

DBMS Notes

Unit- 1 Introduction

PART 1: What is a DBMS?

College Level

A **Database Management System (DBMS)** is software that helps users **store, organize, retrieve, and manage data efficiently**.

Industry View

In companies like Google, Amazon, or Zomato:

- Data is stored in **databases** like **PostgreSQL, MySQL, MongoDB, Oracle**, etc.
 - DBMS helps them **manage users, ensure performance, and maintain data integrity**.
-

DBMS vs File System

Feature	File System	DBMS
Data Redundancy	High	Low (through normalization)
Data Consistency	Hard to manage	Enforced with constraints & transactions
Querying	Manual file reading	Powerful languages like SQL
Security	Weak (OS-level)	Role-based, access-controlled
Backup/Recovery	Manual	Built-in backup, crash recovery

Example:

- In a file system: Student data is saved as .txt or .csv manually.
 - In DBMS: A structured **Students** table with name, roll, course, etc., and accessed via **SQL queries**.
-

PART 2: Database Architecture

Three-Schema Architecture (College Syllabus)

1. **External Level:** How users view the data (e.g., front-end UI shows only user name).
2. **Logical Level:** What data is stored and what relations exist.
3. **Internal Level:** How the data is actually stored in memory (indexes, blocks).

[External Schema]



[Logical Schema]



[Internal Schema]

❖ Industry View

💡 In real-world apps:

- **External** = what the frontend needs.
 - **Logical** = ORM models like User in Django/Node.js.
 - **Internal** = engine-level storage (e.g., PostgreSQL heap files, B-Trees).
-

🔌 PART 3: Data Independence

📘 College Definition

Data Independence means **changes at one schema level don't affect others**.

- **Logical Data Independence:** Change in logical schema (e.g., adding a column) doesn't break applications.
- **Physical Data Independence:** Change in physical storage (e.g., B-tree to hash index) doesn't break schema.

💻 Real-World Relevance

Helps companies:

- Upgrade storage backends without changing code.
 - Add new columns/tables with zero downtime (zero disruption to frontend).
-

🧠 PART 4: ER Model (Entity-Relationship Model)

Let's now shift to **Data Modeling using ER Diagrams**.

◆ **Step-by-Step ER Model Concepts:**

1. **Entity:** Real-world object (e.g., Student, Course)
2. **Attribute:** Property of entity (e.g., name, age)
 - Simple, Composite, Derived, Multivalued
3. **Entity Set:** Collection of similar entities
4. **Relationship:** Association (e.g., Student ENROLLED_IN Course)
5. **Degree:** Number of entity sets in a relationship (Unary, Binary, Ternary)
6. **Cardinality:** 1:1, 1:N, M:N

 **Keys (Important for Normalization)**

- **Super Key:** Any set of attributes uniquely identifying a record
 - **Candidate Key:** Minimal super key
 - **Primary Key:** One chosen candidate key
 - **Foreign Key:** Attribute that refers to another table's primary key
-

 **ER Diagram Example (COLLEGE SYSTEM):**

Entity: STUDENT

- ID (PK)
- Name
- DOB
- Department

Entity: COURSE

- Code (PK)
- Title
- Credits

Relationship: ENROLLS

- Attributes: Grade

- Cardinality: M:N (Many students can enroll in many courses)

We use **diamonds** for relationships, **rectangles** for entities, and **ellipses** for attributes.

⭐ Extended ER Concepts

1. **Generalization:** Combine multiple entities → super entity

Example: Undergrad and Postgrad → Student

2. **Specialization:** Split entity → sub-entities

Example: Student → RegularStudent, ExchangeStudent

3. **Aggregation:** Treat relationship as an entity

Example: Professor supervises a ResearchProject (which is itself a relationship between Student and Project)

✅ What We Covered:

College Topic	Covered
DBMS vs File System	✓
Architecture (3-schema)	✓
Data Independence	✓
ER Model, Relationships, Keys	✓
Extended ER (Aggregation, Gen., Spec.)	✓

Industry-Level Insight	Included?
Modern DB usage (SQL vs NoSQL)	✓
3-schema mapping to real backend layers	✓
ORM and data modeling in practice	✓

✳️ 1. What is an ER Model?

The **Entity-Relationship Model (ER Model)** is a high-level **data model** used for **conceptual design** of a database.

It allows us to:

- Model the **real-world entities**
 - Define **relationships** between entities
 - Design before implementing in SQL
-

◆ 2. Key ER Model Concepts

Term	Meaning
Entity	Object or thing in the real world (e.g., Student, Course)
Entity Set	A group of similar entities (e.g., all Students)
Attribute	Property of an entity (e.g., Student has Name, Age)
Key Attribute	Uniquely identifies an entity (e.g., RollNo)
Relationship	Association between entities (e.g., Student <i>enrolls in</i> Course)

⌚ 3. ER Diagram Notation (with ASCII-style diagram)

[Student] ----- Enrolls ----- [Course]

| |
Name Code
RollNo Title

- **Entities:** Rectangle (Student, Course)
 - **Attributes:** Ovals (Name, RollNo, Title)
 - **Relationships:** Diamond (Enrolls)
-

🔒 4. Keys in ER Model

Key Type	Meaning
Super Key	Set of one or more attributes that uniquely identify an entity

Key Type	Meaning
Candidate Key	Minimal super key (no extra attribute)
Primary Key	Selected candidate key
Foreign Key	Attribute that refers to another entity's primary key

Example:

- Student(RollNo, Name, Email)
 - Super Keys: {RollNo}, {RollNo, Name}, {RollNo, Email}
 - Candidate Keys: {RollNo}
 - Primary Key: RollNo
-

5. Generalization, Specialization, Aggregation

Concept	Meaning
Generalization	Bottom-up approach (combine entities into superclass)
→ Car and Bike → Vehicle	
Specialization	Top-down approach (split entity into subclasses)
→ Employee → Manager, Developer	
Aggregation	Treat a relationship as an entity
→ Project involves Employee and Client as a unit	

6. ER to Table Conversion (Reduction)

ER Concept → SQL Table

Entity	→ Table
Attribute	→ Column

ER Concept → SQL Table

Relationship → Foreign Key (or separate table for many-to-many)

Example:

Student(RollNo, Name)

Course(Code, Title)

Enrolls(RollNo, Code) ← relationship table with FK to both

Industry Relevance

Area	Relevance
<input checked="" type="checkbox"/> ER Diagrams	Used during DB Design (before writing SQL)
<input checked="" type="checkbox"/> Primary/Foreign Keys Foundation of relational databases	
<input checked="" type="checkbox"/> Generalization	Used in object-oriented DBs and ORMs
<input checked="" type="checkbox"/> ER → SQL Mapping	Used by tools like Lucidchart, Draw.io, MySQL Workbench

- Great! Let's now continue **DBMS Unit-I** with the final part:
-

ER Model to Relational Schema Mapping

Once we design an ER diagram, we must **convert** it into a relational schema (tables). This step bridges **conceptual modeling** and **actual database implementation**.

Step-by-Step ER → Relational Table Mapping

ER Component	Relational Mapping Result
--------------	---------------------------

Entity A Table

Attribute Column in table

Primary Key Primary Key of table

Relationship (1:1, 1:N, N:N) Depends on cardinality (see below)

ER Component	Relational Mapping Result
Multivalued Attribute	Create separate table
Composite Attribute	Flatten to atomic columns

Example ER to Table

ER Diagram Snippet:

- **Entity:** Student
 - Attributes: RollNo (PK), Name, Email
- **Entity:** Course
 - Attributes: CourseID (PK), Title
- **Relationship:** Enrolled between Student and Course
 - Attributes: Grade

Mapping:

Table: Student

RollNo (PK) Name Email

Table: Course

CourseID (PK) Title

Table: Enrolled (*relationship table*)

RollNo (FK) CourseID (FK) Grade

- Foreign keys represent the relationship
 - We use a **separate table** for **many-to-many** (Student ↔ Course)
-

Special Cases

1. 1:1 Relationship

→ Add foreign key to either entity (usually the one with total participation)

2. 1:N Relationship

→ Add foreign key in the "many" side

3. N:N Relationship

→ Always create a separate relationship table

4. Multivalued Attributes

→ Create separate table with FK to parent and each value as a row

Industry View

- ER Diagrams are often skipped in **agile startups** (move fast), but used in:
 - **Enterprise projects, banking, telecom, hospital systems**
 - For database design tools (e.g., **MySQL Workbench, dbdiagram.io, Draw.io**)
 - **Reverse engineering** is also common: generating ER diagrams from existing databases
-

Summary

- ER model is **for design**, relational model is **for implementation**
 - Mapping requires understanding **relationship cardinality** and **data dependencies**
 - In real-world projects, teams review ER diagrams before creating schemas
-

Great! Let's begin **DBMS Unit II – Relational Model & SQL** — a *core pillar* for both exams and real-world backend development.

◆ Unit II: Relational Data Model & SQL (College + Industry)

1. What is the Relational Model?

A **Relational Model** organizes data into **tables (relations)** made of:

- **Rows (tuples)** = records
- **Columns (attributes)** = fields

It's the most widely used database model today (used by MySQL, PostgreSQL, Oracle, etc.)

Relational Model Terms

Term	Meaning
Relation	Table
Tuple	Row in a table
Attribute	Column in a table
Domain	Valid set of values for an attribute
Degree	No. of attributes (columns)
Cardinality	No. of tuples (rows)

2. Properties of a Relation

- Each tuple is **unique** (no duplicates)
 - Tuples are **unordered**
 - Attributes are **atomic** (no multivalued/composite fields)
 - Every attribute has a **domain**
-

3. Keys in Relational Model

Type	Purpose
Primary Key	Uniquely identifies a tuple
Candidate Key	Minimal set of attributes to uniquely identify
Super Key	Superset of candidate key
Foreign Key	References primary key from another table

 *Industry:* Keys help enforce **data integrity**. Missing/wrong foreign keys can break joins and cause real-world app bugs.

4. Integrity Constraints (Exam + Real-World)

Constraint Type	Ensures
Domain	Data type & valid range
Entity Integrity	No nulls in primary key
Referential Integrity	Foreign key must match existing PK or be null
User-Defined	Business rules (e.g., age > 18)

 **Industry:** These constraints ensure **data quality** and **business logic** are enforced in DB.

5. Introduction to SQL (Structured Query Language)

SQL = language to **create**, **insert**, **query**, and **update** relational databases.

SQL Components

SQL Category	Description	Example
DDL (Data Definition Language)	Define schema	CREATE TABLE, DROP
DML (Data Manipulation)	Insert/update/delete	INSERT, UPDATE, DELETE
DQL (Data Query)	Retrieve data	SELECT * FROM students;
DCL (Control)	Access control	GRANT, REVOKE
TCL (Transaction)	Manage transactions	COMMIT, ROLLBACK

6. Basic SQL Syntax

```
CREATE TABLE Student (
    RollNo INT PRIMARY KEY,
    Name VARCHAR(100),
    Age INT
);
```

```
INSERT INTO Student VALUES (1, 'Amit', 21);
```

```
SELECT * FROM Student;
```

```
UPDATE Student SET Name = 'Rohit' WHERE RollNo = 1;
```

```
DELETE FROM Student WHERE RollNo = 1;
```

SQL in Industry

- SQL is **essential in 100% of backend jobs**
 - Tools: pgAdmin, MySQL Workbench, DBeaver, Supabase Studio
 - Also used in: **Data analysis (Excel, Power BI, Looker)**
-

Awesome! Let's now cover the remaining **core concepts of Unit II** – this includes **Relational Algebra, Relational Calculus, and Joins** – all important for both college and industry.

◆ 7. Relational Algebra (College + Industry)

Relational Algebra is a **procedural query language** — it tells the *how* to get the result.

Used to:

- Theoretically validate queries
 - Build query compilers
 - Optimize SQL queries internally in databases
-

Basic Operations

Operation	Symbol	Description	Example
Selection	σ	Filters rows	$\sigma_{age > 18}(\text{Student})$
Projection	π	Selects columns	$\pi_{\text{Name}}(\text{Student})$
Union	\cup	Combines results (no duplicates) A \cup B	

Operation	Symbol Description	Example
Set Difference	-	A minus B
Cartesian Product	x	All combinations
Rename	ρ	Renames relation
		$\rho_{\text{New}}(\text{Student})$

Advanced Operations (Built using basic ones)

Operation Description

Join Combines related rows from 2 tables

Intersection Common rows

Division For queries like “all courses taken by all students”

Industry Relevance

Though you don't *write* relational algebra in companies, DB engines (like PostgreSQL or MySQL) **internally optimize queries using it**.

Also helps in **interviews** and **query optimization** topics.

8. Relational Calculus

Opposite of algebra: it's **non-procedural** – it tells **what** to fetch, not **how**.

Two types:

- **Tuple Relational Calculus (TRC)** → works with **tuples**
- **Domain Relational Calculus (DRC)** → works with **values/domains**

Think of it like:

"Get me all students where age > 18" — rather than specifying how to fetch them.

9. Joins (Super Important – College + Real-World Dev)

Joins allow combining rows from two or more tables based on a common key.

Types of Joins

Type	Description	Example
INNER JOIN	Only matching rows	Students enrolled in courses
LEFT JOIN	All from left, matched from right	Students with or without courses
RIGHT JOIN	All from right, matched from left	Courses with/without students
FULL JOIN	All rows from both sides	All students & all courses
CROSS JOIN	Cartesian Product	Every student × every course
SELF JOIN	Table joins with itself	Employee → Manager

Real-World Join Example:

```
SELECT s.Name, c.Title  
FROM Student s  
JOIN Enrolled e ON s.RollNo = e.RollNo  
JOIN Course c ON e.CourseID = c.CourseID;
```

This gives: Which student is enrolled in which course

Joins in Industry

- **Most common SQL operation in production**
 - Used in: Reporting, analytics, microservices, admin panels, dashboards
 - You must master joins for:
 - Backend development
 - Data Engineering
 - SQL interviews
-

Summary of Unit II (College + Industry)

Topic	Usage
Relational Model	Schema design
SQL Basics	CRUD operations
Relational Algebra	Query optimization, exams
Relational Calculus	Theoretical understanding
Joins	Real-world & interview essential

◆ Unit III: Normalization + Query Optimization (College + Industry Level)

◆ 1. What is Normalization?

Normalization = Process of organizing data in a database to:

- Reduce redundancy
 - Avoid anomalies (insert, update, delete issues)
 - Improve data integrity
-

✳ Why Normalize?

Without normalization:

- Repeated data wastes space
 - Changes become inconsistent
 - Relationships between data become messy
-

✳ 2. Normal Forms (NFs)

Normal Form	Rule	Removes
1NF	Atomic values (no lists)	Multivalued attributes

Normal Form	Rule	Removes
2NF	No partial dependency (on part of a PK)	Redundant data for composite PK
3NF	No transitive dependency ($A \rightarrow B \rightarrow C$)	Indirect dependencies
BCNF	Stronger version of 3NF	Any anomaly
4NF	No multivalued dependency	Multiple independent 1-to-many relationships
5NF	No join dependency	Complex joins only for valid combinations

Example (Normalization Step-by-Step):

Unnormalized Table (UNF):

RollNo Name Courses

1 Ayan DBMS, OS

→ Courses column is multivalued → **Not 1NF**

1NF:

RollNo Name Course

1 Ayan DBMS

1 Ayan OS

→ Atomic values 

2NF (Remove Partial Dependency):

If you had a composite PK, you'd split the table like:

- Student(RollNo, Name)
- Course(CourseID, Title)

- Enrollment(RollNo, CourseID)

→ Now no attribute depends on only part of PK 

3NF:

Avoid things like:

"Department depends on Subject, and Subject depends on Teacher"

→ Separate such dependencies into different tables

Real-World Need for Normalization

Context Example

App Backend Avoid storing same user info everywhere

Finance Prevent transactional inconsistencies

SaaS Clear user-subscription-plan separation

E-commerce Product → Variant → Seller breakup

3. Denormalization (Industry Only)

Sometimes in industry, **denormalization** is used to:

- Improve read performance
- Reduce JOINs (at the cost of some redundancy)

Used in:

- Analytics dashboards
- Read-heavy microservices
- Caching systems (Redis, ElasticSearch)

 You normalize to design, denormalize to scale.

4. Query Optimization (Exam + Industry Level)

Relational databases **optimize your SQL queries internally** using:

- Cost-based query planners
 - Heuristics
 - Index usage
 - Join order rearrangement
-

5. How Queries Are Optimized

► In RDBMS:

- Parse → Rewrite → Plan → Execute

► Optimizations include:

Technique	Example
Selection Pushdown	Apply WHERE filters early
Join Reordering	Smallest tables first
Index Usage	WHERE email = 'x' uses index
Subquery Flattening	Convert nested queries to joins
Avoid Full Scans	Use indexes or LIMITs

Query Performance Tips (Industry)

Tip	Description
Use EXPLAIN	Shows query plan
Avoid SELECT *	Fetch only needed fields
Index foreign keys	Speeds up joins
Use pagination	Don't load 10k rows
Write WHERE conditions early	Helps filtering
Avoid nested subqueries when possible	Use JOINs instead

Example:

-- BAD (slow)

```
SELECT * FROM Orders WHERE YEAR(order_date) = 2024;
```

-- BETTER

```
SELECT * FROM Orders  
WHERE order_date BETWEEN '2024-01-01' AND '2024-12-31';
```

Summary of Unit III

Topic	College	Industry
Normalization	 Core design principle	 Clean DB schema
Denormalization	 (not in syllabus)	 Scale, performance
Query Optimization	 Basics needed	 Required daily in backend/data teams

◆ Unit IV: Transactions, Concurrency Control & Recovery

(College + Industry Level)

◆ 1. What is a Transaction?

A **Transaction** is a logical unit of work in a database — a sequence of operations performed as a single unit.

Example:

```
BEGIN;
```

```
UPDATE Bank SET balance = balance - 1000 WHERE acc_no = 101;
```

```
UPDATE Bank SET balance = balance + 1000 WHERE acc_no = 202;
```

```
COMMIT;
```

- This is a **transaction** — money is transferred only if both operations succeed.
-

◆ **2. ACID Properties of Transactions**

Property	Meaning	Industry Use
A – Atomicity	All or nothing	No partial updates
C – Consistency	DB moves from one valid state to another	Maintains rules
I – Isolation	Transactions don't interfere	Avoid dirty reads
D – Durability	Once committed, it's permanent	Crash-proof

💡 **ACID** is a **MUST** for all banking, financial, or critical systems.

◆ **3. Concurrency Control**

Multiple users may access the DB **at the same time** → must ensure correct results!

⚠ **Problems Without Control:**

Problem	Description
Lost Update	Two transactions overwrite each other
Dirty Read	Read uncommitted changes
Non-Repeatable Read	Data changes mid-transaction
Phantom Read	Row appears/disappears on re-query

🔧 **4. Concurrency Control Techniques**

Technique	Description
Locks (most common)	Prevents simultaneous access
Timestamps	Assigns global time to transactions
Optimistic Concurrency	Assume no conflict, verify before commit

Technique	Description
Multiversion Concurrency Control (MVCC)	Read consistent snapshots (used in PostgreSQL, MySQL InnoDB)

Locks in Detail

Type of Lock	Purpose
Shared Lock (S)	Read only
Exclusive Lock (X)	Read + Write
Deadlock	Two transactions waiting on each other

 **Industry Tip:**
Backend developers often use **MVCC** and **lock-free** mechanisms for performance.

◆ 5. Isolation Levels (SQL Standard)

Control **how isolated** transactions are. Higher isolation = fewer anomalies but more locking (slow).

Level	What It Prevents	Allows
Read Uncommitted	 Nothing	Dirty Reads
Read Committed	 Dirty Reads	Non-repeatable Reads
Repeatable Read	 Dirty, Non-repeatable Reads	Phantom Reads
Serializable	 All	Slower but safest

→ You choose based on **system needs**. E.g., banking = Serializable. Analytics = Read Committed.

6. Transaction States

Active → Partially Committed → Committed (or Failed → Aborted → Restart)

7. Recovery Techniques (Very Important for Industry)

What if the system crashes mid-transaction?

Recovery Tool	Description
Logs (Write Ahead Logging)	All changes written to log before DB
Checkpointing	Snapshot of DB taken at intervals
Shadow Paging	Keep old & new versions of data
Undo/Redo Operations	Revert or replay logs during recovery

Real-World Use

Industry Case	Feature Used
Bank Transfer	Atomicity, Durability
Ride Booking (Ola/Uber)	Concurrency, Isolation
Payment Gateways	ACID, WAL logs
Microservices	Distributed Transactions, Eventual Consistency

Summary of Unit IV

Topic	College	Industry
Transactions	<input checked="" type="checkbox"/> Core concept	<input checked="" type="checkbox"/> Required in all apps
ACID	<input checked="" type="checkbox"/> Theory	<input checked="" type="checkbox"/> Real systems
Concurrency Control	<input checked="" type="checkbox"/> Types & problems	<input checked="" type="checkbox"/> Locks, MVCC in practice
Isolation Levels	<input checked="" type="checkbox"/> Names only	<input checked="" type="checkbox"/> Pick levels per use-case
Recovery Techniques	<input checked="" type="checkbox"/> Checkpoints/logs	<input checked="" type="checkbox"/> WAL, failover systems

◆ Unit IV: Transactions, Concurrency Control & Recovery

(College + Industry Level)

◆ 1. What is a Transaction?

A **Transaction** is a logical unit of work in a database — a sequence of operations performed as a single unit.

Example:

BEGIN;

UPDATE Bank SET balance = balance - 1000 WHERE acc_no = 101;

UPDATE Bank SET balance = balance + 1000 WHERE acc_no = 202;

COMMIT;

→ This is a **transaction** — money is transferred only if both operations succeed.

◆ 2. ACID Properties of Transactions

Property	Meaning	Industry Use
A – Atomicity	All or nothing	No partial updates
C – Consistency	DB moves from one valid state to another	Maintains rules
I – Isolation	Transactions don't interfere	Avoid dirty reads
D – Durability	Once committed, it's permanent	Crash-proof

 **ACID** is a **MUST** for all banking, financial, or critical systems.

◆ 3. Concurrency Control

Multiple users may access the DB **at the same time** → must ensure correct results!

Problems Without Control:

Problem	Description
Lost Update	Two transactions overwrite each other
Dirty Read	Read uncommitted changes
Non-Repeatable Read	Data changes mid-transaction
Phantom Read	Row appears/disappears on re-query

4. Concurrency Control Techniques

Technique	Description
Locks (most common)	Prevents simultaneous access
Timestamps	Assigns global time to transactions
Optimistic Concurrency	Assume no conflict, verify before commit
Multiversion Concurrency Control (MVCC)	Read consistent snapshots (used in PostgreSQL, MySQL InnoDB)

Locks in Detail

Type of Lock	Purpose
Shared Lock (S)	Read only
Exclusive Lock (X)	Read + Write
Deadlock	Two transactions waiting on each other

 **Industry Tip:**
Backend developers often use **MVCC** and **lock-free** mechanisms for performance.

◆ 5. Isolation Levels (SQL Standard)

Control **how isolated** transactions are. Higher isolation = fewer anomalies but more locking (slow).

Level	What It Prevents	Allows
Read Uncommitted	 Nothing	Dirty Reads
Read Committed	 Dirty Reads	Non-repeatable Reads
Repeatable Read	 Dirty, Non-repeatable Reads	Phantom Reads
Serializable	 All	Slower but safest

→ You choose based on **system needs**. E.g., banking = Serializable. Analytics = Read Committed.

6. Transaction States

Active → Partially Committed → Committed (or Failed → Aborted → Restart)

7. Recovery Techniques (Very Important for Industry)

What if the system crashes mid-transaction?

Recovery Tool	Description
Logs (Write Ahead Logging)	All changes written to log before DB
Checkpointing	Snapshot of DB taken at intervals
Shadow Paging	Keep old & new versions of data
Undo/Redo Operations	Revert or replay logs during recovery

Real-World Use

Industry Case	Feature Used
Bank Transfer	Atomicity, Durability
Ride Booking (Ola/Uber)	Concurrency, Isolation
Payment Gateways	ACID, WAL logs
Microservices	Distributed Transactions, Eventual Consistency

Summary of Unit IV

Topic	College	Industry
Transactions	<input checked="" type="checkbox"/> Core concept	<input checked="" type="checkbox"/> Required in all apps
ACID	<input checked="" type="checkbox"/> Theory	<input checked="" type="checkbox"/> Real systems
Concurrency Control	<input checked="" type="checkbox"/> Types & problems	<input checked="" type="checkbox"/> Locks, MVCC in practice
Isolation Levels	<input checked="" type="checkbox"/> Names only	<input checked="" type="checkbox"/> Pick levels per use-case
Recovery Techniques	<input checked="" type="checkbox"/> Checkpoints/logs	<input checked="" type="checkbox"/> WAL, failover systems

◆ Unit V: File Organization and Indexing

(College + Industry Level)

◆ 1. File Organization

How data is physically stored on disk — affects performance, search time, insert/update/delete speed.

Types of File Organization

Type	Description	Usage
Heap File	Records stored in any order	Fast insert, slow search
Sequential File	Records stored in sorted order	Good for range queries
Hashed File	Uses hashing to place data	Super-fast for exact-match search

Real-World Analogy:

- **Heap** → Like putting random clothes in a bag
 - **Sequential** → Neatly stacking books in order
 - **Hashed** → Using pigeonhole boxes labeled with student IDs
-

◆ 2. Indexing

Indexing = Data structures that make searching **faster**, like a book's table of contents.

🔍 Why Index?

Without index:

`SELECT * FROM students WHERE roll_no = 5;`

→ Needs to check *every row* (full scan).

With index:

→ Direct jump to correct row!

◆ Types of Indexes

○ Single-Level Index

- A sorted list of keys with pointers to data blocks.
 - Like a **phonebook** – fast lookup.
-

牢固树立 B+ Tree Index (Very Common in SQL)

Feature Description

Balanced Same depth for all leaves

Sorted Good for range queries

Multi-level Fast search in $\log(n)$

Used in MySQL, PostgreSQL, Oracle, etc

→ **Industry Standard for indexing**

⚡ Hash Index

Feature Description

Uses Hashing Direct address lookup

Feature	Description
Super-fast for = queries	Not good for range queries
→ Used in:	In-memory databases like Redis or point queries in MySQL.

Clustered vs Non-Clustered Index

Type	Description	Example
Clustered Index	Sorts the data itself	Primary key
Non-Clustered Index	Separate from data	Secondary keys

Index in SQL (Industry-Level)

```
CREATE INDEX idx_name ON student(name);
CREATE UNIQUE INDEX idx_roll ON student(roll_no);

👉 Helps queries like:
SELECT * FROM student WHERE name = 'John';
```

Optimization Strategy

Goal	Technique
Fast inserts	Heap / Hash
Fast lookups	B+ Tree Index
Range queries	B+ Tree
Random access	Hash Index
Massive data	Clustered index

Real-World Example

E-commerce site

- Products table: indexed by product_id (hash), category (B+ Tree)

- Orders: indexed by user_id, order_date for fast sorting & querying
-

Summary of Unit V

Concept	College	Industry
File Org	<input checked="" type="checkbox"/> Basic theory	<input checked="" type="checkbox"/> Used for DB engine tuning
Heap/Sequential	<input checked="" type="checkbox"/> Must learn	<input checked="" type="checkbox"/> Storage layer optimization
Indexing	<input checked="" type="checkbox"/> B+ Tree	<input checked="" type="checkbox"/> Real query performance boost
Hashing	<input checked="" type="checkbox"/> Concept	<input checked="" type="checkbox"/> Used in key-value stores
SQL Index Use	<input checked="" type="checkbox"/> Often skipped	<input checked="" type="checkbox"/> Must learn as dev

DBMS Completion Recap

- Unit I – Intro + ER Models
- Unit II – Relational Model, SQL, Joins
- Unit III – Normalization + Optimization
- Unit IV – Transactions & Concurrency
- Unit V – Indexing + File Organization

You've now completed both **college syllabus** and **industry-level DBMS** foundations! 