

11/11/19

PAGE NO.  
DATE.

→ We come to DBMS from file system because because it takes more memory and uses many resources which make other processes slow, so that other processes get enough space for its processing in main memory, we store data in secondary memory. This is the main reason we move to DB.

→ Traditional Database Model

→ ER model (Entity-Relationship model)

21/11/19

DB → Empty book (Raw data)

DBMS → Pen

(How to write, what will be its contents, etc)

→ Images, videos, etc form of data → multimedia DB  
→ GIS based data (maps, etc) → spatial DB (postgres)  
→ Real time data (~~not~~ efficient, but we need <sup>or no</sup> SQL in real time) → real time DB.

→ Create conceptual view :-

Before designing a database a model is constructed. E.g., before construction of a building -

- ① A chart is constructed
- ② Its real model

### ③ Conceptual Implementation.

In DBMS 3 steps of modelling is done.

- ① Conceptual view
- ② Representational view
- ③ Logical level or physical data models

#### ① Conceptual or high level view

To explain to naive users, here we used diagram for entities, attributes and relationship between entities.

Eg:- Student-course

Entities → noun

Relationship → verb

#### E-R model

→ It is the most popular & conceptual model because it is easy & clear to naive user & easy to understand.

#### ② Representational view (Implementation view)

→ Relational model

→ It is used by developers

→ It tells how SQL query runs or any query runs on programming language

Here we

→ Generally use tables.

② Logical level or physical data models.

→ It tells how the tables are stored in the memory device.

→ So deal about the structure of records i.e., how to store the data type, how to place the record, how many bytes are required to retrieve the data.

Assignment :-

① Create an E-R diagram for the cookie database.

② Name the entities and relationship along with their attributes using the ER notations.

③ Define the relationships between various entities (one-one, one-many, many-many).

By "entity" one should understand the table names in the cookie db (e.g., BIB FILE, PUBFILE, etc.)

④ Indicate minimum & maximum cardinalities for each entity.

⑤ A company db needs to store info about employees (ssn, salary, phone); departments (dno, budget, dname) and children of employees (name, age). Draw an ER diagram.

Set name → Entity type

Set elements (instance) → Entity

Employee → Entity type

Employee.name → instance (Entity)

Thing → Entity

Properties → attributes

Relationship → Association

① Difference b/w Entity and Entity type

→ Entity type is the schema which is going to represent the type of identity.

→ Entity is a person who is going to be represented in entity type itself.

② → Entity type → Heading

→ Entity → instance of that heading

③ → Entity type → Intension

→ Entity → Extension

④ Entity might get added or deleted

⑤ With respect to set theory,

Set name → Entity type

Set elements → Entity

25/1/19

PAGE NO.

DATE

### Attribute -

→ single valued & multivalued.



Age.



account number

(that can have  
more than  
one  
values)

phn. no.

→ simple v/s composite (which can further be divided).

↳ Age

↳ Gender

↳ Name

Phno.

Address

→ stored & derived



We have to  
store it in  
database



We can derive  
form the values  
stored in DB.

↳ Gender

↳ Age (from DOB).

Q Employee → 2 headquaters

Name → ~~Multivalue~~ Single, Composite, stored

Address → Multivalue, Composite, stored

Phno → Multi, Composite, stored

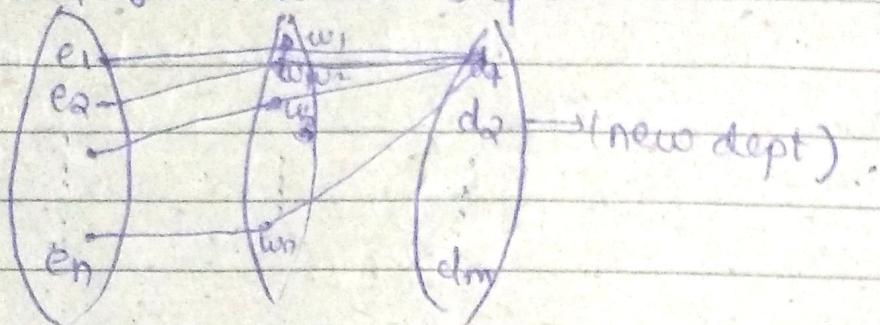
Age → Single, Simple, Derived

DOB → Single, composite, stored

## Relationships :

Q. Every employee works for a Dept. and every dept. can have many employees and new dept. need not have any employee.

Employee      works      Dept



Degree :- Q. (How many entities are there for that relation)

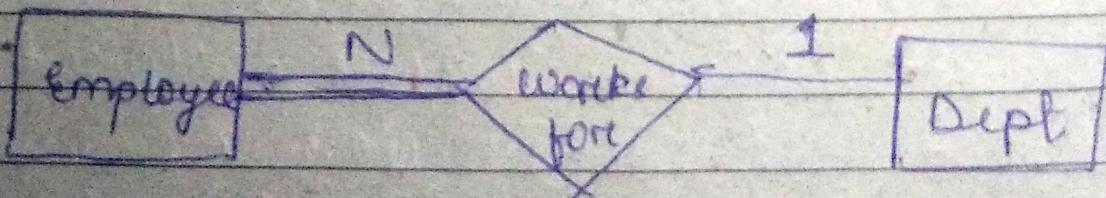
## Cardinality

Employee  $\textcircled{O} N$       { One to many }  
 Dept  $\rightarrow \textcircled{O} 1$       many to one

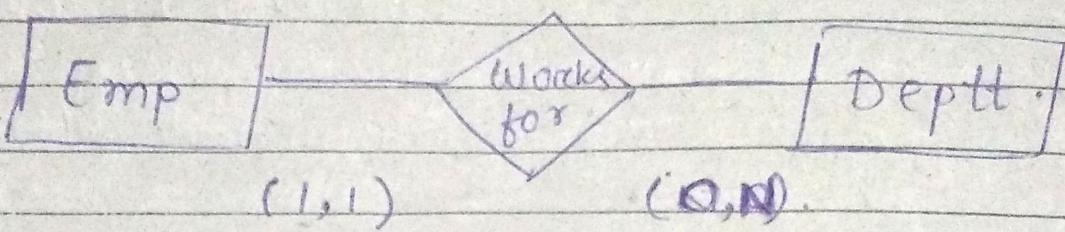
## Participation / Existence / Structural constraint

Employee  $\rightarrow$  total.

Dept  $\rightarrow$  partial.

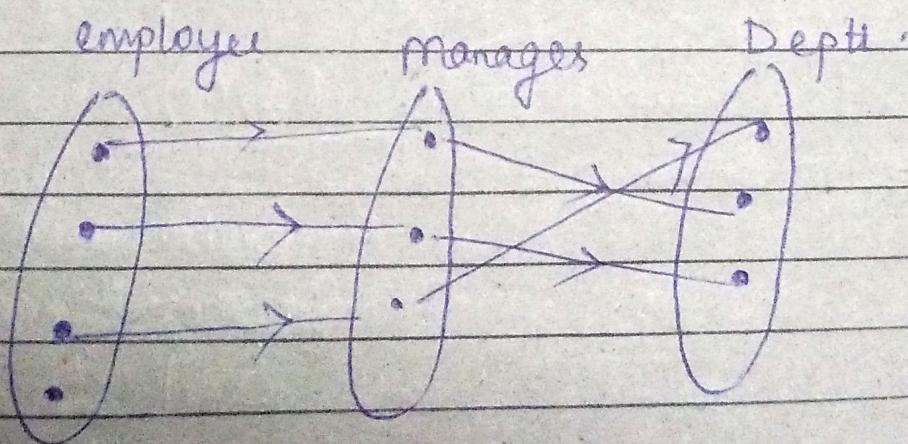
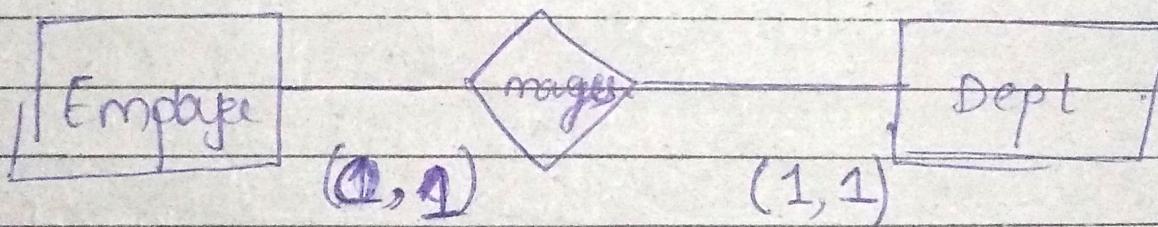


## Min-Max Representation



min max { first participation }  
and cardinality

- Every dept. should have a manager and only one employee manages a dept. and  $\geq 1$  employees can manage only one dept.

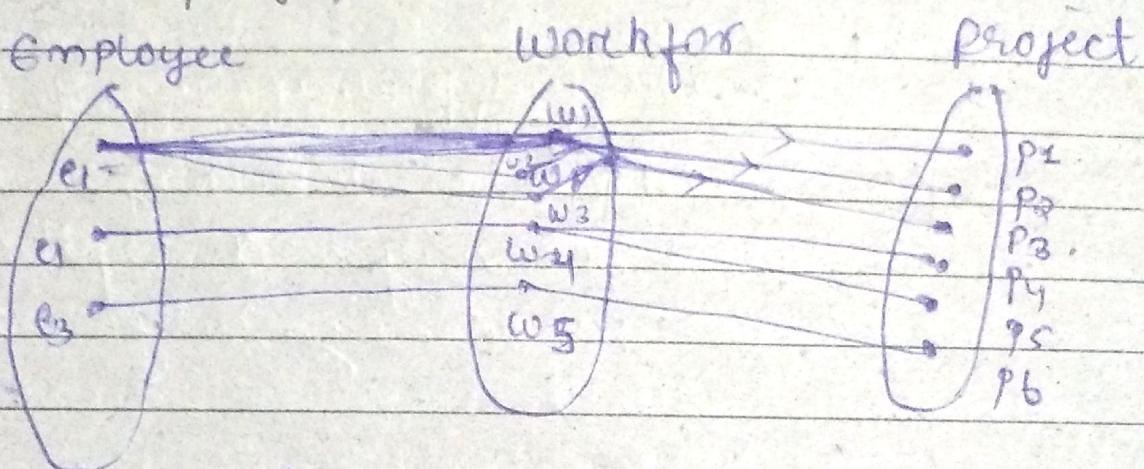


Degree = 2.

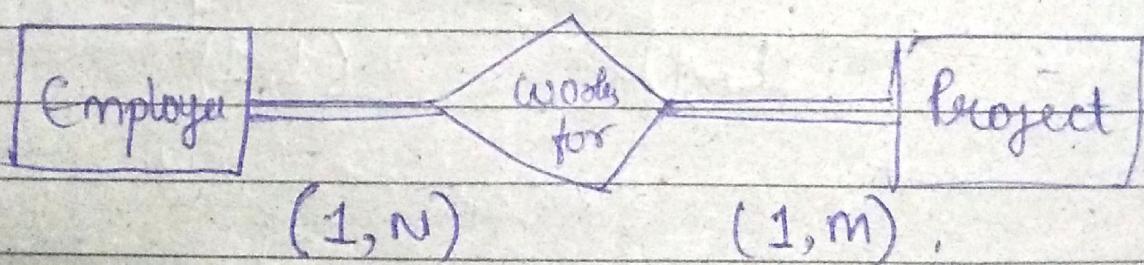
③ Every employee is supposed to work for at least 1 project.

An employee can work on many projects.

A project can have many employees and every project is supposed to have at least one employee.



Many to many

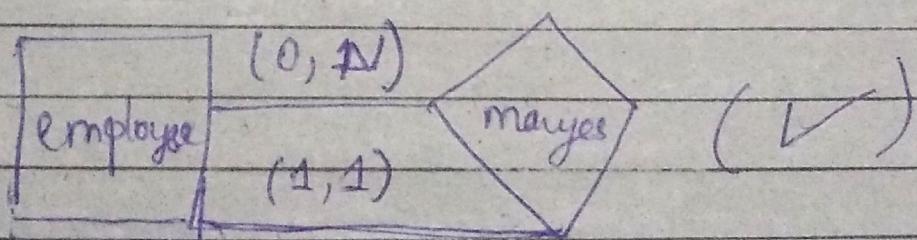
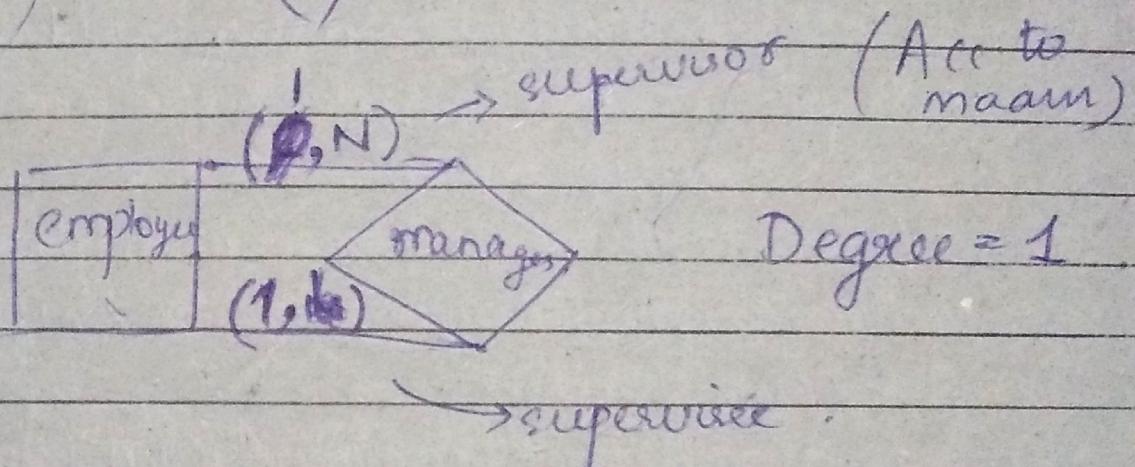
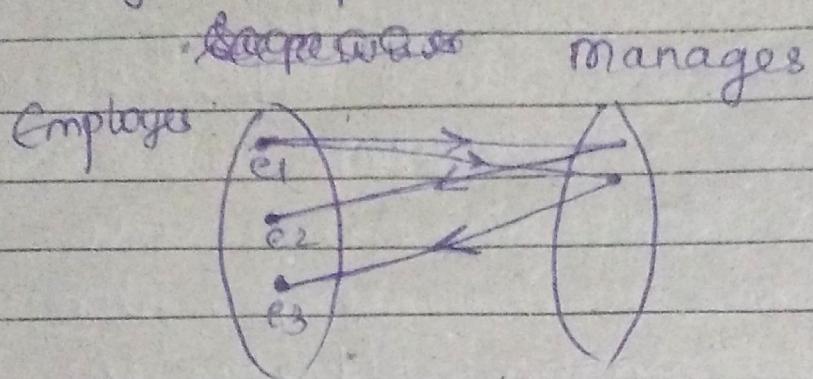


Degree = 2.

{ Both min-max  
 } relationship?  
 and structural  
 representation

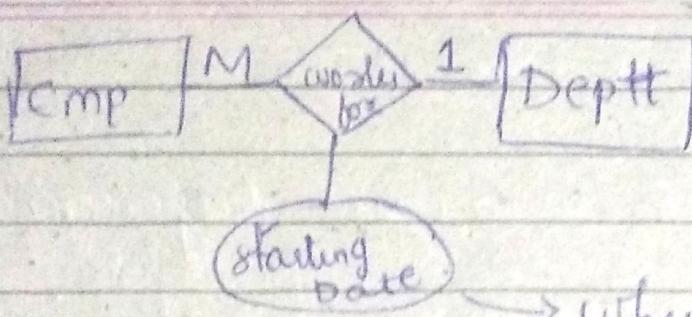
## Recursive Relation

~~Employee~~ Supervisor manages Supervisor.



28/11/19

Whenever we have many to one relationship, we will shift the descriptive attributes towards the many side.



→ where to shift this descriptive attribute

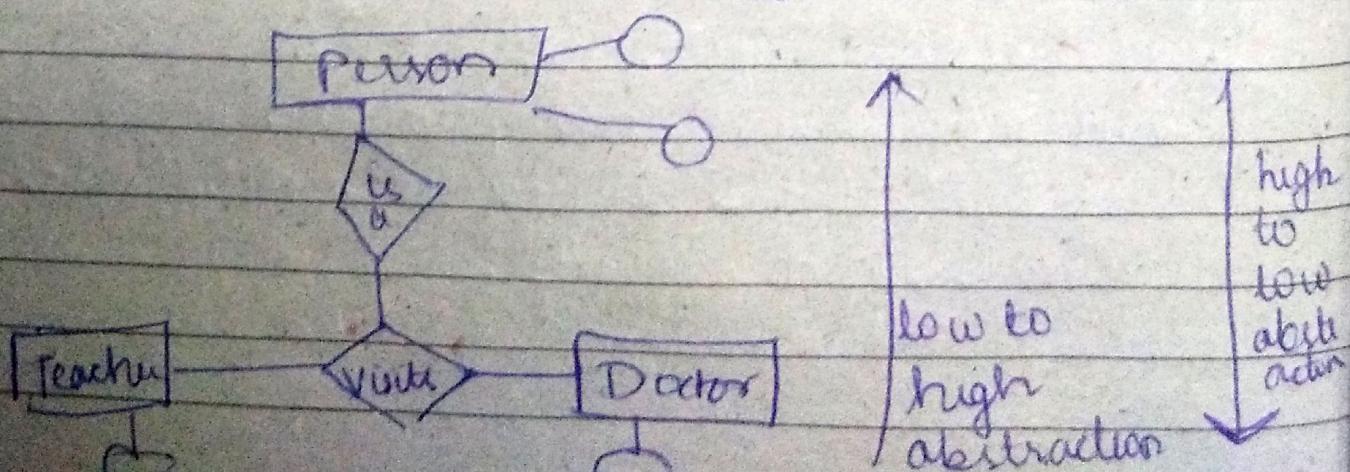
Then we will shift towards many side.  
If we have one-to-one relationship, then we can shift the descriptive attributes towards any side.

If we have many-to-many relationship, then we can't decide which side to shift.

For this we have;

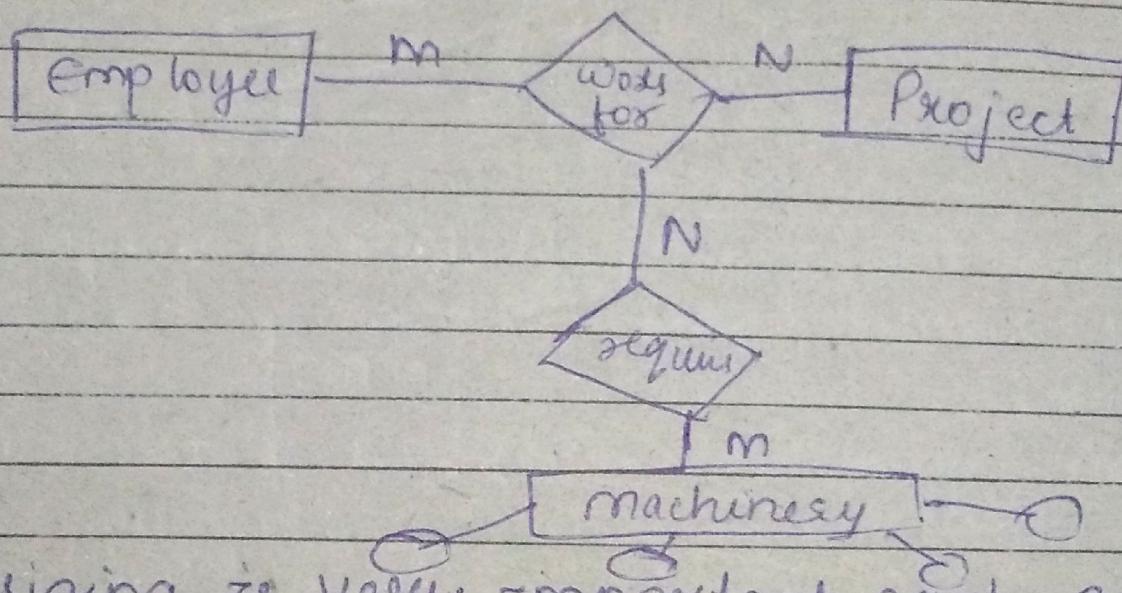
- ① Generalization - (moving from lower abstraction level to high abstraction level)
- ② Specialization
- ③ Aggregation

Specialization - moving from higher abstraction level to lower abstraction level.



## Aggregation

→ To understand better, we further divide relation into attributes and entities, to understand the detail better.



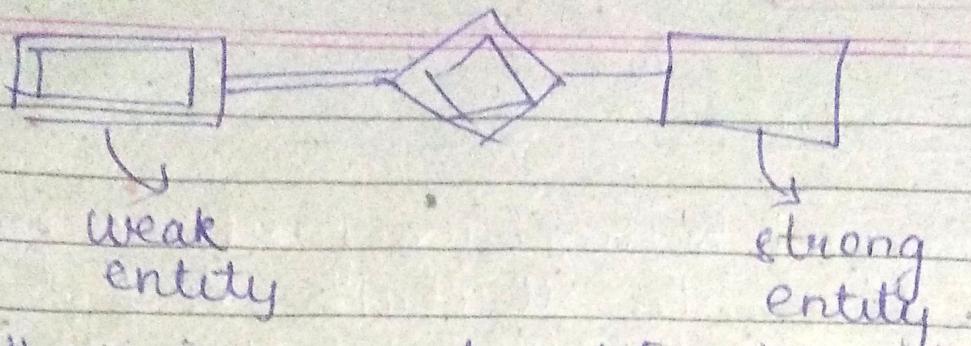
→ Designing is very important part so ER diagram is necessary for this.

→ First we choose database models and then we choose ~~the~~ DBMS (software).

Relational model is chosen most in the market because it is standardized.

so for RDBMS, we have many DBMS.

→ If our model is standardized, then we need not have to train our employee (as if technology changes, model remain same).



- If there is a weak entity, then there will be a total participation ~~see~~ with relation
- But vice-versa is not true.

### Various database models.

- ① Relational model
- ② Object-oriented model
- ③ Hierarchical model
- ④ Network model
- ⑤ ER model.

### Properties of relational model

- ① Tabular form
- ② Rows → tuple      columns - attribute
- ③ Use SQL queries to retrieve, access, change, data
- ④ Composite attributes are not allowed
- ⑤ ~~conceptual~~ implementation level view

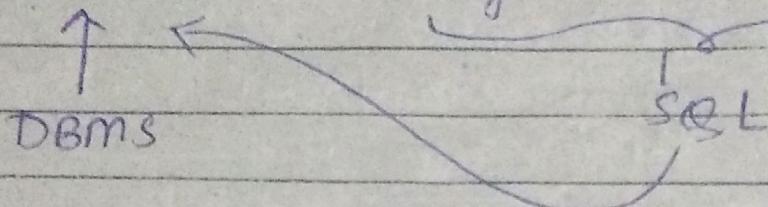
During design time also we can set constraints, domain constraints, key constraints, authorization " ", referential integrity assertion " ", constraint

PAGE NO.  
DATE

Relational database can be commercially used.

## Relational Model

RDB + Rel algebra & calculus.



→ Relation - Table

↔ Tuples | Rows | Records & Attributes |  
fields/cols.

Domain → All attributes are defined over a specific domain.

Relational schema

$$R \rightarrow (A_1, A_2, A_3, \dots, A_n)$$
$$\quad \quad \quad \downarrow \quad \downarrow \quad \downarrow \quad \quad \quad \downarrow \\ D_1 \quad D_2 \quad D_3 \quad \dots \quad D_n$$

Size = Cross product of all the domains  
 $D_1 \times D_2 \times D_3 \times D_4 \times \dots \times D_n$

Constraints → Integrity constraint

Entity  
Integrity  
Constraint

(Primary Key)

Referential  
Integrity  
constraint

(Foreign key)

" Every record is identified by a unique identifier (number)

89/1/19

PAGE NO.

DATE

## Constraints

- Domain constraints : (constraint on domain of attributes)
- Key constraints
- Entity constraints (In the primary key, value should not be null)
- Referential Integrity Constraints (reference should be integral)

## Key constraints

→ There are some attributes with key associated to it.

Super key :- set of all attributes that uniquely define a tuple

Candidate key :- ~~not~~ minimum set of attributes that uniquely define a tuple  
↳ Minimal super key

If we get a key by extraction to smaller units → minimal

If we choose minimum count then it is minimum key

e.g. Sk: A<sub>1</sub> A<sub>2</sub>, A<sub>1</sub> A<sub>3</sub> A<sub>4</sub>

Two superkeys  
minimal → both

minimum → A<sub>1</sub>, A<sub>2</sub>

Semantic constraints during

→ which you apply ~~after~~ the implementation of the database.

→ You cannot apply during the designing of the database.

eg:- Salary should be greater than 1 lakh.

### Violations

→ Insertion

→ Update ( Del + insertion )

→ Delete

→ Insertion can violate domain constraint

( like instead of 'int', it inserts varchar ).

→ Key constraint → Should be unique while inserting.

→ Entity constraint → Should be unique while inserting.

→ Referential Integrity → Should be there in the referred table while insertion.

### Delete

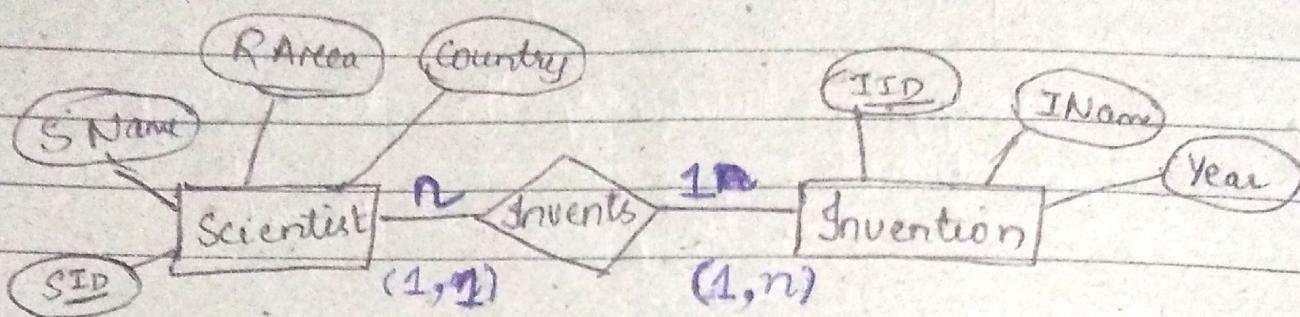
DC → No effect

KC → "

EC → "

RIC → { If we delete value in referred table,  
then we should delete or set null in refer table }

Q. Convert the given ER diagram :-



Include IID of invention as the foreign key of scientist

Scientist (SID, SName, RArea, Country, IID)  
Invention ( IID, IName, Year).

(OR)

Include SID of scientist as foreign key of invention.

Here considering SName as composite attribute  
and RArea as multi

① Many-to-many relationship Considering

- ① Scientist → strong entity
- ② Invention → " " "

Scientist (SID, SName, RArea, Country), IID

Invention ( IID, IName, Year)

Considering many to many relationship.

Scientist (sid, SName, RArea, country)

Invention (IID, IName, Year)

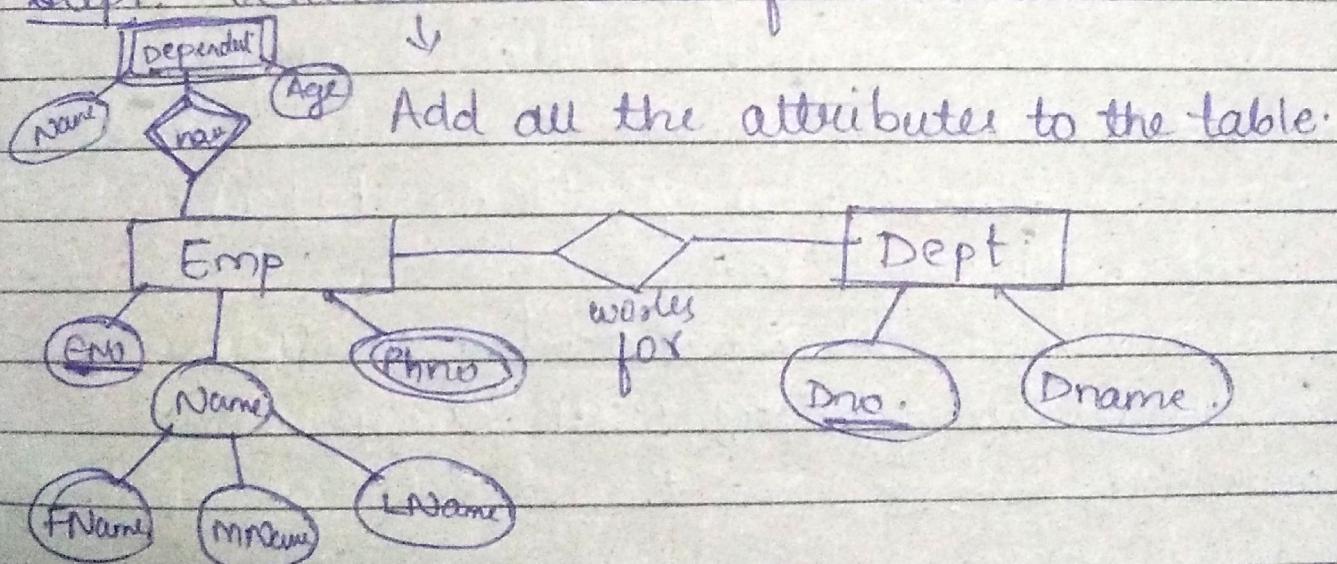
Inventi (IID, sid)

? This is because we are considering many to many relationship.

01/01/19

Converting the ER Diagram/Model to Relational Schema.

Step1:- Create the table for all entities



<u>eno</u>	<u>fname</u>	<u>lname</u>	<u>mname</u>	<u>dno</u>	<u>dname</u>

## Step 2 :- Dealing with weak entities.

foreign key      Dependents

ENo	Name	Relation

Here name is the partial key.

## Step 3 :- Converting the relationships to table;

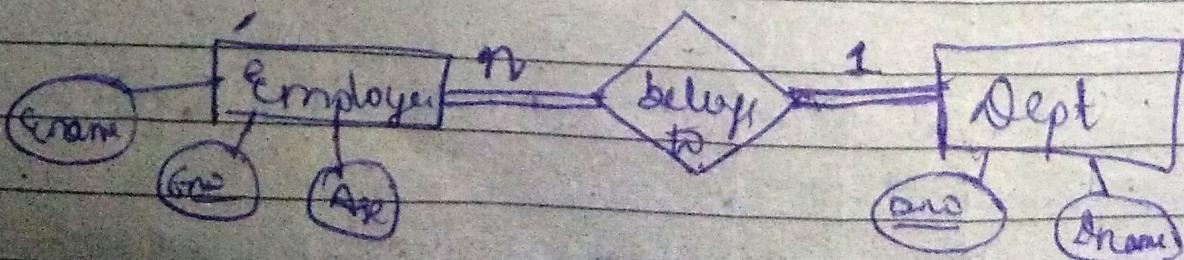
(i) Let's talk about 1:1 relationship.

Eg:- Only one employee can manage <sup>only one</sup> deptt. and each deptt is managed by an employee.

→ Add foreign key towards the total participation side.

→ There is no problem to add towards other side but there will be space wastage.

(ii) 1:m. relationship (One to many relationship)



employee

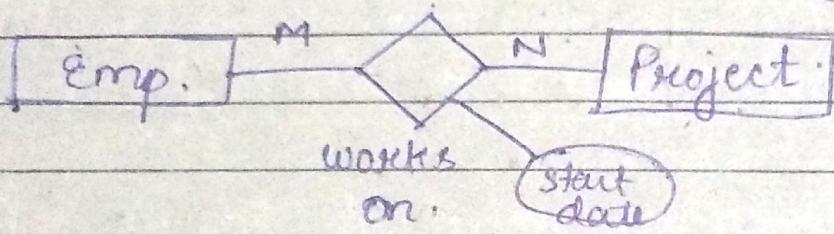
Eno	Ename	Age	Dno

Department

Dno	Dname

We will add foreign key towards the employee side or we can say many side.

### (iii) M:N relationship (many to many relationship)



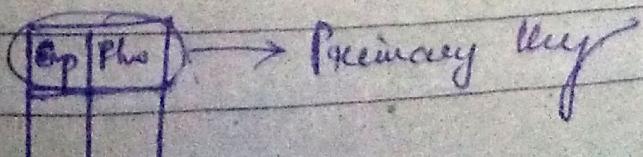
Here we will make a separate table with primary key as Empno and Project no combined.

Eno	Dno	Startdt

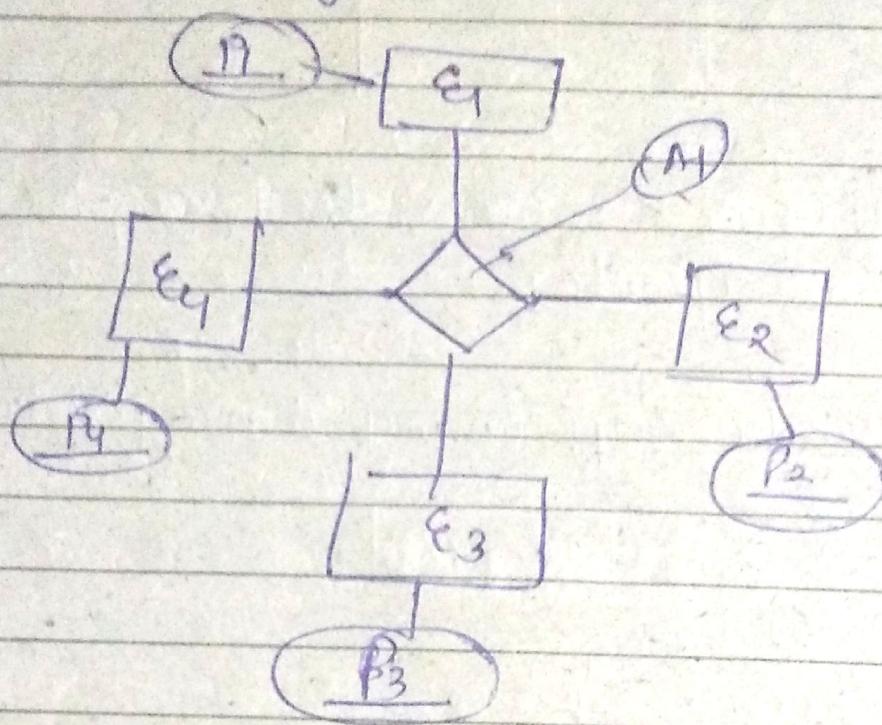
The attributes of the relation will be added to that table.

### Step 4: Multivalued attribute

We will another table for multi-valued attribute

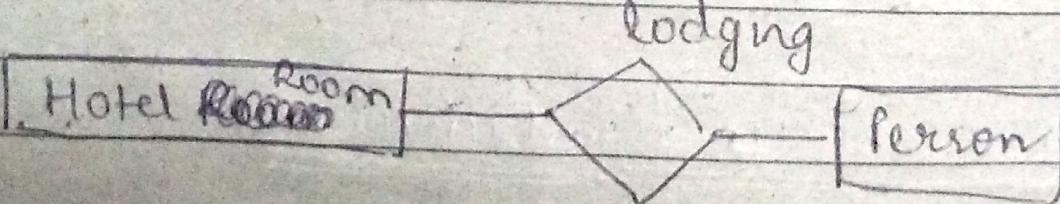


## Step 5: m-many relationships



We will make a separate table with  $P_1, P_2, P_3$  and  $P_4$  and also attribute of that junction.

$P_1$	$P_2$	$P_3$	$P_4$	$A_1$



Lodging is a m:m relationship.

Rent is a payment to be made by a person occupying different hotel rooms,

should be added as an attribute to

- ① Hotel Room
- ② Lodging
- ③ Person
- ④ Null

~~01/02/19~~

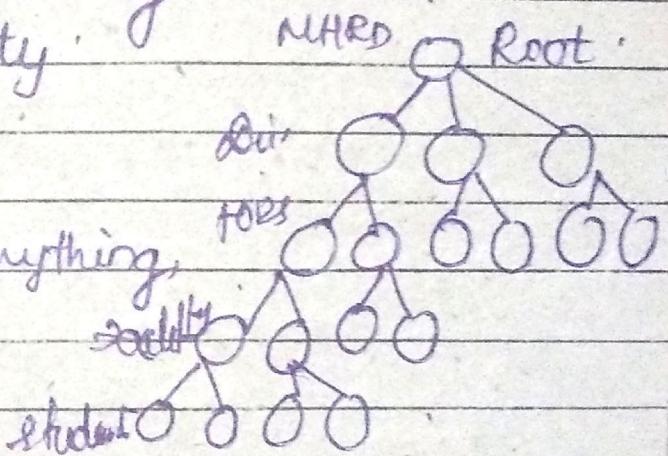
### Network Hierarchical Model

→ Data have been classified <sup>in some kind of</sup> hierarchy system.  
 → This is do because when data becomes large it is very difficult to manage.

→ Adv → Data integrity.  
 → Search is easier.

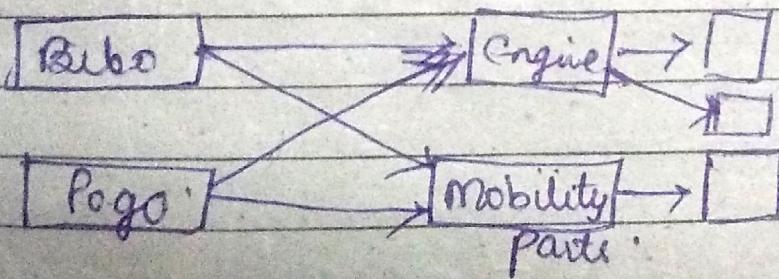
Disadv → Less flexible.

→ We cannot add anything anywhere



### Network model

owner



→ Peer to peer relationship.

Adv → Flexibility got increase

Disadv → It got complicated.

→ Slow down the database.

## Object oriented model

- Extensions to ER model
- Entities and relationships
- Similar kind of objects are put in one class.

~~3/8/19~~

Varchar2 → range is larger than  
varchar

Check → check values that should be in list

Alter → change in schema, we can say  
(not the tuple value)

Restruct → if we are deleting column,  
then others should not be affected  
that's why we use struct

## Auto increment

→ Max value = 10000

→ There should not be cycle,  
after value reaches 10000 the  
value automatically should  
not be incremented to 1

PAGE NO.  
DATE

Data dictionary :- data about data

- All kind of info is stored in data dictionary because schema is entered in data dictionary.
- This is stored in hard disk.

## Aliasing

Select clst.CustomerName As Name from

Customer As Clst

<sup>↳ Alias</sup>  
for table name

<sup>↳ Alias</sup>  
for column  
name

→ AND should be evaluated before OR in any boolean expression.

CHECK  $\implies$  See on your own

when  
Customer IN 'A' OR Customer IN 'B' OR -----

↓

Can be written as:-

when

Customer IN ('A', 'B', 'C', -----)

b

→ Scalar aggregate  
Vector aggregate

→ Views. (Read on ur own)

→ Can there be nested views?

→ Can we create a view from view

→ Can ~~any~~ view reside in database

→ Advantages & Disadv. of views.

→ If we make changes in view, will it be reflected in table?

— X —

08/02/19

## Views

→ To hide some details, we go for views

→ Views does not get saved in the memory, rather query which is processed is saved in memory

→ Materialised view → views which get saved in the memory in some applications  
It is get updated after some period of time.

→ We can create view from view.

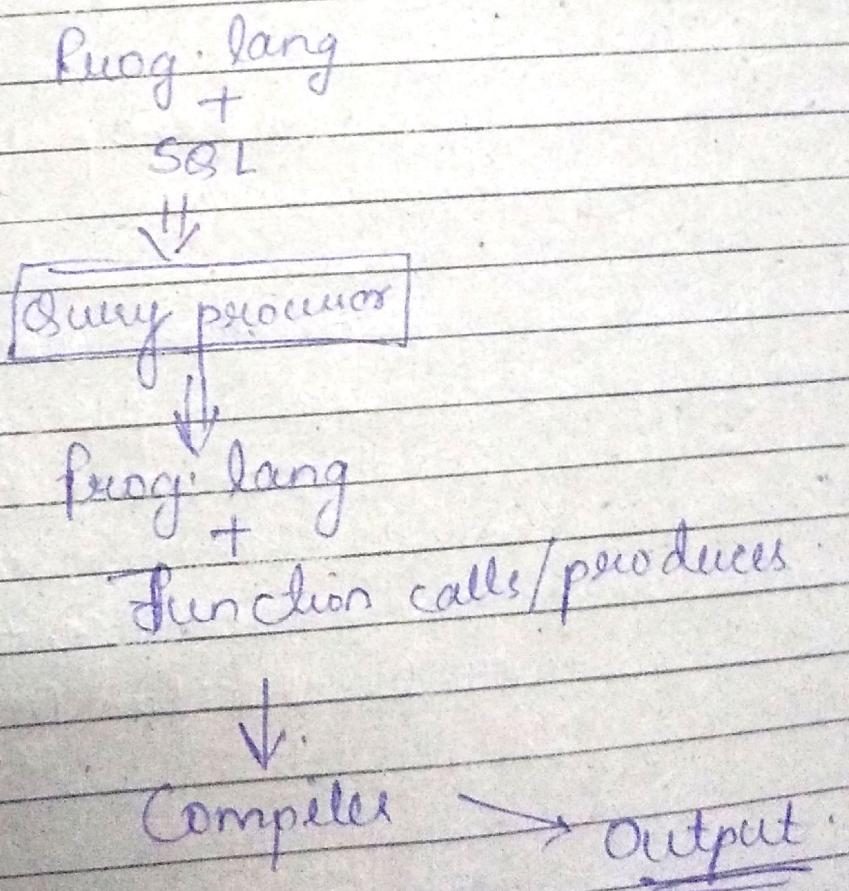
## Embedded SQL

- We apply SQL directly to database.
- General applications require many components :-

### Impediment Mismatch :-

Compiler of our lang. does not or is not able to compile or run our SQL language.

- We introduce our prog lang (C) + SQL input into a preprocessor.



## // Program

EXEC SQL <SQL Query>

→ to execute  
SQL query  
in b/w any  
prog. language

In database, if we want to execute a loop  
or get multiple tuples at a time, this is  
done by cursor

→ Helps to pick values from database  
and then get it to the prog. lang

→ Cursor is bit slower

→ Embedded SQL is an approach where language  
is extended to allow SQL in program blocks.

→ There can be a specialized compiler that  
handles the extended language or, there can be a  
preprocessor that converts  $\text{C} + \text{SQL}$  into  
plain C using call level interface.

$\text{C} + \text{SQL}$

↓

C

## Impediment mismatch

Procedural languages have a strict set of data types. They support mutable variables, loops, if statement, etc. They work with one data at a time, to process n data items it uses loops. Whereas, dB works on many data items at one time. It don't support if stmt or loops, etc so the challenge is to make 2 diff. program models to co-exist together.

- EXEC SQL BEGIN DECLARE SECTION  
-     int presentwith;  
-     char StudioName[50];  
-     char SQLSTATE[6];

↳ Special kind of variable that tells us about the SQL status.

If SQLSTATE returns:

00000 ← no error.

02000 ← No more results to display

- EXEC SBL END DECLARE SECTION

- `StoredProcedure(studioname, "xyz");`

- EXEC SBL SELECT network INTO :presnetwork

FROM

use to fetch  
value ~~as~~ of  
that variable

- FROM Studios JOIN MOVIES ON Presc = att#  
WHERE studios name = :studioname;

cursor

Cursors are used to return the collection  
of results.

Step 1 :-

→ Declare the cursor (i.e, associate the  
cursor to the query)

Step 2 :- Open the cursor. (Connecting the  
query to the DB) Initialization

Step 3 :- To fetch the rows.

Step 4 :- Close the cursor when finished.

SYNTAX:-

→ Declaring the cursor

- EXEC SQL DECLARE cursor-name CURSOR for Query;
- EXEC SQL OPEN Cursor-name;
- EXEC SQL FETCH FROM Cursor-name INTO :var1, :var2;
- EXEC SQL CLOSE Cursor-name;

Eg:-

```
EXEC SQL BEGIN DECLARE SECTION
int curt_number;
char title[100];
int year;
char SQLSTATE[6T];
```

EXEC SQL END DECLARE SECTION;

EXEC SQL DECLARE C-mais CURSOR FOR  
SELECT title, year FROM Movies where  
year = :curt-number ORDER by year.

→ During fetching also we can modify the data

→ first declare the shared variable and then execute the SQL Query in Embedded SQL

Eg:-

```
EXEC SQL OPEN C-movies;
for(;;) {
    EXEC SQL FETCH FROM C-movies INTO
        :title, :year ;
    if (!strcmp(SQLSTATE, "02000"))
        break;
    printf ("%d, %s\n", year, title);
}
EXEC SQL CLOSE C-movies;
```

Modifying the rows :-

11/2/19

PAGE NO.  
DATE

→ All the data are stored in the memory in the form of blocks.

→ File structure :-

→ Ordered way

→ Unordered way

→ we will choose one attribute and will do that sorting according to that attribute

Adv.

→ Insertion is efficient

Disadv.

→ linear searching,  
so it is costlier.

→ Insertion is difficult

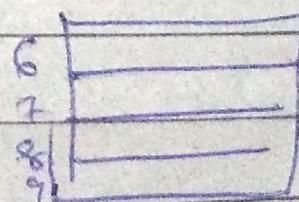
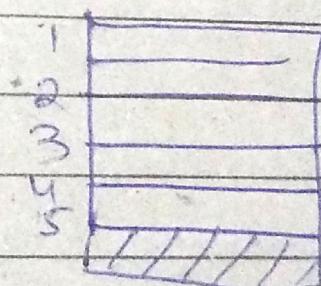
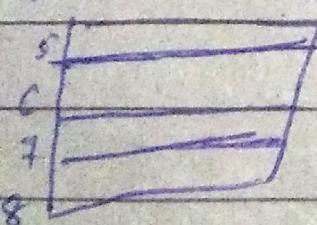
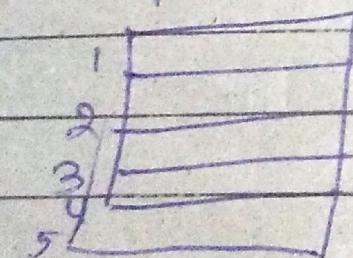
→ Deletion is easy

→ Searching is easy (binary search)

Data can be stored in 2 ways :-

① Spanned

② Unspanned



~~Space concern~~

Memory usage is efficient

→ Memory wastage  
is there

→ Fixed data uses unspanned strategy.  
 Variable length data uses spanned strategy  
 (like content of website)

→ Generally, we convert ~~our~~ variable data into fixed data.

→ In ordered way, complexity of searching is  $\log_2 B + 1$ .

→ Indexing      sparse (Only key value have pointer).

Dense . (Every record has a pointer record)

→ Indexed file must always be ordered and unique.

→ Indexes      → Single level indexed  
                   (Primary, Secondary, clustered)  
     → Multi-level indexing  
                   ( $B$ ,  $B^+$  trees)

→ Primary Index -

    → Ordered + primary key

    Clustered

    → Ordered + non-key

Secondary

    → Considering any key other than PK  
        $\geq \log_2 B + 1$

Q Suppose that we have an ordered file of 30,000 records on disks whereas having block size of 1024 bytes. The record strategy is fixed & unspanned of size 100 bytes. Suppose we have created primary index on key field on file size = 4 bytes and  $BP = 6$  bytes.

$$BF = \frac{1024}{100} = 10$$

$$\text{No. of blocks} = \left\lceil \frac{30,000}{1024} \right\rceil = \frac{30,000}{1024} = 29.375 \rightarrow \underline{\text{Ans}}$$

Without indexing searches =  $\log_2 3000 + 1$

$\rightarrow \underline{\text{Ans}}$

$$\begin{aligned} BF \text{ in index} &= \frac{1024}{15} = 68.2 \\ &= 68 \end{aligned}$$

$$\begin{array}{r} 3 \\ 68 ) 3000 \\ - 27 \\ \hline 32 \\ - 27 \\ \hline 5 \\ - 4 \\ \hline 1 \end{array}$$

$$\text{No. of blocks} = \left\lceil \frac{3000}{68} \right\rceil = \frac{3000}{68} = 45$$

$$\text{Searches} = \log_2 45 + 1 \rightarrow \underline{\text{Ans}}$$

Q Suppose there is an ordered file of 30,000 records,  $BS = 1024$  bytes.

$$RS = 100 \text{ bytes.}$$

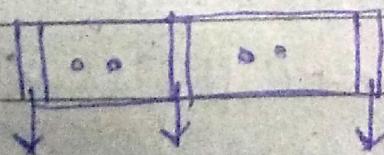
Secondary index on key field of file of size 8 bytes and  $BP = 6$  bytes.

## Multi-level indexing

- Flat file structure in the secondary disk.
  - This is done to reduce the access time and mapping time of the data from secondary memory.
  - B and B+ trees are dynamic in nature.
- B-trees (characteristics)
- In B trees at every level we are going to have key and data pointers and that data pointer actually pointing to either block or record.

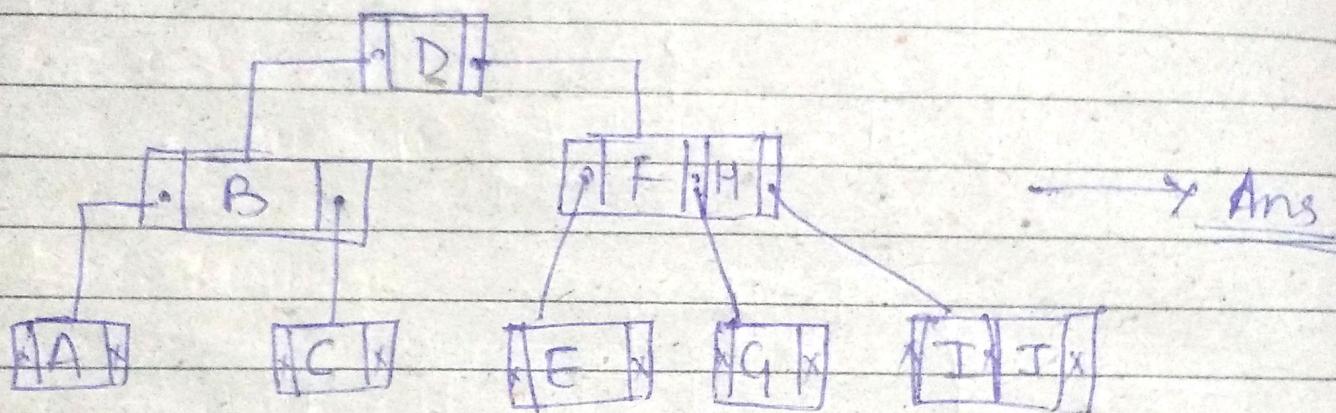
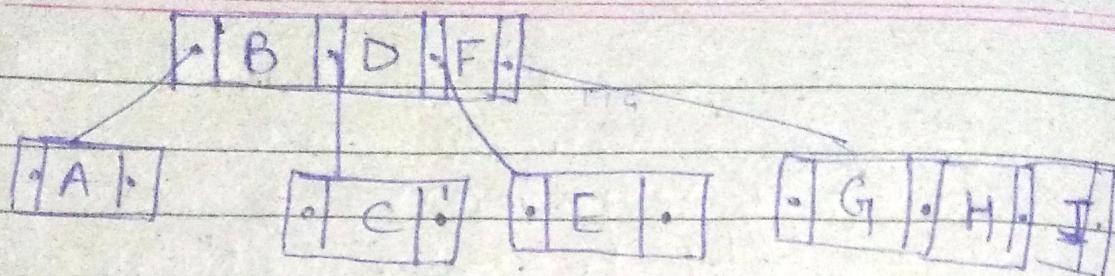
## Properties

- ① A root of Btree can have children between 2 and p, where p is called as order of the tree.
- Order of tree means maximum no' of children a node can have.



$\langle K_1, p_{11} \rangle$     $\langle K_2, p_{12} \rangle$

QP=4 ; ABCDEF<sub>G</sub>HIJ (9insert)



Overflow :- If there is no space for another key; then there is a problem of overflow.  
To overcome this, we split the nodes.

Underflow :- When order is less than desired order then there is a problem of underflow.

To overcome this we do borrow and merge.



## Deletion

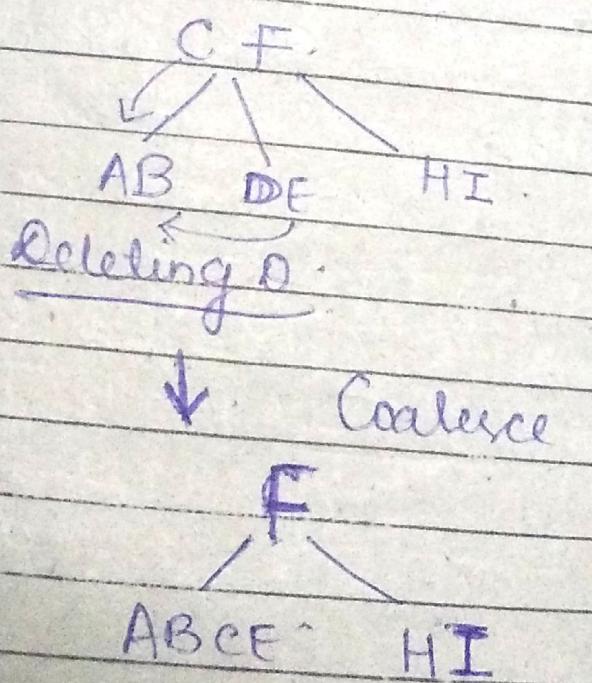
If underflow occurs then we have to adopt two methods.

- ① Borrow.
- ② Coalesce.

### Deletion:-

~~Borrow~~ Removing a key from leaf node leaving  $l$  keys in leaf node. If  $l \geq \text{min no. of nodes}$ , then we can stop else rebalance.

### Coalesce



→ The internal nodes in B+ trees, do not have record pointers unlike B trees.

## Page child node structure

$\langle k_1, p_{n_1} \rangle, \langle k_2, p_{n_2} \rangle \dots \dots \langle k_n, p_{n_n} \rangle \langle p_g \rangle$

— X —

next  
leaf  
node

12/12/19

## ISAM Organization :-

I : Indexed

S : Sequential

A : Access

M : Method

It was introduced by IBM. In this we maintain two levels of indexes. Indexing takes place (2 level) based upon primary indexing and ordered file.

In this case, insertion was costly and searching was easy.

This is static and also costly, so we use dynamic hashing which uses B and B<sup>+</sup> trees.

## Buffer Management :-

We have record that lies in hard disk when a command is to be processed, the pages are transferred from 2<sup>o</sup> memory to 1<sup>o</sup> memory. This work is done by buffer manager.

It divides 1<sup>o</sup> memory into section which is called as memory pool. In this, pages which are more frequently used are kept in buffer pool. When we require a page then that page is brought to the buffer pool.

- We maintain two flags / pointers →
  - count  $\rightarrow (+, -)$
  - dirty  $\rightarrow$  (when page is not in memory)
- When we want to address access small number of addresses from a large amount of address space to avoid wastage of memory, we use hashing.

### Methods of hashing :-

① Division modulo.

② Mid square       $1964 \rightarrow 96 \underline{6543} 6$

19	64	361	4096
		6	543
		65436	

③ Folding.

### Mid Square

Squaring the value that we want to access

$$\text{e.g. } 1964 = \underline{96} \underline{6543} 6$$

If we assume value to be stored as

000 - 999 i.e., 1000 values

i.e., 3 digit number

		000
1964	654	
	779	

Folding

Boundary  
folding

fold shifting

Ex. 123 456 789

if we want to store values in 3 digit  
number address

store at  $(123 + 789) = 912$

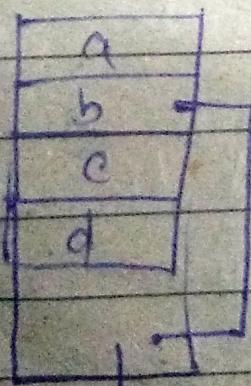
value at	912
	999

Collision Removal Method

① open addressing :-

linear  
probing

quadratic  
probing



insert  
here

→ linear probing

★ Here is a limit on the no. of attempts we can make to find free space.

Hashing → Internal (on Buffer)

External (on  $\mathbb{A}^n$  memory)

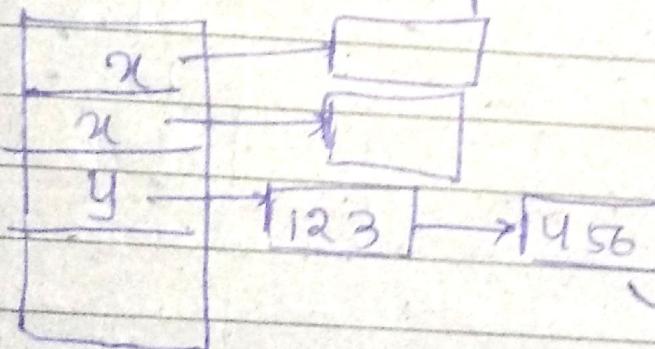
② Double Hashing

→ Apply hashing on the collision.

(Another Best method)

③ Chaining:

In this linked list is created when more than one value occurs at a particular address.



→ It can grow to any level.  
(Here it is 2).

load factor = no. of entries in one place (n)

total memory space (m)

It tells whether the entries are evenly distributed or not. Used to check the amt of collision, if collisions are more then we

Should change our method for hashing.

Example -

10, 25, 48, 20, 55, 69, 41, 99, 108

$$\underline{h(\text{key})} = \text{key mod } m$$

$m = (0-9)$  so total value  $\geq 10$

$$LP(\text{key}, i) = (\text{key} + i) \text{ mod } m$$

$$LP(10, 0) = 10 \text{ mod } 10 = 0$$

10	0
----	---

$$LP(25, 0) = 25 \text{ mod } 10 = 5$$

20	1
----	---

$$LP(48, 0) = 48 \text{ mod } 10 = 8$$

41	2
----	---

$$LP(20, 0) = 20 \text{ mod } 10 = 0$$

99	3
----	---

Already filled. ←

108	4
-----	---

$$LP(80, 1) = 80 \text{ mod } 10 = 1$$

25	5
----	---

55	6
----	---

$$LP(55, 1) = 55 \text{ mod } 10 = 5$$

7	7
---	---

48	8
----	---

$$LP(69, 0) = 69 \text{ mod } 10 = 9$$

69	9
----	---

$$LP(41, 0) = 4 (x) \quad LP(41, 1) = \cancel{0} = 2$$

$$LP(99, 0) = 9(x), 0(x), 1(x), 2(x), 3(x)$$

$$LP(108, 0) = 8(x)$$

$$9(x) 0(x) 1(x) 2(x) 3(x) 4(v)$$

## Q. Quadratic hashing

$$h(k,i) = (h(k) + c_1 i + c_2 i^2) \bmod M$$

$$M = 10, c_1 = 6, c_2 = 3$$

$$h(10,0) = 10 \bmod 10 = 0$$

$$h(25,0) = 5$$

$$h(20,0) = 0 \text{ (already filled)}$$

$$\begin{aligned} h(20,1) &= (20 + 6 + 3) \bmod 10 \\ &= 29 \bmod 10 \\ &= 9 \end{aligned}$$

$$h(55,0) = 5 \times$$

$$\begin{aligned} h(55,1) &= (55 + 6 + 3) \bmod 10 \\ &= 64 \bmod 10 = 4 \end{aligned}$$

$$h(69,0) = 0 \bmod 10 = 0$$

$$h(69,1) = (69 + 6 + 3) \bmod 10 = 78 \bmod 10 = 8 \times$$

$$\begin{aligned} h(69,2) &= (69 + 12 + 12) \bmod 10 \\ &\Rightarrow 93 \bmod 10 = 3 \end{aligned}$$

29  
293  
293  
293  
293  
293  
293  
293  
293  
293

$$h(41,0) = 1$$

$$h(99,0) = 0 \bmod 10 = 0 ; h(99,1) = (99 + 9) \bmod 10 = 8 \times$$

$$h(99,2) = (99 + 24) \bmod 10 = (123) \bmod 10 = 3 \times$$

$$h(99,3) = (99 + 18 + 27) = 144 \bmod 10 = 4 \times$$

## Q Double Hashing

28, 77, 54, 95, 16, 5, 22, 99, 61

$$h_1(k, i) = k \bmod m$$

$$h_2(k) = (1 + (\text{key} \bmod (m-1)))$$

$$DH(\text{key}, i) = (h_1(k, i) + i(h_2(\text{key}))) \bmod m$$

$$DH(95, 1) =$$

$$(5 + 1 \times h_2(\text{key})) \bmod 10$$

↓

$$\begin{aligned} & 1 + 95 \bmod 8 \\ & = 8 \end{aligned}$$

$$(5 + 8) \bmod 10 = 3$$

$$h_1(5, 0) = 5$$

$$h_2(5) = 6$$

$$DH(5, 1) = (5 + 6) \bmod 10 = 1$$

0	
5	1
22	2
95	3
54	4
25	5
16	6
77	7
99	8
	9

$$h_1(61, 0) = 1$$

$$h_2(61) = 6$$

$$DH(61, 1) = 7 \bmod 10 = 7 (X)$$

$$DH(61, 2) = (4 + 12) \bmod 10 = 3 (X)$$

$$DH(61, 3) = (1 + 18) \bmod 10 = 9 (X)$$

$$DH(61, 4) = (1 + 24) \bmod 10 = 5 (X)$$

$$DH(61, 5) = (1 + 30) \bmod 10 = 1 (X)$$

$$DH(61, 6) = (1 + 36) \bmod 10 = 7 (X)$$

$$DH(61,7) = (1+42) \bmod 10 = 3(X)$$

$$DH(61,8) = (1+48) \bmod 10 = 9(Y)$$

Cycle come, so this won't come at any position.

### Codd's Rule:-

① Information Rule :- everything in a database must be stored in a table format.

② Guaranteed Access Rule :-

Every single data element (value) is guaranteed to be accessible logically with a combination of table-name, primary-key (row value), and attribute-name (column name).

③ Systematic Treatment of NULL Values.

The NULL values in a database must be given a systematic and uniform treatment.

This is a very important rule because a NULL can be interpreted as one of the following - data is missing, data is not known, or data is not applicable.

④ Active Online Catalog

The structured description of the entire DB must

be stored in an online catalog, known as data dictionary.

### ⑤ Comprehensive Data sub-language Rule

A DB can only be accessed using a language having linear syntax that supports data def., data manipulation and transaction mgmt operations. If DB allows access to data without any help of this language, then it is considered as a violation.

### ⑥ View Updating Rule

All the views of a DB, which can theoretically be updated, must also be updatable by the system.

### ⑦ High-Level Insert, Update and Delete Rule

This must not be limited to a single row, that is, it must also support union, intersection and minus operations to yield sets of data records.

### ⑧ Physical Data Independence

Any change in the physical structure of the DB must not have any impact on how the data is being accessed by external applications.

### ⑨ Logical Data Independence

The logical data in a DB must be independent of its user's view (application).

### ⑩ Integrity Independence

A DB must be independent of the application that uses it. This rule makes a DB independent of the front-end application and its interface.

### ⑪ Distribution Independence

The end-user must not be able to see that the data is distributed over various locations.

### ② Non-subversion Rule

If a system has an interface that provides access to low-level records, then the interface must not be able to subvert the system and bypass security and integrity constraints.

### ③ The Foundation Rule

For a system that is advertised as, or claimed to be, a RDBMS, the system must be able to manage data bases entirely through its relational capabilities.

19/2/19

PAGE NO.  
DATE

Q) Indices speed query processing, but it is usually a bad idea to create indices on every attribute, and every combination of attributes that is a potential search keys. Explain why?

Ans - ① Every index requires additional CPU time and disk I/O overhead during inserts and deletions.

② Indices on non-primary keys might have to be changed on update, although an index on the primary key might not.

③ Each extra index requires additional storage space.

④ For queries which involve conditions on several search keys, efficiency might not be bad even if only some of the keys have indices on them. Therefore, DB performance is improved less by adding indices when many indices already exist.

Q) When it is preferable to use a dense index rather than a sparse index?

It is preferable to use a dense index instead of a sparse index when the file is not sorted on the indexed field (such as when the index is a secondary index) or when the index file

is small compared to the size of memory.

Q3: What is the difference b/w a clustering index and a secondary index.

Ans - The secondary index may be generated from a field which is a candidate key and has a unique value in every record, or a non-key with duplicate values.

Clustering index is defined on an ordered data file. The data file is ordered on a non-key field.

Q4. Construct a B+ tree for the following set of key values - 2, 3, 5, 7, 11, 17, 19, 23, 29, 31. Assume that the tree is initially empty and values are added in ascending order.

Construct B+ trees for the cases where the number of pointers that will fit in one node is as follows:-

1) 4

2) 6

3) 8

Q5. Explain the distinction between closed and open hashing. Discuss the merits of each technique in database application.

Q6. What are the causes of bucket overflow in a hash file organization? What can be done to reduce the occurrence of bucket overflow?

Q7. Why is the hash structure not the best choice for a search key on which range queries are likely?

Q8. Suppose a relation is stored in a B+ tree file organization. Suppose  $2^0$  indices stored record identifier that are pointer to record on disk, then

@ What would be the effect on the  $2^0$

Index of a node split happen in the file organization.

- What would be the cost of updating all affected record in a secondary index.
- How does using a search key of file organization as logical record identifier solve this prob.
- What's the extra cost due to the use of such logical record identifier?

### Extended Hashing or Dynamic Hashing

e.g. data

Kartik - 00000

Sneha - 00001

Sharam - 01010

Vagun - 10011

Imran - 01111

Ravi - 01100

Buckets are used to store the data. We can store the above data by taking last two or three, etc. The total number of buckets  $2^2$  (check bit) can be used.

→ It eliminates overflow chains by splitting a bucket when it overflow

→ Range of hash function has to be extended to

- accommodate additional bucket.
- Family of hash function based upon  $h_k(v) \text{ mod } 2^k$  (use last  $k$  bits of  $h(v)$ )
  - ↑  
search value
- At any given time a unique hash is used depending upon the number of times buckets have been split.
- Extendable hashing using a directory to accommodate family of hash functions.

### Transaction processing

A - Atomicity

C - Consistency

I - Isolation

D - Durability

Q Show how the extendable hash structures of the following figure changes as a result of each of the following steps :-

1) Delete 11

2) Delete 31

3) Insert 1

4) Insert 15

## Syllabus:-

① Overview of database models.

(only Introduction)

② Conceptual DB design

(General ER diagram concepts)

③ logical DB design

(Relational model (no calculus))

Converting CR to relational model

④ Physical DB design.

(Hashing, B BT tree (insertion & deletion))

⑤ Normalization

⑥ Transactions

(Properties, schedules, serializability,  
consistency, transaction performance)

Eg:- (if trans. is interleaving  
then what will be  
its effect on performance)

22/02/19

DATE \_\_\_\_\_

→ To improve the CPU utilization, we introduce the concept of concurrency

→ Problems →   
      → lost update (write some uncommitted trans.)  
      → dirty read (read some uncommitted trans.)  
      → unrepeatable read.

→  $\frac{(n+m)!}{n! m!}$  = Total schedules if  
      a trans have m operation  
      and other have n.

Total serial schedule =  $n!$

when n is no. of trans.

Q Find total number of schedule with  
 number of operation wret each trans are:-  
 2, 3 and 4.

$$\text{Total schedule} = \frac{(2+3+4)!}{2! \times 3! \times 4!}$$

$$= \frac{9!}{2! \times 3! \times 4!}$$

$$= 9 \times 4 \times 7 \times 5$$

$$\text{Total non-serial schedule} = \text{Total schedule} - \text{Total serial schedule}$$

Here : = 
$$\boxed{9 \times 4 \times 7 \times 5 - 3!}$$

$\textcircled{1} T_1 := \text{read}(P);$   
 $= \text{read}(Q);$   
 $\text{if } P=0 \text{ then } Q=Q+1.$   
 $\text{write}(Q).$

$T_2 := \text{read}(B)$   
 $\text{read}(P)$   
 $\text{if } Q=0 \text{ then } P=P+1$   
 $\text{write}(P).$

What does:

Any non-serial interleaving of  $T_1$  and  $T_2$  for concurrent execution leads to?

①  $T_1$                        $T_2$   
 read(P)  
 read(Q)  
 write(Q)                      read(Q)  
 read(P)  
 write(P).

②                              read B  
 read P.  
 read P  
 read Q  
 $(2+8) \times 2!$               write(P)  
 write(Q)

T<sub>1</sub>T<sub>2</sub>

(3) read P.

read Q

read Q

read P.

write Q

write P.

(4) T<sub>1</sub>  
read (P).

read (Q)

read (Q)

read (P)

write (Q)

write (P)

(5) T<sub>1</sub>  
read (P)T<sub>2</sub>  
read (Q)

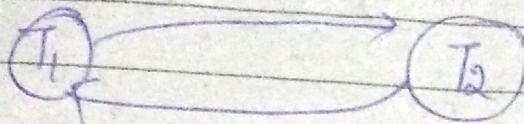
read (Q)

read (P)

write (Q)

write (P)

There can be many more non-serial  
schedulers.



This is not conflict serializable!

S

$S_1 : R_1(A), W_1(A), R_2(A), W_2(A), R_1(B),$   
 $W_1(B), R_2(B), W_2(B)$ .

$S_2 : R_1(A), W_1(A), R_1(B), W_1(B), R_2(A),$   
 $W_2(A), R_2(B), W_2(B)$ .

Whether it is conflict serializable or not.

$S_1$

$S_2$

$T_1$	$T_2$	$T_1$	$T_2$
$R(A)$		$R(A)$	
$W(A) \rightarrow R(A)$		$W(A)$	
$(W(A))$		$R(B)$	
		$W(B)$	
$R(B)$			$R(A)$
$W(B) \rightarrow R(B)$			$W(A)$
$(W(B))$			$R(B)$
			$W(B)$

It is conflict serializable.

Q 11  $\{g_1(x); g_1(y); w_1(u); w_1(z)\}$

$T_2 : - g_2(y); g_2(z); w_2(y)$ .

$T_3 : - g_3(y); g_3(x); w_3(y)$

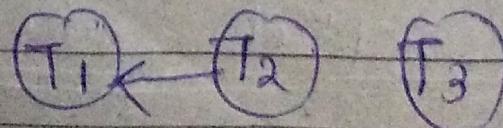
$S_1 : g_1(x); g_3(y); g_3(x); g_2(y)$   
 $g_2(z); w_3(y); w_2(z); g_1(3); w_1(x),$   
 $w_1(z).$

$S_2 : g_1(x); g_3(y); g_2(y); g_3(x); g_1(z)$

$g_2(z); w_3(y); w_1(x); w_2(z); w_1(3)$

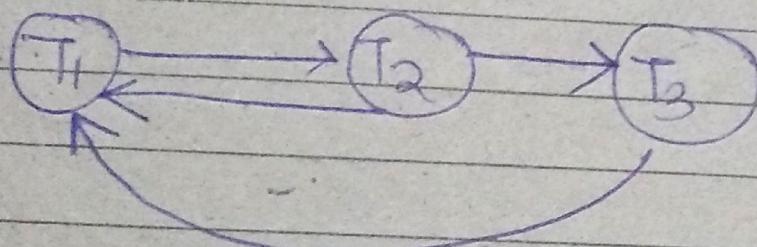
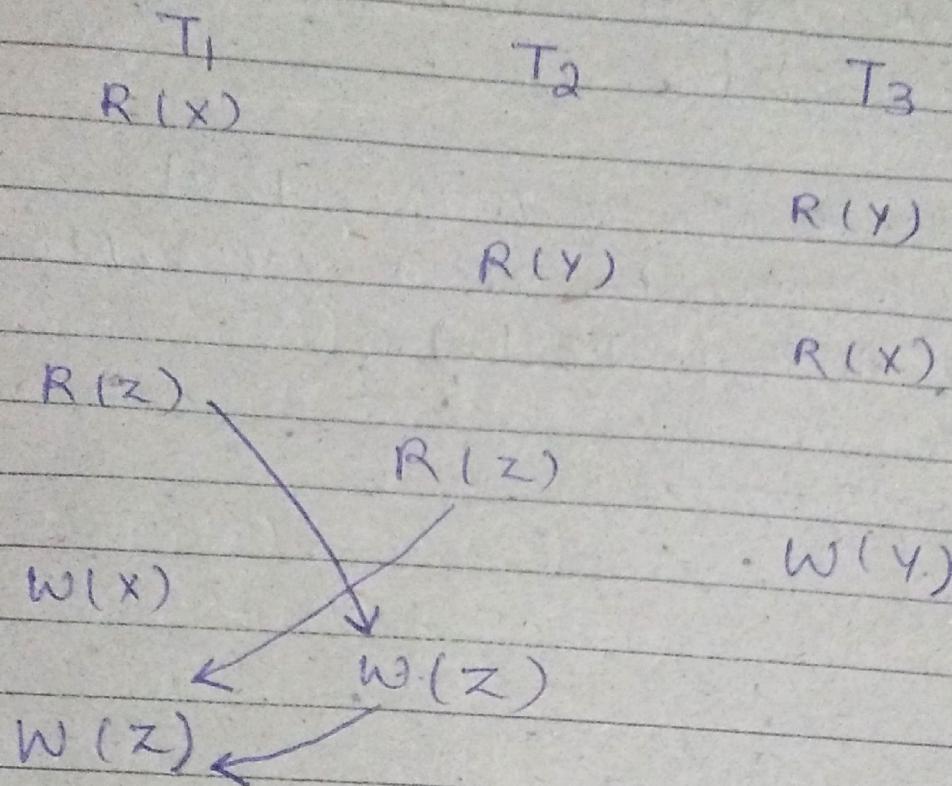
S1 $T_1$  $T_2$  $T_3$  $R(x)$  $R(y)$  $R(x)$  $R(z)$  $R(z)$ 

\*

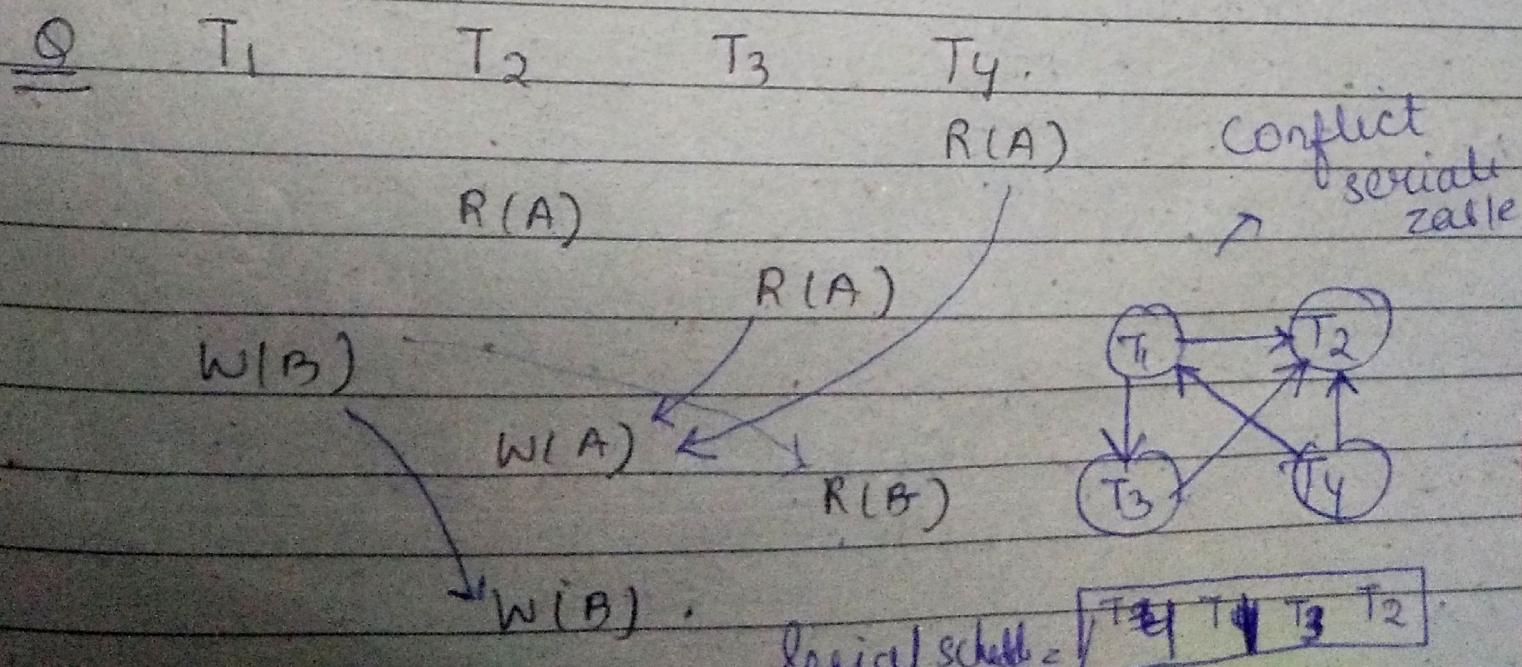
 $w(y)$  $w(z)$  $H(z) \leftarrow$  $w(u)$  $w(x)$ 

It is conflict serializable  
 $T_3 T_2 T_1$  or  $T_2 T_1 T_3$

S2



Cycle exists, it is not serializable



Logical schedule:  $T_1 \rightarrow T_3 \rightarrow T_2 \rightarrow T_1$

- Every conflict serializable is also view serializable.
- ~~empty~~ If a schedule is not conflict serializable, it may be view serializable (in case of blind write).

— X —