

Books → ① Koth ③ Schaum Series  
② Navathe

<u>Topics</u>		<u>Page No.</u>
① SQL	★★★	(3 - 66)
② Normalization	*	(67 - 120)
③ Hashing and Indexing	*	(121 - 166)
④ Transaction	★★	(167 - 206)
⑤ E-R Diagram		(207 - 22)

# Definition: A DBMS is a software whose basic requirement is to provide facility to store data.

\* Any collection of relevant data is database.

### # Other Requirements of DBMS:

(1) It should provide an easy language for data manipulation (SQL)

(2) It must maintain integrity (uniqueness) of data.

\* There can be two types of integrity:

① Entity Integrity

② Referential Integrity

(3) It must provide concurrent access to data i.e. multiple users can access data simultaneously.

(4) It must provide security.

→ —

\* Database is a collection of relevant or meaningful information.

# RDBMS (Relational DBMS) :

basic requirement is to

- \* It stores data in the form of tables or relations

eg:

RollNo	Name	Phone
1	Mayank	100
2	Aman	200
3	Neeraj	300

\* Here, RollNo, Name and Phone are columns / attributes / fields

\* Also, 1 Mayank 100 is a row / tuple / record

\* When data is stored in the form of rows & columns or tuples & attributes or records & fields, this format is known as tabular format or relation.

\* Nowadays, Object Oriented DBMS is being used.



SQL

⇒ Structured Query Language

① Create a table in database :

⇒ create table Student ( RollNo number(6) primary key, name varchar(20), Phone number(10) );

② Insert a record in table :

⇒ insert into Student values ( 1, 'Mayank', 100 );

\* All non-numeric values are written in single codes.

\* NULL is used wherever information is missing.

eg

⇒ insert into Student values ( 2, 'Aman', NULL );

⇒ insert into Student values ( 3, NULL, 300 );

\* NULL is always written without single codes.

→

\* To insert in particular columns ;

⇒ insert into Student ( RollNo, Phone ) values ( 1, 400 );

\* Remaining columns will contain NULL, but primary key must be inserted

⇒ 1 NULL 400

as it can't be  
NYU11

\* where clause works row-wise.

eg update student set RollNo = 1;  $\Rightarrow$  error as RollNo is primary key.  
Page No. \_\_\_\_\_  
(has to be unique)

eg update student set RollNo = 100 update where RollNo = 1; ✓ no error

### (3) Update a record:

\* it can also update primary key.

$\Rightarrow$  update Student set Phone = 500

\* it will set phone of all the students to 500.

$\Rightarrow$  update Student set Phone = 500

where RollNo = 10;

\* If there's no RollNo 10  $\Rightarrow$  no row updated.

\* If there are 2 RollNo 10  $\Rightarrow$  2 rows updated

Ques Write a query to update marks of every students by +10.

Student	RollNo	Name	Marks
	1	A	10
	2	B	
	3	C	15

\* B didn't appear in exams.

$\Rightarrow$  update Student set marks = marks + 10;

Student	RollNo	Name	Marks
	1	A	20
	2	B	
	3	C	25

$\Rightarrow$  NULL + Anything = NULL

#### ④ Delete a record :

⇒ delete Student or delete from student

\* All the records are deleted but the table still exists

\* delete command is used to delete rows not columns.

⇒ delete Student where RollNo = 10 ;

#### ⑤ Drop a table :

⇒ drop table Student

\* entire table is deleted.

#### ⑥ Fetch / Retrieve data from database :

⇒ select \* from student

\* it will display all records of the student table.

⇒	RollNo	Name	Phone
	1	A	100
	2	B	200
	3	C	300
	4	D	400

Ques Write a query to find name & salary of all the employees of dept. 10.

emp : empno	ename	sal	deptno
1	A	100	10
2	B	200	10
3	C	300	20

⇒ select ename , sal from emp where deptno = 10.

⇒ ename	sal
A	100
B	200

Ques Write a query to find name of employees of dept 10 and 20.

⇒ select ename from emp  
where deptno = 10 and deptno = 20 ; X

⇒ select ename from emp  
where deptno = 10 or deptno = 20 ; ✓

OR

⇒ select ename from emp  
where deptno in (10,20) ; ✓

~~Ques~~ Write a query to find name of employees having salary between 10000 to 20000.

⇒ select ename from emp  
where sal > 10000 or sal < 20000 X

⇒ select ename from emp  
where sal > 10000 and sal < 20000 ✓

OR  
⇒ select ename from emp  
where sal between 10000 and 20000 ✓

⇒ select ename from emp  
where sal between 10000 to 20000 X

→

• = : equal to

• != or <> : not equals

\* 'between' is usually inclusive

\* between also works with date

Ques

⇒ select ename from emp  
where ename between 'A' and 'J' ;

⇒  $s_1 = 'ABC'$       } which is greater?  
 $s_2 = 'x'$       }

\* ASCII values are compared in strings.

ABC

65              ⇒ 120 > 65      thus, 'x' is greater than 'ABC'

120

x

⇒ eg 'ABC'

⇒ 'ABC'      65 66      66 > 0 → 'ABC' > 'A'  
 'A'      65 0

⇒ 'ABC'      65      75 > 65 ⇒ 'J' > 'ABC'  
 'J'      75

thus, 'A' < 'ABC' < 'J'

'ABC' will be displayed.

\* It will display all the names b/w A & I &  
only 'J' can be displayed for J.

⇒ eg 'JAI'

⇒ 'JAI'      75       $75 > 65 \Rightarrow "JAI" > 'A'$   
                'A'      65

⇒ 'JAI'      75 65       $65 > 0 \Rightarrow 'JAI' > 'J'$   
                'J'      75 0

\* 'JAI' is greater than A but not smaller than 'J'  
thus, it won't be displayed.

→ X

1. is a wildcard character. It includes any no. of characters including NULL character.

Date: \_\_\_\_\_ Page No: \_\_\_\_\_

Ques Write a query to find the name of employees whose names start from A & end with T.

\* 'like' operator is used for string comparison.

⇒ select ename from emp  
where ename like 'A%T';

Ques Write a query to find name of employees whose names contain 2T. (i) atleast (ii) exactly

(i) ⇒ select ename from emp where  
ename like '%.T%.T%';

(ii) ⇒ select ename from emp where  
ename like '%.T%.T%'  
and ename not like '%.T%.T%.T%';

Ques Write a query to find name of employees such that A is at 1<sup>st</sup> position, anything at 2<sup>nd</sup>, T at 3<sup>rd</sup> & then anything.

⇒ select ename from emp where  
ename like 'A-T%';

%. → any no. of characters  
\_ → exactly one character

\* to search %, %

(11)

⇒ %.%. X /%./.% ✓  
Date \_\_\_\_\_ Page No. \_\_\_\_\_

Ques Write a query to find name of employees whose names contain exactly 3 characters.

⇒ select ename from emp where  
ename like '\_\_\_';

Ques Write a query to find name of employees whose emailid contains an underscore(\_).

⇒ select ename from emp where  
emailid like '%\_-%'; X

\* escape character (/) is used to search -, %, /, etc

⇒ select ename from emp where  
emailid like '%/\_-%';

\* escape character is different for different softwares.

Ques write a query to find name of employees having NULL salary.

⇒ select ename from emp where  
sal = NULL; X

• true ~~and~~ unknown = unknown  
• false and unknown = false

true or unknown = true  
false or unknown = unknown

Date \_\_\_\_\_ Page No. \_\_\_\_\_

\* In SQL, along with true and false, there's one more result of comparison i.e. NULL or unknown anything (+ or - or =) NULL  $\Rightarrow$  NULL

\* thus, 'is' operator is used for NULL comparison

$\Rightarrow$  select ename from emp where sal is NULL;  
 $\rightarrow$

\* To display in order;

\* Display follows order of insertion by default.

$\Rightarrow$  select ename from emp order by ename;

$\Rightarrow$  select ename from emp order by ename desc

Order of Insertion	Ordering	Ordering desc
e2	e1	e4
e4	e2	e3
e3	e3	e2
e1	e4	e1

\* If 1<sup>st</sup> column is same, then order of insertion will be followed.

\* ordering can be done in more than one column.

eName	Job
e1	Manager
e1	Salesman
e2	Manager
e2	Analyst

⇒ select eName, job from emp order by eName, job

eName	Job
e1	Manager
e1	Salesman
e2	Analyst
e2	Manager

\* if 1<sup>st</sup> column is same then ordering is done acc. to 2<sup>nd</sup> column.

\* \* ~~eName~~, job desc ⇒ eName (ascending), job (desc.)

\* eName desc, job desc ⇒ both (descending)

\* eName ~~job~~ desc, job ⇒ eName (desc.), job (asc.)

# Aggregate Functions / Group Value Functions

- (1) sum
- (2) Min
- (3) Max
- (4) Count
- (5) Average

(1) Sum:

⇒ select sum (sal) from emp;

	Name	Sal	
e1		10000	
e2			Sum ⇒ 80000 ✓
e3		30000	⇒ NULL ✗
e4		40000	

\* aggregate functions ignore NULL values.

(2) Min:

⇒ select min (sal) from emp; ⇒ 10000 ✓  
NULL ✗

\* If there's no row in table, except 'count', all can return null.

max, min also works with date.

max → recent date  
min → oldest date

Date \_\_\_\_\_ Page No. \_\_\_\_\_

### ③ Max :

⇒ select max(sal) from emp; ⇒ 40000

### ④ Count :

⇒ select count(sal) from emp; ⇒ 4 X  
3 ✓

\* It also ignores NULL value

\* select count(distinct sal) from emp;

⇒ select count(ename) from emp; ⇒ 4

\* ⇒ count(\*) ⇒ counts no. of rows in database.

\* If there's no row ⇒ it returns zero

(not NULL)\*

### ⑤ Average :

⇒ select avg(sal) from emp;

$$\text{avg(sal)} = \frac{\text{sum(sal)}}{\text{count(sal)}} = \frac{80000}{3} = 26666.66$$

$$\frac{80000}{4} X$$

\* It also ignores NULL value.

\* sum, max, min, avg don't work with strings.

⇒ select ename, avg(sal) from emp; X error

\* ename & avg(sal) are not compatible as column  
ename contains many values while avg has only  
one value.

\* Thus, we cannot use any aggregate function with  
non-aggregate columns.

\* Two or more aggregate functions are compatible  
with each other.

→

ename	sal	deptno
e1	10000	10
e2	40000	10
e3	15000	20
e4	30000	20

Ques Write a query to print sum of salaries dept. wise.

⇒ select sum(sal) from emp  
group by deptno;

\* 2 rows will be printed.

O/P ⇒ 50000  
45000

- \* Select deptno , sum(sal) from emp  
group by deptno; ✓
  - \* A non-aggregate column can't be selected with aggregate function unless that column is grouped.
- O/P →
- | deptno | sum(sal) |
|--------|----------|
| 10     | 50000    |
| 20     | 45000    |
- \* groups are also formed for NULL.i.e if deptno contains a NULL , then one more row would be printed in O/P.
- 
- \* select ename from emp  
group by deptno; X
  - \* If we've done grouping . then we can only use aggregate functions or grouped column ( deptno here )
-

\* select deptno from emp group by deptno;

⇒ 10

20

\* select deptno from emp;

⇒ 10

10

20

20

\* select distinct deptno from emp;

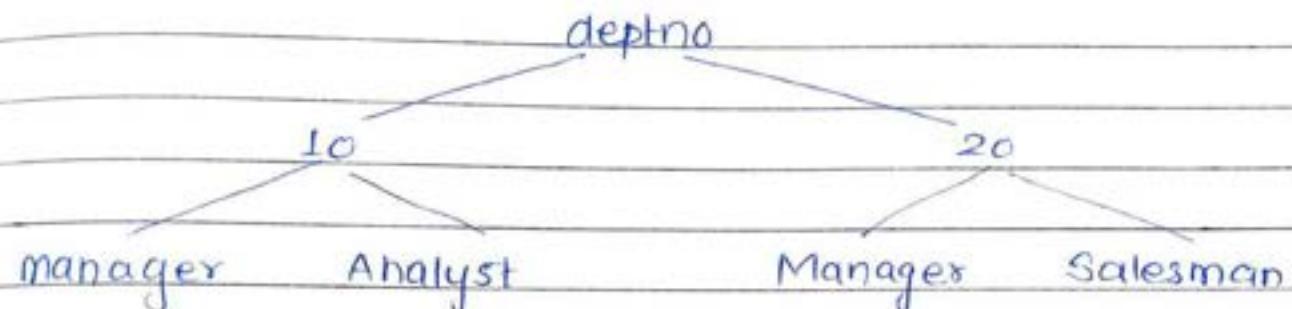
⇒ 10

20

→ —

\* Grouping can also be done on multiple columns.

⇒ select deptno , job , sum (sal) from emp  
group by deptno , job ;



O/P ⇒      deptno      job      sum (sal)

10	Manager	10000
10	Analyst	20000
20	Manager	30000
20	Salesman	40000

→

- \* 'where' clause works row-wise
- \* 'having' clause works group-wise  
thus 'having' clause comes after 'group by' statement.

Ques Write a query to find the deptno where average salary of employee is more than 5000.

⇒ select deptno from emp  
where avg(sal) > 5000    X  
group by deptno;

\* aggregate functions never use 'where' clause , they use 'having' clause. to check conditions.

⇒ select deptno from emp  
group by deptno  
having avg(sal) > 5000 ;

Ques write a query to find the deptno where average salary of managers is greater than 5000.

⇒ select deptno from emp where job = 'Manager'  
group by deptno.  
having avg(sal) > 5000;

\* First, managers are selected then grouping will be done.

\* 'having' clause can be used without grouping. but, Date \_\_\_\_\_ Page No. \_\_\_\_\_  
then complete table will be treated as a single group.  
but we can use only aggregate functions in select.

⇒ select deptno from emp  
group by deptno

having avg(sal) > 5000 and job = 'Manager';

\* job is not unique for a group

⇒ select deptno from emp

group by deptno, job

having avg(sal) > 5000 and job = 'Manager';

→ —

⇒ select deptno from emp where job = 'Manager'  
group by deptno, job

having avg(sal) > 5000 and job = 'Manager';

→ —

• To find deptno where avg. salary is greater than 5000.

⇒ select deptno, asal from (select deptno, avg(sal)  
as asal (wx) group by deptno  
from emp);

where asal > 5000;

$\pi \rightarrow$  project ('select' in SQL)  
 $\sigma \rightarrow$  select ('where' in SQL)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

## # Relational Algebra :

\* It is a theoretical language of query writing & SQL is practical implementation of relational algebra.

$\Rightarrow$  select ename from emp  
 $\Rightarrow \pi_{ename}(\text{emp})$

$\Rightarrow$ ename	SQL	R.A
e1	$\Rightarrow e1$	e1
e2	e2	e2
e3	e3	e3
e3	e3	

\* set notation  $\Rightarrow$  duplication ✗ (default)  
bag notation  $\Rightarrow$  duplication ✓

\* Relational algebra is based on set theory.

$\Rightarrow$  select distinct ename from emp =  $\pi_{ename}(\text{emp})$

eg  $\pi_{ename} \circ \sigma_{\text{sal} > 1000}(\text{emp})$  •  $\pi_{\substack{\text{sal} > 10 \\ 1000}}(\text{emp})$  ✓

eg  $\pi_{ename} \circ \sigma_{\text{sal} > 10000 \wedge \text{deptno} = 10}(\text{emp})$

✳  $\sigma_{\text{sal} > 1000}(\text{emp}) \Rightarrow$  select \* from emp where sal > 1000

## # Nested Query: (query within a query)

Ques Write a query to find the name of employees having highest salary.

⇒ select ename from emp where  
 $\text{sal} = (\text{select max(sal)} \text{ from emp})$

\* first inner query is executed.

Ques write a query to find the name of employees having salary greater than that of e1.

⇒ select ename from emp where  
 $\text{sal} > (\text{select sal from emp where ename} = 'e1');$

\* if there is no e1 ⇒ inner query will return NULL. ✓

\* if there are more than 1 e1 ⇒ error

⇒ Thus, to resolve this,

① ~~>~~ Any

② > All

③ ...  $\text{sal} > (\text{select max(sal)} \text{ from emp where ename} = e1)$   
 or  
 $\text{min(sal)}$

①  $\rightarrow$  select ename from emp where  
 $\text{sal} > \text{Any}$  (select sal from emp  
 where ename = 'e1');

\*  $\text{sal} > 5000 / 10000 \Rightarrow \text{sal} > \underline{\underline{5000}}$

②  $\rightarrow$  select ename from emp where  
 $\text{sal} > \text{All}$  (select sal from emp  
 where ename = 'e1');

\*  $\text{sal} > 5000 + 10000 \Rightarrow \text{sal} > \underline{\underline{10000}}$

Ques write a query to find the name of employees  
 having salary less than that of e1.

$\Rightarrow$  select ename from emp where  
 $\text{sal} < (\text{select sal from emp where ename} = \text{'e1'})$

\* if there is no e1  $\Rightarrow$  inner query will return NULL ✓

\* If there are more than 1 e1  $\Rightarrow$  error

\* Thus, to resolve this,

①  $<$  Any

②  $<$  All

①  $\Rightarrow$  select ename from emp where  
~~sal <~~ Any (select sal from emp where ename = 'e1')

\*  $\text{sal} < 5000 / 10000 \Rightarrow \text{sal} < \underline{10000}$

②  $\Rightarrow$  select ename from emp where  
~~sal <~~ All (select sal from emp where ename = 'e1')

\*  $\text{sal} < 5000 \times 10000 \Rightarrow \text{sal} < \underline{5000}$

$\rightarrow$

\*  $\Rightarrow$  select ename from emp where  
~~sal >~~ Any (select sal from emp where ename = 'e1')

\* if there is no e1  $\Rightarrow$  inner query will return NULL

$\Rightarrow \text{sal} > \text{Any (NULL)}$

$\rightarrow$

$\Rightarrow$  select ename from emp where  
~~sal >~~ All (select sal from emp where ename = 'e1');

\* if there is no e1  $\Rightarrow$  inner query will return 0(zero).

$\Rightarrow \text{sal} > \text{All (0)}$

\* it will display ename of sal > 0.

Ques write a query to find the name of employees having salary equal to that of e1.

⇒ select ename from emp where  
~~sal = ( select sal from emp where ename = 'e1' )~~

\* if there is no e1 ⇒ inner query will return NULL ✓

\* if there are more than 1 e1 ⇒ error

⇒ Thus, to resolve this,

① = Any	}	Both are same.
② = All in		

① ⇒ select ename from emp where

sal = Any (select sal from emp where ename = 'e1')

② ⇒ select ename from emp where

sal in (select sal from emp where ename = 'e1');

↑

emp:

ename	deptno
e1	10
e2	20
e3	20
e4	10

dept:

deptno	dname
10	Research
20	Accounts

Ques Write a query to find the name of employees of Research dept.

→ select ename from emp where  
 deptno = (select deptno from dept where  
 dname = 'Research') ;

Ques Write a query to find the name of employees having 2<sup>nd</sup> highest salary in emp table.

→ select ename from emp where  
 sal = (select max(sal) from emp where  
 sal != (select max(sal) from emp)) ; ✓

OR

→ select ename from emp where  
 sal = (select max(sal) from emp where  
 sal < (select max(sal) from emp)) ; ✓

Que Write a query to find 3<sup>rd</sup> highest salary-employee.

⇒ select ename from emp where  
 sal = ( select max(sal) from emp where  
 sal != ( select max(sal) from emp where  
 sal != ( select max(sal) from emp))) ;      X

\* It will return 1<sup>st</sup> highest salary.

⇒ select ename from emp where  
 sal = ( select max(sal) from emp where      ✓  
 sal < ( select max(sal) from emp where  
 sal < ( select max(sal) from emp))) ;



•  $\Pi_{\text{sum(sal)}}$  X      , For aggregate functions, we can't use project ( $\Pi$ ).

⇒  $G_{\text{sum(sal)}}$  (emp)

• for grouping,  
 (on deptno column)      ⇒ define  $G_{\text{sum(sal)}}$  (emp)

" $\leftarrow$ "  $\Rightarrow$  assignment (in R.A.)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

## # Joins

Ques Write a query to print the employee name with resp. dept.

- \* join is used to fetch the data from more than 1 table
- $\Rightarrow$  select ename, dname from emp, dept.
- \* It is cartesian product of 2 tables. Each row of emp table is 'join' to each row of dept table.
- \* If one table contains m rows & another contains n rows, then their cartesian product will have  $m \times n$  rows.



Ques Write a query to find the name of employees having salary greater than that of e1. (in R.A.)

$\Rightarrow S \leftarrow \pi_{\text{sal}} \circ \text{ename} = 'e1' (\text{emp}) \quad // \text{inner query}$

$\pi_{\text{ename}} \circ \text{sal} \geq s (\text{emp}) \quad // \text{outer query}$



\* select \* from emp, dept  $\Rightarrow$  8 rows & 4 columns

\* equal condition on common field  $\rightarrow$  Natural Join.

Date: \_\_\_\_\_ Page No. \_\_\_\_\_

\* If one table contains 10 rows & other contains no rows ( $\emptyset$ ), then their cartesian product will also be  $\emptyset$ .



# Natural Join:

$\Rightarrow$  select ename, dname from emp, dept  
where emp.deptno = dept.deptno;

\* if columns are common in both the tables, then write it as;  $\rightarrow$  tablename.columnname else there ~~if~~ emp.deptno, dept.deptno. will be error.

$\Rightarrow$  emp:

ename	deptno
e1	10
e2	20
e3	20
e4	10

dept:

deptno	dname
10	Research
20	Accounts

O/P  $\Rightarrow$

e1	Research
e2	Accounts
e3	Accounts
e4	Research

\* This operation is called Natural Join.  
(equality on common column)

\* Alternate way :

⇒ select ename, dname from emp natural join dept;

\* select \* from emp natural join dept;

O/P ⇒ 3 columns (common column is displayed only once (deptno here)).

\* If there is no common column, it will simply perform cartesian product.

\* If there are more than 1 column, it will equate all the columns.

eg R : S:

A	B	A	B
a1	b1	a1	b2
a2	b2	a2	b2
a3	b2	a3	b2
a4	b1	a4	b2

\* select \*  
from R natural  
join S.

A	B
a2	b2
a3	b2

\* But if want to equate only 1 column.

$\Rightarrow$  select \* from R natural join S using f

O/P $\Rightarrow$	A	R B	S B
a1	b1	b2	
a2	b2	b2	
a3	b2	b2	
a4	b1	b2	



\* In Relational algebra, cartesian product  $\Rightarrow \times$

eg  $\pi_{ename, dname} (\text{emp} \times \text{dept})$

\* natural join  $\Rightarrow *$  or  $\bowtie$

eg  $\pi_{ename, dname} \circ \text{emp.deptno} = \text{dept.deptno} (\text{emp} \times \text{dept})$

OR  $\pi_{ename, dname} (\text{emp} * \text{dept})$

OR  $\pi_{ename, dname} (\text{emp} \bowtie \text{dept})$



If R & S are same schema (ie same no. of rows & columns)  
 then,  $R \times S = R \cap S$

(33)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

$\Rightarrow R:$

A	B
a1	b1
a2	b1
a3	b2
a4	b2

S:

A	B
a1	b1
a2	b2
a3	b3
a4	b2

$$\Rightarrow \pi_A(R \times S) \Rightarrow \{a1 \\ a4\}$$

Both are different.

$$\Rightarrow \pi_A(R) \times \pi_A(S) \Rightarrow \{a1 \\ a2 \\ a3 \\ a4\}$$

$\Rightarrow R:$

A	B
a1	b1
a2	
a3	b3

S:

B	C
b1	c1
b2	c2
b3	c4

\* select \*

from R natural  
join S.

$$\Rightarrow O/P = a1 \ b1 \ c1 \\ a3 \ b3 \ c3$$

\* In cartesian prod.,  
all columns are displayed  
but in natural join

common column is  
displayed once.

\* NULL - NULL join  $\neq$   
 eg a2 c4  $\neq$

## # Outer Joins :

\* They are the extension of natural join. They displays unmatched rows also.

- ① Left Outer Join: displays unmatched rows of left table.
- ② Right Outer Join: displays unmatched rows of right table.
- ③ Full Outer Join: displays unmatched rows of both tables.



$\Rightarrow P.:$

A	B	B	C
a1	b1	b1	c1
a2		b2	c2
a3	b3		c4
		b3	c3

$\Rightarrow S.:$

$\Rightarrow$  Left Outer Join

$\Rightarrow$  Right Outer Join

A	B	C	$\Rightarrow$	A	B	C
a1	b1	c1		a1	b1	c1
a3	b3	c3		a3	b3	c3
a2					b2	c2

c4

$\Rightarrow$  Full Outer Join :

	A	B	C
a1	b1	c1	
a3	b3	c3	
a2			
	b2	c2	
		c4	

In SQL,



R  $\rightarrow$  left table  
S  $\rightarrow$  right table

- \* select \* from R left outer join S
- \* select \* from R right outer join S
- \* select \* from R full outer join S.

In R.A.,

- \*  $\pi_{A,C}(R \bowtie S) \Rightarrow$  left outer join
- \*  $\pi_{A,C}(R \bowtie S) \Rightarrow$  right outer join
- \*  $\pi_{A,C}(R \bowtie S) \Rightarrow$  full outer join

$R \Rightarrow$  left table  
 $S \Rightarrow$  right table

In ORACLE,

$\Rightarrow$  select \* from R, S  $\Rightarrow$  natural join  
where  $R.B = S.B$

$\Rightarrow$  select \* from R, S  $\Rightarrow$  left outer join  
where  $R.B = S.B (+)$

$\Rightarrow$  select \* from R, S  $\Rightarrow$  right outer join  
where  $R.B (+) = S.B$

$\Rightarrow$  select \* from R, S  $\Rightarrow$  full outer join  
where  $R.B (+) = S.B (+)$



# Primary Key : unique & not null.

\* only 1 primary key is allowed.

\* All those keys which are capable of being primary are known as candidate keys.

eg rollno, enrollmentno are 2 candidate keys.

⇒ create table Student ( Rollno number(6) primary key,

Empl. varchar(10) unique not null ,

name varchar(10) not null ,

age number(2) check (age between 21 and 28),

branch varchar(6) default 'CS' );

\* check constraint is checked on insertion & updation.

→ →

$\Rightarrow$	Rollno	courseno	grade
	R1	c1	A
	R1	c2	B
	R2	c1	B
	R2	c3	A

\* Here no column is unique, but combination of columns is unique eg. (Rollno., courseno) Thus, here candidate key consists of 2 columns.

\* candidate key is said to be composite candidate key if it is formed by more than 1 column.

$\Rightarrow$  create table Student ( Rollno number(6),  
courseno number(6) , grade varchar(2),  
primary key ( Rollno , courseno));

\* Rollno & courseno cannot be null & (duplicate, as they are part of candidate key / primary key. at same time)

\* Individual attributes can be duplicate but not null.

# Foreign Key :

emp:

ename sal deptno

dept:

deptno	dname
10	Research
20	Design
30	Testing

- ⇒ Can we insert (e2, 1000, 20) ✓  
 ⇒ Can we insert (e3, 2000, 40) X entire row will be rejected.

\* dept table is parent table & emp is child table.

~~or updation~~

- \* Whenever there's insertion in child table, parent table is referred. If that value is present in parent table then only it can be inserted in child table.
- \* Now, deptno. column of emp table is declared as foreign key.
- \* We can declare foreign key in child table only when that column is either primary key or unique in parent table.

→ —

- \* Database ensures entity integrity by giving facility to declare primary key.
- \* Database ensures referential integrity by giving facility to declare foreign key.

→ —

### ① create parent table :

⇒ create table dept (deptno number(2) primary key,  
dname varchar(20));

- \* deptno must be primary or unique.

### ② create child table :

⇒ create table emp (ename varchar(10) primary key,  
deptno number(2) foreign key references dept);

- \* foreign key can accept any no of NULL values.

- \* In dept table , deptno can't be NULL as it's primary key , but it can be NULL in emp table as it's foreign key in emp table.

→ —

emp:

ename	deptno
e1	10
e2	10
e3	20

dept:

deptno	dname
10	Research
20	Accounts
30	Academics

⇒ update dept set deptno = 100 where deptno = 10; X

\* we cannot update / delete any value from parent table when that value is present in child table.

\* But, if you want to delete / update , then there are 2 ways ;

① on delete null : If we delete any value from dept table , that value will become NULL in emp table.

⇒ create table emp (ename varchar(10) primary key, deptno number(2) foreign key references dept on delete null);

→

\* on update null & on update cascade also exists.

⇒ on update cascade , if there's updation in parent table, there will be updation in child table automatic.

② on delete cascade : If we delete any value from dept, that entire row will become NULL in emp table.

→ create table emp (ename varchar (10) primary key, deptno number (2) foreign key references dept on delete cascade);

→

\* parent-child relationship can be there b/w columns within same table.

eg	empno	ename	mgr	
	1	e1		* empno ⇒ parent key
	2	e2	1	
	3	e3	4	* mgr ⇒ child key
	4	e4	3	

\* mgr can accept those values which are present in empno

→

Ques B is foreign key referring to column A of same table. On delete cascade is set.

⇒ delete from R where A = 1 ;

How many rows will be deleted ?

R :

A	B
1	1
2	1
3	2
4	2
5	3
6	4

Ans : All the rows will be deleted as on delete cascade is set.

- \* 'with' clause is used to create temporary table.
- To display employee name drawing highest salary.
- ⇒ With large-sal (msal) as select max(sal) from emp;
- table      column
- select ename from emp, large-sal where sal = msal;
- To display deptno whose largest salary is more than avg. sal.
- ⇒ With large-sal (msal) as select ~~deptno~~, max(sal) from emp group by deptno.
- With avg-sal (asal) as select avg(sal) from emp
- select deptno from large-sal, avg-sal where msal > asal

# Self Join : joining a table to itself.

emp:

empno	ename	mgrno
1	e1	
2	e2	1
3	e3	1
4	e4	2

Ques write a query to print the employee name with manager name.

\* This query can't be solved with one table as we want data from 2 rows. Thus, we need 2 tables

2 copies

$x_1$	$x_2$
empno	empno
ename	ename
mgrno	mgrno
1	1
e1	e1
2	2
e2	e2
1	1
3	3
e3	e3
1	1
4	4
e4	e4
2	2

$\Rightarrow$  select  $x_1$ .ename "employee",  $x_2$ .ename "Manager"  
from emp as  $x_1$ , emp as  $x_2$   
where  $x_1$ .mgrno =  $x_2$ .empno;

$x_1$  &  $x_2$  are known as alias or instances or virtual copies of the table & these are valid till query is under execution.

O/P $\Rightarrow$	employee	Manager
	e2	e1
	e3	e1
	e4	e2

\* This query can be solved by making only 1 copy.

$\Rightarrow$  select emp.ename "employee",  $x_2$ .ename "Manager"  
 from emp, emp as  $x_2$   
 where emp.mgno =  $x_2$ .empho;

In relational algebra, above query will be,

$\Rightarrow \pi_{emp.ename, x_1.ename} \sigma_{emp.mgno = x_1.empho}$   
 $(emp \times f(emp, x_1))$

\*  $f$  is operator for rename.

\* In above query, we rename the table but in relational algebra, we can write query of self-join also by renaming column, it can't be done in SQL.

- \* For natural join, 'using' is used. (Here, 'on' can also be used)
- \* For other joins, 'on' is used.

Date \_\_\_\_\_ Page No. \_\_\_\_\_

⇒ By renaming column :

empno	ename	mgrno	// column names in 1 <sup>st</sup> copy (table)
eno	name	mno	// in 2 <sup>nd</sup> copy

⇒  $\pi_{ename, name} ( emp \bowtie_{mgrno=eno} f(eno, name, mno) \text{emp} )$

→

\* Equal on common columns → Natural Join

\* Equal on uncommon columns ⇒ Equi Join

\* Unequal, greater than, etc. on columns ⇒ Theta Join

eg  $\bowtie$ ,  $\bowtie$ , etc.  
 $mgrno \neq eno$ ,  $mgrno > eno$

\* all joins are theta joins.

→

\* All joins except outer join are inner joins.

eg select x1.ename as "employee", x2.ename as "manager"  
from emp as x1 join emp as x2  
on x1.mgr = x2.empno

theta  
join

# Set Operations:① Union:

credit:	debit:
cname	dname
c1	d1
c2	d2
c3	d3
d1	c1

Ques Write a query to print the names of account holders

⇒ select cname from credit  
union

select dname from debit;

O/P ⇒ 

cname
c1
c2
c3
d1
d2
d3

 \* It will display column name of first table i.e. cname

\* If any value is common in both the tables then these set operations will suppress those duplicate values. (by-default) (eg c1, d1 here)

\* But if we want to display duplicate values too.

⇒ select cname from credit  
union all

select dname from debit;

→

\* types of columns must be same, their size can be different.

e.g: cname varchar (20) , dname varchar (60) ✓

→

\* The union of 2 tables can also be performed if they are union compatible to each other.

⇒ Two relations are said to be union compatible if there exists one to one correspondence between their columns. i.e. if first table contains 3 columns then 2<sup>nd</sup> table must have 3 columns & the types must be like; (size doesn't matter)

- \* type of a1 = type of b1
- \* type of a2 = type of b2
- \* type of a3 = type of b3

\* size of data  
type doesn't  
matter.

A:

a1      b2      c3

B:

b1      b2      b3

eg

credit :

cno	cname
1	c1
2	c2
3	d1

debit :

dno	dname
1	d1
2	d2
3	d1

⇒ select \* from credit  
union

select \* from debit

O/P ⇒ cno cname

1	c1
2	c2
3	d1
1	d1
2	d2

\* if entire tuple  
is duplicate, then  
it will be displayed  
only once.

eg credit :

addr	cname
XYZ	c1
EJK	c2
123	d1

debit :

dname	addr
d1	ABC
d2	DEF
d3	GHI

⇒ select \* from credit  
union

select \* from debit

O/P  $\Rightarrow$

addr	cname
XYZ	c1
IJK	c2
123	d1
d1	ABC
d2	DEF
d3	GHI

$\gamma$

## ② Intersection:

$\Rightarrow$  select \* from credit  
 intersect or intersection  
 select \* from debit

eg credit: debit:

cno	cname	dno	dname
1	c1	1	c1
1	c2	1	c1

O/P  $\Rightarrow$

cno	cname
1	c1

$\Rightarrow$  select \* from credit  
 intersect all or intersection all  
 select \* from debit

O/P $\Rightarrow$	cno	cname	
1	c1		{ twice, not
1	c1		4 times }

 $\rightarrow$ 

### ③ Minus / Except :

\*  $A - B \Rightarrow$  present in A and not in B

$\Rightarrow$  select \* from credit

minus / except

select \* from debit

eg

credit:		debit:	
cno	cname	dno	dname
1	c1	1	d1
2	c2	2	c2
3	c3	3	d3

O/P  $\Rightarrow$

cno	cname
1	c1
3	c3

\* In Relational algebra, for two relations R & S,

① union  $\Rightarrow$   $R \cup S$

② intersection  $\Rightarrow$   $R \cap S$

③ minus  $\Rightarrow$   $R - S$

Ques what is the estimated size of natural join?

R : (child table)

A	B(FK)
a1	b1
a2	b1
a3	b2
a4	b3

S : (parent table)

B(PK)	C
b1	c1
b2	c1
b3	c2
b4	c3

- \* If table R contains n tuples & S contains m tuples then their natural join can have maximum n tuples & minimum 0 tuples as there can be NULL values in column B.

Ques

R :

A.	B
a1	b1
a2	b1
a3	b1
a4	b1
a5	b1
a6	b1

S:

B	C
b1	c1
b1	c2
b1	c3

\* no. Info is given about primary & foreign key.

- \* If table R contains n tuples & S contains m tuples then their natural join can have maximum m \* n tuples & minimum 0 tuples.

# Correlated Subqueries :

Ques write a query to find the name of employees having highest salary in their dept.

① Nested Query:

⇒ select ename from emp where sal in (select max(sal) from emp group by deptno);

ename	sal	deptno
e1	10000	10
e2	20000	10
e3	30000	20
e4	40000	20
e5	20000	20

O/P ⇒ e2  
e4  
e5

⇒ select ename from emp where  
(sal, deptno) in (select max(sal), deptno from emp group by deptno);

O/P ⇒ e2  
e4



## ② Correlated Subquery:

\* In correlated subqueries, inner query is executed whenever 'where' clause is executed & 'where' clause works row wise, thus, inner query is executed per row.

⇒ select e1.ename from emp as e1 where

e1.sal = (select max(sal) from emp  
where deptno = e1.deptno);

\* When inner query contains reference of outer table, it is correlated subquery, not nested query.

\* Outer query is executed per row while the inner query is executed per entire table.

-p-

③ Exists

- \* if inner query returns some rows  $\rightarrow$  true
- \* if inner query doesn't return any row  $\rightarrow$  false

$\Rightarrow$  select x1.ename from emp as x1 where

not exists (select \* from emp where  
deptno = x1.deptno and sal > x1.sal) ;

O/P  $\Rightarrow$  e2  
e4

$\rightarrow$

\* 'exists' can also be used with nested query.

$\Rightarrow$  But, it will display either entire table or nothing.  
as in nested queries

ename	sal	deptno	inner query is run only once.
e1	10000	10	
e2	20000	10	
e3	20000	20	
e4	30000	20	

$\Rightarrow$  select ename from emp exists (select \* from emp where deptno = 10);

$\Rightarrow$  exists (TRUE)  $\Rightarrow$  displays entire table ; as in nested query, inner query is executed only once.

$\rightarrow$

# unique:

customers (cno, cname), purchase (cno, itemno)

$\Rightarrow$  To find the name of customer who has bought only 1 items

$\Rightarrow$  select cname from customers where unique (select ~~cno~~<sup>cno</sup> from purchase where Purchase.cno = Customer.cno).

Purchase:

cno	itemno
-----	--------

1	I1
---	----

$\Rightarrow$  unique (1,1)  $\Rightarrow$  not unique

1	I2
---	----

2	I1
---	----

$\Rightarrow$  unique (2)  $\Rightarrow$  unique

Customer:

cno	cname
-----	-------

1	C1
---	----

2	C2
---	----

O/P  $\Rightarrow$  C2

$\rightarrow$

eg unique (1,2)  $\Rightarrow$  unique

alter table emp add (name number(6)); //add column  
Date \_\_\_\_\_ Page No. \_\_\_\_\_

(53)

\* existing rows will have NULL for this new column.

→ SQL commands are divided into :

① DDL (Data Definition Language)

\* to change structure of table/database

⇒ create, drop, alter, etc.

② DML (Data Manipulation Language)

\* to manipulate data.

⇒ select, insert, update, delete, etc.

③ DCL (Data Control Language) : grant, revoke

④ TCL (Transaction Control Language) : commit, rollback.

⇒ commit, rollback, grant, revoke

\* commit :

\* Whatever changes we perform, are saved in primary memory. To save it permanently in harddisk, 'commit' is used.

\* Autocommit or exiting software correctly saves changes permanently.

- Rollback :

- \* ~~store~~ till last commit, 'rollback' undo the changes we performed.

eg commit

insert into emp values (100, \_\_\_\_\_);

delete t1

create table t2 ( \_\_\_\_\_ );

update student set phone = 9933101601  
where rollno = 2;

select \* from emp

drop table t3

Rollback

they're auto committed

- \* DDL commands can't be rolled back, thus, dropped tables cannot be recovered.

→ → →

- Grant :

users : u1 u2

tables : t1 t2

- \* Users cannot access the tables of other users.  
If u1 wants to access t2 then u2 will run a 'grant' command to give permission to u1 to access t2.

u2  $\Rightarrow$  grant select, insert on t2 to u1

\* u1 cannot run update, delete, etc. commands on t2.

u2  $\Rightarrow$  grant all on t2 to u1

all  $\Rightarrow$  select, insert, update, delete, create

\* Now if u1 wants to give permissions further to other users.

u2  $\Rightarrow$  grant all on t2 to u3 with grant option.

→ —

#### \* Revoke :

\* to take given permission back, 'revoke' is used.

u2  $\Rightarrow$  revoke all on t2 <sup>from</sup> ~~to~~ u1

u2  $\Rightarrow$  revoke select on t2 from u1

→ —

# View : They are virtual tables.

→ create view v1 as  
(select ename, empno from emp)

→ select \* from v1.

\* Views are known as virtual tables because they don't have their own database; they fetch data from table on which they're defined.

→ In grant command, we give permission to access entire table, but if we want other users to access only some selected columns then we can create views.

\* It is possible to insert/update, with the help of views.; provided that view satisfies following conditions;

① view should be single table based. (there should not be any join or nested query in view definition)

② view definition should not contain any aggregate functions.

③ primary key & all 'not null' keys must be the part of view definition.

$$Q_1(x, y) \div Q_2(x) \quad * Q_2 \text{ can project only those columns which are being 'project' by } Q_1.$$

$$\Pi(x, y) \div \Pi_x \quad * x \subseteq (x, y)$$

Date: \_\_\_\_\_ Page No. \_\_\_\_\_

(61)

$\Rightarrow$  I'll select those values of  $x$  in which all values of  $x$  with  $y$  are present

$\Rightarrow$  Create table emp ( empno varchar(10),

deptno number(6) check ( deptno in ( select deptno from dept ));

\* This query is not equivalent to foreign key because, in foreign key, constraint is on both tables while here constraint is only on emp table.

\* There's no constraint on dept table ; if a record is updated in dept table, it won't affect emp table.



\* "For all" type of queries, divide operator is used.

Part: Pno Pcolor Pname			Supplier: Sno Pno	
1	Red	P1	S1	1
2	Green	P2	S1	2
3	Red	P3	S1	3
4	Blue	P4	S2	3
			S2	4
Query to find the suppliers which supply <u>all</u> red color parts.			S3	1
			S3	3

$\Rightarrow \Pi_{Sno, Pno}(\text{Supplier}) \div \Pi_{Pno} \sigma_{Pcolor = 'red'}(\text{Part}) \text{ O/P } \uparrow S1$

R.A. is procedural but Rel. calculus is non-procedural language

Date \_\_\_\_\_ Page No. \_\_\_\_\_

# Relational Calculus: It is also a theoretical language of query writing.

- ① Tuple calculus (Row)
- ② Domain calculus (Column)

\* Relational calculus is based on predicate / first order logic.

# Tuple Calculus:

$\Rightarrow t \mid t \in \text{emp} \Rightarrow \text{select } * \text{ from emp}$

$\Rightarrow t.\text{ename} \mid t \in \text{emp} \Rightarrow \text{select ename from emp}$

~~$t \mid \forall t \in \text{emp} \wedge t = t[\text{ename}] \Rightarrow \text{select ename}$~~

$t \mid \exists d \in \text{emp} (t[\text{ename}] = d[\text{ename}]) \Rightarrow \text{from emp}$

~~$t \mid \exists U \in \text{emp} \wedge (t[\text{ename}] = U[\text{ename}] \wedge U[\text{sal}] > 10000)$~~

$\Rightarrow t.\text{ename} \mid t \in \text{emp} \wedge t.\text{sal} > 10000$

$\Rightarrow \text{select ename from emp where sal} > 10000$

\* Like R.A., R.C. also avoids duplicate values.

Ques Write a query to find the name of employees of Research dept. (in Tuple calculus).

$\Rightarrow t \cdot \text{ename} \mid t \in \text{emp} \wedge \exists d \in \text{dept} \wedge$

$d \cdot \text{dname} = \text{'Research'} \wedge d \cdot \text{deptno} = t \cdot \text{deptno}$

OR

$\Rightarrow t \mid \exists u \in \text{emp} \Delta (t[\text{ename}] = u[\text{ename}]) \wedge$

$\exists d \in \text{dept} (d[\text{dname}] = \text{'Research'} \wedge d[\text{deptno}] = u[\text{deptno}])$

\* if  $\exists$  is replaced by  $\forall$ , this query will only work if dept table ~~is~~ has only one row & the dname of that row is Research.

- - -

Bounded or

\* variable with any quantifier  $\Rightarrow$  quantified variable

\* variable without quantifier  $\Rightarrow$  free variable.

\* A ~~tuple~~ relational calculus query can have only one free variable.

e.g.  $t \mid \exists d \in \text{emp} \dots$   
 free      quantified.

emp

empno	ename	sal
(a)	(b)	(s)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

### # Domain calculus:

$\Rightarrow e | \langle a, e, s \rangle \in \text{emp} \Rightarrow \text{select ename from emp}$

$\Rightarrow \langle a, e, s \rangle | \langle a, e, s \rangle \in \text{emp} \Rightarrow \text{select * from emp}$

$\Rightarrow e | \langle a, e, s \rangle \in \text{emp} \wedge s > 10000$

$\Rightarrow \text{select ename from emp where sal} > 10000$

Ques Write a query to find the name of employees of Research dept. (in Domain calculus)

emp: ename deptno      dept: deptno dname

(e)	(d)	(dno)	(n)
-----	-----	-------	-----

$\Rightarrow e | \langle e, d \rangle \in \text{emp} \wedge \exists \langle dno, n \rangle \in \text{dept} \wedge$

$(n = \text{'Research'}) \wedge d = dno$

- \* The expressive power of Relational Algebra & Relational Calculus is equal for safe queries.
- \* It is not possible to write unsafe queries in relational Algebra, but in Relational calculus:

eg  $\Rightarrow t \mid t \notin \text{emp}$  (select from all tables except emp table)

- \* It's an unsafe query, & can't be written in Relational Algebra & SQL.

$\rightarrow P \leftarrow$

## # Scalar Subquery :

dept:	deptno	dname	emp:	ename	deptno
	10	Research		e1	10
	20	Design		e2	10

deptno	deptno	dname	ename	deptno
10	10	Research	e1	10
20	20	Design	e2	10
			e3	20
			e4	20
			e5	20

- To display no of employees in every dept.

$\rightarrow$  select dname, (select count(\*) from emp where emp.deptno = dept.deptno) as totalemployee from dept;

- \* It should return only one value, thus, it's known as scalar.

dname	deptno	totalemp
Research	10	2
Design	20	3

## Supplier:

Sno	Pno
S1	P1
S1	P3
S2	P2
S2	P4
S2	P1

## Part:

Pno	Pcolor
P1	red
P2	blue
P3	red
P4	green

eg  $\exists s \in \text{supplier} (\exists p \in \text{Part} \wedge s.pno = p.pno \wedge p.color = \text{red})$   
 (there exists) (some red part)

eg (all red parts) (for all)

$\forall t | \exists s \in \text{Supplier}$

$(\forall p \in \text{Part} \wedge (s.pno = p.pno))$

$\wedge ((p.color \neq \text{red}) \wedge (s.sno = p.sno))$

## Supplier:

Sno	Sname	Pno	Pcolor	Sno
S1	A	P1	red	S1
S2	B	P2	blue	S2
S3	C	P3	red	S1
		P4	green	S3

# NORMALIZATION

67

Date: \_\_\_\_\_ Page No. \_\_\_\_\_

## \* Functional Dependency (FD)

⇒ given a rollno, you will get only one name. Thus,  
it's an FD      ( $\text{Rollno} \rightarrow \text{Name}$ )

⇒ given a name, you will get many rollno. Thus,  
it's not an FD.

\* FDs are decided while creating tables.



eg	X	Y	Z	
	x1	y1	z1	Here, $X \rightarrow Y$ holds ✓
	x1	y1	z2	
	x2	y2	z3	$X \rightarrow Z$ X
	x2	y2	z4	

\*  $Z \rightarrow X$ ,  $Z \rightarrow Y$  ✓

\* Here, z is unique, thus  
it'll always decide other  
attributes.



eg

x	y	z
x1	y1	z1
x1	y2	z2
x2	y1	z3
x2	y2	z4

$\Rightarrow \underline{xy} \rightarrow z, z \rightarrow x, z \rightarrow y, \dots$   
 combination is unique.

 $\rightarrow x$ Ques

x	y	z
x1	y1	z1
x2	y2	z2
x3	y2	z2
x4	y1	z1

Ans: what you can conclude on this database;

- (a)  $x \rightarrow y$  &  $y \rightarrow z$       (b)  $y$  doesn't decide  $x$

$\Rightarrow$  On this instance,

- $x \rightarrow y$  ✓
- $x \rightarrow z$  ✓
- $y \rightarrow x$  ✗
- $y \rightarrow z$  ✓
- $z \rightarrow x$  ✗
- $z \rightarrow y$  ✓
- $xy \rightarrow z$  ✓
- $yz \rightarrow x$  ✗
- $xz \rightarrow y$  ✓

 $\rightarrow x$

## # Laws related to FD :

(1) Reflexivity : an attribute always decides itself.

$$X \rightarrow X$$

(2) Transitivity :

$$\text{If } X \rightarrow Y \text{ & } Y \rightarrow Z \text{ then } X \rightarrow Z$$

(3) Pseudo Transitivity :

$$\text{If } X \rightarrow Y \text{ & } YZ \rightarrow W \text{ then } XZ \rightarrow W$$

(4) Augmentation :

$$\text{If } X \rightarrow Y \text{ then } XZ \rightarrow Y \text{ & } XZ \rightarrow YZ$$

<u>eg</u>	rollno	name	add
	1	n1	a1
	2	n2	a2

If  $\text{rollno} \rightarrow \text{name}$   
then,

$$\text{rollno add} \rightarrow \text{name} \text{ & } \text{rollno add} \rightarrow \text{name add}.$$

⑤ Additivity :

If  $X \rightarrow Y$  &  $X \rightarrow Z$  then  $X \rightarrow YZ$

⑥ Projectivity :

If  $X \rightarrow YZ$  then  $X \rightarrow Y$  &  $X \rightarrow Z$

Note : You cannot break LHS.

$AB \rightarrow C \Rightarrow A \rightarrow C$  &  $B \rightarrow C$  X

\* These rules are used to find closure of attributes.

eg  $R(A, B, C, D)$

$A \rightarrow B$

$B \rightarrow C$

$C \rightarrow D$

\* closure of A contains all those attributes which are decided by A directly or indirectly.

$$\{A\}^+ = \{ABC\}, \quad \{B\}^+ = \{BCD\}$$

$$\{C\}^+ = \{CD\}, \quad \{D\}^+ = \{D\}$$

→

\* The attribute whose closure contains all the attributes of the relation, is known as candidate key. (A here)

$R(A, B, C, D, E) \Rightarrow$  if A is candidate key, how many superkeys?

$$\Rightarrow {}^A_1 + {}^4C_1 + {}^4C_3 + {}^4C_4 = \underline{16} \text{ Ans}$$

Date: \_\_\_\_\_ Page No: \_\_\_\_\_

eg  $R(A, B, C, D)$

$$A \rightarrow B \Rightarrow \{A\}^+ = \{ABCD\}$$

$$B \rightarrow C$$

$$C \rightarrow D \Rightarrow \{B\}^+ = \{BCDA\}$$

$$D \rightarrow A$$

$$\Rightarrow \{C\}^+ = \{CDA\}$$

\* All 4 are candidate keys.

$$\Rightarrow \{D\}^+ = \{DABC\}$$

eg  $R(A, B, C, D, E)$

$$A \rightarrow B$$

$$C \rightarrow D$$

$$\{A\}^+ = \{AB\}, \{B\}^+ = \{B\}, \{C\}^+ = \{CD\}$$

$$\{D\}^+ = \{D\}, \{E\}^+ = \{E\}$$

Here, candidate key must be composite.

\* Those attributes which are only in L.H.S. of FDs, will compulsorily be in candidate key (A, C here)

\* Those attributes which are not a part of any FD, will also compulsorily be in candidate key (E here)

$$\Rightarrow \{ACE\}^+ = \{ACEBD\} \Rightarrow \underline{\text{CK}}$$

R(A, B, C, D, E) if A & B are c.k., how many super keys?

$$\Rightarrow 16 + 16 - 8 = 24 \quad (\text{inclusion exclusion})$$

Date \_\_\_\_\_ Page No. \_\_\_\_\_

eg R(S, C, D, G, V)

$$S \rightarrow C$$

$$SD \rightarrow G$$

$$CG \rightarrow V$$

$$\Rightarrow \{SD\}^+ = \{SDCGV\} \Rightarrow \underline{\text{CK}}$$

eg R(A, B, C, D, E, H)

$$A \rightarrow B$$

$$BC \rightarrow D$$

$$E \rightarrow C$$

$$D \rightarrow A$$

$$\Rightarrow \{EH\}^+ = \{EHC\} \times$$

$$\Rightarrow \{AEH\}^+ = \{AEHBCD\} \Rightarrow \underline{\text{CK}}$$

$$\Rightarrow \{BEH\}^+ = \{BEHCD\} \Rightarrow \underline{\text{CK}}$$

$$\Rightarrow \{CEH\}^+ = \{CEH\} \times$$

$$\Rightarrow \{DEH\}^+ = \{DEHACB\} \Rightarrow \underline{\text{CK}}$$

$R(A, B, C, D, E)$ , max. no. of super keys?

For max. no. of super keys, assume all 5 are C.K.

$$+ SC_1 + SC_2 + SC_3 + SC_4 + SC_5 = 2^5 - 1 \text{ (Ans)}$$

Date: \_\_\_\_\_ Page No. \_\_\_\_\_

eg  $R(A, B, C, D, E, F, G, H)$

$$CH \rightarrow G$$

$$\Rightarrow \{D\}^+ = \{D\}$$

$$A \rightarrow BC$$

$$B \rightarrow CFH \Rightarrow \{AD\}^+ = \{ADBCFHEG\} \Rightarrow CK$$

$$E \rightarrow A$$

$$F \rightarrow EG \Rightarrow \{BD\}^+ = \{BDCFHEGA\} \Rightarrow CK$$

$$\Rightarrow \{ED\}^+ = \{EDABCFCHG\} \Rightarrow CK$$

$$\Rightarrow \{FD\}^+ = \{FDEGABCCH\} \Rightarrow CK$$



eg  $R(A, B, C, D, E, F, G)$

$$AB \rightarrow CD$$

$$\Rightarrow \{ABG\}^+ = \{ABG, CDEF\} \Rightarrow CK$$

$$DE \rightarrow F$$

$$C \rightarrow E$$

$$F \rightarrow C$$



\* Trivial FD: An FD is said to be trivial if,  
 $B \subseteq \alpha$

$$\alpha \rightarrow \beta$$

\* trivial FDs always hold.

eg  $AB \rightarrow B, A \rightarrow A, \text{ etc}$



$R(A, B, C, D, E)$ , max. no. of keys that can be C.K. simultaneously

 $\Rightarrow (A, B, C, D, E) \Rightarrow \sum \times ({}^5C_1)$ 
 $\Rightarrow (AB, AC, AD, AE, BC, BD, BE, CD, CE, DE) = \frac{10}{AT} ({}^5C_2) \text{ also, } ({}^5C_3)$ 

$\Rightarrow$  How to find closure of set of FDs ( $F^+$ ) :

$\Rightarrow R(A, B, C)$

• ~~No. of FD's can be possible.~~

$F = \left\{ \begin{array}{l} A \rightarrow B \\ B \rightarrow C \end{array} \right\}$   $\Rightarrow$  How many FD's can be derived from given FD set.

$\Rightarrow A^+ = \{ABC\}$  \* its closure contains 3 attributes, thus, there will be  $2^3 = \underline{\underline{8}}$  FD's in  $F^+$

$\Rightarrow B^+ = \{BC\} \Rightarrow 2^2 = \underline{\underline{4}}$

$\Rightarrow C^+ = \{C\} \Rightarrow 2^1 = \underline{\underline{2}}$

$\Rightarrow \{AB\}^+ = \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}}$

$\Rightarrow \{BC\}^+ = \{BC\} \Rightarrow 2^2 = \underline{\underline{4}}$

$\Rightarrow \{AC\}^+ = \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}}$

$\Rightarrow \{ABC\}^+ = \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}}$

$\therefore$  total FD's in  $F^+$  =  $8 + 4 + 2 + 8 + 4 + 8 + 8 = \underline{\underline{42}}$

+ 1( $\emptyset \rightarrow \emptyset$ )  
 $= \underline{\underline{43}}$

$F^+ = \left\{ \begin{array}{l} A \rightarrow \emptyset, A \rightarrow A, A \rightarrow B, A \rightarrow C, A \rightarrow AB, A \rightarrow BC \\ A \rightarrow AC, A \rightarrow ABC \end{array} \right. \underline{\underline{43}} \left. \right\}$

If a rel. contains 10 attributes, max. how many keys can be R.K simultaneously. (Q5)

$$\Rightarrow \text{max. value} \rightarrow {}^n C_{(n/2)} \Rightarrow {}^{10} C_5 = \underline{\underline{252}}$$

Ans

Date \_\_\_\_\_ Page No. \_\_\_\_\_

e.g.  $R(A, B)$

$$F = \left\{ \begin{array}{l} A \rightarrow B \\ B \rightarrow A \end{array} \right\}$$

$$\Rightarrow A^+ = \{AB\} \Rightarrow 2^2 = \underline{\underline{4}}$$

$$\Rightarrow B^+ = \{AB\} \Rightarrow 2^2 = \underline{\underline{4}}$$

$$\Rightarrow \{AB\}^+ = \{AB\} \Rightarrow 2^2 = \underline{\underline{4}}$$

$$\therefore \text{total FD's in } F^+ = 4 + 4 + 4 = \underline{\underline{12}} + 1 (\emptyset \rightarrow d) = \underline{\underline{13}}$$

→

Ques  $R(A, B, C, D, E)$

$$A \rightarrow B$$

$$A \rightarrow C$$

$$CD \rightarrow E$$

$$B \rightarrow D$$

$$E \rightarrow A$$

Q: Which of the following is not a part of closure of given FD set?

(a)  $CD \rightarrow AC$  ✓ (b)  $BD \rightarrow CD$

(c)  $BC \rightarrow CD$  (d)  $AC \rightarrow BC$

\* For this, don't find entire closure just find closure of LHS of FD & if RHS is present then it'll be part of closure else not.

Q: Which of the following FD is not implied by above FD set?

Another way to ask this question.

$$\Rightarrow CB^+ = \{C, D, E, A\}, BD^+ = \{B, D\}$$

$$BC^+ = \{B, C, D, E, A\}, AC^+ = \{A, C, B, D\}$$

eg R (A, B, C, D)

$$F = \{ A \rightarrow B \}$$

$$\{ C \rightarrow D \}$$

$$\Rightarrow A^+ = \{AB\} \Rightarrow 2^2 = \underline{4}, \Rightarrow B^+ = \{B\} \Rightarrow 2^1 = \underline{2}$$

$$\Rightarrow C^+ = \{CD\} \Rightarrow 2^2 = \underline{4}, \Rightarrow D^+ = \{D\} = 2^1 = \underline{2}$$

$$\Rightarrow \{AB\}^+ = \{AB\} \Rightarrow 2^2 = \underline{4}, \Rightarrow \{BC\}^+ = \{BC\} \Rightarrow 2^3 = \underline{8}$$

$$\Rightarrow \{AC\}^+ = \{ACD\} \Rightarrow 2^4 = \underline{16}, \Rightarrow \{CD\}^+ = \{CD\} \Rightarrow 2^2 = \underline{4}$$

$$\Rightarrow \{AD\}^+ = \{ADB\} \Rightarrow 2^3 = \underline{8}, \Rightarrow \{BD\}^+ = \{BD\} \Rightarrow 2^2 = \underline{4}$$

$$\Rightarrow \{ABC\}^+ = \{ABCD\} \Rightarrow 2^4 = \underline{16}, \Rightarrow \{BCD\}^+ = \{BCD\} \Rightarrow 2^3 = \underline{8}$$

$$\Rightarrow \{ACD\}^+ = \{ACDB\} \Rightarrow 2^4 = \underline{16}, \Rightarrow \{ABD\}^+ = \{ABD\} \Rightarrow 2^3 = \underline{8}$$

$$\Rightarrow \{ABCD\}^+ = \{ABCD\} \Rightarrow 2^4 = \underline{16}$$

$$\begin{aligned} \text{total FDs in } F^+ &= 4 + 2 + 4 + 2 + 4 + 8 + 16 + 4 \\ &\quad + 8 + 4 + 16 + 8 + 16 + 8 + 16 \end{aligned}$$

$$= 120 + 1 = \underline{\underline{121}}$$

→ If a table contains  $n$  attributes, atmost how many FD's are possible?

e.g.  $R(A, B, C)$

$$\begin{aligned} A^+ &= \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}} \\ B^+ &= \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}} \\ C^+ &= \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}} \\ \{AB\}^+ &= \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}} \\ \{BC\}^+ &= \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}} \\ \{AC\}^+ &= \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}} \\ \{ABC\}^+ &= \{ABC\} \Rightarrow 2^3 = \underline{\underline{8}} \end{aligned}$$

$$\Rightarrow (2^3 - 1) * 2^3 = \underline{\underline{56}}$$

Also,  $\emptyset \rightarrow \emptyset, \emptyset \rightarrow A \dots$  exists.  $\therefore (2^3 - 1 + 1) * 2^3 = \underline{\underline{64}}$

$$\therefore 2^n \times 2^n = \underline{\underline{2^{2n}}} \quad \text{Ans}$$

~~+ —~~

→ How many non trivial FD's are possible?

⇒

$\alpha \rightarrow \beta$   
is trivial  
if  $\beta \subseteq \alpha$

## # Normalization :

primary key (composite)

	Rollno	Name	Cno	Cname	grade	Bookname
	1	n1	c1	DBMS	A	Korth
	1	n1	c2	OS	B	Galvin
	2	n2	c1	DBMS	B	Korth
	3	n1	c3	CN	A	Tanenbaum

① Insertion Anomaly : If we want to add a new course to the university, we cannot because Rollno can't be NULL & initially no student is enrolled in the course. Thus, we are unable to insert the data.

② Deletion Anomaly : If we've to delete Rollno 3, it'll delete the details of CN. Thus, on deleting some data, extra data is getting deleted.

③ Updation Anomaly : Due to data redundancy, we've to update data many times.  
(eg Bookname)

④ Data Redundancy : If there are lots of students, we've to insert coursename for all the students. Thus, the unnecessary repetition of data is known as data redundancy.

\* practically, the desired normal form is 3NF. 79

Date: \_\_\_\_\_ Page No: \_\_\_\_\_

⇒ Thus, to remove these problems, we normalize ~~at the~~ relation. (decompose one table into many tables)

\* there exists various normal forms;

### ① 1NF

\* An information representation is said to be in 1NF if it is atomic.

eg	ename	job	phone
It's info. → (not table)	e1	programmer	9801236012
		analyst	8435501286
		leader	
	e2	salesman	9009012368
		manager	7999720116
			8770401972

\* It cannot be directly stored in a table. Thus, it's not atomic & thus, not in 1NF.

⇒ To convert it into 1NF(table), in one row, per job store row for each phone no.; 12 rows will be in table.

\* Every relation is in 1NF.

## (2) 2NF :

\* A relation is said to be in 2NF if all the non-prime attributes fully functionally depends on a candidate key

OR

\* A relation is said to be in 2NF if there exists no partial FD.

\* prime attribute : An attribute is said to be prime if it is a part of candidate key but not the candidate key itself. & those which are not related to candidate key are called non-prime attributes.

eg  $R(A, B, C, D)$

$$AC^+ = \{A, B, C, D\}$$

$$A \rightarrow B$$

Here, AC is candidate key.

$$C \rightarrow D$$

thus, A & C are prime & remaining B & D are non-prime.

eg  $R(S, C, D, G, V)$

$$S \rightarrow C$$

Here, SD is candidate key.

$$SD \rightarrow G$$

thus, prime  $\Rightarrow S, D$

$$CG \rightarrow V$$

non-prime  $\Rightarrow C, G, V$

eg R (A, B, C, D)

A  $\rightarrow$  B

Here, A is candidate key.

B  $\rightarrow$  C

Thus, there exists no prime attributes.

C  $\rightarrow$  D

non prime  $\Rightarrow$  B, C, D.

\* A is neither prime nor non-prime attribute.

\* If the candidate key is not composite then there will be no prime attributes.

eg R (A, B, C, D)

A  $\rightarrow$  B

Here, candidate key  $\Rightarrow$  A, B, C, D

B  $\rightarrow$  C

Thus, A, B, C & D are neither

C  $\rightarrow$  D

prime nor non-prime attributes.

D  $\rightarrow$  A

$\rightarrow$  —

\* partial FD : An FD is said to be partial if its LHS is prime & RHS is non-prime.

eg R (A, B, C, D)

A  $\rightarrow$  B

Here, AC is candidate key,

C  $\rightarrow$  D

prime  $\Rightarrow$  A, C

non-prime  $\Rightarrow$  B, D

$\Rightarrow$  Both are partial FDs. Thus, the relation is not.

in 2NF.

• For partial dependency,  $\alpha \rightarrow B$ ,  $\alpha$  should be proper subset of CK, but not CK itself.

Date \_\_\_\_\_ Page No. \_\_\_\_\_

eg R (S, C, D, G, V)

$S \rightarrow C$  (partial) Here, candidate key  $\rightarrow SD$

CK:  $SD \rightarrow G$  (not partial) prime  $\Rightarrow S, D$

itself  $CG \rightarrow V$  (not partial) non-prime  $\Rightarrow C, G, V$

$\Rightarrow$  Not in 2NF

$\rightarrow -$

eg R (A, B, C, D)

$A \rightarrow B$  (not partial) Here, candidate key  $\rightarrow A$

$B \rightarrow C$  (not partial) thus, there exists no prime attribute

$C \rightarrow D$  (not partial) non prime  $\Rightarrow B, C, D$

\* A is neither prime nor non-prime.

$\Rightarrow$  In 2NF

$\rightarrow -$

eg R (A, B, C, D, E, F, G, H)

$CH \rightarrow G$  (not partial) Here, candidate key  $\Rightarrow DA, DB,$

$DE, DF$

$A \rightarrow BC$

$B \rightarrow CFH$

prime  $\Rightarrow D, A, B, E, F$

$E \rightarrow A$

non prime  $\Rightarrow C, G, H$

$F \rightarrow EH$

$\Rightarrow A \rightarrow B$  (not partial)

$A \rightarrow C$  (partial)

$\Rightarrow$  Not in 2NF

A4e Part - Supplier Table

Part ( Pno, pname, color, sno, sname )

$Pno \rightarrow pname$

Here, Pno is candidate key.

$Pno \rightarrow color$

Thus, there exists no prime attribute.

$Pno \rightarrow sno$

$sno \rightarrow sname$

non prime  $\Rightarrow$  pname, color,  
sno, sname

$\Rightarrow$  In 2NF

\* Pno is neither prime nor  
non prime.

\* If the candidate key is not composite then the relation will always be in 2NF.



\* Here, all non-prime attributes fully functionally depends on candidate key.

( $sno \rightarrow sname$ ,  $pno \rightarrow sno$ )  
transitively,  $pno \rightarrow sname$

\* There exists insertion, deletion, updation anomaly & data redundancy ; all these problems exist due to transitive FDs.

eg

$Pno \rightarrow sno$



$pno \rightarrow sname$

$sno \rightarrow sname$

### ③ 3NF :

- \* A relation is said to be in 3NF if all non-prime attributes non-transitively depends on candidate key.

OR

- \* A relation is said to be in 3NF if in FD  $\alpha \rightarrow \beta$ , either  $\alpha$  should be super key or  $\beta$  should be prime attribute.

eg  $R(A, B, C, D)$

$$\begin{array}{ll} A \rightarrow B & \checkmark \\ B \rightarrow C & \times \\ C \rightarrow D & \times \end{array}$$

Here, candidate key  $\Rightarrow A$   
non prime  $\Rightarrow B, C, D$

A is neither prime nor non-prime

- \* Here, AB can also be candidate key but it's not minimal.
  - \* every candidate key is a super key, but not viceversa
  - \* superset of candidate key is also a super key.
  - \* minimal super key is candidate key.
- $\Rightarrow$  Not in 3NF.

eg R (A, B, C, D)

$A \rightarrow B$  ✓

Here, candidate key  $\Rightarrow A, B, C, D$ .

$B \rightarrow C$  ✓

~~non prime~~  $\Rightarrow A, B, C, D$  are neither prime nor non-prime.

$C \rightarrow D$  ✓

$D \rightarrow A$  ✓

$\Rightarrow$  In 3NF

$\rightarrow$  —

Que R ( Rollno, Name, Cno, Cname, grade )

$Rollno \rightarrow Name$  ✓

Here, candidate keys  $\Rightarrow$

$Name \rightarrow Rollno$  ✓

$Rollno\ Cno, Name\ Cname,$

$Cno \rightarrow Cname$  ✓

$Rollno\ Cname, Name\ Cno.$

$Cname \rightarrow Cno$  ✓

$Rollno\ Cno \rightarrow grade$  ✓

prime  $\Rightarrow Rollno, Name, Cno, Cname$

non-prime  $\Rightarrow grade$

$\Rightarrow$  In 3NF

$\rightarrow$  —

\* If a relation is in 3NF, all anomalies will be vanished but there can be data redundancy.

\* In questions, start checking from BCNF, & then further, only for those FD's which are violating BCNF.

Date \_\_\_\_\_ Page No. \_\_\_\_\_

#### (4) BCNF:

\* A relation is said to be in BCNF if the determinant (L.H.S. of FD) is super key.

\* If  $\alpha \rightarrow \beta$  is an FD, then  $\alpha$  must be super key.

e.g. R (A, B, C, D)

$A \rightarrow B$  ✓      Here, candidate key  $\Rightarrow A, B, C, D$   
 $B \rightarrow C$  ✓  
 $C \rightarrow D$  ✓  
 $D \rightarrow A$  ✓

$\Rightarrow$  In BCNF

$\rightarrow$  —

e.g. R (A, B, C, D, E); What's the ~~again~~ max. normal form this relation is?

$A \rightarrow B$

$BC \rightarrow D$

Here, candidate key  $\Rightarrow AC$

$CD \rightarrow E$

prime  $\Rightarrow A, C$

non-prime  $\Rightarrow B, D, E$

(a) BCNF X

(b) 3NF X

(c) 2NF X

~~1NF~~ 1NF ✓

$\rightarrow$  —

CS-2004

Ques 50. R ( Rollno, Name, Cno, grade )

Rollno (Cno → grade

Here, candidate keys ⇒

Name (Cno → grade

Rollno (Cno, Name Cno

Rollno → Name

Name → Rollno

prime → Cno, Rollno, Name

non-prime ⇒ grade

• BCNF X

• 3NF ✓

→ Y →

Ques R ( A, B, C, D, E, F, G )

AB → CD

Here, candidate key ⇒ AB

DE → F

C → E

prime ⇒ A, B

F → C

non-prime ⇒ C, D, E, F, G

B → G

• BCNF X

• 3NF X

• 2NF X

• 1NF ✓

→ Y →

- \* Till BCNF, normal forms are based on FDs. But after BCNF, normal forms are not based on FDs but on MVDS.
- \* If a relation is in BCNF, all the anomalies & data redundancy due to FD will be vanished.
- \* Normalization is done to decompose tables & denormalization to merge tables.



⇒	empno	ename	deptno	dname	dcity
1	e1	10	Research	Delhi	
2	e2	10	Research	Delhi	
3	e3	20	Finance	Mumbai	
4	e4	20	Finance	Mumbai	

\* All the anomalies & data redundancy exist here.  
Thus, we'll normalize the table.

⇒	emp:		dept:		
empno	ename	deptno	dname	dcity	
1	e1	10	Research	Delhi	
2	e2	20	Finance	Mumbai	
3	e3				
4	e4				

⇒ But, which emp belongs to which dept. X data-loss

\* Thus, on join, we must be able to create original table.

\* Such decompositions where data is lost, are known as lossy join decomposition. & where data is not lost, known as lossless join decomposition.

⇒ Above decomposition is lossy.

$\Rightarrow$  emp.

empno	ename	dcity
1	e1	Delhi
2	e2	Delhi
3	e3	Delhi
4	e4	Delhi

dept

deptno	dname	dcity
10	Research	Delhi
20	Finance	Delhi

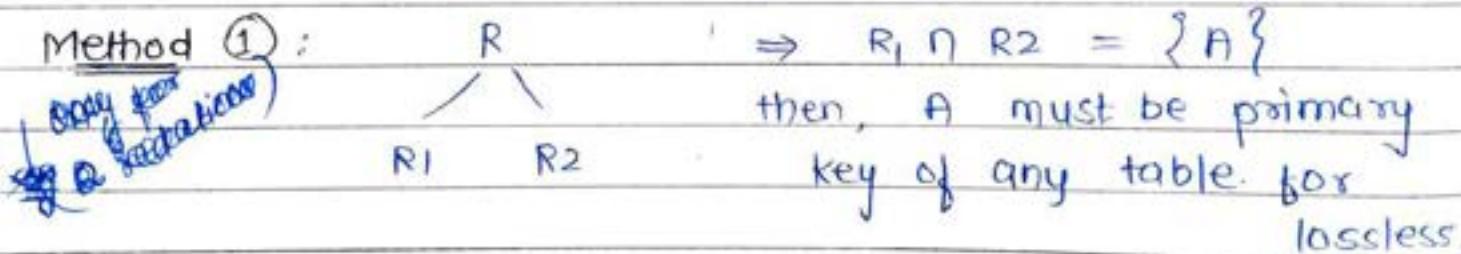
$\Rightarrow$  This decomposition is also lossy.

• We'll get total 8 tuples  
on join.

\* In lossy join, we get some extra tuples known as spurious tuples (cartesian product)

$\Rightarrow$  Identify Lossy or Lossless Decomposition:

Method ① :



eg  $R_1 \cap R_2 \Rightarrow \text{emp} \cap \text{dept} = \underline{\text{dcity}}$

it's not primary key of any table, thus, lossy.

eg  $R_1 \cap R_2 \Rightarrow \text{emp} \cap \text{dept} = \underline{\text{deptno}}$

it's primary key of dept  
thus, lossless.

Apply method ① for only 2 rel. decomposition.

(g1)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

Method ② :

eg  $R(A, B, C, D, E)$

$A \rightarrow B$

$R_1(A, B, C)$

$\Leftrightarrow A \rightarrow B, B \rightarrow C$

$B \rightarrow C$

$C \rightarrow D$

$R_2(B, D, E)$

$\Rightarrow B \rightarrow D$

$B \rightarrow D$

	A	B	C	D	E
R1	✗	✗	✗	✗	
R2		✗	✗	✗	✗

Step ① : Search the column with 2 ✗ (B here)

Step ② : Search the column with 1 ✗ (A, C, D, E here)

Repeat

Now, if  $B \rightarrow C$  then  $R_2(C) = \text{✗}$

if  $B \rightarrow D$  then  $R_1(D) = \text{✗}$

\* If any row can be converted into ✗, then the decomposition is lossless else lossy.

⇒ This decomposition is lossy.



eg previous question with method ①;

$$R_1(A, B, C) \cap R_2(B, D, F) = \underline{\underline{B}}$$

- \* primary key of  $R_1 \Rightarrow A$  it's not primary key of  $R_1$  or  $R_2$ . thus,
  - \* primary key of  $R_2 \Rightarrow BE$  lossy.
- ←

eg  $R(A, B, C, D, E)$

$$A \rightarrow B$$

$$B \rightarrow C$$

$$R_1(A, B, E)$$

$$C \rightarrow D$$

$$B \rightarrow D$$

$$R_2(B, C, D)$$

method ①:  $R_1 \cap R_2 = \underline{\underline{B}} \rightarrow$  it's primary key of  $R_2$ ,  
thus, lossless

- \* primary key of  $R_1 \Rightarrow AE$

- \* primary key of  $R_2 \Rightarrow B$
- ←

FOR  $A \rightarrow D$ , then A & C should have 2 $\alpha$   
& D should have 1 $\alpha$ .

(93)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

②  
method:

	A	B	C	D	E
R1	$\alpha$	$\alpha$	$\alpha$	$\alpha$	$\alpha$
R2		$\alpha$	$\alpha$	$\alpha$	

•  $2\alpha \Rightarrow B$ ,  $1\alpha \Rightarrow A, C, D, E$

\*  $B \rightarrow C \Rightarrow R_1(C) = \alpha$

\*  $B \rightarrow D \Rightarrow R_1(D) = \alpha$

\*

⇒ This decomposition is lossless.



e.g.  $R(A, B, C, D, E)$

$A \rightarrow B$

$R_1(A, B)$

$B \rightarrow C$

$R_2(B, C, D)$

$C \rightarrow D$

$R_3(A, E)$

$B \rightarrow D$

①  
method:  $R_1 \Delta R_2 = B \rightarrow$  its primary key of  $R_2$

$R_2 \cap R_3 = \emptyset$

$R_1 \cap R_3 = A \rightarrow$  its primary key of  $R_1$

Thus, it's lossless.



(2)  
method

	A	B	C	D	E
R1	$\alpha'$	$\alpha'$	$\alpha$	$\alpha$	
R2		$\alpha$	$\alpha$	$\alpha$	
R3	$\alpha$	$\alpha$	$\alpha$	$\alpha$	$\alpha$

- \*  $2\alpha' \rightarrow A, B$ ,  $1\alpha' \rightarrow C, D, E$

~~$R_1 R_2$~~

\*  $B \rightarrow C \Rightarrow R_1(C) = \alpha'$

\*  $B \rightarrow D \Rightarrow R_1(D) = \alpha'$

~~$R_1 R_2 R_3$~~

\*  $A \rightarrow B \Rightarrow R_3(B) = \alpha'$

\*  $B \rightarrow C \Rightarrow R_3(C) = \alpha'$

\*  $B \rightarrow D \Rightarrow R_3(D) = \alpha'$

$\Rightarrow$  This decomposition is lossless.



- \* A decomposition should be lossless & dependency preserving.

\* A decomposition is dependency preserving if

$$F^+ = (F_1^+ \cup F_2^+)$$

e.g.  $R(A, B, C, D, E)$

\* All FD's should exist after normalization.

$$A \rightarrow B$$

$R_1(A, B, E)$

$$B \rightarrow C$$

$$C \rightarrow D$$

$R_2(B, C, D)$

$$B \rightarrow D$$

$$F_1 = \{ A \rightarrow B \}$$

$$F_2 = \{ B \rightarrow C \\ C \rightarrow D \\ B \rightarrow D \}$$

Here,  $F^+ = F_1^+ \cup F_2^+$ , thus, the decomposition is dependency preserving.

→ Identify decomposition is dependency preserving or not

$R(A, B, C, D)$

$$\textcircled{1} \quad A \rightarrow B \quad \checkmark_{R1} \quad R_1(A, B, C)$$

$$\textcircled{2} \quad B \rightarrow C \quad \checkmark_{R1}$$

$$\textcircled{3} \quad C \rightarrow D \quad \cdot \quad R_2(B, D)$$

$$\textcircled{4} \quad B \rightarrow D \quad \checkmark_{R2}$$

$$\emptyset^+ = \emptyset$$

Date \_\_\_\_\_ Page No. \_\_\_\_\_

\* Thus, we've to check only for FD  $c \rightarrow D$  & for this,  
there's an algorithm.

$$R^1 \Rightarrow Z = C \quad (\text{LHS of FD}) \subseteq c \rightarrow D$$

$$Z = Z \cup (R_1 \cap (R_1 \cap Z)^+)$$

$$Z = C \cup (ABC \cap (ABC \cap C)^+)$$

$$Z = C \cup (ABC \cap C^+)$$

$$Z = C \cup (ABC \cap CD)$$

$$Z = C \cup C$$

Also, check for R2.

$\because$  we didn't get RHS of FD, here  
 $Z = C$   $\Rightarrow$  this is not dependency preserving.  
reflex

$$R^2 \Rightarrow Z = C$$

$$Z = C \cup (BD \cap \emptyset^+)$$

$$Z = C$$

-y-

eg  $R(A, B, C, D)$

$A \rightarrow B$  ✓<sub>R1</sub>       $R_1(A, B, C)$

$B \rightarrow C$  ✓<sub>R1</sub>

$A \rightarrow C$  ✓<sub>R1</sub>       $R_2(C, D)$

$B \rightarrow D$  •

$R_1 \Rightarrow Z = B$

$$Z = Z \cup (R_1 \cap (R_1 \cap Z)^+)$$

$$Z = B \cup (ABC \cap BCD)$$

$$Z = BC$$

$R_2 \Rightarrow Z = B$

$$Z = B \cup (CD \cap (CD \cap B)^+)$$

$$Z = B$$

$R_1 \Rightarrow Z = BC$

$$Z = BC \cup (ABC \cap (ABC \cap BC)^+)$$

$$Z = BC \cup BC$$

$$Z = BC$$

\* Here,

value of  $Z$

is updated,

thus check

again for its

new value.

$R_2 \Rightarrow Z = BC$

$$Z = BC \cup (CD \cap (CD \cap BC)^+)$$

$$Z = BC \cup C$$

$$Z = BC$$

$\Rightarrow$  not dependency preserving

\* Repeat this until Z achieves a constant value  
or when we get RHS of FD in Z.

Ques  $R(A, B, C, D)$

$$A \rightarrow B$$

$$A \rightarrow C \quad \checkmark_{R1}$$

$$R1(A, C)$$

$$B \rightarrow C \quad \checkmark_{R2}$$

$$B \rightarrow D \quad \checkmark_{R2}$$

$$R2(B, C, D)$$

$$\stackrel{R1}{\Rightarrow} Z = A$$

$$Z = A \cup (AC \cap (AC \cap A)')$$

$$Z = A \cup AC$$

$$Z = AC$$

$$\stackrel{R2}{\Rightarrow} Z = AC$$

$$Z = AC \cup (AC \cap (AC \cap AC)')$$

$$Z = AC \cup AC$$

$$Z = AC$$

$$\stackrel{P2}{\Rightarrow} Z = A$$

$$Z = A \cup (BCD \cap (BCD \cap A)')$$

$$Z = A$$

$$\stackrel{P2}{\Rightarrow} Z = AC \cup (BCD \cap (BCD \cap AC)')$$

$$Z = AC \cup C$$

$$Z = AC$$

$\Rightarrow$  not dependency preserving

11-2008  
Que $R(A, B, C, D)$  $A \rightarrow B \quad \checkmark_{R1}$  $R1(A, B)$  $B \rightarrow C \quad \checkmark_{R2}$  $C \rightarrow D \quad .$  $R2(B, C)$  $D \rightarrow B \quad \checkmark_{R3}$  $R3(B, D)$ 

$$R1 \Rightarrow Z = BC$$

$$Z = BC \cup (AB \cap (AB \cap C)^+)$$

$$BC = B \cup B \cup B \cup B$$

$$Z = BC$$

$$R2 \Rightarrow Z = C$$

$$Z = C \cup (BC \cap (BC \cap C)^+)$$

$$Z = C \cup BC$$

$$Z = BC$$

$$R3 \Rightarrow Z = C$$

$$Z = C \cup (BD \cap (BD \cap C)^+)$$

$$Z = C \cup B$$

$$R3 \Rightarrow Z = BC$$

$$Z = BC \cup (BD \cap (BD \cap BC)^+)$$

$$Z = BC \cup BD$$

$$Z = BCD$$

 $\Rightarrow$  dependency preserving.

⇒ How to find Minimal cover / Canonical cover of FDs set:

① Convert all FDs to simple.

\* An FD is said to be simple if its RHS contains a single attribute.

eg  $A \rightarrow BC \Rightarrow A \rightarrow B \ \& \ A \rightarrow C$

② Remove extra attributes from LHS of FDs.

eg  $AB \rightarrow C \Rightarrow B \rightarrow C$  (if A is extra)

Here, A will be an extra attribute if the closure of B contains A, & vice versa.

③ Remove extra attributes from RHS of FDs  
OR Remove redundant FDs.

\* An FD is said to be redundant if it can be derived from other FDs.

eg  $A \rightarrow B$   
 $B \rightarrow C$   
 $A \rightarrow C \Rightarrow$  redundant

eg Remove redundant FDs.

- ①  $A \rightarrow B$
- ②  $B \rightarrow C$
- ③  $A \rightarrow C$

\* Find the closure of all FDs without considering <sup>LHSet</sup> itself.

①  $\Rightarrow A^+ = \{AC\} \Rightarrow \text{RHS}(B) \text{ is not present}$   
thus, not redundant

②  $\Rightarrow B^+ = \{B\} \Rightarrow \text{RHS}(C) \text{ is not present}$   
thus, not redundant

③  $\Rightarrow A^+ = \{ABC\} \Rightarrow \text{RHS}(C) \text{ is present.}$   
thus, redundant FD.

\* Once you say that an FD is redundant, then  
don't consider that redundant FD while computing  
redundancy of other FDs.



que Find minimal cover:

- ①  $A \rightarrow B$
- ②  $A \rightarrow C$
- ③  $A \rightarrow D$
- ④  $A \rightarrow H$
- ⑤  $DH \rightarrow A$
- ⑥  $DH \rightarrow B$
- ⑦  $DH \rightarrow C$
- ⑧  $DH \rightarrow G$

Step ①: All FDs are simple

Step ②: There's no extra attribute on LHS of FDs.

Step ③:  $(1) \Rightarrow A^+ = \{A, \underline{B}, \underline{C}, \underline{D}, \underline{H}, \underline{B}, \underline{G}\} \Rightarrow$  redundant

$(2) \Rightarrow A^+ = \{A, \underline{B}, \underline{D}, \underline{H}, \underline{C}, \underline{G}\} \cap \{ADHBCG\} \Rightarrow$  redundant

$(3) \Rightarrow A^+ = \{AH\} \Rightarrow$  not redundant

$(4) \Rightarrow A^+ = \{AD\} \Rightarrow$  not redundant

$(5) \Rightarrow (DH)^+ = \{DH, BC, G\} \Rightarrow$  not redundant

$(6) \Rightarrow (DH)^+ = \{DH, AC, G\} \Rightarrow$  not redundant

$$\textcircled{7} \Rightarrow (DH)^+ = \{ DHABC_G \} \Rightarrow \text{not redundant}$$

$$\textcircled{8} \Rightarrow (DH)^+ = \{ DHABC \} \Rightarrow \text{not redundant}$$

→  $\times$

⇒ previous question from bottom:

$$\textcircled{8} \Rightarrow \{ DH \}^+ = \{ DHABC \} \Rightarrow \text{not redundant}$$

$$\textcircled{7} \Rightarrow \{ DH \}^+ = \{ DHABC_G \} \Rightarrow \text{redundant}$$

$$\textcircled{6} \Rightarrow \{ DH \}^+ = \{ DHAG_B_C \} \Rightarrow \text{redundant}$$

$$\textcircled{5} \Rightarrow \{ DH \}^+ = \{ DHG_I \} \Rightarrow \text{not redundant}$$

$$\textcircled{4} \Rightarrow A^+ = \{ ABCD \} \Rightarrow \text{not redundant}$$

$$\textcircled{3} \Rightarrow A^+ = \{ ABCH \} \Rightarrow \text{not redundant}$$

$$\textcircled{2} \Rightarrow A^+ = \{ ABDH_G \} \Rightarrow \text{not redundant}$$

$$\textcircled{1} \Rightarrow A^+ = \{ ACDH_G \} \Rightarrow \text{not redundant}$$

\* Thus, there can be multiple minimal covers.

⇒ Now, these 2 minimal covers are equivalent or not?

If  $F_1 \subseteq F_2$  &  $F_2 \subseteq F_1$ ,

then  $F_1 \equiv F_2$

$F_1$	$F_2$
$A \rightarrow D$ ✓	$A \rightarrow B$ ✓
$A \rightarrow H$ ✓	$A \rightarrow C$ ✓
$DH \rightarrow A$ ✓	$A \rightarrow D$ ✓
$DH \rightarrow B$ ✓	$A \rightarrow H$ ✓
$DH \rightarrow C$ ✓	$DH \rightarrow A$ ✓
$DH \rightarrow G$ ✓	$DH \rightarrow G$ ✓

Thus,  $F_1 \equiv F_2$ .



\*  $F_1 \subseteq F_2 \Rightarrow$  means, all the FDs of  $F_1$  can be derived from FDs of  $F_2$ .



\*  $DH \rightarrow B$  is in  $F_1$  but not in  $F_2$ ,

$$\Rightarrow \text{In } F_2, \{DH\}^+ = \{DHAGBCD\}$$

thus, it can be derived

\* Do this for all FDs. in  $F_2$



Ques Find minimal cover of following FD set;

- ①  $CH \rightarrow G$
- ②  $A \rightarrow BC$
- ③  $B \rightarrow CFH$
- ④  $E \rightarrow A$
- ⑤  $F \rightarrow EG$
- ⑥  $A \rightarrow CG$

Step ① :  $CH \rightarrow G$

- $A \rightarrow B, A \rightarrow C$
- $B \rightarrow C, B \rightarrow F, B \rightarrow H$
- $E \rightarrow A$
- $F \rightarrow E, F \rightarrow G$
- $A \rightarrow C, A \rightarrow G$

Step ② : there's no extra attribute on LHS of FDs

- Step ③ :
- ①  $CH \rightarrow G$
  - ②  $A \rightarrow B$
  - ③  $A \rightarrow C$
  - ④  $B \rightarrow C$
  - ⑤  $B \rightarrow F$
  - ⑥  $B \rightarrow H$
  - ⑦  $E \rightarrow A$
  - ⑧  $F \rightarrow E$
  - ⑨  $F \rightarrow G$
  - ⑩  $A \rightarrow G$

$$\textcircled{1} \Rightarrow \{CH\}^+ = \{CH\} \Rightarrow \text{not redundant}$$

$$\textcircled{2} \Rightarrow A^+ = \{ACG\} \Rightarrow \text{not redundant}$$

$$\textcircled{3} \Rightarrow A^+ = \{ABG\underline{C}FHE\} \Rightarrow \text{redundant}$$

$$\textcircled{4} \Rightarrow B^+ = \{BFHEGA\} \Rightarrow \text{not redundant}$$

$$\textcircled{5} \Rightarrow B^+ = \{BCHG\} \Rightarrow \text{not redundant}$$

$$\textcircled{6} \Rightarrow B^+ = \{BCFEGA\} \Rightarrow \text{not redundant}$$

$$\textcircled{7} \Rightarrow E^+ = \{E\} \Rightarrow \text{not redundant}$$

$$\textcircled{8} \Rightarrow F^+ = \{FG\} \Rightarrow \text{not redundant}$$

$$\textcircled{9} \Rightarrow F^+ = \{FEAG\underline{BC}\} \Rightarrow \text{redundant}$$

$$\textcircled{10} \Rightarrow A^+ = \{ABC\underline{FHEG}\} \Rightarrow \text{redundant}$$

→ p →

QueF<sub>1</sub>

$$AB \rightarrow C \checkmark$$

$$D \rightarrow E \checkmark$$

$$E \rightarrow C \checkmark$$

F<sub>2</sub>

$$AB \rightarrow C \checkmark$$

$$D \rightarrow E \checkmark$$

$$AB \rightarrow E$$

$$E \rightarrow C \checkmark$$

\*  $F_1 \subset F_2$  but  $F_2 \not\subset F_1$

thus,  $F_1 \neq F_2$

→

⇒ How to convert a relation into 3NF :

Step ① Find minimal cover of FD set.

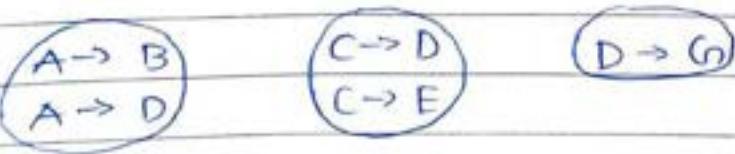
Step ② Find all candidate keys of relation.

Step ③ Create the relations with the attributes of FD having same LHS.

eg       $A \rightarrow B$                    $R(A, B, C, D, E, G)$   
 $A \rightarrow D$   
 $C \rightarrow D$   
 $C \rightarrow E$   
 $D \rightarrow G$

⇒ Step ① : The FDs set is already minimal.

Step ② : candidate key  $\Rightarrow AC$

Step ③ : 

$R_1(A, B, D)$        $R_2(C, D, E)$        $R_3(D, G)$

\* It's a lossy decomposition, thus, there's one more step to make it lossless.

- this step is performed only when we've composite
- Step ④ If candidate key is not a part of any relation then create a relation of that candidate key.

Thus, R<sub>1</sub> ( A, B, D )  
 R<sub>2</sub> ( C, D, E )  
 R<sub>3</sub> ( D, G )  
 R<sub>4</sub> ( A, C )

Now, it's a lossless decomposition.

\* If there are more than 1 CK then create relation for each CK.

\* 3NF decomposition is always lossless join & dependency preserving.

\* BCNF is always lossless but may / may not be dependency preserving.

⇒ How to convert a relation into BCNF :

Step ① Find the minimal cover of FD set.

Step ② Identify the FDs which violates BCNF

Step ③ Create a relation with attributes of FD which violates BCNF. (one by one)

e.g. R (A, B, C, D)

$$A \rightarrow B$$

$$B \rightarrow C$$

$$C \rightarrow D$$

Step ① FD set is already minimal.

Step ② candidate key  $\Rightarrow A$

$$\begin{matrix} B \rightarrow C \\ C \rightarrow D \end{matrix} \quad \left\{ \begin{matrix} \text{violates BCNF} \\ \text{violates BCNF} \end{matrix} \right.$$

e.g.

$$\begin{matrix} B \rightarrow C \\ C \rightarrow D \\ C \rightarrow E \end{matrix} \quad \left\{ \begin{matrix} \text{violates BCNF} \\ \text{violates BCNF} \end{matrix} \right.$$

then create a rel. for (C, D, E)  
not for (C, D) & (C, E)

Step ③ R (A, B, C, D)

\* taking  $C \rightarrow D$  first :

$$R_1 (A, B) \quad \cancel{\&} (C, D)$$

\* cut their RHS from original relation.

R2 (B, C) → create this rel. only if B & C are in original rel.

∴ R1 (A, B), R2 (C, D), R3 (B, C)

\* taking  $B \rightarrow C$  first

$R (A, B, C, D)$

$R_1 (B, C)$

\* Now, we will not make a relation for  $C \rightarrow D$  as  
C is not an attribute of original relation anymore.

$R$  contains  $(A, B, D)$  &  ~~$C$~~

But, now, candidate key of  $R$  is  $AD$ , so, we'll  
check again if some FDs are violating BCNF.

$A \rightarrow B$	}	violates
$B \rightarrow C$		BCNF
$C \rightarrow D$		

$\Rightarrow R_2 (A, B)$  & cut  $B$  from  $R$

(II)\*  $R_1 (A, D)$      $R_2 (B, C)$      $R_3 (A, B)$

$\Rightarrow$  Both decompositions are valid & lossless but  
II is not dependency preserving.

\* Here,  $R_1 (A, D) \rightarrow$  no FD for  $R_1 \rightarrow$  is it in BCNF?  
 $R_2 (B, C) \rightarrow B \rightarrow C \checkmark$  BCNF  
 $R_3 (A, B) \rightarrow A \rightarrow B \checkmark$  BCNF

yes, if all attr  
make CK then

it's in BCNF

OR

Also, it has 2 attr

Thus, individual rel. are in BCNF.

\* Any relation with only 2 attributes will always be in BCNF.

e.g. R (A, B)

\* FDs can be;

$$\begin{array}{lll} \cdot \underline{\underline{A}} \rightarrow B & \text{OR} & \cdot \underline{\underline{B}} \rightarrow A \\ \text{super key} & \text{super} & \text{super} \\ \cdot \underline{\underline{A}} \rightarrow B & \text{OR} & \cdot \underline{\underline{B}} \rightarrow A \\ & & \text{super} \end{array}$$

Here, CK  $\Rightarrow$  A

Thus, it's in BCNF

Here, CK  $\Rightarrow$  B

Thus, it's in BCNF

super

Here, CK  $\Rightarrow$  A, B

$\rightarrow$

Thus, it's in BCNF

\* A relation where no FD exists then in that case, all the attributes together form candidate key & that relation is always in BCNF.

$\rightarrow$

\* If a decomposition is lossless & individual relations are in BCNF then that decomposition is a valid BCNF decomposition (although not returned by our algorithm).

$\rightarrow$

\* In 3NF, an attribute may depend transitively on candidate key.

In 3NF, all non prime attributes non-transitively depend on candidate key but the prime attributes may or may not depend transitively on candidate key.

$\underline{\underline{A}} \rightarrow \underline{\underline{B}}$   
super key / prime attribute

- \* In BCNF, an attribute cannot depend transitively on candidate key.

$$\begin{array}{c} \alpha \rightarrow \beta \\ \text{super key} \\ \longrightarrow \gamma \end{array}$$

- \* If a relation is in 3NF & there exists only one candidate key, then the relation is also ~~also~~ in BCNF.

⇒ If there's no composite candidate key, there will be no prime attributes, then relation will also be in BCNF.

eg R(A, B, C, D)

$$AC \rightarrow D \quad \checkmark_{\text{BCNF}}$$

$$CK = AC, BC$$

$$BC \rightarrow D \quad \checkmark_{\text{BCNF}}$$

$$\text{prime} \Rightarrow A, B, C$$

$$A \rightarrow B \quad \checkmark_{\text{3NF}}$$

$$\text{non-prime} \Rightarrow D$$

$$B \rightarrow A \quad \checkmark_{\text{3NF}}$$

Thus, it's in 3NF, but not in BCNF

- \* In  $\alpha \rightarrow \beta$ , if  $\alpha$  is not super key &  $\beta$  is prime then  $\alpha$  will also be prime attribute but from another candidate key. This type of FDs are possible only when we've more than 1 composite candidate key.

- \* Thus, if there's only 1 key, then in 3NF, all the FD's  $\alpha \rightarrow \beta$ ,  $\alpha$  will be super key, thus, it will always be in BCNF.

## # 4NF :

\* It is based on Multivalued Dependency (MVD).

\* For a given  $\alpha$ , you can have one or more  $\beta$ .

$$\alpha \rightarrow\!\!\! \rightarrow \beta$$

\* every FD is an MVD; but not vice-versa.

\* MVD that is not FD is true or pure MVD.

$\Rightarrow$  Course:

Course	Instructor	Author
DBMS	I1	Korth
DBMS	I2	Navathe
OS	I1	Tanenbaum
NW	I2	Tanenbaum
OS	I1	Galvin
NW	I3	Tanenbaum
NW	I1	Tanenbaum

\* Here, all the attributes combinedly form the candidate key. Thus, it's in BCNF.

$\rightarrow$

$\Rightarrow$	Cname	Instructor	Author
	DBMS	I1	Korth
	DBMS	I2	Navathe
	DBMS	I1	Navathe
	DBMS	I2	Korth

\* there exists some MVDs.;

- ① Cname  $\rightarrow\!\!\!\rightarrow$  Author
- ② Cname  $\rightarrow\!\!\!\rightarrow$  Instructor



\* If  $\alpha \rightarrow\!\!\!\rightarrow \beta$  is an MVD then  $\alpha \rightarrow\!\!\!\rightarrow R - (B \cup \alpha)$  is also an MVD.

e.g. R (A, B, C, D)

if A  $\rightarrow\!\!\!\rightarrow$  B holds  
then A  $\rightarrow\!\!\!\rightarrow$  ABCD - BA  
A  $\rightarrow\!\!\!\rightarrow$  CD will also hold.



\* In FDs,  $A \rightarrow B = A \rightarrow BC$   
 $A \rightarrow C$

In MVDs,  $A \rightarrow\!\!\!\rightarrow B \neq A \rightarrow\!\!\!\rightarrow Bc$   
 $A \rightarrow\!\!\!\rightarrow C$

Name  $\rightarrow$  Author  
Name  $\rightarrow$  Instructor

\* It means that, there's no rel. b/w author & instructor.

(17)

Page No. \_\_\_\_\_

e.g. Name  $\rightarrow$  Author Instructor \* It means that, there's relationship b/w author & instructor.

$\rightarrow$  —

$\Rightarrow$	A	B	C	$\Rightarrow$	A	B	C
	a1	b1	c1		a1	b1	c1
	a1	b2	c2		a1	b2	c2
					a1	b1	c2
• A $\rightarrow\!\!\rightarrow$ BC		✓			a1	b2	c1
• A $\rightarrow\!\!\rightarrow$ B		x					
• A $\rightarrow\!\!\rightarrow$ C		x		• A $\rightarrow\!\!\rightarrow$ BC	x		
				• A $\rightarrow\!\!\rightarrow$ B	✓		
				• A $\rightarrow\!\!\rightarrow$ C	✓		

$\rightarrow$  —

$\Rightarrow$	A	B	C	$\Rightarrow$	A	B	C
	a1	b1	c1		a1	b1	c1
	a1	b2	c2		a1	b2	c2
	a1	b3	c1		a1	b3	c1
					a1	b1	c2
• A $\rightarrow\!\!\rightarrow$ BC		✓			a1	b2	c1
• A $\rightarrow\!\!\rightarrow$ B		x			a1	b3	c2
• A $\rightarrow\!\!\rightarrow$ C		x		• A $\rightarrow\!\!\rightarrow$ BC	x		
				• A $\rightarrow\!\!\rightarrow$ B	✓		
				• A $\rightarrow\!\!\rightarrow$ C	✓		

$\rightarrow$  —

MVD is also known as tuple generating dependency.

- \* A relation is said to be in 4NF if  $\alpha \rightarrow\!\!\! \rightarrow \beta$  is not a trivial MVD then  $\alpha$  must be super key.
- \* <sup>if</sup> there should not be any true MVDs, only FDs are allowed.

$\alpha \rightarrow\!\!\! \rightarrow \beta$  \* if  $\alpha$  is super key,

$\beta$  will have exactly one value.  
i.e. it will be an FD.

- \* An MVD is said to be trivial, if in MVD  $\alpha \rightarrow\!\!\! \rightarrow \beta$

$$\beta \subseteq \alpha \text{ or } \alpha \cup \beta = R$$

eg	Cname	Author
	DBMS	Korth
	DBMS	Navathe

$\Rightarrow Cname \rightarrow\!\!\! \rightarrow Author$  is a trivial MVD.

•  $Cname \cup Author = R$  ( $Cname, Author$ )

→

- \* Trivial MVDs always hold.

# HASHING & INDEXING

Date \_\_\_\_\_ Page No. \_\_\_\_\_

- \* considers a harddisk whose block size is 512 bytes
- & one record takes 100 bytes. in HD storage.

1	ABC	100 bytes
2	DEF	100 bytes
3	GHI	100 bytes
4	JKL	100 bytes
5	MNO	100 bytes
		12 bytes

- \* The file organisation which leads to wastage of storage is unspanned organisation.
- \* The file organisation which doesn't lead to wastage of storage is spanned organisation.
- ⇒ In unspanned organisation, the next record is stored in another block & those 12 bytes will be wasted.
- ⇒ In spanned organisation, a part of next record is stored in same block & the remaining 88 bytes in another block.
- \* Unspanned file organisation is generally used to provide faster access i.e. data will be accessed from one block.

\* For OS, database is one file. The files i.e. tables in database are stored by database itself, it's hidden from OS.



Ques: Consider block size is 512 bytes, record size is 100 bytes. How many blocks are required to store the file if the file contains 4000 records & if file organisation is ; ① spanned ② unspanned

⇒ ① spanned :

$$\begin{aligned} \text{No. of blocks} &= \frac{\text{file size}}{\text{block size}} = \frac{4000 \times 100}{512} \\ &= 781.25 \approx \underline{782} \quad \text{Ans.} \end{aligned}$$

⇒ ② unspanned :

$$\begin{aligned} \text{No. of records} &= \frac{512}{100} = 5.12 = 5 \\ \text{in a block} & \end{aligned}$$

$$\begin{aligned} \text{no. of blocks} &= \frac{4000}{5} = \underline{800} \quad \text{Ans} \end{aligned}$$



- ⇒ select ename from emp where empno = 100;
- \* If linear search is used ; and there are 800 blocks & 1 us is required to read a block , then,
  - \* in worst case ( if record is stored in last block)
 
$$\Rightarrow 800 \times 1 \text{ us} = \underline{800 \text{ us}} \text{ is required.}$$
  - ⇒ Thus , linear search is not feasible in the database.
  - \* Binary search is also not feasible because there are frequent insertions , updations & deletions in the database. (sorting is not possible).
  - ⇒ Thus , to speedup searching within the database , indexing & hashing are introduced.

→

Ques There are records of 200 students , block size is 512 bytes , record size is 100 bytes & file organisation is unspanned . If linear search is performed , how many blocks of harddisk database will read in worst case ?

$$\Rightarrow \frac{200}{5} = 40 \text{ blocks} \quad \text{Ans}$$

→

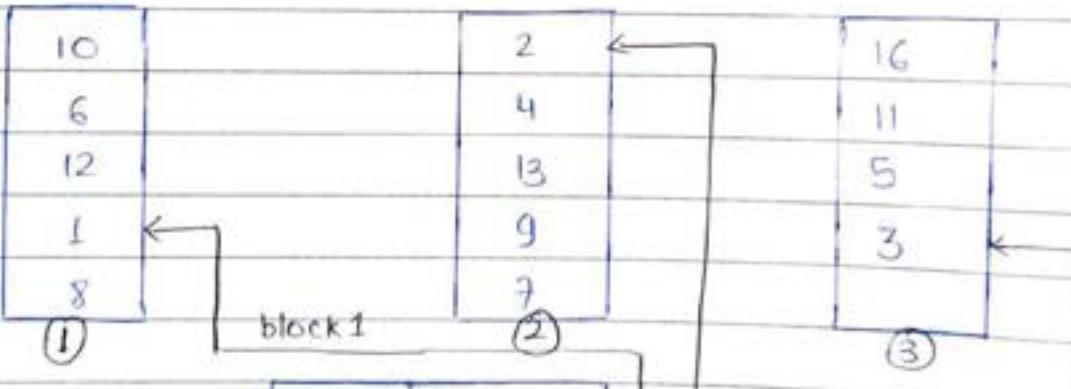
# Indexing:

- \* In indexing, we create a separate file known as index file.
- \* index file has two columns;

6

- |                                       |   |
|---------------------------------------|---|
| ① key <del>field</del><br>primary key | ② record <del>pointer</del><br>address of record<br>in HD |
|---------------------------------------|---|

eg



Index file :    1    ←    block 2

2

3

4

5

block 3

→ block 2

→ block 3

1

1

1

- \* index file is sorted on key field.

\* take unspanned if not given in que

Date \_\_\_\_\_ Page No. \_\_\_\_\_

\* no. of records = no. of records  
in index file in database file

Que Consider a database, record size is 100 bytes, key field is 4 bytes long, record pointer is 4 bytes, unspanned file organisation. How many blocks are required to store database & index file, if block size is 512 bytes & there are 200 records in the database.

$$\Rightarrow \frac{200}{5} = 40 \text{ blocks for database file.}$$

$$\Rightarrow 200 * (4+4) = 1600 \text{ bytes}$$

$$\Rightarrow \frac{1600}{512} = 3.125 \approx 4 \text{ blocks for index file}$$

OR

$$\Rightarrow \frac{\text{no. of records}}{\text{record size}} = \frac{\text{block size}}{\text{record size}} = \frac{512}{8} = 64$$

*(is called blocking factor)*

$$\text{no. of blocks} = \frac{200}{64} = 3.125 \approx 4 \text{ blocks for index file}$$

→ —

\* In previous question,

⇒ Indexing : Best case : 2 blocks have to be read  
Worst case : 5 blocks have to be read

⇒ If indexing was not done;

Best case : 1 block has to be read

Worst case : 40 blocks have to be read



\* If a ~~block~~ record is not present in harddisk i.e.  
Unsuccessful Search:

⇒ Indexing : Best Case : 1 block has to be read.

Worst Case : 4 blocks have to be read.

⇒ If indexing was not done;

Best case: 40 blocks have to be read

Worst case: 40 blocks have to be read.



\* If the index file is sorted, & linear (sequential) search is used then the file org is known as indexed sequential file organisation.

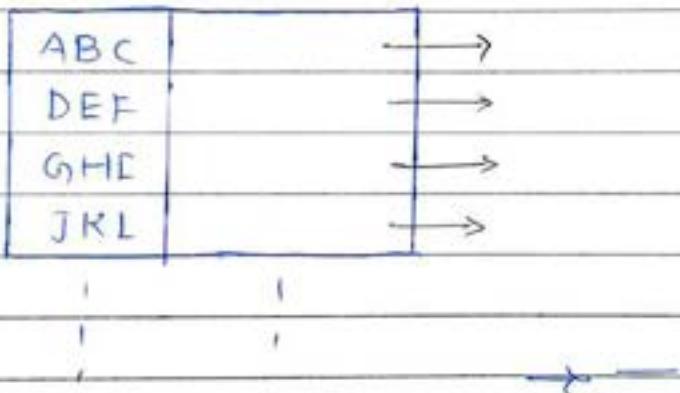
## Drawbacks of Indexing :

①

⇒ select rollno from student where name = "ABC";

⇒ select name from student where rollno = 1;

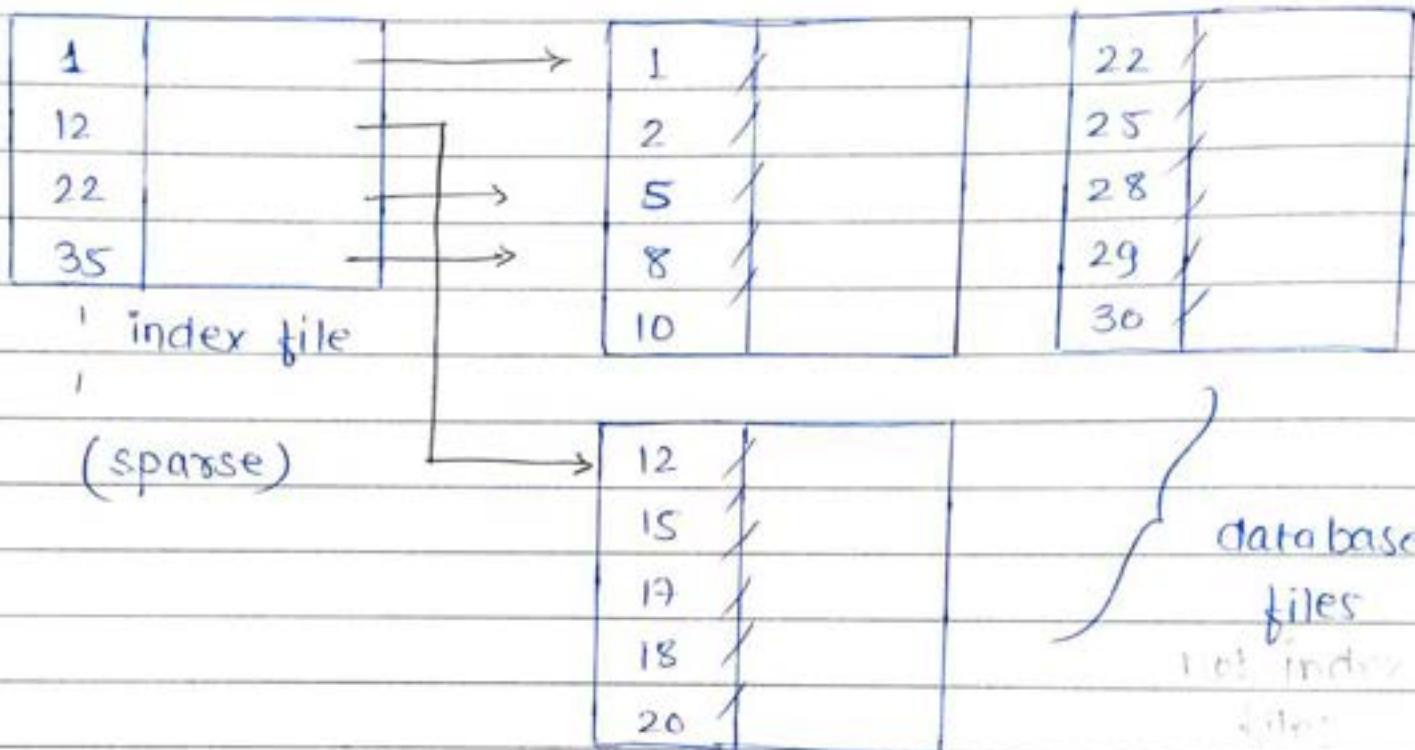
\* Thus, we have to make index files for all attributes.



② If the database is very large then its index file will also contain a lot of records & will be stored in large no. of blocks & in worst case, all these blocks have to be read.

⇒ Thus, how to speedup searching in an index file

## ① Minimize Index file :



- \* But for this, database must be ~~decreasing~~ ~~increasingly~~ sorted &
- \* The index file in which there's an entry for every value (rollno) is known as dense index file.
- \* The index file in which there's an entry for some value (rollno) is known as sparse index file.
- \* database must be sorted for sparse index.

- \* In dense index;

- \* In sparse index;

- \* no. of dense index file = no. of attributes & their combinations

- \* no. of sparse = atmost one ; as database can be sorted on any one column.

-x

If database is not sorted, index will be secondary or non-clustered index (always dense)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

\* Primary Index : The index on a unique field on which the database is sorted is primary index (eg. rollno)

\* It's ~~should be~~ a sparse index as the records are sorted.

\* Cluster Index :

\* If the database is sorted on a non-unique field (eg. deptno)

10		10	
20		10	
40		10	

sparse

cluster index  
Korth (dense)  
Noname (sparse)

↓

30 is not  
present  
in Index.

20		20	
20		20	
20		30	
30		30	

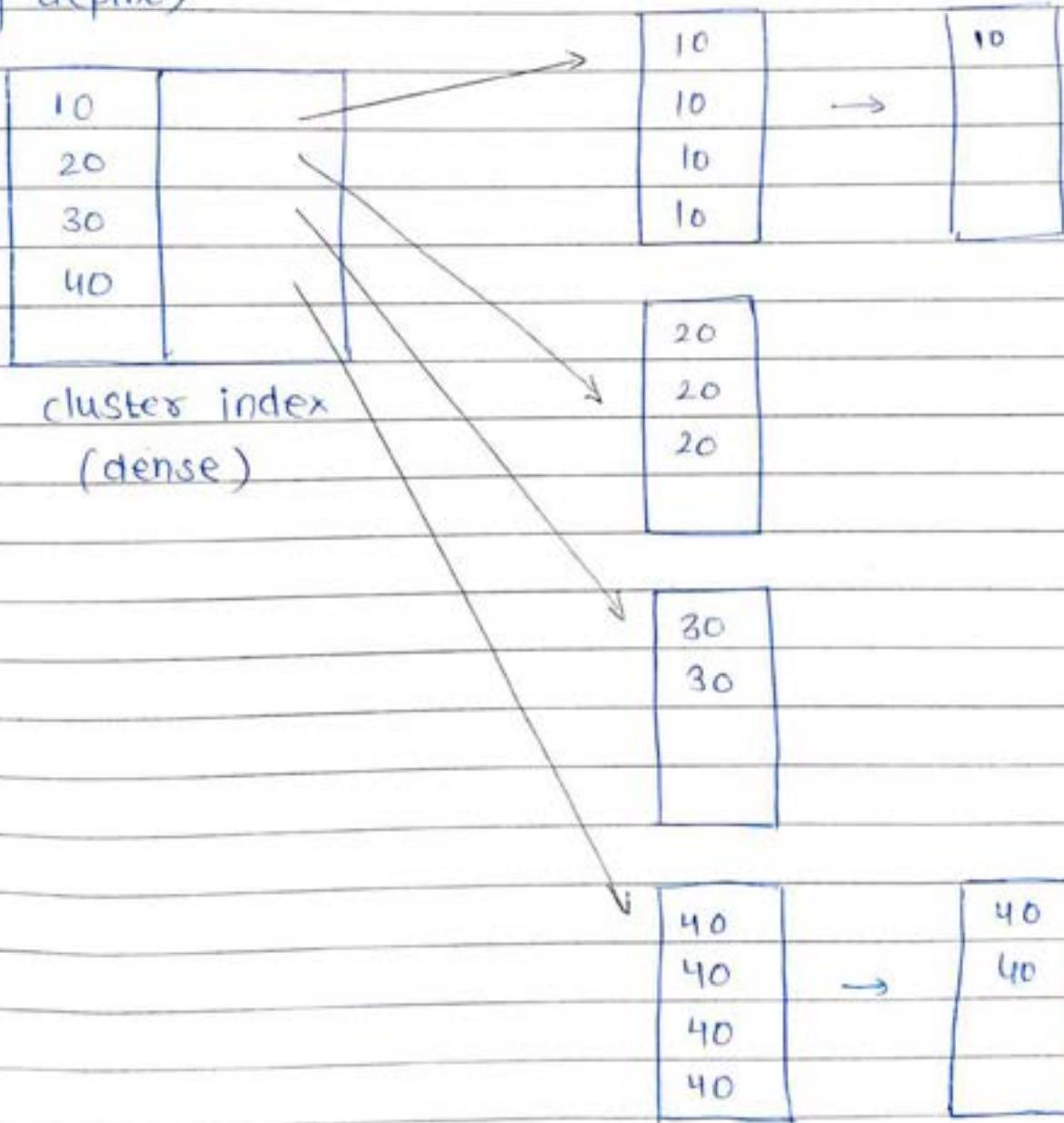
40		40	
40		40	
1		1	
1		1	

problem  $\Rightarrow$  select \* where deptno = 20

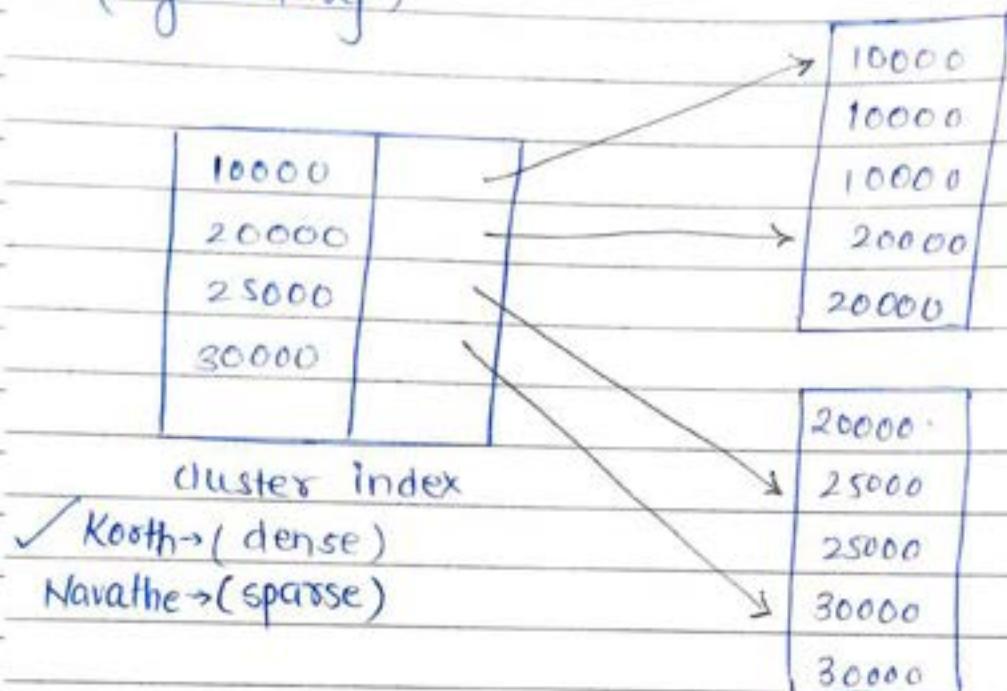
\* only 3 records will be printed, block 1 is not read.

solution:

① If there are more repetitions in column (eg deptno)



(2) if there are less repetitions in column  
(eg salary)



Kothth → (dense)  
Navathe → (sparse)

\*

30000

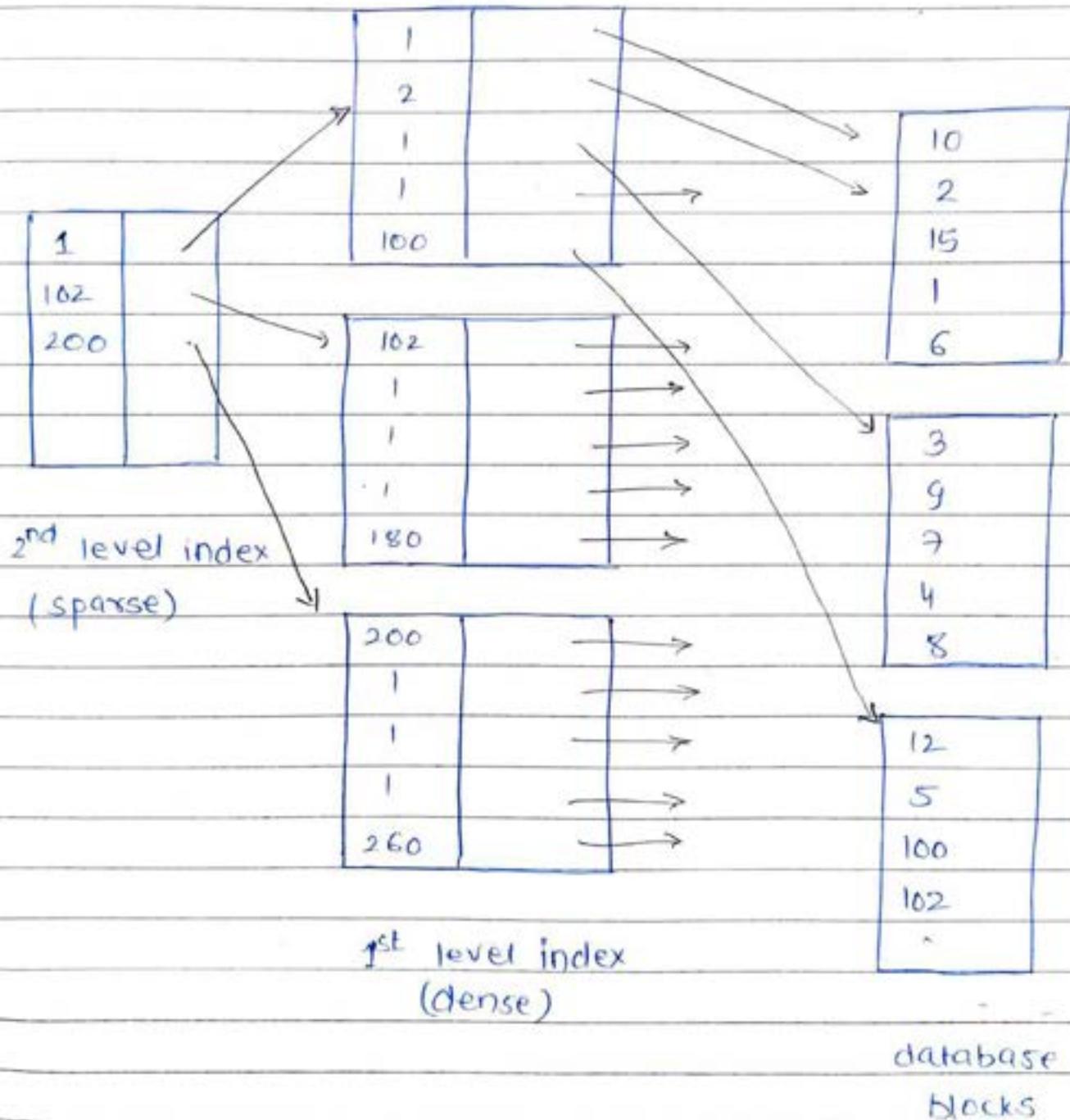
\* These can be only 1 sparse index, the remaining index(s) have to be dense. Thus, how to speedup searching in dense index;

① Multilevel Index

② Tree Index

## # Multi level Index :

- \* First, we create dense index of database file & then well create index to index file.



\* 2<sup>nd</sup> level index will be sparse as 1<sup>st</sup> level index is sorted.

$$\Rightarrow \text{no. of records in } 2^{\text{nd}} \text{ level index} = \text{no. of blocks in } 1^{\text{st}} \text{ level index}$$

$$\Rightarrow \text{no. of records in } 1^{\text{st}} \text{ level index} = \text{no. of records in database}$$

Ques Consider a database file consisting of 4000 records & record size is 100 bytes, 2-level indexing is used where 1<sup>st</sup> level is dense index on rolling & size of rolling is 4 bytes, record pointer is also 4 bytes long. How many blocks are required to store 1<sup>st</sup> level & 2<sup>nd</sup> level index if block size is 512 bytes.

~~$$\Rightarrow \text{no. of records in a block} = \frac{512}{100} = 5$$~~

~~$$\begin{aligned} \text{no. of blocks req. to store database} &= \frac{4000}{5} = 800 \\ * \text{index record size} &= 4 + 4 = 8 \text{ bytes} \end{aligned}$$~~

$$\text{no. of index records in one block} = \frac{512}{8} = 64$$

→ no of blocks req. =  $\frac{4000}{64} \approx \underline{\underline{63}}$   
to store index of  
1<sup>st</sup> level

\* 1 block for 2<sup>nd</sup> level index.

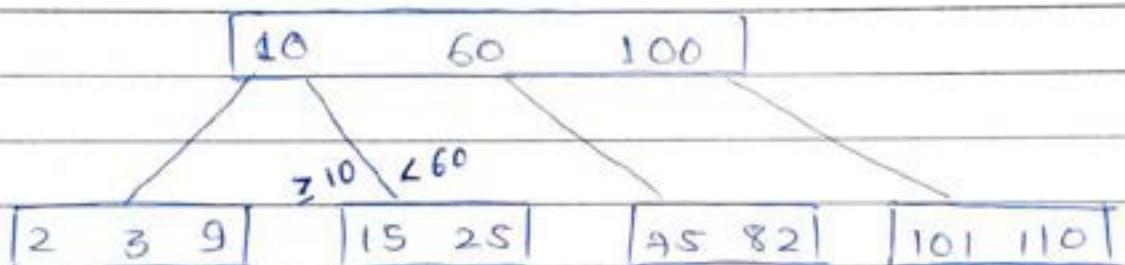


# Tree Index :

# B-trees : They are multiway search tree or n-way search tree in which per node there can be more than 1 key.

- \* within the nodes, keys are sorted.

eg



- \* Before creating tree, order of tree is fixed i.e how many children a node can have at maximum.
- \* if order = n, there will be almost n children & there is n-1 keys



\* B-tree is a multiway search tree in which ;

① all the leaf nodes should be at same level.

② Every non-root node should have atleast  $\lceil \frac{n}{2} \rceil$

children , where  $n = \text{order}$  &

min. keys will be  $\lceil \frac{n}{2} \rceil - 1$ .

eg B-tree of order = 3

max. children = 3 , min children = 2

max. keys =  $3-1 = 2$  , min keys =  $\lceil \frac{3}{2} \rceil - 1 = 1$

eg B-tree of order = 10

max. children = 10 , min children = 5

max. keys =  $10-1 = 9$  , min keys =  $\lceil \frac{10}{2} \rceil - 1 = 4$

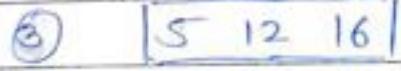
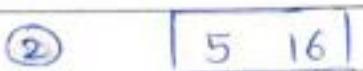
→

Que B-tree of order 3

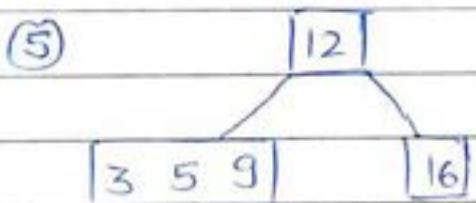
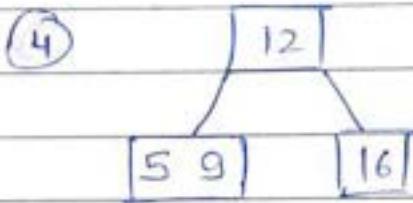
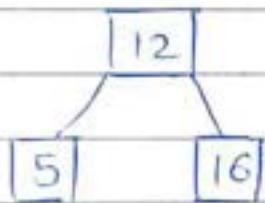
(max. key = 2, min key = 1)

→ 16, 5, 12, 9, 3, 10, 2, 14, 20, 6

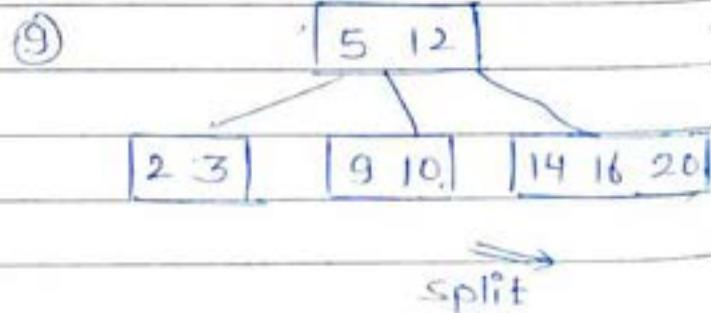
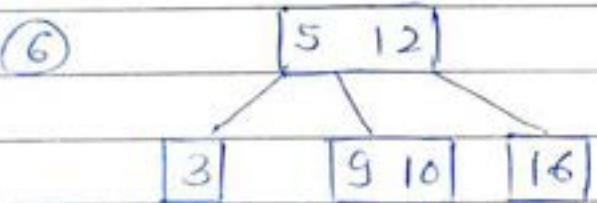
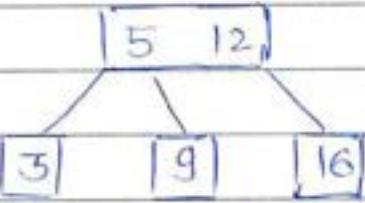
\* Always perform insertion in leaf nodes.

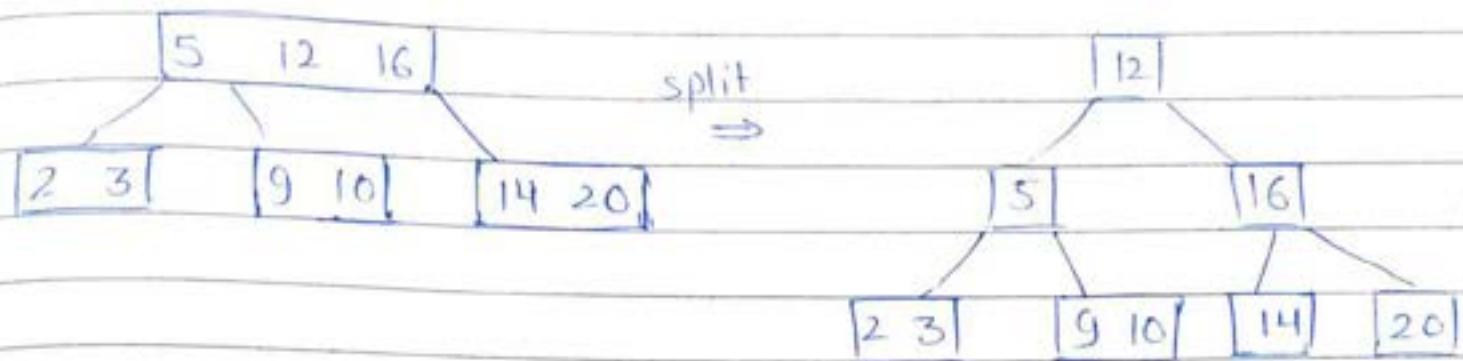


↙ split

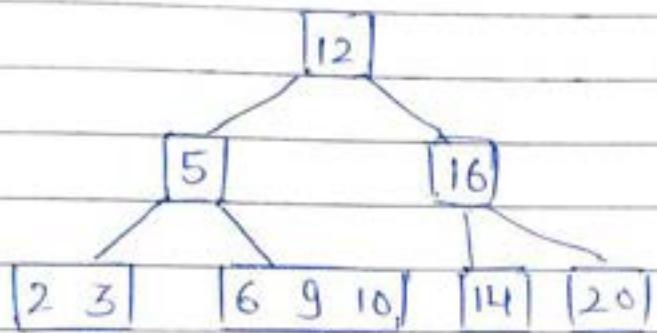
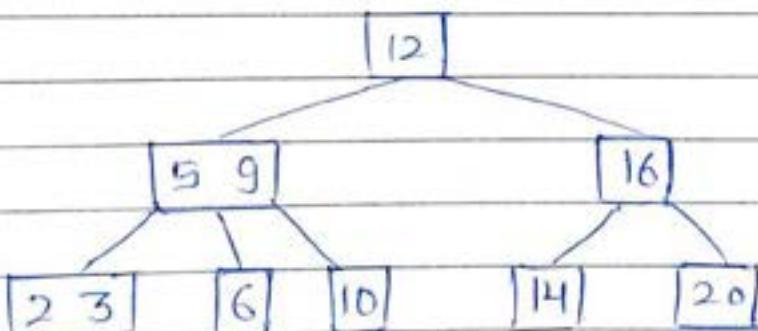


split





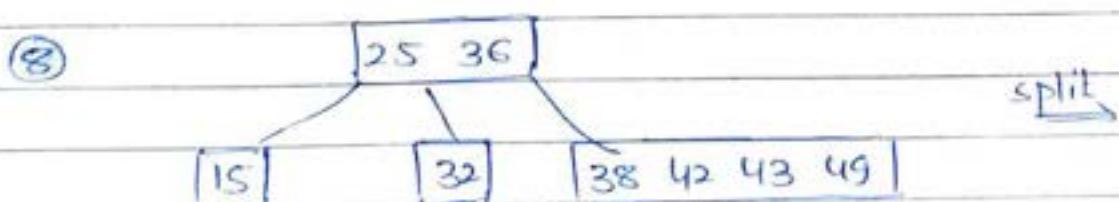
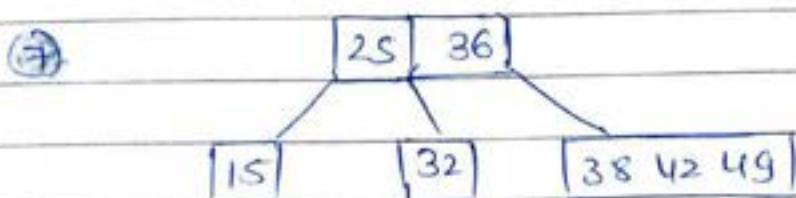
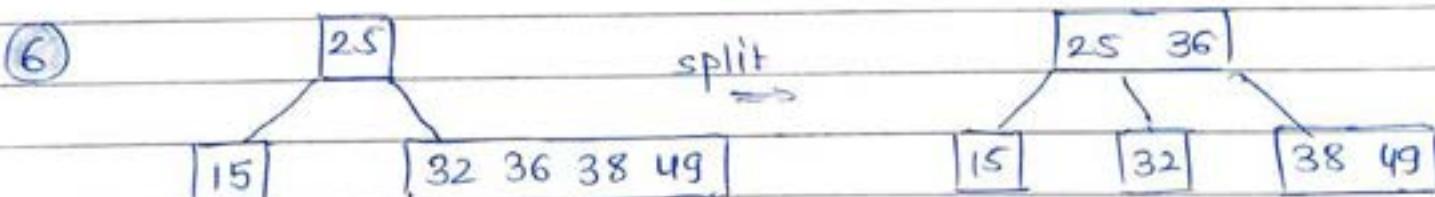
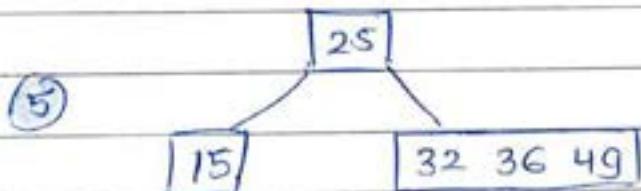
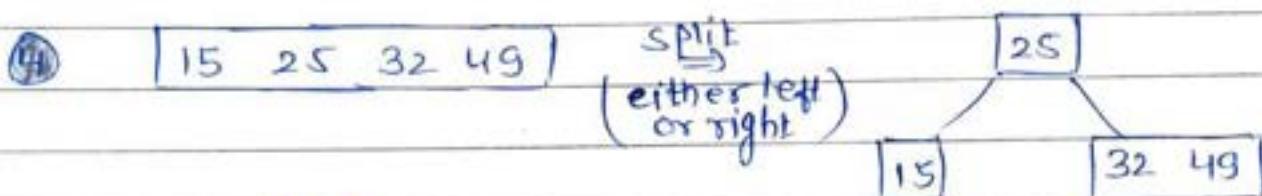
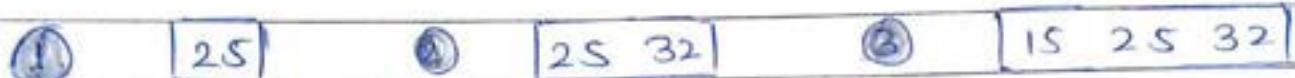
⑥

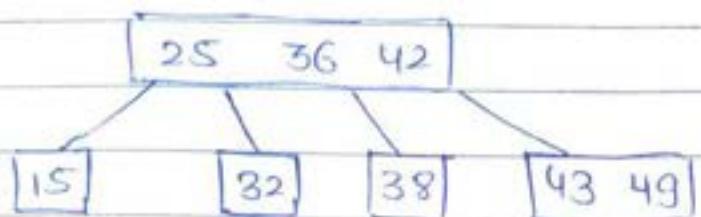
 $\Downarrow$  split $\rightarrow$ 

\* New node is created only when root node is splitted.

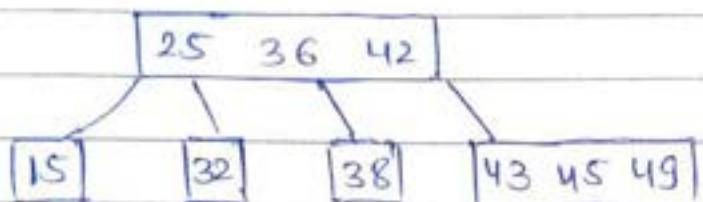
Ques B-tree of order 4

→ 25, 32, 15, 49, 36, 38, 42, 43, 45, 39, 37

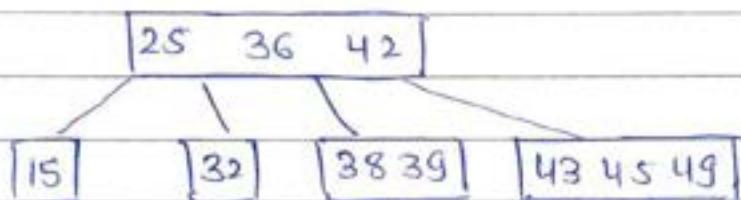




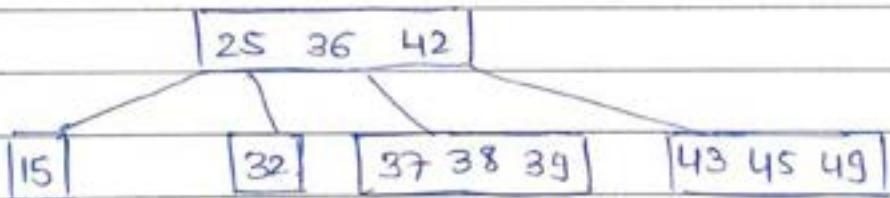
(9)



(10)



(11)



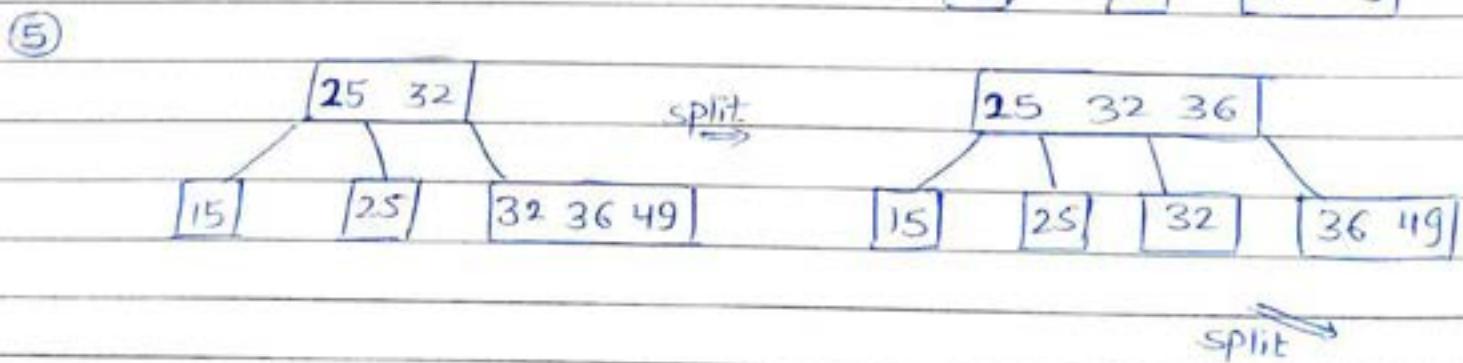
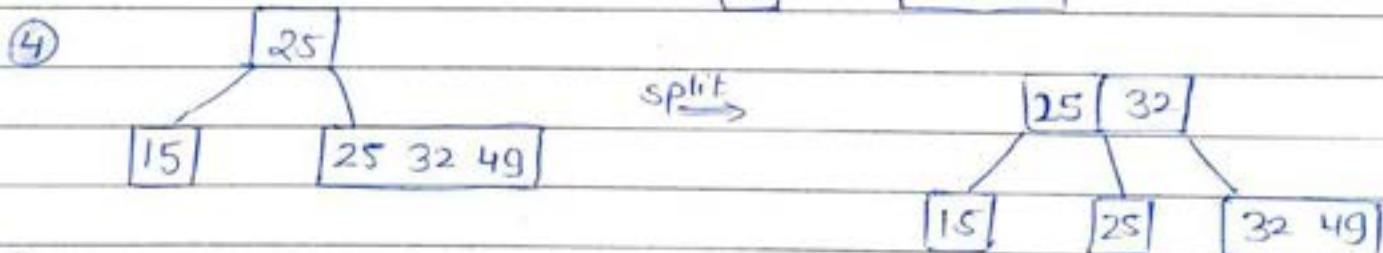
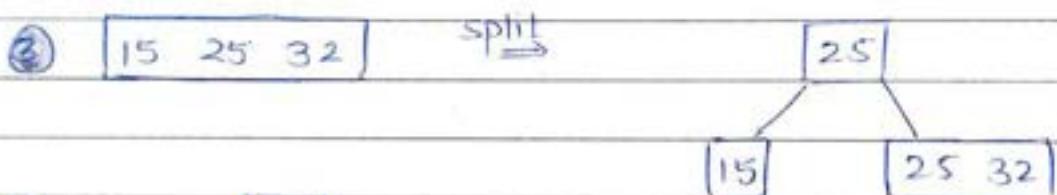
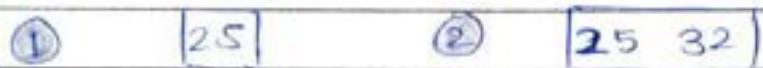
→ —

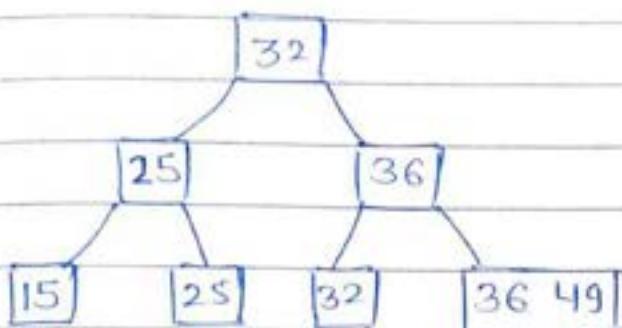
# B<sup>+</sup>-Tree : They are B-trees in which every key must present at leaf node.

\* Perform replication only when leaf is splitted, not for non-leaf node.

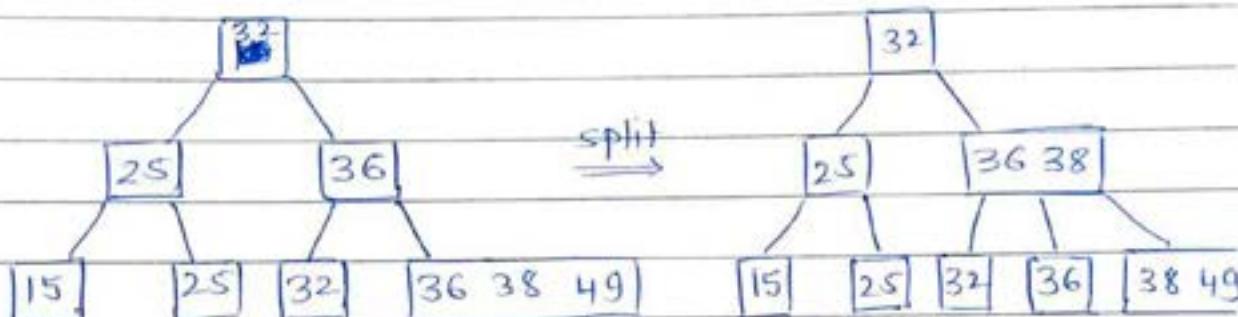
Ques B<sup>+</sup>-tree of order = 3.

→ 25, 32, 15, 49, 36, 38, 42, 43, 45, 39, 37

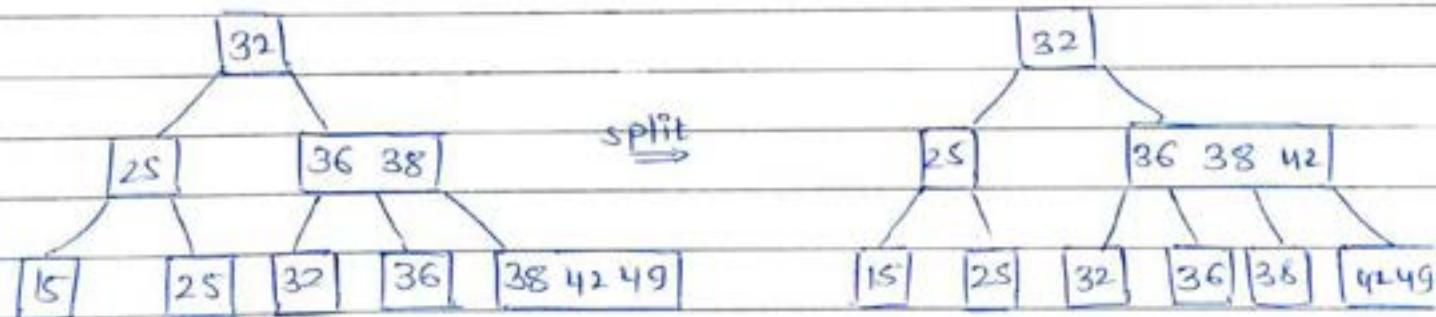




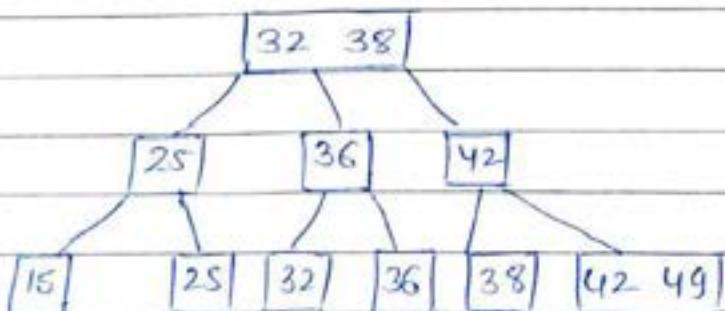
⑤



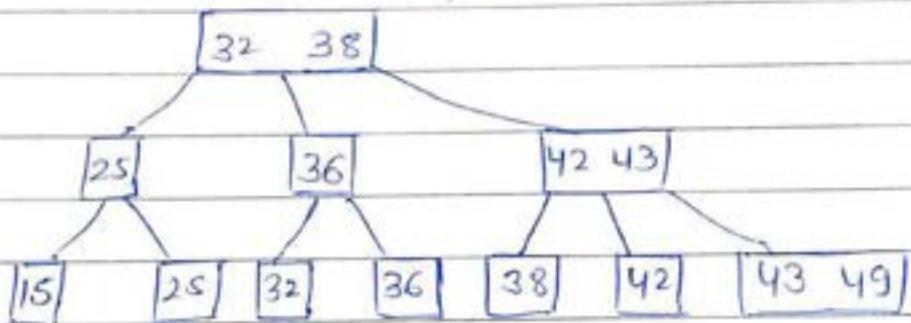
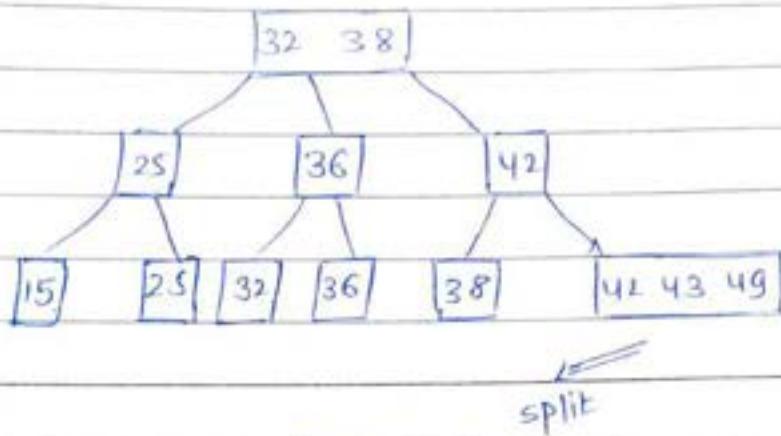
⑥



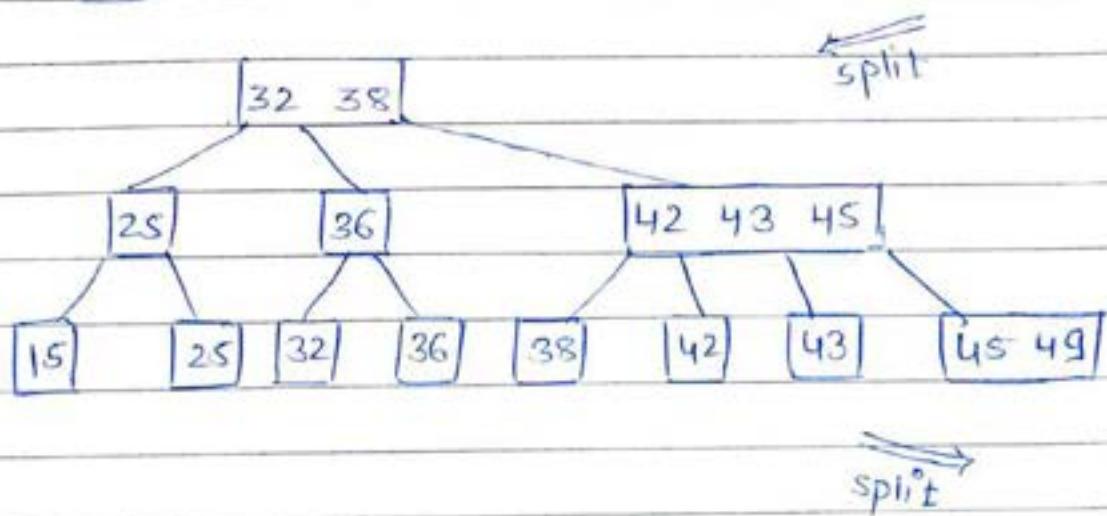
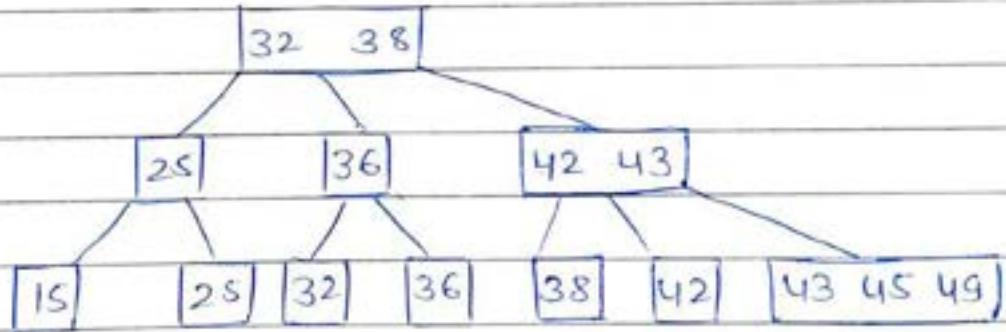
split

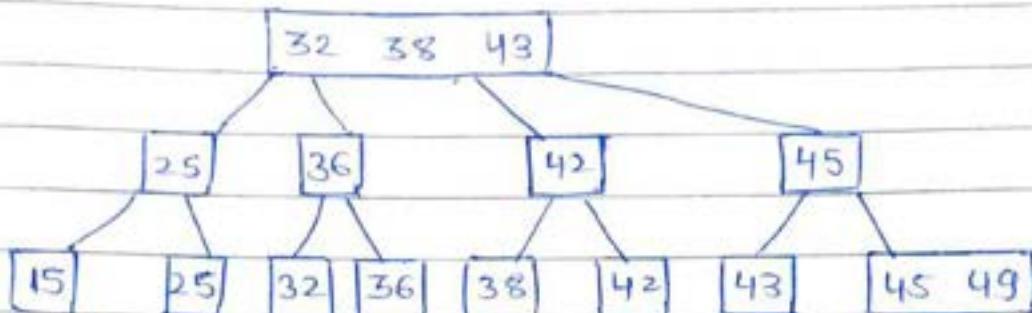


(8)



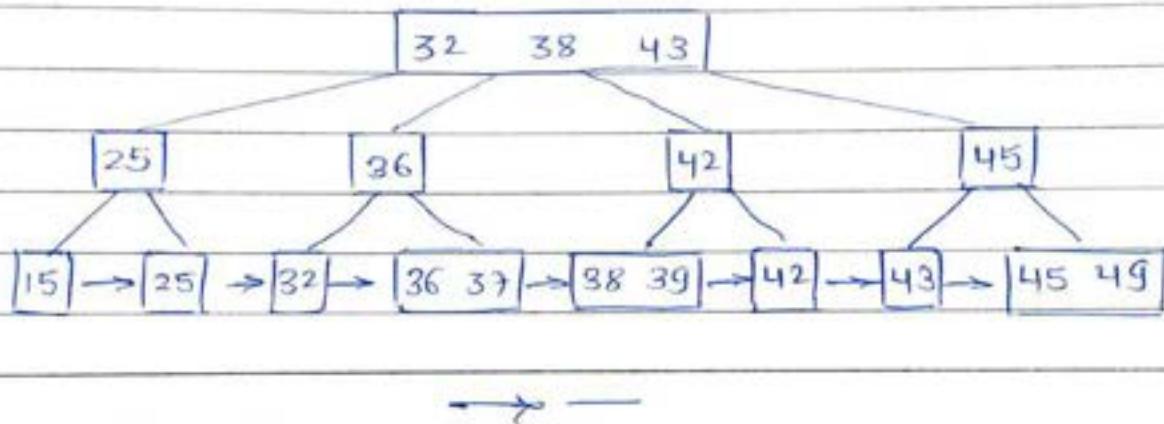
(9)





10

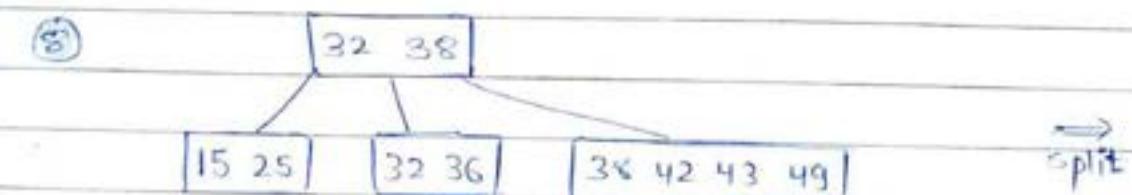
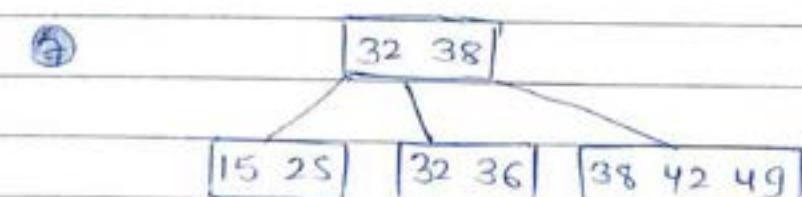
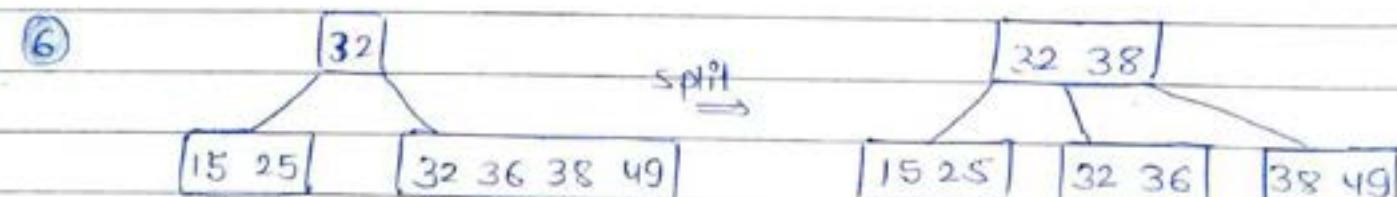
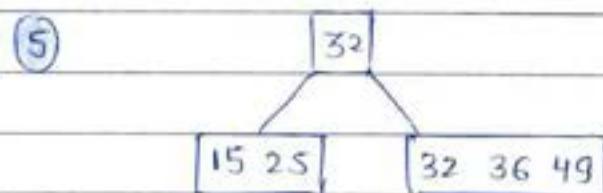
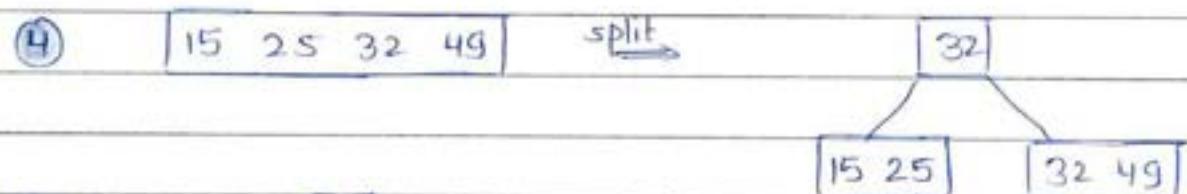
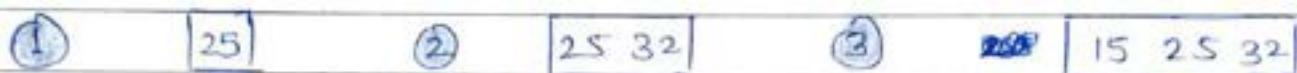
11

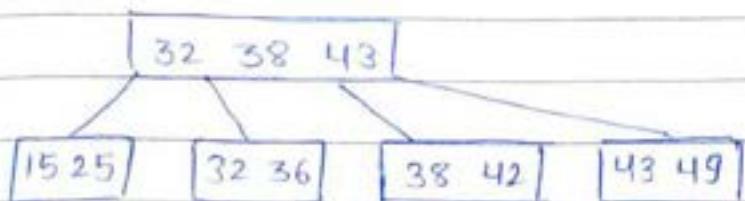


- In B<sup>+</sup> tree, always perform left biasing. [In case of even order]
- Left biasing means left subtree has more elements than right.

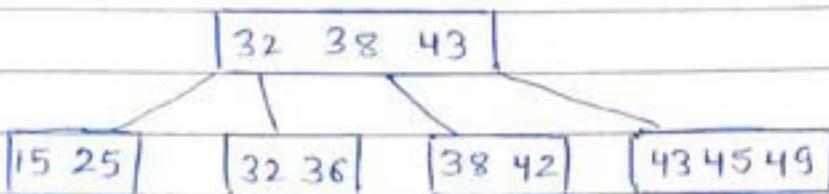
Ques B<sup>+</sup>-tree of order 4

→ 25, 32, 15, 49, 36, 38, 42, 43, 45, 39, 37



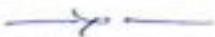
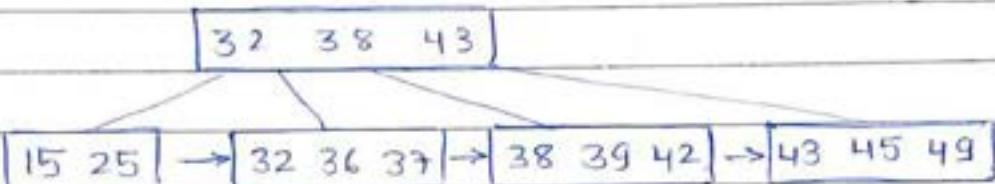


9



10

11



- \* In  $B^+$  tree, all the leaf nodes are connected to each other through linked lists.
- \* In B tree, record pointer is with key itself in every node
- \* In  $B^+$  tree, record pointer is in leaf nodes thus the searching is always upto leaf nodes.
- \* In B tree, non leaf node contains key field, child pointer & record pointer.
- \* In  $B^+$  tree, non leaf node contains key field & child pointer.
- \* order of tree is decided by keeping in mind that single node can be accommodated in single block of harddisk.
- \* In B tree, leaf node also stores child pointer. Null value is stored in their child pointers.
- \* In  $B^+$  tree, leaf node doesn't store child pointer. It stores one sibling pointer.

Ques What is the order of non leaf node of B tree & B<sup>+</sup> tree if key field is 4 bytes, record pointer is 8 bytes, child pointer is 6 bytes & block size is 512 bytes.

⇒ If order = n, n child pointers are required in 1 node

\* One child pointer is of 6 bytes, thus n child pointers will be of  $6n$  bytes

\* For <sup>leaf or</sup> non leaf in B tree: (In B tree, structure of leaf & non leaf are same)

$$\Rightarrow \frac{6n}{n} + 4 + (n-1) * k + (n-1) * r = \text{block size}$$

$$6n + 4n - 4 + 8n - 8 = 512$$

$$18n - 12 = 512$$

$$n = 29.11 \approx \underline{\underline{29}} \quad \text{Ans}$$

— — —

Ans

\* For non leaf in B<sup>+</sup> tree: (In B<sup>+</sup> tree, structures of leaf & nonleaf are different)

$$\Rightarrow n * c + (n-1) * k = \text{block size}$$

$$6n + 4n - 4 = 512$$

$$n = 51.6 \approx \underline{\underline{51}} \quad \text{Ans}$$

Ans

— — —

child pointer is also known as tree pointer

Date \_\_\_\_\_ Page No. \_\_\_\_\_

separate

- \* Sometimes, in  $B^+$  tree, orders of leaf & nonleaf are defined as:
  - \* In  $B^+$  tree, order of non leaf is defined as maximum no. of children a node can have. So,
  - \* if order of non-leaf is  $n$ ; it can have maximum  $n$  children & maximum  $(n-1)$  keys.
  - \* For non-root & non-leaf nodes of order  $n$ ;

minimum  $\lceil \frac{n}{2} \rceil$  children & minimum  $\lceil \frac{n}{2} \rceil - 1$  keys

→ —

- \* In  $B^+$  tree, order of leaf node is defined as, maximum no. of keys a node can have.
- \* if order is  $n$ ; leaf node can have maximum  $n$  keys,  $n$  record pointers & 1 sibling pointer & minimum  $\lceil \frac{n}{2} \rceil$  keys.
- \* For root ; minimum 2 children & 1 key
- \* For root & non-root-non leaf ; maximum  $n$  children.

→ —

(for leaf)

Here, Order = max. no. of keys, a node can have.

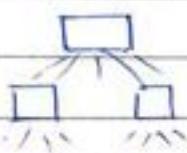
 $\Rightarrow$  previous question for leaf in  $B^+$  tree:

$$\Rightarrow n * k + n * r + c = \text{block size}$$

$$4n + 8n + 6 = 512$$

$$n = 42 \cdot 16 \% \underline{42} \quad \underline{48}$$

Ques Find max. no. of nodes in  $B$  tree or  $B^+$  tree of order = 10 & height = 4  
 (by default height starts from zero).

 $\Rightarrow$  height

nodes

1

$$\Rightarrow 1 + 10 + 10^2 + 10^3 + 10^4$$

10

$$= \underline{\underline{11111}}$$

2

$$10 \times 10$$

3

$$10^3$$

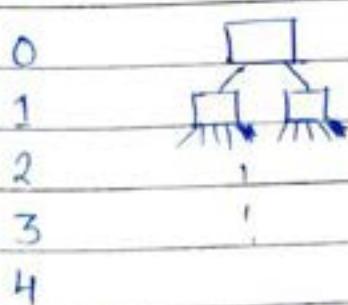
4

$$10^4$$



Ques Find min. no. of nodes in  $B$  tree or  $B^+$  tree of order = 10 & height = 4

- In  $B$  or  $B^+$  tree, <sup>max. or min.</sup> no. of nodes are same but, <sup>max. or min.</sup> no. of ~~children~~ <sup>keys</sup> are different.



1

• min. nodes for root = 1

2

2

• min. key for nonroot =  $\lceil \frac{n}{2} \rceil$ 

3

$$2 \times 5$$

$$= \left( \frac{10}{2} \right)$$

4

$$2 \times 5 \times 5 \times 5$$

$$= \frac{5 \cdot 1}{2} = 4$$

 $\Rightarrow$  For order = n & height = k (starting from zero)

$$\bullet \text{Max. nodes} = 1 + n + n^2 + n^3 + n^4 + \dots + n^k$$

$$\bullet \text{Min. nodes} = 1 + 2 \left[ 1 + \lceil \frac{n}{2} \rceil + \lceil \frac{n}{2} \rceil^2 + \dots + \lceil \frac{n}{2} \rceil^{k-1} \right]$$

Ques Consider a harddisk with block size 1024 bytes, key field 4 bytes, record pointer 10 bytes & child pointer 8 bytes. What should be the order of leaf & non-leaf node of  $B^+$  tree.

$\Rightarrow$  For leaf :

$$\Rightarrow n * k + n * r + 1 * c = 1024$$

$$4n + 10n + 8 = 1024$$

$$n = 92.5 \approx \underline{72} \quad \cancel{73}$$

$\Rightarrow$  For non-leaf :

$$\Rightarrow n * c + (n-1) * k = 1024$$

$$8n + 4n - 4 = 1024$$

$$n = 85.66 \approx \underline{85} \quad \cancel{86}$$

$\rightarrow$

- \* Block pointer is addr. of block

- \* Record pointer is addr. of block & offset within the block.

- \* In multilevel indexing, 1<sup>st</sup> level has record pointer & 2<sup>nd</sup> level has block pointer.

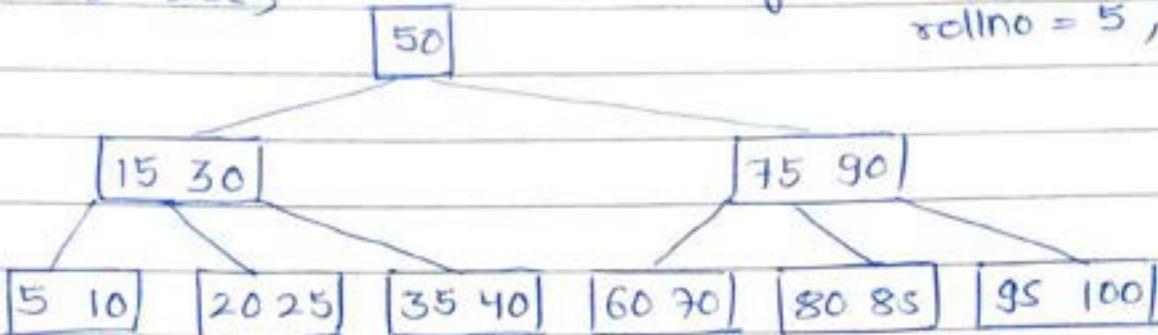
- \* In  $B/B^+$  tree, child pointer is also block pointer as one node of  $B/B^+$  tree equals one block of H.D.

or  $B^+$  tree

$\Rightarrow$  1 node of  $B$  tree = 1 block of harddisk Date \_\_\_\_\_ Page No. \_\_\_\_\_

\* In B tree;

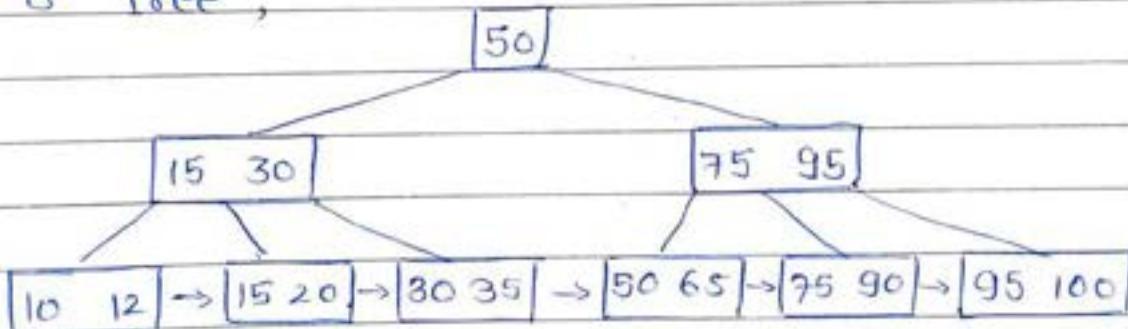
e.g. select name where rollno = 5;



read

- $\Rightarrow$  No. of blocks to search 5  $\Rightarrow \underline{3} + 1$  (database)
- $\Rightarrow$  No. of blocks to search 75  $\Rightarrow \underline{2} + 1$  "
- $\Rightarrow$  No. of blocks to search 50  $\Rightarrow \underline{1} + 1$  "

\* In  $B^+$  tree;



- $\Rightarrow$  No. of blocks to search any key  $\Rightarrow \underline{3} + 1$  (database)  
as record pointer is stored in leaf nodes.

\* But, yet  $B^+$  tree is better because depth of  $B$  tree is more than that of  $B^+$  tree. Thus,  $B^+$  tree is used.

 $\rightarrow$ 

- \* In  $B^+$  tree, non-leaf doesn't store record pointer, they just help to traverse to leaf nodes.

- ⇒ to print values between 10 & 100 ⇒ range query.
- \* In B tree, we have to search all the values between 10 & 100 whether they are present or not.
- \* But in  $B^+$  tree, after searching first element, all the values can be accessed. Thus, range queries are faster ~~than~~ in  $B^+$  trees.

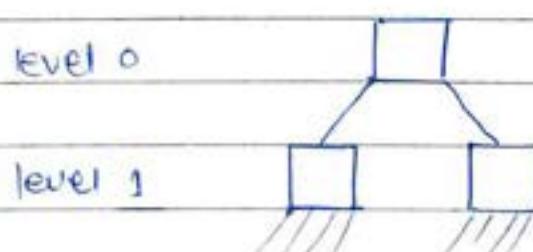
→ —

Ques What is the min. no. of keys a  $B^+$  tree of order 101 can store if level of tree is 2.

$$\Rightarrow \text{min. children} = \left\lceil \frac{101}{2} \right\rceil = 51$$

$$\text{min. keys} = 51 - 1 = 50$$

Root  $\Rightarrow$  min. children = 2, min. key = 1



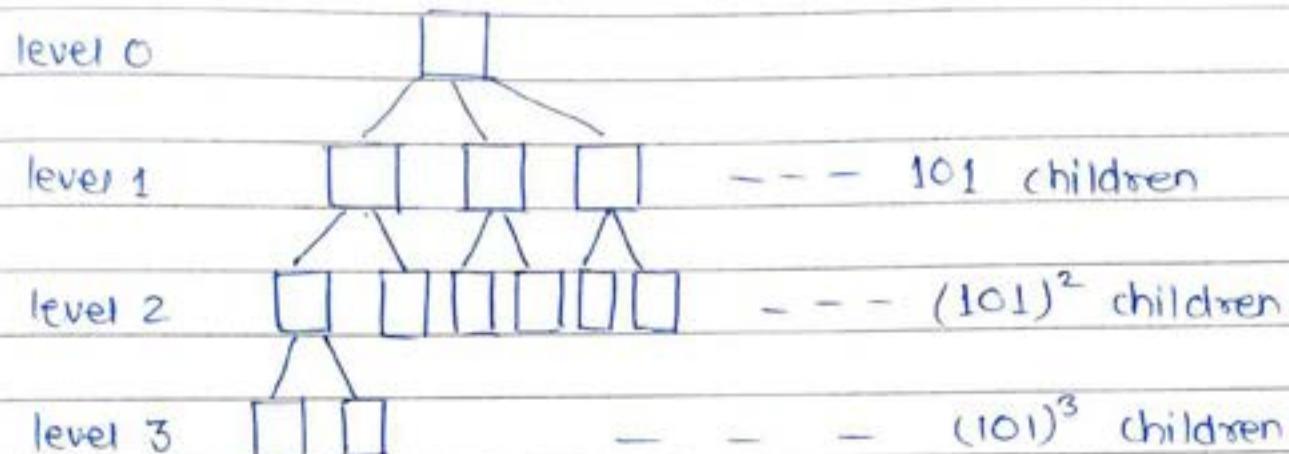
→ here, every children will have 50 keys

$$\text{Thus, } \Rightarrow (51 + 51) * 50 = 5100 \text{ keys (min)}$$

\* It's  $B^+$  tree thus, we've computed for leaf nodes only.

Ques what is the max. no. of keys a  $B^+$  tree of order 101 can store if level of tree is 3.

$\Rightarrow$  max. children = 101, max. keys =  $101 - 1 = 100$

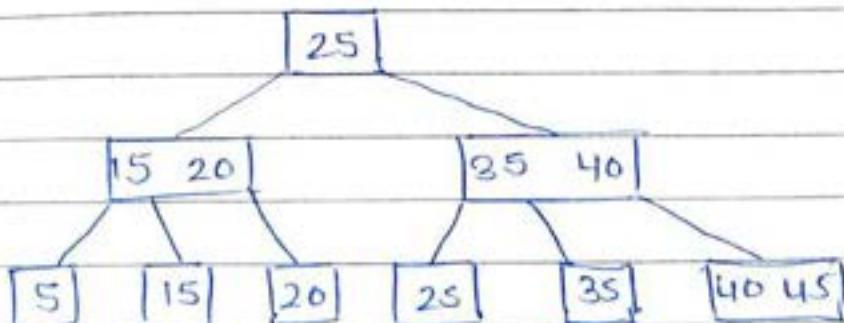


$$\Rightarrow 101 \times 101 \times 101 = (101)^3 \text{ children}$$

$$\Rightarrow (101)^3 \times 100 \text{ keys}$$

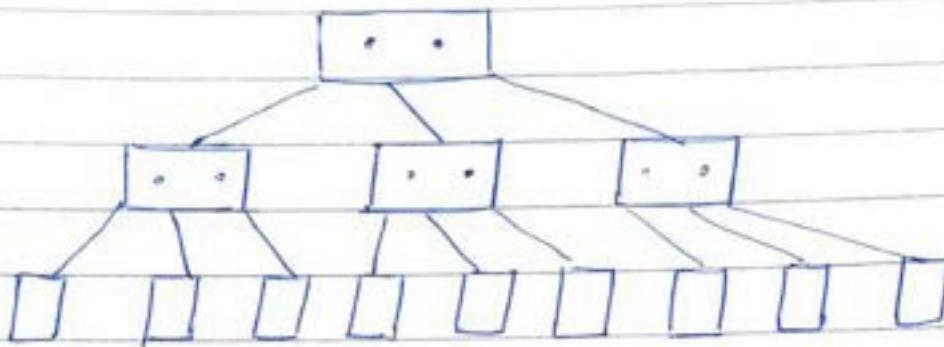
→

Ques what is the max. no. of keys can be inserted in given  $B^+$  tree without increasing level of tree, order = 3



\* min key at root is 1, not  $\left(\left\lceil \frac{D}{2} \right\rceil - 1\right)$ .

Date \_\_\_\_\_ Page No. \_\_\_\_\_



\* total keys at leaf nodes =  $2 \times 9 = 18$

\* keys ~~previously~~<sup>already</sup> inserted = 7

⇒ Thus,  $18 - 7 = 11$  keys can be inserted.

↓ min.

Ques Find max. no. of keys in B tree of order = 10 & height = 4.

⇒ Max keys = max. keys in 1 node \* max. no. of nodes =  $9 \times \text{max nodes}$

$$\cdot \text{Max keys} \Rightarrow (n-1)(1+n+n^2 \dots n^k)$$

$$\cdot \text{Min keys} \Rightarrow 1 + \left[ \left( \left\lceil \frac{n}{2} \right\rceil - 1 \right) 2 \right] \left[ 1 + \left\lceil \frac{n}{2} \right\rceil + \left\lceil \frac{n}{2} \right\rceil^2 \dots \left\lceil \frac{n}{2} \right\rceil^{k-1} \right]$$

→ —

Ques Find max. & min. no. of keys in B<sup>+</sup> tree of order = n & height =  $k$ .

• In B<sup>+</sup> tree, all keys are at leaf level & max nodes at leaf level can be  $(n^k)$  & max. keys at a node is  $(n-1)$

$$\therefore \text{Max. keys} \Rightarrow (n-1)(n^k)$$

$$\text{Min keys} \Rightarrow \left( \left\lceil \frac{n}{2} \right\rceil - 1 \right) \left( 2 * \left\lceil \frac{n}{2} \right\rceil^{k-1} \right)$$

min keys at  
a node.

min nodes at leaf level

# Hashing:

- \* Hashing is a technique in which searching is not required.
- \* In hashing, we choose a hash function which when applied on a key, returns the index.

eg hash function : index = key % 1000

- $4080 \% 1000 = \underline{80}$ , key 4080 will be saved at index 80.
- $3126 \% 1000 = \underline{126}$

- \* The same hash function will be applied in reverse queries. (select name where rollno = 4080)

$$\text{But, } 4080 \% 1000 = \underline{80}$$

$$3080 \% 1000 = \underline{80}$$

- \* This problem is known as hash collision. If a hash function returns same value for 2 different keys.
- \* Thus, a hash function is chosen such that there will be less or no hash collisions.

# Techniques to handle collision :① Rehashing

(Open Addressing)

② Chaining

- ↳ Linear Probing
- ↳ Quadratic Rehashing

\* In rehashing, we change the hash function for the key which collides.

# Linear Probing :

\* In linear probing, we search for free space ~~be~~ linearly.

eg  $3080 \% 1000 = \underline{80} \Rightarrow 3080$  will be stored  
on 80 index.

$4080 \% 1000 = \underline{80} \Rightarrow$  Now, we'll search for  
free space after 80.

eg  $9998 \% 1000 = \underline{998} \Rightarrow 9998$  will be stored  
on 998

$2998 \% 1000 = \underline{998} \Rightarrow$  Now, we'll search for  
free space after 998

& if there's no free space till end (999) then  
we'll search from the beginning.

→ drawback of linear probing :

- ① It may behave as linear search on collision if hash function is not chosen accordingly. (If the hash function returns same index for all the keys.)
- ② Primary clustering :

\* Hash function should distribute data uniformly on array

eg Which hash function should be preferred.

$$H1 : k^2 \bmod 10, \quad H2 : k^3 \bmod 10$$

index	H1	H2	
0	0	0	$0^2 \bmod 10 = 0, \quad 0^3 \bmod 10 = 0$
1	1/9	1	$2^2 \bmod 10 = 4, \quad 2^3 \bmod 10 = 8$
2		8	$3^2 \bmod 10 = 9, \quad 3^3 \bmod 10 = 7$
3		7	etc.
4	2/8	4	* in H1, 4 index(s) are not used while in H2, all are used.
5	5	5	
6	4/6	6	
7		3	
8		2	Thus, H2 will be preferred. as it distributes uniformly.
9	3	9	

\* To remove primary clustering, quadratic rehashing is used.

CS-2010  
Aye 52Hash func.  $h(k) = k \% 10$ 

0	
1	
2	42
3	23
4	34
5	52
6	46
7	33
8	
9	

\* 6 values are stored

 $\Rightarrow$  If linear probing is used,  
which order will be used?

- (a) 46 42 34 52 23 33  
 (b) 34 42 23 52 33 46  
 (c) 46 34 42 23 52 33  
 (d)

Also, how many orders of insertion exist using  
same key values so that you will get same  
resulting hash table.

⇒ Given a hash table, what is load factor:

0		(occupied)
1		
2	42	Load = <u>total keys in slots</u>
3		Factor                    total slots
4	34	= $\frac{3}{10} = 0.3$
5		
6	86	
7		
8		
9		

Ques Given a hash table of size 20, after how many insertions, probability of collisions will be greater than equal to  $1/2$ .

⇒ On 1<sup>st</sup> insertion; prob = 0

⇒ On 2<sup>nd</sup> insertion; prob =  $\frac{1}{20}$

⇒ On 3<sup>rd</sup> insertion; prob =  $\frac{12}{20}$

Thus, on 11<sup>th</sup> insertion; prob =  $\frac{10}{20} = \frac{1}{2}$

⇒ Thus, after 10 insertions, prob of collisions will be greater than or equal to  $1/2$ .

# Quadratic Rehashing:

\* linear probing  $\rightarrow +1, +2, +3, \dots$

\* quadratic rehashing  $\rightarrow +1^2, +2^2, +3^2, \dots$   
 $\Rightarrow +1, +4, +9, \dots$

e.g. 80

$\Rightarrow 80, 81, 84, 89, 96, 105, 116, \dots$

\* Now, there will be no clustering but its drawback is that it stores key values only in selected index despite the array is empty.

## • Quadratic Rehashing.

$\Rightarrow \text{index} = \text{key} \% 1000$

{ if ( $a[\text{index}]$  is not free)

for  $i=1$  to  $(n-1)$

{ if ( $a[\text{index} + i^2] \cdot n$  is free)

// Replace  $i^2$

by  $i$  for  
linear probing

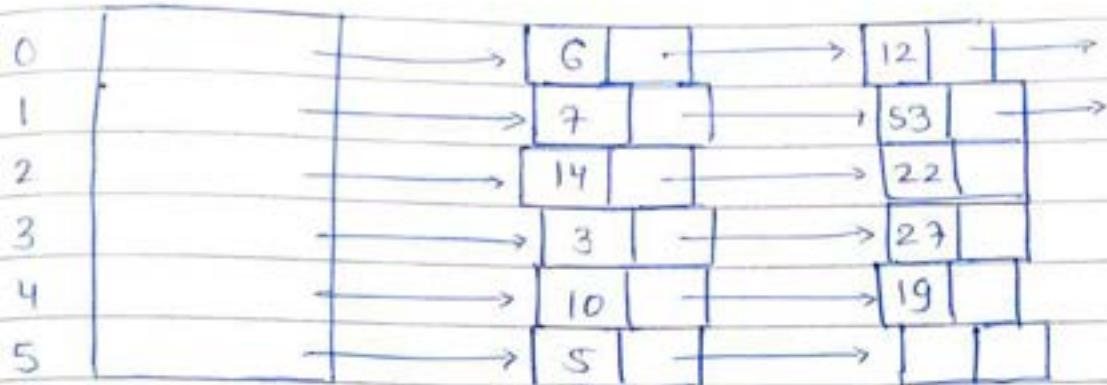
allocate the index;

break;

{

}

## # Chaining :



array of  
pointers  
(stores addresses)

$$\text{Here, load factor} = \frac{11}{6} = 1.83$$

\* if all elements follow a single linked list then the complexity of searching will be  $O(n)$  (linear search). Thus, a hash function is used that distributes uniformly.

⇒ How to perform hashing in harddisk :

① Static Hashing

② Dynamic / Extensible Hashing

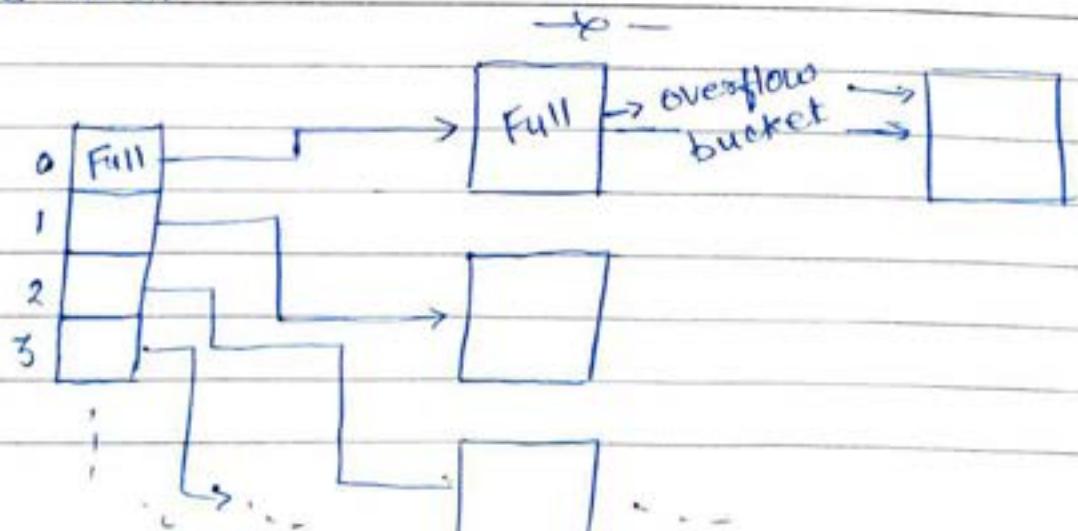
## # Static Hashing :

- \* If bucket (block) gets full, the database will make a new bucket known as overflow bucket.   
 *→ skew*
- \* When a table is created, database cannot know how many rows will be in the table. Thus, it cannot choose hash function accordingly. Thus, hash func. is chosen by some estimation / assumption.

eg key  $y \cdot 4 \Rightarrow$  4 blocks will be reserved in harddisk  
 if more rows are inserted in table,  
 there will be many overflow buckets &  
 linear search has to be applied

key  $y \cdot 1000 \Rightarrow$  1000 blocks will be reserved in harddisk  
 if less rows are inserted in table,  
 it'll lead to wastage of storage.

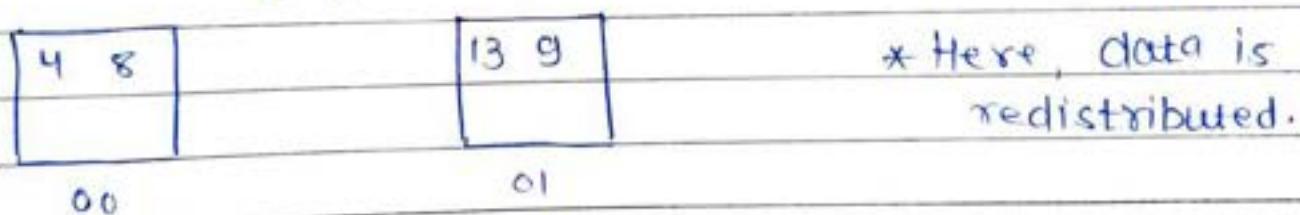
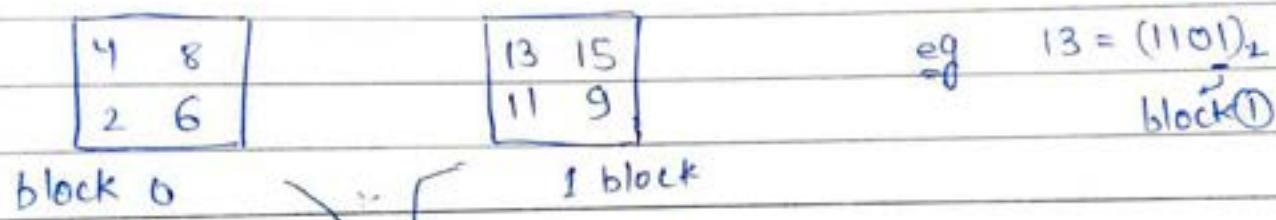
- \* Static Hashing fails if the table grows dynamically. Thus, to remove this problem, dynamic hashing is used.



## # Dynamic Hashing :

\* If hash function is like ; key  $\propto 10^6$ ,  
it won't reserve  $10^6$  blocks in HD.

⇒ First, only 2 blocks will be reserved & when  
they will get full, 4 blocks will be reserved then  
8, 16, 32 - - -



\* If dynamic hashing is implemented successfully,  
searching will be zero.

- \* While using hashing, indexing is also important.
- \* Hashing can not be done in more than one column.
- \* Along with hashing on rolling, indexing will also be done on rolling.
- \* Range queries are better executed in indexing & equal queries are better executed in hashing.

#### o Dynamic Hashing (Navathe):

e.g. Hash func.  $\Rightarrow$  key  $\times .16$

32	16
7	5
17	S

blocks

15	
8	

1

- block 1
- ①  $15 \times .16 = (1111)_2$
  - ②  $32 \times .16 = (0000)_2$
  - ③  $7 \times .16 = (0111)_2$
  - ④  $17 \times .16 = (0001)_2$
  - ⑤  $16 \times .16 = (0000)_2$
  - ⑥  $8 \times .16 = (1000)_2$
  - ⑦  $5 \times .16 = (0101)_2$

\* Initially, only 2 blocks are allocated.

assume, one block stores 5 records.

$\Rightarrow$  Now, 0<sup>th</sup> block gets full, add 2 blocks to it & redistribute

32	16
7	5
17	S

0


local depth = 1

00

Keys  
2 so at  
upto  
total 16  
blocks


local depth = 2

Here, global depth = 4 (max bits in remainder)  
& local depth =

# TRANSACTION

(167)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

- \* Transaction is a set of statements which either committed (executed) fully or shouldn't be committed at all.
  - \* If it's executed partially, there will be some inconsistency.
- ⇒ There are 4 basic properties of transaction :

- ① A : Atomicity
- ② C : Consistency
- ③ I : Isolation
- ④ D : Durability

- ① Atomicity : Transaction should be atomic. i.e either committed fully or nothing.
- ② Consistency : Transaction should lead one valid state to another valid state. eg Money transfer
- ③ Isolation : When a transaction is running, other transactions shouldn't run, especially when both are manipulating same data item.  
eg last seat reservation.

④ Durability : Once the end user is informed about successful completion of a transaction, then later on, it shouldn't inform about its failure.

- \* If a transaction ensures atomicity & isolation then it will definitely ensure consistency, if there's no errors.

- \* If the transactions are executed one by one, then that execution is serial execution & all four properties are ensured in serial execution.

- \* But, the transactions cannot be executed serially, concurrent or parallel execution is required but the effect should be same as in serial execution i.e. all four properties must be ensured.

⇒ A concurrent schedule is serializable or valid if it is equivalent to one of the serial schedules.

→

PTO

↗

eg

T1

T2

S11: update bank

S21: update bank

S12: set bal = 10000

S22: set bal = 20000

S13 where acno = 10;

S23: where acno = 10;

\* if ~~either~~ serial schedule  $\rightarrow$  T1 T2 ✓  
 then, bal = 20000

\* if serial schedule  $\rightarrow$  T2 T1 ✓  
 then, bal = 10000

\* Thus, concurrent schedule should be equi. to (T1 T2) or  
 $(T2 T1)$

$\Rightarrow$  If some concurrent schedule is given, then to  
 check whether it is serializable or not, we've  
 2 tests of serializability;

① Conflict Serializability Test

② View Serializability Test



\* 2 schedules can't be termed as equivalent if  
 their outputs are same (they're just output-equivalent)

\* 2 schedules are equivalent if in both schedules, the  
 values read by T1 & T2, should be same.

## # Conflict Serializability Test:

Consider 3 set of operations to be conflicted;

- ① Read  $x \rightarrow$  write  $x$
- ② Write  $x \rightarrow$  read  $x$
- ③ Write  $x \rightarrow$  write  $x$

$\Rightarrow$  2 operations are conflicted if:  
 ① they belong to diff trans & they manipulate same item  
 ② one of them must be write operation.

$\Rightarrow$  For this test, we draw precedence graph in which transactions are vertices & there will be an edge from  $T_1$  to  $T_2$  if  $T_1$  performs conflicting operation first & then  $T_2$  performs.

Ques Create precedence graph for given schedule;

S1 :  $r_1(x), r_2(x), w_3(x), w_1(y), r_2(y), w_1(x)$

•



\* Just check for next conflicting

$\Rightarrow$  if precedence graph contains cycle, then the schedule doesn't satisfy conflict serializability test.

\* no self loops

Ques S1 :  $\tau_1(x), w_2(y), \tau_2(y), w_1(x), w_2(x), w_2(y)$



$\Rightarrow$  no cycle,  
thus, the schedule  
is conflict serializable.

CS-2009

Q4e S1 :  $\tau_1(x), \tau_2(x), \tau_2(y), w_1(x), w_1(y), w_2(y)$



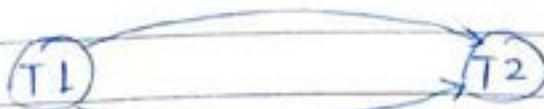
$\Rightarrow$  cycle, thus,  
not conflict  
serializable.

S2 :  $\tau_1(x), \tau_2(x), \tau_2(y), w_1(x), w_2(y), w_1(y)$



$\Rightarrow$  no cycle, thus,  
conflict serializable

S3 :  $\tau_1(x), w_1(x), \tau_2(x), w_2(y), \tau_2(y), w_2(y)$



$\Rightarrow$  no cycle, thus,  
conflict serializable

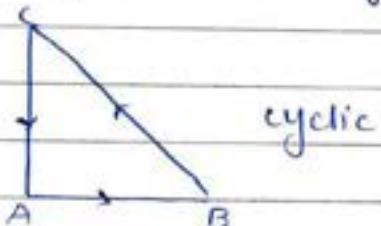
⇒ How to find order of serializability :

(How to find that our concurrent schedule is equivalent to which of the serial schedules)

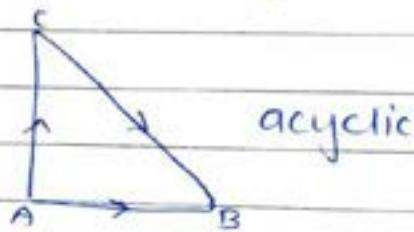
$(T_1, T_2)$  or  $(T_2, T_1)$

# Topological Sort :

- \* It's a graph traversal technique that can only be applied on acyclic graphs (without cycles).



cyclic



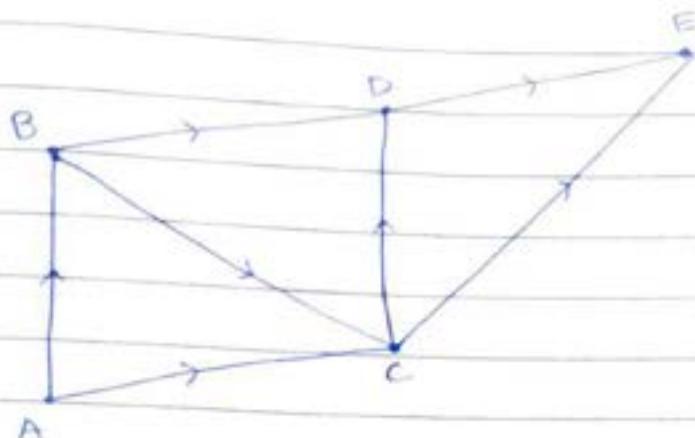
acyclic

- \* If a graph is acyclic, then there will be atleast one vertex with 0 indegree & atleast one with 0 outdegree.

- \* We can start from any 0 indegree vertex.

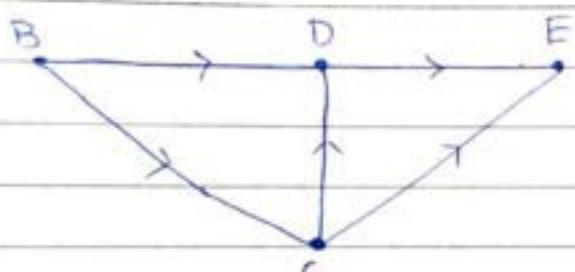


eg

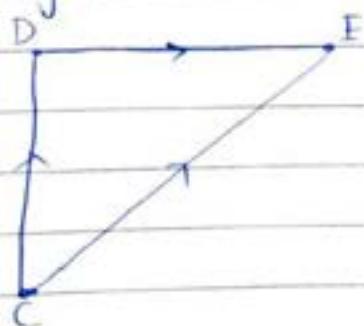


0 indegree vertex  $\Rightarrow$  A

$\Rightarrow$  traverse A & delete edges associated with A.



Now, 0 indegree vertex = B



Now, 0 indegree vertex = C



Thus, order of serializability  $\Rightarrow$  ABCDE

-8 -

\* 2 schedules are conflict-equivalent if their precedence graph is exactly same ie no. of edges should also be eq.

Date \_\_\_\_\_ Page No. \_\_\_\_\_

\* To find the order of serializability, perform topological sort of precedence graph.



→ 0 indegree  $\Rightarrow$  T2

Thus, order  $\Rightarrow$  T2 T1

eg  $\xrightarrow{\text{check}}$  check the serializability of given schedule;

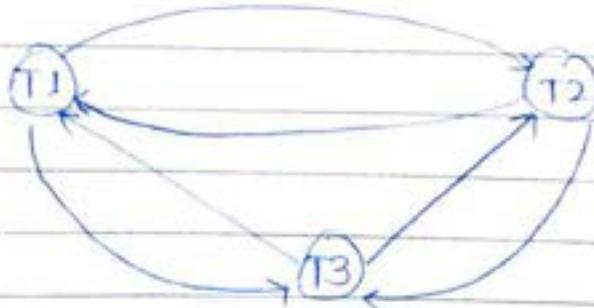
S:  $\tau_1(x) w_2(x) w_3(y) \tau_2(y) w_1(x) w_3(x) \tau_1(y) \tau_2(y)$

→ \* Consider only 2 transactions at a time

T1 & T2  $\Rightarrow$   $\tau_1(x) w_2(x) \tau_2(y) w_1(x) \tau_1(y) \tau_2(y)$

T1 & T3  $\Rightarrow$   $\tau_1(x) w_3(y) w_1(x) w_2(x) \tau_1(y)$

T2 & T3  $\Rightarrow$   $w_2(x) w_3(y) \tau_2(y) w_3(x) \tau_2(y)$

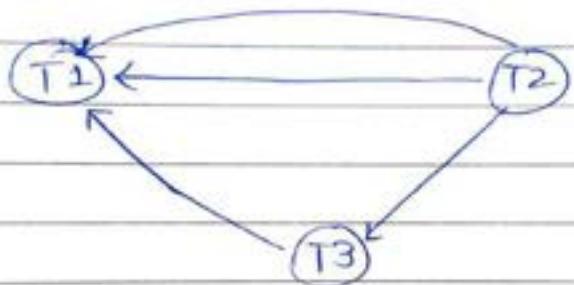


$\Rightarrow$  cycle, thus  
not serializable

\* consider  $\Rightarrow$

27.2.2007  
Ques

S: 2RA, 2WA, 3RC, 2WB, 3WA, 3WC, 1RA, 1RB, 1WA, 1WB

\* 1 & 2  $\Rightarrow$  2RA, 2WA, 2WB, 1RA, 1RB, 1WA, 1WB\* 2 & 3  $\Rightarrow$  2RA, 2WA, 3RC, 2WB, 3WA, 3WC,\* 1 & 3  $\Rightarrow$  3RC, 3WA, 3WC, 1RA, 1RB, 1WA, 1WB $\Rightarrow$  no cycles,  
thus, serializable.\* Order  $\Rightarrow$  T2, T3, T1

\* Conflict Serializability is stricter definition of serializability i.e. there are some schedules which are serializable but doesn't satisfy conflict serializability test



\* If there's no edge in precedence graph then the schedule is serializable & order can be T1T2 or T2T1

T1

T2



## # View Serializability + Test :

T1	T2	T3
$\tau(A)$		
	$w(A)$	
$w(A)$		
		$w(A)$

$\Rightarrow \tau_1(A) \ w_2(A) \ w_1(A) \ w_3(A)$

• It's not conflict serializable as it can't be equivalent to any serial schedule.

\* Here, write of T2 & T3 is not read by any transaction, thus their order doesn't matter.

\* T2 & T3 are doing blind writes here, as they are writing A without reading.

$\Rightarrow$  If there is blind write & the schedule is not conflict serializable, then only you've to check for view serializability test.

- \* Every conflict serializable is view serializable but every view serializable schedule is not necessarily a conflict serializable.
- \* If there is no blind write & schedule is not conflict serializable then it ~~is~~ won't definitely be view serializable.

Serializable Schedules	
Conflict Serializable Schedule	View Serializable schedules

Q Given concurrent schedule is equi to which serial schedule?

- |   |  |   |
|---|--|---|
| T1<br>Read A<br>$A = A + 50$<br>Write A | T2<br>Read A<br>$A = A + 30$<br>Write A<br>Read B<br>$B = B - 30$<br>Write B | ① <u>T1 T2</u> $\Rightarrow$ T1 reads <del>A's initial value</del><br>T2 reads <del>B's updated value</del><br>② <u>T2 T1</u> $\Rightarrow$ T2 reads<br>A's initial value<br>but there T1<br>reads updated value<br>③ Both<br><del>④ None</del> |
|---|--|---|
- Read B  
 $B = B - 50$   
 Write B

This schedule  $[r_1(A), w_2(A), w_1(A), w_3(A)]$  is not conflict serializable but it's view serializable as view checks for order of reads only.

Date \_\_\_\_\_ Page No. \_\_\_\_\_

⇒ How to check View Serializability :

T1	T2	T3	T <sub>f</sub>	
$r_1(A)$				• Add one trans. T <sub>f</sub>
	$w_2(A)$			
		$w_1(A)$		
			$r_3(A)$	

\* T<sub>f</sub> will read all the data items which have been written. (only A, here)

\* View Serializability checks ~~order~~ of read operations only. to check equivalence of schedules.

⇒ T1 reads the initial value of A.

⇒ T<sub>f</sub> reads the value of A written by T3.

Serial schedules:

① T2 T1 T3 T<sub>f</sub> X Here, T3 reads updated value of A, not initial. X

② T1 T3 T2 T<sub>f</sub> X Here, T<sub>f</sub> reads value by T2, not T3

③ T1 T2 T3 T<sub>f</sub> ✓ given this, concurrent schedule is equivalent to T1 T2 T3 T<sub>f</sub> serial schedule.

\* A schedule is view serializable if it's view equi. to any serial schedule.

(179)

Date \_\_\_\_\_ Page No. \_\_\_\_\_

\* View Serializability cannot be checked by finite time algorithm. It's an NP-complete problem.

\* Practically, server cannot check even for conflict serializability. Thus, there are some protocols by which servers ensure serializability;

① Two Phase Locking Protocol

② Time Stamp based Protocol



• How many schedules are possible with  $n$  transactions having  $n_1, n_2, \dots, n_n$  statements respectively.

$$\Rightarrow (n_1 + n_2 + n_3 + \dots + n_n)!$$

$$n_1! \cdot n_2! \cdot \dots \cdot n_n!$$

↓  
total schedules

→ because internal order  
of a transaction can't  
be violated.

→ total schedules -  $\frac{n!}{\prod \text{serial}}$  → concurrent schedules



\* For view equi, check for ① initial reads of all data items used  
② final writes  
③ intermediate reads only

\* If all 3 are same as any serial schedule, then schedule is view serializable.

\* Some trans should possess

## # Concurrency Control Protocols :

### ① Two Phase Locking Protocol (2PL) :

\* Locking : If a transaction is manipulating a data item then another transaction should not interfere unless first transaction completes its execution.

Thus, the first transaction will lock that data item

\* There exists two types of locks ;

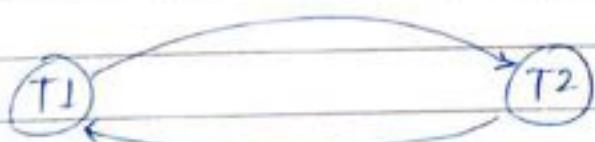
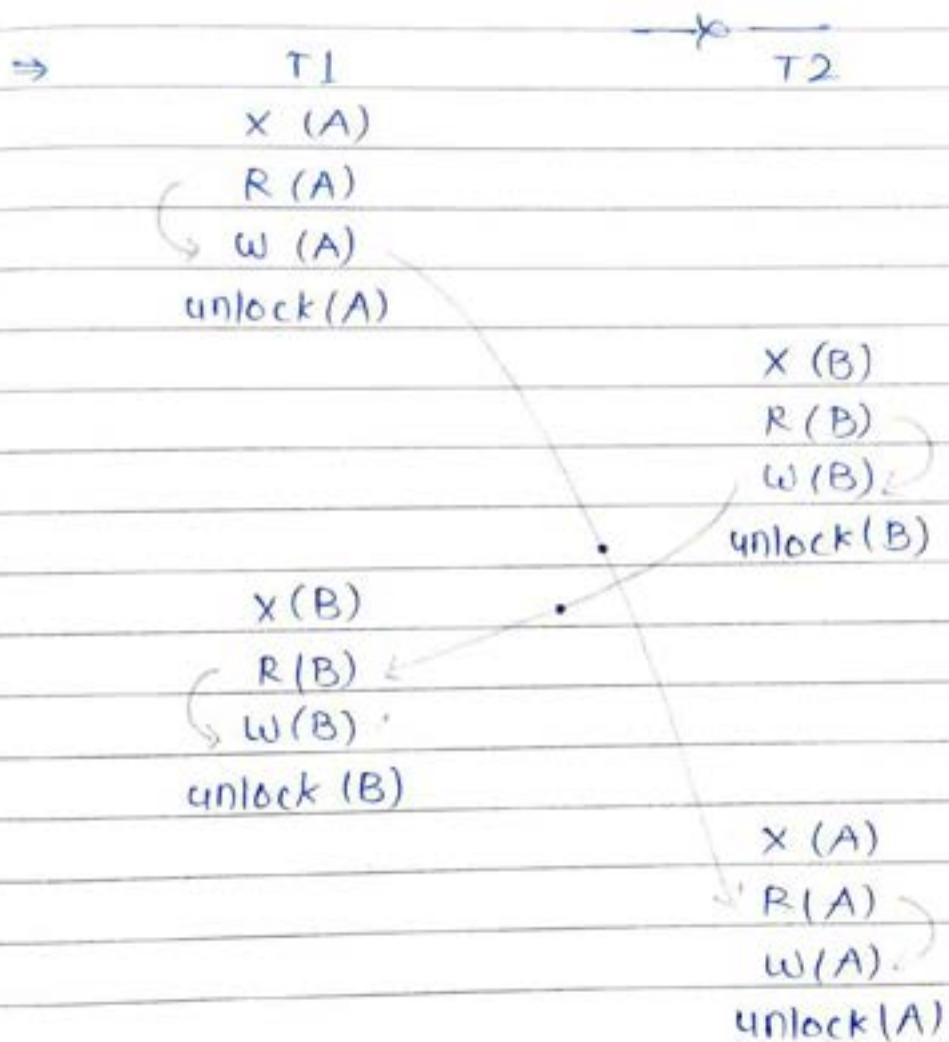
① Shared Lock : It's used when the data items have only to be read.

② Exclusive Lock : It's used when the data items have to be written or manipulated.

\* Many shared locks on a data item can be applied at a time but only one exclusive lock can be applied.

	Unlocked	S	X
S	✓	✓	✗
X	✓	✗	✗

\* Exclusive lock can be applied only when that data item is unlocked.



$\Rightarrow$  cycle, thus not serializable.



• Here, we've performed locking-unlocking, but yet serializability is not ensured.

\* Thus, random locking & unlocking cannot return a serializable schedule. Thus, locking & unlocking should be done according to some rules known as two phase locking protocol (2PL).

⇒ 2PL protocol divides the transaction into two phases;

## ① Growing Phase      ② Shrinking Phase

\* A transaction is in growing phase initially.

- \* A transaction can lock a data item when it's in growing phase ; not in shrinking phase.

- \* A transaction enters in shrinking phase when it unlocks any data item.

- \* Thus, before unlocking a data item, we've to lock data items to be used.

⇒ T1                    T2

X(A)

R(A)

W(A)

X(B)

unlock(A)

X(B)

R(B)

W(B)

X(A)

unlock(B)

R(B)

W(B)

unlock(B)

R(A)

W(A)

unlock(A)



Blocked

only one  
exclusive locks  
can be used at  
a time

thus, T2 will be  
blocked & when

T1 gets chance  
back it executes  
completely. Then

T2 will be  
executed. Thus,  
this schedule  
is no more  
concurrent but

serial, thus,  
serializable.

\* Schedules under 2PL always  
ensure conflict serializability.

\* advantage of 2PL :

\* A schedule returned by 2PL always satisfy conflict serializability. test.

\* drawbacks of 2PL :

- (1) Cascade Rollbacking
- (2) Deadlock

	T1	T2
①	read x	③ read y
②	write x	④ write y

⇒ How many schedules are possible?

\* ② cannot come before ① & ④ can't come before ③.

eg S:  $\gamma_1(x) \quad \gamma_2(y) \quad w_1(x) \quad w_2(y)$  ✓

S:  $\gamma_1(x) \quad \gamma_2(y) \quad w_2(y) \quad w_1(x)$  ✓

S:  $w_1(x) \quad \gamma_1(x) \quad w_2(y) \quad \gamma_2(y)$  ✗

\* If a transaction has m statements & another transaction has n statements then there can be;

$$\Rightarrow (m+n)! \text{ schedules}$$

$$m! \times n!$$

⇒ in above example  $\Rightarrow \frac{(2+2)!}{2! \times 2!} = 6$  schedules

\* no. of transactions = 2, thus  $2! = 2$  schedules will be serial & remaining 4 will be concurrent.

serial : ① ② ③ ④ , ③ ④ ① ②

total - serial

$$\text{* concurrent schedules} = (m+n)! - t! \\ m! \times n!$$

where,  $t$  = no. of transactions

\* If there are 3 transactions, there will be  $3! = 6$  serial schedules.

-  $\downarrow$  -

~~Ques~~ How many concurrent schedules are possible?  
which are conflict serializable?

T1	T2
read x	write x
write y	read z

⇒ order can be T1T2 or T2T1.

- T1T2 ⇒ ①  $r_1(x) w_2(x) w_1(y) r_2(z)$
- ②  $r_1(x) w_2(x) r_2(z) w_1(y)$
- ③  $r_1(x)$
- ④  $r_1(x)$
- ⑤  $r_1(x)$
- ⑥  $r_1(x)$

$$3! = 6$$

$$3!$$

Ques How many conflict serializable schedules are possible?

T1	T2
read(x)	write(y)
write(y)	write(x)

# Cascade Rollbacking :

T1	T2
x (A)	
read (A)	
write (A)	
x (B)	
unlock (A)	
	x (A)
	read (A)
	write (A)
	unlock (A)
read (B)	
if ( B < 5000 )	
rollback	

⇒ Now, T1 will rollback due to which T2 will also have to rollback as T2 is reading value of A written by T1. This is called cascade rollbacking.

T1

x(A)

read(A)

write(A)

x(B)

unlock(A)

T2

x(A)

read(A)

write(A)

unlock(A)

commit

read(B)

if ( $B < 5000$ )

rollback

thus, there will  
be inconsis-  
tency.

$\Rightarrow$  Now, as T2 has committed, it cannot rollback.  
thus, this is known as ~~an~~ irrecoverable schedule.

thus, we prefer transaction to not read the changes by the uncommitted transaction.

Date \_\_\_\_\_ Page No. \_\_\_\_\_

- \* Reading the changes done by uncommitted transaction  $\Rightarrow$  required cascade rollbacking

Ques S:  $r_1(x)$   $r_2(x)$   $w_1(y)$   $r_3(x)$  C2 C3 C1

$\Rightarrow$  T1 & T2 are reading initial value of x but T3 is reading value of x updated by T1.

$\Rightarrow$  Required Cascade Rollbacking & Non-recoverable

Ques S:  $w_1(x)$   $w_2(y)$   $r_1(y)$  C2  $r_3(x)$  C1 C3

$\Rightarrow$  Required Cascade Rollbacking & Recoverable

Ques S:  $r_1(x)$   $w_2(y)$   $r_2(y)$  C2  $r_1(y)$  C1  $r_3(y)$  C3

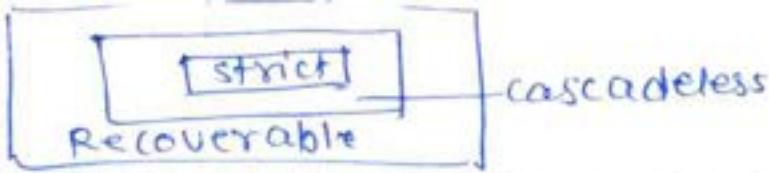
$\Rightarrow$  cascadeless schedule & thus recoverable

Ques S:  $w_1(x)$   $w_2(x)$   $r_3(y)$  C1 C2 C3

$\Rightarrow$  cascadeless schedule & recoverable.

\* For recoverability, check for order of trans. commits

\* the trans. performing dirty read should commit last.



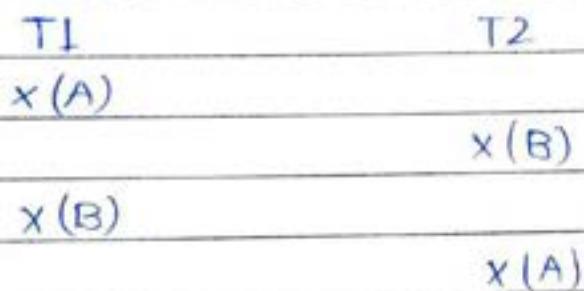
- ⇒ **Strict Schedule:** If value written by trans  $T_1$  can't be read or overwritten by other trans until it (191)  
 (T1) either commits or aborts.

~~Give~~ S:  $w_1(x)$   $w_2(x)$   $r_3(y)$   $c_1(2)$   $c_3$

⇒ Not a strict schedule

- \* Every strict schedule is cascadeless but not viceversa
- To avoid cascade rollbacks in basic 2PL,
- \* Strict 2PL: All the transactions should release their exclusive locks after 'commit'  
 [Thus, dirty reads aren't possible here]
- \* But, it <sup>may</sup> ~~will~~ make our schedule serial.
- \* Rigorous 2PL: All the transactions should release their exclusive & shared locks after 'commit'.
  - In strict schedule, 2PL, shared lock can be upgraded into exclusive, only in growing phase.  
 Also, exclusive lock can be downgraded into shared.  
~~first~~ ~~only~~ & then it will move to 'shrinking' phase.
  - \* In strict & rigorous 2PL, dirty reads are not possible.

## # Deadlock:



- \* T1 is waiting for B & T2 for A. This results in deadlock & there's no solution for this, it's upto programmer.



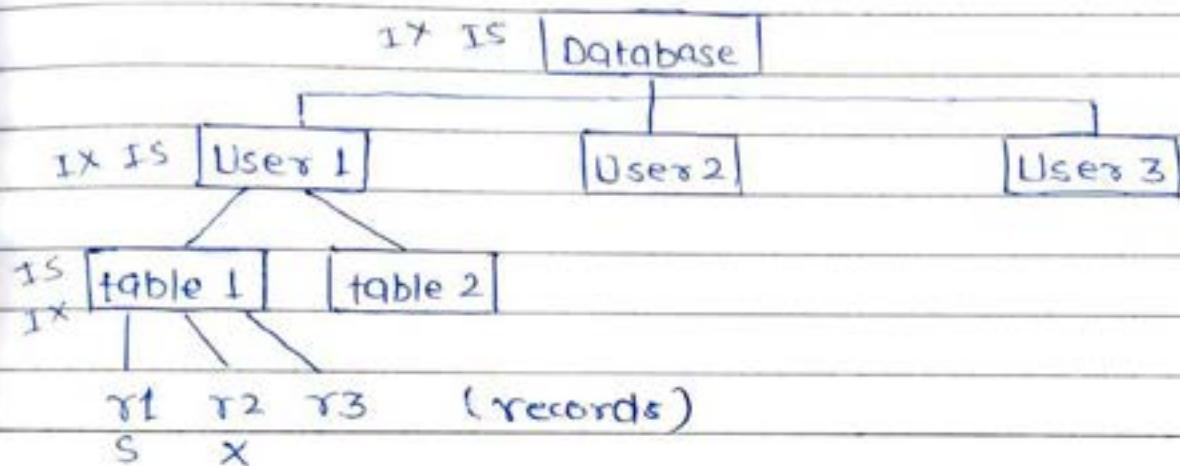
• **Conservative QPL** : A transaction should lock all the data items to be used, in the very beginning.

	Cascade Roll.	Deadlock
① Basic QPL :	✓	✓
② Conserv. QPL :	✓	✗
③ Strict QPL :	✗	✓
④ Rigorous QPL :	✗	✓



## # Granularity of Locking:

- \* Database is collection of users
- \* User is collection of tables.
- \* Table is collection of records.



- \* If this database is to be copied somewhere, but if any record is locked, it cannot be copied.
- \* For a table to be locked, all the records should be unlocked. Thus, how does a table know whether its ~~own~~ records are unlocked or locked.
- \* If a record is shared locked, then its parents, grandparents & so on will be intentionally shared locked.

already locked →

↓ to be locked	IS	IX	S	SIX	X
IS	✓	✓	✓	✓	✗
IX	✓	✓	✗	✗	✗
S	✓	✗	✓	✗	✗
SIX	✓	✗	✗	✗	✗
X	✗	✗	✗	✗	✗

\* locking is pessimistic approach of concurrency control.

## ② Time Stamp Based Protocol :

- \* It ensures serializability without locking.
- \* Every transaction has a time stamp which tells about when it was started.
- \* No two transactions can have same time stamp.
- \* Every data item has 2 time stamp → read time stamp  
write time stamp

eg A → read 11 → it was last read by  
→ write 12 transaction of time stamp 11

↓  
it was last written by  
transaction of time stamp 12.

→



$T_1 \Rightarrow \underline{11}$   
 read A  
 write A

\*  $T_1$  starts on 11, thus,  $T_1$  can read changes done by transactions of time stamp less than or equal to 11.

→ Before reading A,  $T_1$  will first check write time stamp of A, if it is less than 11, then only it will read A else rollback itself & starts with new time stamp.

Now, A → read 9 | 11  
                 → write 9

→ Before writing A,  $T_1$  will check read & write time stamp to be less than 11 then only it will write A else rollback itself & starts with new time stamp

A → read 11  
                 → write 9 | 11

→ → →

\* If time stamps of transactions aren't given assume ;

$T_1 < T_2 < T_3$

Ques Is the given schedule possible?

$\Rightarrow S_1: r_1(A) \quad r_2(A) \quad w_1(A)$        $T_1 \Rightarrow 11$   
           ✓            ✓            ✗       $T_2 \Rightarrow 12$

$A \rightarrow$  read  $10 \cancel{11} 12$

$\searrow$  write 10

\* assume read & write time stamp of A to be less

\* read time stamp of A is not less than time stamp of  $T_1$ , thus, it cannot write A.

Thus, the schedule is not possible.

$\rightarrow$  —

$\Rightarrow S_2: r_1(A) \quad w_2(A) \quad w_1(A)$        $T_1 \Rightarrow 11$   
           ✓            ✓            ✗       $T_2 \Rightarrow 12$

$A \rightarrow$  read  $10 \cancel{11}$

$\searrow$  write  $10 \cancel{12}$

$\Rightarrow$  not possible

$\rightarrow$  —

## # Thomas Write Rule :

- \* The rule for read time stamp is same i.e. check for write time stamp.
- \* When some transaction is writing a data item then read time stamp of data item should be less than or equal to that of transaction and then if write time stamp of data item is less than or equal to that of transaction then <sup>only</sup> transaction should perform write operation.

serial . \*  $r_i(A)$   $w_1(A)$   $w_2(A)$

schedule:  $\hookrightarrow$  useless as no trans. is reading A

$\rightarrow$  if (read time stamp (A)  $\leq T_i$ ) before  $w_2(A)$ .

if (write time stamp (A)  $\leq T_i$ )

}  $T_i$  performs write }

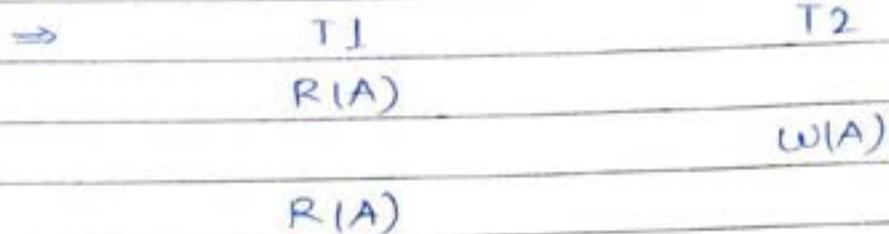
else

}  $T_i$  should ignore its write }

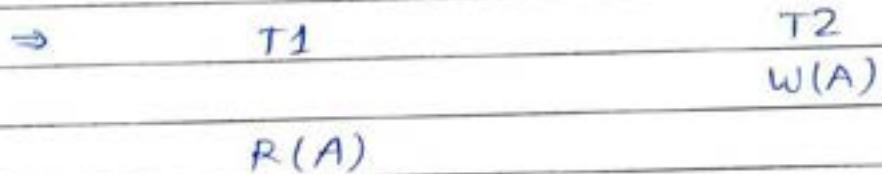
else

}  $T_i$  should rollback

- \* Time Stamp based protocol is deadlock free but cascade rollbacks are possible.
- \* Schedules under basic time stamp based protocol always satisfy conflict serializability test.
- \* Schedules under Thomas write rule in time stamp based protocol may or may not satisfy conflict serializability test but they are always serializable.
- \* concurrent schedules are possible in Thomas rule.

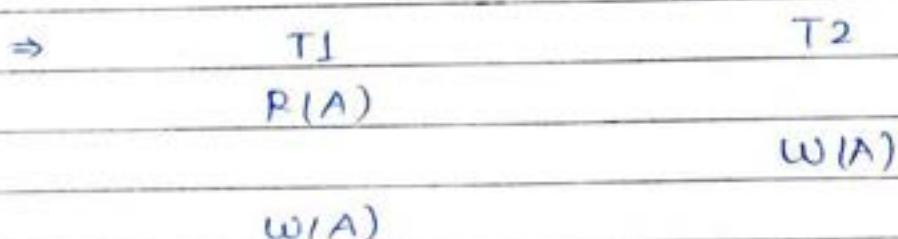


- \* T1 is reading 2 different values of A. Thus, this is known as unrepeatable read.



- \* T1 is reading changes done by uncommitted transaction. It's known as dirty read.

- \* When there's dirty read then the schedule will be required cascade rollbacking.



- \* Here, T2's write is useless i.e. ~~it's lost update~~ A is updated just after by <sup>T1</sup> transaction. It's known as lost update.

# Phantom Phenomenon :

T1

select sum(sal) from emp

T2

insert into emp values  
(100, 'M', 10000)

select sum(sal) from emp



\* row by this statement  
is known as phantom  
or ghost row.

## # Recovery :

- \* To ensure durability, database uses log based recovery.
- \* In log based recovery, database maintains a log file.

log file  $\Rightarrow$  information about operations performed.

```

< T1 start >
< T1 A 100 150 >      // T1 updates A
< T2 start >           from 100 to 150
< T2 B 200 210 >
< T1 C 150 50 >
< T1 commit >
• system crashes
  
```

\* When will system update database <sup>in HD</sup>? On this basis database is of 2 types;

① Immediate update : database is updated immediately.

② Deferred update : database is updated when transaction commits.

- \* Redo changes done by committed transaction in <sup>deferred</sup> update
- \* Undo changes done by uncommitted transaction in <sup>immed.</sup> update
- \* In immediate update, redo is not required.
- \* In deferred update, undo is not required.
- \* In immediate update, new values need not to be stored in log file as they've already been stored in database.
  - ⇒ This log file is too large. thus, database inserts checkpoints between log file.
- \* Those who have committed till a checkpoint, there's no need to undo / redo for them.

→

CS-2015(II)

Ques Consider simple checkpointing protocol;

```

< start T4 >
< write T4, y, 2, 3 >
< start T1 >
< write T1, z, 5, 7 >
< commit T4 >
checkpoint
< start T2 >
< write T2, x, 1, 9 >
< commit T2 >
< start T3 >
< write T3, z, 7, 2 >
.
```

If crash happens now & the system tries to recover using undo & redo operations. What are the content of undo & redo list.

- ⇒ no undo/redo T4 // before checkpoint
- ⇒ redo T2 // committed
- ⇒ undo T1 // not committed
- ⇒ undo T3 //

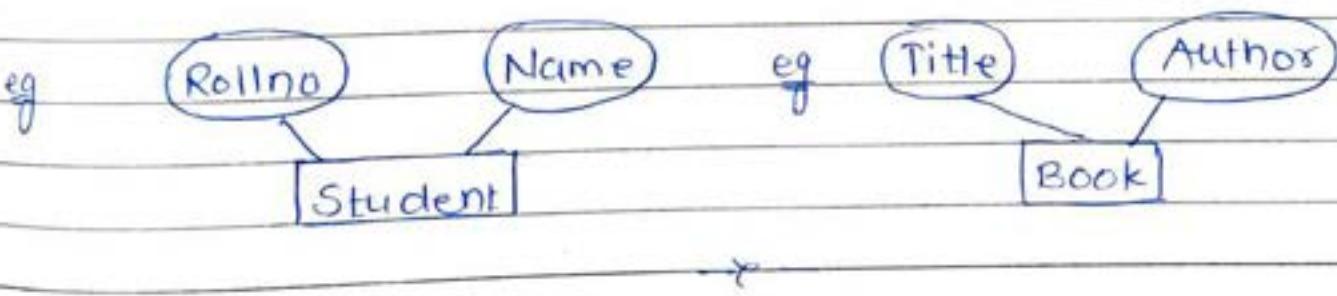
→ —

# ER DIAGRAM

Date \_\_\_\_\_ Page No. \_\_\_\_\_

(207)

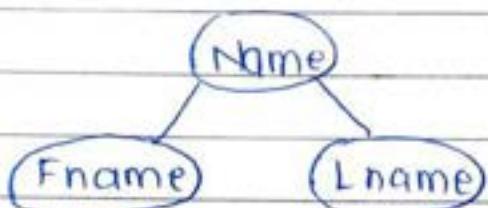
- \* To design a database, we must know how many tables or columns will be present. For this purpose, ER diagram is made.
- \* Any real life object is entity (eg book)
- \* All the books together form a book entity set.
- \* In ER diagram, entity set is represented by rectangle.
- \* Every entity (book) is instance of entity set.
- \* The attributes of entity are represented by ellipse.



## → # Types of attributes:

① Composite attribute : An attribute which can further be divided into subattributes.

eg



② Multivalued attribute : Those attributes which can have multiple values.

eg



\* It's represented by double ellipse.

③ Derived attribute : Those attributes whose value can be derived from other attribute.

eg

age:

\* It's represented by dotted ellipse.

\* age can be derived from DOB.

④ Simple attribute : Those attributes which aren't any of above three are simple.

eg

marks

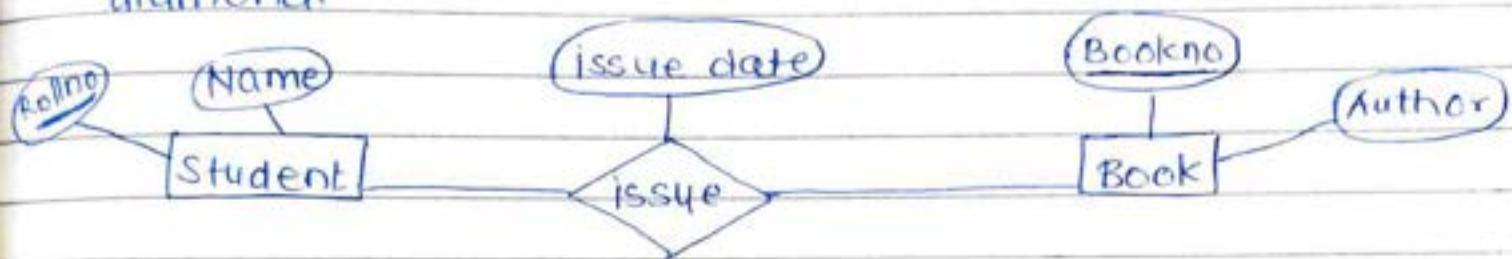
\* unique attribute (primary key) is represented by underline.

e.g.

Rollno



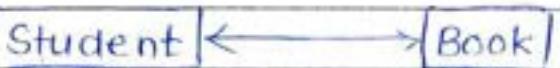
\* Relationship between 2 entities is represented by diamond.



\* issue date is the attribute of issue relationship.

→ There can be 4 types of relationships;

① One to Many : 

② One to One : 

③ Many to One : 

④ Many to Many : 



- \* Initial participation of an entity in relationship is represented by thick or doubled line.

eg

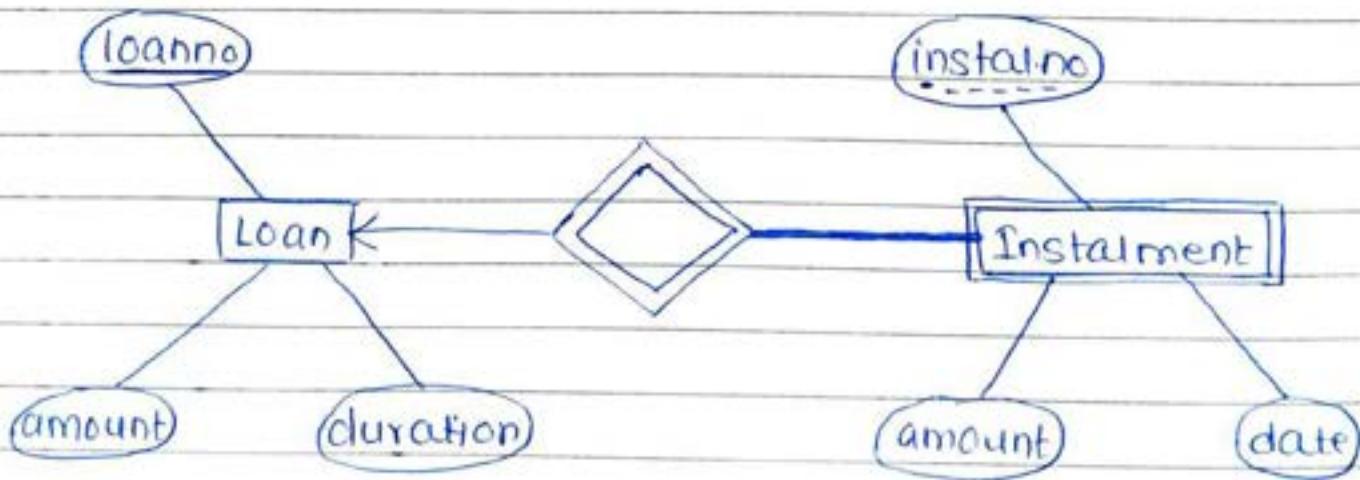


- \* It means every book is issued to some student.



- \* Entities have their own primary key  $\Rightarrow$  Strong entities
- \* Entities don't have primary key  $\Rightarrow$  weak entity  
(don't have even composite)

eg



- \* There is no primary key of instalment entity. Thus, it's a weak entity & represented by double rectangle.
- \* Relationship between a strong & a weak entity is represented by double diamond.
- \* There will always be total participation of weak entity.
- \* There will always be one to many relationship between a strong & a ~~weak~~ weak entity.
- \* That attribute of weak entity which when combined with primary key of strong entity creates a unique key is known as partial key or discriminator & it's represented by dotted line.

eg

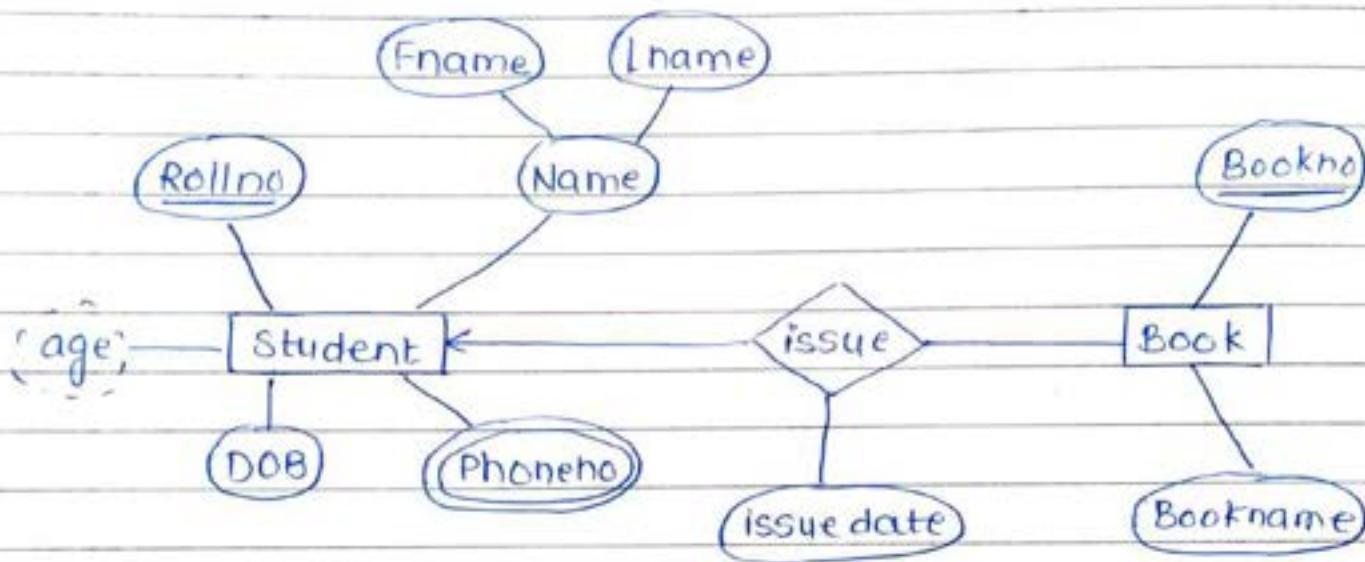
instal no



The relation returned by this method  
will always be in 3NF

Date \_\_\_\_\_ Page No. \_\_\_\_\_

→ How to map ER diagram with a table:



- \* there will be a table corresponding to every entity.
- \* For every simple attribute, there will be a column in table. (eg Rollno, DOB)
- \* For composite attribute, there will be no column but there will be columns for every individual attributes of composite attribute. (eg Fname, Lname)
- \* For derived attribute, there will be no column. (eg age)
- \* For multivalued attribute, there will be no column but a separate table, & in that table, primary key of main table is present along with phoneno.

e.g | Rollno | Phoneno |

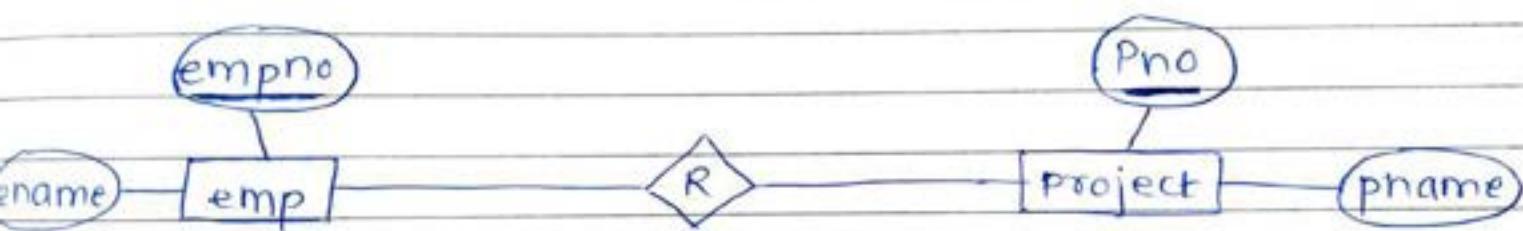
\* In a table of relationship, along with its attributes, primary keys of entities will be present in table

eg | Rollno | issuedate | Bookno |



→ When to create table for relationship :

① Many to many → always create table of relationship



\* In table of relationship, primary key will be composite (empno, Pno)

eg emp:

empno	ename
1	e1
2	e2

Project:

Pno	pname
P1	ABC
P2	DEF

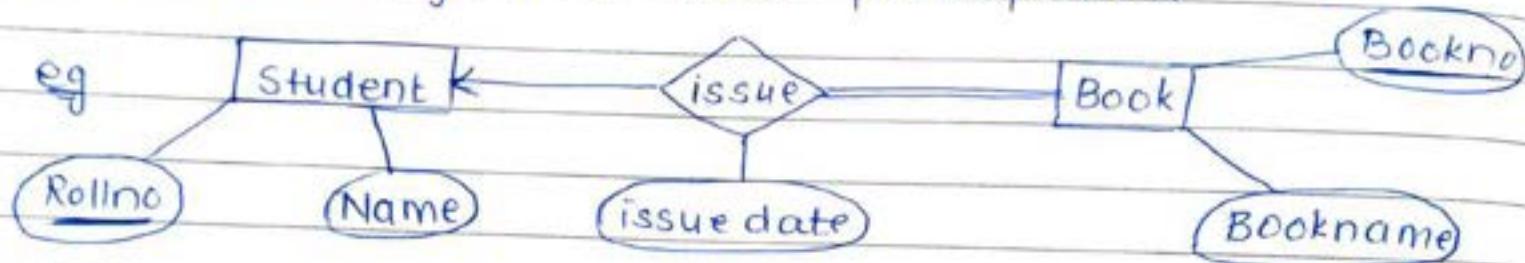
R:

empno	pno
1	P1
1	P2
2	P2

weak entities are always in total participation, thus, the table of their relationship is never made.

② One to Many: create table of relationship in all the cases except;

(a) when many is in total participation.



\* Now, the attributes of relationship will be shifted towards many. i.e. issuedate in table of book.

\* But, there will also be primary key of Student tabl

eg Book:

/ Bookno / Rollno / Bookname / issuedate /

\* If there's no total participation of many, then there will be a lot of books which are not issued to any student due to which large no. of NULL values have to be stored. Thus, we ~~can't~~ create table of relationship in this case.  
1. don't shift the attributes of relationship towards many.

→ But, if question is min. no. of tables, then ans. will be ?  
not 3.

## Book:

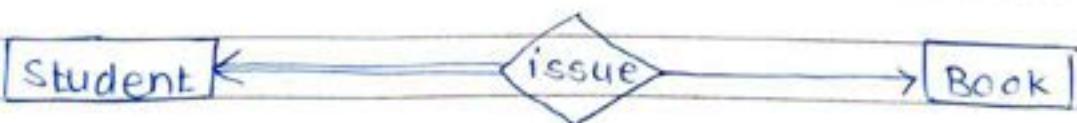
Bookno	Rollno	Bookname	issuedate
b1	1	ABC	20/2/17
b2	2	DEF	22/2/17
b3	NULL	GHI	NULL
b4	NULL	JKL	NULL
b5	NULL	MNO	NULL
1	1	1	1
1	1	1	1
1	1	1	1

\* In one to many, primary key of many will be the primary key of relationship. (eg bookno)

→

③ One to One : Create table of relationship in all the cases except ;

(a) When any 'one' is in total participation.



\* Now, ~~student~~ table the attributes of relationship will be shifted towards table of entity having total participation.

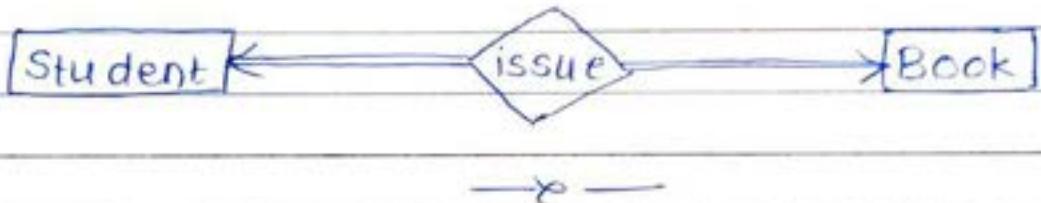
\* But, there will also be primary key of Book table

e.g Student :

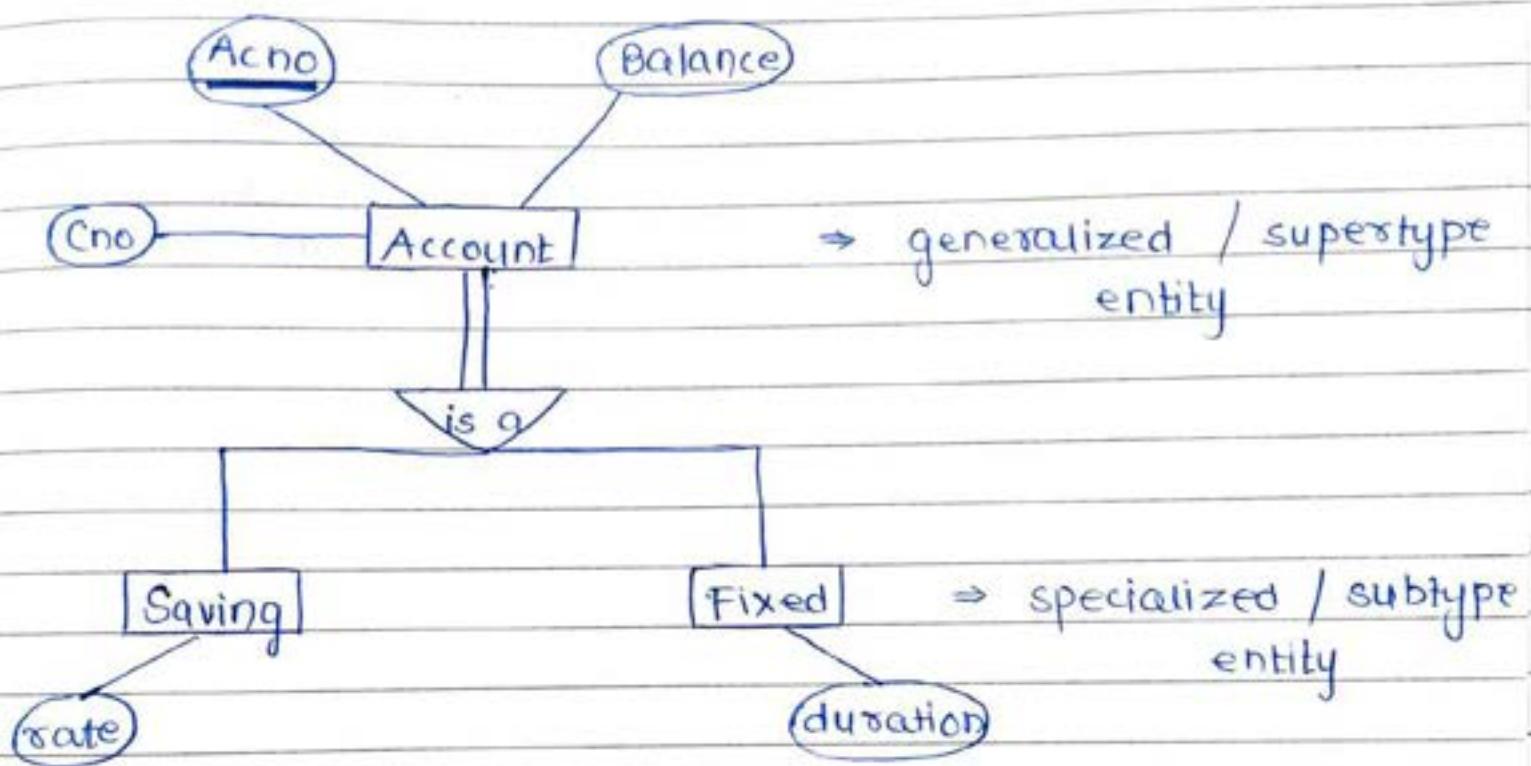
/ Rollno / Name / issuedate / Bookno /

→x→

\* If both have total participation, there will be single table for student & book.



## # Generalisation & Specialization :



- \* interest rate, FD duration, etc. cannot be the attribute of every account.
- \* if every supertype entity belongs to a subtype entity (i.e. if every account is either saving or fixed) then this is known as total specialization.
- \* If there's an account which doesn't belong to subtype (eg. current account) then this is known as partial specialization.

⇒ Mapping into table :

① way : Create tables for supertype & all subtype entities

- \* In ~~super~~<sup>sub</sup> type table , along with its attributes , primary key of supertype will also be there.
- \* every record of subtype will also be present in table of supertype.

OR

② way : Create tables for subtype entities only.

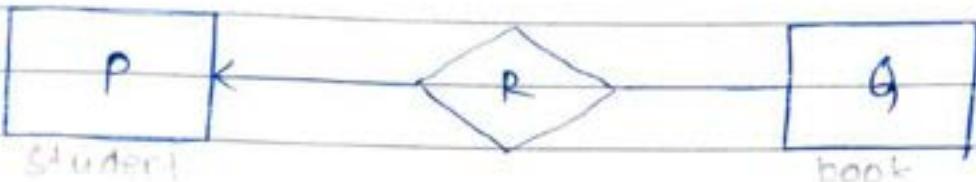
- \* In every table of subtype , along with its attributes , all the attributes of supertype will also be present.
- \* But , it fails when there's partial specialization.

eg Saving:

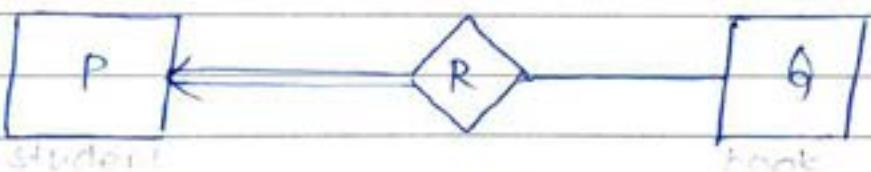
Achno	balance	Cno / rate
-------	---------	------------

Fixed:

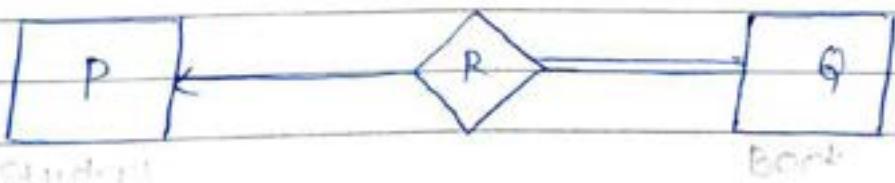
Achno	balance	Cno / duration
-------	---------	----------------

Ques

- (a)  $|P| = |Q|$
- (b)  $|P| \leq |Q|$
- (c)  $|P| \geq |Q|$
- (d) no relationship ✓

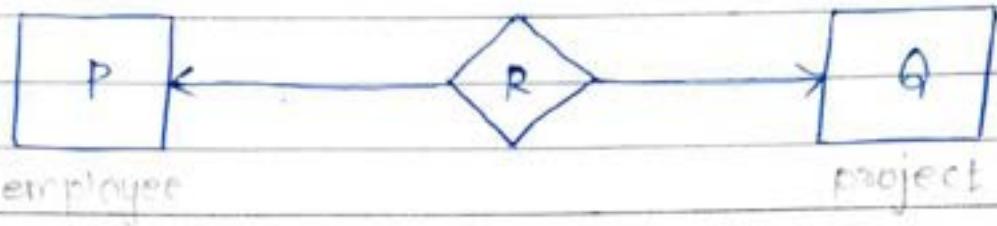
Ans

$$\Rightarrow |P| \leq |Q|$$

Ans

$\Rightarrow$  No relationship

Que what is the max & min cardinality of relationship R?



$$|P| = 10, |Q| = 100, |R| = ?$$

$\Rightarrow$  maximum  $\Rightarrow 10 * 100 = \underline{1000}$  rows in table  
of relationship

$\Rightarrow$  minimum  $\Rightarrow \underline{0}$  rows as there's no total participation.



$$|P| = 10, |Q| = 100,$$

maximum  $\Rightarrow 10 * 100 = \underline{1000}$ , minimum  $\Rightarrow \underline{100}$  (1 emp - 10 proj)



maximum  $\Rightarrow \underline{1000}$ , minimum  $\Rightarrow \underline{100}$  (1 emp - 10 proj)