

Parallel & Distributed Database

Database System Architectures

- Database systems can be centralized, where one server machine executes operations on the database.
- Database systems can also be designed to exploit parallel computer architectures.
- Distributed databases span multiple geographically separated machines.

Centralized Database system

- Centralized database systems are those that run on single computer system.
- Such database systems span a range from single user database systems running on mobile devices or personal computers to high performance database systems running on a server with multiple CPU cores & disks & large amount of main memory that can be accessed by any of the CPU cores.
- It widely used for enterprise - scale applications.
- We distinguish two ways in which computers are used as
 - single user systems ✓
 - Multiuser systems ✓

single user systems

- Smartphones & personal computers are example of single user systems.

- A typical single user system is a system used by a single person, usually with only one processor (with multiple cores) & one or two disks.

Multiuser systems

- A typical multiuser system on the other hand, has multiple disks, a large amount of memory & multiple processors.
- Such systems serve a large number of users who are connected to the system remotely & they are called server systems.

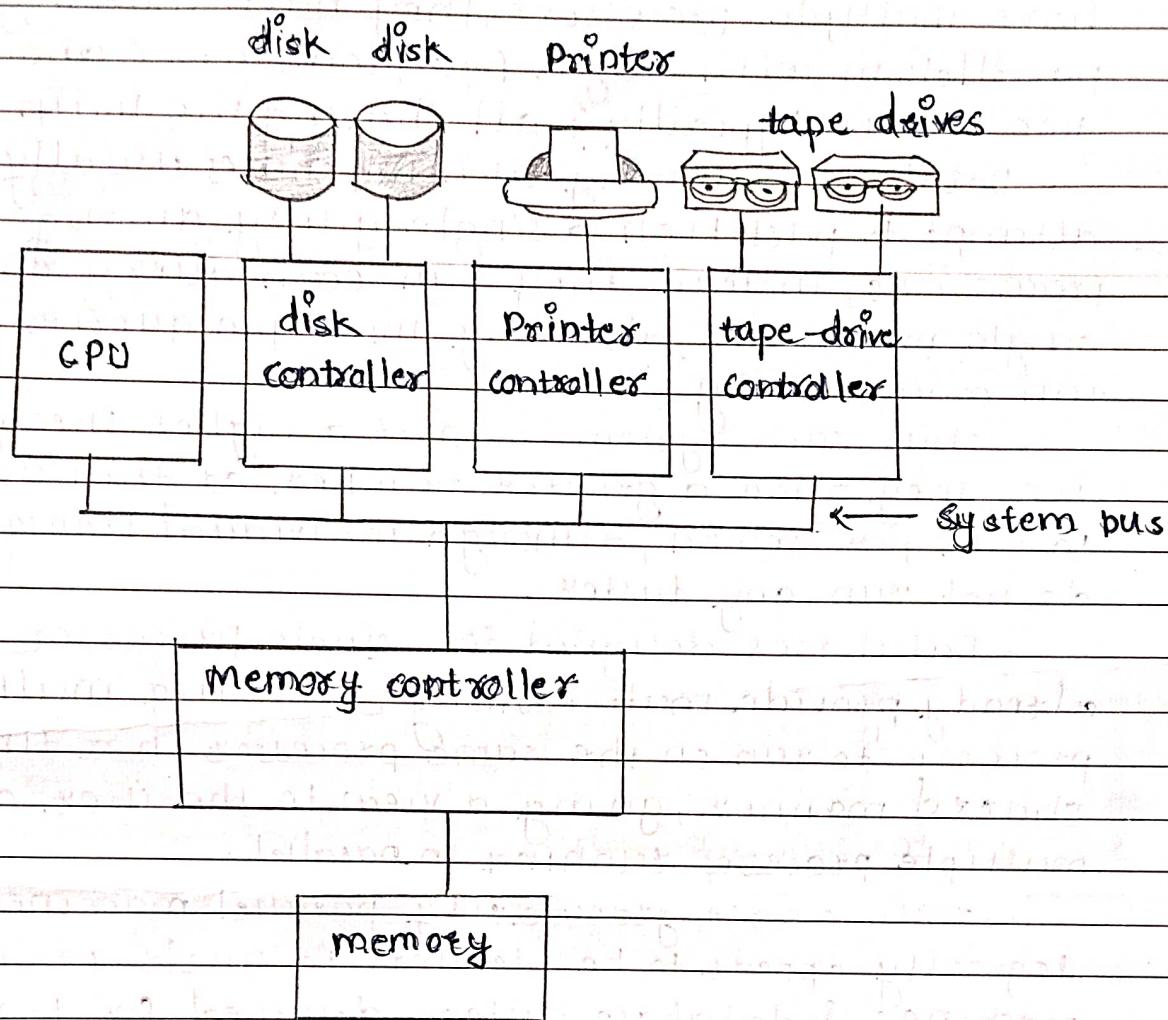
- Database system designed for single user sys usually do not provide many of the facilities multiuser database provides.
- In particular, they may support very simple concurrency control schemes, since highly concurrent access to the database is very unlikely.
- Provisions for crash recovery In such system may also be either very basic or even absent in some cases, such systems may not support SQL & may instead provide an API for data access, such database systems are referred to as embedded databases, they are designed to be linked to a single application program & they accessible only from that application.
- multiuser database systems support full transactional features that are as databases are usually designed as servers, which service requests received from application programs, the request could be in the form of SQL queries or they could be requests for retrieving, storing or updating data specified using API

- centralized database systems are those that run on a single computer system & do not interact with other computer system, such database system running on personal computers to high performance database systems running on high end server systems.

- We distinguish two ways in which computers are used as single user systems & multiuser system. Personal computers & workstations fall into the first category. A typical single user system is a desktop unit used by single person, usually with only one person using the machine at a time.

- A typical multiuser system on other hand has a more disks & more memory, may have multiple CPUs & has a multiuser operating system. It serves a large number of users who are connected to the

system via terminals



centralized computer system

- Database systems designed for use by single users usually do not provide many of the facilities that a multiuser database provides.
- In particular, they may not support concurrency control, which is not required when only single user can generate update.
- Provisions for crash recovery in such systems are either absent or primitive. For example, they may consist of simply making a backup of the database before any update.
- Many such systems do not support SQL & provide a simple query language, such as variant of QBE.
- In contrast, database systems designed for multiple user system support the full transactional features.

- Although general purpose computer systems have multiple processors, they have coarse-grain parallelism, with only a few processors (about two to four typically) all sharing the main memory.

- Databases running such machines usually do not attempt to partition a single query among the processors, instead, they run each query on a single processor, allowing multiple queries to run concurrently.

- Thus, such systems support a higher throughput, i.e., they allow a greater number of transactions to run per second, although individual transactions do not run any faster.

- Databases designed for single processor machines already provide multitasking, allowing multiple processes to run on the same processor in a time shared manner, giving a view to the user of multiple processes running in parallel.

- Thus coarse granularity parallel machines logically appear to be identical to single-processor machines & database systems designed for time shared machines can be easily adapted to run on them.

- In contrast, machines with fine granularity parallelism have a large no. of processors & DB systems running on such machines attempt to parallelize single tasks submitted by users.

Client server system

- As personal computers became faster, more powerful & cheaper there was a shift away from the centralized system architecture.

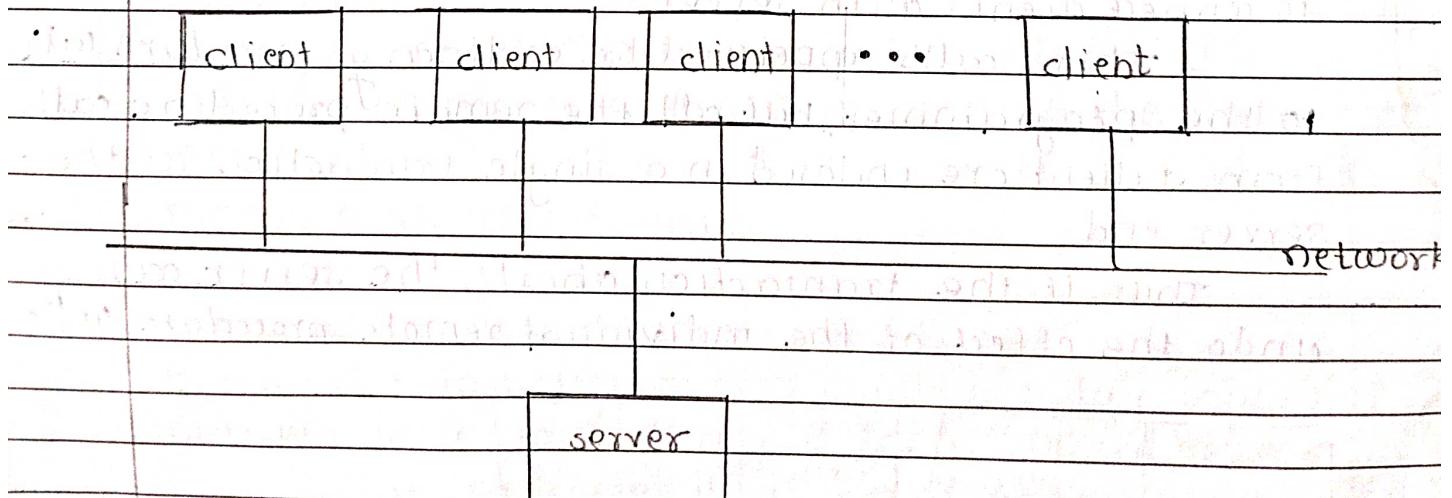
- Personal computers supplanted terminals connected to centralized systems. Correspondingly, personal computers assumed the user interface functionality that used to be handled directly by the centralized systems. As a result, centralized systems today act as server systems that satisfy requests generated by client system.

- Database functionality can be broadly divided into two parts: front end & back end as in figure.

- The back end manages access structures, query evaluation & optimization, concurrency control & recovery.

- The front end of database system consists of tools such as forms, reports writers & graphical user interface facilities.

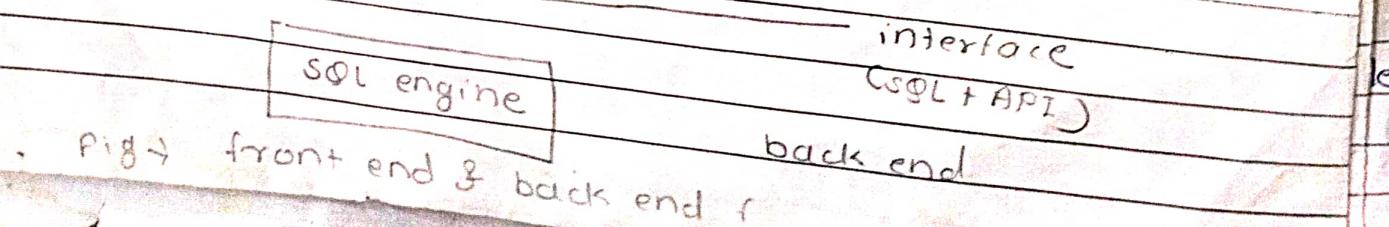
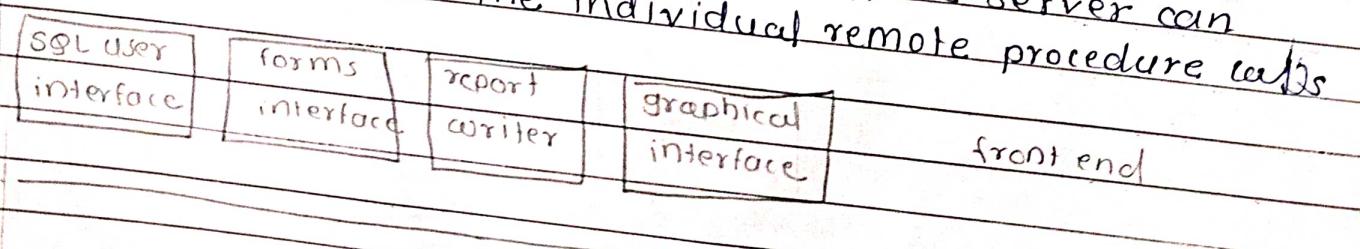
- The interface between the front end & back end is through SQL or through an application program.



client - server system

- standards such as ODBC & JDBC, were developed to interface clients with servers.
- Any client that uses the ODBC or JDBC interfaces can connect to any server that provides the interface
- In earlier generation database systems, the lack of such standards necessitated that the front end & back end be provided by the same software vendor. With the growth of interface standards, the front end uses interface & the back end servers are often provided by different vendors.
 - Application development tools are used to construct user interfaces. They provide graphical tools that can be used to construct interfaces without any programming.
 - Some of the popular application development tools are PowerBuilder, Magic & Borland Delphi. Visual Basic is also widely used for application development.
 - Further certain application programs, such as spreadsheets & statistical-analysis packages use the client-server interface directly to access data from a backend server.
 - In effect they provide front ends specialized for particular tasks. Some transaction processing systems provide a transactional remote procedure call interface to connect clients with server.
 - These calls appear like ordinary procedure calls to the programmer, but all the remote procedure calls from a client are enclosed in a single transaction at the server end.

Thus, if the transaction aborts, the server can undo the effect of the individual remote procedure calls



* Server system Architectures

Server systems can be broadly categorized as transaction servers & data servers.

i) Transaction server system also called query server system provide an interface to which clients can send requests to perform an action in response to which key execute the action & send back results to the client.

- Usually, client machine ship transactions to the server systems where those transactions are executed & results are shipped back to clients that are in charge of displaying the data. Requests may be specified by using SQL or through a specialized application program interface.

ii) Data server systems allow client to interact with the servers by making requests to read or update data, in units such as files or pages.

e.g., file servers provide a file system interface where clients can create, update, read & delete files.

- Data servers for database systems offer much more functionality, they support units of data such as pages, tuples or objects that are smaller than a file.

- They provide indexing facilities for data, & provide transaction facilities so that the data are never left in an inconsistent state if a client machine or process fails.

i) Transaction server process structure

A typical transaction server system today consist of multiple processes accessing data in shared memory.

The processes that form part of the database system include.

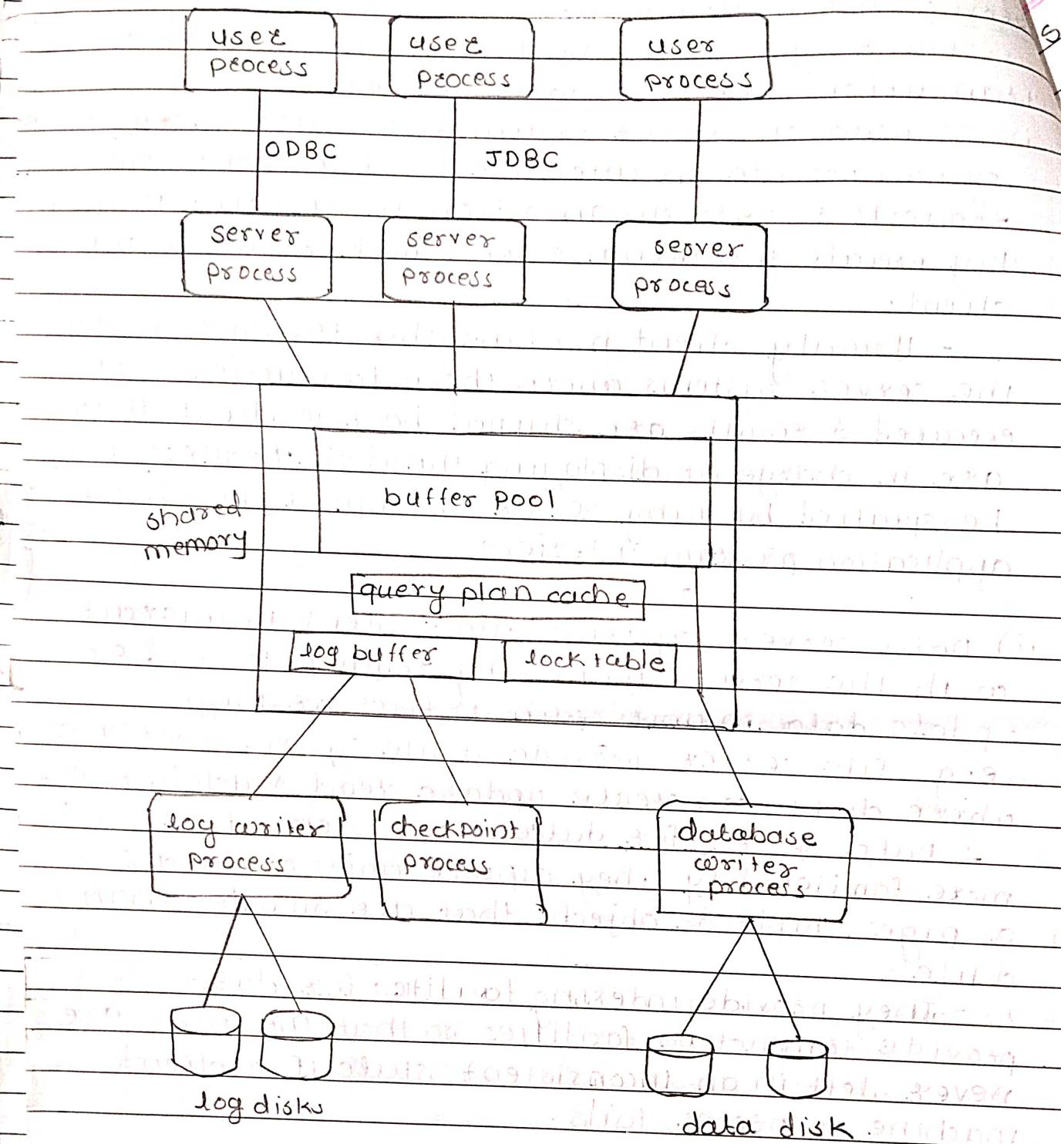


Figure → shared memory & process structure

The processes that form part of the database system include.

i) server processes

These are processes that receive user queries (transactions), execute them & send the results back.

The queries may be submitted to the server processes from a user interface or from a user process running embedded SQL or via JDBC, ODBC or similar protocols.

- Some database systems use a separate process for each user session & few use a single database process for all user sessions but with multiple threads so that multiple queries can execute concurrently.

(A thread is like a process, but multiple threads execute as part of the same process, run in the same virtual memory space. Multiple threads within process can execute concurrently.)

- Many database systems use a hybrid architecture with multiple processes, each one running multiple threads.

ii) lock manager process

This process implements lock manager functionality which includes lock grant, lock release & deadlock detection.

iii) database writer process

There are one or more processes that output modified buffer blocks back to disk on a continuous basis.

iv) log writer process

This process outputs log records from the log record buffer to stable storage. Server processes simply add log records to the log record buffer in shared memory & if a log force is required, they request the log writer process to output log records.

v) checkpoint process

This process performs periodic checkpoints.

vi) Process monitor process

This process monitors other processes & if any of them fails, it takes recovery actions for the process such as aborting any transaction being executed by the failed process & then restarting the process.

The shared memory contains all shared data such as

- Buffer pool

- Lock table

- Log buffer, containing log records waiting to be output to the log on stable storage.

- Cached query plans, which can be reused if the same query is submitted again all database processes can access the data in shared memory. Since multiple process may read or perform updates on data structures in shared memory, they must be a mechanism to ensure that only one of them is modifying any data structure at a time & no process is reading a data structure while it is being written by others.

Data server

- Data server system are used in local area network where there is a high speed connection between the clients & the server, the client machines are comparable in processing power to the server machines & the tasks to be executed are computation intensive.

- In such an env., it makes sense to ship data to client machines, to perform all processing at the client machine & then to ship the data back to the server machine & the tasks to be executed are computation intensive. Note that this architecture requires full back end functionality at the clients. Data server

architectures have been particularly popular in object oriented database system.

- Interesting issues arise in such an architecture since the time cost of communication between client & server is high compared to that of local memory reference.

- Page shipping versus item shipping.

- The unit of communication for data can be of coarse granularity such as a page or fine granularity such as a tuple.

- We use the term item to refer to both tuples & objects if the unit of communication is a single item, the overhead of message passing is high compared to the amount of data transmitted. Instead when an item is requested, it makes sense also to send back other items that are likely to be used in the near future.

- Fetching items even before they are requested is called prefetching. Page shipping can be considered a form of prefetching if multiple items reside on a page, since all the items in the page are shipped when a process desires to access a single item in the page.

Locking

- Locks are usually granted by the server for the data that it ships to the client machines.
- Disadvantage of page shipping is that client machines may grant locks of too coarse a granularity a lock on page implicitly locks all items contained in the page.
- Even if the client is not accessing some items in the page, it has implicitly acquired locks on all prefetched items.
- Other client machines that require locks on those items may be blocked unnecessarily.
- Techniques for lock deceleration have been proposed where the server can request its clients to transfer back locks on prefetched items. If the client machine does not need prefetched items, it can transfer locks on the item blocks to the server & locks can then be allocated to other clients.

Data caching

- Data that are shipped to a client on behalf of a transaction can be cached at the client, even after the transaction completes, if sufficient storage space is available.
- successive transaction completes, if sufficient storage space is available at the same client may be able to make use of the cached data.
- However, cache coherency is an issue: even if a transaction finds cached data, it must make sure that those data are up to date, since they may have been updated by a different client after they were cached.
- Thus a message must still be exchanged with the server to check validity of the data & to acquire a lock on the data.

Lock caching

If the use of data is mostly partitioned among the clients with clients rarely requesting data that are also requested by other clients, locks can also be cached at the client machine.

parallel system

- * Parallel database systems consist of multiple processing & multiple disks connected by a fast interconnection network.
- A coarse grain powerful parallel machine consists of a small no. of powerful processors
- A massively parallel or fine grain parallel machine utilizes thousands of smaller processors.
- In parallel processing many operations are performed simultaneously.
- Parallel computers with hundreds of processors & disks are available commercially.
- There are 2 main measures of performance of DB system
 - 1) Throughput \rightarrow No. of tasks that can be completed in given time interval.
 - 2) Response time \rightarrow The amount of time it takes to complete single task from the time it is submitted.

speedup scaleup

- A fixed sized problem executing on small system is given to a system which is N -times larger measured by:

$$\text{speedup} = \frac{\text{small system elapsed time}}{\text{large system elapsed time}}$$

- speedup is linear if equation equals N .

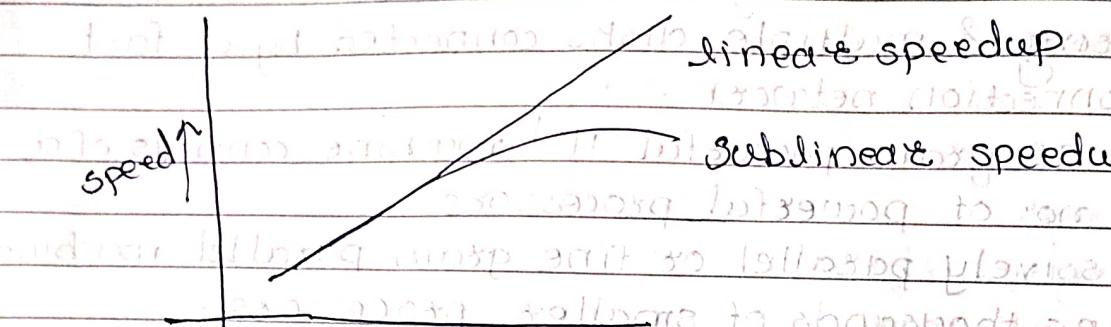
~~Increase the size of both the problem & the system~~

~~It is N times larger system used to perform N -times larger job~~

~~measured by~~

$$\text{scale up} = \frac{\text{small system small problem elapsed time}}{\text{big system big problem elapsed time}}$$

Scaleup is linear if equation equals 2



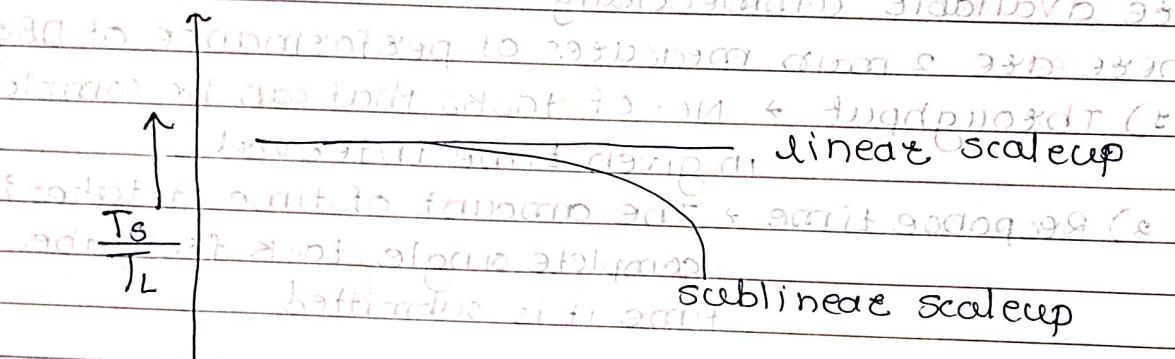
linear & speedup

sublinear & speedup

speed ↑

resources →

② fig. speedup



problem size →

Batch scaleup -

- A single large job, typical of most decision support queries & scientific simulation
- Use an N-times larger computer on N-times larger problem.

Transaction scaleup

Numerous small queries submitted by independent users to shared database, typical transaction processing & time sharing systems.

- N times as many users submitting requests to an N-times larger db on an N-times larger computer.

- Suited to rel DB.

Factors limiting speedup & scaleup

i) startup costs

cost of starting up multiple processes may dominate computation time, if the degree of parallelism is high.

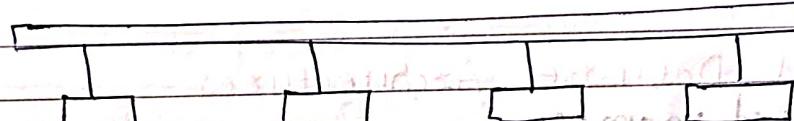
ii) Interference

Processes accessing shared resources compete with each other, thus spending time waiting on other processes, rather than performing useful work.

Interconnection Network Architecture

1) Bus

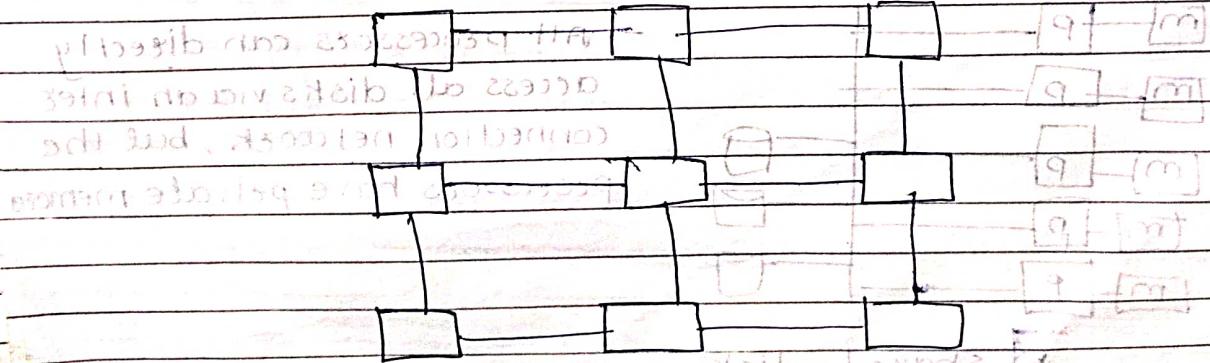
- i) System components send data on & receive data from a single communication bus.
- ii) It does not scale well with increasing parallelism.



a) bus

2) mesh

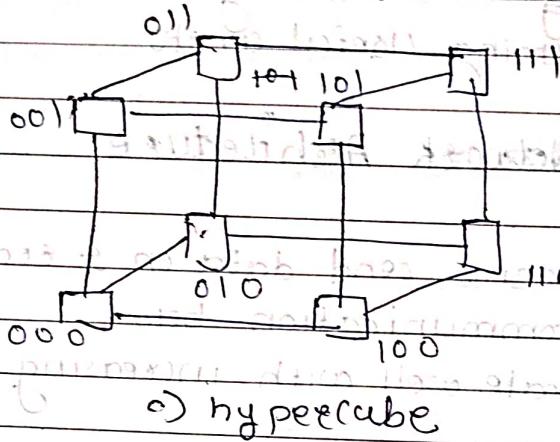
- i) components are arranged as nodes in a grid & each component is connected to all adjacent components.
- ii) communication link grow with growing no. of components & so scales better.
- iii) But may require 2n hops to send message to a node.



b) mesh

Hypercube

- i) components are numbered in binary, components connected to one another if their binary representations differ in exactly one bit.
- ii) n components are connected to log(n) other components & can reach each other via at most log(n) links, reduces communication delays.

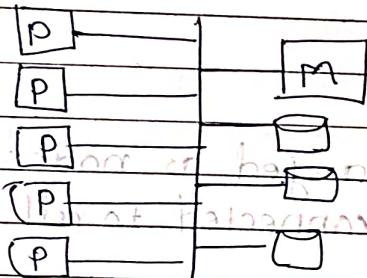


c) hypercube

Parallel Database Architectures

a) shared memory

- Processors share a common memory

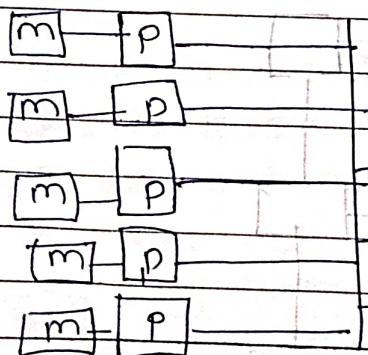


- Process & disks have access to a common memory typically via a bus or through an interconnection network

a) shared memory,

b) shared disk

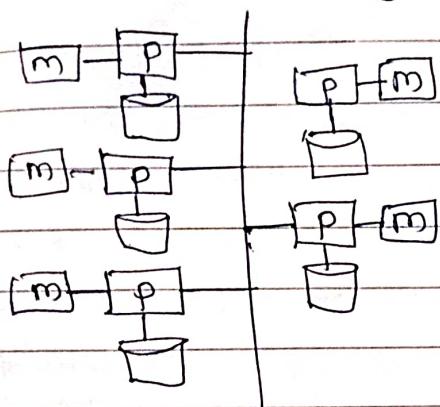
Processors share a common disk



- All processors can directly access all disks via an interconnection network, but the processors have private memory

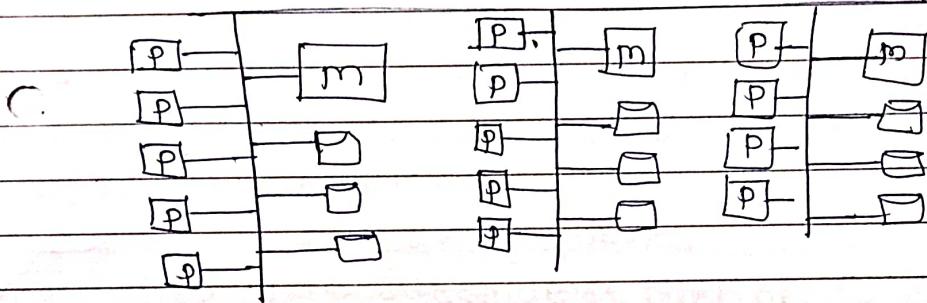
b) shared disk

3) shared nothing



- Processors share neither a common memory nor common disk

4) Hierarchical



- Hybrid of above

- It combines characteristics of shared memory, shared disk & shared nothing arch.

Distributed DB System

- i) It is logically related collection of data that is shared which is physically distributed over computer network on different sites.
- ii) It consists of loosely coupled sites that share no physical components.
- iii) DB systems that run on each site are independent of each other.
- iv) Transactions may access data at one/more sites.

Homogeneous Distributed DB

In an homogeneous distributed DB

- i) All sites have identical software.
- ii) Are aware of each other & agree to cooperate in processing user requests.
- iii) Each site surrenders part of its autonomy in terms of right to change schemas or rules.
- iv) Appears to user as single system.

Heterogeneous distributed DB

- i) Different sites may use different schemas & rules.
- ii) Difference in schema is major problem for query processing.
- iii) Difference in rule is a major problem for transaction processing.
- iv) sites may not be aware of each other & may provide only limited facilities for cooperation in transaction processing.

* Distributed data storage

i) Replication

- i) System maintains multiple copies of data stored in different sites for faster retrieval & tolerance.
- ii) A relation or fragment of relation is replicated it is stored redundantly in two or more sites.
- iii) Fully replication of R is the case where relation is stored at all sites.

* Data fragmentation

Division of relation R into fragments R₁, R₂ ... R_n which contain sufficient info. to reconstruct R.

i) Horizontal fragmentation

Each tuple of R is assigned to one / more fragments.

E.g.

branch-name	acct-no.	bcl
Hillside	A-305	500
Hillside	A-226	800
Hillside	A-155	62

account = 6branch-name = "Hillside" (account)

branch-name	acct-no.	bcl
valleyview	A-177	205
—ii—	A-402	10000
—ii—	A-408	1123
—ii—	A-639	750

account₂ = 6branch-name = "valleyview" (account)

ii) vertical fragmentation

The schema for relation R is split into several smaller schemas.

- a) All schemas must contain a common candidate key or super key to ensure lossless join property
- b) A special attribute, the tuple-id attribute may be added to each schema to serve as a candidate key.

account-no bal tuple-id

A - 805

500

fragment 1

A - 226

836

2

A - 177

205

3

A - 402

10000

4

A - 105

62

5

A - 804

1123

6

A - 302

750

7

$\text{deposit}_2 = \Pi_{\text{account-no}, \text{balance}, \text{tuple-id}} (\text{employment}$

Advantages of fragmentation

- i) Horizontal fragmentation
- ii) It allows parallel processing on fragments of relation
- iii) It allows a relation to be split so that tuples are located where they are most frequently accessed.
- iv) Vertical fragmentation
- v) Allows tuples to be split so that each part of the tuple is stored where it is most frequently accessed.

- ii) tuple-id attribute allows efficient joining of vertical fragments.
- iii) Allow parallel processing on a relation.
- iv) Vertical & horizontal fragmentation can be mixed.
- Fragments may be successively fragmented to an arbitrary depth.

* Data Transparency

- i) It is a degree to which system user may remain unaware of the details of how & where the data items are stored in a distributed system.

Transparency issues in relation to

- i) Fragmentation transparency
- ii) Replication transparency
- iii) Location transparency.

* Naming of data items - criteria

- 1) Every data item must have a system wide unique name.

- 2) It should be possible to find the location of data items efficiently.

- 3) It should be possible to change the location of data items transparently.

- 4) Each site should be able to create new data items of autonomously.

* centralized scheme - Name server

i) Structure

- i) name server assigns all names.
- ii) Each site maintains record of local data items.
- iii) Sites ask name server to locate non local data items.

* Use of Aliases

- Alternative to centralized scheme \rightarrow each site prefixes its own site identifier to any name that it generates i.e. site ID. account.
- Fulfils having unique identifiers & avoids problems associated with central control.
- However, fails to achieve network transparency.
- Solution - create set of aliases for data items, store the mapping of aliases to the real names at each site.
- The user can be unaware of the physical location of data item & is unaffected if the data item is moved from one site to another.

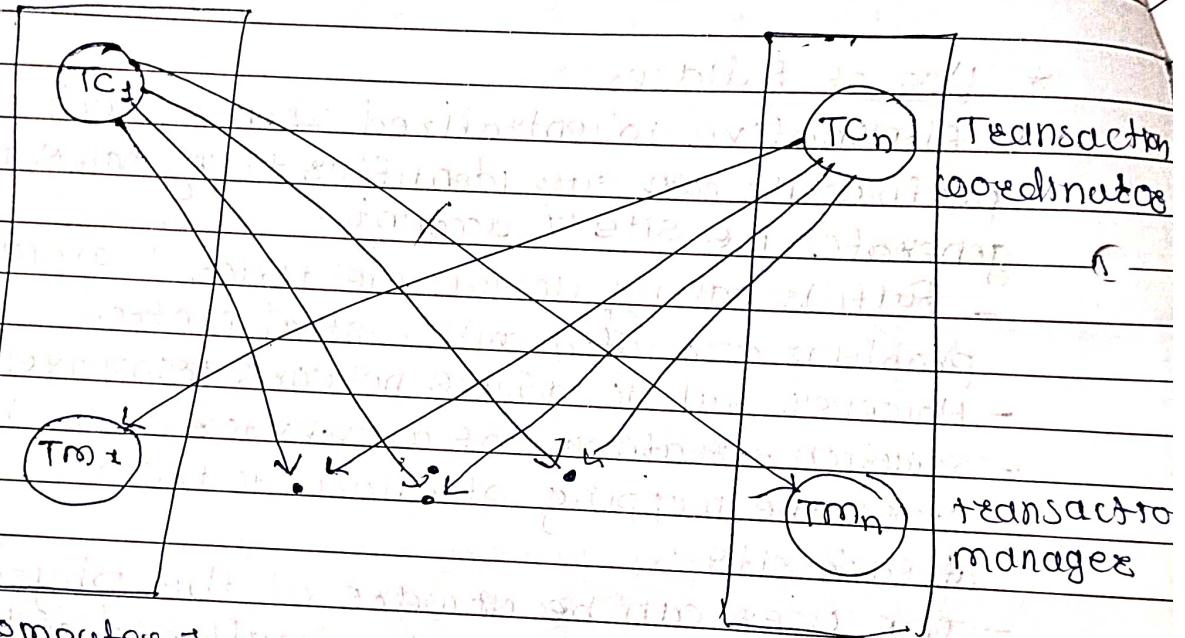
Distributed Transaction

- Transaction may access data at several sites.
- Each site has a local transaction manager responsible for
 - i) maintaining a log for recovery purposes.
 - ii) participating in coordinating the concurrent execution of the transactions executing at that site.
- Each site has transaction coordinator, which is responsible for
 - i) starting the execution of transactions that originate at the site
 - ii) distributing subtransactions at appropriate sites for execution

commit
coms
adv.

- iii) coordinating the termination of each appropriate transaction that originates at the site, which may result in the transaction being committed at all sites or aborted at all sites.

Transaction system Architecture



* System failures Modes

- failures unique to distributed systems

a) Failure of a site

b) Loss of messages

- Handled by new transmission control protocols such as TCP-IP

c) Failure of communication links

- Handled by new protocols by sending msg. via alternative links.

commit protocol

- commit protocols are used to ensure atomicity across sites
- i) A transaction which executes at multiple sites must either be committed at all the sites or aborted at all the sites
- ii) Not acceptable to have a transaction committed at one site & aborted at another
- The two phase commit (2PC) protocol is widely used
- The three phase commit (3PC) protocol is more complicated & more expensive but avoids some drawbacks of two phase commit protocol. This protocol is not used in practice.

1) Two phase commit Protocol (2PC)

- i) Assumes fail stop model - failed sites simply stop working & do not cause any other harm such as sending incorrect msg. to other sites.
 - ii) Execution of the protocol is initiated by the co-ordinator after the last step of the transaction has been reached.
 - iii) The protocol involves all the local sites at which the transaction executed.
 - iv) Let T_j be the transaction initiated at site j & let the transaction co-ordinator at site i be C_i
- phase 1 - obtaining a decision
- i) coordinator asks to all participants to prepare to commit transaction T_j
 - C_i adds the records ($\text{Prepare } T_j \rightarrow$) to the log & forces log to stable storage
 - sends $\text{prepare } T_j$ msg. to all sites at which T_j executed.

- ii) Upon receiving msg, transaction manager at determines if it can commit the transaction
 - if not, add record \langle no T \rangle to the log & send abort T msg. to Ci
 - if the transaction can be committed, then
 - add the records \langle ready T \rangle to the log
 - force all records for T to stable storage
 - send ready T msg. to Ci

Phase 2 → Recording the Decision

- T can be committed if Ci received a ready T message from all the participating sites otherwise T must be aborted.
- Coordinator adds a decision record \langle commit T \rangle or \langle abort T \rangle , to the log & forces record onto stable storage. Once the record stable storage, it is irreversible (even if failure occurs).
- Coordinator sends a msg. to each participant informing it of the decision (commit / abort)
- Participants take appropriate action locally

* Handling of failures - site failure

- When site Si recovers, it examines its log to determine the fate of transactions active at the time of the failure.

- i) Log contains \langle commit T \rangle record : txn had completed, nothing to be done.
- ii) Log contains \langle abort T \rangle record : txn had completed, nothing to be done.
- iii) Log contains \langle ready T \rangle record : site must consult Ci to determine the fate of T
 - If T committed, redo CT, write \langle commit T \rangle record
 - If T aborted, undo CT

- iv) The log contains no log records concerning T:
 - Implies that Sk failed before responding to the prepare T message from Ci.
 - Since the failure of Sk precludes the sending of such response, coordinator Ci must abort T.
 - Sk must execute undo (T).

Recovery & concurrency control.

- In doubt transactions have a \langle ready T \rangle but neither a \langle commit \rangle nor an \langle abort T \rangle log record.
- The recovering site must determine the commit abort status of such transactions by contacting other sites. This can slow & potentially block recovery.
- Recovery algo. can note lock info. in the log.
- Instead of \langle ready T \rangle write out \langle ready T, L \rangle L = list of locks held.

3PC

- i) No new partitioning.
- ii) At any point at least one site must be up.
- iii) At most K sites can fail.

Phase 1 - Obtaining preliminary Decision: Identical to 2PC Phase 1

- every site is ready to commit if instructed to do so.

Phase 2 of 2PC is split into 2 phases, Phase 2 of 2PC & Phase 3 of 3PC

- In Phase 2 coordinator makes a decision as in 2PC & records it in multiple sites.
- In Phase 3 coordinator sends commit msg. to all participants.

concurrency control

- i) Modify concurrency-control schemes for use in distributed env.
- ii) we assume that each site participates in the execution of commit protocol to ensure global transaction automatically.
- iii) we assume all replicas of any item are updated.

2) single lock manager Approach

- System maintains single lock manager. That resides in single chosen site, say S_1 .
- When transaction needs to lock data item it sends a lock request to S_1 & lock manager determines whether the lock can be granted immediately.

- i) If less lock manager sends a msg to the site which initiated the request
 - If Not request is delayed until it can be granted, at which time a msg is sent to the initiating site.

2) Di

2) Distributed lock manager

- i) In this approach functionality of locking is implemented by lock managers at each site
- ii) lock managers control access to local data items but special protocols may be used for replicas.
- iii) work is distributed & can be made robust to failures.

(iv) several variants of this approach

- a) primary copy protocol
- b) majority copy protocol
- c) Biased protocol
- d) Quorum consensus

A. Primary Copy Consensus Protocol

• Primary copy is replicated at all sites

(Primary + backup) = PAND

(Primary + backup + witness) = 3PAND

• "Biased" is branch-based replication

• "Quorum" is primary + witness + backup

• "Majority" is primary + backup + witness

• "Replica" is primary + backup + witness

Distributed Query Processing

- i) For centralized systems, the primary criterion for measuring the cost of a particular strategy is the no. of disk accesses.
 - ii) In a distributed system, other issues must be taken into account.
 - iii) The cost of data transmission over the network.
 - iv) The potential gain in performance from having several sites process parts of the query in parallel.
- consider algebraic queries on fragments,
- it must be possible to construct relation \bowtie from its fragments
 - Replace relation \bowtie by the expression to construct relation \bowtie from its fragments
 - consider horizontal fragmentation of account relation into:

acct₁ = $\sqcap_{\text{branch-name}} \bowtie$ "Baroda" (acct₁)

acct₂ = $\sqcap_{\text{branch-name}} \bowtie$ "Kodoli" (acct₂)

The query $\sqcap_{\text{branch-name}} \bowtie$ "Baroda" (acct₁ \bowtie acct₂)

becomes $\sqcap_{\text{branch-name}} \bowtie$ "Baroda" (acct₁ \sqcup acct₂)

which is optimized into

$\sqcap_{\text{branch-name}} \bowtie$ "Baroda" (acct₁ \sqcup $\sqcap_{\text{branch-name}} \bowtie$ "Baroda" (acct₂))

Database Link

- i) Use the CREATE DATABASE LINK statement to create a database link. It is a schema object in one DB that enables you to access objects on another DB.
- ii) The other DB need not be an Oracle DB system. However to access non-Oracle systems you must use Oracle Heterogeneous Service.
- iii) After you have created a DB link you can use it to refer to tables & views on the other DB.
- iv) In SQL statements you can refer to a table or view on the other DB by appending @dblink to the table or view name. You can query a table or view on the other DB with SELECT statement.
- v) You can also access remote tables & views using any INSERT, UPDATE, DELETE or LOCK TABLE statement.
- vi) To create a private DB link, you must have the CREATE DATABASE LINK system privilege. To create public DB, you must have the CREATE TABLE DATABASE LINK system privilege. Also you must have the CREATE SESSION system privilege. Also on the remote Oracle DB.

Syntax

create_database_link

→ create → shared → public → Database → link → dblink →

Updatable view

- i) It is a special case of deletable view.
- ii) A deletable view becomes an updatable view when at least one of its columns is updatable.
- iii) A column of view is updatable when all of the following rules are true.
 - The view is deletable.
 - The column resolves to column of a table & READ ONLY option is not specified.
 - All the corresponding columns of the operand of a UNION ALL have exactly matching data types & matching default values if the full select of the view include a UNION ALL.
- iv) The following examples uses constant values that can't be updated. However the view is a deletable view & at least one of its columns is updatable. Hence it is an updatable view.

CREATE VIEW updatable_view

C NO., CURRENT_DATE, CURRENT

AS

SELECT NO., CURRENT_DATE, C

FROM weather.forecast

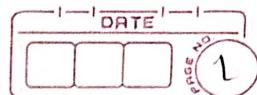
WHERE NO = 300.)

- * When can a view be updated?
- i) The view is defined based on one & only one table
 - ii) View must include primary key of table based upon which view has been created.
 - iii) View should not have any field made out of agg. fun.
 - iv) The view must not have any DISTINCT clause in its def.
 - v) Not group by, having by clause, view not have subquery
 - vi) Any of the selected output fields must not use constants, strings or value expressions.

Syntax

Update {view-name} SET {column} [,...] where {cond}

2. Advanced SQL



Relational Set Operators

SQL provides relational set operators to combine the output of two queries to generate a new relation.

- 1) Union
- 2) Union All
- 3) Intersect
- 4) Minus.

1) Union

i) The union statements combines rows from two or more queries without including duplicate rows.

ii) The syntax of UNION statement is

an ABSTRACT query UNION query .
+-----+
+-----+

iii) UNION statement combines output of two SELECT queries
SELECT statement must be union-compatible.
iv) That is, they must return same no. of attributes
& similar data types.

Example :-

customers 1			customers 2		
CUS-LName	F-Name	Phone	CUS-LName	F-Name	Phone
Mode	Swati	99294	Mode	Swati	99294
Patil	Seema	22244	Patil	Seema	22244
Shinde	Ashu	22245	Shinde	Ashu	22245
Yadav	Trepti	22246	Yadav	Trepti	22246
			Pawar	Aatav	22247
			Kumbhdas	Shreya	22248

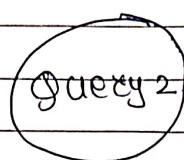
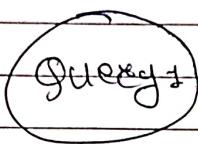
SELECT CUS-LName, CUS-FName, CUS-Phone
From customers 1

UNION

Select, CUS-LName, CUS-FName, CUS-Phone
From customers 2;

2) Union All

similar to Union just that union all returns also the duplicate values



Syntax

```
SELECT * from table 1  
UNION ALL
```

```
SELECT * from table 2;
```

when using UNION & UNION ALL column is SELECT Statements need to match exactly. This would return an error.

```
SELECT column 1 from table 1
```

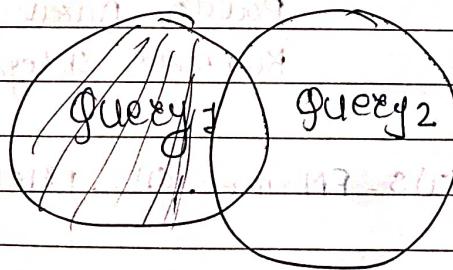
```
UNION ALL
```

```
SELECT * from table 2;
```

3) MINUS

- Minus returns the difference between the first & second SELECT statements.

- It is one where we need to be careful which statement will put it cause we will get only those results that are in the 1st SELECT statement & not in the second



Syntax

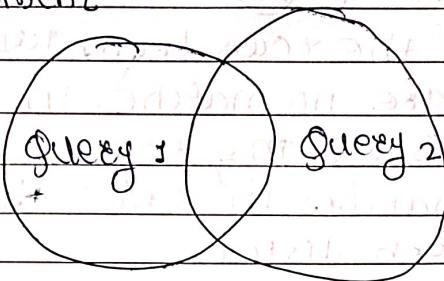
```
SELECT * from table 1  
MINUS
```

```
SELECT * from table 2;
```

④

INTERSECT

- i) It is opposite from MINUS as it returns up the results that are both to be found in 1st & 2nd select statement



SQLE

Syntax

```
SELECT * FROM table1
INTERSECT
SELECT * FROM table2
```

* SQL Join Queries

→ Inner join



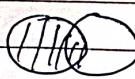
- It returns all rows from both tables where there is match - If there are rows in employees that don't match in orders, those rows will not be listed.

Syntax

```
SELECT field1, field2 FROM 1st table
INNER JOIN second table
```

```
ON first table . keyfield = 2nd table . foreign keyfield
```

2) Left outer join



- It returns all the rows from 1st table, even if there are no matches in the 2nd table (orders)
- If there are rows in employees that do not have matches in orders, those rows also will be listed.

Syntax

Select field₁, field₂ from table 1

Left join table 2

ON table1 · keyfield = table2 · foreign · key field

3) Right outer join

- It returns all the rows from table 2 (orders) even if there are no matches in the 1st table
- If there had been any rows in orders that didn't have matches in tables those rows also would have been listed

Syntax

Select field₁, field₂ from table 1

Right join table 2

ON table1 · keyfield = table2 · foreign · key field

4) Cross join

- It is also known as cartesian join, which return the cartesian product of two or more joined tables.
- It produces table that merges each row from the 1st table with each row of the 2nd table.
- It is not required to include any condition in cross join

Syntax

Select * from table 1 · table 2 · ... · table n

Cross join

table₂ ·

5) Full Outer Join

SQL Functions

- i) SQL functions are very useful tools.
- ii) when you want to list all employees ordered by year of birth or when your marketing dept. wants you to generate a list of all cust. ordered by zip code & first 3 digits of their telephone no., in both cases you will need SQL functions that can be derived from an existing attributes.
- iii) Functions always take numerical, date or string value.
- iv) the value may be part of the command itself or it may be an attribute located in table.
- v) Therefore fun. may appear anywhere in an SQL statement where value or an attribute can be used.
- vi) There are many types of SQL fun. such as arithmetic, trigonometric, string, date & time fun.

Date & time function

i) All date fun. take one parameter & return value.

a) YEAR → It returns four digit year

Syntax → YEAR(date-value)

e.g. → List all employees born in 1996

b) MONTH → It Returns two digit month code

Syntax → MONTH(date-value)

c) DAY → It returns name of day

Syntax → DAY(date-value)

d) DATE() or GETDATE() → It returns todays date

Numeric function

i) Numeric functions can be grouped in many different ways, such as algebraic, trigonometric & logarithmic.

ii) ABS → It returns absolute value of no.
Syntax → ABS (numeric_value)

iii) ROUND → It rounds a value to specified precision no. of digits.

Syntax → ROUND (numeric_value, p)
P = precision

iv) CEIL/CEILING/FLOOR

It returns the smallest integer greater than or equal to no. to returns the largest integer equal to or less than no. respectively,

Syntax → CEIL (numeric_value) — oracle
CEILING (— ") — SQL server
FLOOR (— ")

String function

i) concatenation

|| → oracle
+ — ms Access / SQL server

→ concatenation data from two different character columns & return a single column.

→ Syntax:

string_value || string_value

string_value + string_value

e.g. → List all employees names.

SELECT EmpName || ' ' || Emp_Fname AS NAME
FROM EMPLOYEE

ii) UPPER/LOWER

It returns a string in all capital or all lowercase letters

Syntax → UPPER (string_value)
LOWER (string_value)

e.g - list all employee name in all capital letters

```
SELECT UPPER(EMP_LNAME) || ',' || UPPER(EMP_FNAME) AS NAME
FROM EMPLOYEE;
```

(iii) SUBSTRING

It returns substring or part of given string parameters.

syntax → SUBSTRING (string-value, p, l)

SUBSTR (string-value, p, l)

p = start position

l = length of character

e.g → List all first three characters of all employee

```
SELECT SUBSTR(EMP_PHONE, 1, 3)
FROM EMPLOYEE;
```

iv) LENGTH

It returns no. of characters in string value

syntax → LENGTH (string-value)

e.g → lists all employee last names & length of their names; ordered descended by last name length.

```
SELECT EMP_LNAME, LENGTH(EMP_LNAME) AS NAME_SIZE
FROM EMPLOYEE;
```

conversion functions

i) It allows you to take value of given datatype & convert it to the equivalent value in another datatype.

ii) TO-CHAR & TO-DATE , TO-CHAR function takes date value & returns a character string representing a day, month or year.

iii) TO-DATE fun. takes character string representing data & returns an actual date.

Numeric to character

TO-CHAR

- i) It returns character string from numeric value
- syntax → TO-CHAR (numeric-value, fmt) → Oracle
CAST (numeric AS varchar (length)) ? SQL
CONVERT (varchar (length), numeric) SQL Server

Date to character

TO-CHAR → Oracle

CAST - SQL Server

CONVERT - SQL Server

- i) It returns character string or a formatted character string from date value

ii) syntax

Oracle → TO-CHAR (date-value, fmt)

(SQL Server → CAST (date AS varchar (length))
CONVERT (varchar (length), date))

String to Number

TO-NUMBER

- i) It returns formatted no. from character string using given format

ii) Syntax

Oracle → TO-NUMBER (char-value, fmt)

fmt = format used ; can be

EE (EAN) g = display digit amount . 999 123.12

0 = " " leading zero

> = " " the comma

· = " " decimal point

\$ = " " dollar sign

B = leading blank space (will be filled)

S = leading sign (+ or -) + sign

M = trailing minus sign

PL/SQL

- Procedural language using SQL
- By using SQL we can do many things
- PL/SQL is block structured language
- SQL limitation → does not support checking, looping, branching so that we use PL/SQL
- All statements of block are passed to oracle engine all at once which increases processing speed & decreases traffic
- PL/SQL supports functions, looping, branching & it able to handle exceptions.

PL/SQL blocks

Declare

Begin

End.

PL/SQL have some parts:

- i) stored procedure
- ii) triggers
- iii) cursors

i) stored procedure

Format: CREATE [OR REPLACE] PROCEDURE [proc-name]
[list of parameters] [as] [language]

[TS / AS]

[variable-name data-type [:= initial-value]]

BEGIN

PL/SQL or SQL statement;

END;

- i) In above syntax argument specifies the param that are passed to the stored procedure. A stored procedure could have zero or more arguments or parameters.
- ii) IN/OUT indicates whether the parameter is for input, output or both.
- iii) Data type is one of the procedural SQL data types used in RDBMS, the data types normally match those used in the RDBMS table creation statement.
- iv) Variables can be declared bet. keywords IS & BEGIN, you must specify variable name, its data type & initial value.

Example

```

CREATE OR REPLACE PROCEDURE PRC-PROD-DISCOUNT
AS BEGIN
    UPDATE PRODUCT
    SET P_DISCOUNT = P_DISCOUNT + .05
    WHERE P_QOH >= P_MIN*2;
    DBMS_OUTPUT.PUT-LINE ('* * Update finished * *');
END;
  
```

Here procedure is created.

- In above example, PRC-PROD-DISCOUNT stored procedure uses the DBMS-OUTPUT.PUT-LINE fun. to display a message when procedure executes.

- To execute stored procedure, you must use following syntax:

```
EXEC[procedure-name[(parameter-list)]];
```

* PL/SQL processing with CURSORS

- i) In SQL statement you have used inside a PL/SQL block have returned a single value.
- ii) If the SQL statement returns more than one value you will generate an error.
- iii) By using cursor, you use an SQL statement that returns more than one value inside PL/SQL code.
- iv) Cursor is a special construct used in procedural SQL to hold the data rows returned by SQL query.
- v) Cursor as a reserved area of memory in which the output of the query is stored.
- vi) Cursors are held in reserved memory area in the DBMS server, not in client computer.

There are two types of cursors.

i) Implicit

ii) Explicit

i) Implicit cursor

Implicit cursor is automatically created in procedural SQL when the SQL statement return only one value.

ii) Explicit cursor

- Explicit cursor is created to hold the output of an SQL statement that may return two or more rows. It returns only one row.

Syntax

```
CURSOR cursor_name IS select_query;
```

→ Declaring cursor

```
CURSOR C_customers IS  
select id, name, address, FROM customers;
```

opening cursor

opening cursor allocates memory for cursor & make it ready for fetching rows returned by SQL statement into it.

e.g. OPEN C_customers

* Fetching the cursor

It involves accessing one row at a time for example, we will fetch rows from above opened cursor as follows:

```
FETCH c-customers INTO c-id, c-add;
```

* Closing the cursor

Closing cursor means releasing allocated memory.

```
CLOSE c-customers;
```

Example →

c-id customers : open and close cursor command

i) Open → Opening the cursor executes the SQL command & populates the cursor with data, opening the cursor for processing

- The cursor declaration command only reserves a named memory area for the cursor, it doesn't populate the cursor with data before you can use cursor, you need to open it.

e.g → OPEN cursor-name

ii) FETCH

- once the cursor is opened, you can use FETCH command to retrieve data from the cursor & copy it to the PL/SQL variable for processing.

Syntax

```
FETCH cursor-name INTO var1 [,var2, ...]
```

- The PL/SQL variables used to hold the data must be declared in the DECLARE section & must have data types compatible with the columns retrieved by the SQL command. If the cursor SQL statement returns 5 columns, there must be 5 PL/SQL vars to receive the data from cursor.

iii) CLOSE

The CLOSE command closes the cursor for processing

CURSOR Attributes

1) %ROWCOUNT

It returns the no. of rows fetched so far.

If the cursor is not OPEN, it returns an error, if no FETCH has been done but cursor is OPEN it returns 0.

2) %FOUND

It returns TRUE if the last FETCH returned a row & FALSE if not.

- If the cursor is not OPEN, it returns an error.
- If no FETCH has been done, it contains NULL.

3) %ISOPEN

It returns TRUE if the cursor is open (ready for processing) or FALSE if the cursor is closed.

* PL/SQL Stored functions

i) stored procedures & functions are very similar.

ii) stored function is basically named group of procedural & SQL statements that returns a value i.e. (RETURN)

iii) Syntax :

CREATE FUNCTION fun-name(argument IN datatype,

RETURN datatype [IS] PL/SQL Statement;

BEGIN

PL/SQL Statement;

RETURN (value or expression);

END;

iv) stored fun. can be invoked only from within stored procedures / triggers & can't be invoked from SQL statements.

PL/SQL Stored Functions

i) A stored function is basically named group of procedural & SQL statements that returns a value i.e indicated by a RETURN statement in its program code.

ii) Syntax

```
CREATE FUNCTION fun-name (argument IN data-type, -)
RETURN data-type [IS]
BEGIN
    PL/SQL statements;
    RETURN (value/expression);
END;
```

iii) Stored functions can be invoked only from within stored procedures or triggers & cannot be invoked from SQL statements.

Example

```
DECLARE
    a number;
    b number;
    c number;
```

```
FUNCTION findMax (x IN number, y IN number)
RETURNS number
```

```
CREATE FUNCTION findMax (x IN number, y IN number)
RETURNS number
BEGIN
```

```
IF x > y THEN
```

```
    z := x;
```

```
ELSE
```

```
    z := y;
```

```
END IF;
```

```
RETURNS number
```

```
END IS;
```

```
BEGIN
```

```
a := 23;
```

```
b := 45;
```

```
c := findMax (a, b);
```

```
dbms_output.put_line ('maximum of (23, 45): ' || c);
```

Embedded SQL

i) Embedded SQL is a term used to refer to SQL statements that are contained within an application programming lang. such as visual Basic, .Net, C#, or java.

ii) static/embedded SQL are SQL statements in an application that do not change at runtime & therefore can be hard coded into the appn

* Dynamic SQL

- It is a programming technique that enables you to build SQL statements dynamically at runtime.

- You can create more general purpose, flexible appn by using dynamic SQL because the full text of SQL statement may be unknown at compilation.

Embedded SQL

Dynamic SQL

i) DB access procedure in pre-DB will be accessed can be determined in the statement determine only at run time

ii) statements are compiled at compile time and executed at run time

iii) stat: execute, prepare, used
immediate not used

Used

Programs like soft

when cond" is true
then trigger is automatically
occurred

DATE		
PAGE NO.	16	

when event is occurred & condition

Trigger - is true, execute the action

- i) It is a stored procedure in DB which automatically invokes whenever a special event a in the DB occurs.
- ii) e.g. → A trigger can be invoked when a row is inserted into a specified table or when certain table columns are being updated.

iii) syntax for creating trigger

`create trigger [trigger-name]`

[before | after]

{ insert | update | delete } - types of events

on [table-name]

[for each row] -

[trigger-body] - body of trigger

i) `create trigger[trigger-name]` → creates / replace
an existing trigger with the trigger-name

ii) [Before | after] :

This specifies when the trigger will be executed

iii) { insert | update | delete } -
This specifies the DML operation.

iv) on [table-name]: to bind the trigger to the table

This specifies the name of the table associated with the trigger.

v) [for each row]:

This specifies row level trigger i.e. trigger will be executed for each row being affected.

vi) [trigger-body]:

This provides the operation to be performed as trigger is fired.

Example: `create table empal`

`before|after insert|update on employee`

Triggers

- i) It cannot execute manually it is big difference between stored procedure & trigger
- ii) They only execute in response to user action like INSERT.
- iii) It is used to check the validity of data that are inserted into table.
- iv) If you need a certain piece of code to always be executed in response to an event, the best option is to use triggers.

Example →

Table ① create table Employee
 (id int primary key, name varchar(45),
 sal int, gender varchar(12), deptId int)
 insert into employee values ('12ABC', 10000, Male,
 1) if insertion fails return a message
 Table ② create table empAuditTest
 (id int identity, AuditAction text)

Now, we create trigger that stores the record of each insert operation on temp table into the Employee Audit-Test table.

- Create trigger trigInsertEmployee
 on employee
 for insert

AS

BEGIN

Declare @id int

select @id = id from inserted

insert into EmpAuditTest

values ('New empl with Id = ' + cast(@id as varchar(10)) +
 End ' is added at ' + cast(getdate() as varchar(22)))'

If you could create sequence to automatically assign values to the customer code each time new customer is added & create another sequence

To check sequence →

→ SELECT * FROM USER SEQUENCES;

After creating trigger, we will add following

Insert into employee values (6, 'Sachin', 36000, 'female')

(13) Insert Data

Delete

Create Trigger Delete emp
on Employee
For delete

As

BEGIN

Declare @ Id integer for bit 1

Select @ Id = Id from deleted

Insert into emp-AuditTest
values ('An existing empl with Id = '+CAST
(@Id AS VARCHAR(10))+' is deleted at '+
CAST(GETDATE() AS VARCHAR(22)))

END

After creating trigger

Delete from employee where id = 2

Output after inserting data

EmployeeID Name Gender Age

1 Sachin male 36

2 Rakesh male 30

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

Difference bet. fun stored function & stored procedures

function

Procedure

- | | |
|------------------------------------------------------------------------|----------------------------------------------------------------------|
| i) It must return a value | It is optional, it returns zero |
| ii) Functions can have only input parameters | It can have input or output parameters |
| iii) Functions can be called from procedure | Procedures cannot be called from function |
| iv) It allows only select statement in it | It allows select as well as DML statements |
| v) It can be used in join clause as a result set | It can't be used in join clause |
| vi) Used in SQL statements anywhere in WHERE / HAVING / SELECT section | It can't be used in SQL statements in WHERE / having / select clause |
| vii) Try catch block can't be used | An exception can be handled by try catch block in this |
| viii) We use select command to execute a function | We use EXECUTE Keyword to execute stored procedure |

Oracle SEQUENCE

i) Oracle does not support Autonumber data type or the identity column property, instead you can use sequences to assign values to column on a table. But an Oracle sequence is very different from Access Autonumber data type & deserves close security.

- Oracle sequences are an independent object in the DB.
- Sequence are not data type.
- Oracle sequences are not tied to table or column.
- Oracle sequence have name & can be used anywhere value is expected.

→ Oracle sequences generate a numeric value that can be assigned to any column in any table.

→ The table attribute to which you assigned a value based on sequence can be edited & modified.

→ Oracle sequence can be created & deleted anytime.

Oracle Syntax

CREATE SEQUENCE name [START WITH n]

[INCREMENT BY n] [CACHE | NOCACHE]

where

- name is the name of sequence
- n is the integer value that can be +ve or -ve.
- START WITH specifies initial sequence value
- INCREMENT BY determines the value by which sequence is incremented.
- The CACHE / NOCACHE clause indicates whether Oracle will preallocate sequence no. in memory (by default 20 values).

e.g. → CREATE SEQUENCE CUS-SEQ START WITH 2010 NOCACHE;

If you could create sequence to automatically assign values to the customer code each time new customer is added & create another sequence.

To check sequence.

→ SELECT * FROM USER SEQUENCES;

Subquery

- i) A subquery is best defined as query within query. Subquery enable you to write queries that select rows for criteria that are actually developed while the query is executing at run time.
- ii) It is the use of SELECT statement inside one of the clauses of another SELECT statement.
- iii) In fact subquery can be contained inside another subquery, this is inside another subquery & so forth.
- iv) A subquery can also be nested inside INSERT, UPDATE & DELETE statement. Subqueries must be enclosed within parentheses.
- v) A subquery can be used anywhere where an expression is allowed providing it returns a single value, this means that a subquery that returns a single value can also be listed as an object in a FROM clause listing.
- vi) This is termed inline view because when subquery is used as part of from clause, it is treated like a virtual table or view.
- vii) Subquery can be placed either in FROM clause, WHERE clause or HAVING clause of main query.

Types of subqueries

- i) Single Row subquery
- ii) Multiple row subquery
- iii) Correlated subquery
- iv) Single Row subquery
 - i) Subquery which returns single row output.
 - ii) They make the usage of single row comparison operators like used in WHERE conditions.
- 2) Multiple row subquery
 - i) Subquery returning multiple row output.
 - ii) They make use of multiple row comparison operators like IN, ANY, ALL.
 - iii) There can be sub queries returning multiple columns also.

Correlated sub query

102 23 Oct 2023

- i) It depends on data provided by the outer query.
- ii) This type of subquery also includes subqueries that use the EXIST operator to test the existence of data rows satisfying specified criteria.
- iii) Correlated subquery, the inner query is executed repeatedly once for each row that might be selected by outer query.
- iv) It produces result tables that answer complex management questions.

→ ~~SELECT Employee_ID, salary, Dept_id
FROM employees E~~
→ ~~WHERE salary > (SELECT AVG(salary)
FROM emp T)~~

→ ~~Dept_id = T.dept_id~~

Consider above query, unlike the subqueries previously considered, the subquery in this statement cannot be resolved independently of the main query. Notice that the outer query specifies that rows are selected from the employees table with an alias name of two. The inner query compares the employee dept_no column of the employee table with alias e2 to the same column for the alias table name e1.

∴ ~~Dept_id = T.dept_id~~

→ ~~Dept_id = e2.dept_id~~

∴ ~~Dept_id = e1.dept_id~~

* PostgreSQL

prep due batch 9 & 10

- i) It is an enterprise class open source database management system with standards compliant SQL & JSON for relational & non-relational queries for extensibility of SQL.
- ii) It supports both SQL & JSON for relational & non-relational compliance.
- iii) It supports advanced data types & performance optimization features, which are only available in expensive commercial db like Oracle, MySQL etc.
- iv) It is known as PostgreSQL farm.
- v) It is backed by an experienced community of developers who have made tremendous contributions to make it a highly reliable DBMS system.

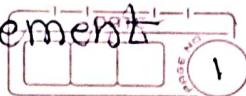
Advantages

- i) It can run dynamic websites & web apps as LAMP stack option.
- ii) It's write-ahead logging makes it a highly fault-tolerant & durable DBMS system.
- iii) Source code is freely available under an open-source license. This allows you freedom to use & modify its implementation as per your business needs.
- iv) It supports geographic objects so you can use it with for location-based scenarios.
- v) To learn PostgreSQL you don't need much learning as it's easy to use.
- vi) Low maintenance & administration for both embedded & enterprise use of PostgreSQL.

Disadvantages

- i) It is not owned by one organization, so it has heavy trouble getting its name out there despite being fully featured & comparable to other DBMS systems.
- ii) Many open source apps support MySQL, but may not support PostgreSQL.

3. NoSQL database Management



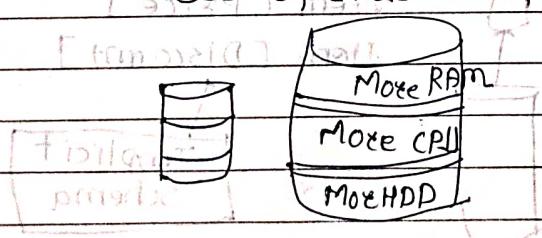
NoSQL: Non-relational or distributed database

- i) NoSQL database is a non-relational database management system, that does not require a fixed schema.
- ii) It avoids joins & is easy to scale.
- iii) The major purpose of using NoSQL database is for distributed data stores with humongous data storage needs.
- iv) NoSQL is used for Big data & real time web apps.
e.g. → Twitter, Facebook & Google collect terabytes of user data every single day.
- v) It stands for "Not Only SQL" & "NoSQL".
- vi) RDBMS uses SQL syntax to store & retrieve data for further insights.
- vii) NoSQL database system encompasses a wide range of database technologies that can store structured, semi structured, unstructured & polymorphic data.

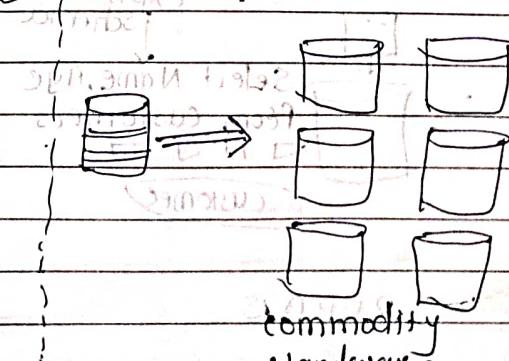
Why NoSQL

- i) The concept of NoSQL database became popular with internet giants like Google, Facebook, Amazon etc. who deal with huge volumes of data.
- ii) The system response time becomes slow when you use RDBMS for massive volumes of data.
- iii) To resolve this problem we could "scale up" our systems by upgrading our existing hardware so this process is expensive.
- iv) The alternative for this issue is to distribute database load on multiple hosts whenever the load increases, this method is known as scaling out.

scale up (vertical scaling)



Scale-Out (horizontal scaling)



→ NoSQL database is non-relational, so it scales out better than relational databases as they are designed with web applications in mind.

History of NoSQL DB

- 1998 - Carlo Strozzi uses the term NoSQL for his lightweight open source relational database
- 2000 - Graph database
- 2004 - Google BigTable is launched
- 2005 - CouchDB is launched
- 2007 - The research paper on Amazon Dynamo is released
- 2008 - Facebook opensources the cassandra project
- 2009 - the term NoSQL was reintroduced.

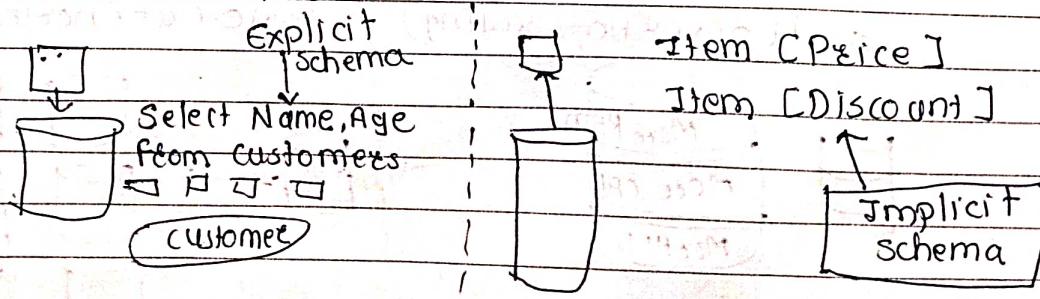
Features of NoSQL

1) Non-relational

- i) NoSQL database never follows the relational model.
- ii) Never provide tables with flat fixed-column records
- iii) Work with self contained aggregates or BLOBs
- iv) Doesn't require object relational mapping & data normalization.
- v) No complex features like query language, query planners, referential integrity, joins, ACID.

2) Schema free

- i) NoSQL databases are either schema free or have relaxed schemas.
- ii) Do not require any sort of definition of the schema of the data.
- iii) offers heterogeneous structures of data in the same domain.



NoSQL is schema free.

3) Simple API

- i) Offers easy to use interfaces for storage & querying data provided.
- ii) APIs allow low level data manipulation & selection methods.
- iii) Text based protocols mostly used with HTTP REST with JSON.

iv) Mostly used no standard based NoSQL query lang.

v) Web-enabled databases running as internet facing services.

4) Distributed

i) Multiple NoSQL databases can be executed in a distributed fashion.

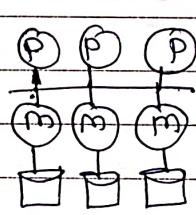
ii) Offers auto scaling & fail over capabilities.

iii) Often ACID concept can be sacrificed for scalability & throughput.

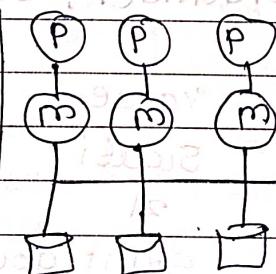
iv) Mostly no synchronous replication between distributed nodes. Asynchronous Multi-Master Replication, peer-to-peer, HDFS Replication.

v) Only providing eventual consistency.

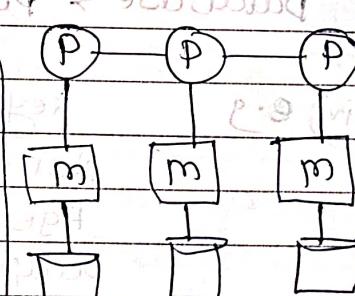
vi) Shared Nothing Architecture. This enables less co-ordinations & higher distribution.



shared mem.
e.g. Oracle 11g



Shared disk
e.g. Oracle RAC



shared Nothing
e.g. NoSQL

Highly available + Scalable

Types of NoSQL Database

i) NoSQL databases are mainly categorized into,

1) Key value pair

2) Column oriented

3) Graph based

4) Document oriented.

ii) Every category has its unique attributes & limitations.

iii) None of the above specified database is better to solve all problems, user should select the DB based on their product needs.

① Key value pair

i) Data is stored in key value pairs. It is designed in such a way to handle lots of data & heavy load.

ii) Key value pair storage database store data as hash table where each key is unique & the value can be JSON, BLOB (Binary Large Objects) string etc.

iii) e.g. → key value pair may contain a key like "Website" associated with a value like "Google".

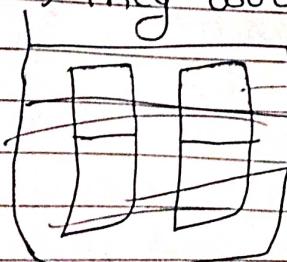
Database → Project Voldemort, Riak, Redis.

iv) e.g.

key	value
Name	Scouti
Age	21
Occupation	Stunt double
Height	175cm
Weight	67kg

v) This kind of NoSQL database is used as collection, dictionaries, associative arrays etc.

vi) They work best for shopping cart contents.



example → Big table, cassandra, Hbase, Hypertable.

Key	Value
Count	100

hashing & framing @

to split partition frame all @

and split online use n.o.p

frame with sh

complaint index vs index of buckets of frame @ add @

key at end of key and bucket @ 2nd last part of

Example → Riak, Tokyo.cabinet, Redis sevee, memcached
+ scalaris.

② column Based	value
row	value
column	value

- i) column oriented database works on document column & are based on Big Table paper by Google.
- ii) Every column is treated separately
- iii) values of single column database are stored contiguously.

iv) e.g. →

column family		
Row	column Name	value
key	key	key
oldid@123456	value	value
oldid@123456	value	value
oldid@123456	value	value

→ The database delivers high performance on aggregation queries like SUM, COUNT, AVG, MIN etc. as the data is readily available in a column.

b) vi) Column based NoSQL databases are widely used to manage data warehouse, business intelligence, CRM, library card catalogs.

vii) HBase, Cassandra, HyperTable are NoSQL query example of column based database.

② Document Oriented

- i) Document oriented NoSQL DB stores & retrieve as a key value pair but the value part is stored as document.
- ii) The document is stored in JSON or XML formats.
- iii) The values is understood by DB & can be queried.

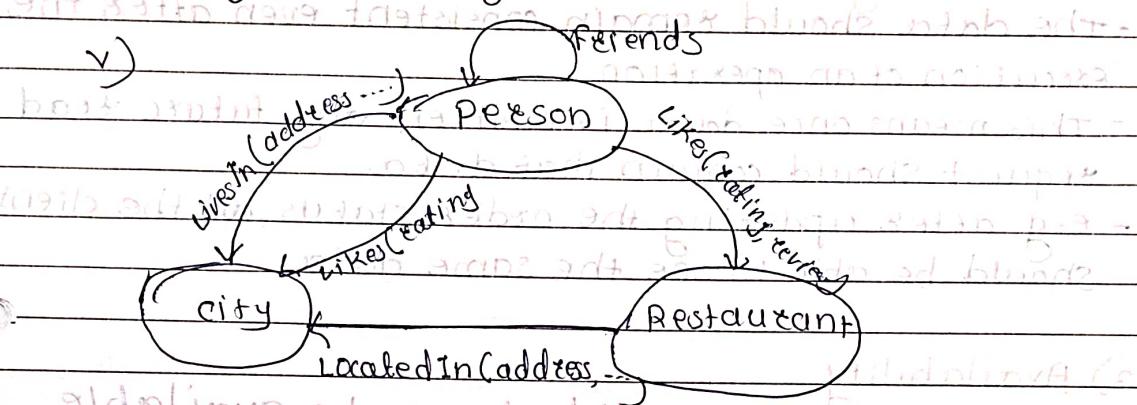
col 1	col 2	Document 1
Data	Data	{ "prop1": data,
Data	Data	"prop2": data,
Data	Data	"prop3": data,
Data	Data	}

Relational Vs. Document

- iv) In above fig. on your left you can see we have rows & columns & in the right we have a document in database which has similar structure to JSON.
- v) In relational db, you have to know what columns you have & so on, however for document database you have data store like JSON object, you do not require to define which make it flexible.
- vi) The document type is mostly used for CMS system, blogging platforms, real time analytics & e-commerce appn.
- vii) It should not use for complex transactions which require multiple operations or queries against varying aggregate structures.
- viii) Amazon simple DB, couch DB, mongoDB, Riak, Lotus Notes are popular document originated DBMS system.

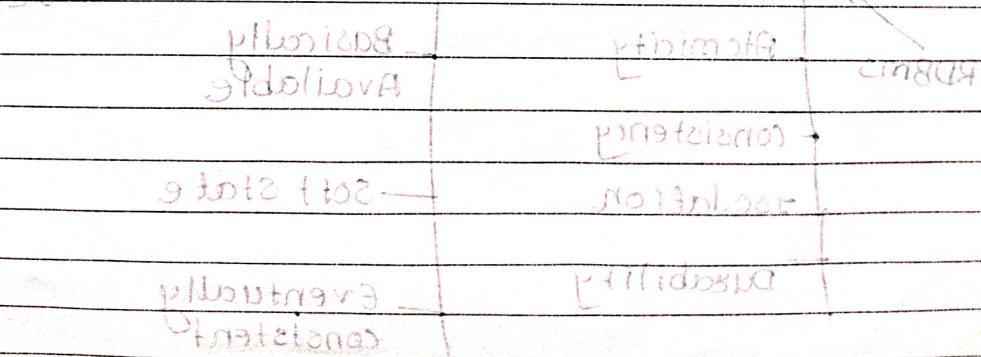
④ Graph based

- i) Graph based database stores entities as well as the relations amongst those entities.
- ii) The entity is stored as a node with the relationship as edges.
- iii) An edge gives relationship between nodes.
- iv) Every node & edge has a unique identifier.



- v) Compared to a relational db where tables are loosely connected, graph database is a multi-relational in nature.
- vi) Traversing relationship is fast as they are already captured into the DB if there is no need to calculate them.
- vii) Graph based database mostly used for social networks, logistics, spatial data.
- ix) Neo4j, infinite graph, OrientDB, FlockDB are some popular graph based database.

QUESTION → ANSWER



CAP theorem (Brewer's theorem)

1) Consistency

2) Availability

3) Partition Tolerance

1) Consistency

- The data should remain consistent even after the execution of an operation.

- This means once data is written any future read request should contain that data.

- e.g. after updating the order status, all the clients should be able to see the same data.

2) Availability

- The database should always be available & responsive.

- It should not have any downtime.

③ Partition Tolerance

The database should always be available & responsive.

- The system should continue to function even if the communication among the servers is not stable.

- e.g. → servers can be partitioned into multiple groups which may not communicate with each other. Here if part of the db is unavailable other parts are always unaffected.

ACID	BASE	Nosql
Atomicity	Basically Available	
Consistency		
Isolation	Soft state	
Durability	Eventually consistent	

1) Basically Available

Database is available all the time as per CAP theorem

2) Soft state

Means even without an input, the system state may change.

3) Eventually consistency

- Means that the system will become consistent over time.

- Eventually consistency means to have copies of data on multiple machines to get high availability & scalability. Thus changes made to any data item on one machine has to be propagated to other replicas.

- Data replication may not be instantaneous as some copies will be updated immediately while others in due course of time, these copies may be mutually inconsistent.

Hence the name eventual consistency.

Q. Explain the various types of consistency.

ANSWER:

1) Strong consistency

Ensures that if a write operation is successful, then all other reads will see the same value.

It is also known as serializable consistency.

It is achieved by using atomic operations or locks.

It is also known as strict consistency.

It is used in distributed systems like replicated databases.

2) Weak consistency

3) Eventual consistency

4) Soft state

5) Basically available

Database in MongoDB

- i) MongoDB is an open source, cross platform document oriented database that provides high performance, high availability & easy scalability.
- ii) It works on concept of Collection & documents.

Database in MongoDB

- i) Database is physical container for collections.
- ii) Each database gets its own set of files on the file system.
- iii) A single MongoDB server typically has multiple db collection.
- i) collection is a group of MongoDB documents.
- ii) collection exists within a single database.
- iii) It do not enforce a schema.
- iv) Documents within collection can have different fields.
- v) All documents in a collection are of similar or related purpose.

Document

- i) A document is a set of key-value pairs.
- ii) It have dynamic schema means that documents in the same collection do not need to have same set of fields or structure & column fields in collections documents may hold different types of data.

RDBMS

Database

Table

Tuple / Row
column

Table join

Primary key

MySQL / Oracle

MySQL / SQLplus

MongoDB

Database

collection

Document

Field

Embedded documents

Primary key (by default generated)

mongod (db server)

Mongo client

Advantages of MongoDB over RDBMS:

i) Schemaless

MongoDB is document database in which one collection holds different documents. No. of fields, content & size of the document can differ from one document to another.

ii) Structure of single object is clear

iii) No complex joins

iv) Deep query ability

It supports dynamic queries on documents using a document based query language that's nearly as powerful as SQL tuning.

v) Ease of scale out

MongoDB is easy to scale.

vi) Conversion / mapping of appn objects to db objects not needed

vii) Uses internal memory for storing the working set;

Installation of MongoDB

1) Pre-Requisite information on windows

1st check your windows 32 / 64 bit

go to cmd

→ C:\ wmic os get osarchitecture

above command shows the window is 32 / 64

2) Download MongoDB on windows from website

select community server

then select

windows [64-bit]

windows [32-bit]

then download

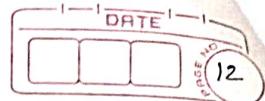
3) Open exe file, click next & accept license

agreement & goes to the next & choose

setup type, there are two setup type

i) complete & ii) custom

you may select complete & click on install & click on finish.



④ Go to the C disk & get into program files & then you will find "MongoDB" folder. Lock disk C → Prog file → MongoDB → server → bin copy above path.

⑤ Go to the computer & props & Advanced then select env. variable & paste above path & edit & then OK, OK, OK.

⑥ Download Robo 8T Setup, & open exe file, then go to next button, agree license & then browse destination folder.

e.g. C:\program files\Robo 8T 1.2.1 & click on next & select Robo 8T 1.2.1 & install then click on finish button.

⑦ MongoDB connection setting
click on create conn. Other one windows is displayed
Type - Direct connection

Name - New Conn.
Address - localhost : 27017

& save it then connect.

Accessing MongoDB

① Open cmd, copy Path upto MongoDB server folder
C → Progfile → mongodb-server-bin

& paste above path on cmd.

Here paste above path /mongo.exe -version

it directly goes to the Mongod command

otherwise you can directly open the mongoDB cmd by using below steps

→ Select C → Progfile → mongodb-server-bin → mongod

click on mongod, it directly open mongoDB cmd.

Next day no 30/08/2018 at 11:47 AM

• Start the file

CRUD operations

① Create Database

- i) MongoDB use `Database_Name` is used to create a database, The command will create a new database, if it doesn't exist otherwise it will return existing db.

Syntax

`use DATABASE_NAME`

example → `> use swati`

switched to db mydb

- ii) To check your currently selected database use command `db`

Syntax → `db database_name`

example → `> db swati`

- iii) If you want to check your database list, then use command `show dbs`

→ `show dbs`

local mydb test testdb testt

swati

② Drop database

Syntax → `dropDatabase c`

(`db: dropdatabase c`)

- i) This will delete the selected database, if you have not selected any database, it will delete default `test` database.

- ii) In mongoDB default database is `test` if you didn't create any database then collection will be stored in `test` database.

② Create documents

i) To create database

→ Use Database name

e.g → Use Swati as database name, MongoDB

switched to db swati & db is selected

ii) To show database

> show dbs

It shows all database that are present in our MongoDB.

& → show dbs

admin 0 : 0006B

config 0 : 0006B

local 0.000GB db < 4. 31GB

Springboot 0.000GB db of 4000 MB free

swati 0.000GB free memory on disk

If database does not contain any document then it doesn't show db, so u can insert atleast one documents.

iii) Insert single Record (single doc)

> db.employee.insert({name: "Swati", email: "abc@gmail.com", age: 28})

inserted 1 document inserted into employee (1 document)

Inserted id : ObjectId ("5f1d1f1f1f1f1f1f1f1f1f1f")

Object id is the unique id assigned by MongoDB at insertion time and it is used for further operations.

iv) Insert many Record (many docs)

db.employee.insertMany([{"name": "Swati", "age": 28}])

* Difference betn couchDB & MongoDB

Apache couch DB

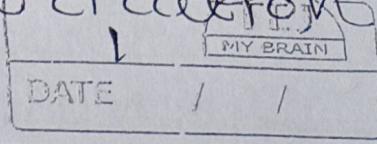
- i) It is developed by Apache software foundation & initially released in 2005.
- ii) It is written in Erlang
- iii) It is open source DB that uses different formats & protocols to store, transfer & process its data.
- iv) It uses JSON to store data, javascript as its query lang. using Map Reduce.
- v) Documents & primary unit of data in couchDB & they also include metadata.
- vi) Document fields are uniquely named & contain values of varying types & there is no set limit to text size or element count.

couch DB

Mongo DB

- | | |
|-------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------|
| 1) Data store in JSON format | Datastores in BSON format |
| 2) The DB contains documents | The DB contains collections |
| 3) It favors availability | It favors consistency |
| 4) It is written in Erlang | It is written in C++ |
| 5) It is eventually consistent | It is strongly consistent |
| 6) MongoDB is faster than couchDB | Mongo DB provides faster read speeds. |
| 7) It follows Map/Reduce query method | Map/Reduce creating a collection & object-based query language |
| 8) couchDB uses map/reduce functions & it will difficult for users with traditional SQL learning experience | Mongo DB is easier to learn as it is closest in syntax to SQL |

4. Database Security & Authorization



Introduction

There are 3 main objectives when designing secure database application:

- 1) Secrecy → Information should not be disclosed to unauthorized users.
e.g. → Student should not be allowed to examine other students grades.
- 2) Integrity → Only authorized users should be allowed to modify data.
e.g. → Student may be allowed to see their grades, yet not allowed to modify them.
- 3) Availability → Authorized user should not be denied access.
e.g. → An instructor who wishes to change a grade should be allowed to do so.

To achieve these objectives, a clear & consistent security policy should be developed to describe what security measures must be enforced..

Access control

- A database for an enterprise contains a great deal of info. & usually has several groups of users.
- Most users need to access only small part of the DB to carry out their tasks.
- User's unrestricted access to all the data can be undesirable & DBMS should provide mechanisms to control access to data.
- A DBMS offers two main approaches to access control:
 - i) Discretionary access control
 - ii) Mandatory access control

i) Discretionary Access control

- It is based on the concept of access rights, or privileges & mechanisms for giving users such privileges.
- A privilege allows a user to access base data object.
- A user, who creates a database object such as a table or view automatically, gets all applicable privileges on that object.
- SQL supports discretionary access control through the GRANT & REVOKE commands.

GRANT

- The GRANT command gives users privileges to base tables & views.
- The syntax of this command is as follows

GRANT privilege ON object TO users [WITH GRANT OPTION]

- For our purposes object is either a base table or a view. Several privileges can be specified, including these.

i) SELECT

The right to access (read) all columns of the table specified as the object including columns added later through ALTER TABLE commands.

ii) INSERT (column name)

The right to insert rows with non null values in the column of the table named as object. If this right is to be granted with respect to all columns, including columns that might be added later, simply use INSERT.

The privileges UPDATE (column name) & UPDATE are similar.

iii) DELETE

The right to delete rows from the table named as object.

REFERENCES (column-name):

- The right to define foreign keys that refer to the specified column of the table object.
- REFERENCES without column name specified denotes this right w.r.t all columns, including any that are added later.

GRANT & REVOKE ON VIEWS & INTEGRITY CONSTRAINTS

- The privileges held by the creator of view change over time as he/she gains or loses privileges on the underlying tables.
- If the creator loses a privilege held with the grant option, users who were given that privilege on the view lose it as well.
- There are some subtle aspects to the GRANT & REVOKE commands when they involve views or integrity constraints.
 - We consider some e.g. that highlight the foll. imp. pair:
 - i) A view may be dropped because a SELECT privilege is revoked from the user who created the view.
 - ii) If the creator of a view gains additional privileges on the underlying tables, he/she automatically gains additional privileges on the view.

The need for & Role of database in organization

1) Data are used by different people in different dept. for different reasons.

DBMS facilitates

- a) interpretation & presentation of data in useful formats by transforming raw data into info.
- b) distribution of data & info. to the right people at right time
- c) Data presentation & monitoring the data usage for adequate periods of time.
- d) control over data duplication & use both internally & externally

- whatever the type of org. the db predominant role is to support managerial decision making at all levels in the org. while preserving data privacy & security

- An org. managerial structure might be divided into 3 levels

- 1) Top → strategic decision
- 2) middle → tactical decision
- 3) operational → daily operational decision

i) Top management derived

- i) provide the info. necessary for strategic decision making strategic planning, policy formulation & goal def'n.
- ii) provide access to external & internal data to identify growth opportunities & to chart the direction of such growth.
- iii) provide a framework for defining & enforcing organizational policies.
- iv) improve the likelihood of the return on investment for the company by searching for new ways to reduce costs & by boosting productivity.
- v) provide feedbacks to monitor whether the company is achieving its goals.

2) Middle management level

- i) Deliver the data necessary for tactical decision & plan.
- ii) Monitor & control the allocation & use of company resources & evaluate the performance of various dept.
- iii) Provide frameworks for enforcing & ensuring the security & privacy of data in db

Privacy → It deals with the rights of individuals & the org. to determine the "who, what, when, where & how" of data usage.

Security → It protecting the data against accidental & intentional use by unauthorized users.

3) Operational mgmt. level

- i) Represent & support the company opers as closely as possible. The data model must be flexible enough to incorporate all required present & expected data.
- ii) Produce query results within specified performance levels. The performance requirements increase for lower levels of mgmt & opern.
- iii) Enhance the company's short term operational ability by providing timely info. for customer support & for appln development mgmt & compute & opern..

* Evolution of DB Administration function

- Data adminstrⁿ has its roots in odd, decentralized world of the file system
- The cost of data & managerial duplication in such file system gave rise to centralized data adminstrⁿ fun. Known as the electronic Data Processing (EDP) or data processing (DP) dept.
- The DP dept. to evolve into an Information system (IS) dept.
- Responsibility of IS
 - i) A service fun. to provide ^{end} users with active data mgmt. support

ii) A production fun. to provide ^{end} users with specific solⁿ for their info. needs through integrated appn or mgn: info. system.

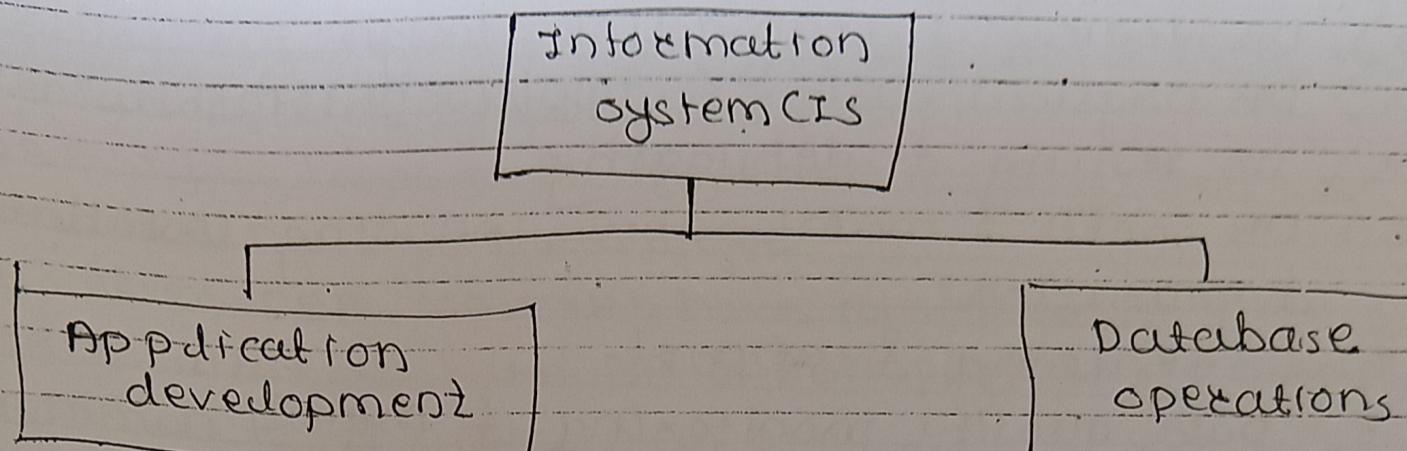
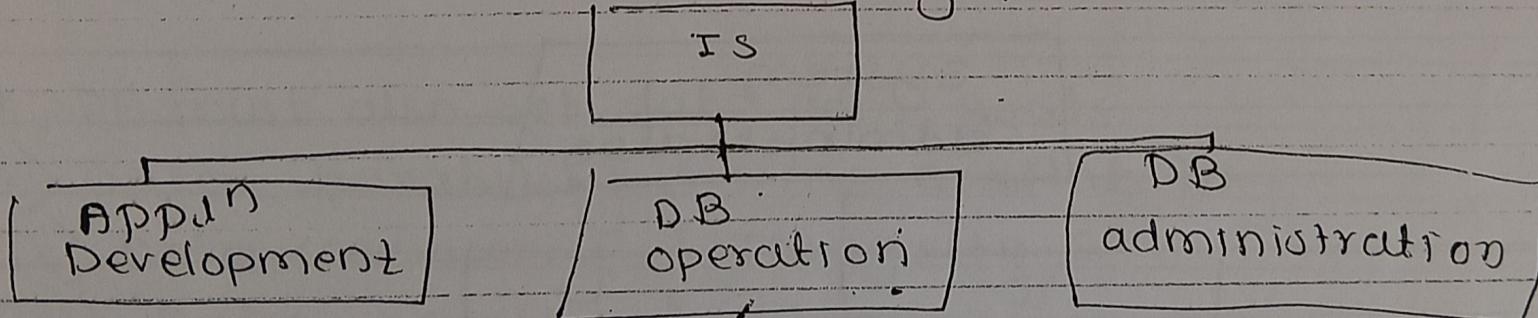


fig → IS dept. internal org.

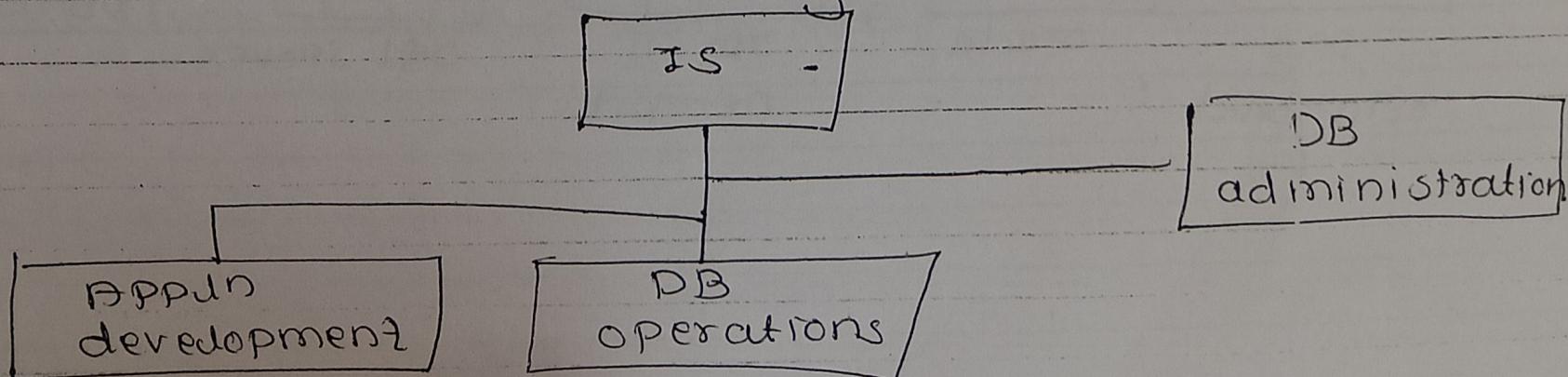
In above fig. as the demand for appn development grew the IS appn development segment was subdivided by the type of supported system: accounting, inventory, marketing & so on. This development meant that DB admin responsibilities were divided. The appn development segment was in charge of gathering DB requirements & logical DB design whereas the DB opern segment took charge of implementing, monitoring & controlling the DBMS opern.

* Placement of DBA Fun.

Line authority position

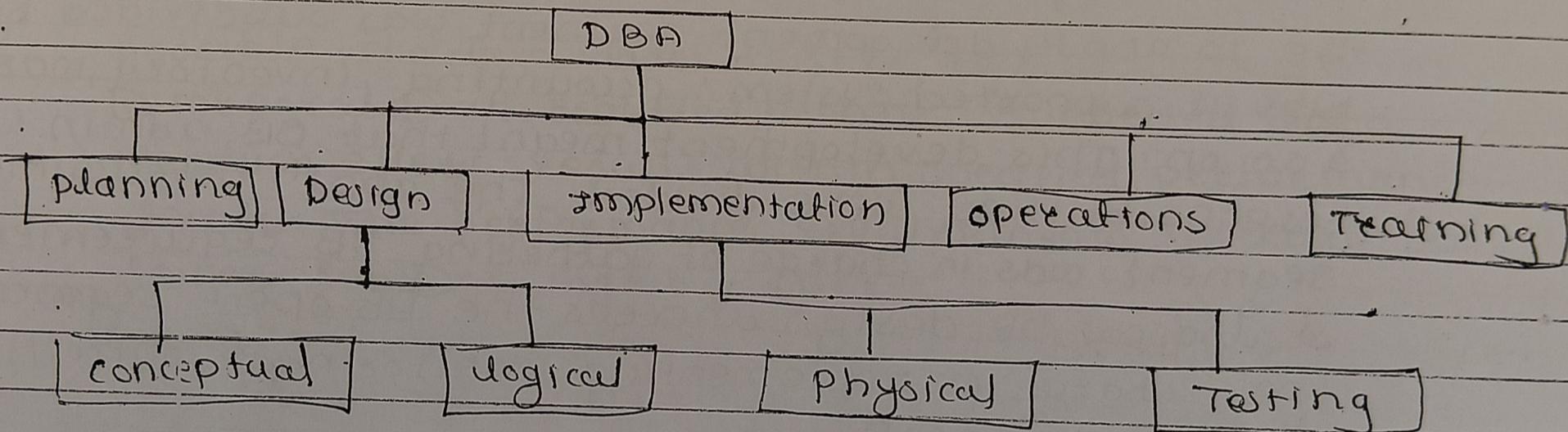


staff consulting position:

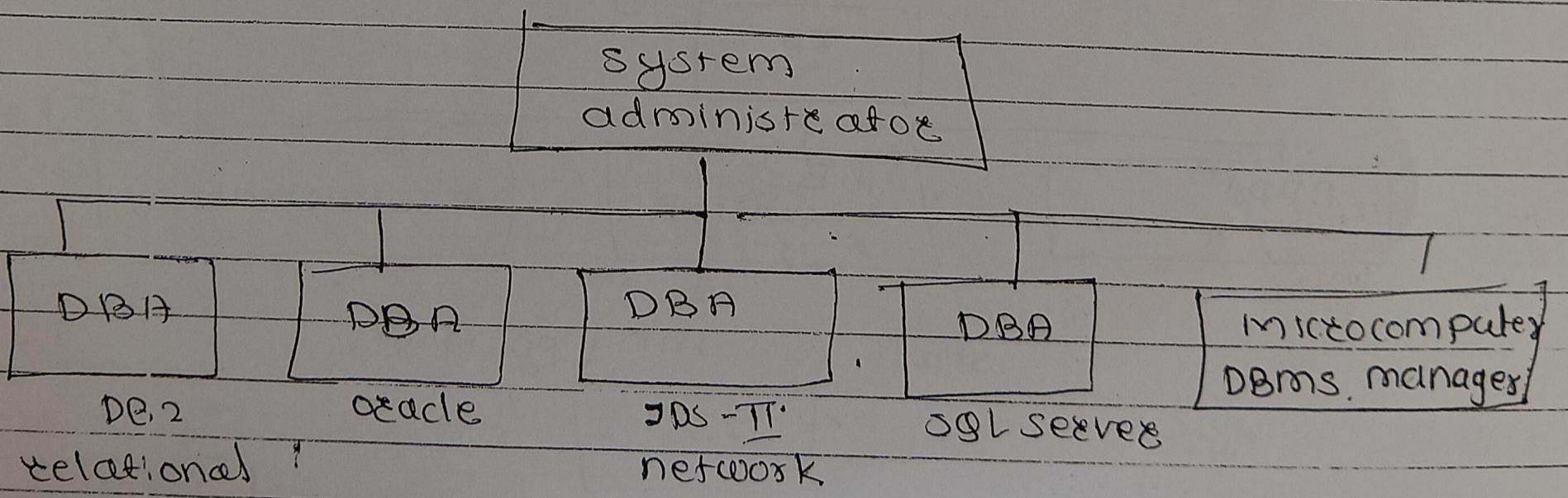


- DBF requires personnel to cover the foll. activities
- 1) DB planning, including the def'n of standards, procedures & enforcement.
 - 2) DB requirements gathering & conceptual design.
 - 3) DB logical & transaction design.
 - 4) DB physical design & implementation.
 - 5) DB testing & debugging.
 - 6) DB opern & maintenance, including installation, conversion & migration.
 - 7) DB training & support.
 - 8) Data quality, monitoring & mgmt.

A DBA functional organization.



Multiple DB admin. in an organization



The DB environments Human components

- i) Effective data admin. requires both technical & managerial skills. e.g. the DA's job typically has strong managerial orientation with company wide scope.
- ii) DBA's job tends to be more technically oriented & has a narrower DBMS specific scope!
- iii) However, DBA, too must have considerable store of people skills after all both DA & DBA perform "people" functions common to all dept. in an org.
 E.g. → Both DA & DBA direct & control personnel staffing & training within their respective dept.

DA	DBA
i) Does strategic planning	controls & supervises
ii) sets long term goals	Executes plans to reach goals
iii) Sets policies & standards	Enforces policies & procedures
iv) Is broad in scope	Enforces programming std Is narrow in scope
v) Focused on the long term	Focuses on short term
vi) Has managerial orientation	Has a technical orientation
vii) Is DBMS independent	Is DBMS specific

* DA must also set data adm. goals

- i) Data sharability & time availability
- ii) Data consistency & Integrity
- iii) Data security & privacy
- iv) Data quality std
- v) Extent & type of data use

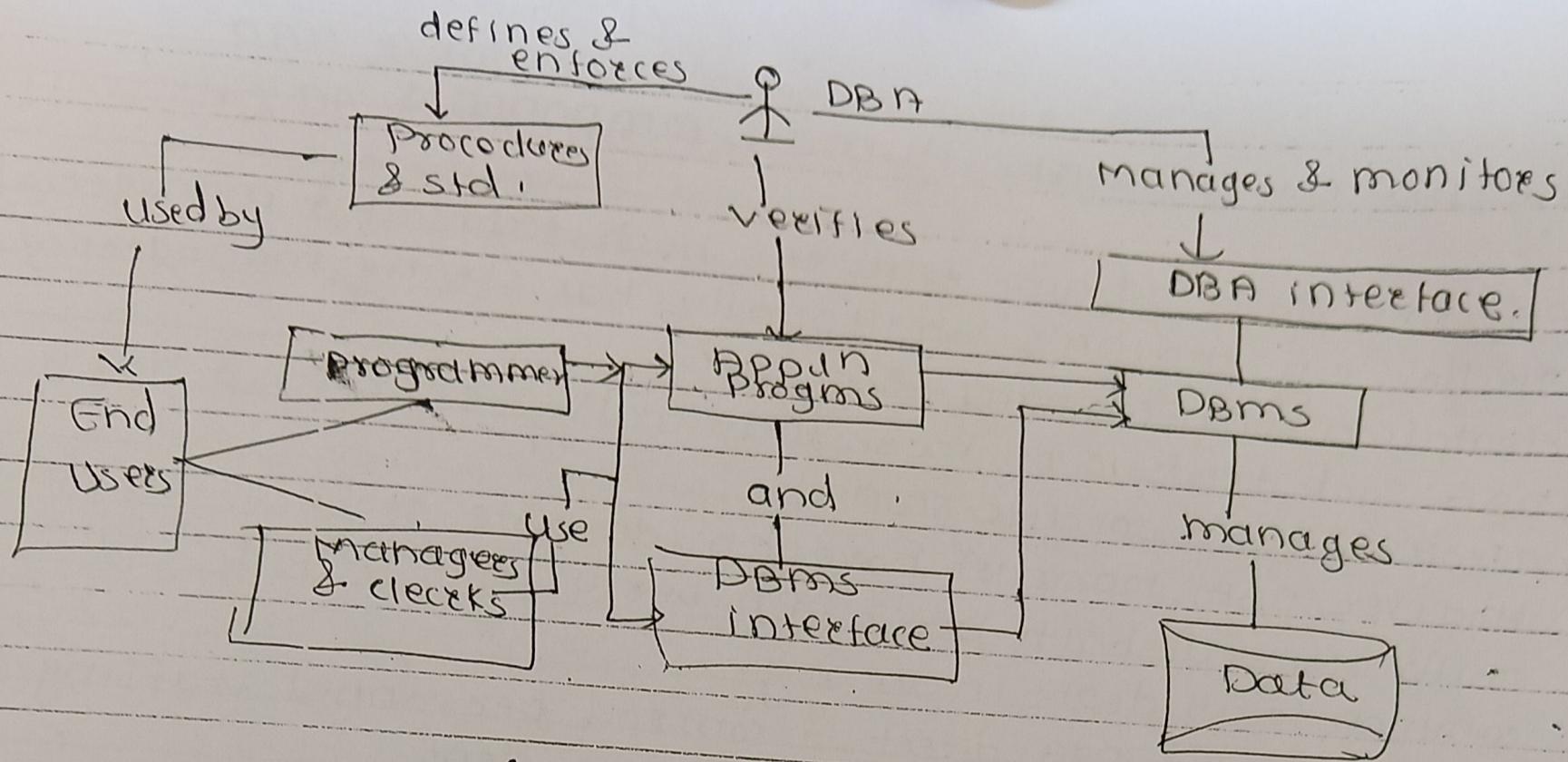
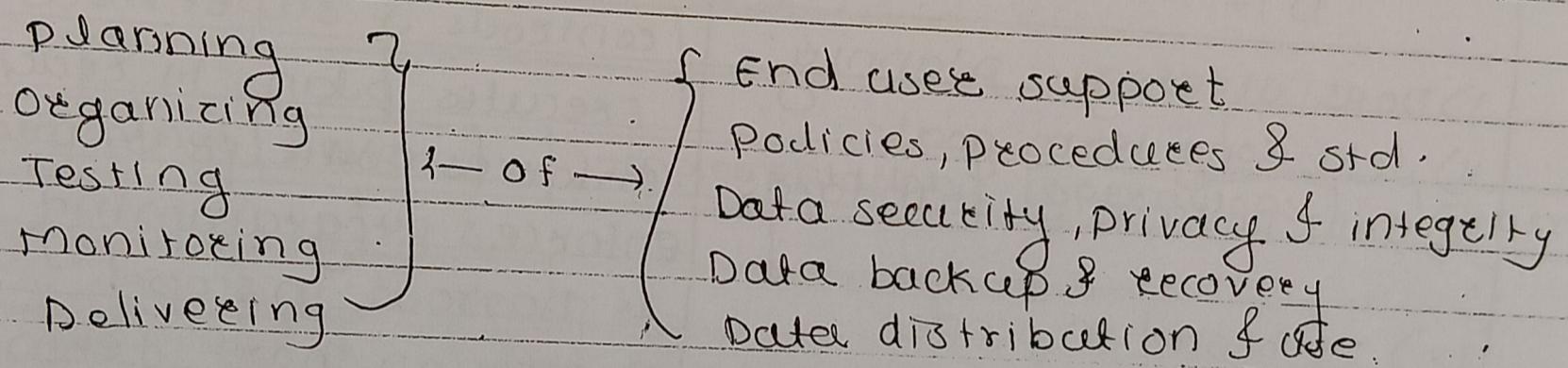


fig. DBA activities

* DBA activities & services



1) End user support

DBA interact with the end user by providing data & info. support services to the org. dept. because end users usually have dissimilar computer backgrounds, end user support services include.

- i) gathering user requirements
- ii) building end user confidence
- iii) resolving conflicts & problems
- iv) finding sol'n to info. needs
- v) ensuring quality & integrity of data & appl'
- vi) managing & training & support DBMS users

2) Policies, procedures & standards

- i) Policies are general statements of direction or action that communicate & support DBA goals

ii) Standards → It describes minimum requirements of given DBA activity. They are more detailed & specific than policies.

e.g. Standards define the structure of appn programs & naming conventions programmers must use.

iii) Procedures

It written instructions that describe series of steps to be followed during performance of given activity

3) Data security, privacy & integrity

The security, privacy & integrity of data in DB are great concern to DBAs who manage current DBMS installation

4) Data backup & Recovery

- when data are not readily available, companies face potentially serious losses, therefore, data backup & recovery procedures are critical in all DB installation
- The DBA must also ensure the data in DB can be fully recovered in case of physical data loss or DB integrity

5) Data distribution & Use

- i) Data are useful only when they reach the right user in timely fashion
- ii) The DBA is responsible for ensuring that the data are distributed to the right people at right time & in right format

* Database administration Tools:

The data dictionary as a data administration tool as well as the DBA's use of computer Aided software engineering (CASE) tools to support database analysis & design.

2) Data Dictionary

- i) The DBMS data dictionary provides the DBMS with its self describing characteristics. In effect data dictionary resembles an X-ray of the company's entire data set & it is crucial element in data administration.
- ii) Two main types of data dictionaries
 - a) Integrated
 - b) Standalone
- iii) An integrated data dictionary is included with DBMS. e.g. all relational DBMS include built in data dictionary or system catalog i.e. frequently accessed & updated by RDBMS.
- iv) Other DBMS other type do not have built in data dictionary instead the DBA may use 3rd party standalone data dictionary system.

Data dictionaries can also classified as active or passive

→ Active data dictionary

It is automatically updated by the DBMS with every DB access, thereby keeping its access info. up to date.

2) Passive data dictionary

It is not updated automatically & usually requires running a batch process.

There is no standard format for the info. stored in the data dictionary, several common features are

- 1) Data elements that are defined in all tables of all database. It stores the names, data types, display format, internal storage format & validation rules.

- 2) Tables defined in all databases, it stores the name of table creator, date of creation, access authorisation & no. of columns.
- 3) Indexes defined for each database, table, for each index DBMS stores at least index name, attribute used, location, specific authorisations & no. of columns, index characteristics & creation date.
- 4) Defined DB. This include who created each DB, when & where DB is located, who the DBA.
- 5) End users & admin. of the DB
- 6) Programs that access the DB. This include screen formats, report formats, appn programs & SQL queries
- 7) Access authorisation of for all users of all DB
- 8) Relationship among data elements.

Examples of Data Dictionary Usage

- List the names & creation dates of all table created by user swati in the current DB.

```
→
SELECT NAME, CTIME
FROM SYSTABLES
WHERE CREATOR = 'SWATI';
```

(2) Case Tools

- i) CASE is the acronym for Computer Aided Systems Engg.
- ii) It provides an automated framework for the system development Life cycle (SDLC).
- iii) CASE uses structured methodologies & great powerful graphical interfaces
- iv) It play an increasingly important role in Info. system development.
- v) CASE tools are usually classified according to the extent of support they provide for SDLC
e.g. → front end CASE tools that provide support for the planning & design phases, back end CASE tools provide support for the coding & implementation phases

CASE Tools include

- i) Reduction in development time & costs
- ii) Automation of the SDLC
- iii) Standardization of systems development methodology
- iv) Easier maintenance of appln systems developed with CASE tools.

- One of the CASE tools most important components is an extensive data dictionary which keeps track of all objects created by the systems designer.
- e.g., CASE data dictionary stores data flow diagrams, structure charts, descriptions of all external & internal entities, data store, data items, report formats & screen formats.

A typical CASE tool provides 5 components

- 1) Graphics designed to produce structured diagrams such as data flow diagrams, ER diagrams, class diagrams & object diagrams.
- 2) Screen painters & report generators to produce the info. system I/P & O/P format.
- 3) An integrated repository for storing & cross referencing the system design data. This repository includes a comprehensive data dictionary.
- 4) An analysis segment to provide a fully automated checks on system consistency, syntax & completeness.
- 5) A program documentation generator.

CASE Tools

casewise	corporate modeler suite
compute & associates	ERwin
Embarcadero Technologies	ER/ studio
microsoft	VISIO
oracle	Designer
IBM	System Architect
sybase	Power Designer
visible	visible Analyst

* Developing a Data Administration Strategy.

- i) The database administration strategy must not conflict with IS plans. After all, the IS plans are derived from detailed analysis of company's goals, it's condition of situation & business needs.
 - ii) Several methodologies are available to ensure the compatibility of data admin & IS plans & to guide the strategic plan development.
 - iii) Information Engg. (IE) allows for the translation of company's strategic goals into data & app-in that will help company achieve those goals.
 - iv) The output of the IE process is an IS architecture (ISA) that serves as the basis for planning, development & control of future IS.

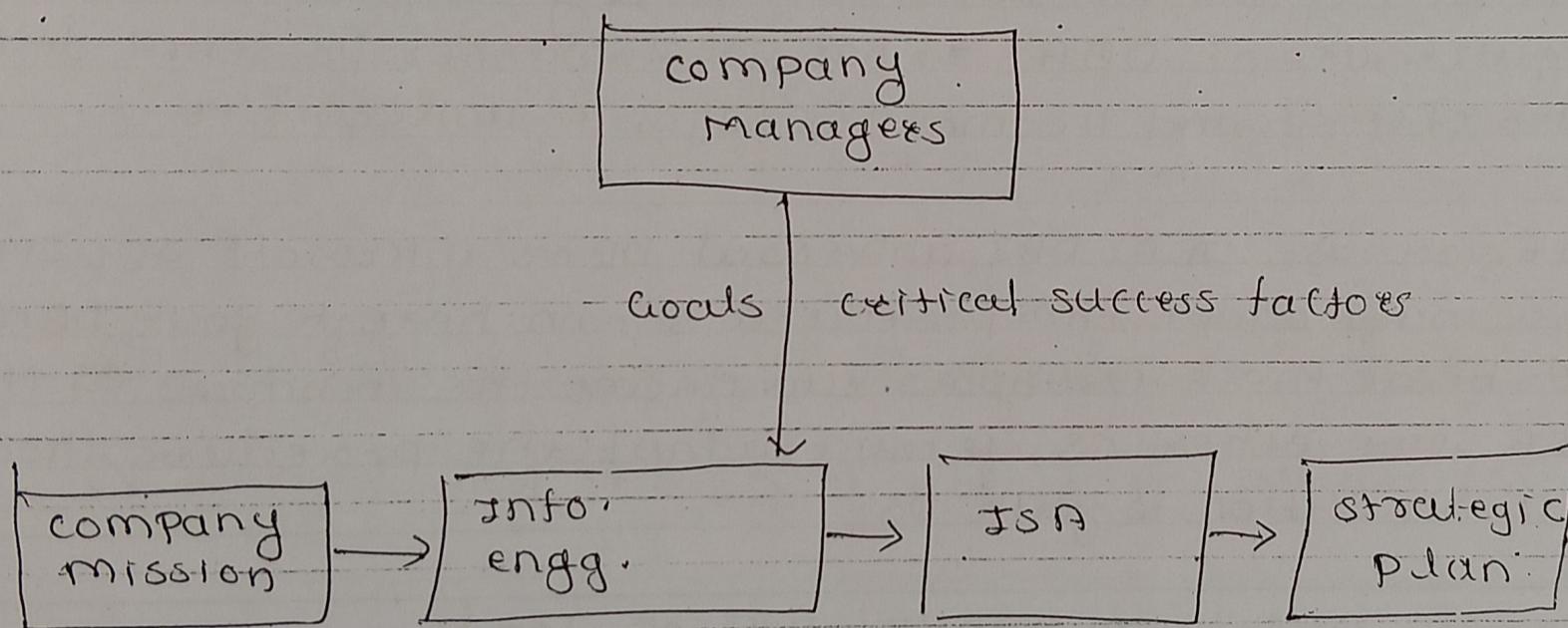


Fig. Forces affecting development of the ISA

- Implementing IE methodologies is an org. is a costly process that involves planning, commitment of resources, mgmt.ability, well defined objectives, identification of critical factors & control.

- ISIA provides framework that includes the use of computerized, automated & integrated tools such as DBMS & CASE tools.

The DBA at work: Using Oracle for DB administration
 How DBA might handle the following technical tasks in specific DBMS

- 1) Creating & expanding database storage structures
 - 2) Managing DB object such as tables, indexes, triggers & procedures.
 - 3) Managing end user DB env. including type & extent of DB access.
 - 4) Customizing DB initialization parameters
- All DBMS vendors provide set of programs to interface with DB & to perform wide range of DB adminstrn tasks
- We choose oracle 11g for windows to illustrate the selected DBA tasks bcoz it is typically found in org. that are sufficiently large & have sufficiently complex db env. to require use of DBA. It has good market presence & it is also often found in small colleges & universities
- If you use IBM DB2 universal DB or microsoft sql server, you must adapt the procedures shown here to your DBMS & bcoz these examples run under the windows OS if you use some other OS, u must adapt the procedures shown in this section to your OS.

Oracle Database Administration Tools

- 1) All DB vendors supply a set of DBA tools. In oracle you perform most DBA tasks via oracle enterprise manager interface
- 2) The default login. → To perform any administrative task you must connect to the DB using username with DBA privileges. By default oracle automatically creates SYSTEM & SYS user id's that have administrative privileges with every new DB you create.
- 3) Ensuring an automatic RDBMS start
 One of the basic DBA tasks is to ensure that your DB access is automatically started when you turn on the computer

- A database is logically composed of one or more tablespaces.
- A tablespace is logical storage space. & are used primarily to group related data logically.

- ⑤ - Oracle automatically creates the tablespaces & datafiles
- i) The SYSTEM tablespace is used to store data dictionary data
 - 2) USERS → is used to store table data created by end users
 - 3) TEMP → to store temporary tables → indexes created during the execution of SQL statements
 - 4) UNDOTBS → - used to store DB transaction recovery info.

⑥ Managing the DB objects, tables, views, triggers & procedures

- Oracle Enterprise Manager gives the DBA graphical user interface to create, edit, view & delete DB objects in DB.

⑦ Managing Users & Establishing security

- ⑧ Customizing the DB initialization parameters

5. Business Intelligence & Data Warehouse

ability to apply acquired knowledge

The Need for Data Analysis

- Managers must be able to track daily transactions to evaluate how the business is performing.
- Strategies should be developed to meet organizational goals using operational databases.
- Data analysis provides info. about short term tactical evaluations & strategies.

Business Intelligence

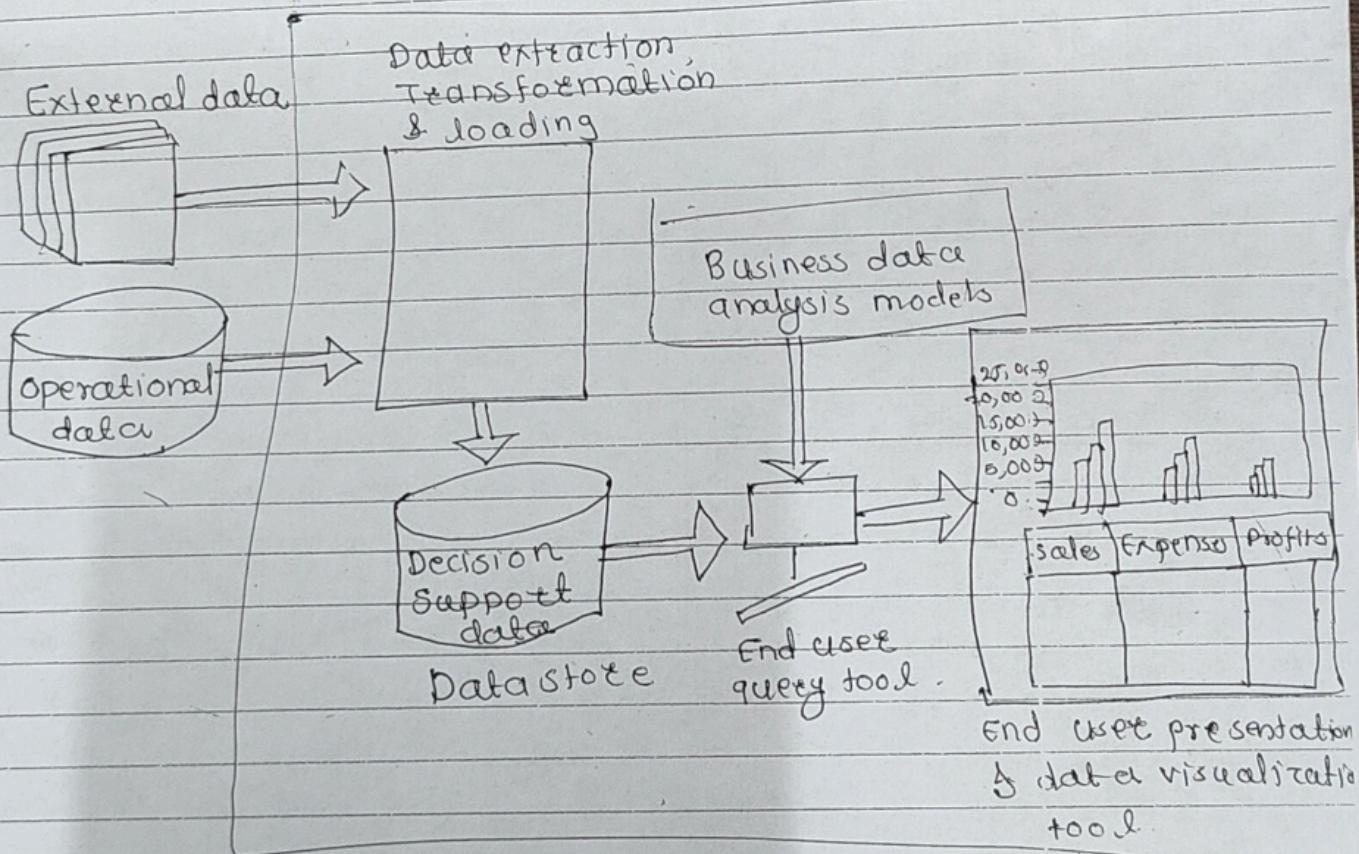
Definition

Comprehensive, cohesive, integrated tools & processes to capture, collect, integrate, store & analyze data & generate info. to support business decision making.

- It includes BI tools & tech to transform raw data into meaningful & actionable info.
- Framework that allows a business to transform
 - Data into info.
 - Info into knowledge
 - Knowledge into wisdom
- Implementing BI captures not only business data (internal & external) but also metadata.
- BI involves the following general steps.
 - 1) collecting & storing operational data
 - 2) Aggregating the operational data into decision support data
 - 3) Analyzing decision support data to generate information
 - 4) Presenting such info. to the end user to support business decision
 - 5) Making business decisions which in turn generate information more data is collected, stored etc. (restarting the process)

c) Monitoring result to evaluate outcomes of the business decisions (Providing data to be collected, stored etc)

- i) composed of data, people, processes, technology & management of components.
- ii) Focuses on strategic & tactical use of information
- iii) Multiple tools from different vendors can be integrated into a single BI framework.
- iv) KPI (Key Performance indicators) measurements that assess company's effectiveness & success in reaching goals.
- v) Master Data management
A collection of concepts, techniques & processes for the proper identification, definition & management of data elements within an organization



① Decision support data

i) Operational data

- mostly stored in relational database
- optimized to support transactions representing daily operations.

ii) Decision support data

It differs from operational data in three main areas:

- Time span
- Granularity
 - drill down & roll up to different levels of aggregation
- Dimensionality

iii) Operational data have narrow timespan, low granularity & single focus, such data are usually presented in tabular format in which each row represents a single transaction. This format often makes it difficult to derive useful information.

iv) DSS data focus on broader timespan tend to have high level of granularity & can be examined in multiple dimensions, for example

- sales by product, region, agent etc.
- sales for all years or only few selected years.
- sales for all products or only few selected products.

characteristics

operation Data

DSS Data

Data currency current operations
Real time data

Historic data, snapshot
of company data, Time
component (week / month/
years)

Granularity

Atomic detailed
data

summarized data

summation level

Low, some aggregate
yields

High, many aggregate
levels.

Data model

Highly normalized
mostly relational DBMS

Non-normalized,
complex structures
some relational but
mostly multidimensional
DBMS

Transaction Type

Mostly updates

Mostly query

Transaction volumes

High update volumes

periodic loads &
summary calculations

Transaction speed

Updates are critical

Retrievals are
critical

Query activity

Low to medium

High

Query scope

Narrow range

Broad range

Query complexity

Simple to medium

Very complex

Data volumes

Hundreds of megabytes
up to gigabytes

Hundreds of gigabytes
up to terabytes.

Decision support DB requirements

- Specialized DBMS tailored to provide fast answers to complex queries.

i) Fox main requirement

- Database schema
- Data extraction & loading
- End user analytical interface
- Database size

ii) Database schema

- Must support complex data representations.
- Bitmap indexes, data partitioning, non-normalized
- Must contain aggregated & summarized data
- Queries must be able to extract multidimensional time slices.

iii) Data extraction & filtering

- Should allow batch & scheduled data extraction
- Support different data sources.
- Flat files, Hierarchical, network & relational databases, multiple vendors.

iv) Data filtering

- Must allow checking for inconsistent data.
- Advanced data integration, aggregation & classification.
- Must solve data-formatting conflicts.

v) End user analytical interface

- One of most critical DSS DBMS components
- Permits user to navigate through data to simplify & accelerate decision making process.

vi) Database size

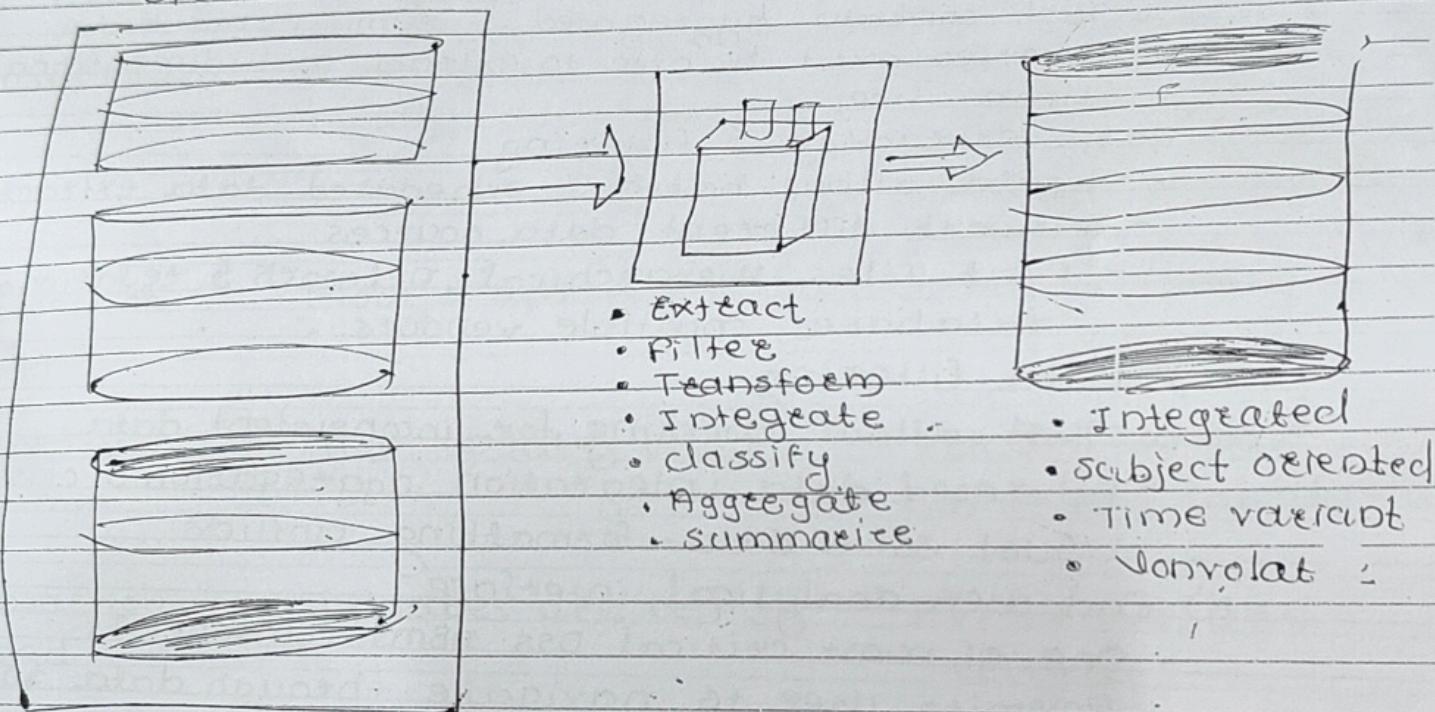
- DBMS must support very large databases e.g.
- Might be required to use advanced hardware such as disk arrays, symmetric multiprocessors or massively parallel processors.

Data warehouse

- Integrated, subject oriented, time variant & nonvolatile collection of data
- It provides support for decision making
- Usually a read only database optimized for data analysis & query processing.
- Requires time money & considerable managerial effort to create.

Creating data warehouse

Operational data



Data mart

- Small, single subject, data warehouse subset
- More manageable data set than data warehouse
- Provides decision support to small group of people
- Typically lower cost & lower implementation time than data warehouse

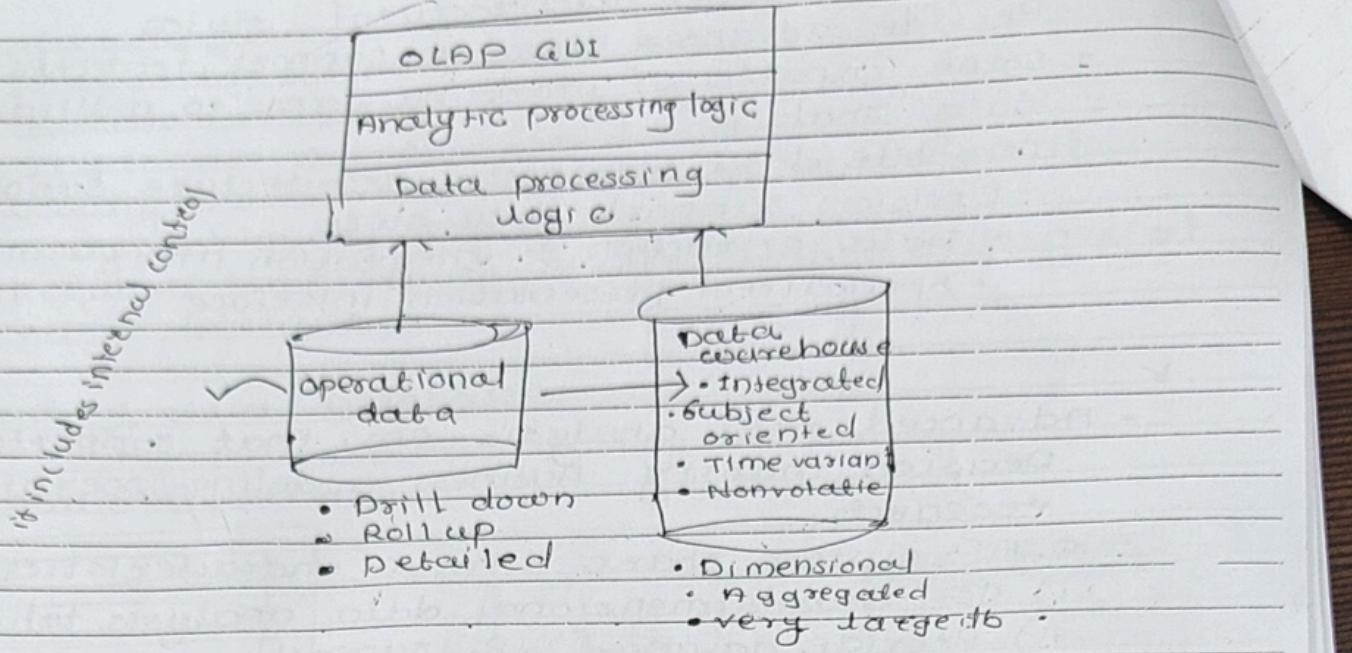
Decision support architectural styles

- Provide advanced decision support features.
- Some capable of providing access to multidimensional data analysis
- Complete data warehouse architecture supports
 - Decision support data store
 - Data extraction & integration filter
 - Specialized presentation interface

- X
- Advanced data analysis env. that supports decision making, business modeling, operations research.
 - OLAP system share 4 main characteristics.
 - i) Use multidimensional data analysis tech.
 - ii) Provide advanced DB support
 - iii) Provide easy to use end user interface
 - iv) Support client/server architecture
 - Easy to use end user interface
 - Advanced OLAP features more useful when access is simple
 - Many interface features are borrowed from previous generations of data analysis tools.

OLAP Architecture

- Operational characteristics 3 main modules.
 - Graphical User Interface (GUI)
 - Analytical processing logic
 - Data processing logic
- Designed to use both operational & data warehouse data.
- In most implementations, data warehouse & OLAP are interrelated & complementary
- OLAP systems merge data warehouse & datamart approaches.



1) multidimensional

- i) In multidimensional analysis, data are processed & viewed as part of multidimensional structure.
- ii) This type of analysis is particularly attractive to business decision makers because they tend to view business data as data that are related to other business data.
- iii) multidimensional data analysis tech. are augmented by the following fun.

a) Advanced data presentation functions.

- 3-D graphics, pivot tables, crosstabs, data rotation & 3 dimensional cubes. Such facilities are compatible with desktop spreadsheets, statistical packages & query & report packages.

b) Advanced data aggregation, consolidation & classification fun.

→ These allow the data analyst to create multiple data aggregation levels, slice & dice data. f drill down & rollup data across different dimension & aggregation levels.

e.g. → Aggregating data across the time dimensions by week, month, quarter & year allow data analyst to drill down & rollup across time dimensions.

c) Advanced computational functions.

- These include business oriented variable, financial & accounting ratio & statistical & forecasting function. These functions are automatically & end user does not need to redefine their components each time they are used.

d) Advanced data modeling functions.

These provide support for what if scenarios, variable assessment, variable contributions to outcome, linear programming & other modelling tools.

② Advanced Database support.

- i) Access to many different kinds of DBMS, flat files & internal & external data sources.
- ii) Access to aggregated data warehouse data as well as to detail data found in operational databases.
- iii) Advanced data navigation features such as drill down & roll up.
- iv) Rapid & consistent query response times.
- v) The ability to map end user requests, expressed in either business or model terms, to the appropriate data source & then to the proper data access language.
- vi) It support for very large databases.

③ Easy to use End user Interface.

- i) OLAP tool vendors learned this lesson early & have equipped their sophisticated data extraction & analysis tools with easy to use graphical interface.
- ii) Many of the interface features are "borrowed" from previous generation of data analysis tools that are already familiar to end user. This familiarity makes OLAP easily accepted & readily used.

OLAP Architecture

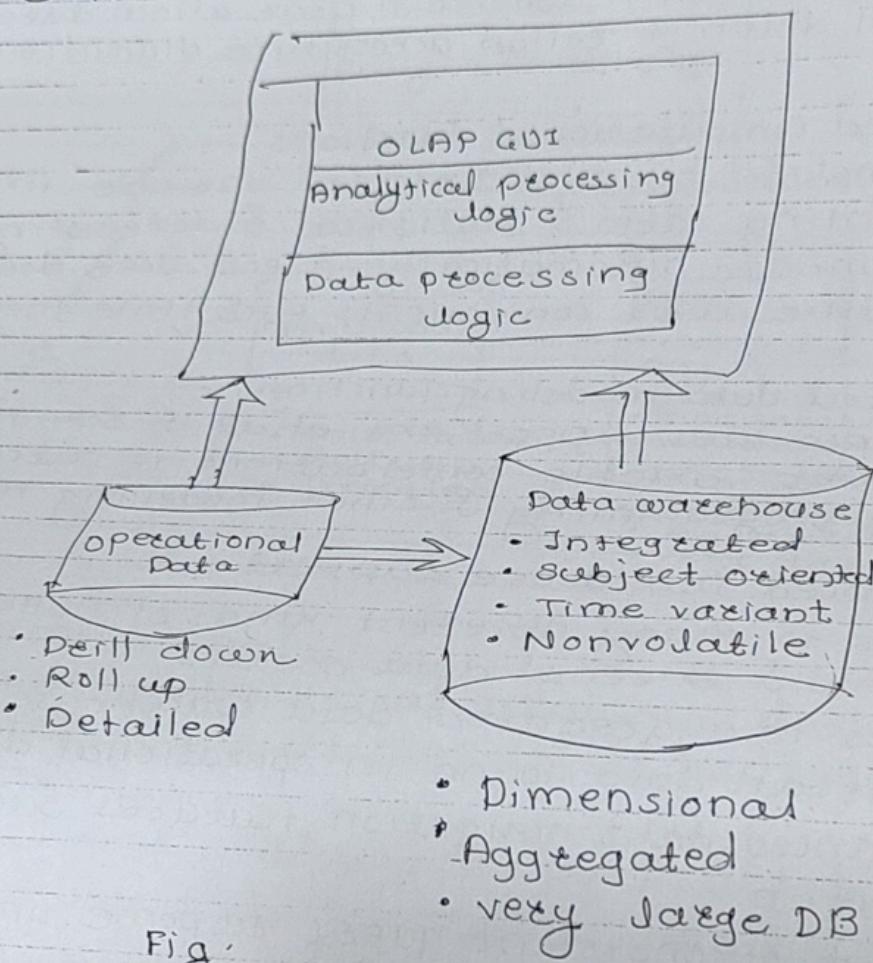


Fig.

OLAP can be divided into three main modules

- i) GUI (Graphical User Interface)
- ii) Analytical processing logic
- iii) Data processing logic

- * OLAP systems are designed to use both operational & data warehouse data.
- * Fig. shows OLAP system components located on single computer but this single user scenario is only one of many.

Data warehousing

- i) Data warehouse is the foremost repository for the data available for developing BI architectures & decision support systems.
- ii) The data warehousing indicates the whole set of integrated activities involved in designing, implementing & using a data warehouse.
- iii) It is possible to identify 3 main categories of data feeding into a data warehouse.
 - a) Internal data.
 - b) External data
 - c) Personal data

a) internal data

- i) Internal data are stored for the most part in db., referred to as transactional systems or operational systems, that are backbone of an enterprise info. system.
- ii) Internal data are gathered through transactional appn that routinely preside over the opers of a company such as administration, accounting, production & logistics.
- iii) This collection of transactional db. appn is termed enterprise resource planning (ERP).
- iv) The data stored in the operational systems usually deal with main entities involved in a company processes, namely customers, products, sales, employees & suppliers. These data usually come from different components of the info. system.

d) back office system

That collect basic transactional records such as orders, invoice, inventories, production & logistic data.

b) front office systems

That contain data originating from call center activities, customer assistance, execution of marketing campaigns.

c) web based systems

That gather sales transactions on e-commerce sites, visits to websites, data available on forms filled out by existing & prospective customers.

2) External data

- i) There are several sources of external data that may be used to extend the wealth of info. stored in the internal DB.
- ii) e.g. - some agencies gather & make available data relative to sales, market share & future trend predictions for specific business industries as well as economic & financial indicators.
- iii) Source of external data is provided by geographic information system (GIS) which represent a set of appn for acquiring organizing & storing & presenting terriro-
cal data.

3) Personal data

In most cases, decision makers performing BI analysis also rely on info. & personal assessments stored inside worksheets or local db located in their computers. The retrieval of such info. & its integration with structured data from internal & external sources is one of the objectives of knowledge mgmt. system.

* Data Warehouse architecture

Figure includes the following major components:

- 1) The data warehouse itself, together with additional data marts, that contains the data & fun. that allow the data to be accessed, visualized & perhaps modified.
- 2) Data acquisition appn also known as extract, transform & load (ETL) or back end tools which allow the data to be extracted, transformed & loaded into the data warehouse.

- 3) Business intelligence & decision support appn. which represent the front end & allow the knowledge workers to carry out the analyses & visualize the results.

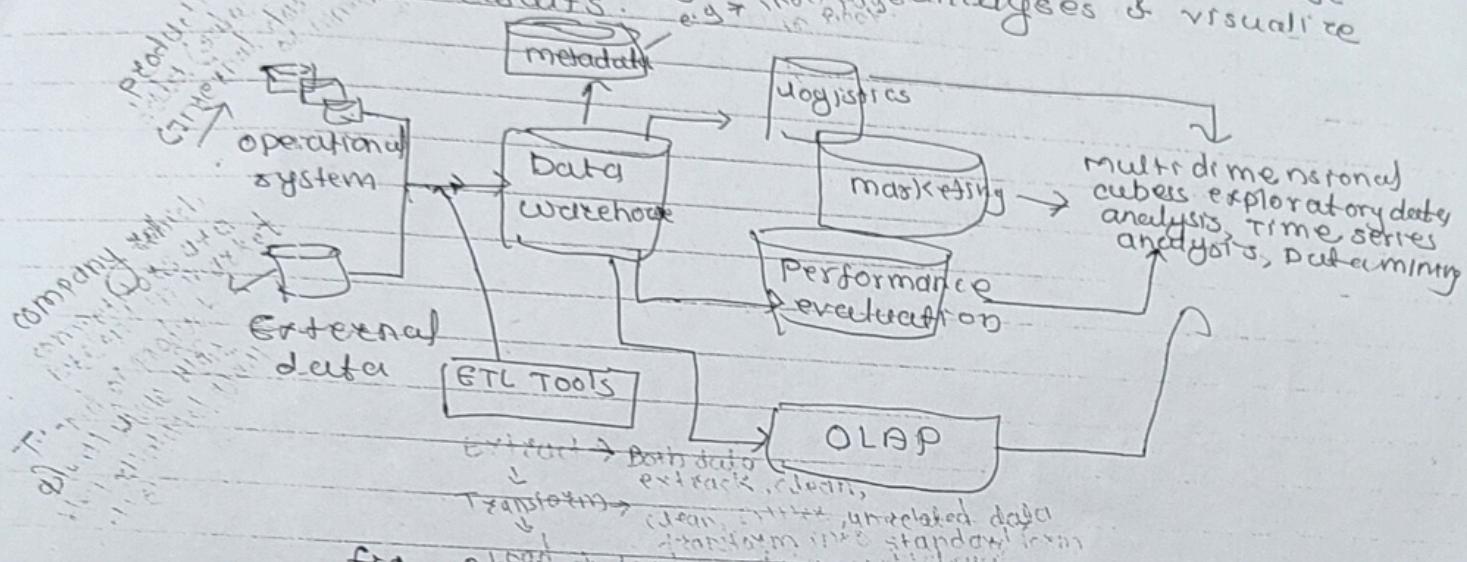


fig. Architecture & Fun. of data warehouse

These level distinction applies to the architecture shown in above fig.

- 1) The level of the data sources & related ETL tools that are usually installed on one or more servers.
- 2) The level of data warehouse & any data mart possibly available on one or more servers as well & separated from those containing the data sources. This second level also includes the metadata documenting the origin & meaning of the records stored in the data warehouse.
- 3) The level of the analyses that increase the value of info. contained in data warehouse through query, reporting & possibly sophisticated decision support tools.

The data warehouse may be implemented according to different design approaches.

- a) Top down → It is based on the overall design of the data warehouse & is therefore more systematic.
- b) Bottom up → It is based on the use of prototypes, therefore system extensions are made according to step by step.

* ETL Tools

3) Extraction

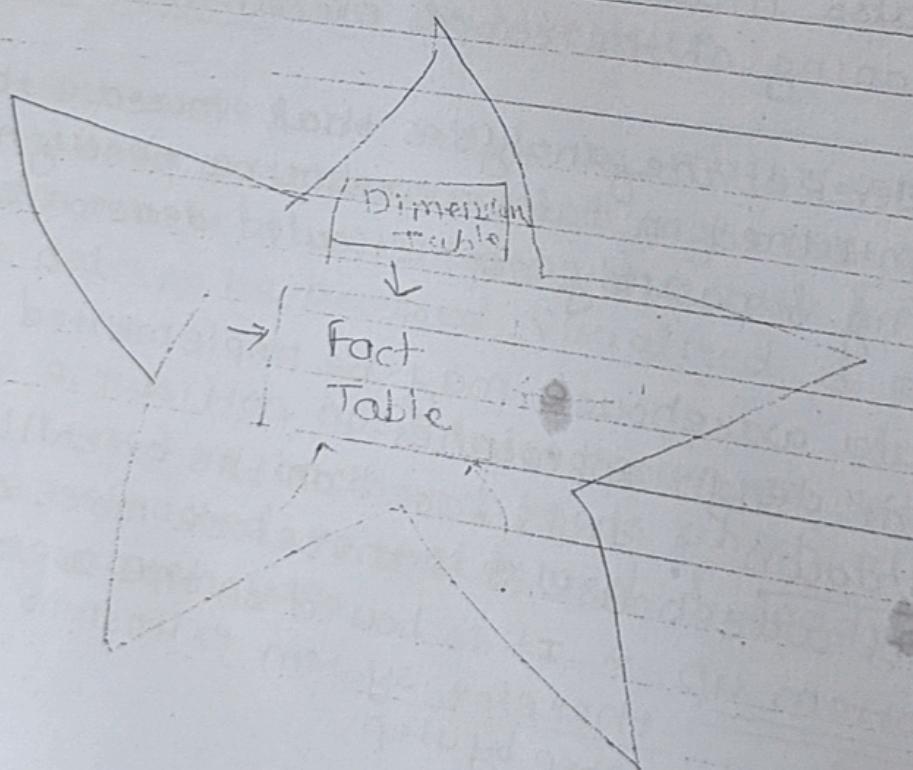
During the 1st phase, data are extracted from the various internal & external sources.

2) Transform

After extracting the data, then clean & filter. data means unrelated data can be removed & make standard format of data.

3) Load

store & load into data warehouse



Star schema

- i) star schema is the fundamental schema along among the data mart schema & it is simplest.
- ii) This schema is widely used to develop or build a data warehouse & dimensional data marts.
- iii) It includes one or more fact tables indexing any no. of dimensional tables.
- iv) It is necessary cause of snowflake schema.
- v) In star schema, business process data, that holds the quantitative data about business is distributed in fact tables & dimensions which are descriptive characteristics related to fact data.
- vi) Sales price, sale quantity, distant, speed, weight, & weight measurements are few e.g. of fact data in star schema.
- vii) If having multiple dimensions is termed as centipede schema.

Advantages

- 1) simple queries
- 2) simplified Business Reporting logic
- 3) Feeding cubes

* SQL extension for OLAP

- i) OLAP requires almost invariably data aggregations & SQL does support such aggregations through its Group-by instruction.
- ii) However during OLAP, each data analysis session usually requires repeated aggregations of similar type over the same data.
- iii) To facilitate the execution of OLAP queries & data aggregation, SQL introduced three extensions to the GROUP BY, CUBE & ROLLUP.

- GROUP BY. son
- iv) The CUBE operator computes a union of GROUP BY on every subset of the specified attribute types.
Its result set represents a multidimensional cube based upon the source table.
 - v) The grouping () function returns 1 in case NULL value is generated during the aggregation & 0 otherwise. This distinguishes generated NULLs & possible real NULLs stemming from the data.
 - vi) - The ROLLUP operator computes the union on every prefix of the list of specified attribute types from most detailed up to the grand total.
 - It is especially useful to generate reports containing both subtotals & totals.
 - The key difference b/w the ROLL UP & CUBE operator is that the former generates a result set showing the aggregates for hierarchy of values of the specified attribute types.

* Materialized view

- i) When the results of view expression are stored in DB system they are called materialized views.
 - ii) SQL does not provide any standard way of defining materialized view, however some DBMS provides custom extension to use materialized views.
 - iii) The process of keeping the materialized views updated is known as view maintenance.
- * DB system uses one of the 3 way to keep the materialized view updated.
- 1) Update the materialized view as soon as the relation on which it is defined is updated.
 - 2) Update the materialized view every time the view is accessed.

3) Update the materialized view periodically.

- Materialized view is useful when view is accessed frequently as it saves the computation time as the result are stored in DB before hand.

6: Data Analysis & Exploration

* Data Mining

- i) Data mining refers to the overall process consisting of data gathering & analysis, development of inductive learning models & adoption of practical models & methods, decision & consequent actions based on the knowledge acquired.
- ii) Data mining process is based on inductive learning methods, whose main purpose is to derive general rules starting from a set of available examples consisting of past observations recorded in one or more DB.
- iii) The purpose of data mining analysis is to draw some conclusions starting from a sample of past observations & to generalise these conclusions with reference to the entire population in such way that they are as accurate as possible.
- iv) A further characteristic of data mining depends on the procedure for collecting past observations & inserting them into a DB. Indeed these records are usually stored for purposes that are not primarily driven by data mining analysis.
- v) Data mining activities can be subdivided into two major investigation streams according to the main purpose of the analysis:
 - 1) Interpretation
 - 2) Prediction

Interpretation

- i) The purpose of interpretation is to identify regular patterns in the data & to express them through rules & criteria that can be easily understood by experts in the appn domain.
- ii) The rules generated must be original & non-trivial in order to actually increase the level of knowledge & understanding of the system of interest.

i) The purpose of prediction is to anticipate the value that random variable will assume in the future or to estimate the likelihood of future events.

e.g. → mobile phone provider may develop a data mining analysis to estimate for its customers the probability of churning in favor of some competitor.

ii) In different context, retail company might predict the sales of given product during subsequent weeks.
iii) Most data mining tech. derive their predictions from the value of set of variables associated with the entities in a database.

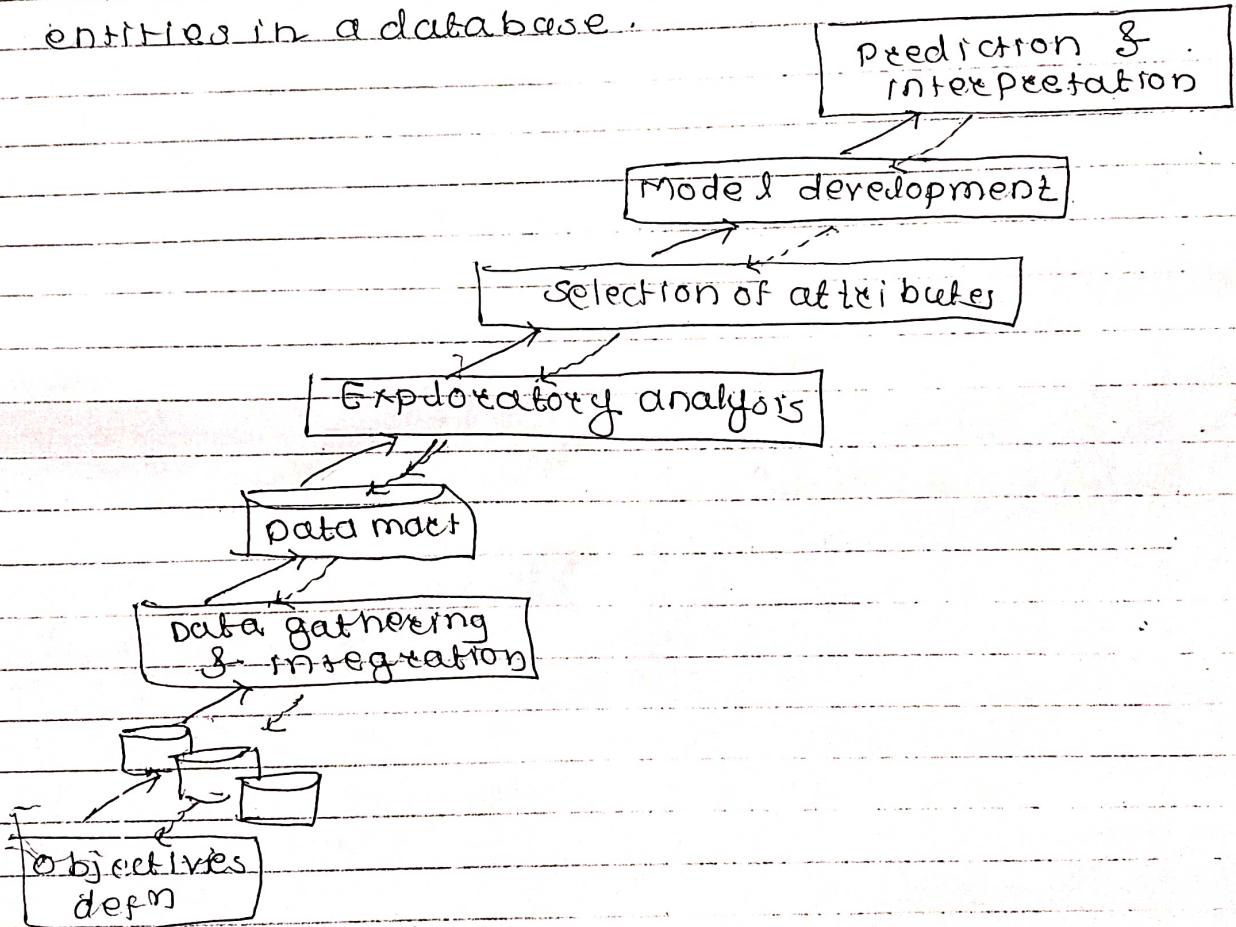


fig. Data mining process.

- * Models & methods for data mining
 - i) there are several learning methods that are available to perform a different data mining tasks.
 - ii) A no. of techniques originated in the field of computer science such as classification trees or association rules & are referred to as machine learning or knowledge discovery in db.
 - iii) In most cases an empirically based approach tends to prevail within this class of techniques & other methods belong to multivariate statistics such as regression or Bayesian classifiers & are often parametric in nature but appear more theoretically grounded.
 - iv) statistical learning theory which are based on solid theoretical foundations & place themselves at the crossroads of various disciplines among which probability theory, optimization theory & statistics.

* Applications of data mining

1) Relational marketing

- identification of customer segments that are most likely to respond to targeted marketing campaigns such as cross selling & up selling.
- identification of target customer segments for retention campaigns.
- prediction of the rate of responses to marketing campaigns.
- interpretation & understanding of the buying behavior of the customers.
- analysis of the products jointly purchased by customers known as market basket analysis.

2) Fraud detection

Fraud may affect different industries such as telephony, insurance & banking.

3) Risk evaluation

The purpose of risk analysis is to estimate the risks connected with future decisions, which often assume to dichotomous form.

e.g. - using the past observations available, bank may develop a predictive model to establish if it is appropriate to grant a monetary loan at home loan based on characteristics of the applicant.

4) Text mining

Data mining can be applied to different kinds of texts which represent unstructured data, in order to classify articles, books, documents, emails & web pages.

E.g. - web search engines or the automatic classification of press releases for storing purpose.

5) Image Recognition

The treatment & classification of digital images both static & dynamic is an exciting subject both for its theoretical interest & great no. of applications it offers.

6) Web mining

They may prove useful for the analysis of e-commerce sites in offering flexible & customized pages to suffice in caching most popular pages or in evaluating the effectiveness of an e-learning training course.

7) Medical diagnosis

- Learning models are an invaluable tool within the medical field for the early detection of diseases using clinical test results.

- Image analysis for diagnostic purpose is another field of investigation i.e currently burgeoning:

Data Preparation

- i) Data preparation is used for proper business data analysis.
- ii) The data preparation process involves collecting, cleaning & consolidating data into a file that can be further used for analysis.

why data preparation is necessary

- i) To filter unstructured, inconsistent & disordered data
- ii) connecting data from real time multiple data sources
- iii) For quick reporting of data.
- iv) To handle data collected from scraped file like Pdf doc.

steps of data preparation

Gather

|
Discover

|
Cleanse

|
Transform

|
Develop

|
Store

① Gather data

This is initial process for each business, in this phase it is necessary to collect data from various sources. The sources can be of any type such as from catalogs or ad hoc can be added.

2) Discover data

The next step is discovering the data, here it is imp. to understand the data & categorize it into different dataset. This step might take long time to filter because of the huge coll'n of datasets.

3) Cleaning & validating data

- i) Removing unnecessary data & outliers
- ii) Use the appropriate patterns for refining all the data.
- iii) Use the lock to protect your sensitive data.
- iv) Fill the empty space for data flow.

4) Transforming the data

It defines maintaining the format of value entered in order to meet well define O.P. & can clearly understand the wider audience.

5) Storing data

This is final step after going through all the above processes. Once data is cleaned it is ready to offer 3rd party tools such as BI tools for analysis.

Data exploration / Exploratory data analysis (EDA)

- i) It provides simple set of exploration tools that bring out the basic understanding of real time data into data analytics.
- ii) The outcomes of data exploration can be powerful factor in understanding the structure of data, values distributions & interrelationships.
- iii) It is the 1st data analytic.

Data exploration Method

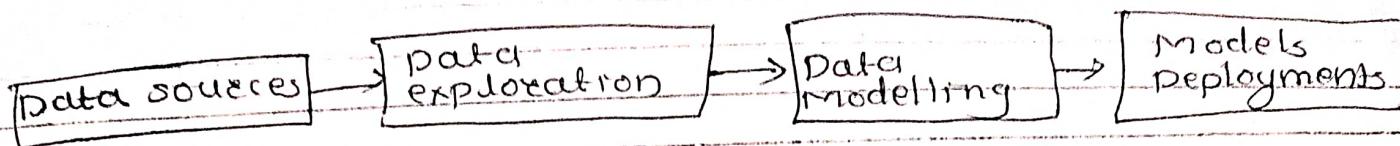
There are two formats of data exploration.

- 1) automatic
- 2) manual

Mostly analysts preferred automated methods such as data visualization tools because of their accuracy & quick response.

Manual data exploration on other hand, methods include filtering & drilling down into data in Excel spreadsheets & writing scripts to analyze raw data sets.

stages of data mining



There are several tech. for analyzing data such as

1) Univariate analysis

It is the simplest form of analyzing data.

Univariate means that there is only one variable in your data.

2) Bivariate analysis

It is simplest form of quantitative analysis. It includes analysis of two variables. It includes analysis of 2 variables (x, y) used for calculating the empirical relation between two variables.

3) Principal components analysis

The analysis & conversion of possibly correlated variables into smaller no. of uncorrelated var.

4) Multivariate analysis

It can be used to refer to any analysis that involve more than one variable.

The next step after data exploration is data discovery. In this phase BI tools are used to inspect trends, sequences & events & create visualization to present to business managers.