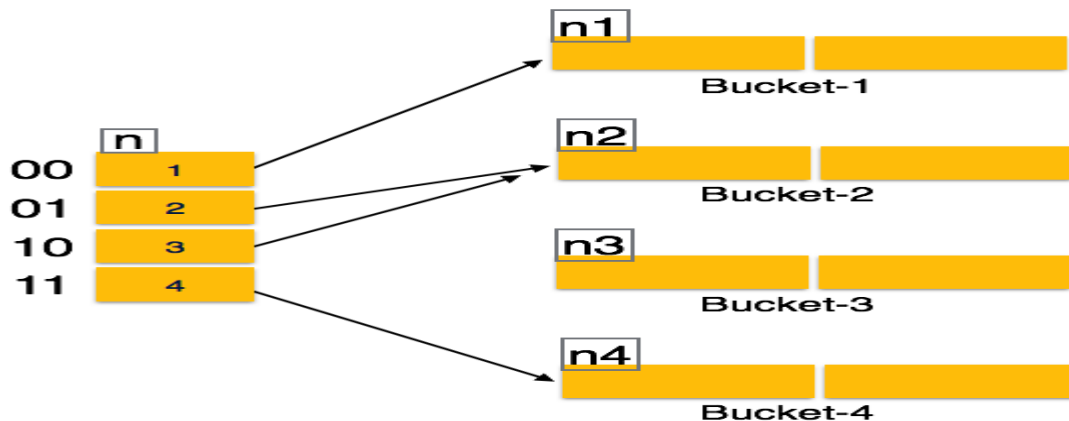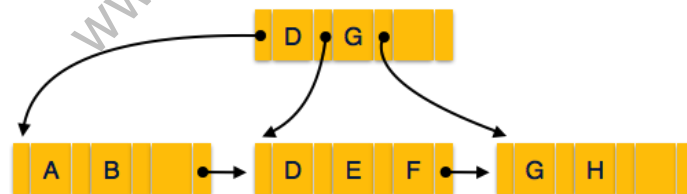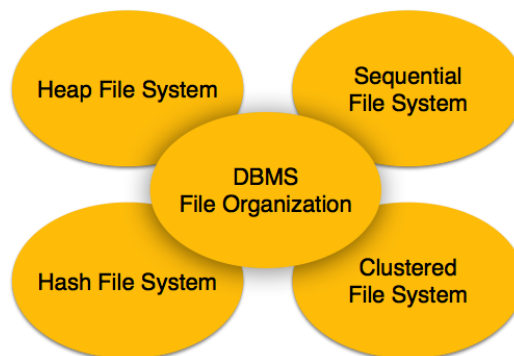