

## UNIT-I :-

DBMS:- A s/w package / system to facilitate the creation and maintenance of a computerized database.

### Database Systems:

The DBMS s/w together with the data itself. Sometimes the applications are also included.

### DBMS Functionality:

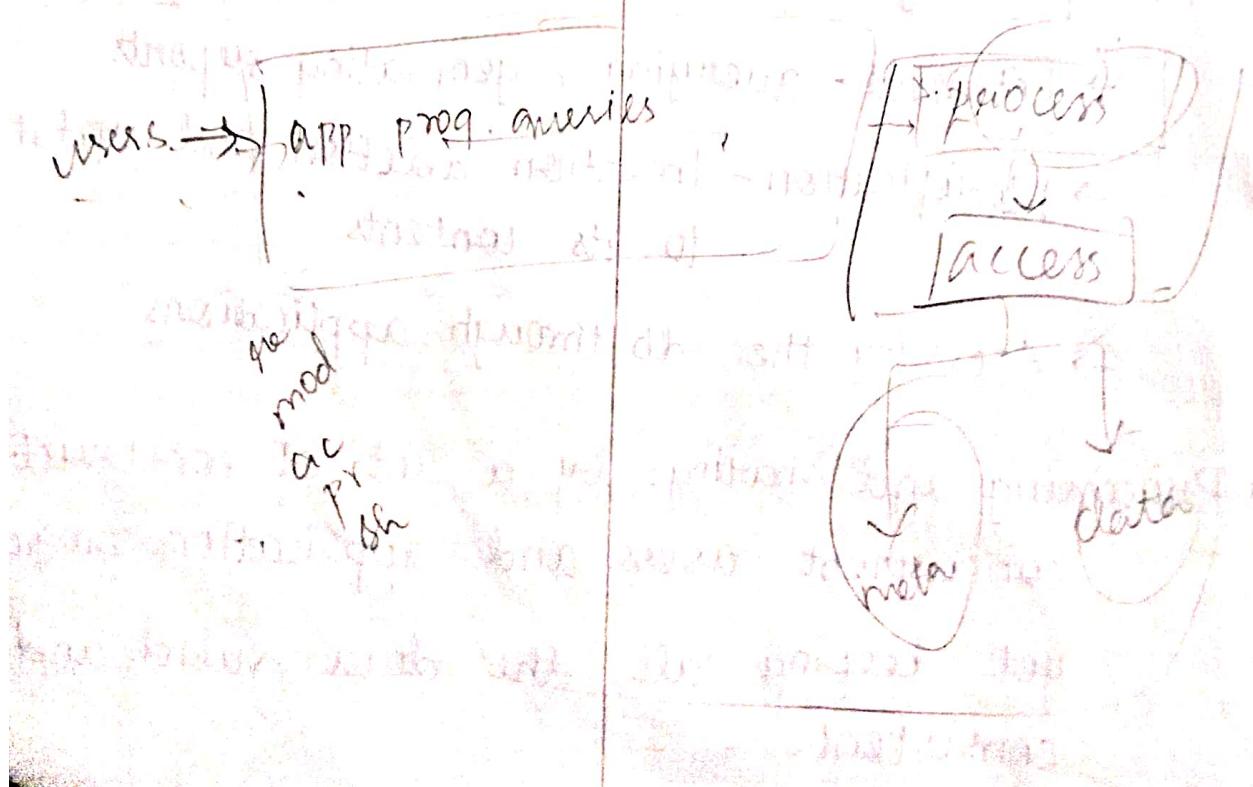
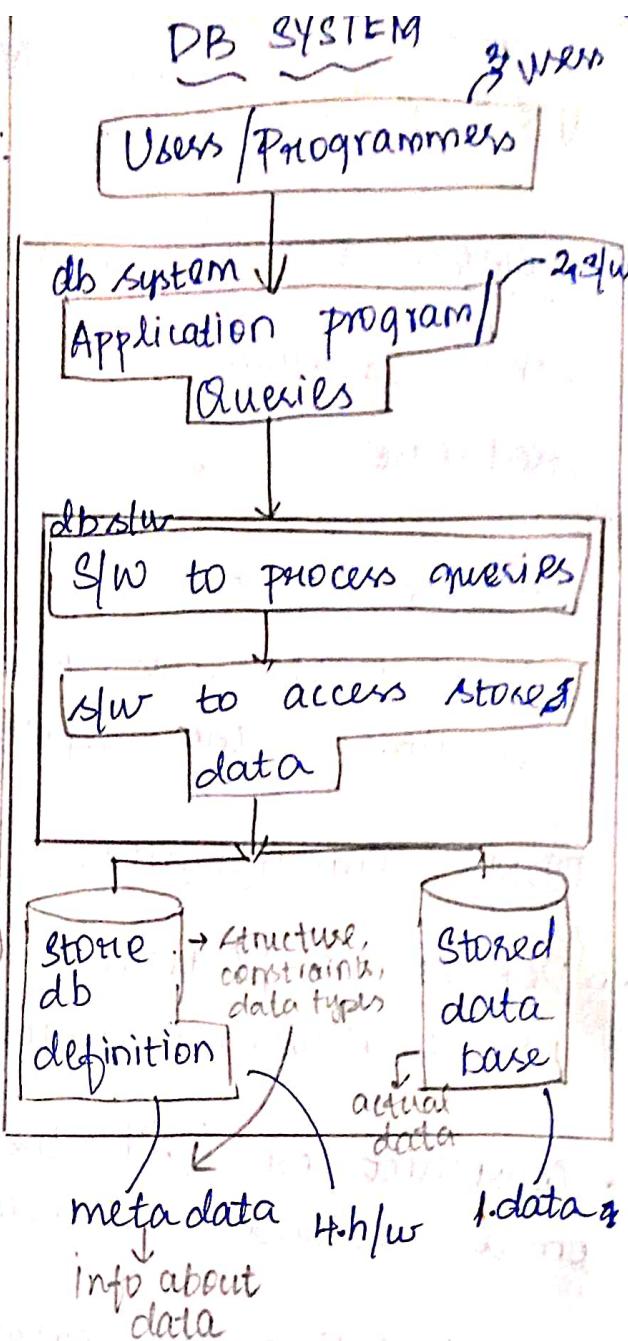
- To Define a particular database in terms of its data types, structures and constraints.
- Construct or load the initial database content on a secondary storage medium.
- Manipulating the database
  - Retrieval - querying, generating reports
  - Modification - Insertion, deletion and update to its contents
  - Accessing the db through applications.
- Processing and sharing: by a set of concurrent users and application programs yet keeping all the data valid and consistent.

## Database Implementation:

- Data Types
- Structures
- Constraint
- Sharing the data
- Querying

## Database Applications:

- Numeric, Textual (notepad)
- Multimedia dbs (YouTube, audio, video)
- GIS (Geographical info system)
- Data warehouses
- Realtime and active db



Data Models: How we store data.

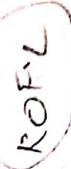
↳ visual representations of an enterprise's data elements and the connections b/w them.

↳ FMS: File management system

→ can store any no. of data

→ stored in sequential order

draw back → retrieval also begins from the begin



↳ HM: Hierarchical model

→ based on one to many relations → tree-like structure

→ less time consumption than FMS

→ can't jump from one data to another  
(has to follow branches of a tree)

↳ NM: Network model

→ faster than both because we use pointers to process data

→ Too many pointers might increase the complexity in searching data. → so can't use for large datasets

1970: EF, Kod

↳ RM: Relational model

→ Table form (data storage)  
↳ rows and columns.

eff:

Any database satisfying more of 12 rules  
of COD → perfect RDBMS

↳ Oracle satisfies 11/12 rules

### ORACLE LAB:

↳ Oak Ridge Arithmetic Calculating Logical  
Electronic device

↳ got popular in 1971

DDL: Create, Alter, Drop, Truncate, Rename

(data definition lang)  
add → modify → delete data  
alone not structure

### SYNTAX:

- Create table Stu (Sno number,  
                            Sname varchar(20),  
                            S.sex char(1));
- Alter table Stu drop column Sno;
- Alter table Stu add (Saddr varchar(30));
- Alter table Stu modify (Saddr varchar(10));
- Rename Stu to Student;
- Truncate table Student; → removes all  
                            rows but the  
                            structure, col,  
                            const, index  
                            remain same
- Drop table Student;

~~DATE~~ eff. ~~10/10/2018~~ DML : Insert, Update, Select, Delete, Merge

- Insert into Student values (1, 'Ram', 'M', 'Hosur');

- For multiple insertions:
- Insert into Student Values (2 Sno, 'R.Sname',  
'R.Ssex', 'R.Saddr');

Insert into all

Insert into Student ( );

Insert into Student ( );

~~Select \* from dual;~~

↓  
dummy table

- Update Student Set S-name = "krishna" where S-name = "Ram";

- Delete from Student; → delete all rows

- Delete from Student where Sname = "Krishna";

now find out the last inserted id

PCL  
Data control commands

at now after running the above query  
Delete contents of a column

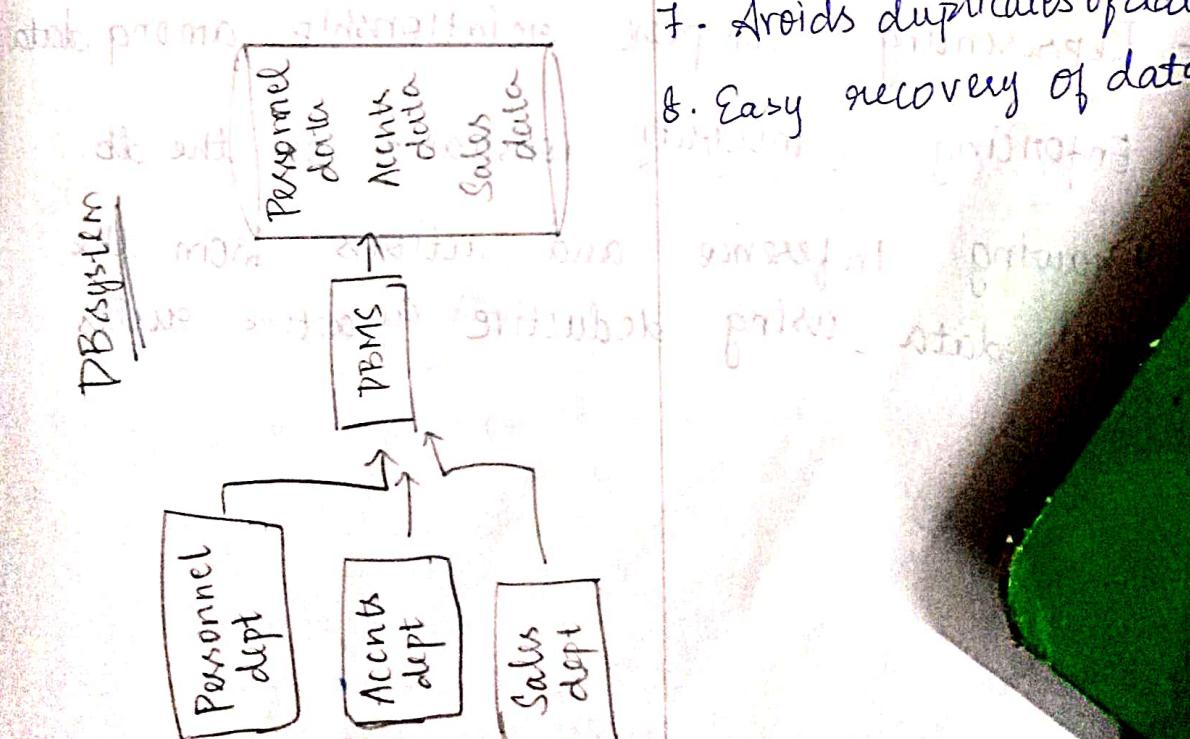
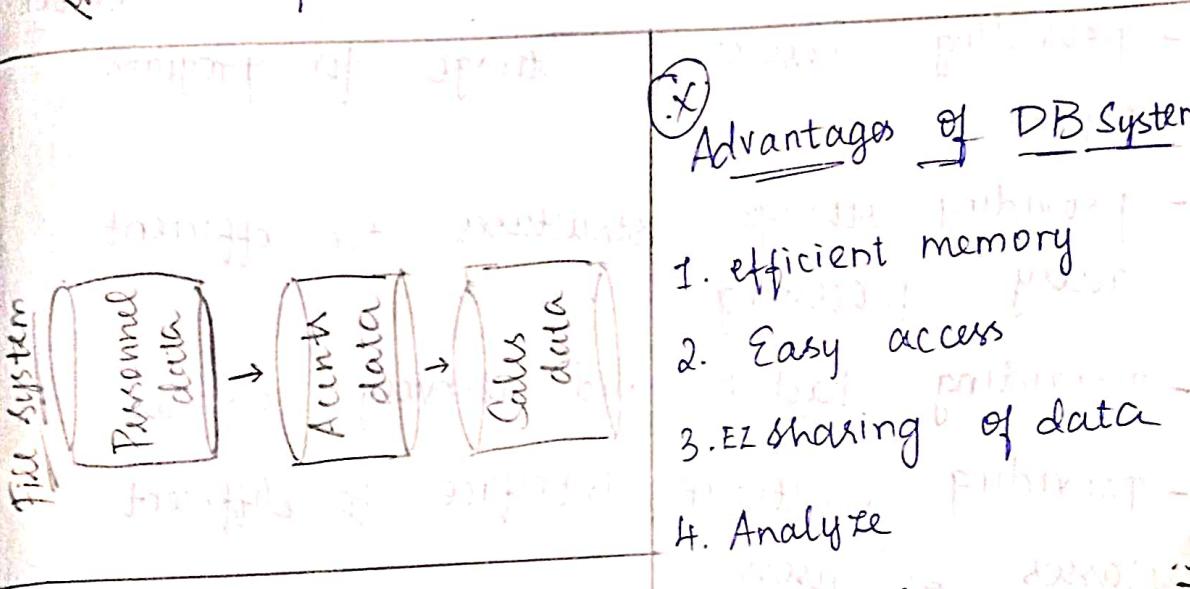
(f)

eff:

## Main characteristics of the db approach.

- self describing nature of a db system
  - A DBMS catalog stores the description of a particular db (eg. types, structures, constraints)
    - how to view
    - the description is called meta-data
- Insulation between programs and data
  - Called program-data independence
  - allows changing data structures and storage organization without having to change the DBMS access program.
- Data abstraction
  - A data model is used to hide storage details and present the users with a conceptual view of the db.
- support of multiple views of the data
  - Each user may see a different view of the db, which describes only the data of interest to that user.
- sharing of data and multiuser transaction processing
  - allows a set of concurrent users to retrieve from and to update the db.

- Concurrency control with the DBMS guarantees that each transaction is correctly executed or aborted.
- Recovery subsystem ensures each completed transaction has its effect permanently recorded in the db.



## Advantages of using of the db approach:

- controlling redundancy in data storage and in development and maintenance effort
  - sharing of data among multiple views
- Restricting unauthorized access to data <sup>→ PCL</sup> <sub>Grant, Revoke</sub>
- providing persistent storage for program objects
- providing storage structures for efficient query processing.
- providing backup and recovery services
- providing multiple interface to different classes of users
- + Refreshing
- Representing complex relationship among data
- Enforcing integrity constraints on the db.
- Drawing inference and actions from the stored data using deductive or active rules

## Data Abstraction :-

Level of abstractions:

→ Physical level / Internal level:

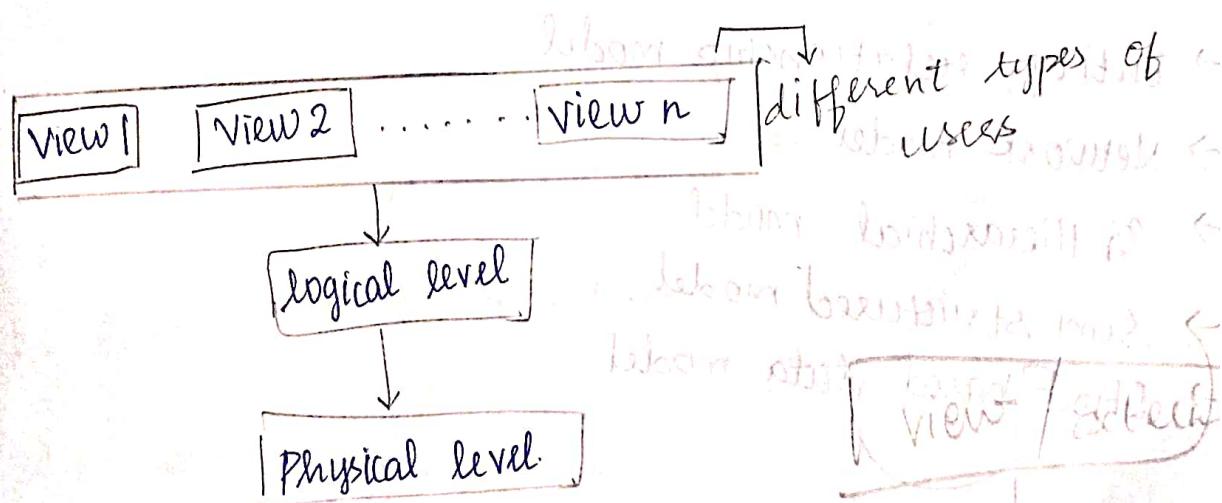
- describes how the data are actually stored

→ Logical level / Conceptual level:

- Describes what data are stored in the db and what relationship exists among the data

→ View level / External level:

- Highest level of abstraction describes only part of the entire db.



Level of data abstraction.

Logical/conceptual

Physical/internal

Instance:

the collection of info stored in the db at a particular moment is called an instance. → the value of variables @ a particular moment.

Schema:

The overall design of the db is called db schema

↳ blueprint of a database which describes how the data may relate to other tables / data models

Application program: to exhibit physical data independence if they do not depend on the physical schema

↳ is a comprehensive, self-contained program that performs a particular function directly for the user.

Data models: → describes the structure of the db

→ relational model

→ entity - relationship model

→ Network model → star, hub, ring

→ Hierarchical model

→ semi structured model

→ Object based data model

comprehensive, self-contained program that performs a particular function for an application

LAB:

1. Select \* from tab; → contains all tables available.

Constraints:

- 1. Key constraints → Primary key → to maintain uniqueness (col)
- 2. Default constraints
- 3. Not Null constraints → not null always
- 4. Check constraints → allows null values
- 5. Unique constraints → should be PK of another table

→ these come under DDL

2. Alter table <table-name> modify <col-name>

notnull;

3. <col-name> varchar(10) notnull; → while defining

4. Alter table ~~product~~ <table-name> add constraint

P.info-Pk primary key  
 ↓  
 (model-no)  
 col name

only for  
 key constraint  
 use  
 gr name  
 of  
 constraint.

5. Alter table PC add constraint PC-Pk primary key  
 (model-no)

6. Alter table pc add constraint <sup>ref:</sup> Pc-FK foreign key  
 (model\_no) references product\_info (model\_no);

7. Alter table laptop add constraint Lap-Fk  
 foreign key (model\_no) references product\_info (model\_no)  
on delete cascade;

if this col is deleted from parent table it

will be deleted from this table also

8. Alter table printer add printcode varchar(10);

9. Alter table laptop modify HDD number(4);

10. ~~Select \* from user-constraints;~~

11. Select table-name, constraint-type, constraint-name  
 from user-constraint where table-name = 'pc';

eff:  
→ referential constraint (foreign key)

↗ PC      R      PC-FK

PC      P      PC-PK

PC      C      PC-CK

One table can  
have only  
one primary  
key

12. Alter table PC ~~modify~~ add speed default 12;

13. Alter PC add ~~column~~ constraint PC-CK check  
price > 0;

14. Alter table product-info add constraint  
Pinfo-ck check (type in ('LP', 'PC', 'PR'));

15. Create table printerinfo as select \* from printer;  
modifying speed not null;  
modifying speed unique;  
modifying speed default 2;

16. Create table printerinfo as select \* from printer  
where 1>2;

17. Delete from printerinfo;  
only data deleted

18. drop table printerinfo; → including structure

Q1:

19. Update PC set price = price + 10;

20. Update PC set price = price + 10 where  
~~model~~ model-no = 'PC111';

21. Delete from PC where model-no = 'PC112';

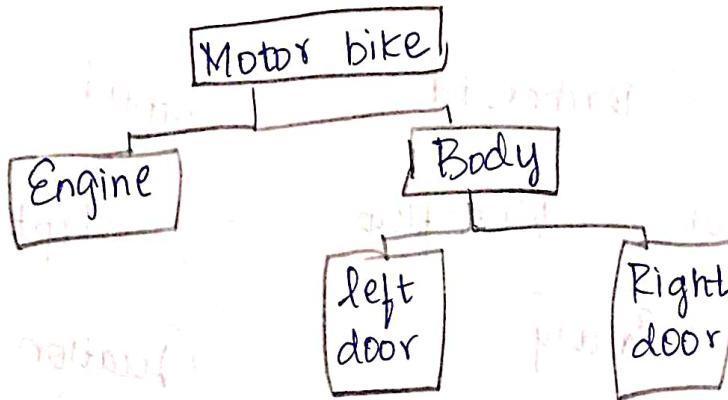
22. Select \* from product-info where  
 model-no = (Select model-no from PC  
~~where price >= 1000~~ where price = 100);

## DATA MODELS :-

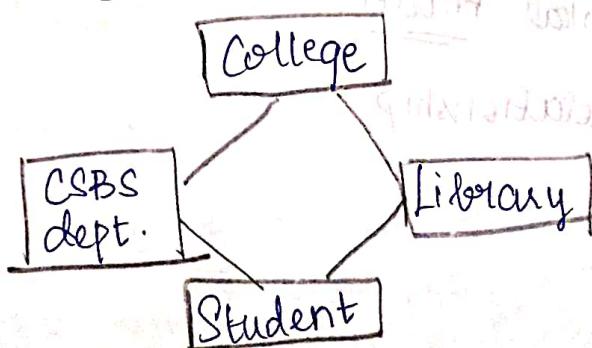
Relational model: → tabulated form

Empid	Ename	Salary	Did
101	Ram	45000	10
102	Renu	42000	11

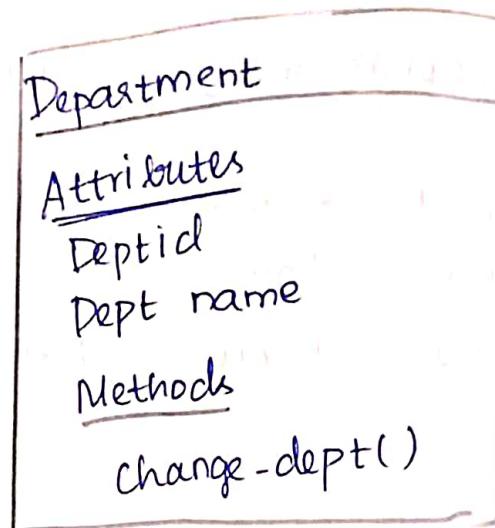
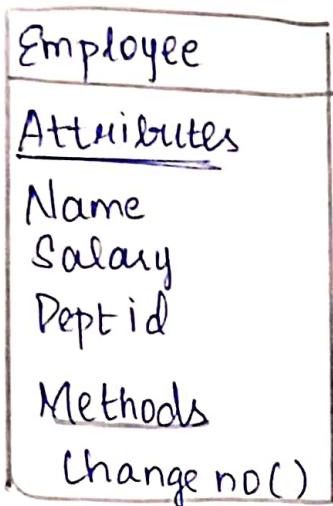
Hierarchical mode: → tree structured



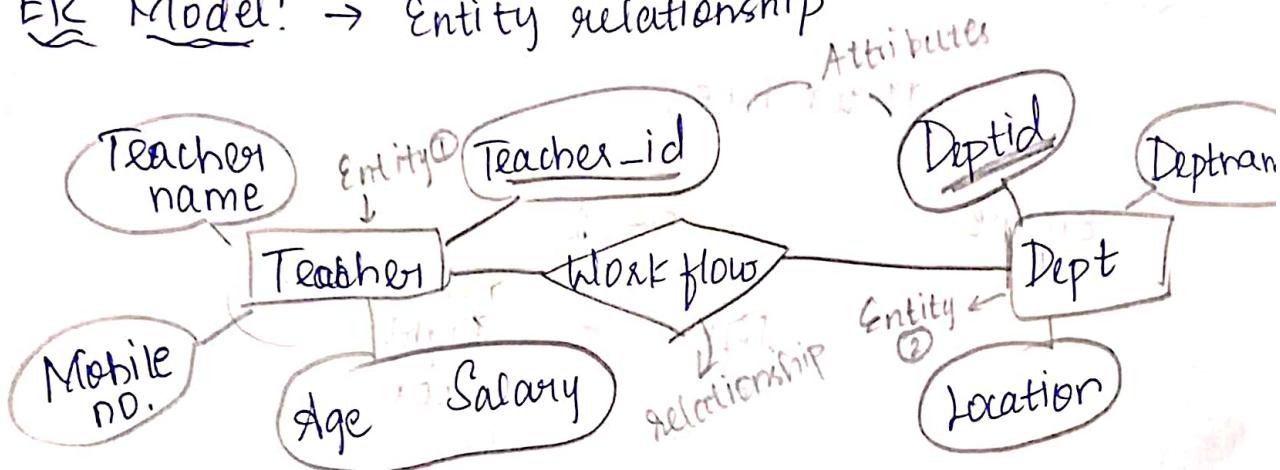
Network model: → diff types of topology



## Object Oriented data model:



ER Model: → Entity relationship



## Features Of Hierarchical model:

- \* one-to-many relationship
- \* parent-child

## Advantages:

- very simple & fast to traverse through a tree like structure
- Any change in the parent node is automatically reflected in the child node

### Disadvantage:

- complex relationships are not supported
- parent node is deleted, the child also gets deleted.

### Features of Network model:

- ability to merge more relationship
- many paths
- circular linked lists

### Advantages:

- the data can be accessed faster as compared to hierarchical model
- Parent-child relationship

### Disadvantages:

- change like update, deletion and insertion very complex

### Features of ER:

Entities: real-world things

Attributes: real-world property

Relationship: Tells how 2 attributes are related

DBMS vs RDBMS

## Advantages:

- Simple
- Effective communication tool
- Easy conversion to any model

## Disadvantages:

- hidden info
- No industry standard for notation

## features of Relational model:

- Tuples
- attribute or field

## Advantages:

- simple
- Scalable
- structural Independence

## Disadvantages:

- H/W overhead
- ✓ requires

### ① Presentation layer:

- client layer
- purpose is to communicate with the application layer

### ② Application layer:

- business logic layer
- middle layer
- transmit partially processed data

### ③ Database layer:

- stored info
- insert, update
- ↑ delete

→ it's a diagram of

Features of DB model: → describes interrelated things of interest in a specific domain of knowledge

- is an abstract of our DB
- high level or conceptual model

## Relation:

- format and relations in a way a db could understand
  - Implementation of our model

→ Implementation of our model  
Two-tier architecture: ex: client/server

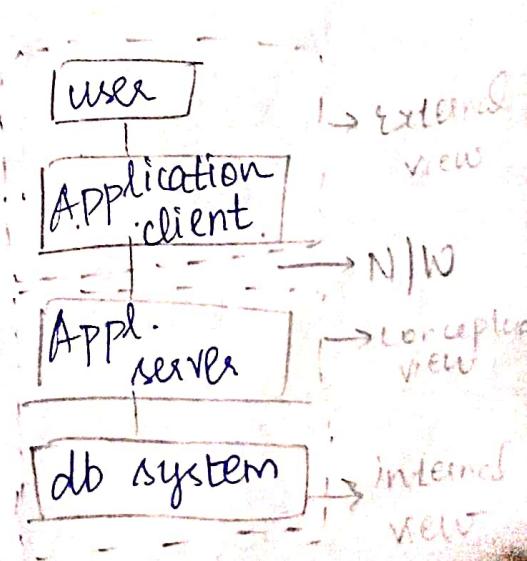
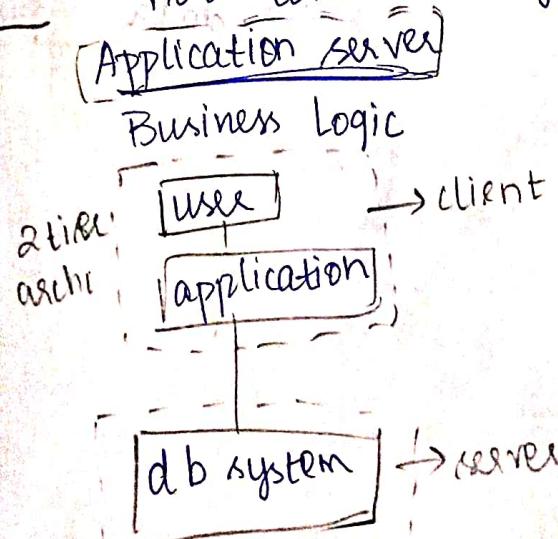
Two-Tier Architecture

→ the application resides at the client machine, where it invokes db system functionality at the server m/c through any language.

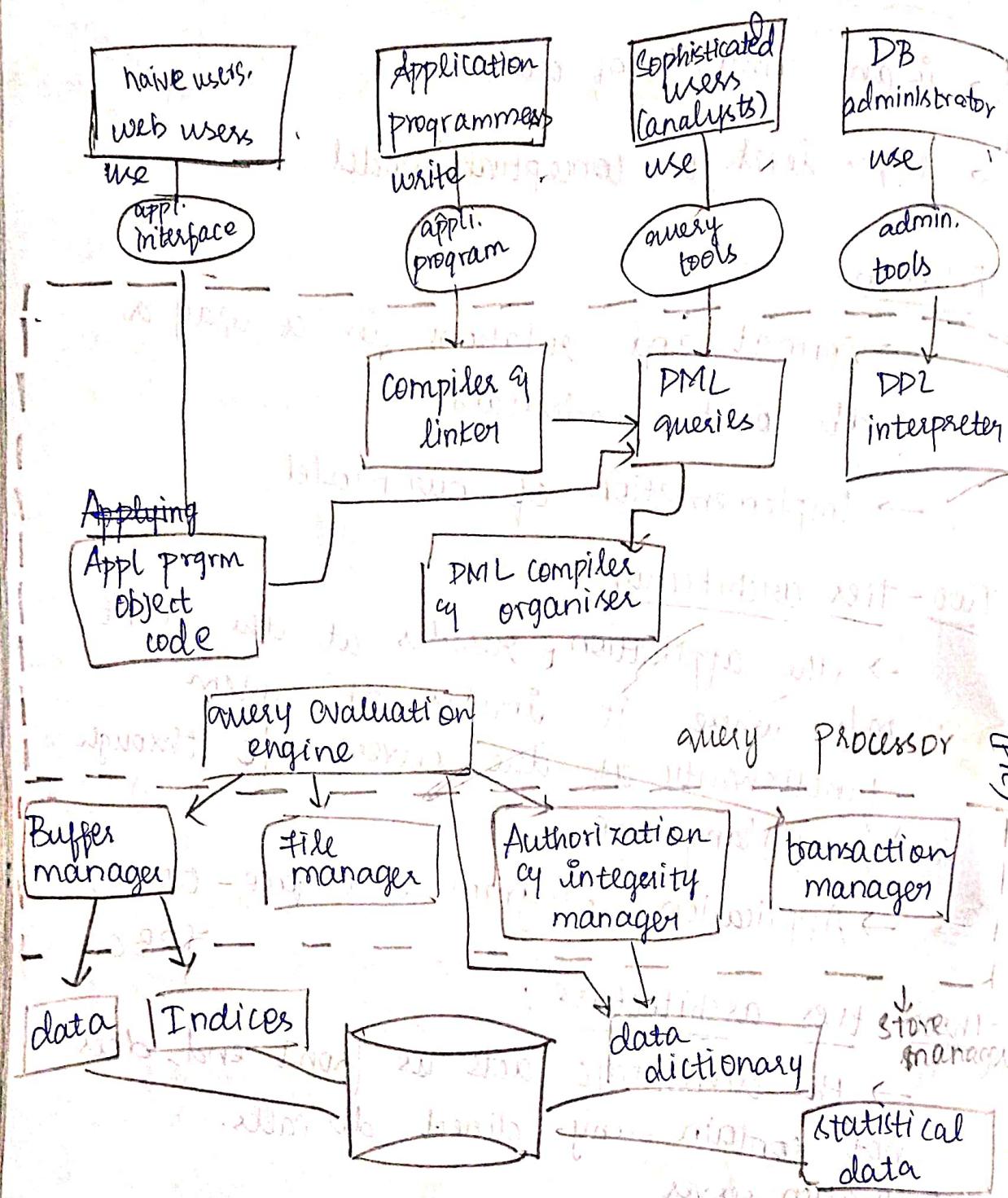
→ Application program interface - ODBC, TDBC

## Three tier architecture

→ the client m/c acts as front end, does not contain any direct db calls.



## Database System architecture:



## Query Processor:

### ⇒ DDL interpreter:

- which interprets DDL statements and record the definitions in the data dictionary.

### ⇒ DML Compiler:

- which translates DML statements in a query language into an evaluation plan consisting of low level instructions in the query evaluation engine understands.

### ⇒ Query evaluation Engine:

- which executes low-level instructions generated by the DML compiler.

### ⇒ Data files:

- which store the db itself

### ⇒ Data dictionary:

- which stores the metadata about the structure of the db, in particular the schema of the db.

## ⇒ Storage manager:

### → Authorization and integrity manager:

- which tests for the satisfaction of integrity constraints and check the authority of user of access data

### → Transaction manager:

- The db remains in a consistent (correct) state despite system failures and the concurrent transaction executions proceed w/o conflicting

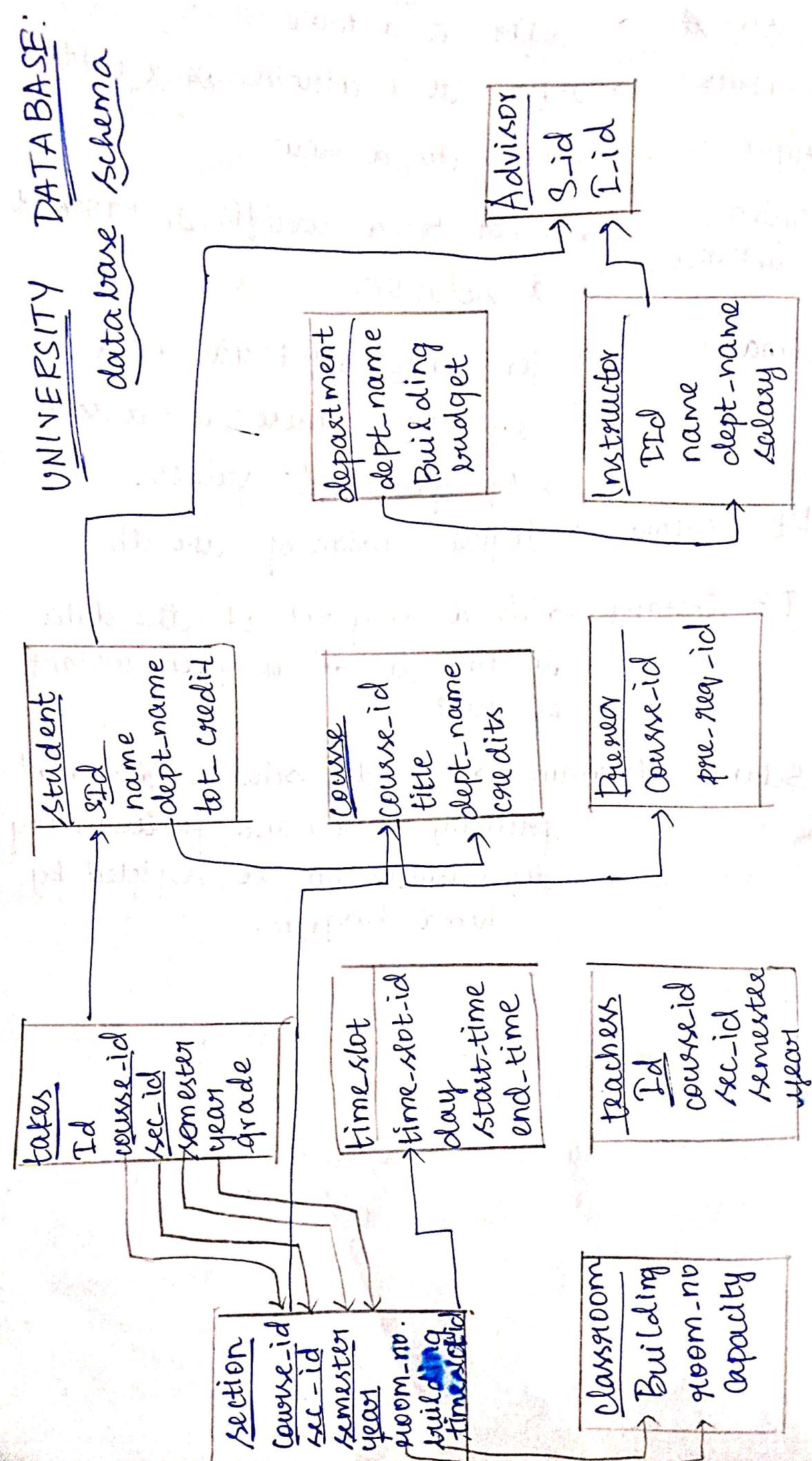
## → File manager:

- manages the allocation of space on disk storage and data structures used to represent info stored on disk

## ⇒ Buffer manager:

- responsible for fetching the data from disk storage units into main memory and deciding what data to access in main memory

- Relation ~~at~~ → refer to a table
- attribute → refer to a column of a table
- tuple → refer to a row
- relation instance → refer to a specific instance of a relation
- domain → for each attribute of a relation there is a set of permitted values.
- DB Schema → logical design of the db
- DB instance → is a snapshot of the data in the db at a given instant of time.
- Schema diagram → A db schema along w/ primary key and foreign key dependencies can be depicted by schema diagram.



## KEY CONSTRAINTS:

- Primary key → unique / not-null
- foreign key → PK of another table
- Candidate key → can be PK (minimal id of superkey)
- Super key → combination of columns in a table to maintain uniqueness
- Unique key → unique PK

## ER model:

### Entities:

- ex: Employee, Department, Project

### Attributes:

- Employee: Name, SSN, DOB, Address

### Types of attributes:

#### → Simple:

- Single atomic value for the attribute

#### → Composite:

- Several components

Ex: • Address (Apt #, House #, St., City)

• Name (first name, lastname, middle name)

#### → Multivalued:

- Multiple values for that attribute

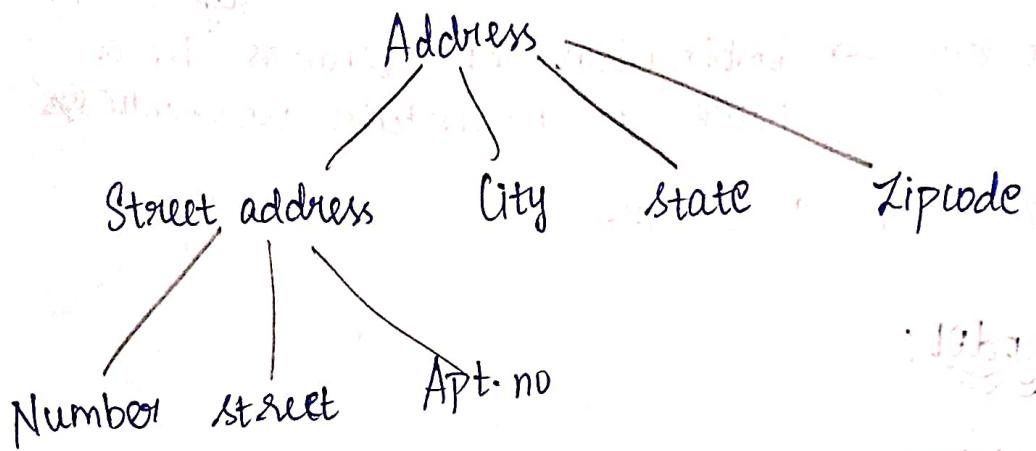
Ex: • color of a car

• previous degrees of a student

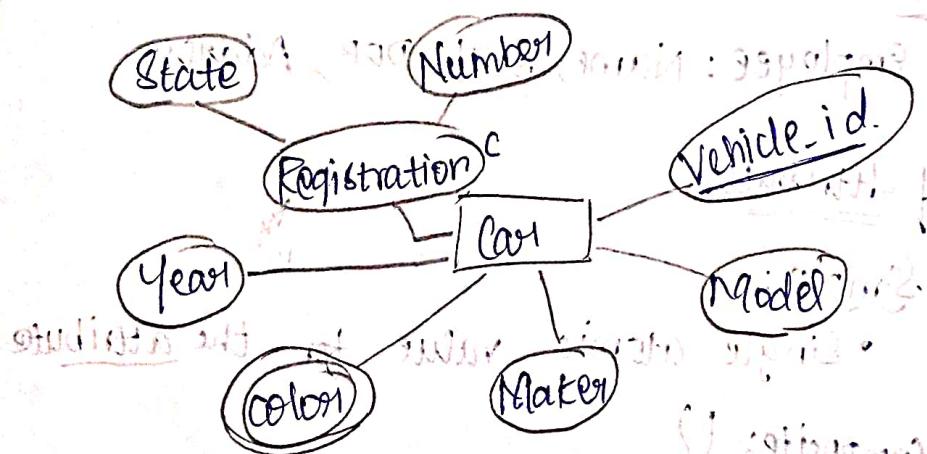
→ Composite valued multivalued:

{previous degrees (college, year, degree, field)}

→ Composite:



Vehicle tag:



Entity set:

DATA: Registration (Number, state), vehicle-id, make, model, year, h\_color.

CAR 1:

(ABC123, Texas), TK629, FORD, H door, 2004,  
{red, black})

CAR 2:

(NSY720, Texas), TD729, NISSAN, H door, 2002,  
{blue, red})

Relationship:

- when an attribute of one entity type refers to another entity type.
- Relationship type R among 'n' entity types  
 $E_1, E_2 \dots E_n$
- Defines a set of associations among entities from these entity types

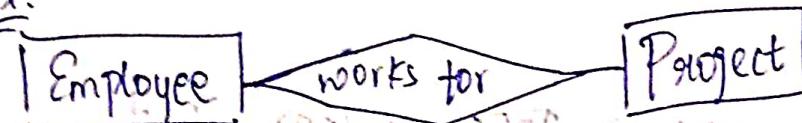
Relationship instances r:

Each  $r_i$  associates n individual entities.  
 $(e_1, e_2 \dots e_n)$  Each entity  $e_j$  in  $r_i$  is a member of entity set  $E_j$ .

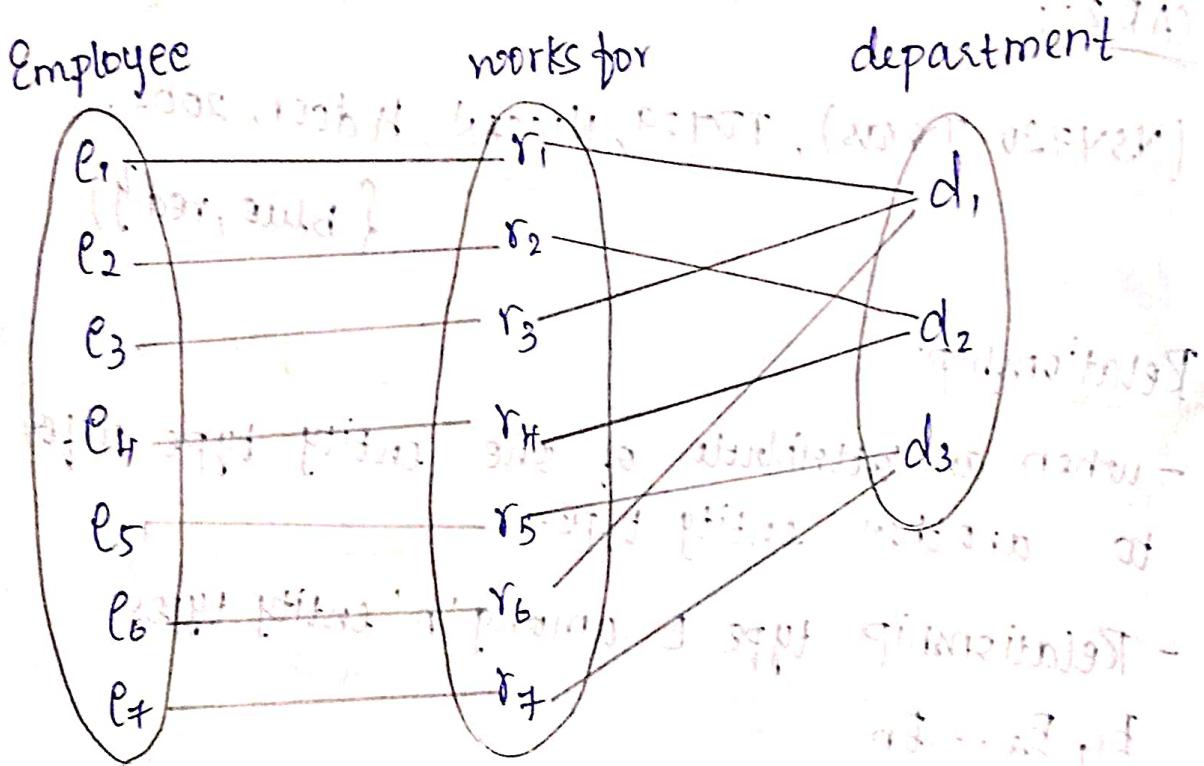
Degree of relationship type:

No. of participating entity types: binary, ternary

Ex:



Relationship instances:



Some instances of the work for relationship b/w employee and dept.

relationship between employee and department is that all the employee works for one department only (one-to-one relationship).

- SAL, EPdia, Archi, UML, Rational Rose, Visio, Microsoft Visio, Microsoft Project
- terminal (private)
- no. of types > 1
- in a present @ any instance
- name (str)
- cont
- (intension, extension)
- (DBMS, RDBMS)
- or delete cascade - DI
- weak attribute

## LAB:

## Arithmetic, Logical, Set Operations:

$\downarrow$  = 3703 OR AND  
+  
 $\downarrow$  Union all  
AND union  
Intersection minus

$$A = \{1, 2, 3, 4\}$$

$$B = \{3, 4, 5, 6\}$$

$$A \cup B = \{1, 2, 3, 4, 5, 6\}$$

A union all B = {1, 2, 3, 4, 3, 4, 5, 6}

$$\text{Intersection } B = \{3, 4\}$$

$$A \text{ minus } B = h(1, 2) \setminus \{ (1, 2) \} + f(3, 4)$$

$$B \quad " \quad A = \{5, b\} \quad (\text{ss})$$

⇒ Select name from worker > union select name  
from worker skill,  
(Intersect / minus)  
union all " ;

⇒ In operator:

→ select \* from Lproduct\_info where type in ('PC', 'PR')  
→ " " " where type = 'PC'  
→ " " " type not in ('PC', 'PR');

→ update Lpc > set price =  $\pi$  Price + 1000  
where model-no = 'Pc112'

→ Select \* from LPC.

→ select \* from Laptop;

→ select \* from pc where price > 5000 and  
price < 10000;

→ select \* from <pc> where price between  
5000 and 10000;

→ Select \* from <Laptop> where model-no = 'LP112'  
or model-no = 'LP113';

⇒ Aggregate functions:

Sum, average, max, min, count.

→ Select max(price) from <PC>;

→ Select min(salary) from <Emp>;

→ Select count(\*) from <product-info>;

⇒ Group By function:

→ Select dept-id, count(\*) from <Student>;

→ Select dept-id, sum(salary) from <Employee>;

→ Select dept-id, sum(salary) from <Employee>;  
Group by dept-id;

→ Select RAM, count(\*) from <PC> group by RAM;

⇒ Group By - having :- condition for group by:

→ Select RAM, count(\*) from <PC> group by RAM  
having count(\*) > 1;

- ⇒ Order by:
- Select \* from LPCS order by price;
- .. " " desc;

Translate  $(3 * [2+1] / [8-4], [ ] 83, () 0)$

gives  $3 * (2+1) / (8-4)$

- Primary key      SuperKey
- unique
  - Not null
- It has all the capabilities to be a PK

Composite key

◦ A combination of attributes to maintain uniqueness

Superkey

◦ Work of maintaining uniqueness

◦ Helped in defining MAT rules

◦ If it has primary key

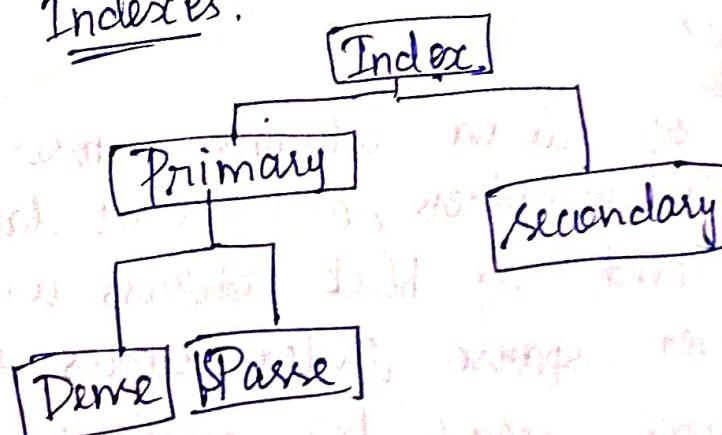
## Indexing:

- is a data structure technique which allows you to quickly retrieve records from a db file.
- An index is a small table having only two columns:
  - 1st column comprises a copy of the primary (or) candidate key of a table
  - 2nd column contains a set of pointers for holding the address of the disk block where that specific key value is stored.

## Advantages:

- \* takes a search key as input
- \* efficiently returns a collection of matching records

## Types of Indexes:



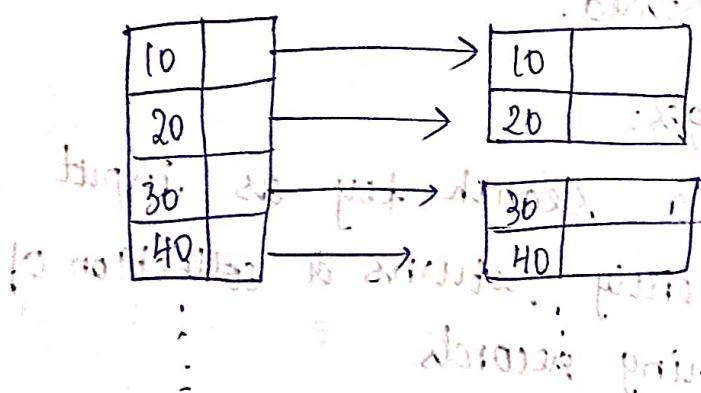
## Primary Index:

→ is an ordered file, which has fixed length & size with two fields. The first field is the same as PK and second field pointed to the specific data block - one to one relationship.

### → Two type:

- \* Dense
- \* Sparse

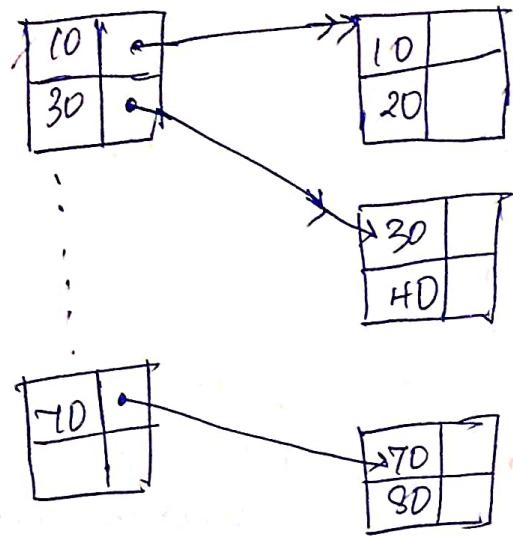
→ Dense index: a record is created for every searchable key in the db, this helps you to search faster but needs more space index records.



### → Sparse Index:

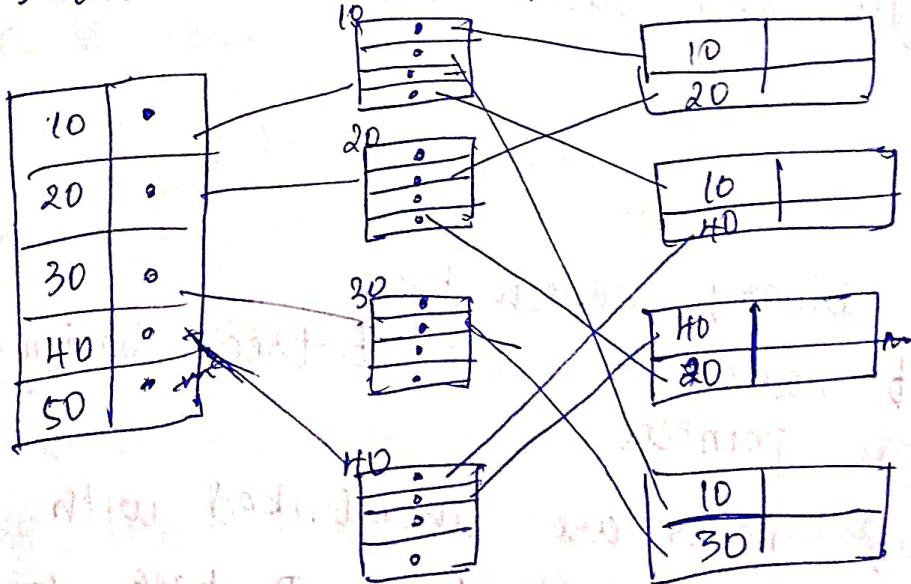
A range of index columns stores the same data block address, and what data needs to be retrieved, the block address will be fetched. However sparse index stores records for only some search key values. It needs less space, less maintenance over read.

It is slower compared to the dense index, for locating records <sup>eff.</sup>



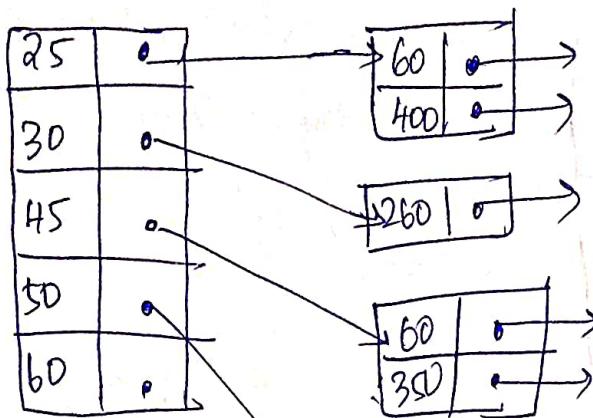
Secondary index: Generated by a field which has a unique value for each record and it should be a candidate key called non-clustered index.

→ Two levels db indexing technique is used to reduce the mapping size of the first level.

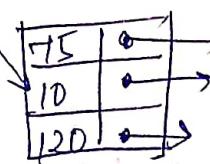


B  
is

## Multilevel Index



Pointer to data



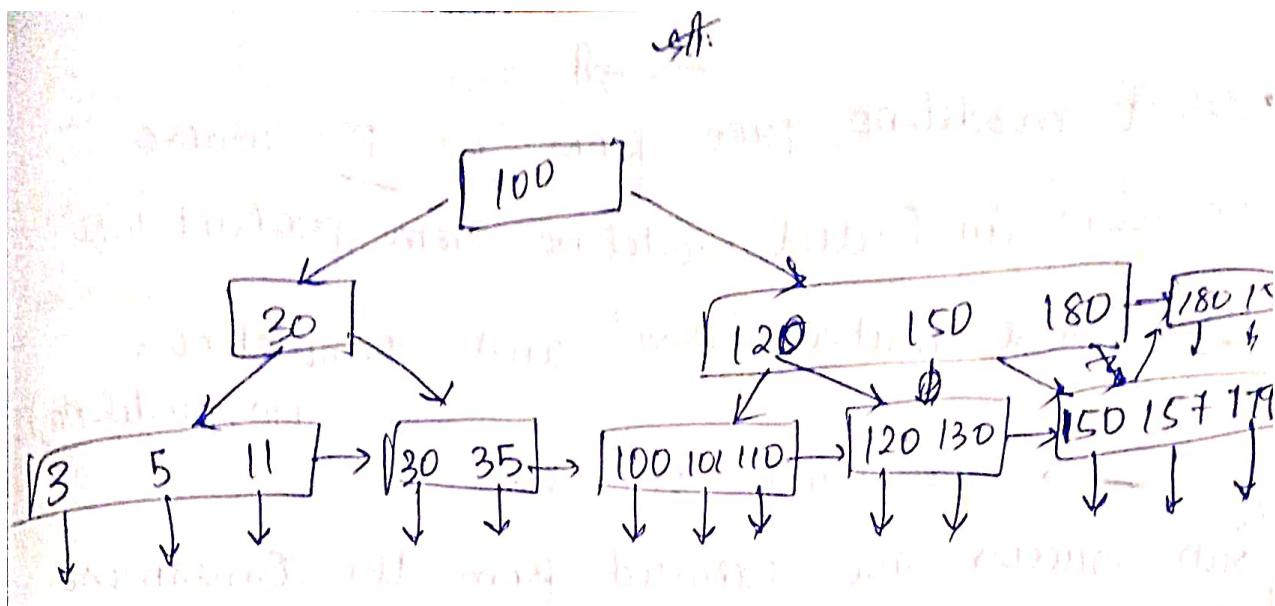
The outer blocks are divided into inner blocks which in turn point to the data disk.

Multi level indexing in db is created when a primary index does not fit in memory.

You can reduce the no. of disk accesses to sort any record and kept on a disk as a sequential file and create sparse base on that file.

## B-Tree:

- Balanced binary search tree
- All leaf nodes of the B-tree signify actual data pointers.
- All leaf nodes are interlinked with a linked list, which allows a B-tree to support both random & sequential access.



Create index in\_pc\_mo on pc (model-no);

### LAB:

#### Sub Query / Nested Query:

→ Outer Query

→ Inner Query

Select salary from employee  
where salary > (Select salary from employee where empid = 700);  
Select salary from employee where salary > (Select avg(salary) from emp);  
Select model-no, price from pc where model-no  
in (select model-no from product\_info where  
maker = 'IBM');

- Select salary from employee where salary >  
(select salary from employee where empid = 700);
- Select salary from employee where salary >  
(select avg(salary) from emp);  $\Rightarrow$  These 2 return single value. So relational
- Select model-no, price from pc where model-no  
in (select model-no from product\_info where  
maker = 'IBM');  
  
referential integrity works here  $\Rightarrow$  This has to return more than a set of values so 'in'.  
1 table involved

- Select model-no, price from pc P where model-no in (select model-no from product-info pin where maker = 'IBM' and p.model-no = pin.model-no)
  - alias name for table.

(Sub queries are executed from the innermost query.)

### Subquery:

→ Correlated subquery

→ Non-correlated subquery.

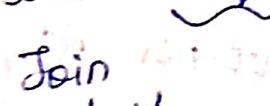
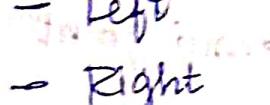
Non-correlated: The inner query is executed only once and then sends it to the outer loop.

Correlated: The inner query is executed for every record of the outer query.

### Examples:

- Select salary from employee where salary > (select avg(salary) from employee),
- Select empid from employee where empid in (select mgr-id from employee),

JOINS: 

- Inner Join → equijoin → Non equi join 
- Outer Join →  
  - Left
  - Right
  - full

For n tables

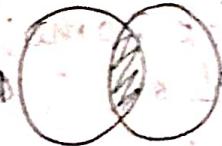
we can have

a max of  $\frac{n(n-1)}{2}$  joins

$n-1$  joins

Select model-no from pc P, product-info Pin

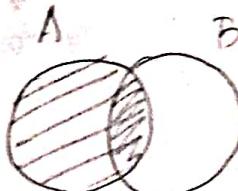
where p.model-no = pin.model-no; inner join

  $A \cap B$   $\rightarrow$  Inner Joint condition

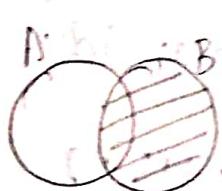
Select \* from pc; product-info; Pin

Full join  $\rightarrow$  Cartesian product

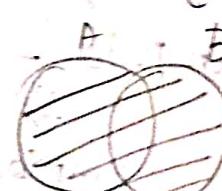
Select \* from pc; product-info; Pin where A.model-no = B.model-no (no condition)



Left  
Outer join



right  
Outer  
join



Full outer  
join

## SYNTAX:

Select w.name, wk.name from worker w, workerskill wk  
 where w.name = wk.name; → equi joint

Select name from worker w, workerskill wk  
 where w.name = wk.name;

Select w.name, wk.name from worker w,  
 workerskill wk where w.name <= wk.name;  
 non equi join not equal or !=

Select w.name, wk.name from worker w,  
 workerskill wk where w.name = wk.name and  
 w.age >= 30; → non-equi joint

Select e1.emp-id, e2.mgr-id from employee e1, e2  
 where e1.emp-id = e2.mgr-id;  
 ↓  
 employees working  
 as manager

Select w.name, wk.name from worker w,  
 workerskill wk where w.name <= wk.name(+);  
 (or)

Select w.name, wk.name from worker left outer  
 join workerskill ~~where~~ on w.name = wk.name;



Select w.name, wk.name from worker w,  
workerskill wk where wk.name(f) = wk.name  
right outer  
join

Views: is a logical rep of physical table

- Since DBMS is a multiuser accessible model

ID name no. . . . (25 columns)

- - - - -  
- - - - -  
- - - - -

View 1

View 2

ID name  
- - - - -  
- - - - -  
- - - - -

ID no. - - -  
- - - - -  
- - - - -  
- - - - -

→ changes in the view will be reflected in  
the actual table and vice-versa.

→ There is no separate memory for view  
table.

Create

view view-emp as ~~dept\_id~~, ename

Select (emp-id, ename) from employee;

Create view emp-view-sum as select

sum(salary), avg(salary) from employee where  
dept-id = 7;

Delete from view-emp where l-name = 'RAM';

Drop view emp-view;

view-emp as select

Indexes:

Create index index-emp on employee (ename);

Create index index-dept on dept (dname);

Create synonym <empl> for employee;

alias  
name

Create index ie  
on emp/empl

Sequence:

create sequence <sequence-name>

increment by n

start with n

max value n

min value n

cycle / no cycle

cache / no cache;

increment by n

start with n

max value n

min value n

Ex:

left:

- create sequence dept-deptno
  - increment by 1
  - start with 91
  - max value 100
  - no cycle
  - no cache;
- select sequence\_name, minvalue, max-value  
increment-by, last-number from user-sequences;
- insert into dept (deptno, dname, loc) values  
(dept-deptno.nextval, 'marketing', 'New delhi');
- select dept-dept no. curval from dept; dual;
- alter sequence dept-deptno
  - increment by 1
  - max value 9999
  - No cache
  - No cycle;
- drop sequence dept-deptno;

PL | SQL: Procedural language:-

→ An extension to SQL with design features of programming languages (procedural by pop)

→

PPT

## Advantages:

- ⇒ Avoids network traffic  $\rightarrow$  n no. of SQL statements can be executed @ a time
- ⇒ Pre Compiled
- ⇒ Portable (can work in any platform)

## Declaration:

- ⇒ V-first-name VARCHAR(35);
- ⇒ V-count NUMBER:=0;

## Executable Section:

```

BEGIN
  SELECT first-name, last-name
    INTO V-first-name, V-last-name
   FROM student
  WHERE student_id = 123;
  DBMS_OUTPUT.PUT_LINE
  ('Student name:' || V-first-name || V-last-name);
END;
  
```

## Exception Handling Section:

- ⇒ These statements are executed when a runtime error occurs

### EXCEPTION

WHEN NO\_DATA\_FOUND THEN

DBMS\_OUTPUT.PUT\_LINE ("Nope studen");

END;

## Reference Variables:

- It directly references specific database column or row
- It assumes data type of associated column or row
- Y.TYPE:

variable-name tablename.colname Y.TYPE;

- Y.ROWTYPE:

variable-name tablename Y.ROWTYPE;

DECLARE

~~DEFINITION~~ Sname Sailors.sname Y.TYPE;

type of column

Sname in  
Sailors

DECLARE

Reserves\_rec. Reserves Y.ROWTYPE;

accessing fields in Reserves\_rec:

Reserves\_rec.sid := 9; → (will it change?)

Creating a Record: → a user defined data type like struct

TYPE Sailor\_record\_type IS RECORD

(sname VARCHAR(10),

sid VARCHAR2(9),

age NUMBER(3),

sailor\_record Sailor\_record\_type);

BEGIN

Sailor\_record.sname := 'peter';

Sailor\_record.age := 45;

To get user input:

Set Serveroutput on

Accept P-Price Prompt 'Enter:'

Comments: /\* \*/

## Normalisation:

①	Stu-id	Proj-id	Stu-name	Proj-name

②	Stu-id	St-name	City	ZIP

③	Sid	Prog-lang	Stud-age
	101	Compiler N/w	20
	101	JAVA	20
	102	DBMS	20
	103	S/w Engg	21
	103	Compiler	21

All these tables are in first normal form

Convert to ① 2NF ② 2NF ③ 2NF

keya hoga? convey and communicate? what happens?

①	Stu-id	Stu-name

Proj-id	Proj-name

②	Stu-id	Stu-name	Zip

already in 2NF  
candidate key

Zip	City

(3)

E-id	S-age
------	-------

-sft.

S-id	Program
------	---------

B

INF → repeating ~~values~~ groups are non-existent

dNF → nonkey attributes are fully dep. on key attributes

3NF → transitivity is removed

BCNF → Every determinant is a candidate key

4NF →

5NF →

Emp-id	Emp-nationality	emp-dept	dept-id	dept-no-emp
1001	UK	Production	D001	200
1002	USA	Sales	D002	300

### ① Student

	Sname	Flag	Team-name	Flag
S1	Alex	A	Fighters	A
S2	Henry	B	WARRIORS	B
S3	Michel	B	Supercool	C
S4	Charlie	D	Beast	D
S5	Ace	D		
S6	HInge	C		

### Formal Lang Query:

select Salary from emp  
 ↓  
projection      where empid = 123  
 ↓  
selection

$\sigma_{A=B \wedge D \geq 5}(r)$  : Select Operation :

- $\sigma_p(r)$  where p is the selection predicate.
- $\wedge$  (and),  $\vee$  (or),  $\neg$  (not)

Cartesian join vs natural join (H10) CLA 2

2nd maximum Salary: select max(salary) from details where salary < max(salary)  
 Ummm.

P ↗ left:

Select  $\max(\text{salary})$  from employee where  
 $\text{Salary} \in (\text{select salary from employee}$   
 $\text{where salary} \geq \max(\text{salary}))$

Selection

Select:  $\sigma$

Project:  $\Pi$

Union:  $\cup$

Sel diff: -

Cartesian:  $\times$

Rename:  $\rho$

$\sigma_P(r) \rightarrow \text{relation}$

$\downarrow$   
 Selection  
 Predicate

⑥  $\Pi_{\text{dept}}(r)$   $\rightarrow$   $\text{relation}$

replace  $\{\text{dept}\}$   $\rightarrow$   $\{2000\}$

translate  $\{\text{dept}\}$   $\rightarrow$   $\{\text{dept}\}$

visit:  $\{\text{dept}\}$   $\rightarrow$   $\{\text{dept}\}$

$\Pi_{\text{course}}(r)$

$\Pi_{\text{course}}(r)$

$\Pi_{\text{course}}(r)$

## The last Update Problem:

↳ Some operation is interleaved  
 ↳ it makes the value of some db  
 inc

## The Temporary Update:

↳ failure after update A then  
 another trans accesses it before  
 rollback.

### Incorrect Summ

↳ one trans is agg another  
 its update  
 misc out on  
 the updated values

$$\text{Sum} = 0$$

read A

$$\text{Sum} = \text{Sum} + A$$

read X  
 $X := X - N$   
 write X

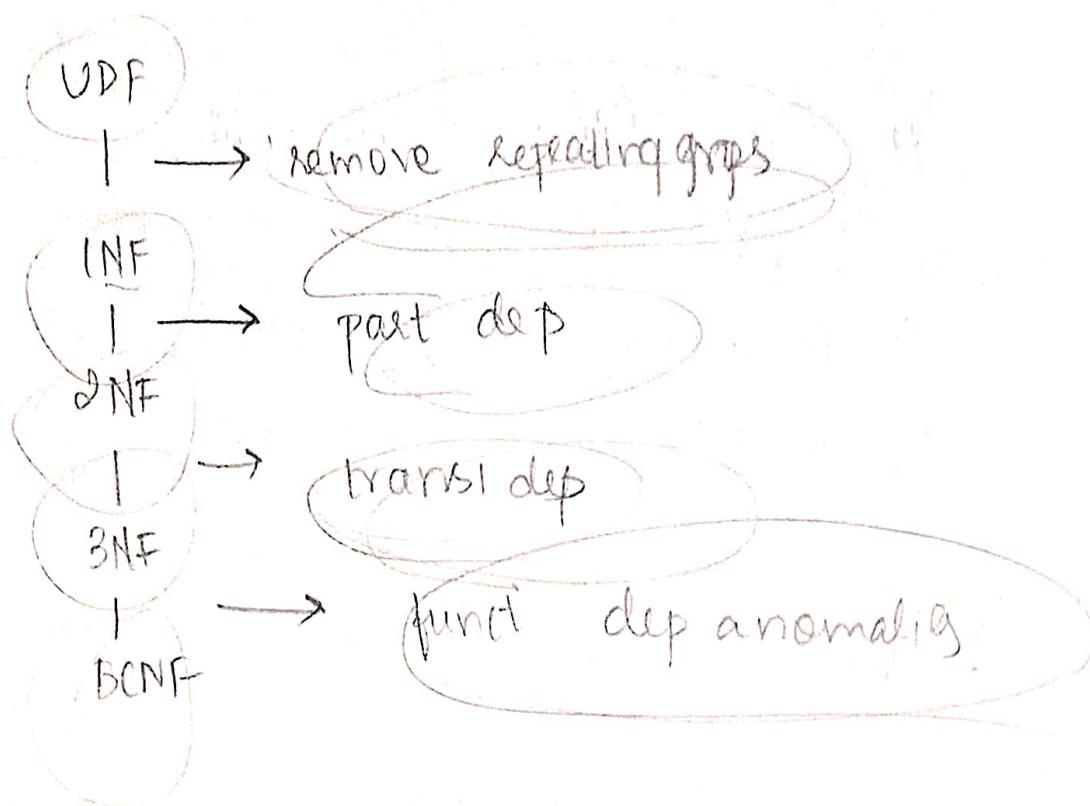
read X

$$\text{Sum} = \text{Sum} - X$$

read Y

$$\text{Sum} = \text{Sum} + Y$$

read Y  
 $Y := Y + N$   
 write Y



Select w.name, ~~wk.name~~ from worker w, workcenter wk  
where w.name = wk.name.

round (price, 2)

length (name)

substr (phone, 9, 13)

instr (name, 'a')

soundex (city) - so

select \* from painter where  
price in (select price from  
printl where model-no in  
(select modelno from  
select make  
where make

## Serializability:

A schedule is serializable if it is equivalent to a serial schedule.

1. Conflict serializable
2. View

## Conflict Serializable:

⇒ if a schedule  $S$  can be transformed into a schedule  $S'$  by a series of swaps of non-conflicting <sup>good</sup> instructions, we say that  $S$  and  $S'$  are conflict equivalent.

⇒ Schedule  $S$  is conflict serializable if it is conflict equiv. to a serial schedule.

## View Serializable:

A schedule is called view serializable if it is view equal to a serial schedule.

## Conflicting operation:

- \* They belong to diff. trans
- \* They operate on the same data item
- \* Atleast one of them is a write operation.

## Precedence Graph for testing Conflict Serializability

Graph  $(V, E)$  consisting of a set of nodes  $V = \{T_1, T_2, T_3, \dots, T_n\}$  and set of directed edges  $E = \{e_1, e_2, \dots, e_n\}$ .

- the graph contains one node for each transaction  $T_i$ .
- An edge  $e_i$  is of the form  $T_j - T_k$ , where  $T_j$  is the starting node of  $e_i$  and  $T_k$  is the ending node of  $e_i$ .
- An edge  $e_i$  is constructed between node  $T_j$  to  $T_k$  if one of the operations on  $T_j$  appears in the schedule before some conflicting operation in  $T_k$ .

### Algorithm:

1. Create a node in the graph for each participating transaction in the schedule.
2. For the conflicting operation read-item and a write-item if a transaction  $T_j$  executes a read-item(x) after  $T_i$  executes a write-item(x) draw an edge from  $T_i$  to  $T_j$  in the graph.
3. For the conflicting operation write-item(x) and read-item(x). if a transaction  $T_j$  executes a write-item(x) after  $T_i$  executes a read-item(x) draw an edge from  $T_i$  to  $T_j$  in the graph.

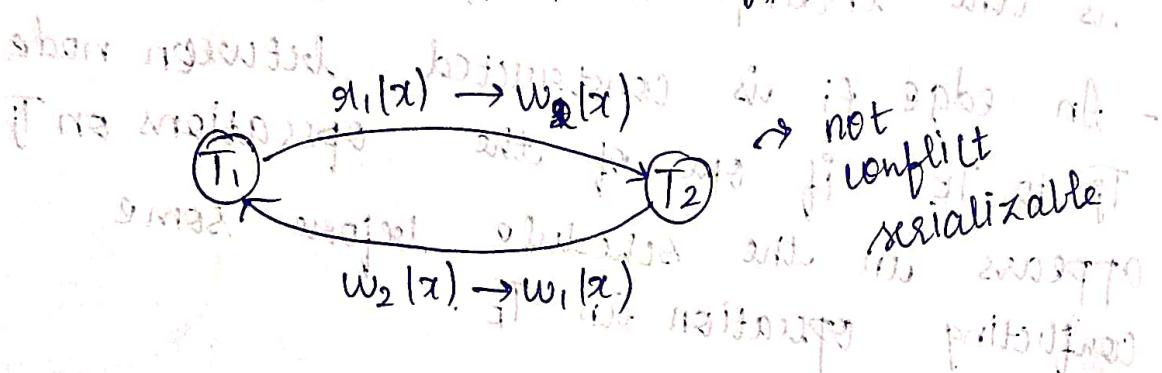
that is so cool & wow!

4. for each conflicting operation  $\text{write-item}(x)$  and  $\text{write-item}(x)$  of a transaction  $T_j$  executes a  $\text{write-item}(x)$  after  $T_i$  executes a  $\text{write-item}(x)$ , draw an edge from  $T_i$  to  $T_j$  in the graph.

5. The schedule is serializable if there is no cycle in the precedence graph.

Example: If  $\text{read}(x)$  is opto up -  
st mode is to share patients. write in

1. S:  $r_1(x) r_1(y) w_2(x) w_1(x); r_2(y)$



2. S:  $r_1(x); r_3(y) w_1(x), w_2(y) r_3(x) w_2(x)$

• Schedule S:  $T_1, T_2, T_3$  is not serializable

because  $r_1(x)$  waits for  $w_1(x)$  before  $r_3(x)$

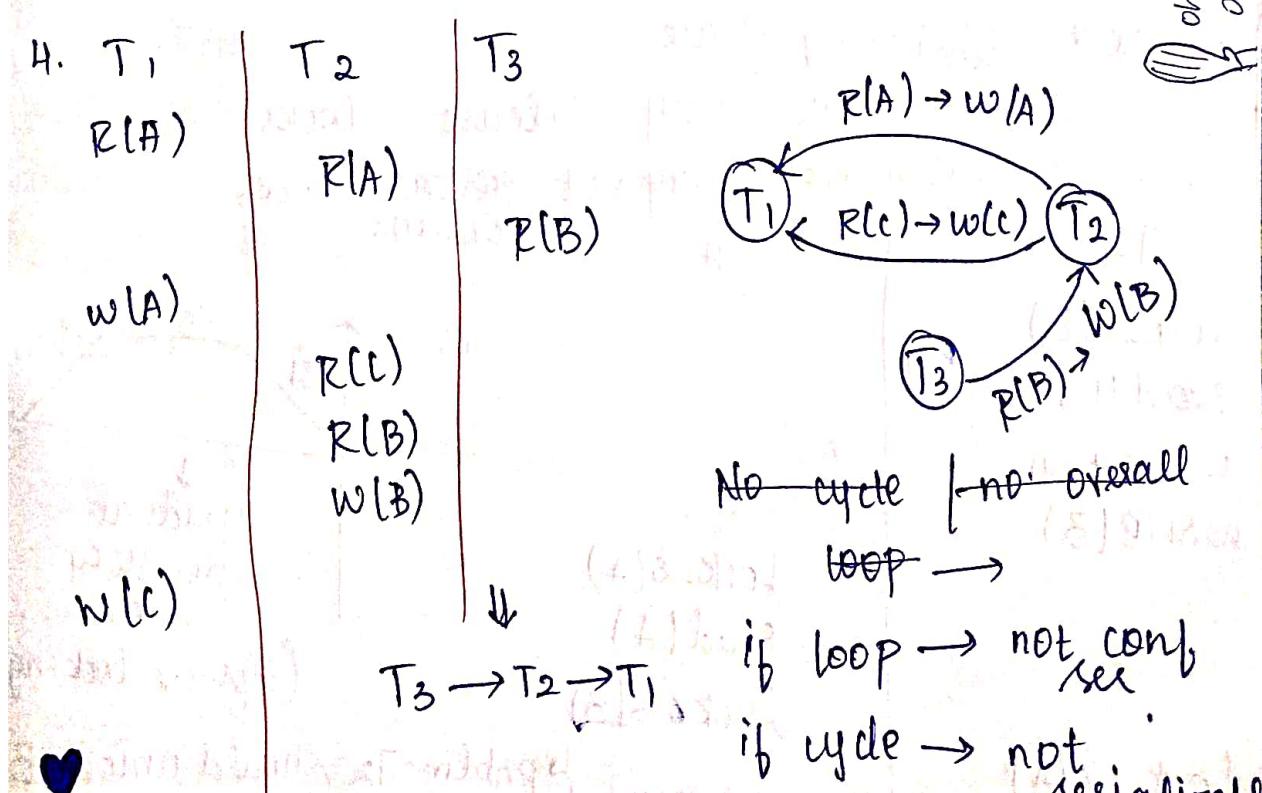
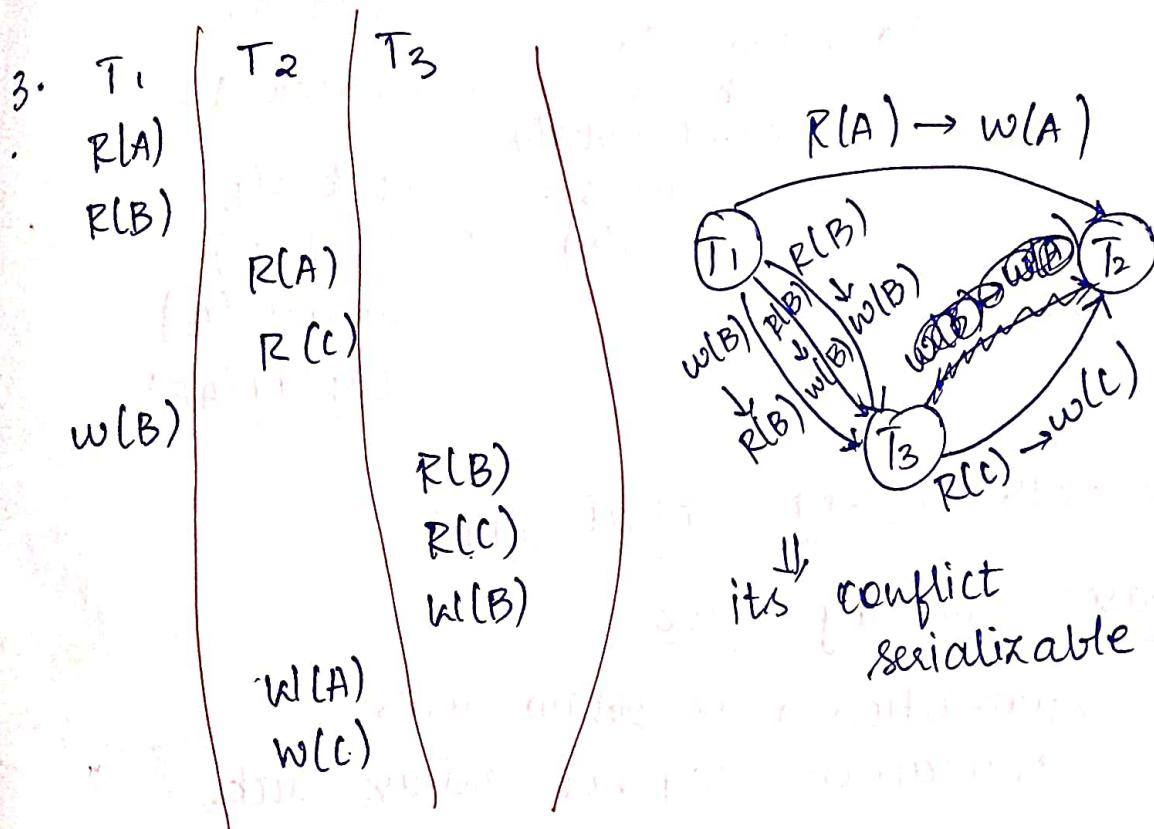
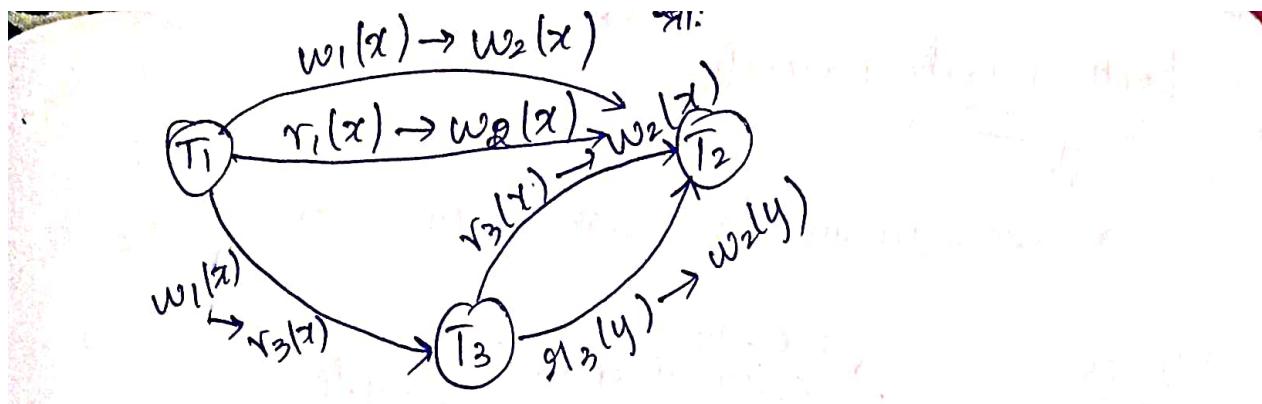
and  $w_1(x)$  waits for  $r_3(x)$  before  $w_2(x)$

and  $w_2(x)$  waits for  $r_3(x)$  before  $w_2(x)$

so  $r_1(x)$  waits for  $w_1(x)$  before  $r_3(x)$

and  $w_1(x)$  waits for  $r_3(x)$  before  $w_2(x)$

and  $w_2(x)$  waits for  $r_3(x)$  before  $w_2(x)$



## Lock based Protocols:

a mechanism to access data item

### Modes:

- shared lock even tho its locked other users can access (ex: read)
- exclusive locks unless the lock is released nobody can access (ex: write)

### Example:

Lock s(A)  
read (A)  
unlock (A)  
  
Lock-s(B)  
read (B)  
unlock (B)  
display (AB),

### Two phase locking protocol (2PL)

#### Phase 1: growing phase

- transaction may obtain locks
- transaction may not release locks

#### Phase 2: shrinking phase

- transaction may release locks
- transaction may not release locks

lock-x(B)

read (B)

B := B + 50

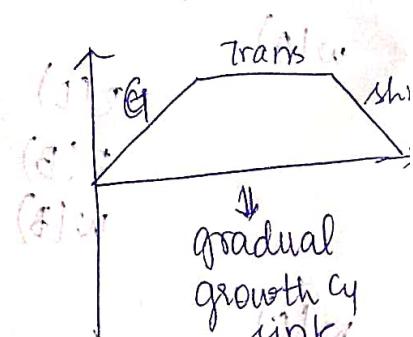
write (B)

Lock-s(A)

read (A)

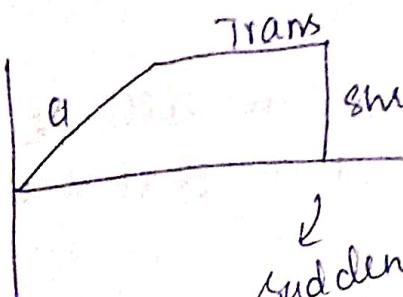
Lock-s(B)

lock-x(A)



(2 phase locking

problem: T<sub>3</sub> should unlock B



$\rightarrow$  (strict 2PL)

suddenly  
releases  
all the locks

the problem in  $T_3$  &  $T_4 \rightarrow$  deadlock

One transaction  
waiting for another  
transaction to release  
a lock on an  
item.

(Prevention is disad<sup>r</sup> of deadlock)

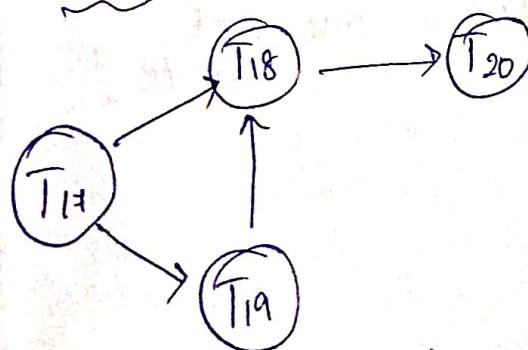
Deadlock handling:

deadlock prevention:

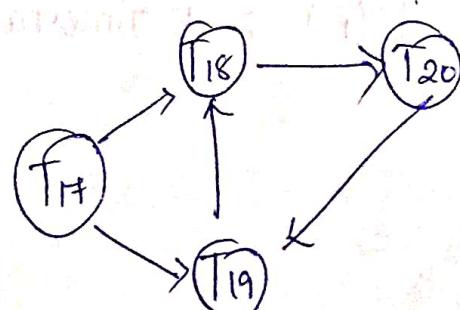
$\rightarrow$  requires each transaction to lock all its items before it begins executions.

$\rightarrow$  impose partial ordering of all data items and require that a transaction can lock data items only in the order specified by partial order (graph based).

deadlock detection algorithm:



wait for graph w/o cycle



wait for graph  
with a cycle

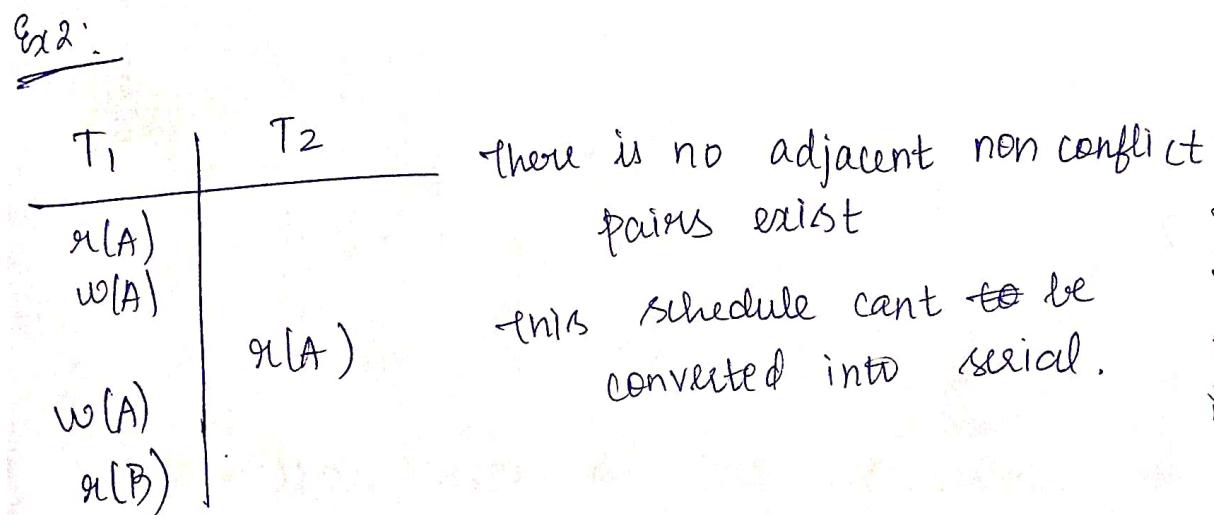
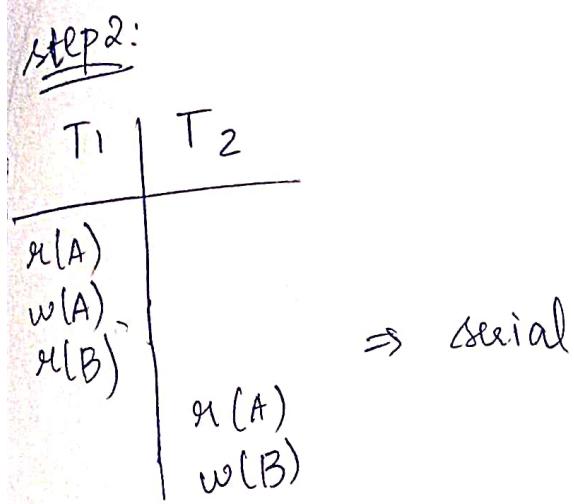
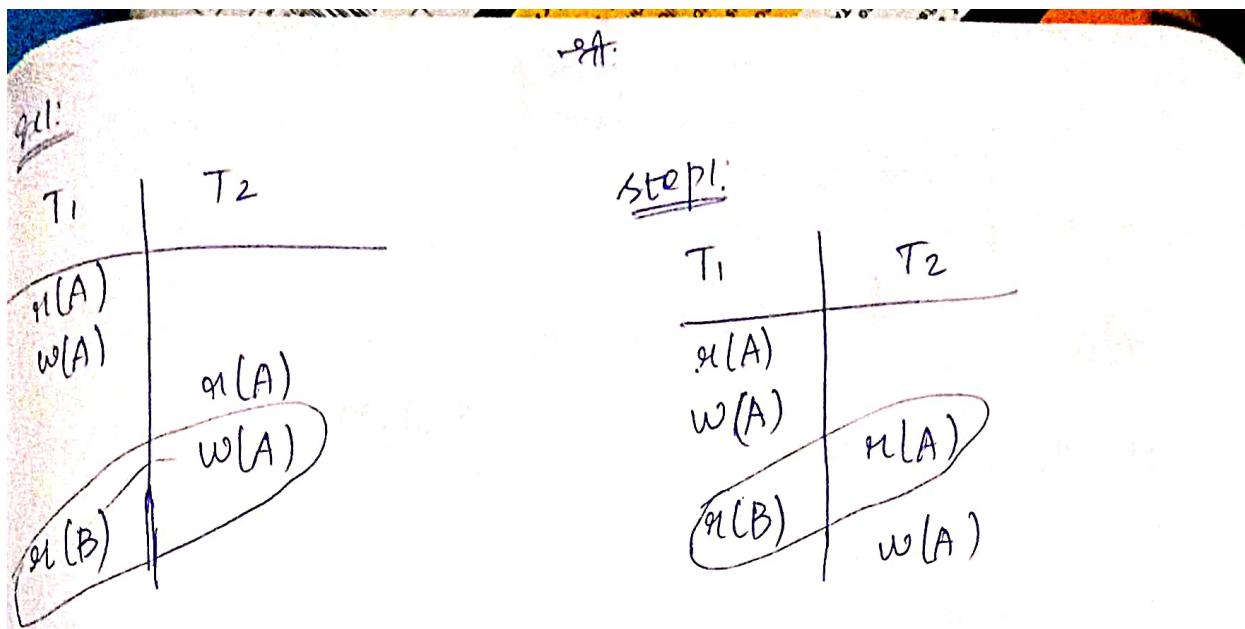
## left:

### Deadlock recovery:

- Some transactions will have to rollback to break deadlock. Select the victim that will incur the minimum cost.
- Rollback - determine how far to rollback transaction
- Total rollback: Abort the transaction and then restart it.

### Timestamp based protocol:

- Each transaction is issued a timestamp when it enters the ~~timestamp~~ system. If an old transaction  $T_i$  has timestamp  $TS(T_i)$  a new transaction  $T_j$  is assigned timestamp  $TS(T_j)$  such that  $TS(T_i) < TS(T_j)$



eff:

## Data Warehousing and Data mining:

⇒ Decision support systems: used to make business decisions often based on data collected by on-line transaction - processing systems.

⇒ Data analysis, Statistical analysis, Data mining

↓

discover meaning,

knowledge

automatically  
from data

in the form of  
statistical rules or  
patterns

⇒ Data Warehouse: is a repository or an archive of information gathered from multiple sources, stored under a unified schema.

⇒ Imp. for large businesses that have various data inputs and outputs

⇒ Greatly simplifies querying,

⇒ permits study of historical trends

⇒ shifts decision support query load away from transaction processing systems.

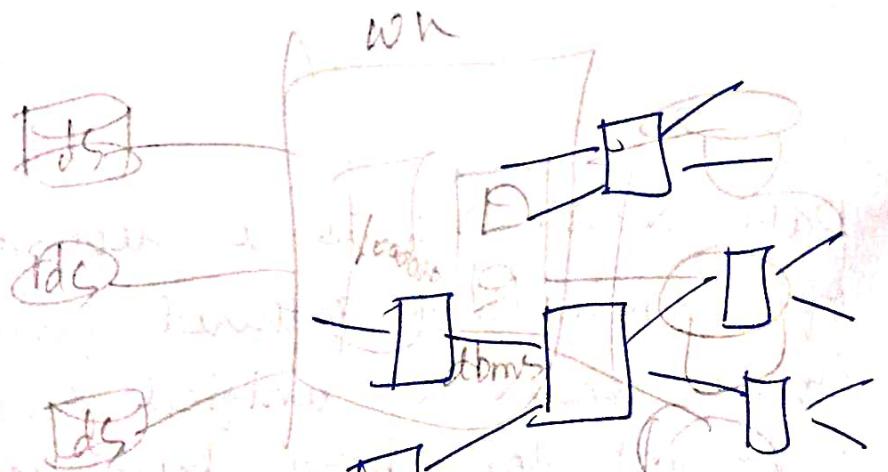
## Design Issues:

left:

- When and how to gather data
  - \* Source driven architecture (source to wh period)
  - \* Destination " " (copy from wh to source)
- keeping wh exactly synchronised with up-to-date data is costly.
- What schema to use
- Data cleansing
- What data to summarize

## WH Schemas:

- star schema
- snowflake "
- constellation "



Data mining: process of semi automatically analyzing large databases to find useful patterns.

- predictions based on past history
  - \* classification
  - \* Regression

- Descriptive patterns

- \* Associations (causations)
- \* Clusters

Classification rules:

- Summarised using a decision tree

Decision Tree:

- each internal node partitions the data into groups based on a partitioning attributes by a condition for the node
- each leaf node:
  - \* all items @ the node belongs to the same class
  - \* all attributes have been considered by no further partitioning is possible.

Best Split:

$$\text{Gini}(S) = 1 - \sum_{i=1}^K p_i^2$$

no. of classes =  $K$   
no. of instances =  $|S|$

fraction of  
inst in class

- When all instances are in same class  $= 0$
- reaches maximum if each class has the same no. of instances.

$$\text{Entropy}(S) = - \sum_{i=1}^K p_i \log p_i$$

When a set  $S$  is split into multiple sets  $S_1, i=1, 2, \dots, r$  we can measure the purity of the resultant set of sets as:

$$\text{purity}(S_1, S_2, \dots, S_r) = \sum_{i=1}^r \frac{|S_i|}{|S|} \text{purity}(S_i)$$

Information gain due to particular split of  $S$  into  $S_i, i=1, 2, \dots, r$

$$\text{Info gain} = \text{purity}(S) - \text{purity}(S_1, S_2, \dots, S_r)$$

[Information gain ratio =  $\frac{\text{Info gain}}{\text{Info content}}$ ]  $\rightarrow \uparrow \rightarrow$  best split

$$\text{Info content} = - \sum_{i=1}^r \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

## other types of classifiers:

→ Bayesian classifiers

→ Regression

→ clustering

clustering: → unsupervised learning

→ inter cluster  $\leftrightarrow$  intra cluster distance

→ nested clusters as a tree

→ hierarchical and partitional sets  $\rightarrow$  non overlapping subsets

→ Agglomerative  $\leftrightarrow$  divisive clustering

↓  
small cl

to big

↓  
big cl to

small

## Distributed Databases:

→ consists of loosely coupled sites that share no physical component

→ 2 or more database files @ diff locations in the network.

→ They run independently but consistently

$$1 - \sum_{i=1}^K p_i^2$$

loosely coupled sites

eff:

Homogeneous :-

- ⇒ All sites have identical s/w
- ⇒ Are aware of each other & agree to cooperate in processing user requests
- ⇒ Each site surrenders a part of its autonomy
- ⇒ appears as a single system to user

Heterogeneous :-

- ⇒ diff sites may use diff schemas & s/w
- ↓  
problem for  
query  
processing
- ↓  
problem for  
trans.  
processing
- ⇒ may not be aware of each other & may provide only limited facilities for cooperation.

Distributed Data Storage:⇒ Replication:

- maintains multiple copies of data, stored in diff sites, for fault tolerance

→ Adv:

- Availability
- Parallelism
- Reduced data transfer

dis:

- Increased cost of updates
- ↑ complexity of concurrency control

## Fragmentation:

- division of relation  $r_1$  into fragments  $r_1, r_2, \dots, r_n$  which contain suff information to reconstruct relation  $r_1$
- Horizontal  $\leftrightarrow$  vertical
  - ↓  
each tuple is assigned to one/more fragments
  - ↓  
split into several smaller schemas
  - all schemas must contain a common candidate key (or super key) to ensure lossless join property
  - a tuple-id attribute may be added to each schema to serve as cand. key

→ Adv:

Hori:

- \* PP on fragments of a relation
- \* tuples are located where they are most frequently accessed

diff b/w distributed & centralized

Vert:

- \* PP on a relation
- \* each part of tuple is located where they are most frequently used

dis: \* very high speed access sometimes

- \* recursive fragm → very expensive
- Availability
- parallelism
- reduced data transfr

## Web databases: p67

↳ access the database on the internet

Information Retrieval: locates relevant documents on the basis of user input such as keywords or example documents.

Ex: Web search engines

## Diffs from database systems:

- IR systems don't deal w/ transaction updates
- DBMS deals w/ structured data
- IRS deals with some querying issues not generally addressed by DBMS
  - \* Approx searching by keywords
  - \* Ranking of retrieved ans by estimated degree of relevance.

Keyword Search: done by first calculating relevance

- Term frequency
- Inverse document frequency
- Hyperlinks to documents

## Relevance Ranking:

- TF-IDF :  $n(d)$  → no. of terms in the document  
 $n(d,t)$  → no. of occurrences of term t in the document d

Relevance of a document  $d$  to a term  $t$ :

$$TF(d, t) = \log \left( 1 + \frac{n(d, t)}{n(d)} \right)$$

Relevance of a document to query  $Q$

$$r(d, Q) = \sum_{t \in Q} \frac{TF(d, t)}{n(t)}$$

- Words that occur in title, author list, section headings are gr imp.
- Words whose first occurrence is late in the document  $\rightarrow$  less imp.
- 'a', 'an', 'the', 'it'  $\rightarrow$  eliminated  $\rightarrow$  stop words
- Proximity: If keywords in query occur close together in the document,  $\rightarrow$  high imp

Similarity based Retrieval: retrieve documents similar to a given document

\* Relevance feedback: user selects a few documents from those retrieved by keyboard query  $\rightarrow$  returns similar docs

\* vector space model:

- vector for document  $d$  goes from origin to a point whose  $i$ th coordinate is  $TF(d, t)/n(t)$

- cosine of angle b/w vectors  $\rightarrow$  measure of similarity

## Relevance using Hyperlink: uphi!

- connections to social networking theories
  - Hub of authority based ranking
- $\downarrow$   
what      what

## PB security:

- ⇒ Types of security issues:
  - Legal and ethical issues
  - policy issues
  - system related issues
  - The need to identify multiple security levels

### ⇒ Threats to databases:

- Loss of integrity
- Loss of availability
- Loss of confidentiality.

### ⇒ Countermeasures:

- Access control → by creating user acc by pass
- Inference "
- flow "
- Encryption

### ⇒ Types of db security mechanisms:

- Discretionary "
- mandatory "

- Access Control:

- for restricting access to unauthorised users
- done by creating passwords by user accounts to manage logins

- Inference Control:

- controlling the access to statistical database that provide statistical data by summaries
- to safeguard such dbs we use inference control measures

- flow control:

- way to prevent information from flowing in such a way that it reaches unauthorised users.

→ channels that are pathways for info to flow implicitly in ways that violate the security policy of an organisation are called covert channels.

- data encryption:

→ used to protect sensitive data transmitted via some ~~type~~ comm. network

→ if data is encoded using some encoding algorithm,

- unauthorised users will find it difficult to decipher it.
- authorised users will be given decoding algorithms

DBA: • Database administrator:  
 → central authority for managing a db  
 \* granting privileges  
 \* classifying users

## Page ②

- DBA has a DBA account in the DBMS.  
 Called a system or Superuser acc
  - Acc creation → access control
  - Privilege granting → discretionary revocation
  - security level assignment → mandatory

- ① account creation
- ② system log
- ③ db audit

## Discretionary Access Control Based on Granting & Revoking Privileges

- Account level: privilege that each account holds independently of the relations in the db.
- relation level: privilege to access each individual relation or view in the db.

~~Eff:~~  
Types of privileges: (account level)

→ create schema / create table

→ create view

→ Alter

→ drop

→ modify

→ select

~~Pg 19~~

(relation model).



follows access matrix model

rows → subjects (users, acc, prgrms)

columns → objects (relation, record, col, views)

each position  $M(i,j)$  ↴

type of privilege

that i holds over j

Pg(20)