

DBMS

- Data

raw bytes

can be in form of integer, float, int etc
unorganized.

- Information

processed data is information.

Information helps in decision making.

When data gets a meaning, it is information.

- Database

electronic system that stores the data.

It can be bank or anything.

- DBMS

Data base Management System.

Made of two things.

- Database

Management System → It provide functionality
to modify data.

Add, update, delete.

- Need for DBMS

DBMS helps us with file system flows as.

i) Data redundancy / inconsistency.

If one thing gets changed or updated then it is reflected to every place. This is called consistency.

Storing same data on multiple places causes duplicates. This is redundancy.

ii) Accessible.

Data can be accessed efficiently and quickly by a simple query.

iii) Atomicity.

Either a yes or no. There is nothing in between.
e.g. credit and debit.

iv) Security

The security is maintained. There are methods to grant or take access of data.

v) concurrent access.

It means multiple queries can run together without any trouble.

vi) Isolation of data.

Data can be in different formats. It leads to data manipulation issues.

vii) Integrity

There are multiple constraints that help in functionality. e.g. balance can not be below 10,000.

• Abstraction

It is a technique used to hide data from different users.

It follows three level architecture.

i) External level

Highest level of abstraction.

Also called view schema.

This is used by end users.

There are multiple schema for different user groups.

These schema are called subschema.

Eg. view for logistics team ≠ view for customer service.

Security measures can also be implemented.

ii) Conceptual / logical level

It is also called logical schema.

It tells us what data and its relationship.

It does not focus on how data is stored.

It is used by DBA.

iii) Internal level / Physical level

It is the place where data is stored.

It tells us how data is stored.

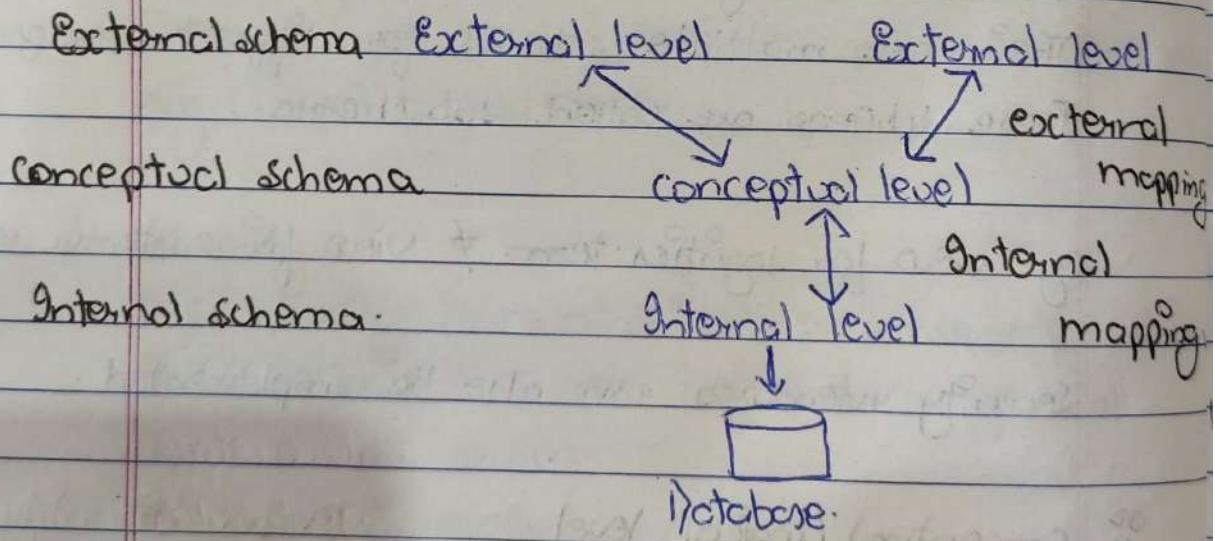
Change in physical level does not affect logical level.

It has physical schema that describe physical storage structure.

Its goal is to make data access efficient.

Note Converting data from physical to logical level is called mapping of data. Internal mapping.

Architecture



• Schema

Blueprint in DBMS that defines how data is organized in database.

It describes tables, fields, columns, views etc.

In short, logical structure of database.

• Instance

Total data points stored in database at a point of time is instance of DB.

i) DBA

Database Administrator.

person who has the control of the system.

It has some functions like:-

- 1 schema definition.
- 2 storage structure and access method
- 3 Authorization control
- 4 Routine maintenance

backups

security patches

Upgrades.

• DBMS System Architecture

i) Tier 1 Architecture.

It states that client, server and Database, all three are stored in single place.

e.g. while learning SQL

ii) Tier 2 Architecture.

It has two components.

Application sides sends query to the database system.

User can directly access the database.

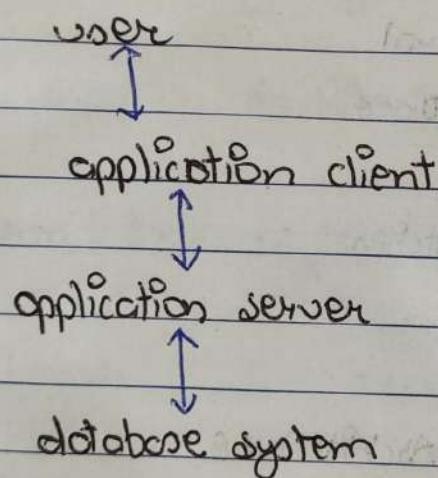
API like ODBC and JDBC are used.

Tier 3 Architecture

It has three parts:

client machine has frontend and communicate to application server.

Application server acts as a mediator between client side and database system.



These architecture are used for large scale system like amazon.

Advantage:

Scalability -> multiple client and application server can be set up.

Integrity -> It does not affect security as client can not connect to database directly.

Data models

These are the models that defines the structure of data, relationships, constraints.

ER Model

It stands for entity relationship model

entity

It is any real world object that has some attributes of itself.

e.g. student is an entity.

It is denoted as

student

entity set

This is a collection of similar type of entities.

e.g. There can be multiple students.

Attributes

These are the features of the entity. e.g.

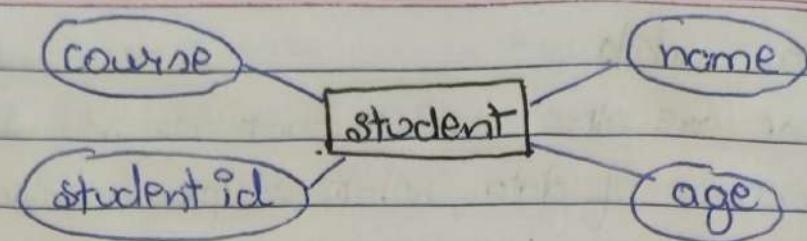
Student has several attributes i.e. student name, age, id etc.

These are different for different students.

These are denoted as

name

∴



Note The attributes used to uniquely define an entity is called primary key.
e.g. student id is primary key
It is denoted as

student id

Primary key is unique for all and can not be null.

Types of Attributes.

• Simple

These are single attributes that can not be further broken down.

e.g. Aadhar number.

• composite

These can be further broken down.

e.g. name broken into fname, lname.

• single valued.

These attributes can have only one value.

e.g. student id.

- Multi valued

These can have multiple values for one entity.

e.g. phone no.

Phone no

- Derived.

These values can be derived from other attributes.

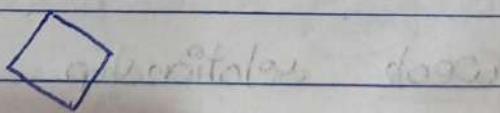
e.g. age derived from date of birth.

age

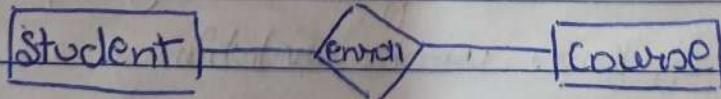
Relationships

These are used to define how two entities are related to each other.

It is denoted by.



e.g.



Strong entity

These entities can be uniquely defined.

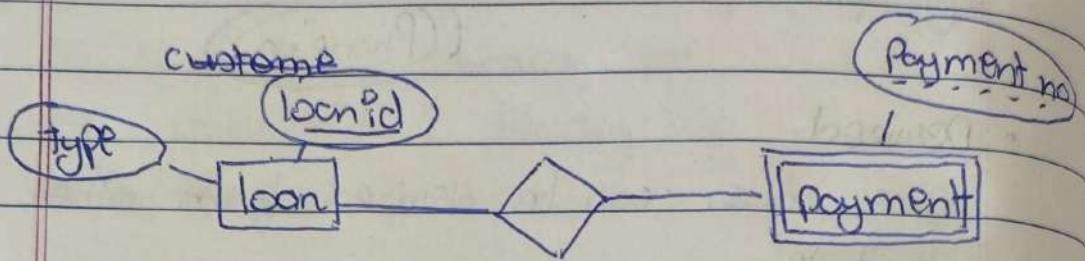
They are independent.

e.g. student can be defined by studentid

weak entity.

These are dependent on other entities for existence.

e.g.



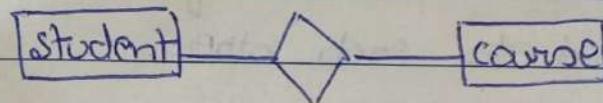
Here, payment depend on loan.

Therefore, they are defined as double rectangle.

strong Relationship

Both entities exist independently.

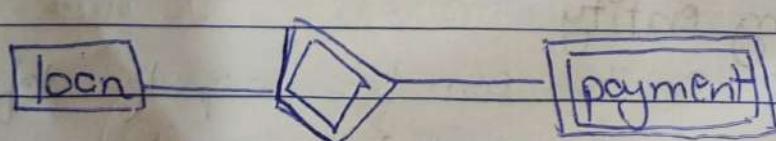
e.g.



weak relationship

When one of the entity depends on other, then it forms weak relation.

e.g.



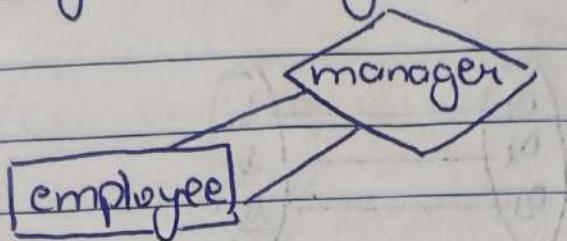
It is denoted by double rhombus.

Types of Relationship

- Unary relation

It has only one entity.

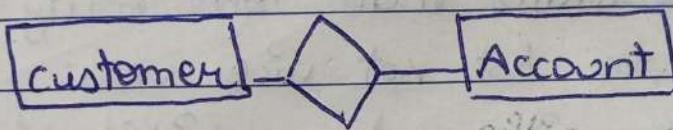
e.g.



- Binary Relation

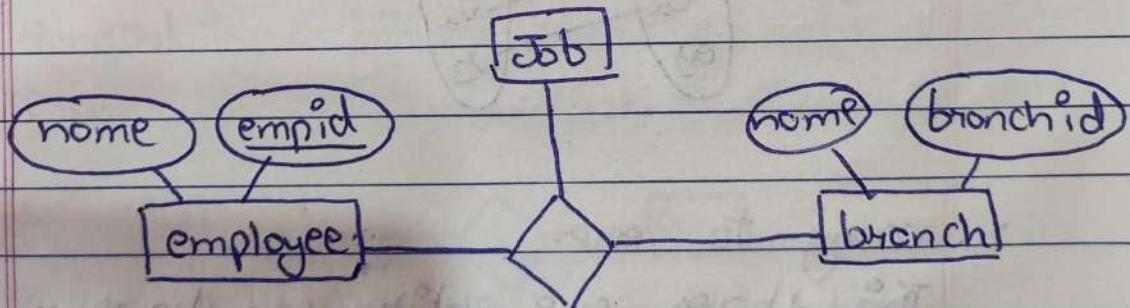
It has two entity.

e.g.



- Ternary relation.

It has three entities.



Relationship constraints

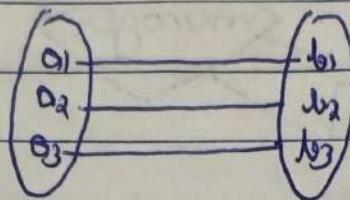
- Mapping cardinality.

This tells us how many entities in an entity set are connected to no. of entities of other entity set.

- one to one.

This states that one entity with exactly one entity.

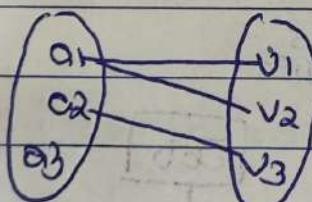
e.g. citizen has address.



- One to many

This states that one entity can have multiple values but not vice versa.

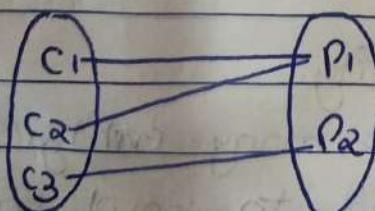
e.g. citizen has vehicle.



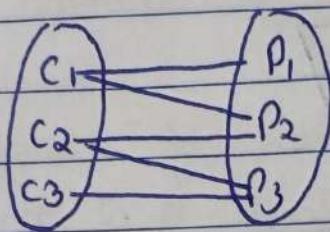
- many to one.

This shows one entity can have multiple values but not vice versa.

e.g. course by professor.



- many to many.
n entities of A can have relation to n entities of B.
e.g. customer buy product



\Rightarrow Participation constraint.

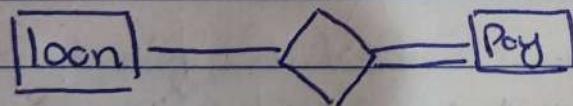
It shows us the minimum cardinality constraint.

• Total participation.

This happens when all entities of entity set participate.

denoted by =====

e.g.



• partial participation.

This shows that not all entities participate.

e.g.



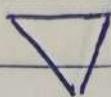
- Specialization

Top-down approach.

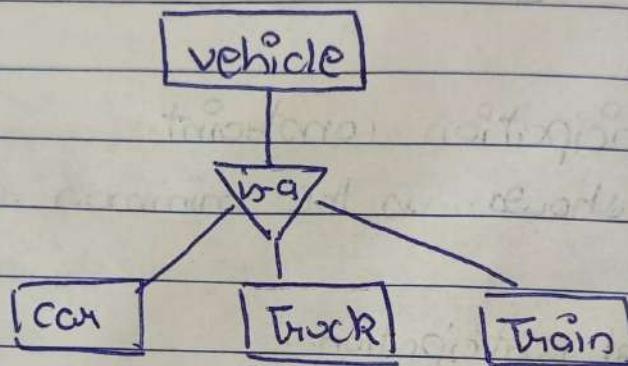
We break down larger entities into smaller.

It has 'is-a' relation.

It is denoted



e.g.



why?

To remove redundancy.

It is like inheritance.

Scenario.

Some attributes are common in above example for all three entities like speed.

∴ we remove speed from them and add it to vehicle.

This helps to remove data redundancy.

- Generalization

Bottom-up approach

We group sub entities into larger entities.

It has 'is-a' relation.

It is opposite of specialization.

Note ER diagram for both are same, just the difference is of thinking.

If smaller group first then generalization, else specialization.

- Attribute inheritance.

Both specialization and generalization have it.

Common attributes are inherited to through parent entity into child entity.

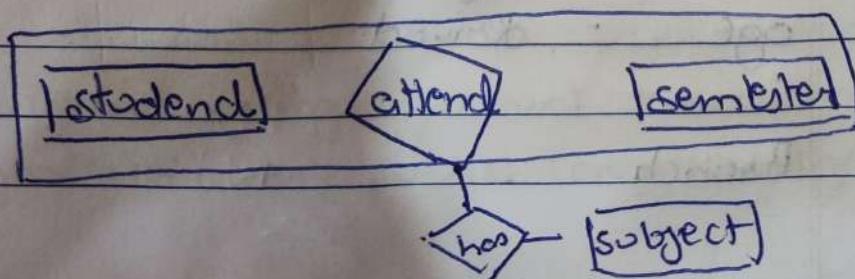
- participation inheritance.

If parent entity involved in a relation then its child also gets involved.

- Aggregation

Combining entities relation to a higher entity

This introduces abstraction.



How to draw er diagram.

There are three steps to draw er diagram

- Identify the entities
- Identify Attributes and their types.
- Identify relationship and constraints

Mapping

Participation

1 entities

Banking System

1 entities

customer, branch, employee, saving acc,
current account, loan, payment

2 Attribute and types.

• customer

name - multivalue composite

cust id - Primary key, single

address - composite

phone - multivalue

DOB - single

age - derived

• Branch

branch id - Primary Key

accts

liabilities

name - single

city - single

- employees

name - composite

emp_id - Primary key

phone - single

year of service - derived.

dependent name - multivalued

start date - single.

- saving acc -

account no - Primary key

balance - single

daily withdrawal limit - single

interest - single.

- current acc

account no - Primary key

balance - single

per transaction charge - single.

Note: current and savings account have some common attributes so we can generalize them.

- generalized entity - Account

account no - account no

balance - single.

- loan

loan id - Primary key

amount - single.

- payment

payment number

payment date

amount

Note loan and payment will have weak relationship.

3 Relationships.

- Customer deposit account

M : N

- customer has loan

M : N

=

- loan has payment

M : N

=

• employee manage employee.

$N : 1$

loan originated branch

$N : 1$

=

• Relational Model

storing the data in the form of tables.

entity are converted to table names.

Attributes converted to columns.

row/tuple define a single entity.

column/Attribute contain a domain specific information.

Note degree of table is no. of columns

cardinality is no. of rows.

steps to draw a database system.

i) draw er model

ii) convert to relational model

iii) implement RDBMS to relational model.

RDBMS : MySQL, Microsoft access, Oracle

- Properties

unique name for all tables.

unique name for all columns in a table

sequence of columns does not matter.

table must follow consistency constraints.

values should be atomic

each row should be unique.

Keys

- Super Key → All the combination of columns that can be used as primary key.

e.g.

custid	name	email	contact

Super key :- cust id, [name email],
[name, contact], [email], contact

etc.

- Candidate Key.

subset of super key.

cannot be null

do not have columns with redundant type.

e.g. → cust id, [email] (contact)

\therefore we removed name as it is redundant.

- Primary Key

It is a candidate key with least no. of attribute.

It can not be null.

It is unique.

e.g. cust id.

Note every primary key is a candidate key but every candidate key is not a primary key.

- Alternate Key

All combination of candidate key left after choosing primary key.

$$\therefore \{CK\} - PK = AK.$$

- Foreign Key

It is a primary key that is used to join two tables.

every foreign key is a primary key of other table.

The table whose primary key is used as foreign key is called Parent/referenced table.

The table who takes primary key as foreign key is called child / referencing table.

- composite key

Primary key formed by at least 2 attributes

- compound key

Primary key formed by 2 foreign keys

- surrogate key

used to join two tables where nothing can be used in place of Pk.

usually holds an integer.

can be used as Primary key

Integrity Constraints

These are important to maintain consistency in database when CRUD operations are performed.

- Domain constraints

To have a check on domain specific column

e.g. id should be int, age ≥ 18

- entity constraint

every table should have a Pk.

Pk != null

- referential constraint

- insert constraint

cannot insert in child relation table if not present in parent relation table.

- delete constraint

cannot delete from parent relation table if entry present in child relation table.

Can we not delete then from parent?

Yes we can by using On delete cascade.

This delete the entry related to parent table and child table at once.

e.g. if we delete custId = 3 from parent relation then it gets deleted from child relation.

can the foreign key be ~~zero~~ ^{Null}?

Yet Yes, by using on delete Null.

This do not deletes entry in child table instead set the foreign key value to Null.

Key constraints

These are the constraints that are used to manipulate the data.

i) Not null

ensures that value in a column is not null.
By default value of a column is null.

cust_id int NOT NULL

ii) Unique.

ensure every value is unique in the column.
There can be multiple unique constraints.

cust_id int,

unique (id)

iii) default

This defines the default value of a column.
It is added if no value is specified.

iv) check.

It is used to apply conditional constraints.

e.g.

age int,

check age >= 18

v) Primary Key

Set of attributes or single attribute that uniquely define a table.

There can only be one Primary Key.

i) Foreign Key.

It is used to join two tables.

The Foreign key must be primary key of other entity set.

e.g.

Primary key (order ID),

Foreign key (cust ID) Referencing
customer (cust ID)

Transformation from ER To RM

i) strong entity

convert to a separate table.

ii) weak entity

convert to a table.

introduce FK ie primary key of strong entity.
Primary key of this table is FK + PK.

iii) composite attribute.

Create separate table.

All attributes included.

separate column for every possible attribute

e.g. address : streetno, streetname, city, state.

iv) multivalued attribute.

Create new table with same name as attrib.

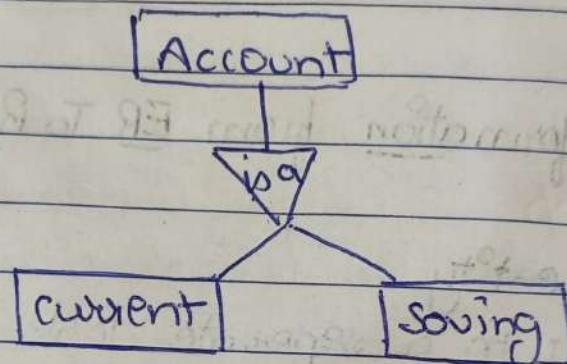
introduce FK ie primary key of main table
FK + Pk distribute in Pk.

ii) Generalization

• Method 1

separate tables for all entity.

e.g. if



Then three tables.

account : acc_no, bal

curr_acc : acc_no, . . .

Scv_acc : acc_no, . . .

• Method 2

create only two tables for lower entity.

curr_acc : acc_no, bal, . . .

Scv_acc : acc_no, bal, . . .

Problem with method 2.

- if an account is both savings and current then redundancy.

- if the account is neither current nor saving
then no way to fetch.

v) Aggregation

convert the relation into table. Take look on
big of aggregation.

All the primary keys of all entities become
attribute.

The PK is combination of all attributes ie a
composite key.

manager (emp-id, manager-id, job-id, branch-id)

PK \rightarrow combination of all attributes.

vii) Unary relation.

add another column, the FK ie the primary
key of some table.

This is for 1:N

• 1:1

Some like 1:N

but primary key is only one attribute ie
PK itself.

• M:N

second table of relationship.

include both id ie both are primary key of main table.

These both will act as PK.

Here the PK will be combination of both PK compound key.

Normalization

Functional dependency.

When one attribute can be determined by other then it is functional dependency.

e.g. $A \rightarrow B$

$\{ \text{emp_id} \} \rightarrow \{ \text{emp_name} \}$

Here A is determiner

B is dependent

Types.

• Trivial

It states that if $A \rightarrow B$ has Trivial FD if $B \subseteq A$.

$\{ \text{emp_id}, \text{emp_name} \} \rightarrow \{ \text{emp_name} \}$

- non trivial

it states if $A \rightarrow B$ but $B \not\subseteq A$

e.g.

$$\{\text{emp_id}\} \rightarrow \{\text{emp_addr}\}$$

Armstrong's Axioms Rules

- Reflexive

if A is set of attribute and $B \subseteq A$ then
 $A \rightarrow B$.

e.g. $A \rightarrow \{a, b, c, d\}$, $B = \{b, c\}$

$$\therefore A \rightarrow B$$

- Augmentation.

if $A \rightarrow B$ then

$$A \cup Z \rightarrow B \cup Z$$
 is also true.

- Transitivity.

if $A \rightarrow B$ and $B \rightarrow C$ then $A \rightarrow C$

Normalization

It is a technique to remove the redundancy in the data to make it efficient.

Why Normalization?

To remove redundancy.

What if not normalization?

If we do not normalize data then we can encounter anomalies.

- insertion.

We cannot insert data until the all data is fetched.

- deletion

Cannot delete data properly.

Deletion of data may cause loss of other data.

- updation.

Data updation will require repeating steps again and again.

It also increase size of database.

∴ Decompose tables until SPP is achieved.
Single responsibility Principle.

Types

1NF

It states that there should only atomic values
no multivalued attributes allowed.

2NF

Table should be in 1st normal form.

There should not be partial dependency.

All non-prime attributes should be fully dependent
on prime attribute.

Partial dependency.

Non prime attribute cannot be determined
by part of prime attribute.

Eg. If R1 {A,B,C,D}

AB \rightarrow Primary key, then

AB \rightarrow C ✓

B \rightarrow C X Partial dependency.

3NF

must be in 2NF

No transitive dependency.

Non prime cannot determine non prime.

e.g. if $R_1(A, B, C)$ then,

$A \rightarrow B$ ✓

$B \rightarrow C$ ✗

- BCNF

Boyce - codd Normal form.

Must be 3NF

prime attribute cannot be on right side.

if $A \rightarrow B$ then, A must be a super key.

Non prime cannot bind prime or non prime.

ACID Properties And Transactions

Transaction

Amount of work done that includes multiple steps executed in a sequence.

This work should either be successful or not successful.

e.g.

A [1000] O [2000]

if A wants to send 100\$ to O then

read (A) $\rightarrow A = A - 100 \rightarrow$ write (A) $\rightarrow 900$

read (B) $\rightarrow B = B + 100 \rightarrow$ write (B) $\rightarrow 3100$

\therefore When all these steps are completed, we say the transaction is successful.

There are few properties that ensures integrity.

- **Atomicity.**

Transaction should either be completed or not completed. If system crashes in middle then option to roll back to old state.

- **Consistency.**

The consistency should be maintained. Therefore the balance sum of A and B should be equal to amount before transaction and after transaction.

$$\therefore (A+B) \text{ before } T_1 = (A+B) \text{ after } T_1$$

- **Isolation.**

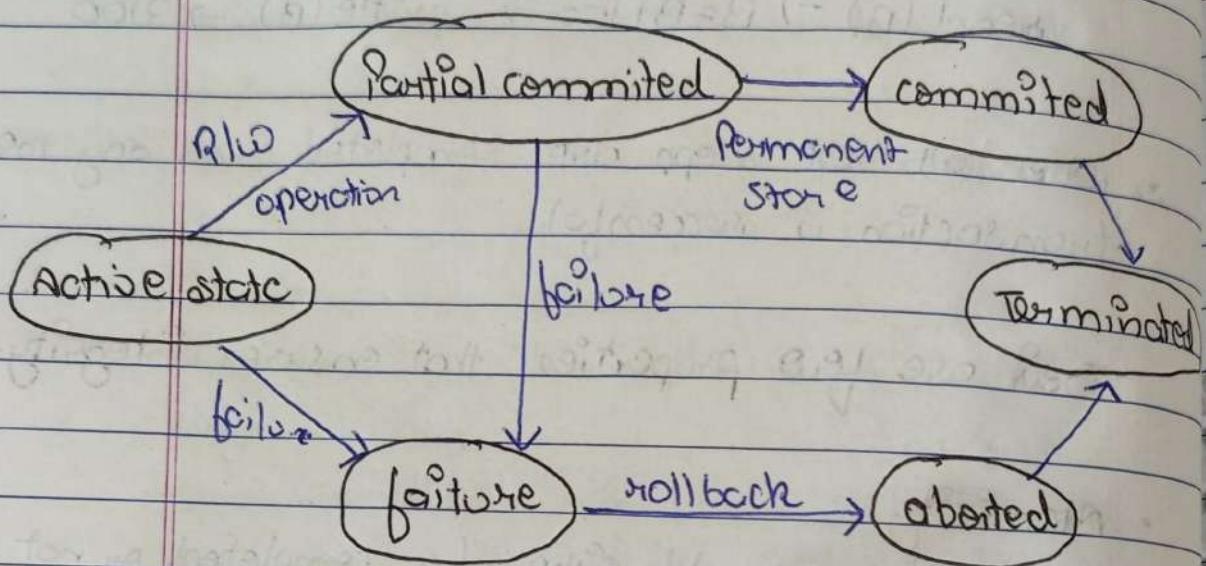
There are N number of transaction.

All transactions should be executed one after other.

- **Durability.**

Changes in DB should be permanent.

Transaction states



- Active state.

First state of transaction

R/W operations perform.

If executed then move to Partial committed

If not then failure.

- Partial committed

Here the changes are stored in buffer of main memory.

If changes permanent then to committed state
else to failed state.

- Committed state.

Here changes are made permanent and then move to terminated state.

- Failed state

This is reached if any error is encountered.
Then the steps are rolled back to original state.

- Aborted state

If roll back is achieved then aborted state is reached.

- Terminated state

It is reached either after committed or aborted state.

How to implement Atomicity?

- Shadow copy scheme:

We create a copy of data and perform operations on the copy.

While the db pointer still points on to the old data.

If all the operations are successfully performed then db pointer points to new copy.

If error then new copy is destroyed.

- For next transaction the new copy act on the old copy.

Log recovery Method

In this method we write logs side by side with each step.

Logs are written before the execution of statement.

- deferred db method

We commit afterwards in this method.

e.g.

read(A)

<to start>

A = A - 50

<to, A, 950>

write(A)

<to, B, 2050>

<commit>

If transaction fails before commit then logs are destroyed.

If after commit then these are used to redo.

- immediate db method.

In this method we directly write on disk and logs are written like

<to start>

<to, A, 1000, 950>

<to, B, 2000, 2050>

Here old value is also stored to redo in case of error.

Indexing

It is a technique to optimize the searching in the table.

It is used to increase the efficiency.

We create an index file that has a sorted key or non key attributes and entries are divided into batches.

Note Index file is always sorted.

It is not necessary that data should be sorted only on PK. non prime attributes can also be used.

Types

- Primary indexing.

It states data should be sorted on any attribute whether primary or non prime.

e.g. sorted on basis of roll no.

Index file

I B₁

II B₂

III B₃

Dense index

When every unique value is entered in the index file then it is called dense index.
Also known as clustering index.

Clustering index

It is also a primary index but here non prime attribute is used as index.

Sparse index

When only one value is listed in index table representing other values of a block it is called sparse index.

It is usually used with primary attribute.

- Multi-level index

Index with two or more levels.

It is done when index table gets too large.
It is sparse index.

- Secondary index (non-clustering index)

It is used on unsorted data.

It is a dense index example.

It can be done on prime or non prime attribute.

If values are not unique then we can maintain linked list of those pointers.

Index file

1 $B_1 \rightarrow B_3 \rightarrow B_9$

2

3

It is helpful as index file is always sorted.
We cannot apply binary search on data file
but can apply on index file.

NoSQL (Not only SQL)

It stores data other than in tables.

One famous example in JSON format.

MongoDB uses NoSQL.

Why NoSQL?

- **flexible schema:**

It provides flexible schema as it does not have any structure.

We can add information that is not complete.

- **Horizontal scaling (scale out)**

It has important feature of horizontal scaling, i.e. dividing the data and storing it in multiple servers.

Note SQL uses vertical scaling (scale up)

Vertical scaling

When upgrades are made on single system
to increase the space.

Hardware, RAM are upgraded.

Fast data retrieval

It helps in fast data access as all the data
can be stored in single file, ie no normalization.
Insert and read operations are fast.

Delete and update operations are time taking

Note SQL can use Horizontal scaling but it is not
practical.

High availability

Data is stored in multiple servers due to its
auto replication feature. Thus if one server
crashes then other server can be used.

Cloud

More suitable for cloud servers

Caching Mechanism

Stores frequently accessed data in a fast access
storage.

Eg. chats in social media app.

When to use NoSQL?

Huge amount of data
fast development

when scale out feature is required

Storage of unstructured data

Types

i) Key-value

They have key value pairs. One pair of key and value.

e.g. cards in amazon, flipkart

ii) C-store

store data in column format

[Dev | Kochipuri | 30 | Mutt | Pune | 20] → noSQL store

[Deve | Mutt | Kochipuri | Pune | 30 | 20] → column store

This is suitable for analytics

Insertion is slow but data access is fast

iii) Document based

MongoDB is an example.

Similar to json format

Multiple values can be stored

e.g. e-commerce platform.

↳ graph based.

It is suitable for storing relationships without using joins.

Uses graph structure to store data.

e.g. Fb friends structure.

1) Disadvantages of NOSQL

- Data redundancy. ∴ storage high.
- Slow delete and update and costly.

Note: MongoDB supports acid properties.

Types of Database

• Relational databases.

• NOSQL databases.

• Object oriented databases.

based on object oriented programming.

Used when establishing relationship gets too complex.

It supports encapsulation, inheritance etc.

Stores information in a class and object is used to access it.

advantage

Data retrieval is fast.

can handle complex data relations.

friendly with real world problems.

Disadvantage

Highly complex

Less community support

Does not support views like relational database.

e.g. objectDB, GemStone etc.

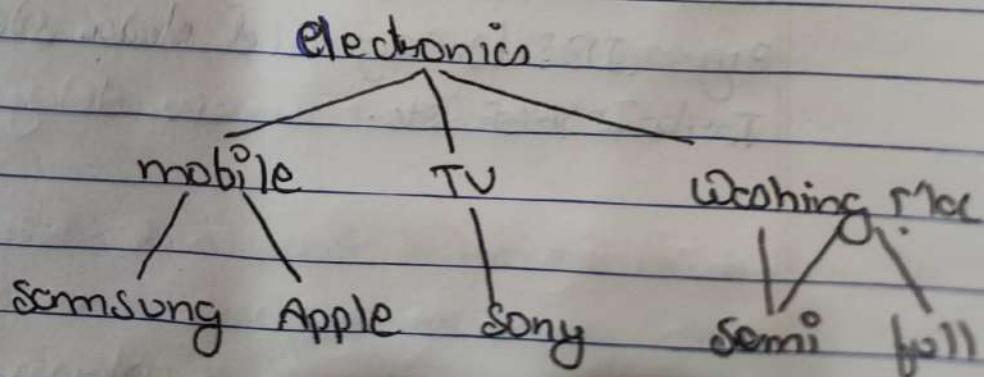
- Hierarchical Databases

Tree like structures.

Used when top-down data retrieval is required in specific use case.

One child can only have one parent.

e.g. Electronics shop.



It is easy to use.

easy to update and insert.

Disadvantage

not flexible; can only have one parent.
not much community support.

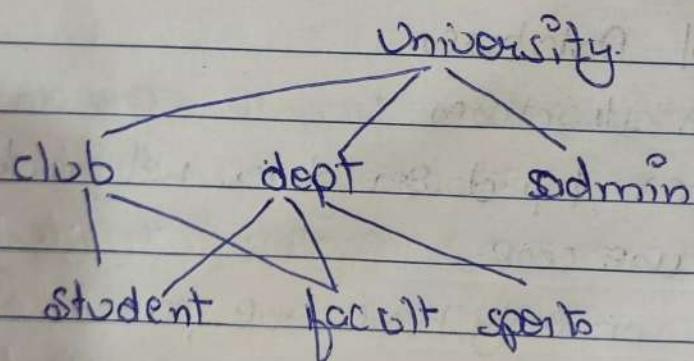
• Network Databases.

Extension to Hierarchical databases.

can have more than one parent.

Graph like structure.

e.g.



These can handle complex relations.

Not much community support.

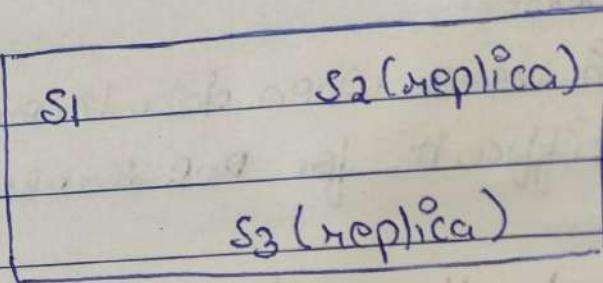
e.g. IDS (Integrated data store), IORMS, TurboIMAGE etc.

clustering

It is also called replica set.

It is a process of combining more than one server that connects a single database.

e.g. There is a server S1 with info about 1000 users, then



clustering.

∴ All three servers contain some information.

Advantage

- Data redundancy.

If one system crashes then the request can be fulfilled by other server.

It also provides a layer of abstraction.

- Load balancing.

We introduce a load balancer in order to balance the load of incoming request.

It also improves the scalability.

- High availability.

Due to clustering, the website is always working and request can be fulfilled instantly.

- Partitioning

It is a process to divide data into multiple servers.

It is done when data becomes huge and it is difficult for one server to manage it.

What all can be done for optimization?

- scale up

vertical scaling can be done

- but it increases the cost.

- clustering

- It is managed by a main server and replica sets.

updates are first on main server then to replica sets that causes delay.

- Partitioning

It allows two types-

- vertical partitioning - column wise

- horizontal partitioning - row wise.

Distributed database.

single logical database spread across multiple servers and logically interconnected by network.

• Sharding

It is a technique to implement horizontal partition.
It requires a routing layer to help in sending request to appropriate shards (servers).

Advantage

Scalability

Availability.

Disadvantage.

Increase complexity cause of routing layer.

Non-uniformity

Not suitable for analytical queries like aggregation.

Note Advantages of Partitioning

Parallelism

Manageability

Reduce cost of vertical scaling

Availability.

Performance.

CAP Theorem

It is one of most important concept in distributed databases.

- consistency

Value of a record must be same across nodes.

- Availability

The system should remain operational all the time.

- Partition tolerance.

When partition occurs, the communication between the nodes is lost.

If system is partition tolerant then it does not fail even after partition.

Note CAP theorem states that distributed system can only provide two of three properties at a time.
∴ There is a trade off between consistency and availability when there is a partition.

CA Database: These provide both consistency and availability but there is no partition.
e.g. MySQL, PostgreSQL.

CP database: Consistency, Partition tolerance.

When a partition occurs, the system turns off inconsistent node that contradicts the availability.
e.g. MongoDB, Banking System.

AP database: Availability, Partition tolerance.

When partition occurs, these systems allow both nodes to run which leads to loss of consistency.

Though, these systems ensure consistency after the partition is fixed and nodes get sync.

Eg → Cassandra, DynamoDB.

App like FB

Master Slave Architecture

There is a master database and many slave databases.

Master DB handles all write operations.

Slave DB handles all read operations.

Once write is performed in master DB, then it is replicated to slave DB.

How replication takes place?

- **Asynchronous.**

This means that if write query is processed in master db at $t=0$, then it will be replicated to slave db in some time.

Let's say it updates in slave 1 at $t=10$ and slave 2 at $t=12$.

- ∴ if any read query arrives between $t=0-t=12$ then then it will show old value.

- **Synchronous.**

This states that after the value is updated in all the slave DB, only after that the request is accepted.

Advantages

Scale out read request

Availability

Parallelism

Backup if master DB fails.