

Data → raw facts

Information → processed data with meaningful interpretation

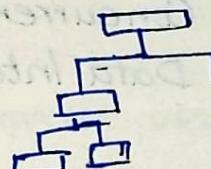
Database : Structured collection of interrelated data  
in form of tables ex: Banking, Management Systems  
Airline Reservation

file system : structured form used for management of data used by an OS

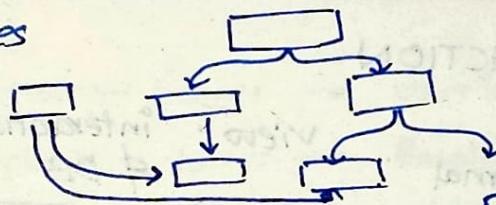
CONS

- Data Redundancy  
• Poor Memory utilization  
• Inconsistency  
• Not secure

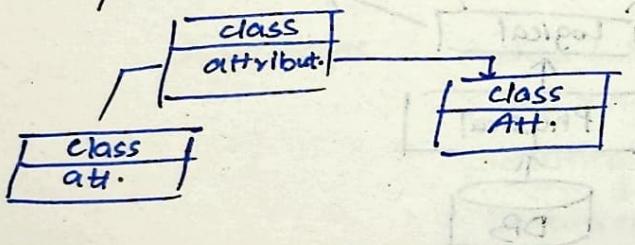
## Types of Databases | \* Hierarchical



## \* Network Databases



## \* Object Oriented



## \* Cloud Databases

# BigQuery Azure SQL

## \* Centralized

## \* Persona

- \* NoSQL  
unstructured data  
Ex: MongoDB

## WHY?

Efficient Management  
No Redundant data

Consistency

Data Security

Data Independence

Scalable

Abstraction

## WHY NOT?

Cost

Not suitable for  
smaller projects

Vendor Lockin

## APPLICATIONS

Organisation

Search and Retrieval

Access Control

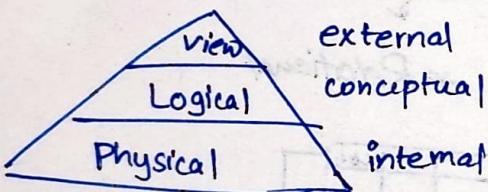
Concurrency control

Data Integrity

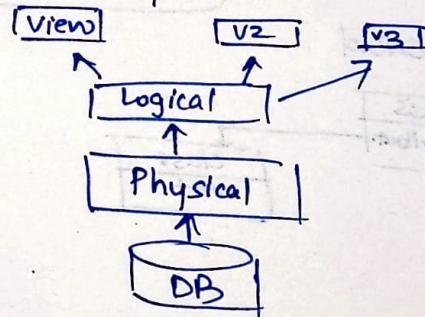
Management Systems

E-comm.

## LEVELS OF ABSTRACTION



View: interaction layer b/w users & DB



### Physical

upto the DBA

deals with actual storage

Org , allocation

### Logical

Simpler commands are used  
to perform ops

usually used by programmers

### External → user view

ex: customers

## Database Architecture

Schema : logical containers

defines structure

### Types :

Physical Schema

Data storage in underlying hardware

High risk and efforts

Logical Schema

No concern to physical storage

↳ Conceptual Schema

Overall view HL structure

↳ External Schema

External user specific view

## 3 Tier Architecture

1. Presentation Layer → user interface
2. Application Layer → Business Logic
3. Data Layer → Manage storage and processing

Advantages : Reduced Redundancy, Inconsistency

Simple Access

Security

Concurrency control

## DATA MODELS : Abstract way to represent the schema

### Types :

- \* Hierarchical
- \* Network
- \* Relational
- \* Entity Relationship
- \* OO Data Model

### \* NoSQL

document (MongoDB)

keyvalue (Redis)

Column (Cassandra)

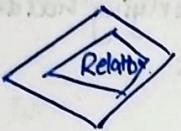
Graph (Neo4J)

## Entity

Strong: Non dependent unique

Weak: Dependent on some primary

Strong



weak



attribute  
multivalued attribute

Strong Relationship  $\Rightarrow$  No independent existence of entities  
weak  $\rightarrow$  inter-related entities.

\* Specialization: Top down

\* Generalization: Bottom up

\* Aggregation: combination

## Relational Data base Management

structured data in form of rows & columns

[Intension]: schema of a database, access, indexes etc. are defined.

Extension: data stored at any given time. (dynamic)  
includes insertion, deletion, updation

\* Data Definition Language (DDL): create, alter, drop, remove, truncate.

\* Data Manipulation Language (DML): insert, update, delete, select.

\* Data Control Language (DCL): Grant, Revoke.

\* Transaction Control Language (TCL): commit, Rollback, Savepoint

# DATABASE NORMALIZATION

↳ To Improve efficiency consistency accuracy

WHY?

Remove Anomalies →  
Improve Efficiency  
Remove Redundant data

↳ Insertion  
↳ updation  
↳ Deletion

Normal Forms

1NF  $\Rightarrow$  single valued attribute  
i.e one cell one value

WHY? To remove duplicates

2NF  $\Rightarrow$  Each non key attribute is dependent on key

WHY? To ensure consistency

3NF  $\Rightarrow$  All non key attributes are independent of each other

Boyce Codd NF  $\Rightarrow$

Be in 3NF

X should be a superkey for every functional dependency  $X \rightarrow Y$

4NF  $\Rightarrow$  No multivalued dependency

5NF  $\Rightarrow$  No join dependency, No further lossless decomposition

Functional Dependency : one set determines the other

$\rightarrow X \rightarrow Y \leftrightarrow$  Dependent  
Determinant

- Full functional Dependency : Removing any of the two attributes of  $X$  the dependency fail.

eg:  $(IFSC + Acc. No) \rightarrow Account\ details$  Removing either will result in failure

- Partial Dependency : depends on a part of the composite key,

eg:  $(Amount + AccNo.) + Name \rightarrow Deposit$  Removing Name bares no issue.

- Trivial functional Dependency : dependent is part of the determinant

$\{RollNo, Name\} \rightarrow \{Name\}$

- Transitive Dependency : Indirect dependency

Withdrawal  $\rightarrow$  BankBalance  $\rightarrow$  Bank  
Withdrawal  $\rightarrow$  Bank

- Multivalued Dependency  $A \rightarrow \{B, C\}$

## \* DENORMALISATION

combining normalized tables

### Reasons:

1. Performance
2. Simplification of queries
3. Reducing Joins
4. Enhancing Retrieval of data

### CONS:

- Problems:  
High Maintenance  
Inc Redundancy  
 $\uparrow$  Inconsistency

Techniques :

1. Add Redundant data
2. Add precomputed tables.
3. Store derived values.

## RELATIONAL DATA MODEL

Structured storage of data in form of tables

Tables = Relations

Relational Schema  $\Rightarrow$  Table structure — Relation, attributes  
constraints, keys

Relational Instance  $\Rightarrow$  set of data at  
a particular time in a relationship

Attribute  $\Rightarrow$  Columns of a table / properties of a relation

Domain of an attribute : Possible set of values for an attribute

Tuple  $\Rightarrow$  Each row for a relation

Cardinality  $\Rightarrow$  No. of distinct values in a column

Degree  $\Rightarrow$  No. of attributes also called Arity

Relational Algebra | Procedural language

1. Union : Avoids duplicates takes common values as well as non common
2. Intersection : Takes only common values
3. Difference :  $(A - B) \rightarrow$  values in A that are not in B
4. Cartesian Product  $(A * B) \rightarrow m * n$  relation
5. Selection ( $\sigma$ )

## KEYS

ensure organization, accuracy and accessibility

1. Super Key : set of one or more attributes that can uniquely identify a tuple.
2. Candidate Key : set of attributes that can uniquely identify a relation record  
It can be composite or simple
3. Primary Key : Candidate key chosen as primary
4. Alternate Key : keys that can be primary but were not chosen
5. Foreign Key : An attribute that points to the primary key of another table
6. Composite Key : Two or more attributes grouped together
7. Unique Key : Contains unique value with one NULL value

<u>Consistent</u>	<u>ACID</u>	- A transaction is a single logical unit of work
Transaction manager	<u>Atomicity</u>	: Entire transaction either takes place or doesn't. It's OR
Application programmer	<u>Consistency</u>	: Valid status before and after any transaction
Concurrency controller	<u>Isolation</u>	: Multiple transactions independent of each other concurrently
Recovery Manager	<u>Durability</u>	: Persistence of changes even after system failure

<u>High Availability</u>	<u>BASE</u>	$\Rightarrow$ flexibility and fluidity	(NOSQL)
			Ex: Cassandra, MongoDB
	<u>Basically Available</u>	: Availability of data over consistency	
	<u>Soft State</u>	: Change in data overtime	
	<u>Eventually Consistent</u>	: will be consistent eventually it takes time to converge	

CRUD  $\rightarrow$  Create Read Update Delete

Truncate → Removes all data from the table

Delete → Can delete specific rows

Drop → delete tables & database.

## OPERATORS

1. Arithmetic (+, -, \*, /)
2. Comparison (=, <, >, ≠, ≤)
3. Logical (AND, OR, NOT)

## AGGREGATE

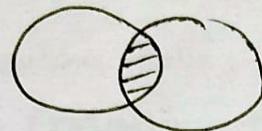
1. count
2. sum
3. Avg
4. Max
5. Min

## CLAUSES

1. WHERE → used to filter
2. GROUP BY → often used with aggregate functions
3. HAVING → filter the rows created by Group By  
GROUP BY → HAVING
4. ORDER BY → sort DESC, ASC
5. LIMIT : limits the results

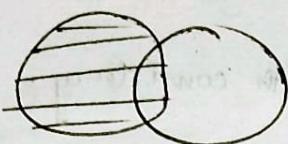
## JOINS

1. INNER :



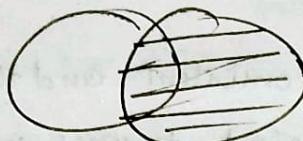
Common value

2. LEFT JOIN:



All records from left table  
+ matched records in from right

3. RIGHT



All records from right & matched records from left

4. CROSS



All records from both tables.

5. CROSS JOIN : Cartesian product of both tables.

6. SELF JOIN: Join a table with itself.

## UNIONS

combine result of two or more SQL queries. (UNIQUE result)  
including ALL gets duplicate data too

VIEWS virtual table to display any retrieval.

they can be updated

Inline views : subqueries inside the FROM clause.

Materialized View: saves the view physically

### PROS:

- \* Abstraction
- \* Data security
- \* Limiting access

### CONS:

- performance overhead.

## Scheduling : organizing and executing transaction's

\* Serial Schedules: No interleaving transactions

\* Non serial schedules: has concurrency with consistency

### a) Serializable schedules:

checks if a database is consistent and the interleaved execution is equivalent to serial execution.

#### ① Conflict serializable

if it can be transformed into a serial schedule by swapping non-conflicting operations.

~~All~~ operations <sup>are</sup> conflicting if

- ← Belong to diff. transactions
- + Operate on the same item
- + At least one operation is write

#### ② View serializable: view is equal to a serial schedule (No overlapping transactions)

### b) Non serializable

#### I) Recoverable schedule:

Transactions commit only after all transactions they read commit.

e.g.: if TA reads from TB the only after TB commits

TA commits

Recoverable  $\rightarrow$  i) Cascading Schedule

also called Avoids cascading Aborts (ACA)

if one transaction fails, all the depending are either rolled back or aborted. (cascading abort/rollback)

ii) Cascadeless Schedule : Transactions only read after all transactions whose changes are to be read are committed.

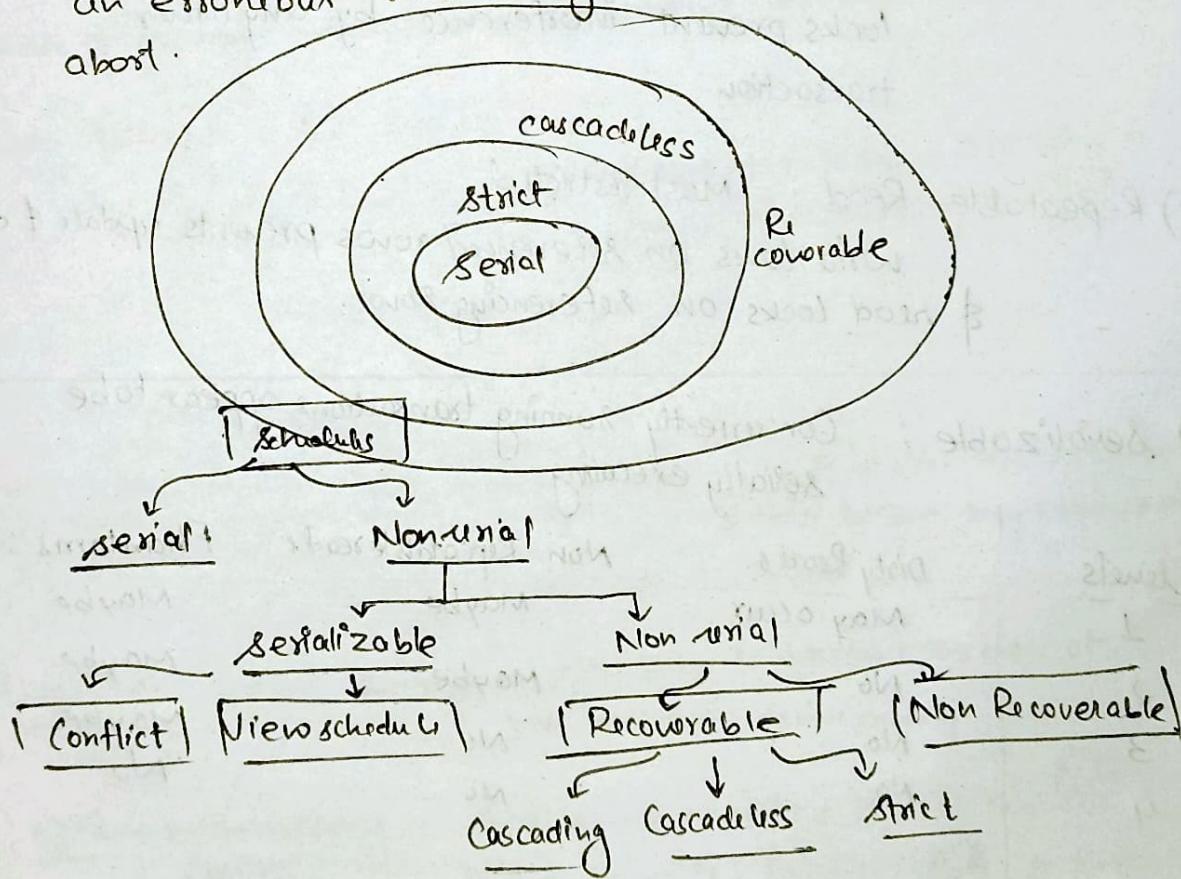
why?  $\Rightarrow$  Avoids the cascading aborts and rollbacks

iii) Strict Schedule:

A can read or write updated or written value of B only after B commits / aborts

## II) Non Recoverable Schedule

value is read before commit / abort leading to an erroneous value being read in case of future abort.



# TRANSACTION ISOLATION

Dirty Read : A transaction reads data that has not been committed

Non Repeatable read : A reads B twice and gets diff values both times, as C has updated B

Phantom Read → two reads of the same query return diff results, by two diff transactions.

## 4 Levels of Isolation :

i) Read Uncommitted : Reading uncommitted data, i.e. a dirty read is allowed.

No isolation at this level.

ii) Read Committed : Any read data is committed. (transactions hold locks)  
locks prevent interference by any other transaction.

iii) Repeatable Read : Most restricted.  
Write locks on referenced rows prevents update & delete  
& read locks on referencing rows.

iv) Serializable : Concurrently running transactions appear to be serially executing.

levels	Dirty Reads	Non Repeatable reads	Phantoms
1	May occur	Maybe	Maybe
2	No	Maybe	Maybe
3	No	No	Maybe
4	No	No	No

## Additional Levels & Features

**Snapshot Isolation** : Allows transactions to work with a snapshot of the database at any specific time.

**Multi-Version Concurrency Control** : Maintains multiple versions of data for better concurrency.

## SERIALIZABILITY AND CONCURRENCY CONTROL

↳ ensures the results are same as if the queries were performed serially instead of concurrent execution.

types: i) Conflict : can be converted into serial schedule  
ii) View : preserves the same view as serial

CONS: Performance Overhead , Complexity

Methods to manage simultaneous operations without any conflicts

| Primary focus: ! to avoid :

i) Lost update : if two transactions are trying to update simultaneously a single value.

ii) Dirty Read

iii) Non Repeatable Read .

Techniques : i) Lock Based Protocol: A transaction must acquire a lock before any ~~any~~ operation. Locks can be shared or exclusive.

ii) Timestamp-Based : Assigns timestamp to determine the order of execution. This maintains serializability.

iii) Optimistic Concurrency Control : Assumes no conflicts and checks only before commits. Rollback if conflicts are found.

4. Multi-version Concurrency Control: Versions of data are maintained.

(MVCC)

PROS

- Data integrity
- Inc. throughput
- Improved User exp.

### APPLICATIONS

- Banking
- Ecom
- Social Media
- Reservations

## LOCKING PROTOCOLS

## LOCKING PROTOCOLS

To control the access of data to ensure consistency and accuracy.

1. Shared Locks
2. Exclusive lock

Type of Locks!

Lock is a simple variable associated with the data accessibility at any given time.

SHARED LOCK : also called the read' lock

(S-lock)      only  
Multiple transactions can read at the single instance.  
ensures no write operations occur on the given data.

EXCLUSIVE LOCK: (X-lock)

only a single transaction can read & write the data  
thus single transaction holds the data

### Lock Based PROTOCOLS

i) Simplistic Lock Protocol : acquire lock before read/ write access  
uses single lock type.

PROS: data consistency is maintained through no concurrent access  
simple

CONS: Performance bottlenecks due to limited concurrency  
No deadlock prevention

ii) Pre-Claiming Lock Protocol : Need to declare all the locks and obtain them before access.

PROS: Avoids cyclic wait

CONS: Reduced concurrency  
Inefficient due to unnecessary holding of locks.

### iii) Two Phase Locking (2PL)

Transaction → Growing Phase : cannot release locks

→ Shrinking Phase : cannot acquire locks

Guarantees serializable transactions.

PROS: Consistency

Ensure no interference

CONS: Deadlocks

Performance overhead

### iv) Strict Two Phase locking (S2PL)

Must hold all exclusive locks until commit/abort.

Shared locks can be released

PROS: Avoids cascading aborts

Simplifies recovery as uncommitted changes are not visible to other processes.

CONS: Deadlock

Reduced concurrency.

## Database Recovery Management

Types :

Transaction Failure	↳ logical	} software
System failure	↳ OS	

Media failure	↳ Physical	} hardware
Natural Disaster		

### Techniques:

i) Backup and Restore : restoring from backup's

ii) Log-Based Backup : Redo completed transactions

& undo failed transactions

3. Shadow Paging: a shadow copy is maintained of the database page.  
In case of update the other stays unchanged acting as a point of restoration in case of failure.
4. Checkpointing: snapshots of the database are created at particular intervals.

### 5. Rollback & Rollforward

↓  
undoes incomplete transactions

→ Applies changes by completed transaction

### Best Practices

- 1. Regular Backups
- Offsite storage: 3:2:1 rule
- Test recovery processes.
- Automate processes.

## QUERY OPTIMISATION

1. Indexing: pointers to speed up retrieval of rows  
Have a overhead thus periodic check for unused index is essential.
2. Avoiding select \*
3. using EXISTS instead of IN or COUNT  
↳ searches for first occurrence only
4. WHERE instead HAVING  
↓  
filters first  
↳ groups first
5. USE JOINS over Subqueries
6. Using LIMIT or TOP
7. Avoid wildcards at begin like patterns.

g. using proper datatypes

10. Avoid functions over indexed columns

\*\*

## INDEXING

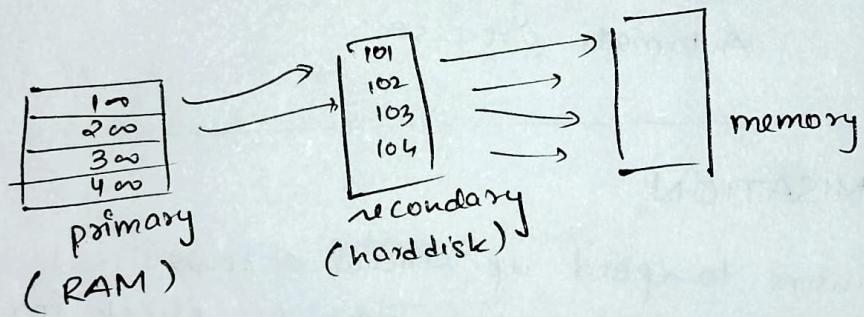
PROS : faster queries

CONS: storage & computation overheads.  
throughput hit

i) Ordered Indices  
increases with nested index

ii) Clustered Index : grouping

iii) Secondary Index : Additional indexing to reduce size of mapping in sparse indexing



iv) Primary Index: Made on the primary key of the table

Sparse Index: records for only some of the values  
ex: for 1000 records at every 100

Dense Index: index for every search key value.

## B - Trees

All leaf nodes must be at the same level

Balanced data structure

Each node has  $m$  children where  $m$  is the tree order  
( $m-1$ ) keys

Non root and Non leaf nodes have atleast  $\lceil m/2 \rceil$  children

Nodes are only added at leaf nodes.

## B + Trees

Internal nodes act as guides while data is stored at leaf nodes.

Leaf nodes are linked | self balanced

PROS

- Fast access and efficient
- Support for large datasets
- Self balancing
- Dynamic size handling

### APPLICATIONS

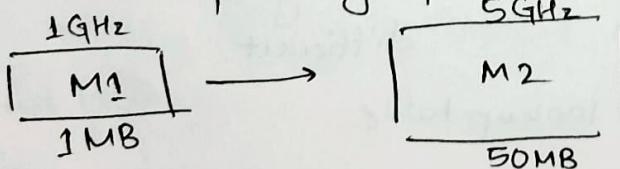
- File systems
- DB
- Search engines
- Telecom Networks
- Geospatial DB

## Vertical and Horizontal Scaling

changing system size in regards to needs.

VERTICAL : Adding additional resources to the system { scale up/down }

i.e powering up the current machine



- \* Single Node
- \* Less time & expertise
- \* Easier to maintain

### PROS

CONS : \* Higher risk of data loss in case of failure

- \* Downtime
- \* Limited upgradeability

## HORIZONTAL SCALING

Changing the number of nodes (scaling in/out)

Distributed workload

increased complexity  
Expensive, more effort  
to maintain

Low downtime

Load balancing

SHARDING : Dividing a large database into a smaller chunk that more manageable.  
The shards can be hosted separately

Components : i) Shards : logical partitions of data  
ii) Shard key : column that determines the distribution

iii) Shard Nothing Architecture :  
Independent op. of shards & unaware of others.

## Methods

i) Range Based : specified Range

\* Easy to implement

\* Can lead to uneven data distribution

ii) Hashed Sharding : using a hash function & key

\* Even distribution

\* Reassignment is difficult.

iii) Directory Sharding : uses a lookup table

iv) Geo Sharding : partition based on geo loc.

## Optimizing sharding

- factors:
- Cardinality: unique value ↑ ↑ shards
  - Frequency: Avoid high freq. to reduce hotspots
  - Monotonic change: Rate of change of keys

## ROLE BASED ACCESS CONTROL

- components:
- 1) Role: could be a job function within an org.
  - 2) Permissions: access rights & privileges
  - 3) Users: individuals

User Assignment, Role Permissions, Role Hierarchies

PROS: simplified Administration  
security  
scalability  
compliance

CONS:

Complexity  
Inflexible  
Delegation  
Granularity

## DATA MASKING : obfuscating data (Hiding)

### i) Static Data Masking

changing words when data is at rest.

usually before transfers commonly during development

and testing

### ii) Dynamic Data Masking

Both for full / partial masking

could be done during data transfers.

### iii) Deterministic Data Masking

value is replaced by similar value in the very same row

### iv) On the fly:

Similar to dynamic data M. but done one value at a time.

## TECHNIQUES

i) Substitution : replace with a fake yet realistic value

↑ Makes realistic data

↓ Not applicable to large data that is unrelated

ii) Averaging : can be used with numeric data

iii) Shuffling

iv) Encryption

v) Null out / Deletion

vi) Redaction : replacing with generic values.

vii) Date Aging

## DATA ENCRYPTION

plaintext → ciphertext

: same key on both sides (private encryption)  
eg: AES

i) Symmetric key

, diff keys public & private

ii) Asymmetric key  
they are mathematically linked

eg: RSA , DSA