

Database Management System (8M)

Database Design:

1. E-R model
2. Relational model
3. Functional dependency, Normalization

Implementation:

4. Relational Algebra, Relational Calculus
5. SQL

Maintainance:

6. Transaction Management
7. File Organization

Indexing

Data: Data is a raw fact (Abstract in nature). 90

Information: Meaningful / processed data.

Data is stored in the form of records

Records:

Collection of logically related data.

Database:

Collection of records

Management:

It is done using set of programs.

DBMS:

Data Structure Used

Graphs

Trees

Tables / set

object

objects & tables

DBMS

Network DBMS

Hierarchical DBMS

Relational DBMS (RDBMS)

object oriented DBMS (OODBMS)

object - Relational DBMS (ORDBMS)

Examples of RDBMS

ORACLE, DB-2, SQL-Server, MYSQL, Sybase

Logical Database Design using Relational model:

→ Relation consists of records (tuple).

Each record has certain no. of attributes.

→ Every relation is described using 2 things.

(i) Relation schema

(ii) Relation instance

* Relation schema specifies structure of relation.

Relation schema specifies name of relation, attributes, domain of attribute, conditions on attributes (Eg: primary key)

constraints

* Relation instance specifies the tuples present at a particular time.

Degree of relation: It is no of columns i.e., no of attributes.

Cardinality of relation: It is, no of tuples.

Integrity constraint: It is a condition specified on a database

schema.

There are 5 types of integrity constraints:

- i) Primary key (unique, not null)
- Table constraints (if they are specified on a table)
 - ii) Unique (no duplicates, but null is allowed)
 - iii) not null
 - iv) Check (User defined, conditions) (Eg: age ≥ 18)
 - v) Foreign key (Used to show relation b/w 2 tables)
- referential integrity constraint.

key constraint:

key is a set of attributes which uniquely determines each tuple in a relation.

There are 3 types of key constraints: i.e. candidate key

i) Candidate key:

It is minimal set of attributes which uniquely identifies each tuple in a relation.

i.e. If S is a candidate key then no proper subset of S is a key.

Eg: Consider

student (Roll-no, name, father-name, passport-no)

Here : to find which attributes can be key

(Roll-no) \leftarrow key

(name, father-name) may or may not be key

(Passport-No) may or may not be key

because some students may not have passport.

(Roll-no, name) \leftarrow key

(Roll-no, passport-no) \leftarrow key

Now

(Roll-no) \leftarrow candidate key

(Roll-no, name) is not a candidate key

(ii) Super key:

It set of attributes which uniquely identifies a tuple in a relation.

(or)

Set of attributes to which a candidate key is subset

(or)

Superset of a candidate key

Eg: $(\text{Roll-no}) \leftarrow \text{superkey}$

$(\text{Roll-no, name}) \leftarrow \text{superkey}$.

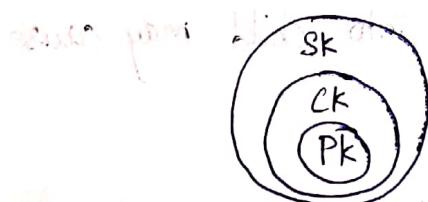
→ Every candidate key is a super key but reverse need not

→ To be true. Here $\{\text{name}\}$ is candidate key and $\{\text{name}\}$ is not superkey.

(iii) Primary key:

Among various no. of primary candidate keys, we choose one ck as primary key.

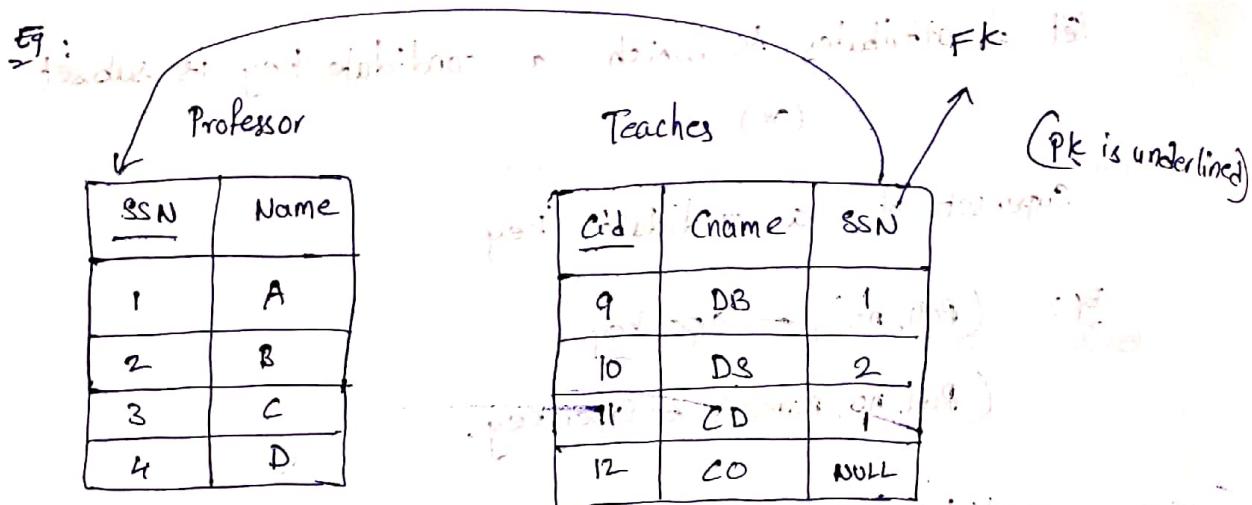
But general design convention pk will less number of attributes is chosen as primary key.



Whenever it is mentioned key, we consider it as candidate key

Foreign key Constraints (Referential integrity constraint) :-

Foreign key is set of attributes used to refer set of attributes of another relation or same relation.



→ FK may have both duplicates & null values.

→ Here Teaches is called Referencing relation (child table). Parent

→ Here Professor is called Referenced relation (parent table). Child

→ Every value under fk column of the child table must be present in PK under PK column of the parent table.

But reverse need not to be true. (NULL is not considered a value, NULL is just a value's absence)

and update

→ Deletion from parent, Insertion into child may cause violations of above rule.

on delete cascade / on update cascade:

→ It is a condition set on FK.

when a violation of deletion / ~~update~~ occurs in parent table, the corresponding rows of child table are deleted / updated.

Set null on deletion violation:

→ when this condition is set on Fk of child table, the Fk values of corresponding rows are set to NULL on deletion.

Set default

→ Same as above condition but a default value is given rather than a null value.

But if this default value is not present under Pk column of parent table, a violation will occur.

→ Fk may also refer to ~~same~~ Pk of the same relation.

Eg:

| Employee - Manager | | |
|--------------------|-------|--------|
| E-No | Ename | MgrENo |
| 1 | Ajay | |
| 2 | Sai | (1) |
| 3 | Sita | 2 |
| 4 | Rama | 2 |

Recursive Relationship

→ NULL

Q1
4-05

Consider below relation R where A is ~~foreign~~ and C is Fk referring to A with on delete cascade.

If $\langle 2,4 \rangle$ is deleted then what are the extra tuples deleted to preserve the referential integrity?

- a) $(5, 2)$ and $4 \neq 5 \Rightarrow 5 > 4$ no violation
- b) $(3, 4)$ is the only pair of numbers
- c) $(5, 2), (7, 2), (9, 5)$
- d) $(3, 4), (4, 3), (6, 4)$

| A | C |
|---|---|
| 2 | 4 |
| 3 | 4 |
| 4 | 3 |
| 5 | 2 |
| 7 | 2 |
| 9 | 5 |

Solution: It shows that there is no violation of the condition.

del $\langle 2, 4 \rangle$

relation now in soft order

del $\langle 2, 4 \rangle$ will ~~cascade~~ for $\langle 2, 4 \rangle$ function with $b_1 = 2$

now $\langle 5, 2 \rangle, \langle 7, 2 \rangle$ didn't have to go

|

$\langle 9, 5 \rangle$

addition of $\langle 9, 5 \rangle$ to the end of the soft order

$\therefore \langle 5, 2 \rangle, \langle 7, 2 \rangle, \langle 9, 5 \rangle$

10/11/20

Q2
a-01 Consider a relation geq which represents "greater than or equal to" that is, $(x, y) \in geq$ only if $y \geq x$

{

create table geq

lb integer not null

ub integer not null

primary key lb

(foreign key (ub) references geq on delete cascade)

} which of the following is possible if a tuple (x, y) is deleted?

- a) A tuple (z, w) with $z > y$ is deleted
- b) A tuple (z, w) with $z > x$ is deleted
- c) A tuple (z, w) with $w < x$ is deleted
- d) The deletion of (x, y) is prohibited

sd: $\langle lb, ub \rangle$

Consider $\langle 8, 10 \rangle$ is to be deleted

then every tuple of form $\langle a, 8 \rangle$ is deleted

then every tuple of form $\langle b, a \rangle$ is deleted

and so one

Here $b \leq a \leq 8 \leq 10$

$$\therefore b < 10$$

$\langle 2, w \rangle$ is deleted

$$\Rightarrow 2 \leq w \leq 8$$

\therefore opt C

Here ~~is~~ is wrong
is correct but
not $b \leq 8$
Ex: $\langle 2, 3 \rangle$
 $\langle 1, 3 \rangle$
Selection of $\langle 2, 3 \rangle$
doesn't cause
deletion of $\langle 1, 3 \rangle$

Note:

In on delete cascade, ~~one delete update~~ on update cascade,
on delete set null, on delete set default, the child table
tuple are modified first.

P/B

If the question is like

s is foreign key referencing p in table T₁ with
on delete cascade and on update cascade then

$\langle 3, 8 \rangle$ will try to delete $\langle 8, 3 \rangle$
(T₁) (T₂)

now $\langle 8, 3 \rangle$ (T₂) will try to delete $\langle 8, 8 \rangle$ of T₁
 $\langle 5, 8 \rangle$
 $\langle 9, 8 \rangle$

so this infinite loop continues.

Conceptual Database Design using E-R model:

→ ER model: Entity Relationship model

→ we convert concept (requirements) to E-R model.

E-R Components:

i) Entity: Real world object. (In general entities are nouns)

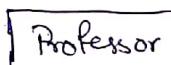
Ex: Professor, Course

ii) Entity set:

Set of entities

Ex: set of professors.

→ Represented as rectangle



iii) Relationship:

It is association b/w two entities or more entities.

Ex: Professor teaches a course.

↓
relationship

(Generally relationship is identified by verbs)

iv) Relationship set:

set of similar relationships.

represented using rhombus



v) Attribute:

Every entity is associated with set of attributes.

It is represented using oval.



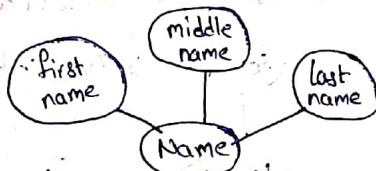
Classification of Attributes:

(i) Simple Attribute:

An attribute is said to be simple, if it cannot be divided into further attributes.

(ii) Composite Attribute:

It can be divided into set of simple attributes.



Based on no of values associated:

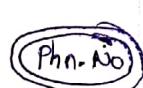
(i) Single valued attribute:

The attribute can take maximum of one value only.

(ii) Multi valued attribute:

The attribute can take more than one value.

It is represented using double ellipse.



Based how the attribute value is stored:

(i) Stored attribute:

An attribute which is stored in database and its value cannot be derived from other attributes.

Derived attribute:

An attribute is said to be derived if can be derived from other attributes.

It is represented using dotted oval.

Eg:

Age

Graduation

Key attribute:

It is the attribute which identifies every entity in the entity set.

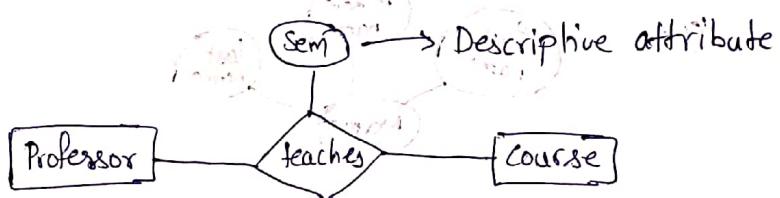
Key attribute is represented as underlined attribute.

Roll.no

Descriptive attribute

Descriptive attribute is used to describe a relation.

Eg:

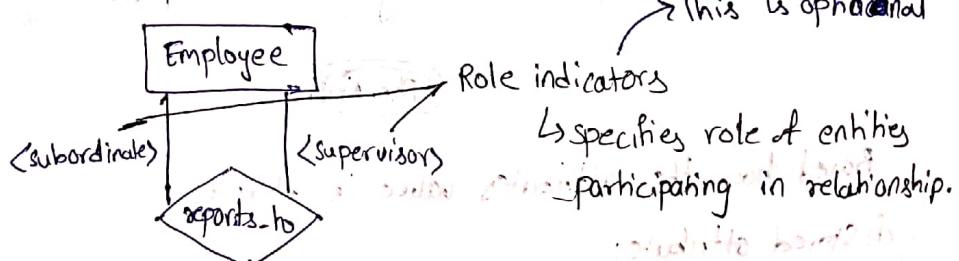


Relationship Degree:

Degree of a relationship is no of entities associated with a relationship.

i) Unary relationship / Recursive Relationship

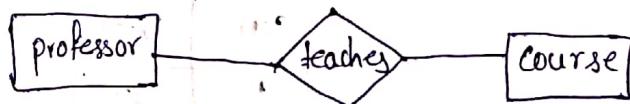
Relationship with degree 1.



This is also known as recursive relationship.

(ii) Binary relationship:

Relationship of degree 2.



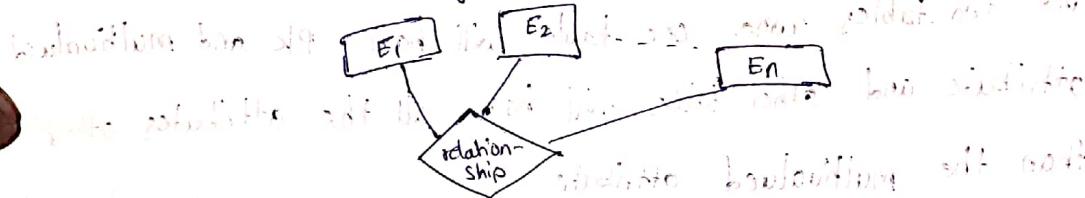
(iii) Ternary relationship:

Relationship of degree 3.

~~n-ary~~

n-ary relationship:

Relationship of degree n.

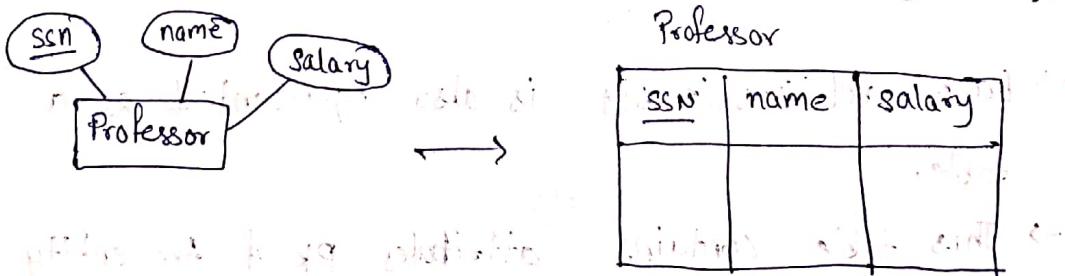


11/11/20

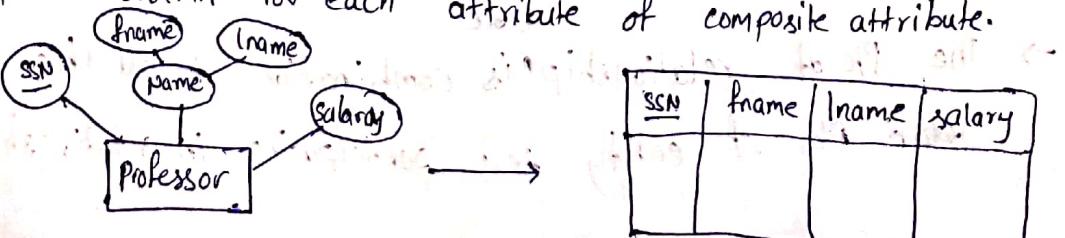
ER to relational model

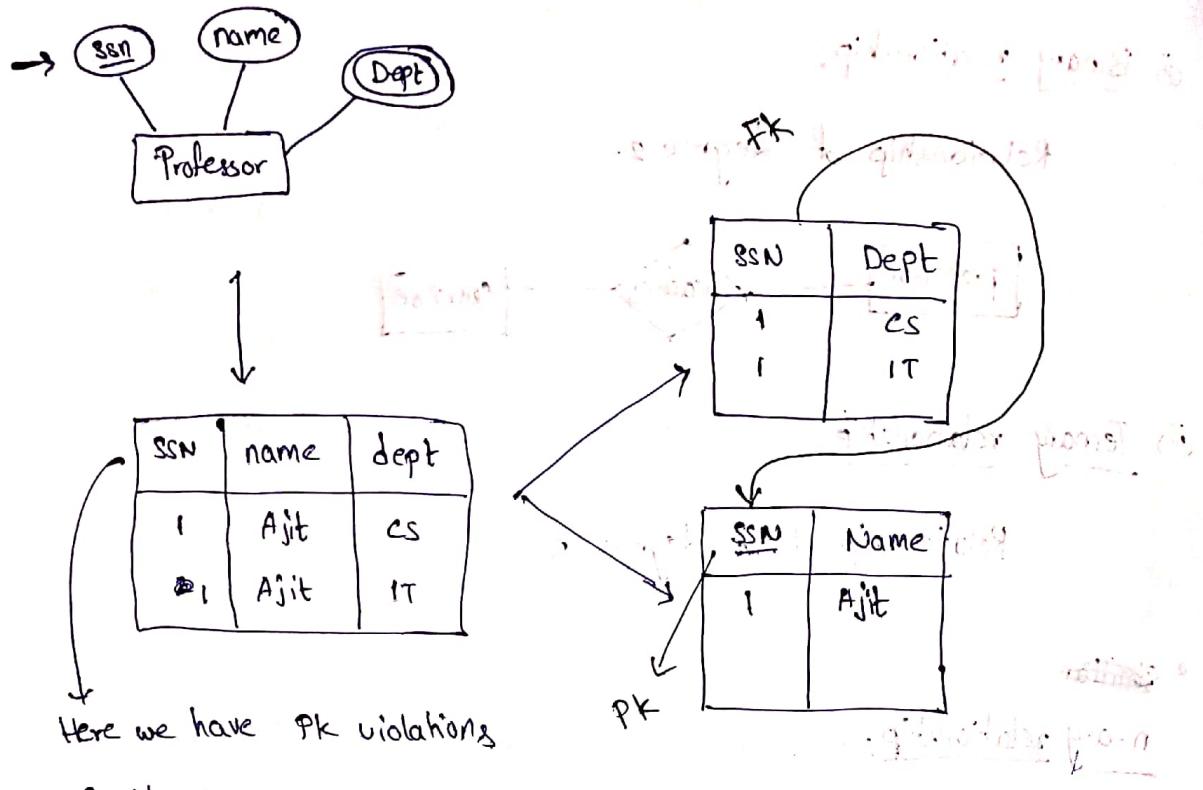
→ Every entity set is represented as one table.

→ The attributes of entity set becomes attributes of table.



→ If we have composite attribute in ER model, then we create separate column for each attribute of composite attribute.





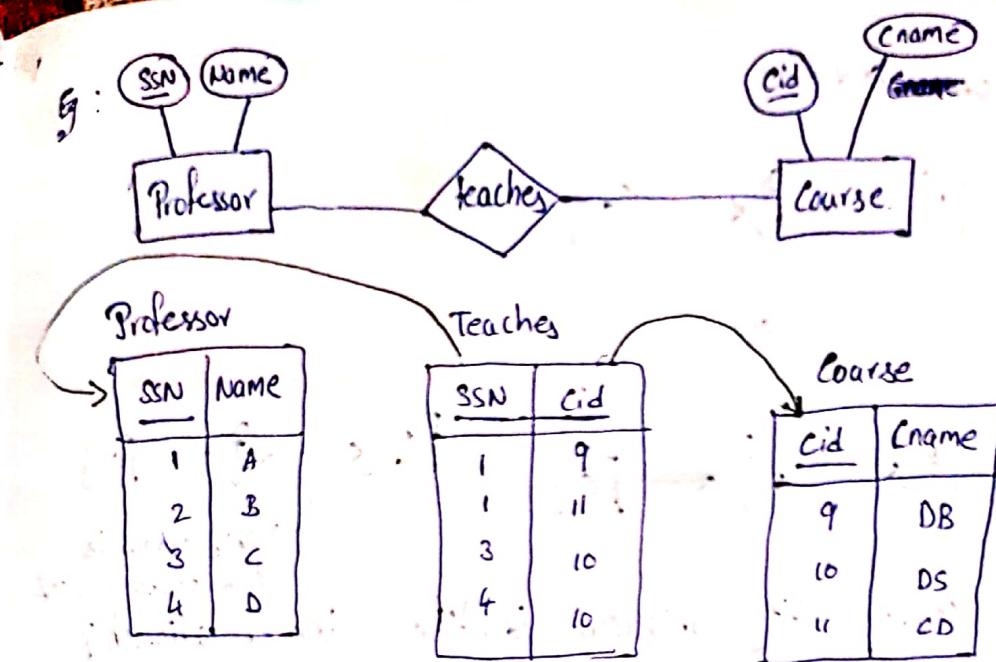
Here we have PK violations

So the table is divided into two tables where one table will have PK and multivalued attribute and other table will have all the attributes other than the multivalued attribute

- If an entity set has multivalued attribute then it has to be represented using 2 tables.

Converting relationship sets into relational model

- Each relationship set is also represented as a table.
- This table contains attributes PK of the entity sets participating in relationship as FK, and other attributes specific to the relationship table.
- The PK of relationship is combination of all the PKs of entity sets participating in relationship.



PK
SSN → FK
Cid → FK
 $\{SSN, Cid\} \rightarrow PK$

key Constraints in E-R model:

If for each entity in the entity set there is atmost one relationship in the relationship set then the ~~entt~~ entity set is said to be a key constraint.

The key constraint is represented with arrow (from entity set to relationship set)

~~one more way is~~

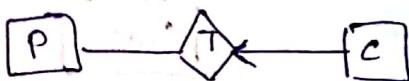
Eg: Each course is taught by atmost one professor.

Here course is the key entity set.

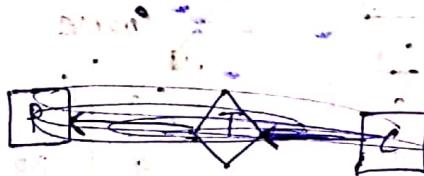
Professor Teaches Course

| Professor | | Teaches | | Course | |
|-----------|------|---------|-----|--------|-------|
| SSN | Name | SSN | Cid | Cid | Cname |
| 1 | A | 1 | 9 | 9 | DB |
| 2 | B | 1 | 11 | 10 | DS |
| 3 | C | 3 | 10 | 11 | CD |
| 4 | D | 4 | 12 | 12 | CO |

It is represented as



(or)



But sometimes numbering may introduce null values which will lead to wastage of space.

But combining can make certain operations (involving both tables) execute faster.



This arrow is from course to prof

- Here each course entity uniquely identifies the relationship.
- Thus the key of relationship table is cid.
- Thus we combine teaches & course tables.
- The resultant tables are:

| Professor | | Course-teacher | | |
|-----------|------|----------------|-------|-----|
| SSN | Name | cid | Cname | SSN |
| 1 | A | 9 | DB | 1 |
| 2 | B | 10 | DS | 3 |
| 3 | C | 11 | CD | 4 |
| 4 | D | 12 | CC | |

Participation Constraints:

total participation:

If every entity of an entity set participates in the relationship then it is called total participation.

Eg: Each course is taught by exactly one professor

So here we make ssn column as not null in course-teachey table (in above example)

It is represented using ~~\leftrightarrow~~

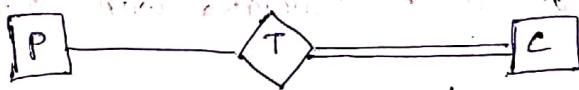
~~see for the relationship~~

~~Directed from entity A to entity B~~

(none of which) \rightarrow



total participation with key constraint



total participation

Partial Participation:

If ~~some~~ some entities in the entity set are not participating in a relationship then it is called partial participation. It is represented using —

Note:

atmost one :

(every course taught by atmost one prof)



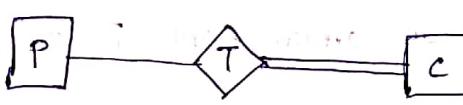
exactly one :

(every course taught by exactly one prof)



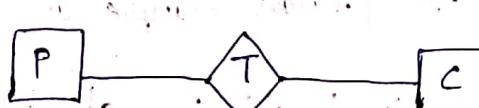
atleast one :

(every course taught by atleast one prof)



0 (or) many :

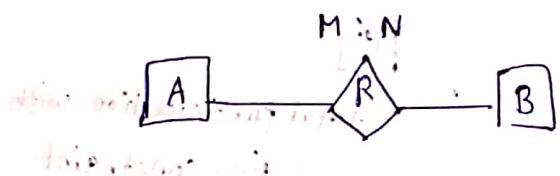
(every course is taught by 0 (or) More prof)



Course entity is not a key in this case.

Cardinality ratios:

i) M:N (many to many)



This means:

→ Each entity of A is associated with 0 (or) more entities of B.

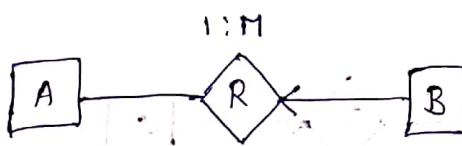
→ Each entity of B is associated with 0 (or) more entities of A.

→ To represent this in relational model,

we need three tables: ~~A, B~~, A, R, B

PK of R is {PK of A, PK of B}

ii) 1:M (one to many)



→ Each entity of A is associated with 0 or more entities of B.

→ Each entity of B is associated with 0 (or) 1 entity of A.

So we draw an arrow from B to R.

→ It is similar to the previous example of

each course is taught by atmost one professor

→ This represented using 2 tables.

↳ (A & R-B)

PK of R-B is → PK of B

(iii) M:1 (Many to one)



→ It requires 2 tables. (A-R & B)

pk of A-R is pk of A

A-R has pk of B as FK

Note:

In 1:M & M:1 → is from 'M' side.

(iv) 1:1



⇒ Each entity in A is associated with atmost one entity of B.

→ Each entity in B is associated with atmost one entity of A.

→ It is represented using 2 tables.

either (AR, B) or (A, RB)

\downarrow
pk = pk of A
fk = pk of B

→ we cannot represent this using one table because for some key value of A, corresponding B can be null and vice versa.

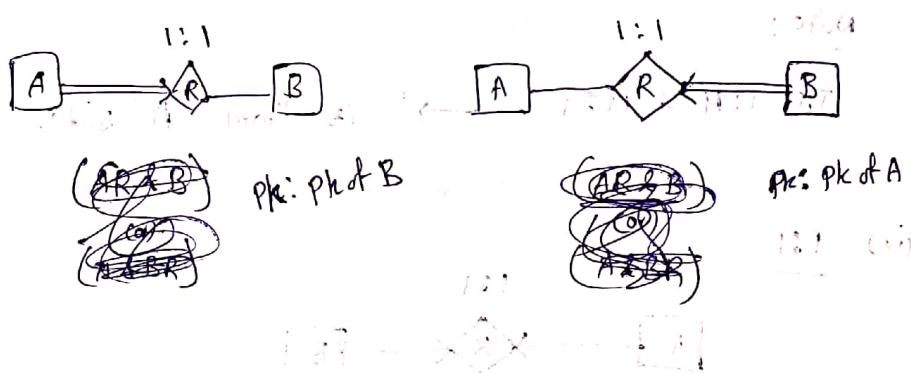
In this case we cannot choose primary key.

Note:

→ In 1:1 cardinality, if both the participations are total then one table will be enough.



→ In 1:1 cardinality, if total participation is only from one side then we need ~~two tables~~ minimum one table

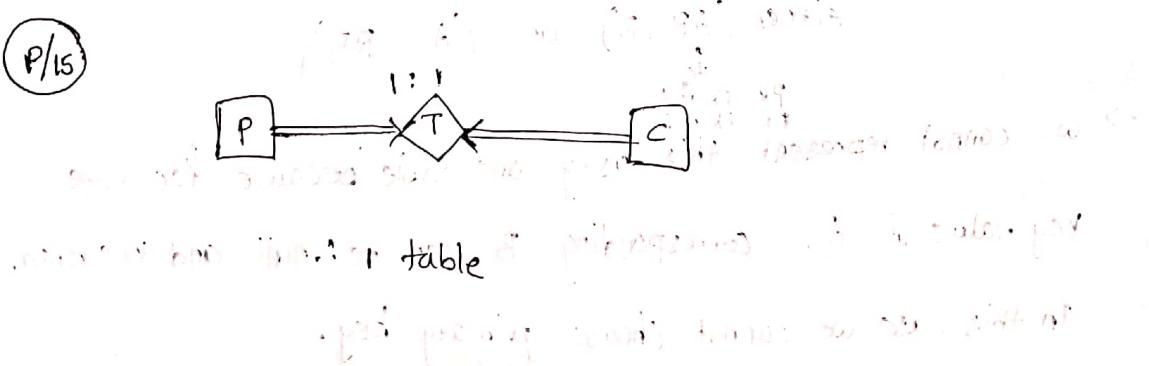


P/5

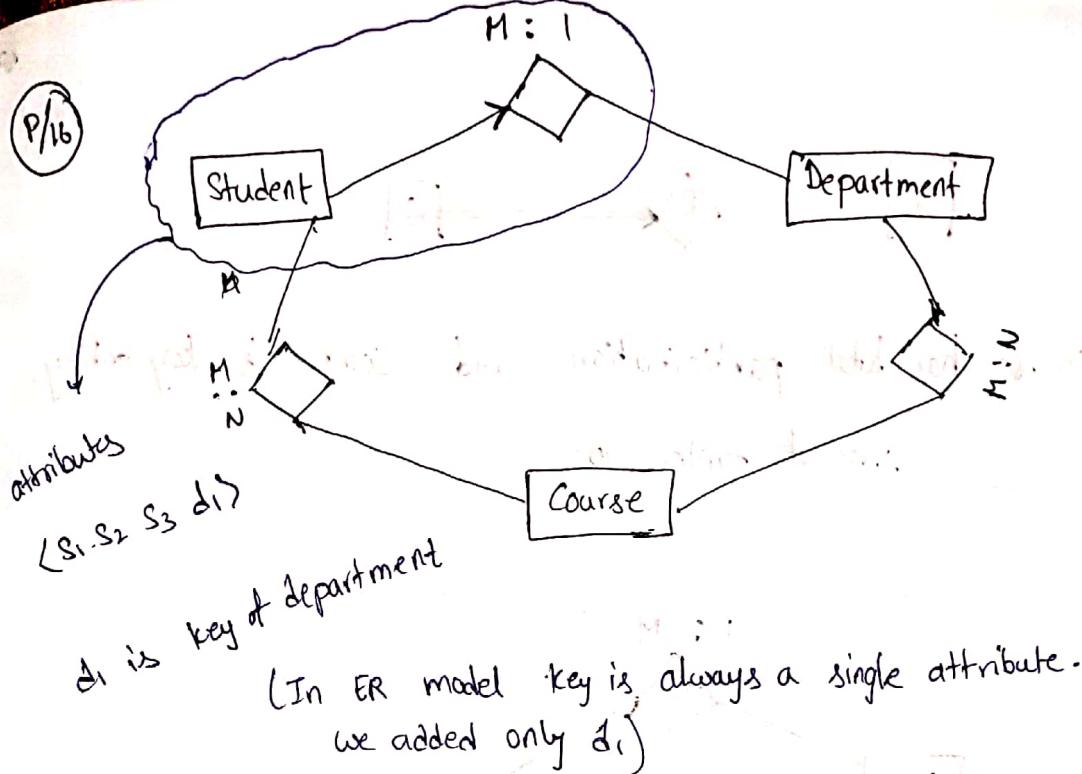
arrow is present from car to purchase. \therefore car is key in the relationship.

\therefore Each car has almost one dealership; or vice versa.

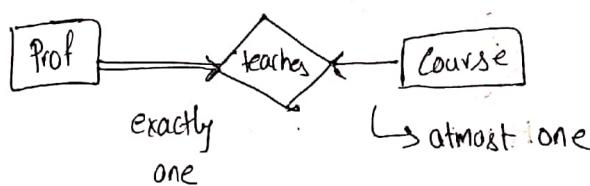
P/15



P/16



P/17



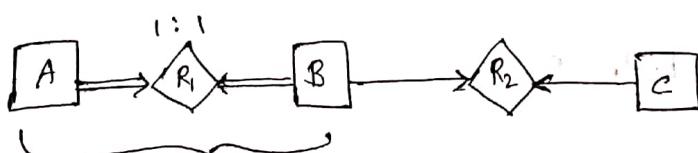
| ssn | Name | cid | cname |
|------|------|-----|-------|
| 1 | A | 9 | DB |
| null | null | 10 | CO |

cid column can never have null.

\therefore cid is pk.

12/11/20

P/18



can be represented
using one table



Now **R₂** can be merged with **AB** or **C**.

\therefore 2 tables.

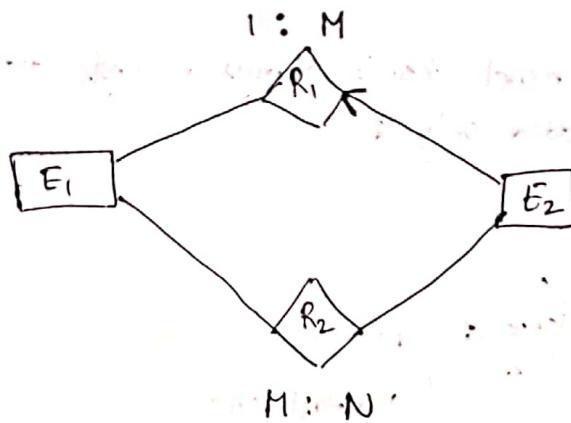
P/19



course has total participation and course is key entity.

∴ no of tuples = 100

P/20



~~R_{E2} can be combined~~

Initially we start with 4 tables

E1, E2, R1, R2

we can combine E2 & R1

∴ we finally have 3 tables

E1, E2R1, R2

P/21

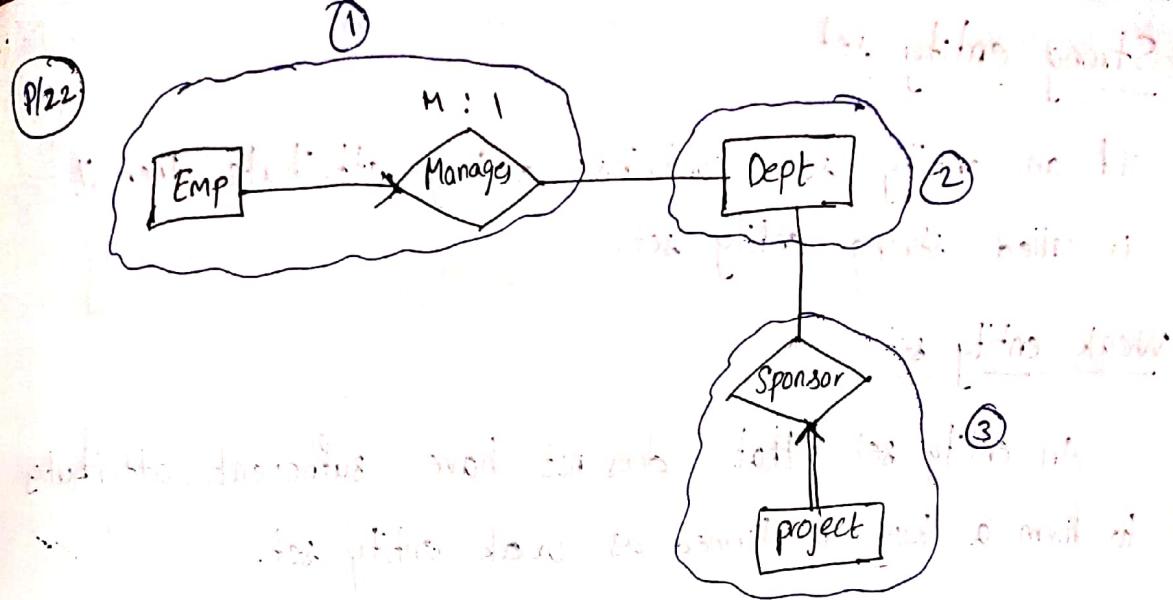


E2 has a multivalued attribute.

∴ It requires 2 tables.

we can combine E2 & R.

So we have 3 tables.



$\therefore 3$ relations

P/23

a)



two diff. relations & 2 beginning tables shown in 1 table
2 tables

b)

M:N

3 tables

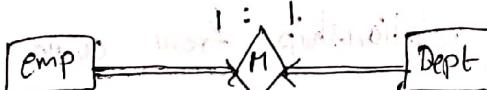
three diff. relations before putting into 1 table
3 tables

c)

1:M

2 tables

d)



$\therefore 1$ table

two diff. relations & 2 beginning tables shown in 1 table
1 table

the best way to do this is to make it 1 table
1 table

the best way to do this is to make it 1 table
1 table

Strong entity set:

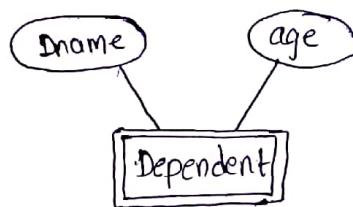
If an entity set contains a key attribute, then it is called strong entity set.

Weak entity set:

An entity set that does not have sufficient attributes to form a key is termed as weak entity set.

weak entity set is represented using double rectangle.

E:



- A weak entity set participates in a relationship and it is identified using key of another entity set (strong entity set).

This strong entity set is called owner entity.

This relationship is called weak relationship/identifying relationship.

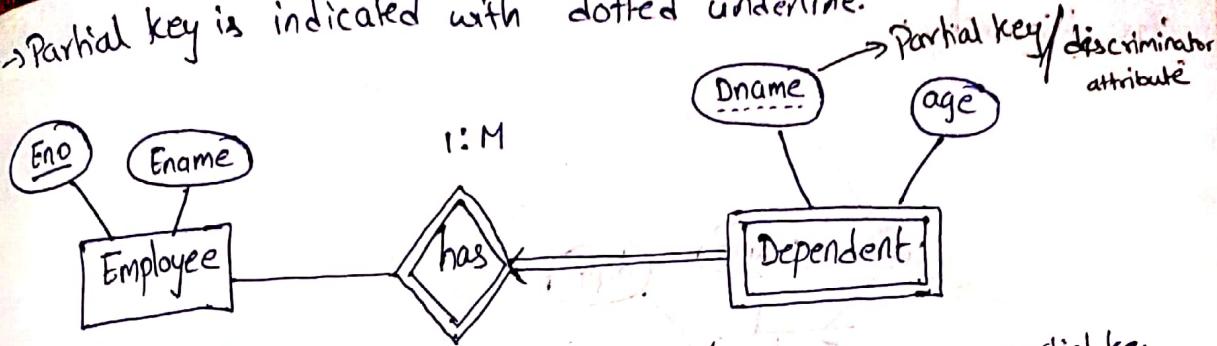
→ The cardinality of relationship from owner to weak is always 1:M

→ Weak entity must have total participation.

→ The key of weak entity set is {key of owner + partial key}

→ The attribute of weak entity set which when combined with key of owner behaves as key of weak entity set is called partial key.

→ Partial key is indicated with dotted underline.



| <u>E#</u> | Ename |
|-----------|-------|
| 1 | P |
| 2 | Q |
| 3 | R |
| 4 | S |
| 5 | T |

| <u>E#</u> | Dname |
|-----------|-------|
| 1 | A |
| 1 | B |
| 2 | A |
| 2 | B |
| 3 | C |
| 4 | C |
| 4 | C |

| <u>E#</u> | Dname | age |
|-----------|-------|-----|
| 1 | A | 12 |
| 1 | B | 6 |
| 2 | A | 10 |
| 2 | B | 9 |
| 3 | C | 3 |
| 4 | C | 3 |

Foreign key

{E#, Dname}

↓
key
↓
partial key
owner

| <u>E#</u> | Dname | age |
|-----------|-------|-----|
| 1 | A | 12 |
| 1 | B | 6 |
| 2 | A | 10 |
| 2 | B | 9 |
| 3 | C | 3 |
| 4 | C | 3 |

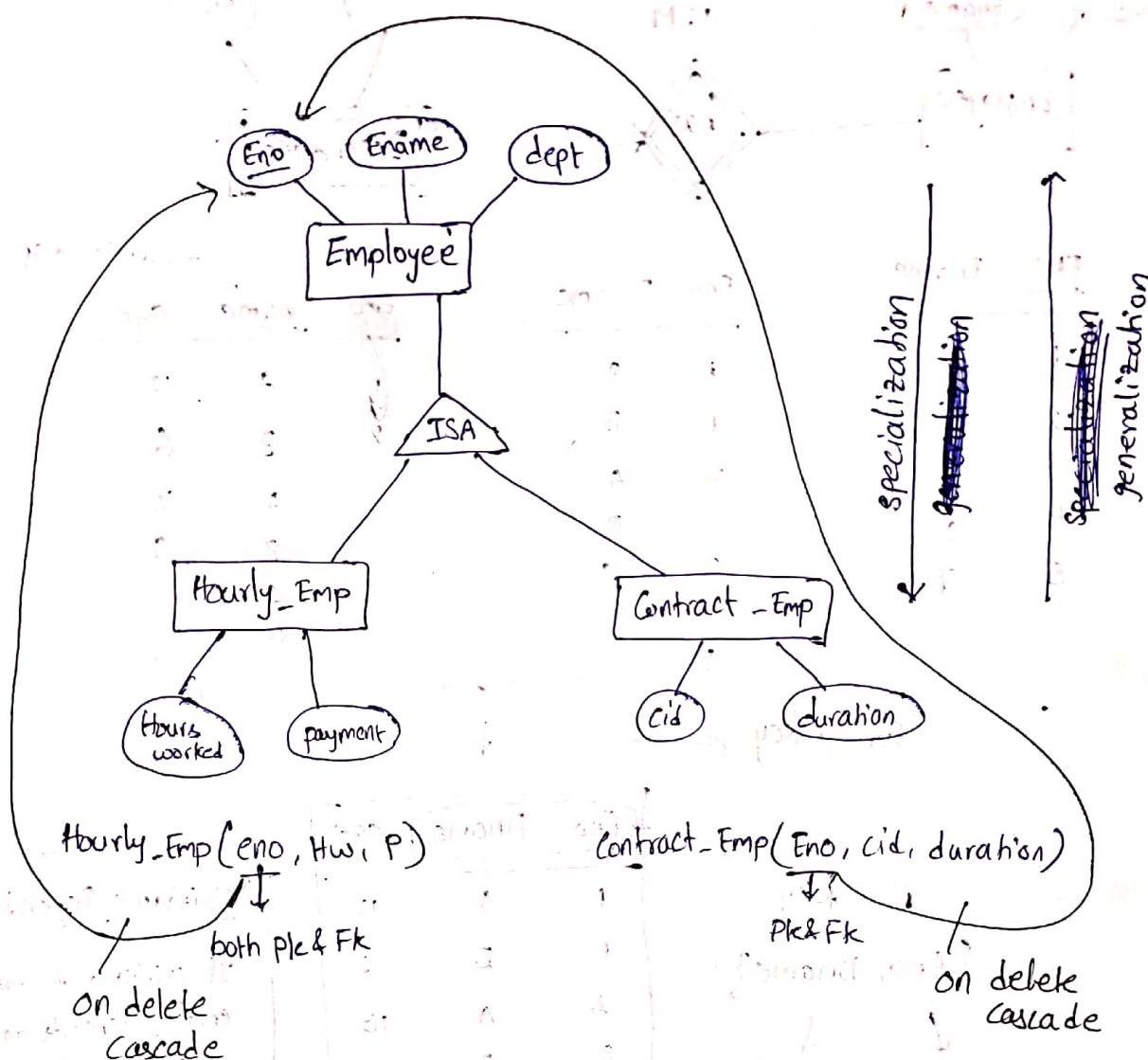
Existence Dependency
If existence of one entity depends on exists of other entity, then it is called existence dependency

→ Existence of a weak entity is possible & meaningful only when its owner entity exists.

∴ The Fk : key 'E#' in weak entity set 'has' has non delete cascade ~~on~~ Fk constraint.

∴ Fk of weak entity set always has non delete cascade relationship.

Class hierarchy:



→ whenever an entity of employee is deleted, the corresponding entity must be deleted from hourly_emp or contract.emp.

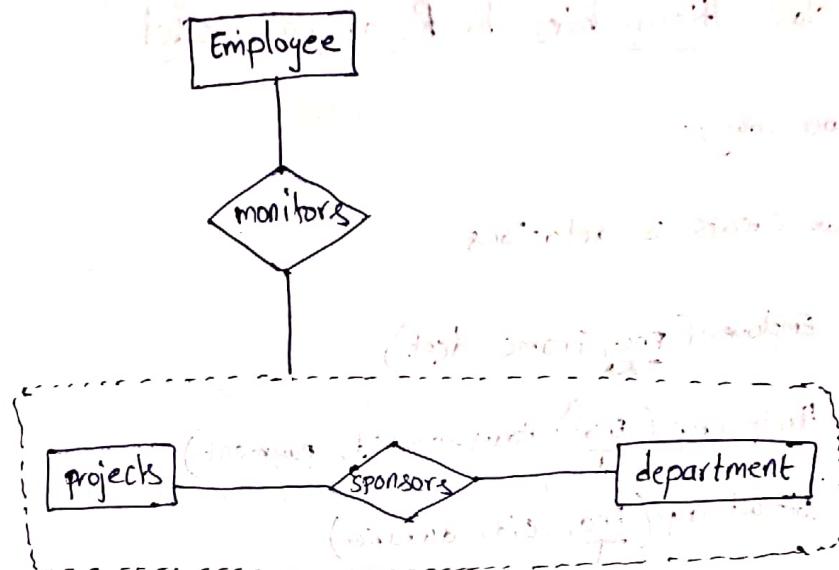
∴ FK constraint is on delete cascade

Aggregation:

Aggregation is relationship over a relationship.

Eg: Consider

Employee monitors the sponsorship of projects by the department.



→ This is represented with 5 tables.

P/24

3 tables

(i) MR₁

(ii) P₁ → {P₁, P₂, P₃, P₄}

(iii) R_{2N}

P/25

MR₁ → {M₁, M₂, M₃, P₁}

P → {P₁, P₂}

R_{2N} → {P₁, N₁, N₂}

∴ opt(a)

In ISA hierarchy, we have 2 constraints
 (i) overlap: says whether two subclasses
 can have same entity.

(ii) Covers: says whether the ~~to~~ all
 subclasses must form all the
 entities of generalized entity set.

→ If overlapping is allowed it has to
 be specified explicitly

→ covering constraints also has to be explicitly
 specified

Schema Refinement:

Problems due to redundancy:

(i) Repeated storage

(ii) Insertion problem: It may not be possible to store certain information unless some other unrelated info is stored as well.

(iii) Updation problem: If one copy of repeated data is updated, inconsistency is created unless all the copies are updated.

(iv) Deletion problem: It may not be possible to delete certain information without losing other unrelated information.

(i) Repeated storage

Assume we have a requirement that

"Employees of same seniority should have same salary".

| Eno | Ename | Seniority | salary |
|-----|--------|-----------|--------|
| 1 | Ajith | 5 | 50k |
| 2 | Raju | 5 | 50k |
| 3 | Kamal | 5 | 50k |
| 4 | Harsha | 10 | 1L |
| 5 | Sharma | 10 | 1L |

(ii) Insertion problem:

If we need to insert

$\langle 6 \text{ shastri} \ 10 \ 1.2L \rangle$

this causes violation to customers requirement.

(iii) Update problem:

Assume we need to change salary of Ajith to 60k.

This ~~will~~ causes violation.

(iv) Deletion problem:

If tuples 4 & 5 are deleted, the association (10-1L) will be lost.

→ The problem's ~~is~~ ~~this~~ solution is decomposition.

We form below two tables.

Emp:

| Eno | Ename | seniority |
|-----|--------|-----------|
| 1 | Ajith | 5 |
| 2 | Raju | 5 |
| 3 | Kamal | 5 |
| 4 | Harsha | 10 |
| 5 | Sharma | 10 |

Emp-Salary:

| Seniority | Salary |
|-----------|--------|
| 5 | 50k |
| 10 | 1L |

Now all the 4 problems are solved.

→ Here we say salary is functionally dependent on seniority.

Functional Dependencies

If A & B are attributes of a relationship and if for each value of A there exactly one value of B then we say "B is functionally dependent on A" or "A functionally determines B".

It is denoted as $A \rightarrow B$

If $A \rightarrow B$ then

| A | B |
|-------|-------|
| x_1 | y_1 |
| x_1 | y_1 |
| x_2 | y_1 |
| x_2 | y_1 |
| x_3 | y_2 |
| x_3 | y_2 |

valid

| A | B |
|-------|-------|
| x_1 | y_1 |
| x_1 | y_2 |

invalid

Eg: Consider below relation

| A | B | C |
|-------|-------|-------|
| a_1 | b_1 | c_1 |
| a_2 | b_1 | c_1 |
| a_1 | b_2 | c_1 |

which of the following FD's are satisfying

- a) $A \rightarrow B$ ✗
- b) $B \rightarrow C$ ✓
- c) $C \rightarrow A$ ✗

- d) $C \rightarrow B$ ✗
 e) $A \rightarrow C$ ✓
 f) $B \rightarrow A$ ✗
 g) $AB \rightarrow C$ ✓
 h) $AC \rightarrow B$ ✗
 i) $BC \rightarrow A$ ✗
 j) $AC \rightarrow C$ ✓

Trivial FD:

The FD in which RHS attributes is subset to LHS attributes.

otherwise it is called non-trivial FD.

13/11/20

Properties of Functional Dependencies:

(i) Reflexivity:

If $Y \subseteq X$, then

$X \rightarrow Y$ is an FD

Eg: $AB \rightarrow B$

The FDs under this property are trivial FDs

(ii) Augmentation:

If $X \rightarrow Y$ then $XZ \rightarrow YZ$ is a valid FD

(iii) Transitivity:

If $X \rightarrow Y$ & $Y \rightarrow Z$ are two valid FDs then
 $X \rightarrow Z$ is also a valid FD.

(iv) Union:

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

(v) Decomposition:

If $X \rightarrow YZ$ then

$X \rightarrow Y$ and $X \rightarrow Z$

$\{X, Y\} \rightarrow Y$

$\{X, Z\} \rightarrow Z$

Closure set of FD's:

If F is a set of function dependencies,

closure of F , (F^+) is set all valid FDs possible from given set of dependencies.

Ex: Consider relation $R(ABC)$, and its FDs

$F: \{A \rightarrow B, B \rightarrow C\}$

Computing F^+ equal to $\{A \rightarrow B, B \rightarrow C, A \rightarrow C\}$

$A \rightarrow B, B \rightarrow C$

$A \rightarrow C$ (transitivity)

$A \rightarrow BC$ (union)

$AC \rightarrow BC$ (Augmentation)

$AC \rightarrow B, A \rightarrow C$ (Decomposition)

⋮

Closure set of attributes:

If X is a set of attributes, then closure of X , (X^+) is set of all attributes that can be determined by X .

Eg: Consider

$$R(ABC) \text{ s.t. } F = \{A \rightarrow B, B \rightarrow C\}$$

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

$$B^+ = \{B, C\}$$

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

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

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

$$(BC)^+ = \{B, C\}$$

$$(ABC)^+ = \{A, B, C\}$$

→ If closure of X includes all the attributes of relation, then X is a key.

∴ In above example A, AB, AC, ABC are keys and A is a candidate key.

AB, AC, ABC are superkeys but not ck.

No of CKs: 1

No of SKs: 4

Functional dependencies

(partial dependency)

Independent functional dependencies

Eg: Consider $R(ABCDE)$ with dependencies $\{AB \rightarrow C, B \rightarrow D, A \rightarrow E\}$

$$F: \{AB \rightarrow C, B \rightarrow D, A \rightarrow E\}$$

Find

i) AB^+

$$\begin{aligned} AB^+ &= \{A, B, C, D\} \\ &\downarrow \quad \swarrow_{B \rightarrow D} \\ &\therefore AB \rightarrow C \end{aligned}$$

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

ii) B^+

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

iii) E^+

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

iv) BE^+

$$\cancel{B \rightarrow E} \quad E \rightarrow A$$

$$AB \rightarrow C$$

$$B \rightarrow D$$

$$\therefore BE^+ = \{A, B, C, D, E\}$$

Applications of attribute closures: To find out if a dependency is possible or not.

i) Additional FDs can be determined.

Given a set of dependencies and if it is asked to check if $X \rightarrow A$ is possible or not.

If $A \subseteq X^+$ then we say $X \rightarrow A$ is possible.

Q3
G-05

Consider R(ABCDE) the given set of FDs: {A → B, A → C, CD → E, B → D, E → A}.

A → B; A → C; CD → E; B → D; E → A

which of the following FDs is not implied

- a) CD → AC
- b) BD → CD
- c) BC → CD
- d) AC → BC

Sol:

- a) CD → AC

$$(CD)^+ = \{C, D, E, A, B\}$$

- b) BD → CD

$$(BD)^+ = \{B, D\}$$

∴ false

- c) BC → CD

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

- d) AC → BC

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

(ii) Determining equivalence of two FD sets:

Given two sets of FDs, say, F & G, it is to be determined if

F & G are said to be equivalent iff $F^+ = G^+$

→ But finding equivalence in this process is tedious.

→ So we find if F covers G and G covers F.

If F covers G and G covers F, then F & G are said to be equivalent.

→ we say F cover G , if G is contained in F (ii)

if all the dependencies of G can be derived from F .

$$\therefore \{(F \text{ cover } G) \wedge (G \text{ covers } F)\} \Leftrightarrow F \cong G$$

(iii)

Eg: Consider all the dependencies of G to be same as F

$$F: \{A \rightarrow B, AB \rightarrow C, D \rightarrow AC\}$$

conditions

$$G: \{A \rightarrow BC, D \rightarrow A\}$$

(i) ~~if $A \rightarrow B$ is to be found, going to $D \rightarrow A$ is better~~
and ~~if $AB \rightarrow C$ is to be found, going to $D \rightarrow A$ is better~~
~~if $D \rightarrow A$ is to be found, going to $AB \rightarrow C$ is better~~

Checking F covers G ~~if G is contained in F~~

start with $A \rightarrow B$ to be found, going to $D \rightarrow A$ is better.



(i) $A \rightarrow BC$ to be found, and the dependencies of G

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

(ii) $D \rightarrow A$ to be found, and the dependencies of G

(ii) $D \rightarrow A$ to be found, and the dependencies of G

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

(iii) $D \rightarrow A$ to be found, and the dependencies of G

$\therefore F$ covers G .

Checking G covers F ~~if F is contained in G~~

(i) $A \rightarrow B$

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

(ii) $AB \rightarrow C$

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

(iii) $D \rightarrow A$

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

$\therefore G$ covers $F \quad \therefore F \cong G$

Q4) F: $\{A \rightarrow C, AC \rightarrow D, E \rightarrow AD, E \rightarrow H\}$

G: $\{A \rightarrow CD, E \rightarrow AH\}$

Equivalent

(iii) Finding minimal set of FDs of canonical form:

A set of FDs F is minimal if it satisfies the below conditions:

- Every dependency in F has a single attribute on its right side.
- We cannot replace any $X \rightarrow A$ in F with $Y \rightarrow A$ where Y is a proper subset of X and still have equivalent set to F.
- No dependency can be removed such that the set is equivalent to F.

Procedure to find canonical cover:

- Minimize LHS of each FD:

Consider $AB \rightarrow C$

A' can be deleted if B^+ contains A

Eg: F: $\{AB \rightarrow C, B \rightarrow A\}$

Canonical cover, G: $\{B \rightarrow C, B \rightarrow A\}$

A set may even have more than one minimal sets possible

for domain of

Eg: Consider

$$F: \{AB \rightarrow C, B \rightarrow A, A \rightarrow B\}$$

A^+ contains B

B^+ contains A

$$G_1: \{A \rightarrow C, B \rightarrow A, A \rightarrow B\}$$

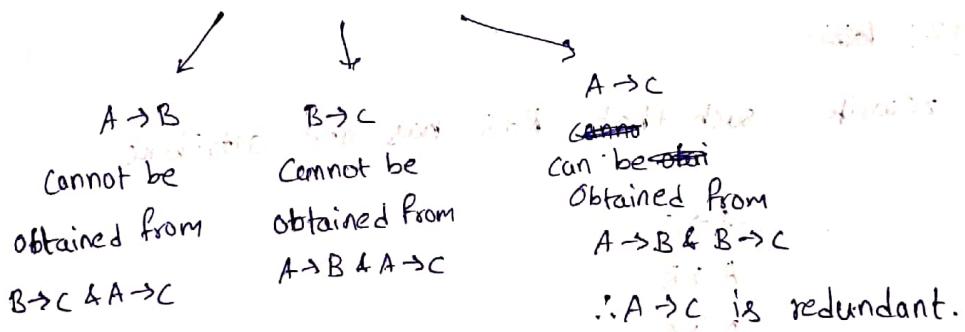
$$G_2: \{B \rightarrow C, B \rightarrow A, A \rightarrow B\}$$

$$\text{Here } F \cong G_1 \cong G_2$$

(iii) Minimize RHS of each FD:

First make sure that right hand side has only one attribute

Let $F: \{A \rightarrow B, B \rightarrow C, A \rightarrow C\}$ (using decomposition)



Thus F can be reduced to

$$\{A \rightarrow B, B \rightarrow C\}$$

→ Thus we find canonical cover performing above two steps.

Q5 F: $\{AB \rightarrow C, A \rightarrow B, D \rightarrow AC\}$

Find minimal set

Sol:

Step 1: Minimize LHS

$$(i) AB \rightarrow C$$

$$A^+ = \{A, B, C\} \Rightarrow B \text{ can be deleted}$$

$$B^+ = \{B\} \Rightarrow A \text{ can't be deleted}$$

A^+ contains B

So remove B

$$A \rightarrow C$$

Step 2: ~~Minimize R.H.S.~~

rewrite such that RHS has one attribute

$$A \rightarrow C$$

$$A \rightarrow B$$

$$D \rightarrow A$$

$$D \rightarrow C$$

• Check if $A \rightarrow C$ is derivable from other dependencies

i.e., ~~for~~

Compute A^+ from $A \rightarrow B, D \rightarrow A, D \rightarrow C$

$$A^+ = \{A, B\}$$

$\therefore A \rightarrow C$ can't be removed

• $A \rightarrow B$

$$A^+ = \{A, C\}$$

$A \rightarrow B$ can't be removed

• $D \rightarrow A$

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

$\therefore D \rightarrow A$ can't be removed

• $D \rightarrow C$

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

$\therefore D \rightarrow C$ can be removed

Final \therefore Final minimized set is

$$F: \{A \rightarrow C, A \rightarrow B, D \rightarrow A\}$$

$\therefore F: \{A \rightarrow BC, D \rightarrow A\}$ is

minimal set / canonical set / irreducible set.

(Q2)

$F: \{AB \rightarrow C, A \rightarrow B, B \rightarrow C\}$. Find minimal set

Soln:

Step 1:

$$\cancel{A^+} \quad A^+ = \{A, B, C\} \quad \{A \rightarrow A, B \rightarrow A, C \rightarrow A\}$$

$$B^+ = \{B, C\}$$

$$\{B \rightarrow B, C \rightarrow C, B \rightarrow C, C \rightarrow B\}$$

$$\therefore F: \{A \rightarrow C, A \rightarrow B, B \rightarrow C\} \quad \{A \rightarrow C, A \rightarrow B, B \rightarrow C\}$$

Step 2:

$A \rightarrow C$ can be obtained from $A \rightarrow B$ & $B \rightarrow C$

$$\therefore F: \{A \rightarrow B, B \rightarrow C\}$$

P/15

$$V \rightarrow W$$

$$VW \rightarrow X$$

$$Y \rightarrow VX$$

$$Y \rightarrow Z$$

$$\{V, W, X\} \text{ is}$$

$$\{V, W, X\}$$

Step (i)
 $VW \rightarrow X$

$V^+ = \underline{\{W\}}$ V^+ contains w
 $w^+ = \{w\}$

\therefore remove w

$V \rightarrow W; V \rightarrow X; Y \rightarrow Vx; Y \rightarrow Z$ $\{V, x, Y, Z\} \rightarrow$

Step (ii)

$V \rightarrow W$

$V \rightarrow X$

$Y \rightarrow V$

$Y \rightarrow X$

$Y \rightarrow Z$

$Y \rightarrow X$ can be removed $\{V, x, Y, Z\} \rightarrow$

$\therefore Y \rightarrow V \neq V \rightarrow X$

\therefore still (some) (no) (no) to be eliminated

\therefore reduced set is

$V \rightarrow W$

$A \rightarrow X$ minimum bnf. $\{A \rightarrow B, A \rightarrow C, A \rightarrow BC, B \rightarrow C\} \rightarrow$

$Y \rightarrow V$

$Y \rightarrow Z$

(P/11)

$\{A \rightarrow B, AB \rightarrow C, A \rightarrow BC, B \rightarrow C\} \rightarrow$ $\{A\} \rightarrow A$ \rightarrow

$\{A \rightarrow B, A \rightarrow C, A \rightarrow BC, B \rightarrow C\} \rightarrow$

$\{A \rightarrow B, A \rightarrow C, B \rightarrow C\} \rightarrow$

$\{A \rightarrow B, AB \rightarrow C\}$

\rightarrow $\{A \rightarrow B, A \rightarrow C, B \rightarrow C\}$ \rightarrow $\{A\} \rightarrow A$

(P/12)

$\{A \rightarrow B, CD \rightarrow A, CB \rightarrow D, CE \rightarrow D, AE \rightarrow F, AC \rightarrow D\}$

$C^t = \{A / C\}, D^t$

$D^t = \{D\}$

Step (i): $\{A \rightarrow B, CD \rightarrow A, CB \rightarrow D, CE \rightarrow D, AE \rightarrow F, AC \rightarrow D\}$

Step 1: ~~For 3rd step minimization and removal of dependencies~~
no minimization can be done in step(i)
~~because all dependencies are~~

Step 2:



$AC \rightarrow D$ can be removed

$\therefore (AC)^+ = \{A, B, C, D\}$ (from remaining 5 dependencies)

$\therefore 5$

14/11/20

(P/13)

$BCD \rightarrow A$: 0 no 3 (Hence 3 dependencies left)

$C \rightarrow E$

$E \rightarrow B$

minimization on left size (left)

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

self-loop on C, B, A, D, no minimization on left size

$\therefore B$ can be eliminated from $BCD \rightarrow A$

$CD \rightarrow A$

$C \rightarrow E$

$E \rightarrow B$

minimization on right size

(P/14)

$BC \rightarrow A$

$EF \rightarrow A$

$B^+ = \{B, C, A\}$: 2 (left) 3 (right)

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

$B \rightarrow D$

$C^+ = \{C, F\}$ (no minimization on left)

$D \rightarrow EF$

$E^+ = \{E\}$ (no minimization on right)

$C \rightarrow F$

$F^+ = \{F\}$ (no minimization on right)

\therefore no minimization on left size

$BC \rightarrow A$:

$(BC)^+ = \{B, C, D, E, F\}$: 5
 $\therefore BC \rightarrow A$ can be removed

P/16

$A^+ = \{A, B, C, G, D\}$ B can be removed from $ABE \rightarrow CDGH$
 \therefore no minimization left side.

$A \rightarrow B$ — essential

$A \rightarrow C$ — essential

$A \# E \rightarrow C$ — not needed

$A \# E \rightarrow D$ — not needed

$A \# E \rightarrow G$ — not needed

$A \# E \rightarrow H$ — needed

$C \rightarrow G$ — not needed

$C \rightarrow D$ — needed

$D \rightarrow G$ — needed

$E \rightarrow F$ — needed

$\therefore \{A \rightarrow B; A \rightarrow C; A \# E \rightarrow H; C \rightarrow D; D \rightarrow G; E \rightarrow F\}$

(iv) Finding the keys of a relation:

If X is set of attributes and if X^+ contains all the attributes then X is called super key.

If X is minimal then we call it candidate key.

Ex: Consider

$R(ABC)$ $F: \{A \rightarrow B, B \rightarrow C\}$

Find closure of LHS of all dependencies

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

$$B^+ = \{B, C\}$$

now 'A' is ck

CK: A

SK: A, AB, AC, ABC

$\therefore 4$ super keys

Eg: Find CKs and SKs for relation R with F

$R(ABCD)$ i.e. $F: \{AB \rightarrow C, A \rightarrow D\}$ minimal pair wise fd

$$(AB)^+ = \{A, B, C, D\}$$

$$A^+ = \{A, D\} \quad \text{since } A \rightarrow D$$

$\therefore AB$ is ek

$$\underline{CK}: AB$$

$$\underline{SK}: AB, ABC, ABD, ABCD$$

Eg:

$R(ABCD)$

$F: \{A \rightarrow B, AB \rightarrow C, B \rightarrow D\}$

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

$$(AB)^+ = \{A, B, C, D\}$$

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

$$\underline{CK}: A$$

$$\underline{SK}: \{A\}, AB, AC, AD, ABC, ABD, ACD, ABCD$$

(8 Combinations of B, C, D are possible)

Eg: $R(ABCD)$

$$F: \{AB \rightarrow C, BC \rightarrow D\}$$

$$(AB)^+ = \{A, B, C, D\}$$

$$(BC)^+ = \{B, C, D\}$$

$$\underline{CK}: AB, BC$$



$$\therefore \text{no of SKs} = 4 + 4 - 2 = 6$$

In these kind of situation we compute no of sks by combining (union) two (cks) and find no of possibilities and subtract it from

$AB \rightarrow 4$ possibilities $BC \rightarrow 4$ possibility

$ABC \rightarrow 2$ possibility

$$\text{Ans} : 4+4-2 = \underline{\underline{6}}$$

Ex: $R(ABC)$ F: $\{A \rightarrow BC, BC \rightarrow A\}$

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

$(BC)^+ = \{A, B, C\}$

CK: A, BC

Sks:

$A \rightarrow 4$ ways $A \rightarrow BC$

$BC \rightarrow 2$ ways $(-) ABC \rightarrow 1$ way

$\therefore (4+2)-1 = 5$

Ex: $R(ABC)$ F: $\{A \rightarrow B, B \rightarrow C, C \rightarrow A\}$

CK: A, B, C

\Rightarrow : Since all attributes are key every combination of key attributes is sks.

\therefore no of sks = ~~2^3~~ $2^3 - 1 = 7$

Eg: R(ABCD)

F: {A → B, C → D}

Find CKs. (for 2nd normal form or 3NF)

$$A^+ = \{A, B\} \quad \text{↳ A determines B}$$

$$C^+ = \{C, D\} \quad \text{↳ C determines D}$$

$$(AC)^+ = A^+ \cup C^+ = \{A, B, C, D\} \subseteq S \Rightarrow AC \text{ is CK}$$

∴ AC is a candidate key

Eg: R(ABCD)

F: {A → BC, B → C}

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

$$B^+ = \{B, C\}$$

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

∴ since A determines BC and B determines C
∴ AD is CK

Eg: R(ABC)

F: {AB → C, C → A}

$$(AB)^+ = \{A, B, C\} \quad \text{↳ AB determines C}$$

$$C^+ = \{A, C\} \quad \text{↳ C determines A}$$

Here B is independent attribute

Since i.e., no attribute determines B

∴ CK: AB, BC

∴ CK: AB, BC

* Eg: $R(ABCD)$

* ~~For $AB \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$, $C \rightarrow D$~~

$$F: \{AB \rightarrow C, A \rightarrow D, D \rightarrow A, C \rightarrow D\}$$

Here B is independent attribute

$$(AB)^+ = \{A, B, C, D\}$$

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

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

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

AB is one ck.

A, B are prime attribute.

To obtain other cks check if there is any dependency such that prime attribute on RHS. If any is found, then replace it with its LHS in ck.

$$CK \rightarrow AB$$

In $D \rightarrow A$, A is present on RHS

$\therefore DB$ is also ck

In $C \rightarrow D$, D is present on RHS

$\therefore CB$ is ck

\therefore cks: AB, BC, BD

Finding cks:

$$AB \rightarrow 4 \quad BC \rightarrow 4 \quad BD \rightarrow 4$$

$$ABC \rightarrow 2 \quad ABD \rightarrow 2 \quad BCD \rightarrow 2$$

$$\therefore (4+4+4) - (2+2+2) + 1 = 12 - 6 + 1 = 7$$

Eg: R(ABCDE)

F: $\{AB \rightarrow C, BD \rightarrow E, DE \rightarrow B\}$

$(AB)^+ = \{A, B, C\}$ (A, B, C are prime attributes)

$(BD)^+ = \{B, D, E\}$ (B, D, E are prime attributes)

$(DE)^+ = \{B, D, E\}$

A, D are independent attributes.

$\Rightarrow \underline{(ABD)^+} \quad (AD)^+ = \{A, D\}$

$(ABD)^+ = \{A, B, C, D, E\}$

$(ADE)^+ = \{A, B, C, D, E\}$

Cks: ABD, ADE

$\therefore A, B, D, E$ are prime attributes.

∴ A & D are independent attributes so it cost

∴ no of Cks = 2

Eg: R(ABCDE)

F: $\{A \rightarrow B, BC \rightarrow D, DE \rightarrow A\}$

$A^+ = \{A, B\}$ (A, B are prime attributes)

$(BC)^+ = \{B, C, D\}$ (B, C, D are prime attributes)

$(DE)^+ = \{D, E, A, B\}$ (D, E, A, B are prime attributes)

Here C, E are independent attributes.

$(CE)^+ = \{C, E\}$

* Independent attributes must exist in ck.

So we try adding different combinations to CE

$$(ACE)^+ = \{A, B, C, D, E\}$$

$$(BCE)^+ = \{B, C, E, D, A\}$$

$$(DCE)^+ = \{D, C, E, A, B\}$$

\therefore CKs: ACE, BCE, DCE

$\{A, B, C\} \in F(O)$

P/24

R(ABCDEFH)

$$(CH)^+ = \{C, H, G, I\}$$

$$A^+ = \{A, B, C, F, H, E, G\}$$

$$B^+ = \{B, C, F, H, E, G, A\}$$

$$E^+ = \{E, A, B, C, F, H, G\}$$

$$F^+ = \{F, E, G, A, B, C, H\}$$

Here D is independent attribute & $D^+ = \{D\}$

AD, BD, ED, FD are CKs

\therefore 4 CKs

(Ans)(c) 4

P/19

$$C^+ \subseteq \{C, F\}$$

Here C, F are independent attributes.

CF must be present in key of R(O)

\therefore opt (a) {C, A, B, D} from

P/23

E, H are independent attributes and must be part of key

$$(EH)^+ = \{E, H, C\}$$

from options

opt (d) {C, E, H, D} is correct

P/17

$$AB \rightarrow 2^4 = 16 \quad AE \rightarrow 2^4 = 16$$

$$ABE \rightarrow 2^3 = 8$$

$$\therefore 16 + 16 - 8 = 24$$

P/18

$$BC \rightarrow 8$$

$$CD \rightarrow 8$$

$$BCD \rightarrow 4$$

$$\therefore 8 + 8 - 4 = 12 \text{ Sks}$$

P/18

$$AB \rightarrow 8$$

$$BC \rightarrow 8$$

$$CD \rightarrow 8$$

$$ABC \rightarrow 4$$

$$BCD \rightarrow 4$$

$$ABCD \rightarrow 2$$

$$ABCD \rightarrow 2$$

$$(8+8+8) - (4+4+2) + 2 = 16$$

P/20

A, B, C, D, k are independent attributes.

\therefore opt (d)

P/21

A, B are independent

~~opt (a)~~ b) $(AB)^T = \{A, B, C\}$

d) $(ABD)^T = \{A, B, C, D, E\}$

\therefore opt (d)

P/22

H is independent

\therefore opt (b)

P/25

$$(AB)^+ = \{A, B, C, D, E, F\}$$

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

$$C^+ = \{C, A\}$$

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

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

AB, F are CKs

A, B, F are prime attributes

In $C \rightarrow A$,

C is on RHS

\therefore In AB , we can replace A with C

$\therefore CB$ is CK

A, B, C, F are prime attributes

$\therefore 3$ CKs (AB, CB, F)

P/26

$$A^+ = \{A, C, E\}$$

$$(AB)^+ = \{A, B, C, D, E\}$$

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

$$C^+ = \{C, E, A\}$$

$$E^+ = \{E, A, C\}$$

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

AB is CK

A, B are prime attributes

AD is CK

BE is CK

now A, B, D, E are prime attributes

∴ BC is key

∴ BA is key

AB, AD, BC, BE

replace A with E

EB, ED, AB, AD, BC

∴ BA is key

∴ BE is key

P/27

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

$$(CD)^+ = \{C, D, E, A, H, B\}$$

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

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

$$(ABH)^+ = \{A, B, H, D, C, E\}$$

$$(DH)^+ = \{B, C, D, H, A, E\}$$

if next, A is part of the primary key then

∴ CK → D, ABH

$$A \rightarrow B$$

With both forms of consideration, ABH is not CK

AH is CK

After this, forming $B \oplus D, AH$ is not CK

∴ 2

equivalent positions in

$$\{BC, A\} \oplus \{B, C\}$$

Properties of Decomposition:

(i) Lossless join decomposition

(ii) dependency preserving decomposition

i) Lossless join decomposition:

The decomposition of R into R_1 & R_2 is said to be lossless iff ~~the~~ the join of R_1 & R_2 gives R .

otherwise it is called lossy decomposition.

→ When R is divided in R_1 & R_2 , for the decomposition to be lossless there must be a common attribute between R_1 & R_2 .

The common attribute ^{must} be either key of R_1 or key of R_2 .

This common attribute if it is key of R_1 , then it acts as Fk in R_2 .

Now we say parent child relationship is established b/w R_1 and R_2 .

The decomposition is lossless iff parent child relationship is maintained properly.

Eg: Consider

$$R(ABC) \quad F: \{A \rightarrow B\}$$

is decompose into $R_1(AB)$ $R_2(BC)$.

Check if this decomposition is lossy or lossless.

$$R_1 \cap R_2 = B$$

The common attribute has to be key of R_1 or key of R_2 .

~~Note But it is impossible to determine which~~

But B is key for neither of the tables.

∴ It is lossy decomposition.

$$\text{Eg: } R(ABC) \quad F: \{A \rightarrow B\}$$

$$R_1(AB) \quad R_2(AC)$$

$$R_1 \cap R_2 = \{A\}$$

$$A^+ = \{A, B\}$$

A determines attributes of R_1 if suppose A

∴ R_1 is parent

R_2 is child

∴ lossless decomposition.

$$\text{Eg: } R(ABCDE) \quad F: \{AB \rightarrow C, B \rightarrow D, D \rightarrow E\}$$

Let R be decomposed into $R_1(ABC)$ $R_2(BD)$ $R_3(DE)$

Consider joining R_1 & R_2

$$R_1 \cap R_2 = B$$

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

∴ R_1 & R_2 a parent child relationship exists

∴ lossless join

$$R_1 R_2(ABCD) \quad R_3(DE)$$

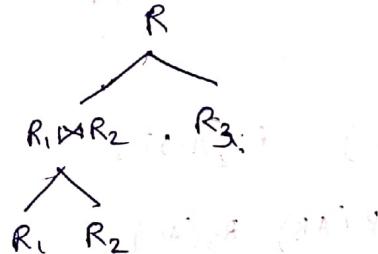
$$R_1 R_2 \cap R_3 = D$$

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

D can determine all attributes of R_3

$\therefore R \xrightarrow{R_1} R_2 \xrightarrow{R_3}$ is lossless decomposition

The decomposition is



dis Dependency Preserving decomposition:

A decomposition is said to be dependency preserving

$$\text{iff } F^+ = (F_1 \cup F_2)^+$$

F is dependencies of R

F_1 & F_2 are dependencies of R_1 & R_2

$$\text{Ex: } R(ABC) \quad F: \{A \rightarrow B, B \rightarrow C, C \rightarrow A\}$$

is decomposition $R_1(AB)$ $R_2(BC)$ dependency preserving

$$R_1(AB) - F_1: \{A \rightarrow B\}$$

$$R_2(BC) - F_2: \{B \rightarrow C\}$$

~~C in F~~ can't be obtained from ~~F1 & F2~~

\therefore the decomposition is not dependency preserving.

$$B^+ = \{B, C, A\} \quad C^+ = \{C, A, B\}$$

\therefore we can add $B \rightarrow A$ to F_1 and

$C \rightarrow B$ to F_2

$F_1: A \rightarrow B$

$B \rightarrow A$

$F_2: B \rightarrow C$ (new paragraph)

$C \rightarrow B$

Now we see if $C \rightarrow A$ can be obtained from $F_1 \& F_2$

Find C^+ in $F_1 \& F_2$

$\therefore C^+ = \{e, B, A\}$

$\Rightarrow C \rightarrow A$

\therefore The decomposition is dependency preserving

Procedure:

i) Write all direct dependencies

ii) Now check for obtain indirect dependencies in F .

we do this by writing all indirect dependencies to

F_1, F_2, F_3, \dots

Now we see if given indirect dependency is possible or not.

Note:

\rightarrow A decomposition is said to be good if it is both
lossless and dependency preserving.

(P/01)

a) $R_1(BC)$ $R_2(AD)$

no common attribute

as it is lossless & dependency preserving

$$\frac{R_1(BC)}{F_1: B \rightarrow C} \quad \frac{R_2(AD)}{F_2: D \rightarrow A}$$

\therefore Dependency preserving

but it is not a good decomposition

3) lossless

$$R_1(ACD) \quad R_2(BC)$$

$$F_1: C \rightarrow A \quad F_2:$$

$$A^t = \{A\}$$

$$C^t = \{A, C, D\}$$

$$D^t = \{D\}$$

$$B^t = \{B\}$$

$$C^t = \{C\}$$

we need to check for $AB \rightarrow C$

$$(AB)^t = \{A, B\}$$

so no disjoint in projection from R to B, so not

\therefore not dependency preserving

3) lossless decomposition

$$R_1(ABC) \quad R_2(AD)$$

$$F_1: A \rightarrow BC$$

$$A^t = \{A, B, C, D\}$$

$$C^t = \{C, A, D\}$$

$$C \rightarrow AB$$

$$A^t = \{A, B, C, D\}$$

$$D^t = \{D\}$$

$$\therefore A \rightarrow D$$

$$F_1 \cup F_2 = \{A \rightarrow BC, C \rightarrow AB, A \rightarrow D\}$$

now we need check if $C \rightarrow D$ is possible from $F_1 \cup F_2$

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

\therefore ~~dependency~~ dependency preserving

\therefore this is good decomposition

4) lossy decomposition

$$\underline{R_1(AB)}$$

$$F_1: A \rightarrow B$$

$$\underline{R_2(AD)}$$

$$F_2:$$

$$\underline{R_3(CD)}$$

$$F_3: C \rightarrow D$$

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

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

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

$$\therefore A \rightarrow D$$

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

$$\therefore F_1 \cup F_2 \cup F_3 = \{A \rightarrow B, A \rightarrow D, C \rightarrow D\}$$

$$B^+ = \{B\}$$

\therefore not dependency preserving

(P2)

$$\underline{R_1(AB)}$$

$$A \rightarrow B$$

$$\underline{R_2(BC)}$$

$$B \rightarrow C$$

$$\underline{R_3(BD)}$$

$$C \rightarrow B$$

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

\therefore dependency preserving

~~dependency preserving~~

~~can't be obtained~~
~~not dep. preserving~~

$R_1 \bowtie R_2$

B is common

$$B^t = \{C, D, B\}$$

$\therefore R_1$ is child R_2 is parent

~~R~~ $(R_1 \bowtie R_2) \bowtie R_3$

B is common

$$B^t = \{C, D, B\}$$

R_3 is parent

\therefore lossless join

\therefore opt (a)

15/11/20

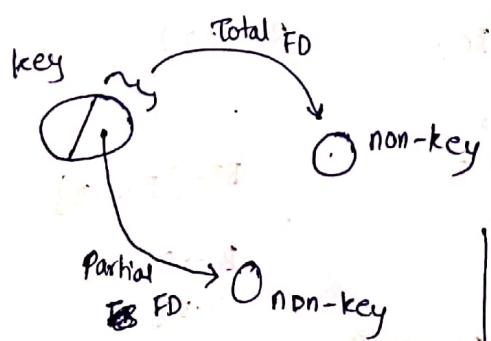
Normal forms:

1NF:

- 1NF doesn't allow multivalued and composite attributes.
- 1NF allows only atomic values
- Every relation in relational model is by default in 1NF.

2NF:

- 2NF doesn't allow partial FDs. but allows full FD.



| |
|---|
| PFD: key |
| part of key \rightarrow non-key attribute |

Ex: If A is in R(ABC)

($A \rightarrow C$) is called FD

($A \rightarrow A$) is called full FD

Ex: Consider, R(ABCD) and AB is CK

$AB \rightarrow C$ is called full FD
 $A \rightarrow D$ is called partial FD

→ Partial dependencies lead to redundancy

| Patient.No | Drug.No | No.of.Units | Patient.Name | Drug.Name |
|------------|---------|-------------|--------------|-----------|
| P_1 | d_1 | 10 | Ajith | Crocin |
| P_1 | d_2 | 15 | Ajith | dolo |
| P_2 | d_2 | 20 | Rama | dolo |

(The redundancies are circled).

FDS:

$\{Patient.No, Drug.No\} \rightarrow no.of.units$ (FFD)

$Patient.No \rightarrow Patient.Name$ (PFD)

CK

$Drug.No \rightarrow Drug.Name$ (PFD)

Now we decompose the relation to eliminate redundancy as

shown below

| Patient.No (FK) | Drug.No (FK) | No.of.Units |
|-----------------|--------------|-------------|
| P_1 | d_1 | 10 |
| P_1 | d_2 | 15 |
| P_2 | d_2 | 20 |

| Patient.No | Patient.Name |
|------------|--------------|
| P_1 | Ajith |
| P_2 | Rama |

Patient.No \rightarrow Patient.Name

| Drug.No | Drug.Name |
|---------|-----------|
| d_1 | Crocin |
| d_2 | dolo |

Drug.No \rightarrow Drug.Name

$\{Patient.No, Drug.No\} \rightarrow no.of.units$

→ Generally in $R(ABCD)$

If AB is CK and $F: \{AB \rightarrow C, A \rightarrow D\}$,

for every PFD we create a separate table.

for all FFD we put them in the same table

Decomposition is:

$R_1(ABC)$, $R_2(AD)$

Every decomposition of
this kind is lossless
and dependency preserving

Eg: Consider $R(ABCDE)$

$F: \{AB \rightarrow C, A \rightarrow D, B \rightarrow E\}$

CK: AB

Now $A \rightarrow D$, $B \rightarrow E$ are PFDS

R is in INF

now we decompose it is

PK
 $R_1(ABC)$ $R_2(AD)$ $R_3(BE)$

(0.77) how to do it \rightarrow first pass, one factoring

(0.77) smart thinking \rightarrow defining

Eg: $R(ABCD)$

$F: \{A \rightarrow B, C \rightarrow D\}$

CK: AC

decomposition is

PK
 $R_1(AC)$

$R_2(AB)$

$R_3(CD)$

wanted result

* Here if we decompose into two tables only $R_1(AB)$, $R_2(CD)$
then the decomposition will be lossy.

* Eg: $R(ABCDE)$

* $F: \{AB \rightarrow C, A \rightarrow D, D \rightarrow E\}$

$$(AB)^F = \{A, B, C, D, E\}$$

$A \rightarrow D$ is a PFD

\therefore Decomposition is

$$R_1(A \underset{\uparrow}{D} E) \quad R_2(A \underset{\uparrow}{B} C)$$

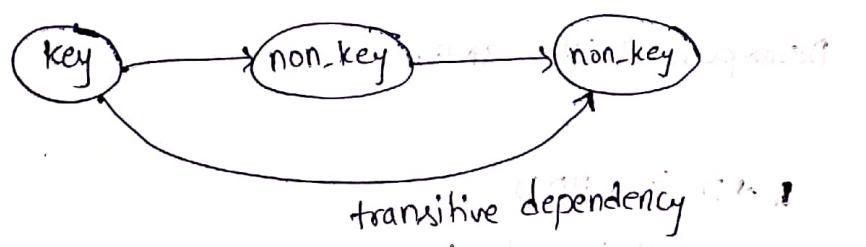
Note:

A relation with key of one attribute is always in 2NF.

3NF:

\rightarrow 3NF should first be in 2NF

\rightarrow 3NF not allows transitive dependencies.



Transitive dependency exists whenever there is
non-key \rightarrow non-key dependency. (or) $\left\{ \begin{array}{l} \text{Part of key} \\ \text{non-key} \end{array} \right\} \rightarrow \text{Non-key}$

Eg: Employee(Eno, Ename, seniority, salary)

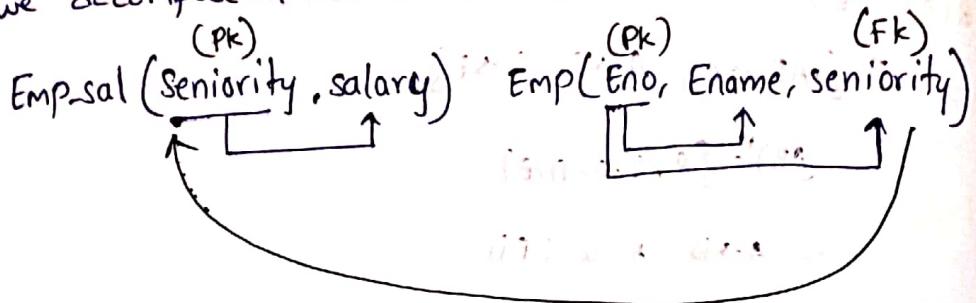
$Eno \rightarrow Ename, seniority$
 $seniority \rightarrow salary$

CK: Eno

It has no PFDs.

\therefore It is in 2NF but not in 3NF.

so we decompose the relation as



This decomposition is also lossless always.

Eg : $R(ABC)$

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

and R is in 2NF but not in 3NF

\therefore It is decomposed as

$$R_1(AB) \quad R_2(BC)$$

Eg : $R(ABCDE)$

$$F = \{AB \rightarrow C, B \rightarrow D, D \rightarrow E\}$$

Decompose R into 3NF.

Sol :

$B \rightarrow D$ is PFD

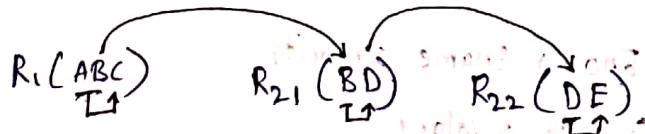
$\therefore R$ is not in 2NF

Decompose into 2NF

$$R_1(ABC) \quad R_2(BDE)$$

R_1 is in 3NF and R_2 is in 2NF but not in 3NF

\therefore Decomposition is



* \rightarrow Most database use 3NF and it is considered adequate.

Algorithm to decompose into 3NF: (from 1NF to 3NF)

This decomposition is lossless and dependency preserving.

Step 1: Find minimal set of F called G.

Step 2: Take each FD into a separate relation

This step ensures that the decomposition is dependency preserving

Step 3: If no decomposed relation contains key then take key attributes into separate table (to ensure lossless decomposition)

Step 4: If there are any redundant relations, then eliminate them.

$$F: R_1(AB) \quad R_2(ABC)$$

Here R_1 is redundant because attribute A, B are already present in R_2 .

Ex: $R(ABCD)$

$$F: \{AB \rightarrow C, B \rightarrow CD\}$$

Step 1: Find minimal

ck: AB

Step 2: Create individual relations for each FD.

$$\therefore R_1(ABC) \quad R_2(BD)$$

Step 3: R_1 has key

Step 4: No redundant relations

$\therefore R_1(ABC) \quad R_2(BD)$ is in 3NF

Eg: $R(ABCD)$

F: $\{A \rightarrow B, C \rightarrow D\}$

Ck: AC

Step 1: F is minimal

∴

Step 2: $R_1(AB) R_2(CD)$

Step 3: There is no relation with key attributes.

So we create one

and then we get $R_3(AC)$, because on R_1 there

Step 4: $R_1(AB) R_2(CD) R_3(AC)$

and obtaining no redundant relations, so now it is in 3NF

∴ 3NF

Eg: $R(ABCDE)$

F: $\{AB \rightarrow C, B \rightarrow D, D \rightarrow E\}$

Ck: AB

Step 1:

F is minimal

Step 2:

$R_1(ABC) R_2(BD) R_3(DE)$

Step 3: R_1 has key attribute

Step 4: no redundant relations

∴ $R_1(ABC) R_2(BD) R_3(DE)$ is in 3NF.

is in 3NF.

P/3

Ex. $(AB)^+ = \{A, B, C, D, E, F, G, H, I, S\}$

AB is key

Decomposing into 3NF: $R_1(A), R_2(B), R_3(C), R_4(D), R_5(E), R_6(F), R_7(G), R_8(H), R_9(I), R_{10}(S)$

$R_1(ABC)$ $R_2(ADE)$ $R_3(BF)$ $R_4(FGH)$ $R_5(DIJ)$

R_1 has key attributes

Also there are no redundant relations

\therefore It is in 3NF

2NF decomposition:

$R_1(ABC)$ $R_2(ADEIJ)$ $R_3(BFGH)$

(P/4)

Step(i):

F is minimal key of R as it is not redundant in R .

Step(ii):

~~$R_1(ACF)$ $R_2(AB)$ $R_3(BE)$ $R_4(CD)$~~

CK: $(AC)^+ = \{A, B, C, D, E, F\}$

Step(iii):

key attributes present in R ,
no redundant attributes, so no anomalies in R .

Step(iv): no redundant attributes

$\therefore R_1(ACF)$ $R_2(AB)$ $R_3(BE)$ $R_4(CD)$

(P/5)

$A^+ = \{A, F, C, D\}$

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

$\therefore CK: AB$

$R_1(AFC)$ $R_2(CD)$ $R_3(BE)$

no relation has all key attributes

$\therefore R_1(AFC)$ $R_2(CD)$ $R_3(BE)$ $R_4(AB)$

Eg: $R(ABC)$

$F: \{AB \rightarrow C, C \rightarrow A\}$

$\overbrace{C}^{\text{prime transitivity}} \rightarrow$ Prime transitivity

$(AB)^t = \{A, B, C\}$

$C^t = \{A, C\}$

$\underline{CK}: AB, CB$

~~There is PFI~~

• There is no PFD

• There is no transitive dependency

\therefore This is 3NF

Note:

→ If a relation schema R is in 3NF then for every dependency $X \rightarrow A$ ~~one~~ at least one of below two conditions holds

- X is super key
- A is prime attribute.

\therefore If all attributes are prime attribute then the relation is in 3NF

Eg: $R(ABCD)$

$F: \{AB \rightarrow C, B \rightarrow D, D \rightarrow B, C \rightarrow A\}$

Find normal form of R

$$(AB)^t = \{A, B, C, D\}$$

$$\underline{CK}: AB, AD, CB, CD$$

all are prime attributes

\therefore 3NF

BCNF (Boyce-Codd Normal Form) :-

→ The problem with 3NF is prime transitivity.

Eg: Consider $R(ABC)$

$$F = \{AB \rightarrow C; C \rightarrow A\}$$

CK: AB, CB

R is in 3NF

Considering the information from FDs alone, BCNF has redundancy



Here prime attribute A is depending on CK AB

through another prime attribute c.

(Additional point) This is called prime-transitivity

This prime transitivity is not allowed in BCNF.

$AB \rightarrow C$ $C \rightarrow A$
 prime attribute A is depending on CK AB
 $\{c_1, a_1\}$ $\{c_1, a_1\}$ Redundancy due to prime
 Redundancy due to prime
 transitivity

→ A relation schema 'R' is in BCNF, if whenever a nontrivial FD of form $X \rightarrow A$ holds,

X must be superkey

If R is a relation and
 $X \rightarrow A$ is dependency
 violating BCNF property
 then we decompose into
 $R_1(R-A)$ $R_2(XA)$ and
 continue this process
 recursively

Eg: $R(ABC)$: $F: \{AB \rightarrow C; C \rightarrow AB\}$

CK: AB; C

In $AB \rightarrow C$, AB is superkey

$C \rightarrow AB$ - C is superkey

∴ R is in BCNF

Eg: R(ABC) \rightarrow F: $\{AB \rightarrow C, C \rightarrow A\}$

Ck: AB, CB

In $C \rightarrow A$, C is not superkey

R is not in BCNF but in 3NF

So we need to decompose

since $C \rightarrow A$ is violation

create R(AC)

to make sure it is lossless we

create another relation R(AB) or R(BC)

partial and missing (B is remaining attribute)

in $R_1(AC)$ (or) $R_1(AC)$ partial and missing

$R_2(AB)$

$R_2(BC)$

In Both these are not dependency preserving

Note:

All key relation is a relation formed by attributes of one candidate key

Every all key relation is always in BCNF

The relation $R_2(BC)$ is in BCNF

Sometimes lossless and dependency preserving decomposition

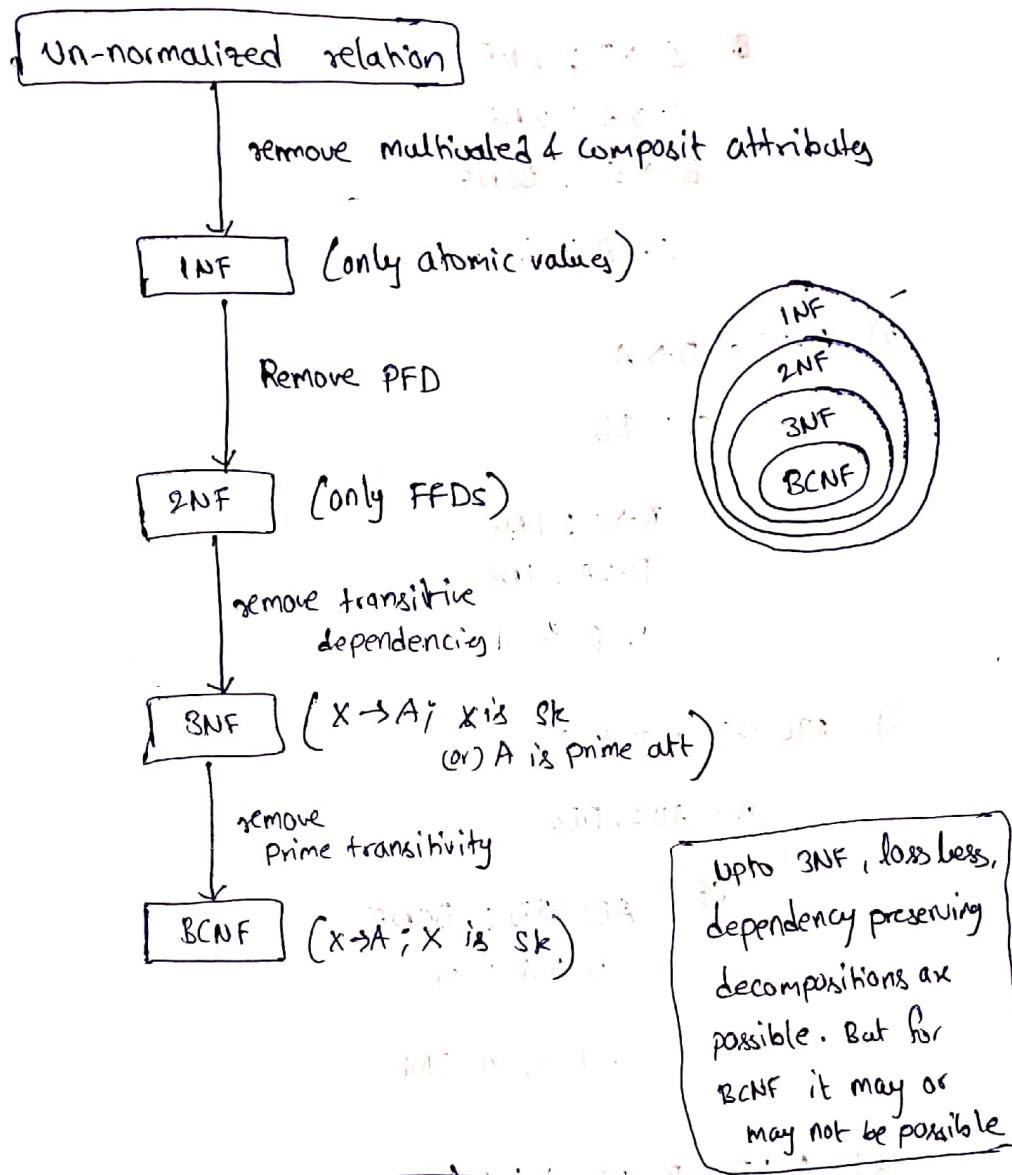
may not be possible for every relation.

Thus 3NF is mostly used.

Also for higher normal forms more no of relations are formed and thus accessing data may be slow.

→ Any relation with two attributes is in BCNF.

Note:



Note:

→ If there is (part of key → non-key) then it is not in 1NF

→ else if there is

Note:

* When asked to find normal form of a relation we start from backwards

i) If every dependency is of form

SK → any : BCNF

SK → PA : 3NF

Nk → Nk : 2NF

Part of key → Nk : 1NF

P/6

R(ABCD)

1) $C \rightarrow D$; $C \rightarrow A$; $B \rightarrow C$

ck: B

④ $C \rightarrow D$: 2NF

$C \rightarrow A$: 2NF

$B \rightarrow C$: BCNF

$\therefore R$ is in 2NF

2) $B \rightarrow C$; $D \rightarrow A$

ck: BD

$B \rightarrow C$: 1NF

$D \rightarrow A$: 1NF

$\therefore R$ is in 1NF

3) $ABC \rightarrow D$; $D \rightarrow A$

ck: ABC, DBC

④ $ABC \rightarrow D$: BCNF

$D \rightarrow A$: 3NF

$\therefore R$ is in 3NF

4) $A \rightarrow B$, $BC \rightarrow D$, $A \rightarrow C$

ck: A

$A \rightarrow B$: BCNF

$BC \rightarrow D$: 2NF

$A \rightarrow C$: BCNF

$\therefore R$ is in 2NF

5) $AB \rightarrow C$, $AB \rightarrow D$, $C \rightarrow A$, $D \rightarrow B$

ck: AB, CB, BAD, CD

all are prime attributes

$\therefore 3NF$

$C \rightarrow A$: 3NF

$\therefore R$ is in 3NF

(P7)

R.(ABC)

$$F_1: \{A \rightarrow B, B \rightarrow C, C \rightarrow AB\}$$

CK: A, C, B

$A \rightarrow B$: BCNF

$B \rightarrow C$: ~~BCNF~~

$\therefore R_1$ is in ~~BCNF~~

(P8)

$$F: \{QR \rightarrow S, R \rightarrow P, S \rightarrow Q\}$$

Y is $\{QR\}$. Lower dependencies are $S \rightarrow Q$

$$F_1: R \rightarrow P$$

$$P^+ = \{P\}$$

$$R^+ = \{R, P\}$$

$$F_2: QR \rightarrow S$$

$$S \rightarrow Q$$

$$Q^+ = \{Q\}$$

$$S^+ = \{S, Q\}$$

$$R^+ = \{R, P\}$$

Y is in BCNF

$$F_2: \{QR \rightarrow S, S \rightarrow Q, SR \rightarrow Q\}$$

CK: QR, SR

$S \rightarrow Q$: 3NF

$\therefore Z$ is not in BCNF

R is common attribute

$$R^+ \text{ in } F = \{R, P\}$$

Y(PR)

\therefore lossless decomposition

$$F_1 \cup F_2: \{R \rightarrow P, QR \rightarrow S, S \rightarrow Q, SR \rightarrow Q\}$$

\therefore dependency preserving

\therefore opt (a)

P/9

$\text{client_id} \rightarrow \text{f.name}, \text{l.name}, \text{city}, \text{zip_code}$

$\text{order_id} \rightarrow \text{order_data}, \text{city}, \text{zip_code}$

$\therefore \text{PFD is present}$ PFDs present

$\therefore \text{relation in INF}$

$\therefore \text{violating 2NF}$

→ when asked to find a normal form of relation we find highest normal form.

when asked to find violating normal form we find lowest normal form.

P/10

$\text{cid}, \text{sid} \rightarrow \text{bill}$

$\text{sid} \rightarrow \text{store_loc}$

↳ PFD

$\{ \text{bill} \} \rightarrow \text{bill}$

$\{ \text{store_loc} \} \rightarrow \text{store_loc}$

$\therefore \text{Take it to separate table and sid pk}$
and sid Fk in
original table

P/11

$\underline{\underline{F_1, F_2}} \not\subseteq F_1, F_2, F_3$

$F_1, F_2 \rightarrow \text{independent}$

$\therefore (F_1 F_2)^\dagger = \{ F_1, F_2, F_3, F_4 \}$

CK: $F_1 F_2$

writer: $F_1 \rightarrow F_3$: INF

$F_2 \rightarrow F_4$: INF

$F_1 F_2 \rightarrow F_5$: BCNF

$\therefore \text{opt(a)}$

17/11/20

(P/13) a) $B \rightarrow C$
 $AB \rightarrow C$
CK: B, AB

$\therefore \text{3NF}$

b) $AC \rightarrow B$
 $B \rightarrow C$

CK: AC, AB

all are prime

$\therefore 3NF$

c) $A \rightarrow B$

$B \rightarrow C$

CK: A

$A \rightarrow B : 3NF$

$B \rightarrow C : 2NF$

$\therefore 2NF$

d) $C \rightarrow AB$

$B \rightarrow C$

CK: C, B

Writing in above BCNF check off to each bit true

(P/15)

$AB \rightarrow D$

$BC \rightarrow D$

$A \rightarrow C$

$C \rightarrow A$

$(AB)^+ = \{ABC, D\}$

CK: AB, CB

$\therefore 3NF$

(P/16)

I. $BCNF$

II. $CK: \{rno, cid\}, \{email, cid\}$

$\therefore 3NF$ but not $BCNF$

Relational Algebra:

→ procedural language that give query evaluation plan.

1) Selection (σ) :

This operator is used for ~~section~~ selection of rows.

2) projection (π) :

This operator is used for selection of columns

Eg: Consider below student relation

| RNo | Name | marks |
|-----|-------|-------|
| 1 | Ajay | 15 |
| 2 | Shiva | 16 |
| 3 | Ajay | 13 |
| 4 | Ramu | 9 |

Query: find name of the students whose marks are greater than 10.

$\pi_{\text{Name}} (\sigma_{\text{marks} > 10} \text{student})$

Select name
from student
where marks > 10;

O/P: Name
Ajay
Shiva

(Here Ajay is printed only once)

Note:

- RA & RC has duplicate elimination implicitly
- But SQL gives duplicates also
- So we need to use distinct keyword in SQL.

Set manipulation operators:

U - Union

Π - Intersection

(-) - Difference

To perform any of these operations b/w two sets, the sets must be compatible.

- Two relational instances are said to be compatible iff
- (i) same no of columns and
 - (ii) corresponding attributes must have same domain

Ex: Consider

| Depositor | |
|-----------|-------|
| Cname | Acno. |
| Ajay | 1 |
| Balu | 2 |
| Charan | 3 |

| Borrower | |
|----------|---------|
| Cname | Loan no |
| Devi | 9 |
| Ajay | 10 |
| Sita | 11 |

Query: Find name of the customers who have an account (or) loan

(or) both

$$\{ (\Pi_{cname} \text{ Depositor}) \cup (\Pi_{cname} \text{ Borrower}) \}$$

Here we don't need to select any rows. So we don't need

to use Π.

Q1P: Cname

Ajay

Balu

Charan

Devi

Sita

Query: Find customer names who are depositors but not borrowers.

$$\left\{ \left(\pi_{Cname} \text{ Depositor} \right) - \left(\pi_{Cname} \text{ Borrower} \right) \right\}$$

Cross Product: (x)

It is useful when we need to process data from more than one tuple.

Consider

| R | | S | |
|---|---|---|---|
| A | B | C | D |
| 1 | 2 | 5 | 6 |
| 3 | 4 | 7 | 8 |



| R x S | | | |
|-------|---|---|---|
| A | B | C | D |
| 1 | 2 | 5 | 6 |
| 1 | 2 | 7 | 8 |
| 3 | 4 | 5 | 6 |
| 3 | 4 | 7 | 8 |

Eg: Borrower \sqcap Loan

Order of attribute must be ABCD
only in SR it is CDAB

| Cname | Lno |
|-------|-----|
| Ajay | 9 |
| Sita | 10 |

| Lno | Branch | Amount |
|-----|--------|--------|
| 9 | ABIDS | 10k |
| 10 | Sec | 9k |

Query: Find name of the customer who have a loan in the branch abids.

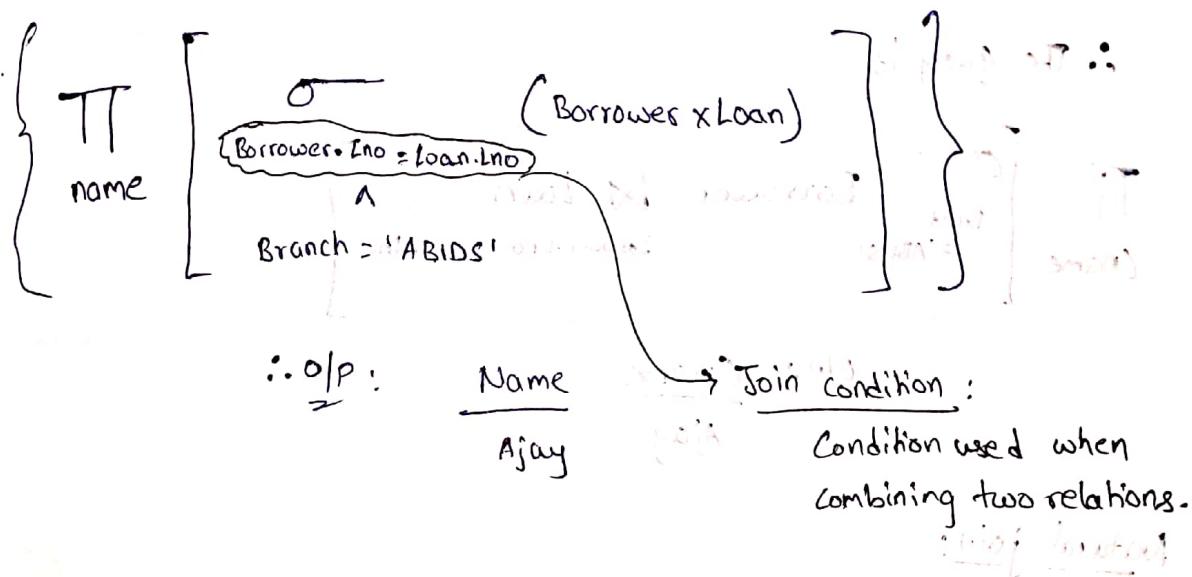
→ For queries involving 2 relations we perform Cartesian product.

Borrower x Loan

| Cname | Lno | Lno. | Branch | Amount |
|-------|-----|------|--------|--------|
| Ajay | 9 | 9 | ABIDS | 10k |
| Ajay | 9 | 10 | SEC | 9k |
| Sita | 10 | 9 | ABIDS | 10k |
| Sita | 10 | 10 | SEC | 9k |

Actually when a name conflict occur while performing cross product, no names will be given and we need to refer them by position

Here combinations are where Lno are equal



18/11/20
 Shortest path algorithm, brief overview of

JOIN: (\bowtie) predicate over relations remains to be explained

It is ~~do~~ a cartesian product followed by a selection and a projection

Consider

$$\frac{R}{A \ B} \quad , \quad \frac{S}{B \ C}$$

Here $R \bowtie S = \Pi_{A, B, C} (R.B = S.B \ R \times S)$

Consider the previous query of finding customer who have loan in ABIDS

Borrower \bowtie Loan

Borrower.Lno = Loan.Lno

| Cname | Lno | Branch | Amount |
|-------|-----|--------|--------|
| Ajay | 9 | ABIDS | 10k |
| Sata | 10 | Sec | 9k |

∴ The query is

$\Pi_{\text{Cname}} \left[\begin{array}{l} \sigma_{\text{Branch} = 'ABIDS'} \\ \text{Borrower} \bowtie \text{Loan} \\ \text{Borrower.Lno} = \text{Loan.Lno} \end{array} \right]$

Op: Cname

Ajay

Natural join:

In natural join, the join condition is by default equality on all common attributes (i.e., attribute with same name)

R
 $\frac{}{A \ B}$

S
 $\frac{}{B \ C}$

$$R \bowtie S = R \bowtie S = \Pi_{A, B, C} \left(\sigma_{R.B=S.B} (R \times S) \right)$$

↓
Natural join

Now the previous query can be written as

$\Pi_{\text{Cname}} \left[\begin{array}{l} \sigma_{\text{Branch} = 'ABIDS'} \\ (\text{Borrower} \bowtie \text{Loan}) \end{array} \right]$

However, we can write this query in efficient way.

$$\Pi_{\text{cname}} \left[\text{Borrower} \bowtie \left(\sigma_{\text{Branch} = 'ABIDS'} \text{Loan} \right) \right]$$

Here no of tuples joined will be reduced by a large factor.

This query requires less buffer size and less no of comparisons.

Normal query : Join & select

Optimized query : select & join

→ Also while joining, performing join on smaller tables first is more efficient

(P/8) If bldg is first and borrowing is not required

Here we have

bldg is subject to an explicit join on

branch \bowtie (account \bowtie depositor)

To optimize we perform select first and then join

$$\left(\sigma_{\substack{\text{branch} \\ \text{b-city} \\ = 'Agra'}} \right) \bowtie \left[\left(\sigma_{\substack{\text{account} \\ \text{balco}}} \right) \bowtie \text{depositor} \right]$$

To optimize further we can first joint branch & account

∴ account & depositor are bigger relations

$$\therefore \left(\sigma_{\substack{\text{branch} \\ \text{b-city} \\ = 'Agra'}} \bowtie \left(\sigma_{\substack{\text{account} \\ \text{balco}}} \right) \right) \bowtie \text{depositor}$$

But from the given options, opt@ is the best.

P/4

~~R(A, B, C) S(B, D, E)~~

For every tuple in S , we have a tuple in R such that $S \cdot B = R \cdot B$

In $R \times S$ we get 100 correct combinations

for $A \cdot B$ i.e. 100 tuples in $R \times S$

This maximum possible

i.e., if there is $R \cdot B$ for every $S \cdot B$

P/5

Here we have exactly 2000 since R in g_1

is first ref. column in g_1 & therefore minimum value of $R \cdot V$ is 2000

Note:

whenever join is performed b/w parent & child, the no of tuples = no of tuples in child.

P/11

$$\pi_X \left(\sigma_{P.Y = R.Y} (P \times R) \right) \text{ is number of minima of}$$

$$R.V = 2^N \text{ i.e. } \{ (x_1, y_1, z_1, v_1) \} \text{ i.e. } \{ (x_1, y_1, z_1) \}$$

$$\pi_X \left(\sigma_{Q.Y = R.Y} (Q \times R) \right) = \{ (x_2, y_2, z_2, v_2) \}$$

$$\text{number of minima of } Q.V = 2^M \text{ i.e. } \{ (x_2) \}$$

$$\therefore \{ (x_2) \} - \{ (x_1) \} = \{ (x_2) \}$$

(9/9) $\Pi_{cid} \left[enroll - \left\{ \left(\Pi_{sid} \left(\sigma_{sex='female'} (studinfo) \right) \right) \times \Pi_{cid} (enroll) \right\} \right]$

gives all ids of female students

gives all courses enrolled by atleast one student

gives all possible combinations of female sids and cids

Now removing this from enroll has no effect on tuples involving males.

\therefore It gives all courses in which male students enrolled.

(P/10) $\Pi_{name} \left(\sigma_{sex='female'} (student) \right)$

: gives names of all female students.

$\Pi_{name} \left[(student) \setminus \left(\begin{array}{l} sex=female \\ Ax = male \\ marks \leq m \end{array} \right) \cup \Pi_{n,x,m} (student) \right]$

: give names of all female students whose marks are less than or equal to atleast one boy.

$\therefore \{ \text{all female} \} - \{ \text{female whose marks} \leq \text{atleast one boy's marks} \}$

\therefore All females whose marks greater than all the boys.

Division (\div)

This used to compare a tuple of a relation with all the tuples of other relation.

Eg : Consider two relation $A(x,y)$ $B(y)$

$A \div B$ contains all x value such that x is related to ~~every~~ every value of y in B .

Eg :

| A | | B | | $A \div B$ | |
|-------|-------|---|-------|------------|--|
| X | Y | | Y | X | |
| x_1 | y_1 | | y_1 | x_1 | |
| x_1 | y_2 | | y_2 | | |
| x_1 | y_3 | | y_3 | | |
| x_2 | y_1 | | | | |
| x_2 | y_2 | | | | |
| x_3 | y_1 | | | | |

Eg : Consider below relations

| Student | | Enroll | | Course | |
|---------|------|--------|-----|--------|------|
| Sid | Name | Sid | Cid | Cid | Name |
| 1 | Ajay | 1 | 9 | 9 | DB |
| 2 | Rama | 2 | 10 | 10 | ids |

Query : find name of student who has enrolled for all the courses

Sol :

Comparing with the previous example here we have

X : Name Y : Courseid

So we need to create two relation $\Pi_{name}(\text{Student})$ and $\Pi_{cid}(\text{Enroll})$ and perform division.

$$\therefore \Pi_{(name, cid)}(\text{Student} \bowtie \text{Enroll}) \div \Pi_{cid}(\text{Course})$$

| name | cid |
|------|-----|
| Ajay | 9 |
| Ajay | 10 |
| Rama | 10 |

| cid |
|-----|
| 9 |
| 10 |

After performing division we get student Ajay

to satisfy the condition of student Ajay

| cid | name |
|------|------|
| Ajay | 9 |

the same

One more of writing query is

student who is not enrolling to

$$\Pi_{name} \left\{ \text{Student} \bowtie \left[\text{Enroll} \div \left(\Pi_{cid}(\text{Course}) \right) \right] \right\}$$

old student satisfies to

be enrolled in student

student to be old

student which is

youngest student

youngest student

$$\{ (name, cid) | \text{Student} \bowtie \text{Enroll} \}$$

$$\{ (cid) | ? \}$$

but student

which satisfies

student

satisfies

student who is not enrolling to

student who is not

student who is not enrolling to

student who is not

student who is not enrolling to

student who is not

student who is not enrolling to

student who is not

student who is not enrolling to

student who is not

Relational Calculus:

- It is a non-procedural language (i.e., it doesn't give the process)
 - It is called declarative language.
 - It describes result of query but not the step by step procedure.
 - two types:
 - (i) Tuple Relational calculus (TRC)
 - (ii) Domain Relational calculus (DRC)
- In TRC, result is described as set of tuples.
- In DRC, result is described as set of values of an attribute.

| TRC | DRC |
|--|---|
| → result is described as set of tuples | → result is described as set of values of an attribute |
| → set of tuples are represented as tuple variable | → set of values are represented as domain variable |
| → form of TRC query $\{T / P(T)\}$ ↓ tuple variable formula that describes T | → form of DRC query $\{\langle d_1, d_2, \dots, d_n \rangle / P(d_1, d_2, \dots, d_n)\}$ ↓ domain variables formula that describes domain variables |

Ex: Consider below relations

student (Rno, name, marks)

Enroll (Rno, Cno)

Course (Cno, cname, duration)

In TRC, for a variable R, if we don't write $\exists R$ relation, then the convention is that R has all the fields of relations that appear in the formula

Q1: Find name of the student whose marks > 10;

TRC: $\{ T \mid \exists S \in \text{student} (S.\text{marks} > 10 \wedge T.\text{name} = S.\text{name}) \}$

tuple
variable

This attribute can be given any name (for display)

DRC: $\{ \langle N \rangle \mid \exists R, M [\langle R, N, M \rangle \in \text{student} (M > 10)] \}$

Domain
variable representing
name

Q2: Find name of the student who enrolled a course, with ~~eno~~ eno = 9.

TRC: $\{ T \mid \exists S \in \text{student} (T.\text{name} = S.\text{name}) \}$

\wedge

$\exists E \in \text{Enroll} (E.\text{cno} = 9 \wedge S.\text{Rno} = E.\text{Rno}) \}$

DRC: $\{ \langle N \rangle \mid \exists R, M [\langle R, N, M \rangle \in \text{student}] \}$

\wedge

$\exists C [\langle R, C \rangle \in \text{Enroll} (C = 9)] \}$

Here R need not to be mentioned again

Here we need value of 'c' to be 9. So this query can also be written as

$$\{ \langle N \rangle \mid \exists R, M [\langle R, N, M \rangle \in \text{student}] \\ \wedge [\langle R, 9 \rangle \in \text{Enroll}] \}$$

Q3: find name of the student and name of the course who enrolled a course of duration 3 months.

$$\text{TRC: } \{ T \mid \exists S \in \text{student} (T.\text{name} = S.\text{name}) \}$$

$$\wedge \exists C \in \text{Course} (C.\text{duration} = 3 \wedge T.\text{name} = C.\text{name}) \\ \wedge \exists E \in \text{Enroll} (E.rno = S.rno \wedge E.cno = C.cno) \}$$

$$\text{DRC: } \{ \langle N, C \rangle \mid \exists R, M [\langle R, N, M \rangle \in \text{student}] \\ \wedge \exists C, \text{duration} [\langle C, \text{name}, \text{duration} \rangle \in \text{Course}] \}$$

$$\{ \langle N, C \rangle \mid \exists R, M [\langle R, N, M \rangle \in \text{student}] \}$$

$$\wedge \exists cno, c, 3 [\langle cno, c, 3 \rangle \in \text{Course}]$$

$$\wedge \exists [\langle R, cno \rangle \in \text{Enroll}]$$

Eg: What does below query return

$$\{T \mid \exists P \in \text{professor} (T.\text{name} = P.\text{name})\}$$

$$\wedge \exists C \in \text{Course} (C.\text{cname} = \text{'DBMS'} \wedge P.\text{ssn} = C.\text{ssn})\}$$

Sol:

It returns names of professors who teach DBMS

Unsafe Query:

A query language than can express all the queries possible with relational algebra, is said to be relationally complete

$$\{T \mid T \in \text{ER}\}$$
 returns finite set no of tuples

which may be absent and result may show error

$$\{T \mid \sim(T \in \text{ER})\}$$
 return infinite no of tuples

→ A query that returns infinite no of tuples is called unsafe query.

→ Relational algebra has only safe queries.

Relational algebra \cong safe relational calculus



Q-07 Consider the relation employee(name, sex, supervisorName),

what does below query produce.

$$\{e.name \mid \text{employee}(e) \wedge$$

$$\forall x [\text{employee}(x) \vee x.\text{supervisorName} \neq e.name, \\ \vee x.sex = "male"] \}$$

Now we write it as

{ e.name / employee(e)

^

{ e.name / employee(e) \wedge $\exists x [employee(x) \wedge x.supervisorName = e.name \wedge x.sex = "male"] \} = "Frank"$

Structured Query Language

Order of evaluation:

(i) from clause

If we write more than one table in from clause,
Cartesian product is performed.

(ii) where clause is evaluated after from clause
selection of rows is done

(iii) group by

divide rows into groups

(iv) HAVING

select the formed groups

(v) Expressions in select clause are evaluated

(vi) distinct in select

(vii) order by

sort the rows

After this the o/p is produced

Simple Select:

- > select name from student;
- > select name, branch from student;
- > select * from student;
- > select rno, marks + 5 from student; // Mathematical expressions are also possible.
- > select distinct name from student;

Select with where:

where is used for selecting rows based on some conditions.

→ The conditions may be connected by below operators if needed

→ and, or, not

→ in, not in

→ between .. and

→ not between .. and

→ is null, is not null

→ <, <=, >, >=, <>, ~~=~~, =
↳ not equal

→ like '%', like '_'

'!= is not there
in SQL'

Eg: Consider a question for selecting data :-

Student(Rno, name, branch, year, marks, passport)

> select *
from student
where branch = 'CSE' and year = '1';

Strings in SQL are enclosed
within single quotes

Q/P:

It displays all students of CSE studying 1st year.

The evaluation is done row by row starting from the 1st row

$\Rightarrow \underline{\underline{=}}$
 where not branch = 'CSE'; } gives all student records who
 (01)
 $\underline{\underline{=}}$
 where branch \neq 'CSE'; } include mark & marks
 in operator; (~~set~~ comparison operator)
 where branch = 'CSE' or branch = 'IT' or branch = 'ME'
 the above can be written as

\Rightarrow where branch in ('CSE', 'IT', 'ME');
 where branch not in ('CSE', 'IT', 'ME');
 Opposite of in operator
 where branch not in ('CSE', 'IT', 'ME');

* gives students of branches other than CSE and IT and ME

between...and:
 where marks between 10 and 20;
 inclusive

\approx where marks ≥ 10 and marks ≤ 20

Similarly

where marks not between 10 and 20; exclusive

\approx where marks < 10 or marks > 20

is null & is not null

To check if the value is present or not.

where passport is null; \rightarrow students who don't have passport

where passport is not null; \rightarrow students who have passport.

like :

\rightarrow used for pattern ~~match~~ matching

'%o' \rightarrow matches to 'o' or more characters

'-' \rightarrow matches to one character.

Eg: where name like 'R-%'; // student name starting with R

where name like '%H'; // student names ending with H

where name like '%S%'; // containing S

where name like '___'; // names of length 3

~~water~~

* 'R---%' // names of length atleast 3 and starting with R.

Note :

> select *
from student
where 1=2;

This kind of queries are valid and it returns 0 rows.

> where 1<2 \rightarrow all rows are printed in O/P

> select *
from student
where .NULL.=0;

NULL is a keyword and it is not zero
 \therefore 0-rows will be returned

> select
where NULL<>0; \rightarrow return 0 rows
(Comparing null with any value is always false) (\because null is unknown)

> where NULL=NULL; \rightarrow returns 0 rows.

- > where NULL is NULL ; ... returns all the rows.
- > where NULL \neq NULL ; ... returns 0 rows
- Comparison of NULL with any value using relational operator is false.
- Comparison of NULL with 'is NULL' is only true.

Aggregate Functions:

~~SQL has~~

- Aggregate functions take a set of values and return one value.
- SQL has 5 aggregate functions

- i) Min
- ii) Max
- iii) Sum
- iv) Avg
- v) Count

Eg: Student

| Name | M1 | M2 |
|------|------|------|
| A | 10 | 5 |
| B | null | 10 |
| C | 20 | 15 |
| D | null | 20 |
| E | 30 | null |
| F | null | 25 |

Note:

Aggregate functions ignore null value

i) select min(M1) from student;

O/P: 10

ii) select max(Name) from student;

O/P: F

COUNT
MIN, MAX can take
input either text or
number

(iii) Select sum(m₁) from student;

O/P: 60

(iv) select sum(m₁, m₂) from student;

This gives error

→ aggregate func takes only one set as i/p

* (v) select sum(m₁+m₂) from student;

*

~~error~~

m₁+m₂ is single set. So no error.

Any mathematical operation with null giving null

null + any value is null

(∴ unknown + any value is unknown)

(15+null+35+null+null)

(makes the result unknown so error)

(∴ O/P: 50 at this time didn't printing out)

(vi) select sum(m₁) + sum(m₂) from student

∴ sum(m₁) = 60

sum(m₂) = 75

∴ O/P: 135

Aggregate functions with DISTINCT

i) COUNT(DISTINCT A):

no of distinct values in column A.

ii) SUM(DISTINCT A):

sum of all unique values

iii) AVG(DISTINCT A)

Avg of all DISTINCT values

* Using DISTINCT with MIN & MAX doesn't make any sense

(vii) select count(m₁) from student;

O/P: 3

(viii) select count(*) from student;

O/P: 6 ∴ i.e. no

(ix) select avg(m1) from student;

$$Q.P: 20 \quad (\text{i.e., } \frac{60}{3})$$

$$\rightarrow \text{avg}(m1) = \frac{\text{sum}(m1)}{\text{count}(m1)}$$

(x) select sum(name) from student;

avg & sum works only
only on numerical data

(xi) select name from student

→ where m1 = max(m1);
~~where m1 = max(m1);~~

~~where~~:
Name is under group function

This gives error.

(\because where is a row function (i.e., it selects row) and
max is a group function (i.e., max considers all column)
so performing this is difficult for query processor)

Note:

→ Using aggregate functions in where clause give error.

→ However we can use aggregate function in 'Having' class.

Set manipulation Operators

→ we have 3 set manipulation operators

R { Union
 intersect
 except }

difference

Set manipulation operators
eliminates duplicates.

R & S must be compatible.

Order by:

used to sort the rows

Eg:

Company



| Cno | Name | turnover |
|-----|------|----------|
| 1 | B | 5Cr |
| 2 | A | 3Cr |
| 3 | B | 4Cr |
| 4 | C | 6Cr |
| 5 | B | 5Cr |

Query: Display all details of company in order of their name

> select *

from company

order by name

ASC;

for descending order we use DESC

Query: Display company details in descending order of their

name and increasing order of their turn over

> select *

from company

order by name desc, turnover asc;

Q2 :

| Cno | cname | turnover |
|-----|-------|----------|
| 4 | C | 6Cr |
| 3 | B | 4Cr |
| 1 | B | 5Cr |
| 5 | B | 5Cr |
| 2 | A | 3Cr |

While storing if

the attribute values are same, then ordering will be based on insertion i.e., first inserted records will come first

Group by and Having clause:

→ Group by used for grouping rows
and 'having' specifies grouping conditions.

Eg:-

Student

| Rno | Name | Branch | Year |
|-----|------|--------|------|
| 1 | A | CS | I |
| 2 | B | CS | I |
| 3 | C | CS | II |
| 4 | D | IT | I |
| 5 | E | IT | I |
| 6 | F | ME | II |

Q: Find no of students in each branch.

> select branch, count(*) -- ③
from student -- ① } order of execution
group by branch; -- ②

Op:-

| Branch | count(*) |
|--------|----------|
| CS | 3 |
| IT | 2 |
| ME | 1 |

Q: Find no of student in each branch and year.

> select branch, year, count(*)
from student
group by branch, year;

O/P:

| branch | year | count(*) |
|--------|------|----------|
| CS | I | 2 |
| CS | II | 1 |
| IT | I | 2 |
| ME | II | 1 |

Note:

Consider below query happens with distinct counts prints

> select branch, year, count(*)

from student

group by branch

This will result an error

Note: condition like and having conditions are applied to

Every attribute of select ~~clauses~~ must be either present in group by

clause or other aggregate functions can also appear in

select clause. but group by clause can't be present in

Q: Find no of students in each branch if the branch contains atleast two students.

> select branch, count(*)
 from student
 group by branch
 having count(*) ~~>~~ ≥ 2

Conditions of
group by are
specified using
having clause

O/P:

| branch | Count(*) |
|--------|----------|
| CS | 3 |
| IT | 2 |

Q: Find the no of students in each group other than cs.

> select branch, count(*)

from student

group by branch

having branch <> 'cs';

Note :

Having clause contains either aggregate functions or attributes that appear in group by.

Where vs Having

→ where selects rows

→ having selects groups

→ aggregate functions cannot be used in where

aggregate functions can be used in having.

20/11/20

* Q: What is the o/p of the following query based on the above table.

above-student table.

> select count(*)

from student

having count(*) >= 3;

a) error

b) 6

c) 0

d) unpredictable o/p.

Sol:

Since group by is not specified, entire table is considered one group.

$$(count*) = 6 \geq 3$$

$\therefore \text{OPT} : \frac{\text{count(*)}}{6}$

if all group by columns are present
then it will be considered as a single group

$\therefore \text{opt(b)}$

in case of having multiple group by
columns consider each of them as a group

Note:

→ A query containing having clause without group by, treats whole table as a single group.

(P/7) A, B and aggregate function can appear in select clause

$\therefore \text{opt(b)}$

(P/8) If a is key

$\text{count}(a) > 1$ will always be false
and no O/P will be given if a is key

$\therefore \text{opt(c)}$

JOIN:

Consider

professor(ssn, name, dept)

course(cno, cname, ssn)

Q: find name of the professor teaching database.

> select name

from professor, course → Cartesian Product

where professor.ssn = course.ssn and cname = 'Database',
↓ join condition ↓ selection condition

→ Here we are first performing join and then selecting rows. This is inefficient (for large tables)
i.e., we are specifying join condition and selection condition in the same where clause.

join...on :

→ To remove the previous discussed inefficiency
~~we use join...on~~

→ Now the same query is written as

> select name

from professor join course on professor.ssn = course.ssn
where cname = 'Database';

This query, in terms of efficiency and evaluation is same as that of the previous query, but it is more readable.

Natural join :

→ Equality is implicitly specified on ~~equal~~ same named fields

Now the same query can be written as

> select name

from professor natural join course
where cname = 'Database';

Database table by name loan_records is given below

A
G-11

| Borrower | Bank Manager | Loan amount |
|----------|--------------|-------------|
| Ramesh | Sunderajan | 10000 |
| Suresh | Ramgopal | 5000 |
| Maresh | Sunderajan | 7000 |

What is the o/p of the following SQL query?

SELECT count(*)

from

Select borrower, bank_manager from loan_records) as S

Natural join is to consider bank manager field as common field.

(Select bank_manager, loan_amount from loan_records) as T.

O/P of 1st sub query O/P of 2nd sub query

| | |
|--------|------------|
| Ramesh | Sunderajan |
| Suresh | Ramgopal |
| Maresh | Sunderajan |

| | |
|------------|-------|
| sunderajan | 10000 |
| Ramgopal | 5000 |
| Sunderajan | 7000 |

∴ The natural join is

| | | |
|--------|------------|-------|
| Ramesh | Sunderajan | 10000 |
| " | " | 7000 |
| Suresh | Ramgopal | 5000 |
| Maresh | Sunderajan | 10000 |
| " | " | 7000 |

∴ 5

Outerjoin:

It is used to join one relation (all rows) even if there is no matching row.

Eg: Consider

| Professor | | Course | | |
|-----------|------|--------|-------|------|
| SSN | Name | c_id | Cname | SSN |
| 1 | A | 9 | DB | 1 |
| 2 | B | 10 | DS | 3 |
| 3 | C | 11 | CD | 1 |
| 4 | D | 12 | CO | NULL |

Q: Find name of all the professors and name of the courses that they teach if any.

i.e., even if a professor doesn't teach any course, we need to display his name.

In situations like this we use outer join.

→ Based tuple of which relation we need completely.

Outer join is of 3 types:

i) Left outer join (keyword: left join)

ii) right outer join (keyword: right join)

iii) Full outer join (keyword: full join)

Consider

> select name, cname
from professor natural join course;

O/p:

| name | cname |
|------|-------|
| A | DB |
| A | CD |
| C | DS |

Here we don't get professor names who ~~teach~~ doesn't teach any course.

> So we use left-outer join, as shown below

> select name, cname

from professor natural left join course;

O/p:

| name | cname |
|------|-------|
| A | DB |
| A | CD |
| C | DS |
| B | null |
| D | null |

Note:

The same above query can be written as shown below

> select name, cname

from professor, course

where professor.ssn = course.ssn(+);

(+) says course.ssn is optional.

E: Consider

| <u>R</u> | |
|------------------------------------|----------------|
| A | B |
| a ₁ | b ₁ |
| a ₂ | b ₁ |
| a ₁ | b ₃ |
| a ₂ | b ₄ |
| <u>a₃ b₅</u> | |

| <u>S</u> | |
|------------------------------------|----------------|
| B | C |
| b ₁ | c ₁ |
| b ₂ | c ₂ |
| b ₃ | c ₂ |
| b ₄ | c ₁ |
| <u>b₅ c₁</u> | |

Find o/p for each of the below queries

(i) > select *

from R natural join S;

| A | B | C |
|----------------|----------------|----------------|
| a ₁ | b ₁ | c ₁ |
| a ₂ | b ₁ | c ₁ |
| a ₁ | b ₃ | c ₂ |
| a ₂ | b ₄ | c ₁ |

(ii) > select *

from R natural left join S;

A B C
a₁

no of rows
 $4+1=5$

| A | B | C |
|----------------|----------------|----------------|
| a ₁ | b ₁ | c ₁ |
| a ₂ | b ₄ | c ₁ |
| a ₁ | b ₃ | c ₂ |
| a ₂ | b ₄ | c ₁ |
| a ₃ | b ₅ | |

(iii) > select *

from R natural right join S;

A B C
a₁

then

| A | B | C |
|----------------|----------------|----------------|
| a ₁ | b ₁ | c ₁ |
| a ₂ | b ₄ | c ₁ |
| a ₁ | b ₃ | c ₂ |
| a ₂ | b ₄ | c ₁ |
| a ₃ | b ₅ | |

| A | B | C |
|----------------|----------------|----------------|
| a ₁ | b ₁ | c ₁ |
| a ₂ | b ₁ | c ₁ |
| | b ₂ | c ₂ |
| a ₁ | b ₃ | c ₂ |
| a ₂ | b ₄ | c ₁ |

no of rows
 $= 4 + 1 = 5$

(iv) select *

from R natural join S;

| A | B | C |
|----------------|----------------|----------------|
| a ₁ | b ₁ | c ₁ |
| a ₂ | b ₁ | c ₁ |
| a ₁ | b ₃ | c ₂ |
| a ₂ | b ₄ | c ₁ |
| | b ₂ | c ₂ |
| a ₃ | b ₅ | |

no of rows
 $= 4 + 1 + 1 = 6$

Note:

→ no of rows in natural join

$$= \text{no of matching rows}$$

→ no of rows in left outer join

$$= \text{no of matching rows} + \text{no of non-matching rows for left table}$$

→ no of rows in right outer join

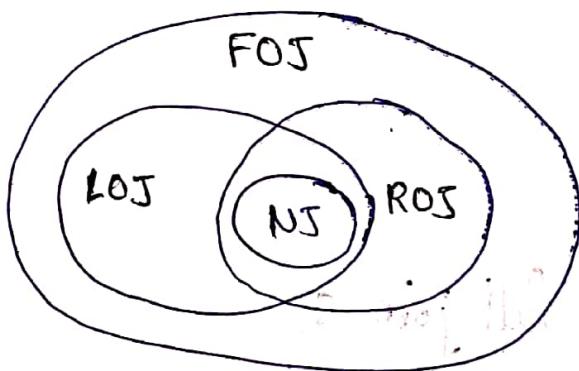
$$= \text{no of matching rows} + \text{no of non-matching rows for right table}$$

→ no of rows in full outer join

= no of matching rows +

no of non-matching rows of left table &

no of non-matching rows of right table.



Set of tuples returned
by NJ

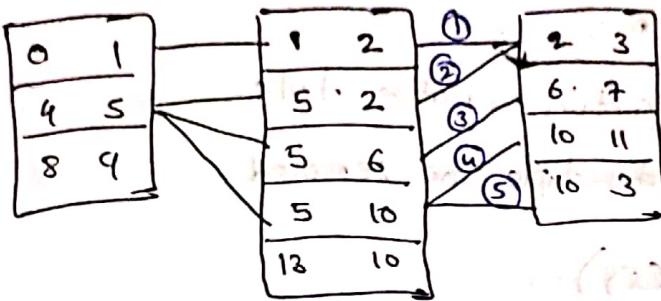
\subseteq Set of tuples returned
by LOJ

Set of tuples returned
by NJ

\subseteq Set of tuples returned by ROJ

$$\left\{ \begin{array}{l} \text{Set of tuples by LOJ} \\ \cup \\ \text{Set of tuples by ROJ} \end{array} \right\} = \left\{ \text{Set of tuples of FOJ} \right\}$$

P/9



$\therefore 5$ rows output because each pair

has 5 pairs.

P/13

R

| A | B | C |
|---|---|---|
| 1 | 5 | 7 |
| 3 | 7 | * |
| 4 | | 9 |

∴ 3 rows output because each pair

\therefore opt C (optimal)

P/16

 $\pi(A, B)$ $\sigma_{B < 5}(S)$

\sqsubset

This may have unmatched tuples for π_1 since we are

Selecting tuples of S under condition $\sigma_{B < 5}$.

But these unmatched tuples are not displayed.

a) $\sigma_{B < 5}(\pi_1 \bowtie S)$

\sqsubset contains only matching tuples

\sqsubset selection

Now this doesn't contain unmatched tuples.

b) $\sigma_{B \in S} (g_1 \text{ LOS})$.

↳ contains unmatched tuple

↳ unmatched tuples are removed

c) $g_1 \text{ LOS} (\sigma_{B \in S} (g_1))$

This may have unmatched tuples

∴ opt ②

Sub Queries:

Nested query: It is a query within a query.

(Query1
 (Query2))

Query 1 → outer query

Query 2 → inner query

| P | S | A |
|---|---|---|
| P | S | A |
| F | E | B |
| F | E | C |
| F | E | D |
| F | E | F |
| F | E | G |
| F | E | H |
| F | E | I |
| F | E | J |
| F | E | K |
| F | E | L |
| F | E | M |
| F | E | N |
| F | E | O |
| F | E | P |
| F | E | Q |
| F | E | R |
| F | E | S |
| F | E | T |
| F | E | U |
| F | E | V |
| F | E | W |
| F | E | X |
| F | E | Y |
| F | E | Z |

→ The result of inner query is given to outer query.

Ex: Consider

| Professor | |
|-----------|------|
| SSN | Name |
| 1 | A |
| 2 | B |
| 3 | C |

| Teaches | | Course | |
|---------|-----|--------|-------|
| SSN | Cid | Cid | Cname |
| 1 | 9 | 9 | DB |
| 2 | 10 | 10 | DS |
| 2 | 10 | 10 | DS |

Q: Find name of the professor who teaches the course with Cid 9.

→ One way to write a query for this is using this which requires join of 2 tables.

> select name

from Professor P, Teaches T,
where P.ssn = T.ssn and ~~T.cid = 9~~ T.cid = 9;

This method is inefficient for larger tables.

Here condition & where clause is checked on $2 \times 3 = 6$ tuples.

→ If prof has 2000 tuples & teaches has 3000 tuples,

then join would result in 6000,000 tuples and select
condition has to be checked on 60 lakh tuples.

→ Rather than this we can first select tuples with cid=9
from teachers and then perform join.

This can be done using nested queries.

> select name

from professor

where ssn in (select ssn

from Teachers

where cid = 9);

This is evaluated by first inner query.

> select name

from professor

where ssn in (i) 2000 tuples + 3000 tuples

∴ O/P: $\frac{\text{Name}}{A}$

→ Here condition is checked only on 3000 tuples and hence this more
efficient. So no of tuples involved = $2000 + 3000 = 5000$

21/11/20

Ques. 1. Write a query

Q: Find name of the professor who teaches the course 'Database'.

Ans: Select name
from professor
where ssn in

(Select ssn

from teaches
where cid in

(Select cid

from course

where cname = 'DB');

O/P: name

Ans: A is final ans since both rows

→ If the same query is written using join, then would require lot of space.

Q: Student

| Name | Marks |
|------|-------|
| A | 90 |
| B | 80 |
| C | 70 |
| D | 60 |

Find name of the student with highest marks

Ans: Select name

from student

where marks = (Select max(marks)

from student);

O/P: name
A

→ Here inner query gives a set but since '=' operator is used the set is converted into value implicitly.

Q: find name of student with 2nd maximum marks

∴ > select name

from student

where marks = (select max(marks)

from student

where marks < > (select max(marks))

limit 1 offset 1 from student);

To make above query simple we use limit clause.

Q: find name of student with 3rd maximum marks

Syntax:

limit x offset y

It skips first y entries and returns next x entries

∴ > select name

from student

order by marks desc

limit 1 offset 2;

O/P:

| name |
|------|
| c |

Correlated subquery:

~~Outer~~ (outer
(inner))

for every tuple of outer query, inner query is executed once.
(i.e., every row of inner query)

→ Here both outer & inner are dependent on each other and hence the name correlated.

(row of outer query → Complete inner query → based on o/p of inner query) outer query

→ Now consider finding student name with third highest marks in below table.

| Student | |
|---------|-------|
| Name | marks |
| A | 90. |
| B | 70 |
| C | 80 |
| D | 60 |

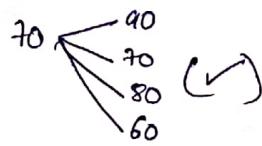
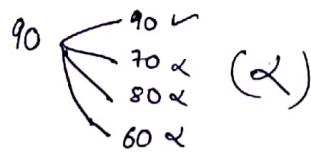
> select s1.name

from student s

where \langle select count(*) \rangle

from student s2

where $s1.name = s2.marks) = 3$,



Op. Name
B

However there is flaws with this query that it doesn't work with duplicate marks.

e.g:

| student | |
|---------|-------|
| Name | marks |
| A | 90 |
| B | 70 |
| C | 80 |
| D | 60 |
| E | 70 |

out requirement is 3rd highest marks, but not

3rd highest rank. So the expected o/p is

| o/p: | Name |
|------|------|
| 3rd | B |
| and | E |

So the appropriate query can be written using DISTINCT

► select s1.name

from student s1

where ~~s1.marks~~ (select count(distinct s2.marks)

from student s2

where s1.marks <= s2.marks) = 3;

Note:

→ Syntactically correlated subquery and nested subquery looks similar.

So whenever inner query uses reference of outer query we say query is correlated subquery.

Q
G-05

The relation book(title, price) contains the titles and prices of different books. Assuming that no two books have the same price, what does the following SQL query list?

```
Select title  
from book as B  
where (select count(*)  
from book as T  
where T.price > B.price) <= 5;
```

- a) titles of 4 most expensive books
- b) titles of the fifth most inexpensive book
- c) titles of fifth most expensive book
- d) titles of five most expensive books

Sq:

Here we print B.title

if there are 4 or less than four books with price greater than B.price.

If there are 4 books whose price is greater than B.price then ~~and~~ B is 5th expensive book.

If there are 3 such books then B is 4th expensive book.

2 --- 3rd

1 --- 2nd

0 --- 1st

∴ It gives 5 most expensive books

Exists:

This is used to compare with an empty set.

i.e., to find whether inner query has returned any row or not.

$\text{exists}(\text{empty set}) = \text{false}$

$\text{exists}(\text{non empty set}) = \text{true}$

Eg: consider

| student | |
|---------|------|
| Rno | name |
| 1 | A |
| 2 | B |

| enroll | |
|--------|-----|
| Rno | Cno |
| 1 | 98 |
| 2 | 99 |

Q: find name of the student who enrolled a course with cno 98

> select name

from student

where EXISTS(

select *

from enroll

where student.Rno = enroll.Rno

and Cno = 98)

op of inner query

point to this { (1, 98) in q } (1, 98)
(2, 99) X }

2 { (1, 98) X } -
(2, 99) X }

$\therefore \text{O/P: NAME}$

A

Set Comparison Operators

Any / some :

any item in the set satisfies the condition

$$x < \text{any}(5, 10, 15)$$

$$\approx (x < 5) \text{ OR } (x < 10) \text{ OR } (x < 15)$$

ALL :

$$x < \text{all}(5, 10, 15)$$

$$\approx (x < 5) \text{ AND } (x < 10) \text{ AND } (x < 15)$$

Eg: Consider

| Instructor | | | |
|------------|------|------|--------|
| Iid | Name | Dept | Salary |
| 1 | A | CS | 50k |
| 2 | B | CS | 60k |
| 3 | C | ME | 90k |
| 4 | D | EC | 55k |
| 5 | E | CE | 40k |

Q: Find name of the instructor whose salary is greater than some instructor of the department CS

> select a.name from instructor a

from instructor b

where salary > Any (select salary

from instructor b

where b.dept = 'CS');

| Name |
|------|
| A |
| C |
| D |

similarly replacing 'any' with 'all' in the above query;

we have O/P: Name
C

Note:

→ all(empty) is always true

$x \left\{ \begin{array}{l} \leq \\ = \\ \geq \\ \geq \end{array} \right\}$ all(empty) is true.

→ any(empty set) is always false

$x \left\{ \begin{array}{l} \leq \\ = \\ \geq \end{array} \right\}$ all any(empty) is false.

→ $x = \text{any}(s) \cong x \in \text{all}(s)$

Eg: select I₂.name

from instructor I₂

where dept <> 'cs'

and salary = all(select I₁.salary

from instructor

where dept = 'IT');

O/P: 3 rows

Eg: select I₂.name

from instructor I₁

where dept <> 'cs'

and salary = any(select I₁.salary

from instructor

where dept = 'IT');

O/P: 0 rows

All evaluates until it finds false and once it finds false it returns false. For empty set it doesn't find any false value so it returns true

22/11/19

Multicolumn Subquery

Consider

Student (sid, ...)

Takes (sid, cid, sem, year)

| | | | | |
|------|------|-----|----|------|
| John | 2001 | 99 | I | 2005 |
| John | 2001 | 100 | II | 2006 |
| John | 2002 | 99 | I | 2005 |
| John | 2002 | 99 | I | 2006 |

Instructor (Id, ...)

Teachers (Id, cid, sem, year)

| | | | |
|----|----|----|------|
| 9 | 99 | I | 2005 |
| 10 | 98 | II | 2006 |
| 9 | 99 | II | 2006 |

Q: Find the total no of distinct students who have taken some course in a ^{some} sem and year, taught by an instructor with id = 9.

→ We need to compare multiple attributes.

i.e., pairwise comparison.

→ we can do this using IN and this type of query is called multicolumn subquery

> select count(distinct sid)

from takes

where (cid, sem, year) in (select cid, sem, year

↓
pairwise
comparison

from teachers

where id = 9);

O/P: count (distinct cid)

Unique:

The unique construct tests whether a subquery has any duplicate tuples in its result.

- The unique construct evaluates to "true" if a given subquery contains no duplicates otherwise false.

Eg. Consider

| Course | |
|--------|-------|
| Cid | Cname |
| 9 | DB |
| 10 | DS |
| 11 | PL |

| offerings | | |
|-----------|-----|------|
| cid | sem | year |
| 9 | I | 2009 |
| 9 | II | 2009 |
| 10 | I | 2009 |
| 11 | II | 2010 |

Q: Find all courses that were offered almost once in 2009.

select cname
from course
where unique(select cid
from offerings
where year=2009);

where course.cid = offerings.cid
and year=2009);

O/P: Cname

DS

With:

The 'with' clause provides a way of defining a temporary relation whose definition is available only to the query in which the with clause occurs.

→ 'with' clause helps write complex queries ~~in a simple way~~ in a simple way.

Q: Consider

| Instructor | | | |
|------------|------|------|--------|
| Id | Name | Dept | Salary |
| 1 | A | CS | 20k |
| 2 | B | CS | 20k |
| 3 | C | CS | 30k |
| 4 | D | IT | 30k |
| 5 | E | IT | 30k |
| 6 | F | ME | 20k |

Q: Find all departments whose ~~total~~ total salary is greater than average salary at all departments.

$$\text{Total salary of CS} = 70k$$

$$\text{Total salary of IT} = 60k$$

$$\text{Total salary of ME} = 20k$$

$$\text{Avg total salary at all depts}$$

$$\frac{70 + 60 + 20}{3} = 50k$$

> with total-sal (dept, total) as
 (select dept, sum(salary)
 from instructor
 group by dept)
 dept-avg (total2) as
 (select avg(total)
 from total-sal)

with
clause

| total-sal | |
|-----------|-------|
| dept | total |
| CS | 70k |
| IT | 60k |
| ME | 20k |

→ Now we find dept-avg

| dept-avg | |
|----------|------|
| total2 | dept |
| 50k | CS |

Here total-sal, dept-avg

are temporary relations
 created using with clause

select dept
 from total-sal, dept-avg
 where total > total2;

∴ Dept
 CS
 IT
 Having found that dept = CS
 sum of all (dept) distinct dept

→ The above query can also be written as shown below
 without using with clause

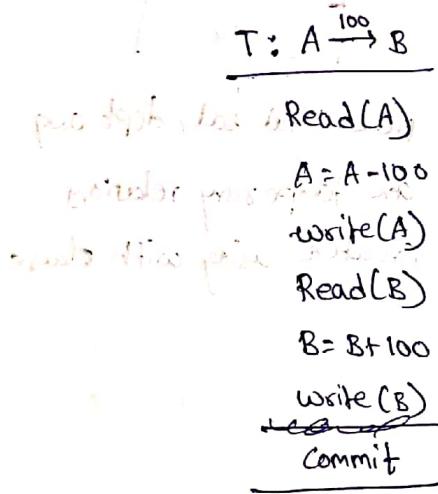
> select dept
 from instructor
 group by dept
 having sum(salary) >
 (select sum(salary)/count(distinct dept))

Transactions & Concurrency Control

Transaction:

It is a collection of operations that performs a single logical unit of work.

Eg: Consider transferring Rs. 100 from account A to B.



A [500] B [1000]

A [400] B [1000]

A [400] B [1100]

→ values will be saved permanently and rollback (undo) is no more possible.

→ Once transaction is committed it cannot be rolled back.

Properties of a transaction (ACID Properties):

i) Atomicity:

A transaction has to execute either all of the operations or none of the operations.

(~~if one fails, then all fails~~)

→ Transaction manager (component of DB) checks if this property is maintained or not.

iii) Consistency:

→ State of system has to be correct always.

(For every transaction, we may have a consistency requirement.)

Eg: In $T_1: A \xrightarrow{100} B$

$A+B=1500$ is a consistency requirement

→ This consistency is ensured by user or app. programmer.

iii) Isolation:

Each transaction is unaware of other transaction that are executing.

i.e., Each transaction is executing independently.

→ Concurrency control manager supports this isolated execution of transactions.

iv) Durability:

All updations done by a transaction should become permanent.

Schedule:

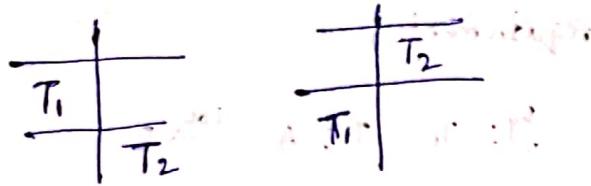
Schedule represents the order in which the operations of transactions ~~execute~~ execute. Schedule is of 2 types:-

(i) Serial schedule

(ii) Concurrent schedule

i) Serial Schedule:

Each transaction is executed one after other without ~~inter~~ intervention.



For example, if T_1 and T_2 have 3 operations each.

→ A serial schedule always ensures consistency.

→ If T_1 and T_2 have 3 operations each.

ii) Concurrent Schedule:

→ Here, intervention is possible.

for e.g., T_1 has 3 operations & T_2 has 3 operations.

Let T_1 & T_2 be 2 transaction with 3 operations in each.

| T_1 | T_2 |
|-------|-------|
| 2 | 1 |
| 3 | 2 |
| 3 | 1 |
| 2 | |

→ Concurrent schedule may lead to inconsistent state.

| T_1 | T_2 |
|--------------|--------------|
| R(A) | |
| A = A + 1000 | R(A) |
| w(A) | A = A + 1000 |
| | w(A) |

→ Concurrent schedules minimize waiting time, response time.

Note:

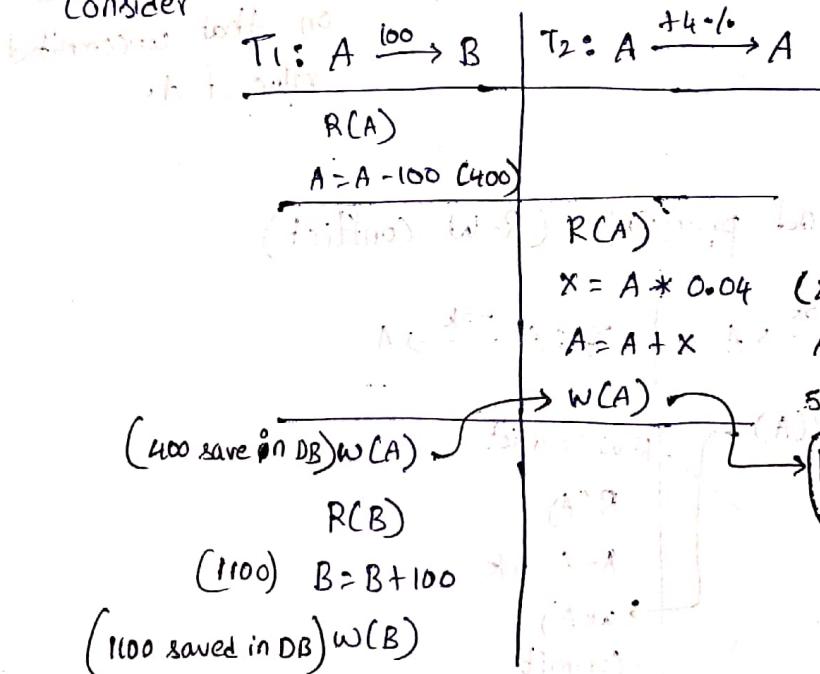
- no of serial schedules possible with n transactions = $n!$
- no of concurrent schedules possible with n transaction
 - where each transaction has m_1, m_2, \dots, m_n operation(s)

$$\frac{(m_1 + m_2 + \dots + m_n)!}{m_1! * m_2! * \dots * m_n!}$$

Problems due to concurrent schedule:

ii) Lost update problem (W-W conflict)

Consider



lost update

Assume initially $A=500$ & $B=1000$

After operation $A = \cancel{500} \quad B = 1100$

but correct state is $A=420$ & $B=1100$

Defn: If update of one transaction is overwritten by the update of another transaction without reading it, then it causes lost update problem.

(ii) Dirty Read Problem (W-R conflict)

→ Reading uncommitted data is called **dirty read**.

Consider

| $T_1: A \xrightarrow{100} B$ | $T_2: A \xrightarrow{+4\%} B$ |
|-----------------------------------|---|
| $R(A)$ $A = A - 100$ $w(A)$ | $R(A)$ (reading uncommitted data) $x = A * 0.04(16)$ $A = A + x$ $w(A)$ |
| $R(B)$ $B = B + 100$ $w(B)$ | If T_1 rolls back after this read then value of A is undone by T_1 , but still T_2 performs operations on that uncommitted value of A . |

(iii) Unrepeatable read problem (R-W conflict)

| $T_1: A \xrightarrow{-10k} A$ | $T_2: A \xrightarrow{-5k} A$ |
|---|--|
| $A = 12k$ $(12k)$ <i>unrepeatable reads (i.e., reads with diff values) (7k)</i> | $R(A)$ <i>R-W conflict</i> $R(A)$ $A = A - 5k$ $w(A)$ <i>commit</i> |
| $A = A - 10k$ $w(A)$ | $R(A)$ $A = A - 10k$ $w(A)$ |

After unrepeatable read, roll back is performed.

If we allow T_1 to continue even after the unrepeatable read, assume a case where T_1 fails and needs to roll back.

Since T_1 is depending on T_2 , it also requires T_2 to roll back which is not possible. cuz, T_2 has committed already.

Note:

Checking for above 3 problems in a given schedule.

i) Lost update:

$w_i(A)$ followed by $w_j(A)$

T_j performs $w_j(A)$ without reading value v_i produced by T_i from $w_i(A)$

ii) Dirty read:

$w_i(A) \rightarrow R_j(A)$

$R_j(A)$ comes before Commit of T_i

iii) Unrepeatable read:

$R_i(A) \xrightarrow{w_j(A)} R_i(A)$ (old value is read) (new value is read)

Eg: Consider schedules

S1: R₁(A) R₂(B) W₂(B) R₁(A) W₁(A) R₂(A) W₂(A) C₁ C₂

S2: R₁(A) R₂(B) W₂(B) R₂(A) W₂(A) C₂ R₁(A) W₁(A) C₁

S3: R₁(A) R₂(B) W₂(B) R₂(A) W₁(A) W₂(A) C₁ C₂

i) which schedule has unrepeatable read?

~~S1: R₁(A) R₂(B) W₂(B) R₁(A)~~

(A) (B) (C)
unrepeatable read

Ans: S2

ii) which schedule has dirty read problem?

| S1: | T ₁ | T ₂ |
|-----|----------------|----------------|
| | R(A) | |
| | | R(B) |
| | | W(B) |
| | R(A) | |
| | | W(C) |
| | | R(A) |
| | | W(C) |
| | C ₁ | |
| | | C ₂ |

∴ Ans: S1

iii) which schedule has lost update problem?

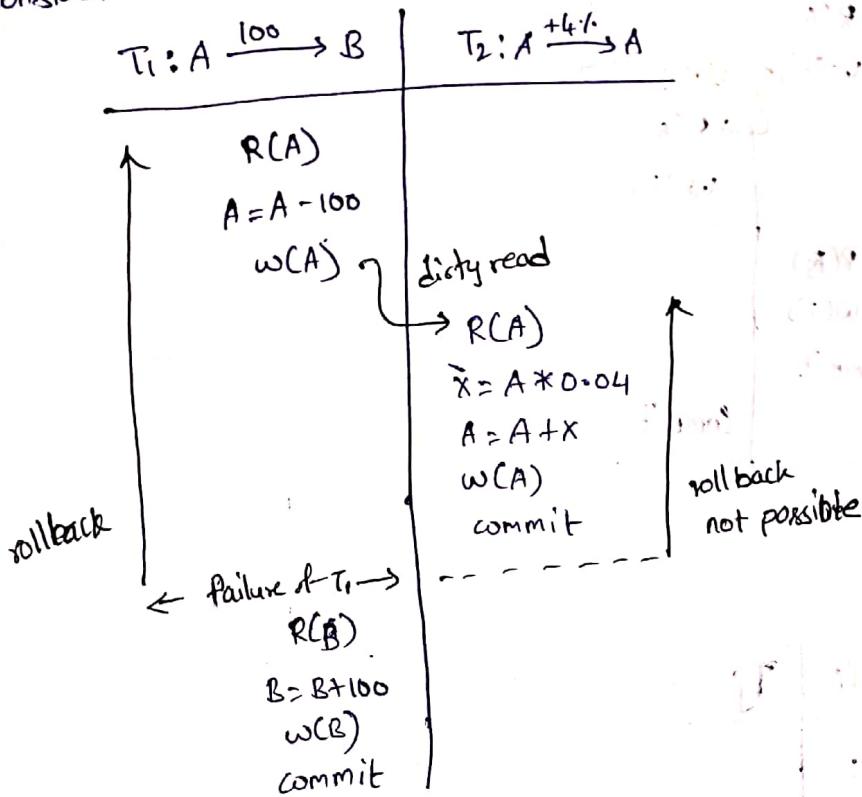
Ans: S3

lost update
(A) (B)

Types of schedules based on recoverability

i) Non-Recoverable Schedule

Consider



→ Here ~~the~~ operation of T_2 is dependent on T_1 .

But $\rightarrow T_2$ has committed before T_1 and hence if

T_1 needs to rollback then also should T_2 but

this is not possible because T_2 has committed.

This is known as non-recoverable schedule.

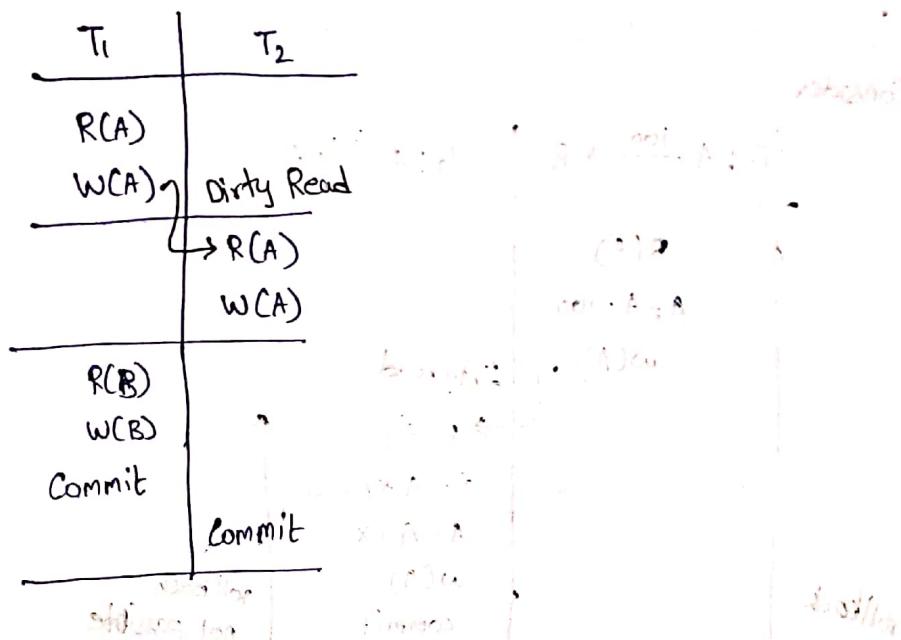
→ So the schedule will be recoverable if T_1 commits before T_2 .

→ Non-recoverable \Rightarrow Dirty read exists.

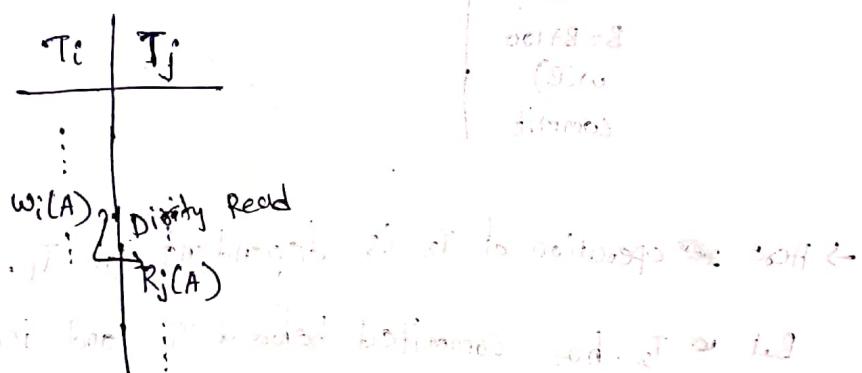
→ No Dirty read \Rightarrow Recoverable

(ii) Recoverable:

→ Same as previous case where T_2 commits after T_1 .



Note:



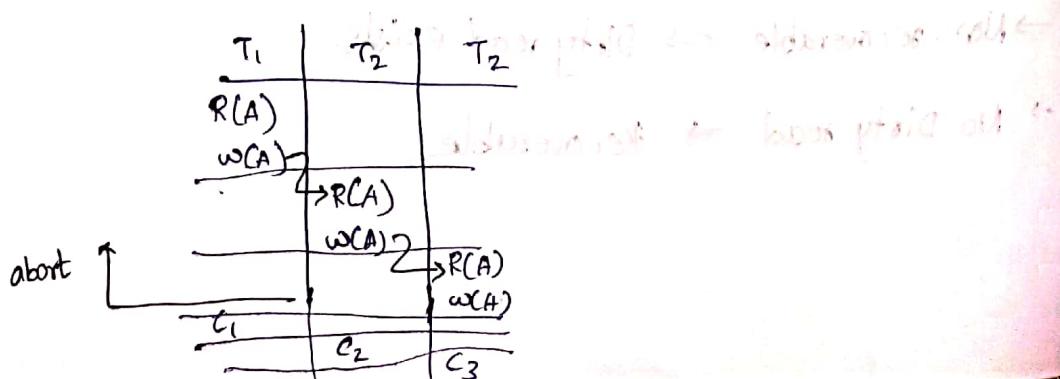
→ T_j commits before $T_i \Rightarrow$ Non-Recoverable

→ T_j commits after $T_i \Rightarrow$ Recoverable

24/11/20

(iii) Cascading Rollbacks / Cascading Aborts:

Consider below schedule



If T_1 needs to abort & rollback, it would cause T_2 to abort & rollback which further require abortion of T_2 . This is called ~~as~~ cascading rollbacks / cascading aborts.

Defn: A schedule is called cascading rollbacks, if abortion of one transaction leads to abortion of other transactions.

(iv) Cascadeless Rollbacks:

Consider below schedule,

| T_1 | T_2 |
|--------|-------|
| R(A) | |
| W(A) | |
| | R(CB) |
| | W(CB) |
| R(C) | |
| W(C) | |
| Commit | |

\rightarrow have placed A & B in T_2 because T_1 has committed

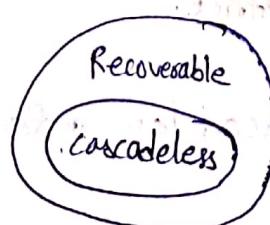
\rightarrow when T_1 fails, T_2 will not fail

$R(A)$, \rightarrow committed read
 $W(A)$'s previous operation in T_1
~~Commit~~ Commit -> dirty reads kept

Here failure of T_1 doesn't cause T_2 to abort.

So this schedule is called cascadeless rollbacks.

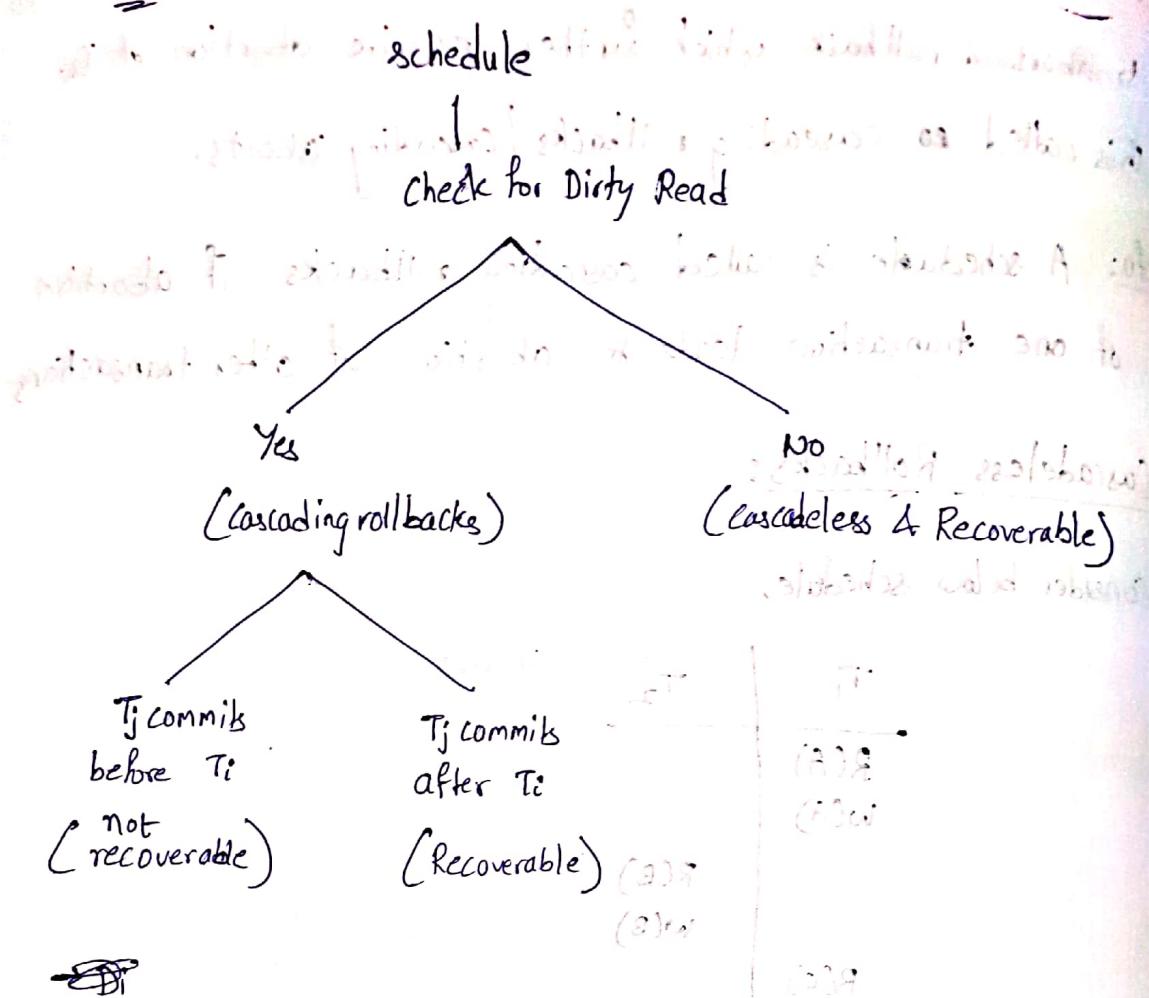
\rightarrow Every cascadeless schedule is recoverable but not vice versa



Cascading rollback
 \Rightarrow Dirty read
 No dirty \Rightarrow No CR

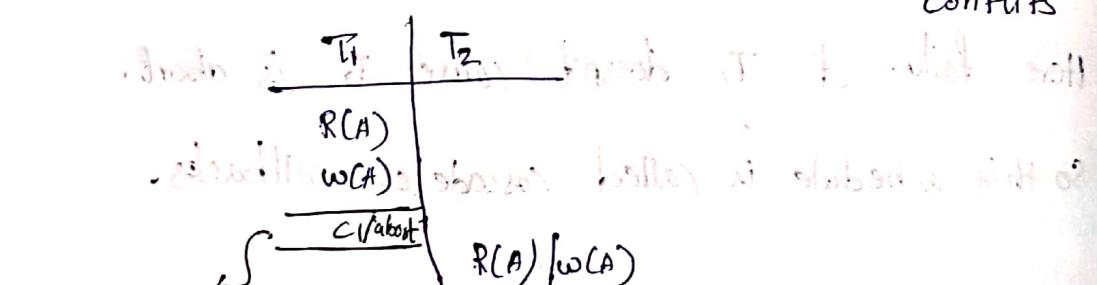
\rightarrow ~~Dirty Read \Rightarrow Cascading Rollbacks~~

Note:

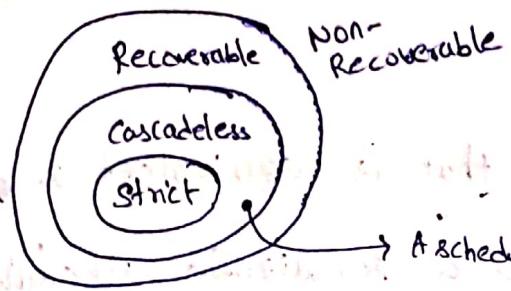


(v) Strict Schedule:

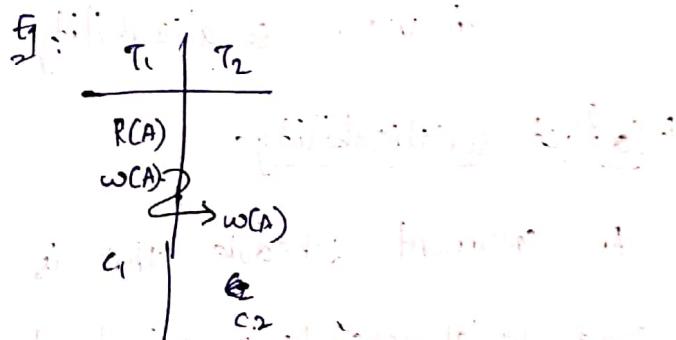
If it is schedule which doesn't have dirty read & lost update problem. i.e. no ~~w-wt-w-r~~ conflicts



→ Strict schedule is cascadeless & Recoverable (\because no w-r)



A schedule may be cascadeless & recoverable but not strict. This is possible by ww conflicts.



(P4) serial $\rightarrow 2! = 2$

$$\text{concurrent} \rightarrow \frac{(8+3)!}{5! \cdot 3!} = \frac{8 \times 7 \times 6}{3 \times 2 \times 1} = 56$$

$$\text{non-serial} \rightarrow 56 - 2 = 54$$

(P5) T₂ performed dirty read & committed first

If it is not recoverable, T₁ cannot be rolled back and only some part of T₁ is executed.

These bugs are observed due to not ensuring atomicity.

Serializability:

Any concurrent schedule that is equivalent to some serial schedule is said to be a serializable schedule.

two types:

(i) Conflict serializability

(ii) View serializability

(i) Conflict serializability:

Any concurrent schedule that is conflict equivalent to some serial schedule is said to be a conflict serializable schedule.

→ Two schedules are said to be conflict equivalent

iff both the schedules have same set of conflicts.

Conflict operations:

$W_1(A)$ $W_2(A)$

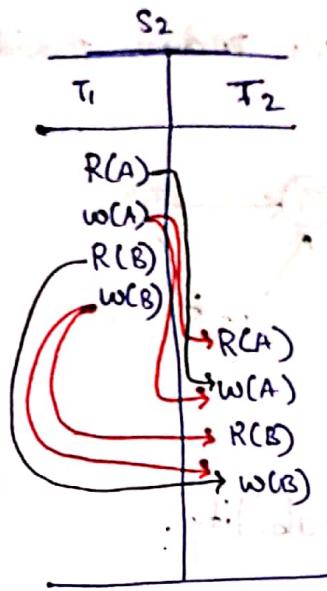
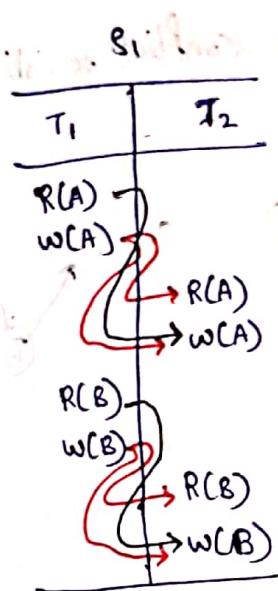
$W_1(A)$ $R_2(A)$

$R_1(A)$ $W_2(A)$

→ If two schedules have same set of conflicts then the result of both the schedules will be same.

So a conflict serializable schedule has equal results as that of a serial schedule which is consistent.

Q: Consider two schedules



Conflicts in S₁ \cong Conflicts in S₂

\therefore S₁ and S₂ are conflict equivalent.

\therefore S₁ and S₂ can be interleaved.

Since S₂ is a ~~conflict~~ serial schedule, S₁ is also serial.

S₁ is said to be view serializable.

Test for conflict serializability:

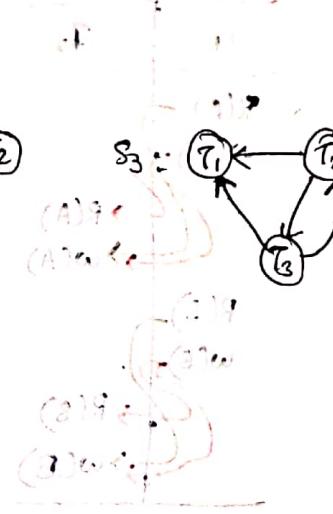
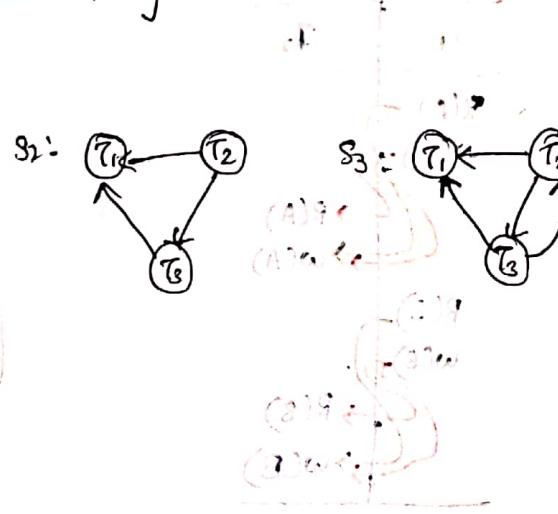
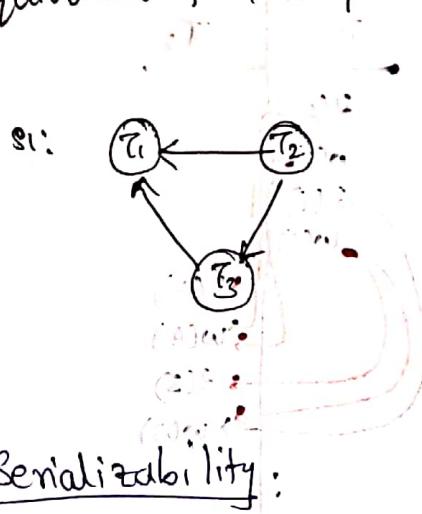
(i) Construct a precedence graph in which each transaction is a vertex and each conflict operation from T_i to T_j is a directed edge from T_i to T_j.

(ii) If the resultant graph contains no cycles, then the schedule is conflict serializable.

(iii) The order of serializability is determined based on topological sort of the precedence graph.

P/17

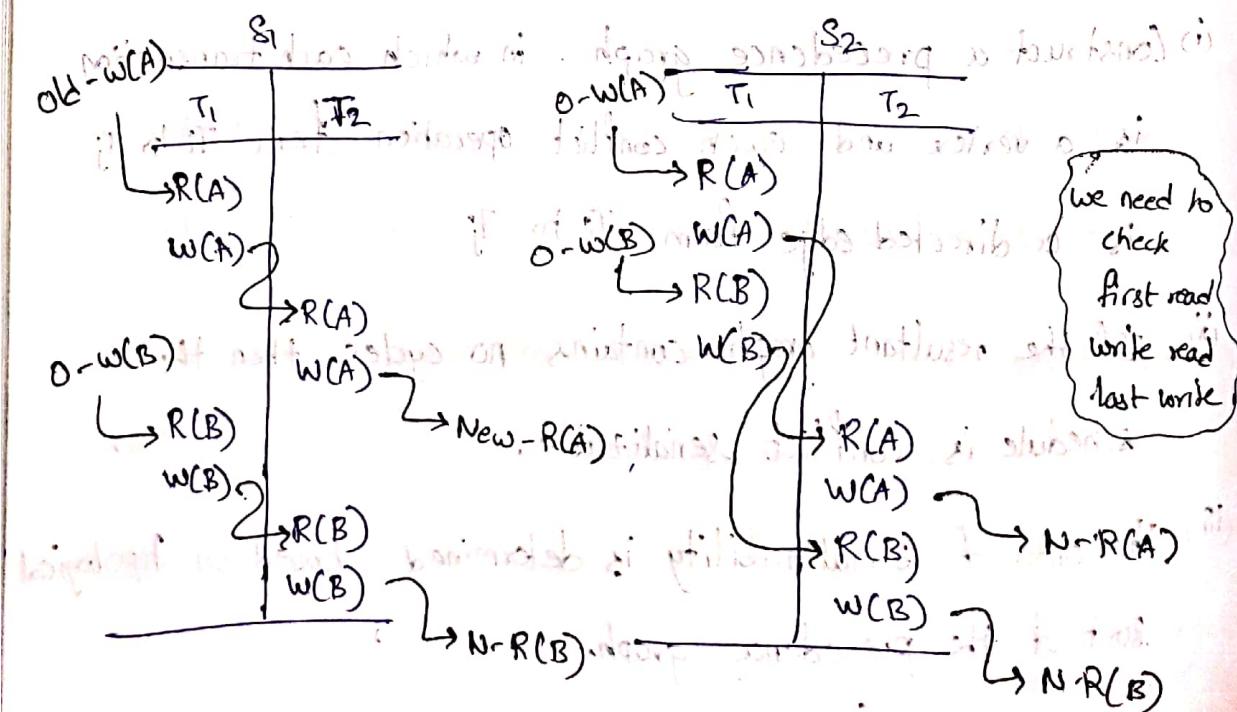
schedules with same precedence graphs are conflict equivalent provided they are conflict serializable.



View Serializability:

- A concurrent schedule is said to be view serializable if it is view equivalent to some serial schedule.
- Two schedules S_1 & S_2 are said to be view equivalent iff all the views $(W_i(Q) \rightarrow R_j(Q))$ are same in both the schedules.

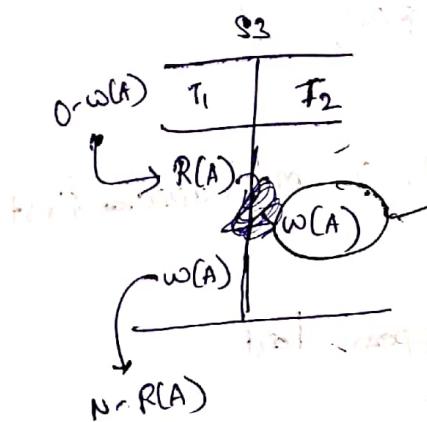
Eg: Consider below schedules.



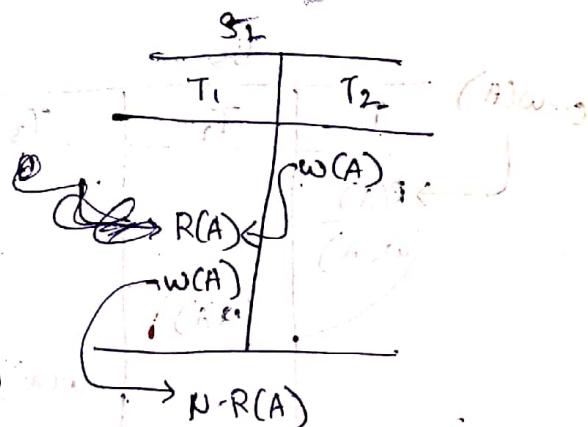
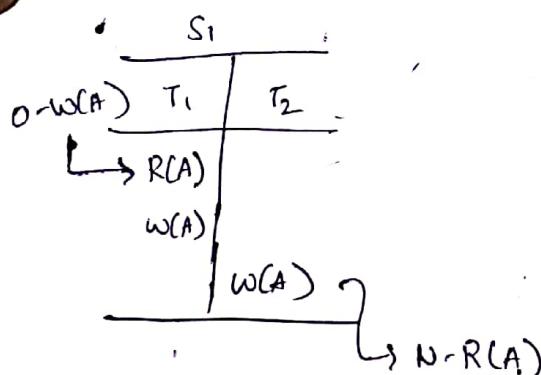
$\rightarrow s_1$ & s_2 has same views

$\therefore s_1$ is view equivalent to s_2 and s_2 is a serial schedule
 $\therefore s_1$ is view serializable.

Eg.: Check whether below schedule is view serializable or not.



Blind write.



$\therefore s_3$ is neither equivalent to s_1 nor to s_2

$\therefore s_3$ is not view serializable.

Note:

\rightarrow Testing view serializability is ~~called~~ NP hard problem.

Blind write:

A write performed without read is called blind write.

Ex: Check whether below schedule is view serializable or not.

| Schedule \$4 | | | |
|--------------|----------------|----------------|--------------------|
| | \$4 | | |
| | T ₁ | T ₂ | T ₃ |
| 0-W(A) | | | |
| R(A) | | | |
| W(A) | | W(A) | |
| | | | w(A ₃) |
| | | | N-R(A) |

Since R(A) is from 0-W(A), T₁ must appear first since W₃(A) is to N-R(A)
So T₃ must appear last

Consider serial schedule T₁, T₂, T₃

| Schedule \$5 | | | |
|--------------|----------------|----------------|--------------------|
| | \$5 | | |
| | T ₁ | T ₂ | T ₃ |
| 0-W(A) | | | |
| R(A) | | | |
| W(A) | | W(A) | |
| | | | w(A ₃) |
| | | | N-R(A) |

The order of write operations in sub \$5 (encircled) are

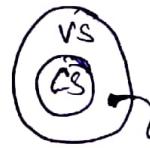
different but still \$4 is view equivalent to \$5.

\$4 is view serializable but not conflict serializable,

~~∴~~ view serializability is not as powerful as conflict serializability.

Note:

→ Every conflict serializable schedule is also view serializable



blind writes
are present

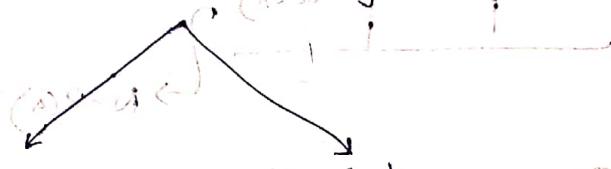
→ ~~blind writes~~ are possible

→ A schedule is VS but not CS \Rightarrow Blind writes exist.

Schedule



Draw precedence graph



(not CS) Cycle present

No Cycle

\therefore CS & VS

check for
blind writes

Yes

No

(Not CS &
Not VS)

Check view rules



(Not CS & Not VS)

Order of commits

doesn't matter for
& checking serializability

(~~not~~ established)

not
satisfied

(Not CS
but VS)



P/19

| T_1 | T_2 | T_3 |
|--------|--------|--------|
| d_1 | d_3 | d_5 |
| d_2 | d_4 | d_6 |
| $w(A)$ | $w(A)$ | $w(A)$ |

Here only order of $w(A)$ matters

Consider

| T_1 | T_2 | T_3 |
|--------|--------|--------|
| $w(A)$ | | |
| | $w(A)$ | |
| | | $w(A)$ |

$\rightarrow N-R(A)$

This view is equivalent to

| T_1 | T_2 | T_3 |
|-------|-------|-------|
| | | |

\therefore opt 2

25/11/20

Concurrency Control

Locking protocols

two types of lock modes:

i) shared lock : used for reading data (Lock-S)

ii) exclusive lock : used for both read & write (Lock-X)

Compatibility of locks

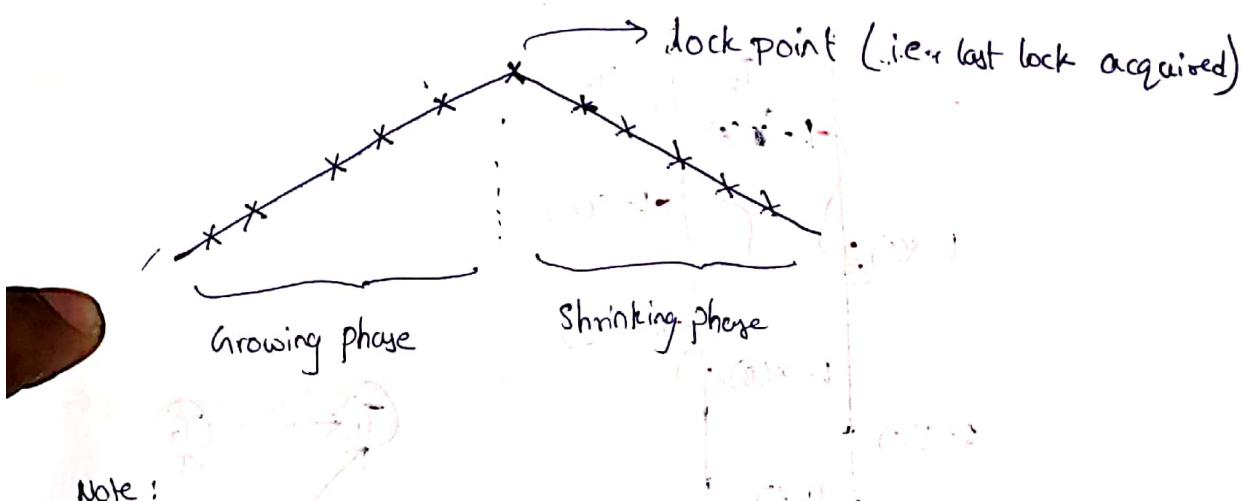
| | S | X |
|---|---|---|
| S | ✓ | ✗ |
| X | ✗ | ✗ |

Two-phase locking protocol (2PL):

Locking and unlocking is done in 2 phases

i) Growing phase: Transactions can acquire the lock but can't release the locks.

ii) Shrinking phase: Transactions can release the locks but can't obtain the locks.



Note:

Every schedule possible under 2PL always ensures conflict serializability. (if no deadlock)

Q: Consider

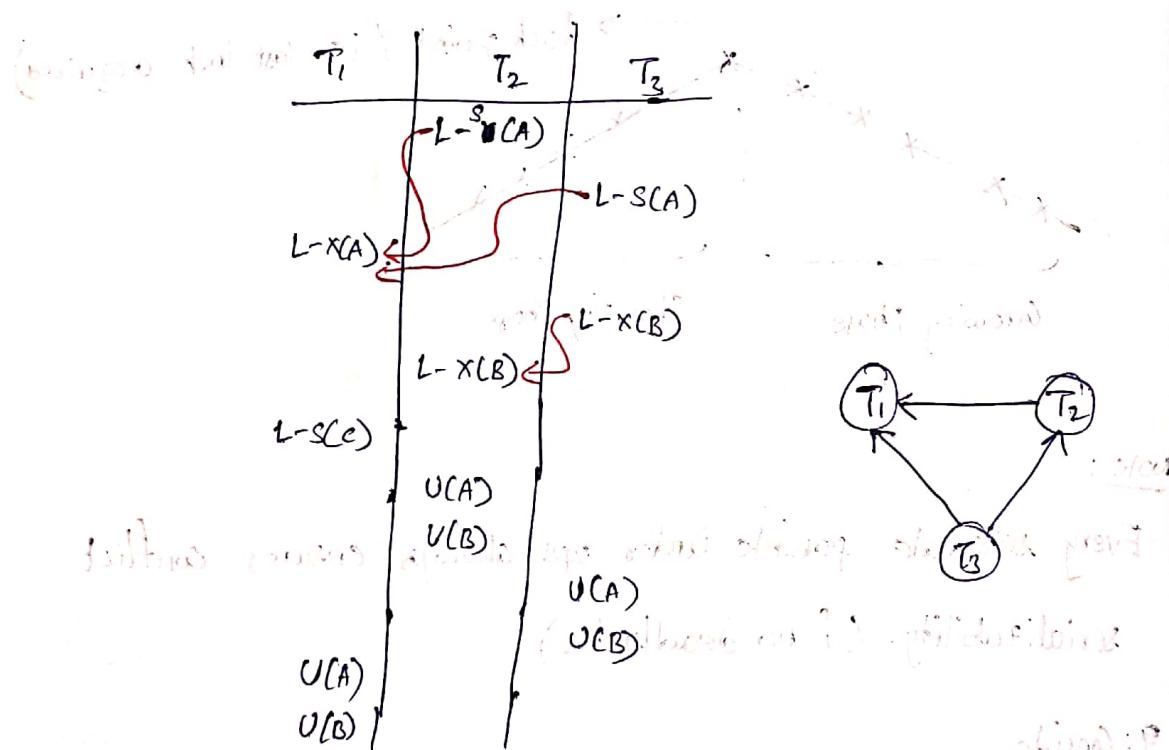
| | T_1 | T_2 |
|----|-------------------------|---------------------------------|
| G | L-X(A), R(A) W(A) | |
| G1 | | L-SC(B) \downarrow G1 R(B) |
| G2 | L-X(B) | R(B) |
| S | R(B) W(B) | U(B) \downarrow S |
| | | U(A) U(B) |

→ Here we draw precedence graph where edges represent conflicting locks.



∴ order of serializability: T_2, T_1

Eg: find the serializability order of the following serializable 2PL schedule



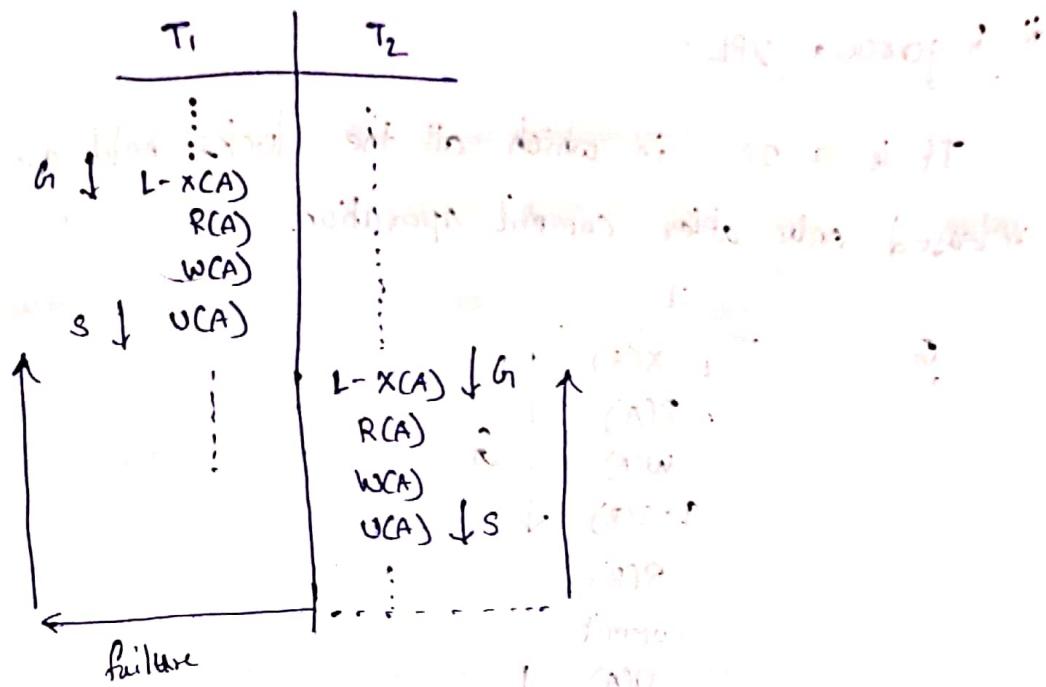
∴ order of serializability: T_3, T_2, T_1

→ If there is a deadlock, the precedence graph will have cycle and it will be conflict serializable.

Note:

Cascading rollbacks may be possible under 2PL.

Eg: Consider below example where cascading rollbacks are possible



If T_1 fails, T_2 should also roll back since there is a dirty read.

→ To avoid cascading rollbacks there are 2 modifications of 2PL

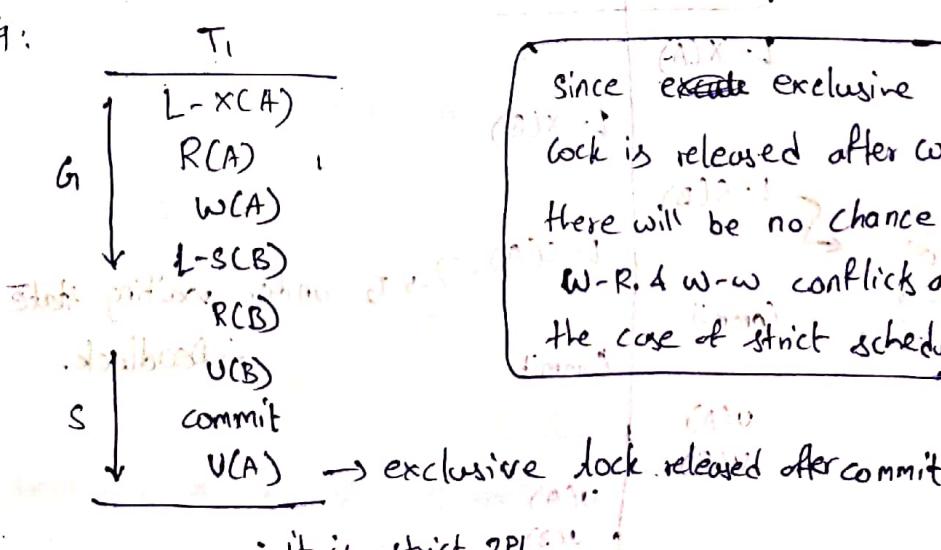
i) Strict 2PL

ii) Rigorous 2PL

(i) strict 2PL:

~~It is~~ It is 2PL in which all the exclusive locks held are released only after commit.
(Now check above example with this modification)

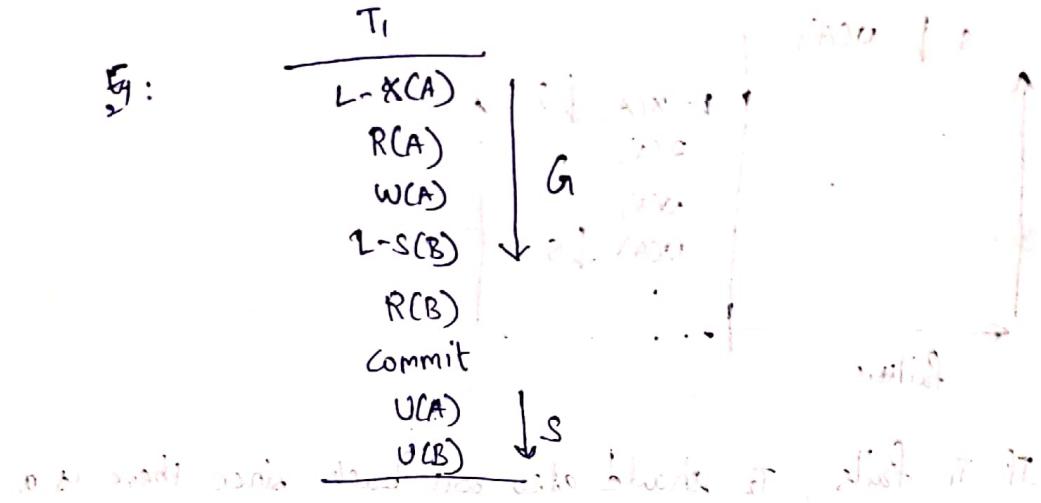
Eg:



Since ~~execute~~ exclusive lock is released after commit, there will be no chance of W-R & W-W conflict as in the case of strict schedule.

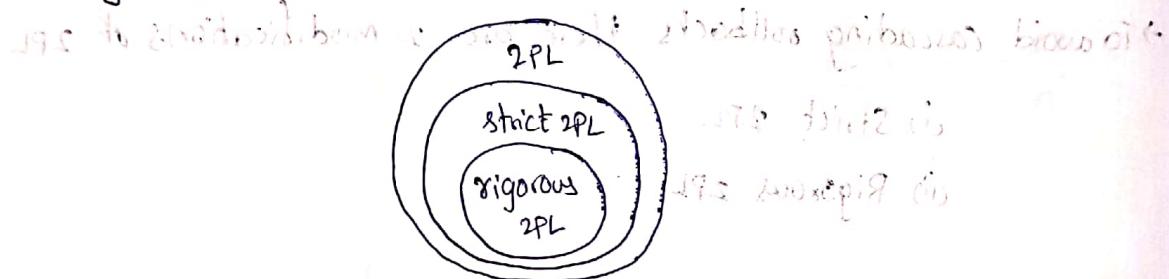
(ii) Rigorous 2PL:

It is a 2PL in which all the locks held are released only after commit operation.



Note:

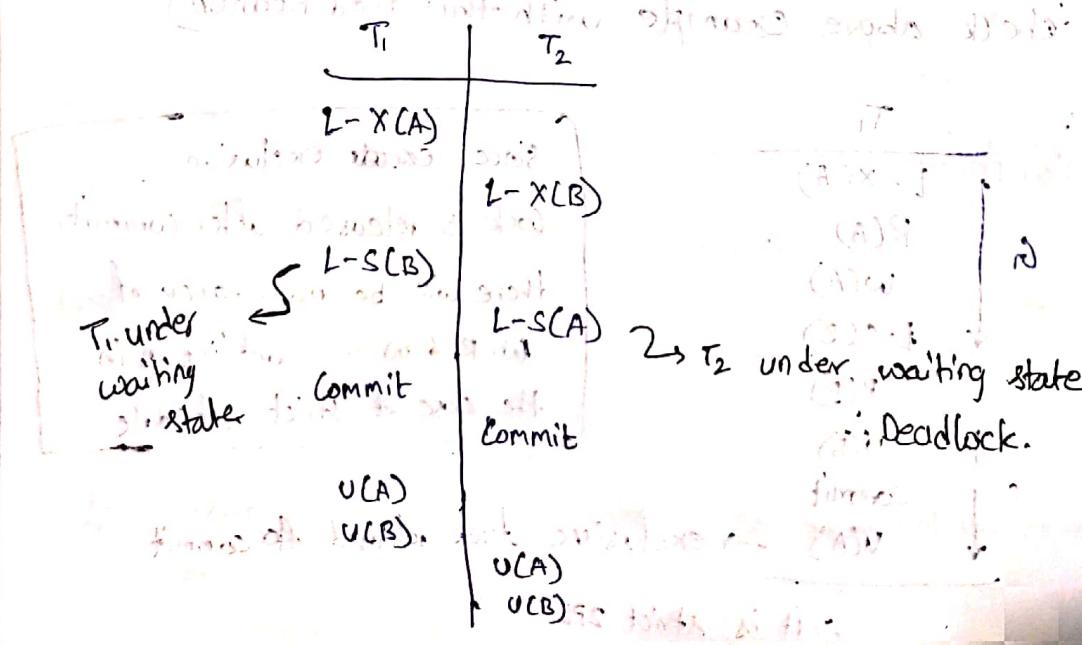
→ Every rigorous 2PL is strict 2PL



Note: Deadlocks will also arise in 2PL as it is.

Deadlocks are possible under 2PL because zero hi

(Additional information: rigorous 2PL does not allow)



Note: Basic 2PL ensures conflict serializability.

- Basic 2PL ensures conflict serializability.
- Strict & Rigorous 2PL ensures conflict serializability, ~~concerning~~ ~~concerning~~ cascadeless schedule, recoverable.

Deadlock Handling:

(i) Deadlock Prevention:

It ensures that deadlock never occurs.

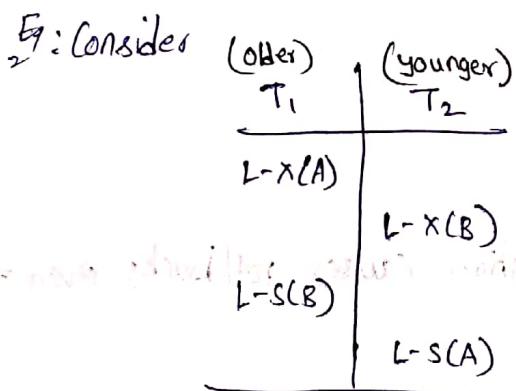
- There are 2 approaches for deadlock prevention.

a) Wait-Die: older waits; younger dies.

Every Tx has a timestamp; older Txs have lower time-stamps.

$$\text{Eg: } TS(T_i) < TS(T_j)$$

$\Rightarrow T_i$ is older process.



Here since T_1 is older, T_1 waits.

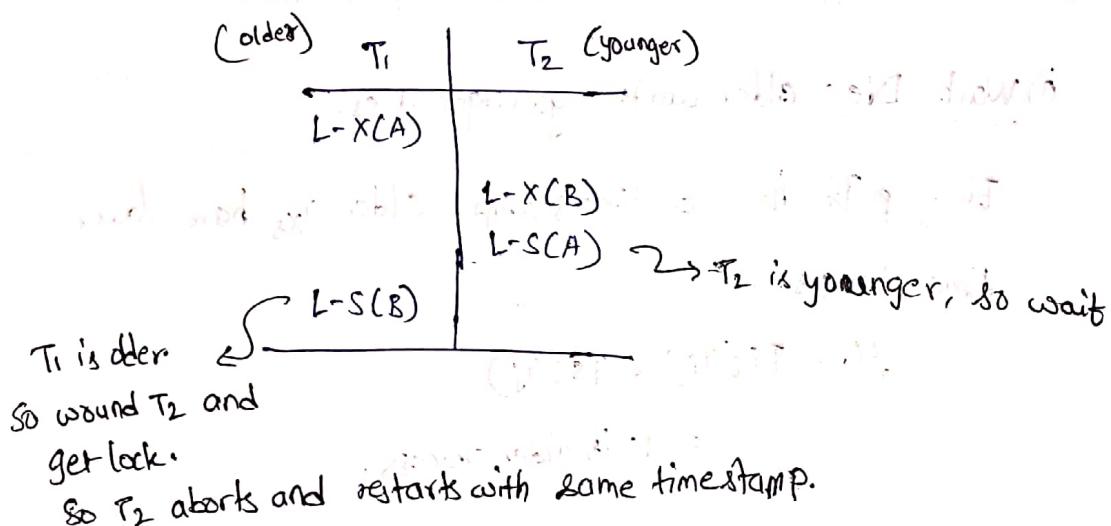
Later when T_2 request lock on T_1 , T_2 must die since T_2 is younger than T_1 .

So T_2 is aborted & rolled back and restarted with same time stamp. (To avoid starvation.)

Once T_2 dies T_1 's wait is over and continues with execution.

b) wound-wait: older wounds, younger waits.

Eg: Consider



27/11/20

(ii) Deadlock detection:

Sometimes deadlock prevention causes rollbacks even if there is no deadlock.

So when probability of occurrence of deadlock is less, it is better to go for deadlock detection.

→ This detection is done using wait-for graph.

| T_1 | T_2 | T_1 | T_2 |
|----------|----------|----------|----------|
| $L-X(A)$ | $L-X(B)$ | $L-X(A)$ | $L-X(B)$ |
| $L-S(B)$ | $L-S(A)$ | $L-S(B)$ | $L-S(A)$ |
| $U(B)$ | $U(A)$ | $U(A)$ | $U(B)$ |
| $U(A)$ | $U(B)$ | | |

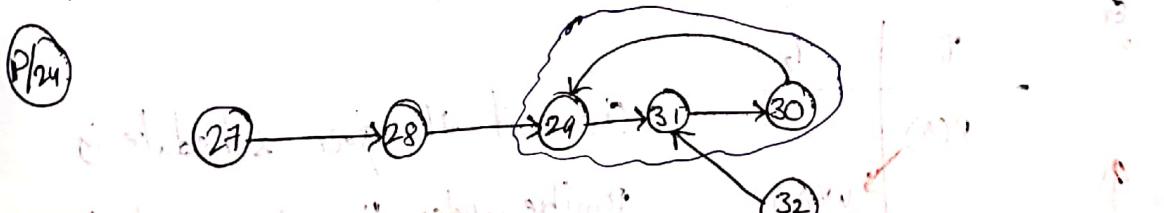


A cycle is present in the graph. No cycle exists in the graph. \therefore No deadlock. \therefore Deadlock has occurred.

\rightarrow Edge from T_1 to T_2 indicates that T_2 is waiting for a resource held by T_1 .

\rightarrow In 2nd example after $U(B)$ is executed by T_2 , the edge from T_1 to T_2 is deleted.

\rightarrow When the system detects deadlock, it selects a transaction from the cycle and aborts it.



\therefore Deadlock has occurred.

Time stamp based protocols:

Here each T_x is given a unique timestamp.

T_x with less timestamp value are more older.

- This protocol say, every schedule must be ~~equivalent~~ equivalent to serial schedule T_i, T_j such that $TS(T_i) < TS(T_j)$

→ This doesn't use any locks.

∴ This protocol is deadlock free protocol.

→ Here the order of serializability is in the order of the timestamps.

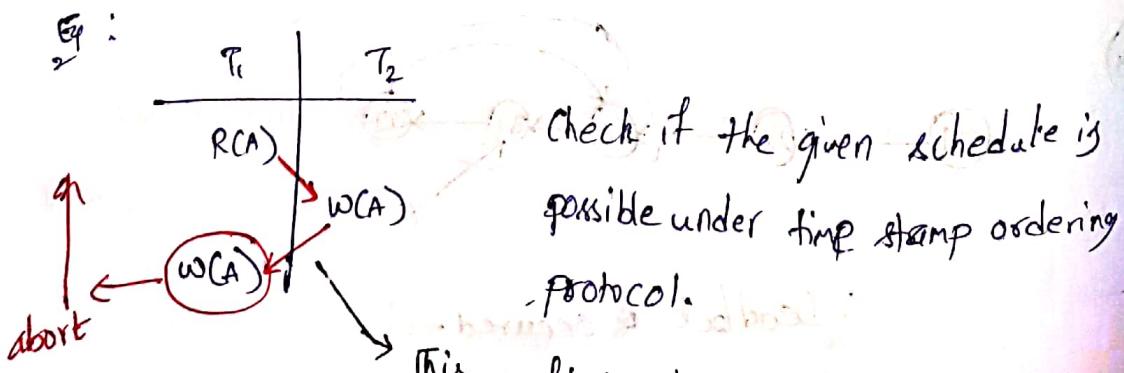
Q 2 types of time stamp based protocols:

i) Time Stamp ordering protocol.

ii) Thomas write rule.

iii) Time stamp ordering protocol.

Here all the conflicting operations must be executed in the order of their time stamps.



This conflict doesn't follow the required order of serializability. So abort T_1 .

To fall under time stamp ordering protocol it must be conflict equivalent to serial schedule T_1, T_2 . But he the precedence graph is



\therefore This schedule is not possible under time stamp ordering protocol.

Eg: Is the below schedule possible under time stamp ordering protocol?

| T_1 | T_2 |
|----------------|-------|
| R(A) | |
| W(A) | |
| R(B) | |
| W(B) | |
| (A) is written | |
| | R(B) |



$$TS(T_1) < TS(T_2)$$

order of serializability: T_1, T_2 .

\therefore possible under TSO.

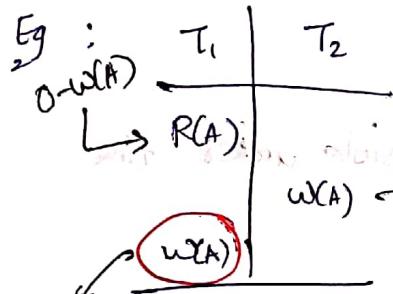
Note:

- TSO ensures CS.
- It ensures freedom from deadlock.
- Aborting & restarting of T_p (with new TS) may cause starvation.

(ii) Thomas write rule:

→ It is a modification of TOP that minimizes the problem of starvation by avoiding the absolute writes.

↳ useless write

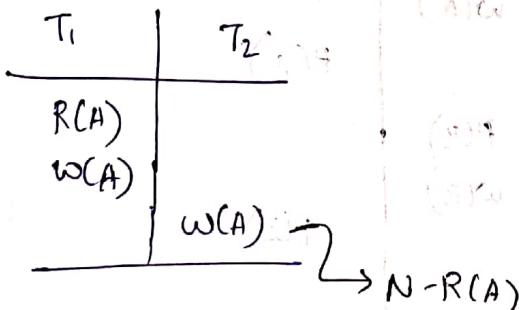


absolute write. $TS(T_1) < TS(T_2)$

~~∴ Out~~ ∵ This should be equivalent to serial schedule T_1, T_2
expected & needed

$w_2(A)$ is already performed.

So we ignore $w_1(A)$.



→ So the above schedule is allowed under Thomas write rule.

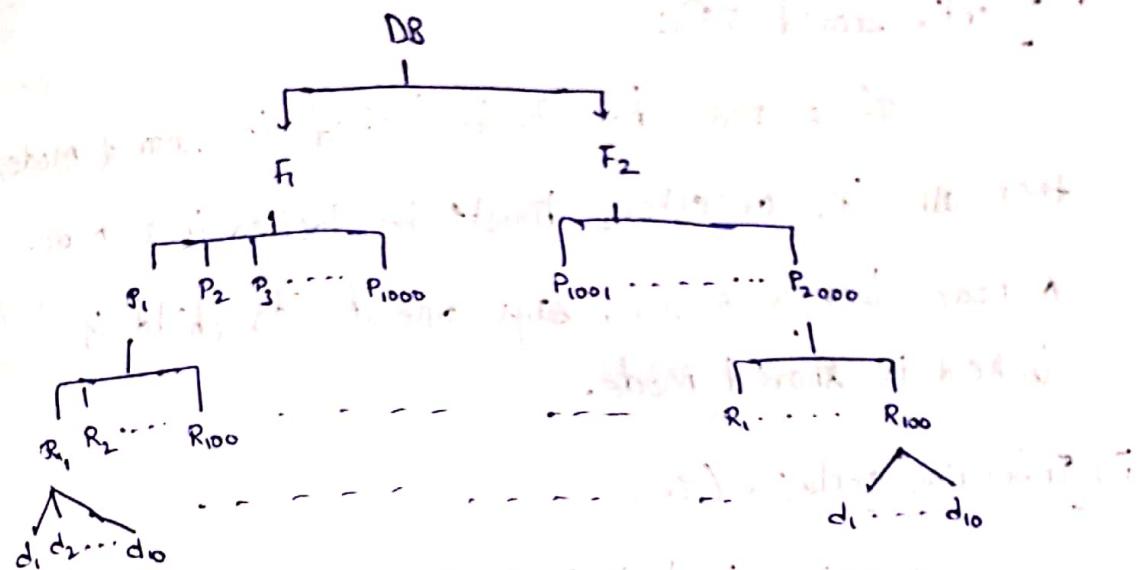
→ T.W.R ensures view serializability by ignoring absolute write.

Note:

→ Every schedule that is possible under TOP is possible under TWR



Multiple Granularity Locking protocol:



→ If a lock is needed on all 10 data items of a record, a lock is requested on record rather than each data item individually and thus it reduces no. of locks to be requested.

→ If P_i is locked in shared mode, it means all the descendants of P_i are locked in shared mode.

~~E~~ Consider 2 transactions

T₁: Read (R_i) of P_i

T₂: update (F_j)

→ If a child node is locked, then its parent cannot be locked in conflicting mode. and vice-versa.

This system uses 3 additional lock modes

(i) Intention shared (IS):

If a node has to be locked in shared mode, then all its ancestors must be locked in IS mode. A node locked in IS says one of its child is locked in shared mode.

(ii) Intention exclusive (IX):

A node is locked in IX says one of its child node is locked in exclusive mode.

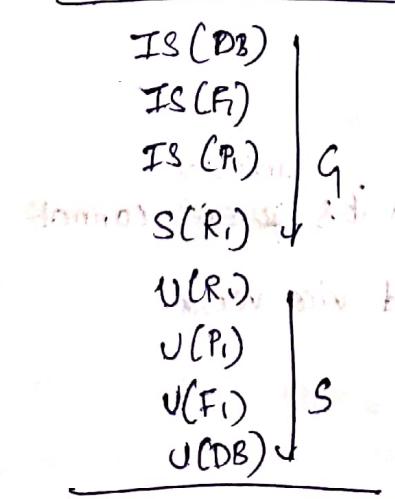
(iii) Shared and Intention Exclusive (SIX):

A node locked in SIX says that the node

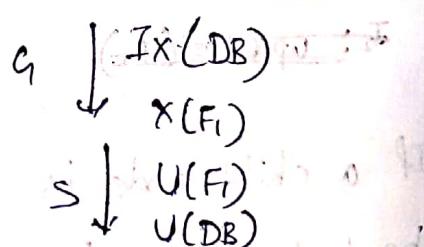
is locked in shared mode and one of its descendants is locked in exclusive mode.

Eg: Consider two transactions

T₁: Read of R₁ of P₁



T₂: Update (F₁)



Multiple granularity 2PL (ensures cs)

- Multiple granularity 2PL ensure CS.
- locking is done from root to leaf.
- Unlocking is done from leaf to root.
- Here deadlock may occur.

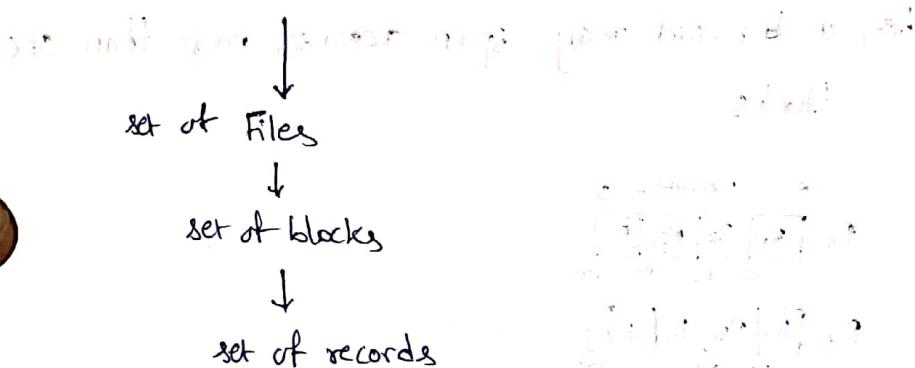
28/11/20

File Organization & Indexing

File Organization:

Database organization starts from bottom to top :-

Data Base



Eg: Consider a file with 1000 records and size of each record is 25 bytes. Let block size be 100 bytes.

$$\text{no of records per block} = \frac{100}{25} = 4$$

$$\text{no of blocks required} = \frac{1000}{4} = 250$$

Blocking factor: It is the avg no of records that can be stored in a block.

$$B_f = \frac{\text{block size}}{\text{record size}}$$

Eg: Consider

Block size = 100 bytes; R = 30 bytes.

$$B_f = \frac{100}{30} = 3.33 \text{ records/block. (spanned strategy).}$$

(or)

= 3 records/block (unspanned strategy).

→ For storing records into memory, we have 2 strategies

i) Spanned strategy:

A part of record is possible to store in a block.

i.e., a record may span across more than one blocks

| | ← 100 bytes → | | | |
|----|---------------|----|----|-----|
| B1 | R1 | R2 | R3 | R4 |
| | 30 | 60 | 90 | 100 |
| B2 | R4 | R5 | R6 | R7 |
| | 20 | 30 | 80 | 100 |
| B3 | R7 | R8 | R9 | R10 |
| | 10 | 40 | 70 | 100 |

Adv: No wastage of memory

Disadv: More no. of block accesses may be needed for one record

Block access: Moving block of records from disk to memory.

→ This technique is suitable for variable length records.

(ii) Unspanned strategy

→ No record is stored across multiple blocks.

| | | | |
|----------------|----------------|----------------|----------------|
| B ₁ | R ₁ | R ₂ | R ₃ |
| | 30 | 60 | 90 |

| | | | |
|----------------|----------------|----------------|----------------|
| B ₂ | R ₄ | R ₅ | R ₆ |
| | 30 | 60 | 90 |

| | | | |
|----------------|----------------|----------------|----------------|
| B ₃ | R ₇ | R ₈ | R ₉ |
| | 30 | 60 | 90 |

adv: Less no of block accesses

disadv: Wastage of memory.

→ Suitable for fixed length records.

→ For storing records into a file, we have two strategies

(i) Ordered file organization

(ii) Unordered file organization

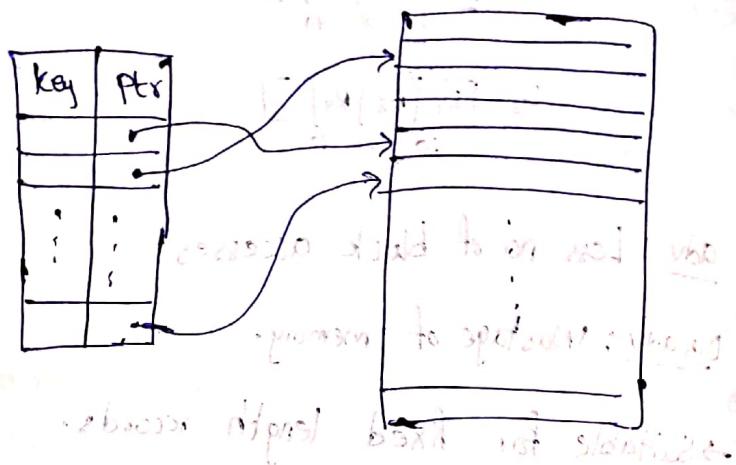
| Ordered file organization | Unordered file organization |
|--|--|
| → Records are stored in the order of some key value. | → Any record can be stored anywhere in the file. |
| → Use binary search | → Uses linear search. |
| → Searching is efficient | → Storing records is easy. |
| → Storing records is difficult | → Searching records is difficult. |

→ Most databases use ordered file organization but searching based on other fields is not efficient.

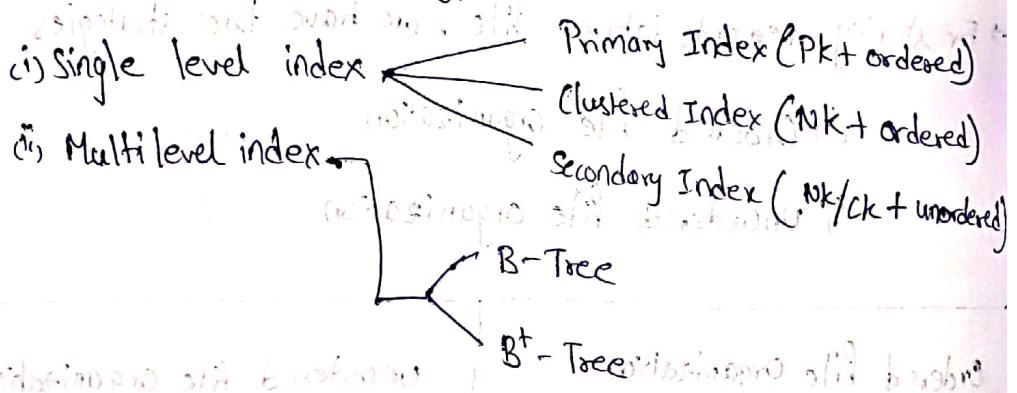
Index: Index is used to speed up the searching process.

→ Index is an ordered file with 2 fields (key & pointer)

key
ptr
Here pointer points to the block that contains the record.



Index is classified into 2 types:

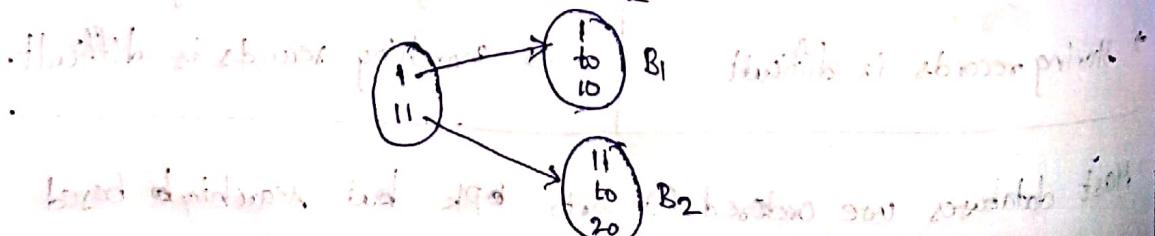


Other classification of index:

(i) Sparse index: No pointers

Dense index: Many pointers

Index record is created only for some key values.



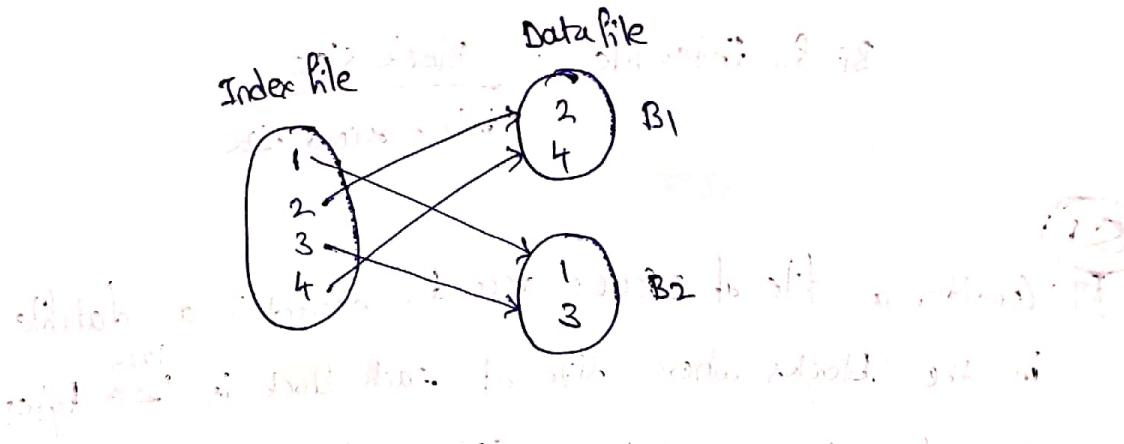
→ This is possible when data file is ordered.

→ no of index records = no of blocks.

→ Here pointer used is block pointer.

(ii) Dense index:

Index record is created for each record of datafile.



→ This is used when datafile is unordered.

→ no of index records = no of records in datafile.

→ Here pointer used is record's pointer.

Single level index:

(i) Primary Index: (Pk + ordered)

→ Primary index is created on pk if the pk is ordered.

w.r.t pk.

→ Type of index: sparse index

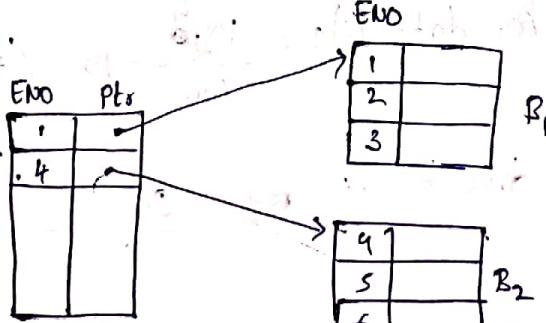
With indexing

→ B_i - no of index blocks

→ Binary search

→ avg no of block accesses

$$= \log_2 B_i + 1$$



Without indexing

→ B - no of datablocks

→ Binary search

→ avg no of block accesses

$$\approx \log_2 B$$

- (ii) Index record is created for each block (sparse index).
 (iii) no of records = no of blocks in datafile

$$BF \text{ for data file} = \frac{\text{Block size}}{\text{data record size}}$$

$$BF \text{ for index file} = \frac{\text{Block size}}{\text{index record size}}$$

Q1

Eg: Consider a file of 300000 records stored in a datafile in the blocks where size of each block is $\frac{1024}{1024}$ bytes.

size of each record is 100 bytes, of which key field is of size 9 bytes and both block and record pointer takes 6 bytes. If we create primary index on the key field of the file, then find the average no of block access to access a record.

Sol:

$$N=30000; B=1000 \text{ bytes}; R=100 \text{ bytes};$$

$$K=9 \text{ bytes}; B_{ptr} = R_{ptr} = 6 \text{ bytes}$$

$$BF \text{ for data file} = \frac{\log N}{R} = 10 \quad (\because \text{fixed length records})$$

$$\text{no of data blocks req} = \frac{30000}{10} = 3000$$

$$\Rightarrow \text{no of index records} = 3000$$

size of each index record = $9 + 6 = 15$ bytes

Blocking factor for index file = $\left\lceil \frac{1000}{15} \right\rceil = 68$

(\because fixed length record)

(iii) no of index blocks req

$$= \left\lceil \frac{3000}{68} \right\rceil = 45$$

(iv) avg no of block accesses = $\lceil \log_2 45 \rceil + 1$

avg no of block accesses without index

$$= \lceil \log_2 3000 \rceil = 12$$

Clustered index: (Nk + ordered) \rightarrow all records

\rightarrow Clustered index is created on non-key. (duplicate values are possible) when data file is ordered wrt Nk

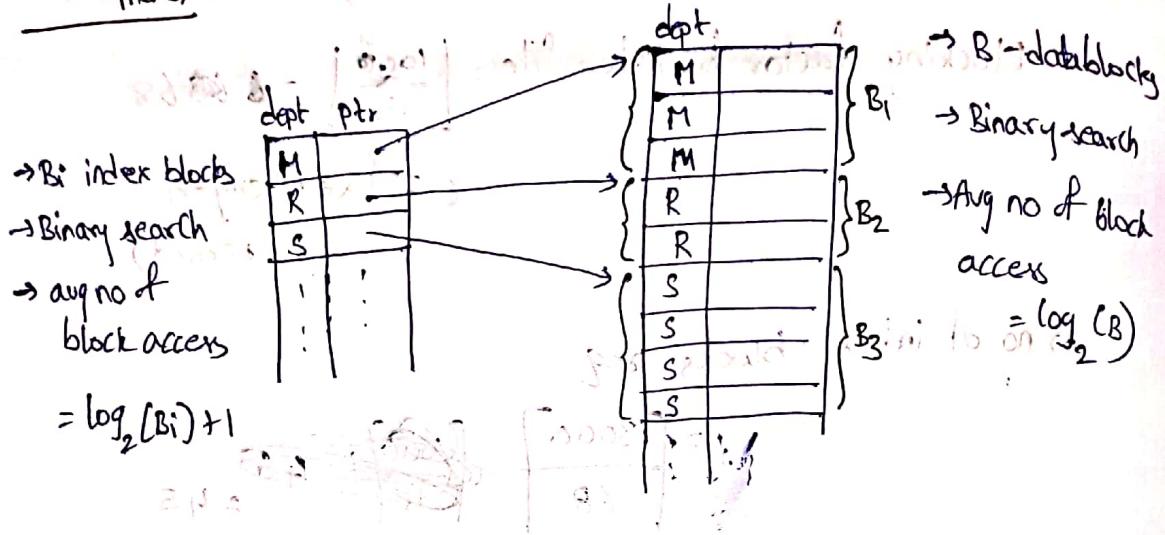
\rightarrow All the records with same Nk value forms a cluster.

\rightarrow Each cluster is assumed as one block.

(i) Index record is created for 1st record of each cluster.

(ii) No of index records = no of clusters
 \therefore Sparse index,

with index



E9: Suppose if we create a clustered index on datafile of Q1

Q2 then find the avg no of block accesses. no of clusters in the data file is 1000. Assume size of nk is 9 bytes.

Sq: ~~Exhibit function, pass into formula~~

$$\text{no of index records} = 1000$$

$$\text{BF for index file} = \left\lceil \frac{1024}{615} \right\rceil = 68$$

$$\text{no of index blocks} = \left\lceil \frac{1000}{68} \right\rceil = 15$$

The avg no of block accesses = $\log_2 15 + 1$
~~Set of new blocks in sorted order (adding 1)~~
 $= 5$

without indexing,

$$\text{no of block accesses} = \log_2(1000) = 10$$

Note:

→ Data block size is variable (depending on cluster size).

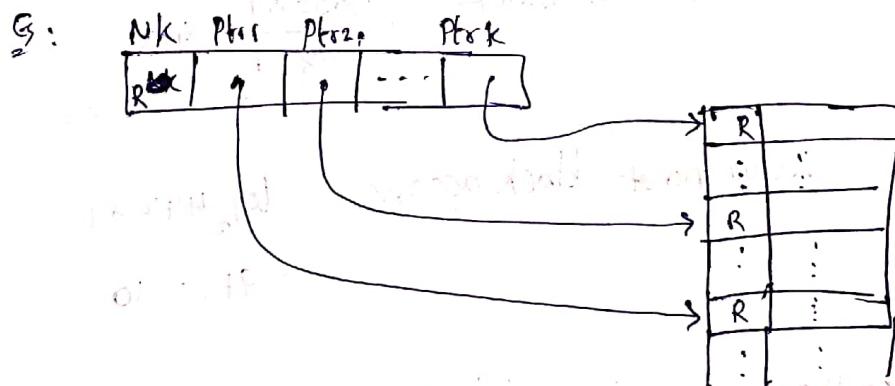
→ Index block size is fixed.

Secondary index (Nk/Ck + unordered)

→ Secondary index is created on Nk/Ck when pk is already present and the file is unordered w.r.t. Nk/Ck.

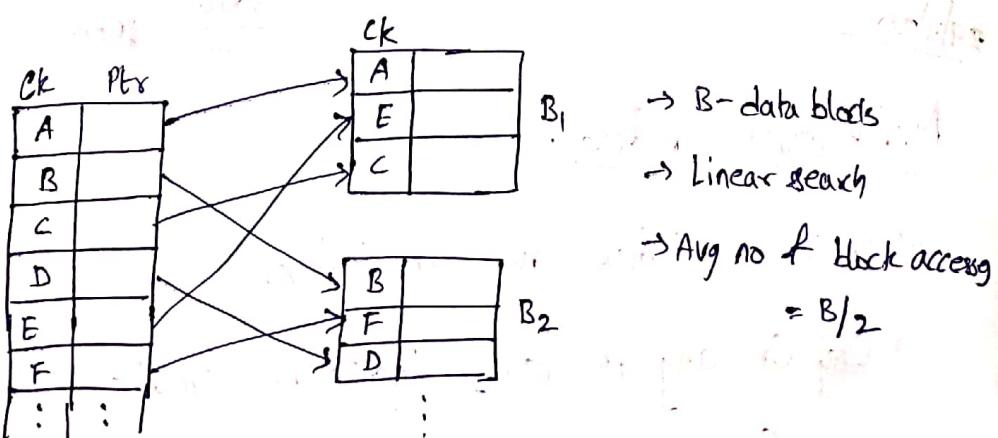
→ Secondary index is possible on Nk, but it is not recommended.

It is because an index record needs to have many pointers and in some cases size of index record could be larger than data record.



→ Thus it is recommended to have secondary index on Ck.

- Bi-index blocks
- Binary search
- avg no of block accesses
- $= \log_2 B_i + 1$



(i) Index record is created for each record (dense index)

(ii) no of index records = no of records in data file

Eg: Suppose if we create a secondary index on the data file,

Q2 then find the avg no of block accesses

Sol:

$$\text{no of index records} = 30000$$

$$\text{size of each index record} = 15$$

$$\text{no of index blocks} = \left\lceil \frac{30000}{15} \right\rceil = 442$$

$$\therefore \text{avg no of block accesses} = \log_2 442 + 1 \\ = 9 + 1 = 10$$

without secondary indexing, no of block accesses

$$\text{no of block accesses} = \frac{\text{no of data blocks}}{2}$$

solution

$$= \frac{3000}{2} = 1500$$

Multi-level Index:

→ It is an index over an index.

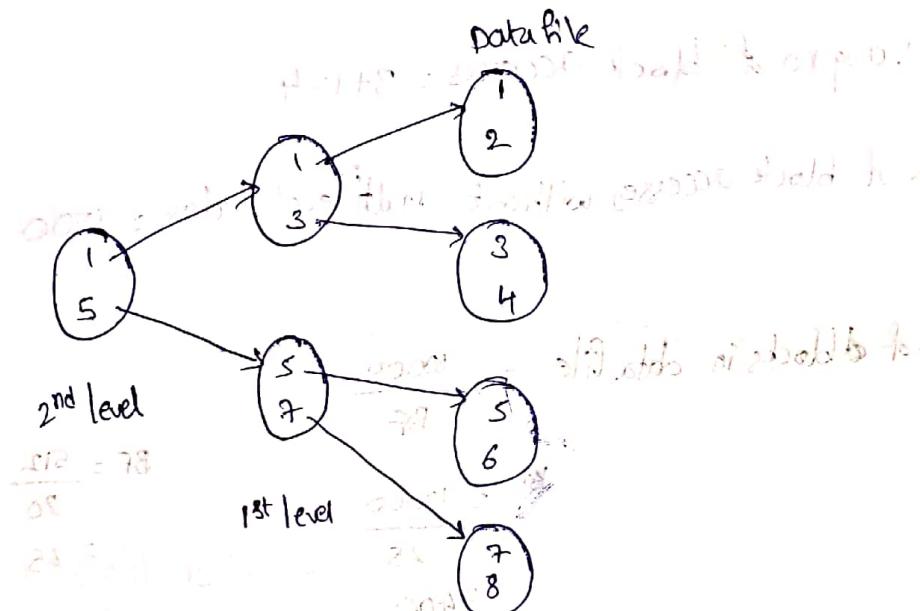
→ On any 1st level index we create a primary index

(∴ Index file is always ordered without any duplicates)

→ we keep on creating multiple indices until the last level index file fits in one block.

→ Here no of index records = no of index records in n^{th} level \leq no of index records in $(n-i)^{th}$ level $(n > i)$

→ The avg no of block accesses = no of levels + 1



→ 1st level may be a sparse index or a dense index.

2nd and higher levels are always sparse indices.

Eg: Q4 Suppose if we convert the secondary index created in the Q3 into multilevel index
Q4 b) find the no of block accesses

Sol:

$$\text{no of index blocks} = 442$$

$$\Rightarrow \text{no of index records in 2nd level} = 442$$

$$\text{size of 2nd level index} = 442 * 15$$

$$\text{no. of index blocks in 2nd level} = \lceil \frac{442}{68} \rceil = 7$$

$$\Rightarrow \text{no. of records in 3rd level} = 7$$

$$\Rightarrow \text{no. of index blocks in 3rd level} = \lceil \frac{7}{68} \rceil = 1$$

\therefore 3-level index. Total no. of block access = 4 + 1 + 1 = 6

$$\therefore \text{avg no. of block accesses} = 3+1=4$$

No. of block accesses without multi-level index = 1500

(P2)

$$\text{no. of blocks in data file} = \frac{10000}{BF}$$

$$BF = \frac{512}{20} \\ = 25 \\ = \frac{10000}{25} \\ = 400$$

$$BF = \frac{512}{20} \\ = 25$$

$$\Rightarrow \text{no. of index records} = 400$$

(P3)

$$n=16384; R=32 \text{ bytes}; k=6 \text{ bytes}; \text{no. of direct blocks} = 16384$$

ordered on Nk; Unspanned;

Block size = 1024; BP = 10 bytes.

Secondary index: 16384 blocks form off part

$$\text{BF}_{(\text{Datafile})} = \frac{1024}{32} = 32$$

$$\text{no. of blocks in datafile} = \frac{16384}{32} = 512$$

$$\text{no. of records in index file} = \frac{\text{no. of records in data file}}{32} \\ = 16384$$

size of index record = $10 + 6 = 16$

$$BF \text{ of index file} = \frac{1024}{16} = 64$$

$$\text{no of index blocks} = \frac{16384}{64} = 2^8 = 256$$

no of records in 2nd level index = 256

$$\text{no of blocks in 2nd level} = \frac{256}{64} = 4$$

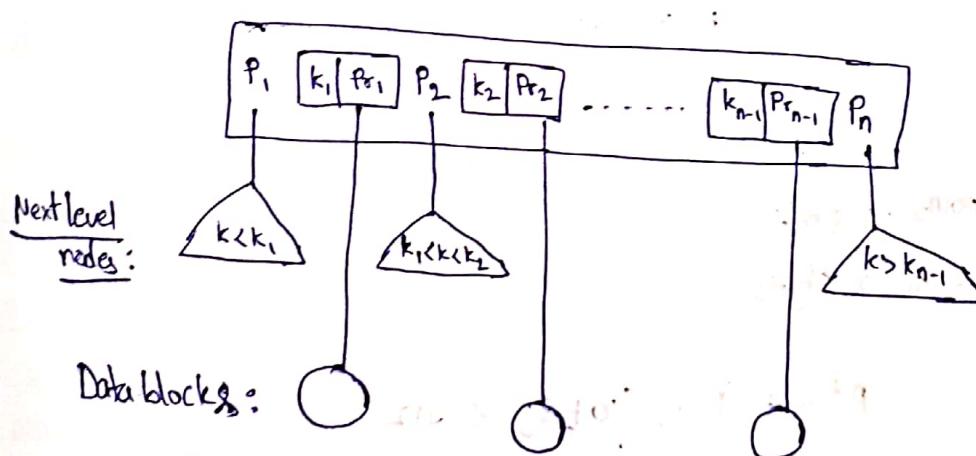
Implementation of Multilevel index:

→ Multilevel index can be implemented using: B-Tree & B⁺-Tree.

B-Tree:

- It is a height balanced search tree.
- Each block in the multilevel index is assumed as one node in B-Tree

Structure of B-Tree node:



B-Tree provides direct access (\because If key found at some level, we can directly move to data block)

In the structure of node

$P_1, P_2, \dots, P_n \rightarrow$ Block pointers

$P_{n+1}, P_{n+2}, \dots, P_{2n-1} \rightarrow$ record / data pointer

$k_1, k_2, \dots, k_{n-1} \rightarrow$ keys

Here order of B-Tree = n

Order of B-Tree:

\rightarrow Order is max no of tree pointers in a node.

If n is order,
no of tree pointers = n

no of index records = $n-1$

(key + record ptrs)

Q

~~if~~ If B is size of block

$$n * p + (n-1) * (k + p_r) \leq B$$

$p \rightarrow$ size of block ptr

$p_r \rightarrow$ size of record ptr

$k \rightarrow$ size of key

$n \rightarrow$ order

(P/4)

data
block pointers \rightarrow 8 bytes

block ptrs \rightarrow 8 bytes

$$5 * 8 + (5-1)(10 + 8) \leq 512$$

$$5P + 18P - 18 \leq 512$$

$$23P \leq 530 \Rightarrow P \leq 23.04$$

$\therefore P = 23$

Insertion:

Note :

A B-Tree of order n must have a minimum of

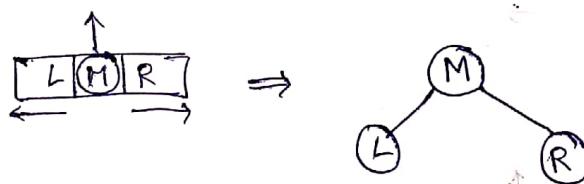
$\lceil \frac{n}{2} \rceil$ pointers

$\lceil \frac{n}{2} \rceil - 1$ keys

But requirement need not to be satisfied by root node.

Insertion:

→ If a node is full then split the node into two nodes and all the key values less than middle key move to left subtree and greater than middle key move to right subtree and middle node to upper level.



Eg: Insert following keys in a B-Tree of order 3

Keys : 20, 40, 30, 15, 45, 5, 60

given order : 3

Max ptrs = 3 min ptrs = $\lceil \frac{3}{2} \rceil = 2$

max nodes = 2 min nodes = $2 - 1 = 1$

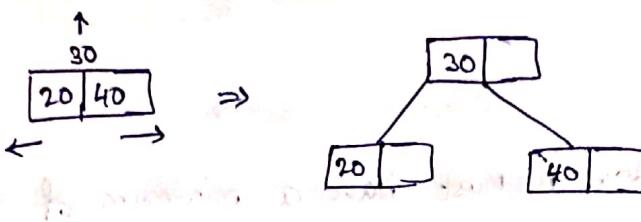
20:

20 → Insert in a new node, balanced & a

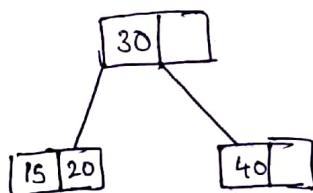
40:

20/40 → no merging yet as both ptns < 1

30:

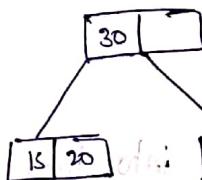


15:



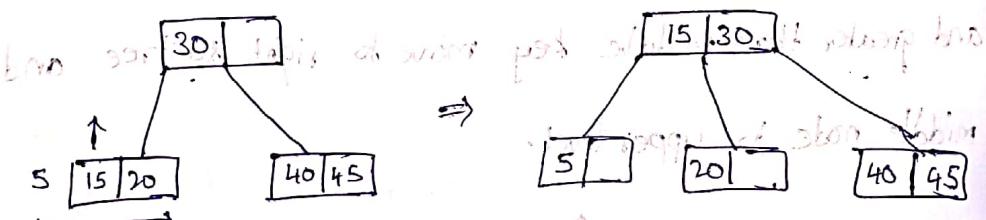
above form of insertion is not for tree formations but

45:

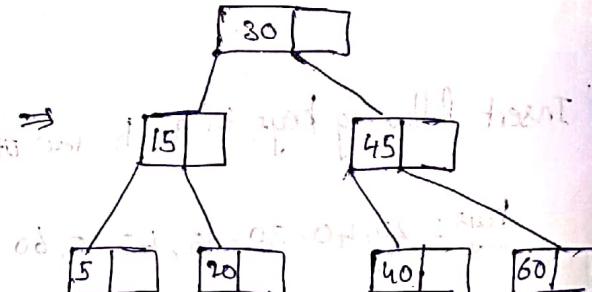
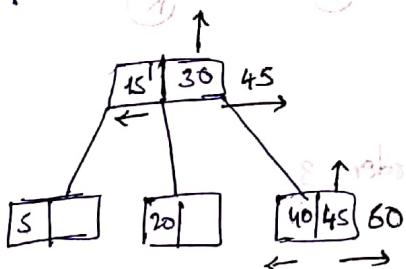


below form of insertion is not for tree formations but

5: first of more part of btree part 2nd is after part 3rd is



60:



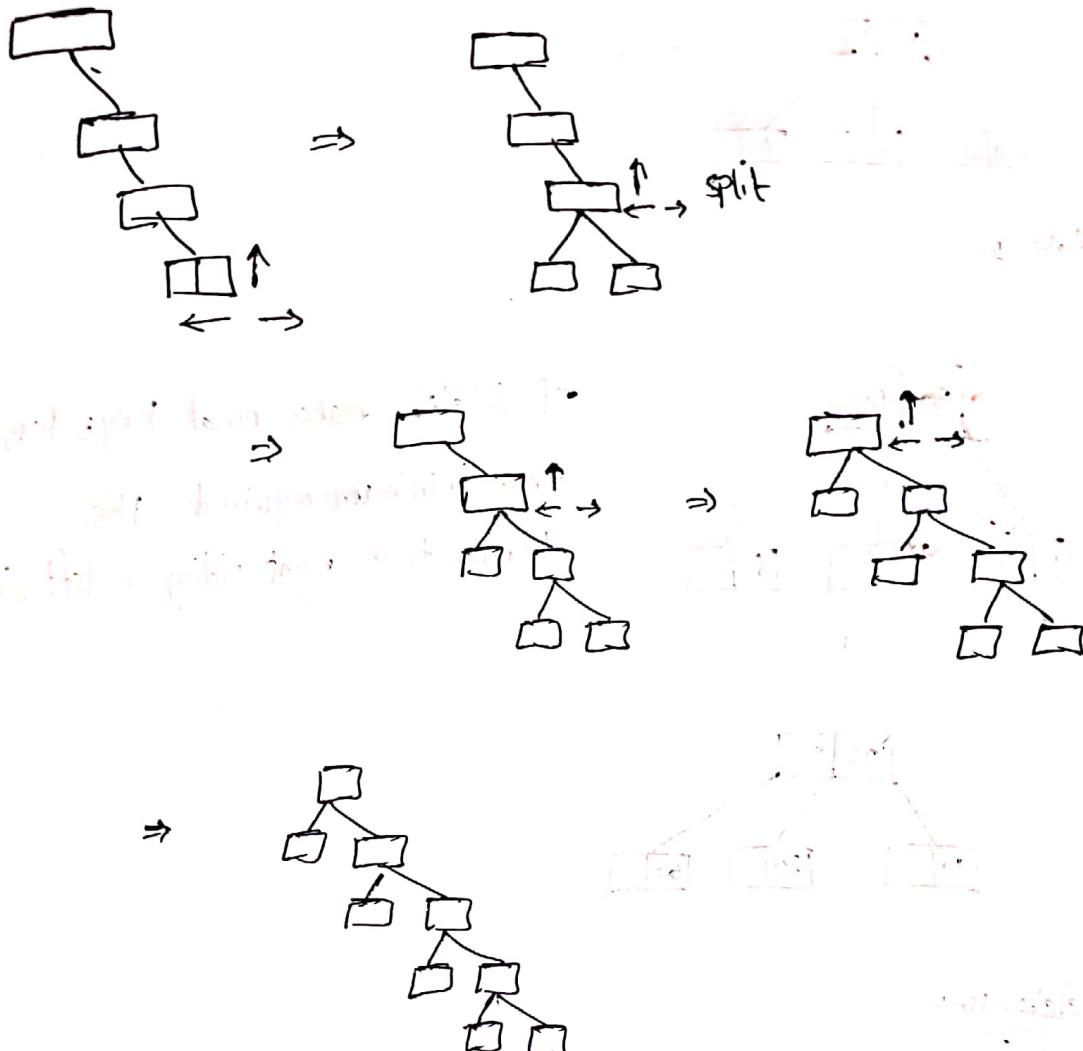
(Q)
G-OS

A B-Tree used as an index for a large database table has

four levels including the ~~no~~ root node. If a new node is

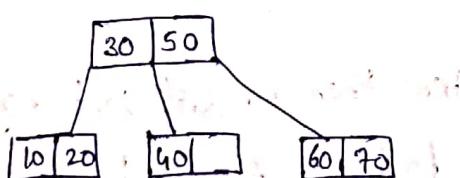
to be inserted, then maximum no of nodes that could be newly created in the process are _____

81:



Deletion:

Consider

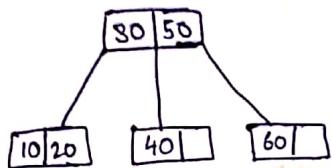


order = 3

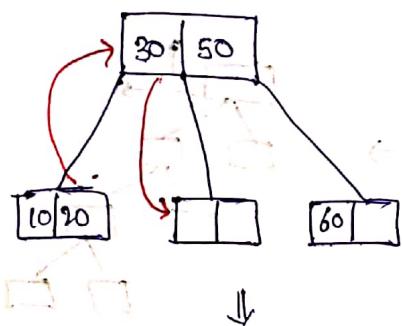
max: 3P, 2k

min: 2P, 1k

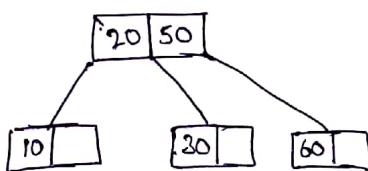
Delete 70:



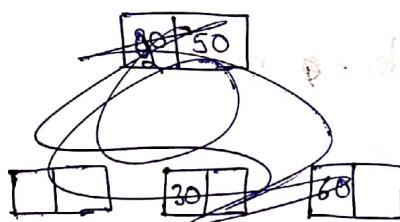
Delete 40:



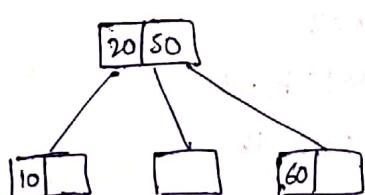
If deletion makes no of keys less than minimum required, then borrow from right sibling or left sibling.



~~Delete 10:~~

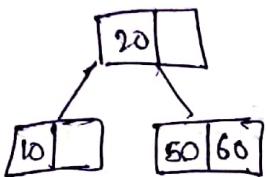


Delete 30:

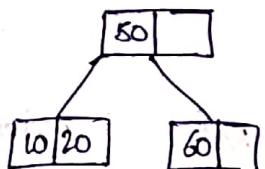


Here there are no sufficient keys to borrow from left or right sibling. Thus we borrow from parent and merge the node from left sibling or right sibling.

Borrowing 50



Borrowing 20



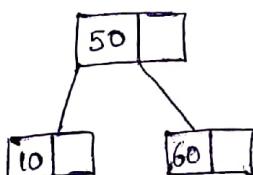
Delete 20:

when an internal node is to be deleted, replace it

either with largest key of LST or smallest ^{key} of RST.

Here LST ^{node} has only one key, so we have to ~~replace~~ ^{key}.

replace 20 with smallest node of RST.



that will be your final answer to given question.

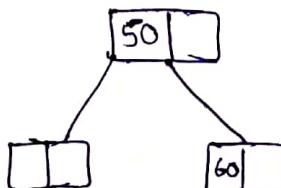
Delete 10:

when 10 is to be deleted, we need to check if borrowing

from left sibling or right sibling is possible. But it is not

possible here. Also borrowing from parent is not possible.

so we merge the nodes

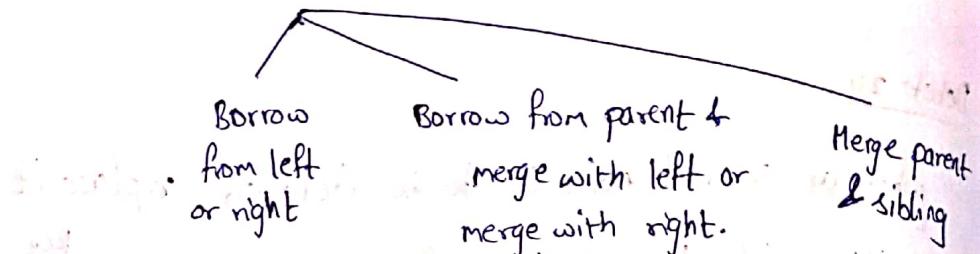


Note:

→ Deletion of leaf

(i) sufficient nodes → no modification

(ii) Insufficient



→ Deletion of internal node

Replace key value with largest of LST or smallest of RST



Note:

→ Leaf nodes of B-Tree doesn't have any children but still we allocate space for tree pointers in leaf node.

This is wastage of space.

→ This drawback is overcome by B+Tree. It is most efficient for doing search operation.



B⁺-Tree:

→ It is a modification of B-Tree

* Leaf nodes ~~don't~~ don't have tree pointers

* Internal nodes don't have record pointers.

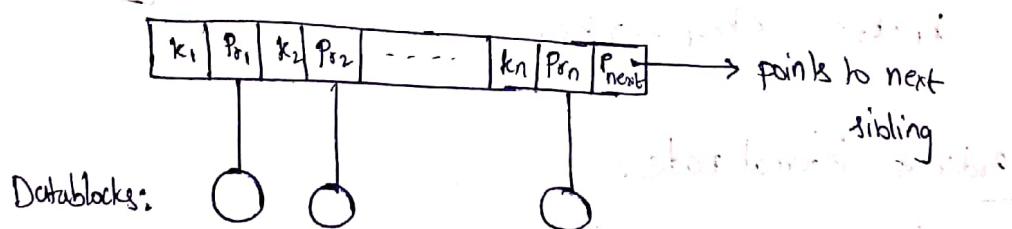
leaf node: (key + record ptr) internal node: (key + tree ptr)

→ B⁺-Tree doesn't provide direct ~~access~~ access.

It provides ordered access.

→ This modification reduces height ~~of~~ and increases width (i.e., no of keys in each node)

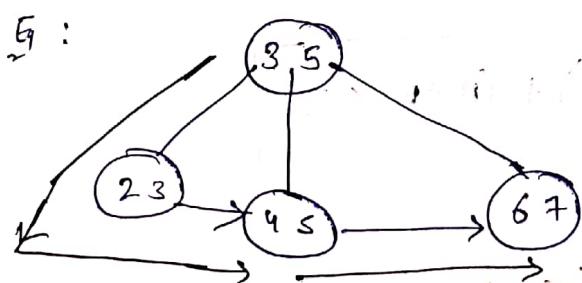
Structure of leaf node:



k_1 to k_n ← keys ; $P_{r1}, P_{r2}, \dots, P_{rn}$ → record / data ptrs

The pointer P_{next} is used for range queries

P_{next} → block ptr
(or)
index ptr
(or)
tree ptr



Query: obtain nodes where $k > 2 \wedge k \leq 6$

no of block access = 4

without P_{next} it would be 6

order of leaf node:

If order is n ,

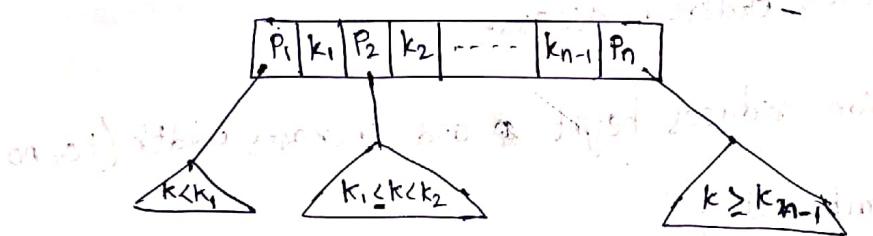
then max no of index records possible = n

If B is size of block then

$$n * (k + p_r) + p_{next} \leq B$$

Order of leaf node may be different from order of internal node

Structure of internal node:



P1, P2, ... Pn → Tree/block/index ptrs

k < k1 < k2 < ... < kn-1 → key

when key is equal, it has to move to left or right. But same convention has to be followed for entire tree

Order of internal node:

If order is n ,

then max no of free ptrs = n

If B is size of block then

$$(n * p) + (n - 1) * k \leq B$$

Q/11

Finding order of internal node:

$$n * (6) + (n - 1)(9) \leq 512$$

$$15n \leq 521$$

$$\Rightarrow n = 34$$

no of ptrs in 0th lvl $\rightarrow 34$

" " " 1st " $\rightarrow (34)^2$

" " " 2nd " $\rightarrow (34)^3$

\Rightarrow no of nodes in 3rd level $= (34)^3$ [leaf nodes]

order of leaf node:

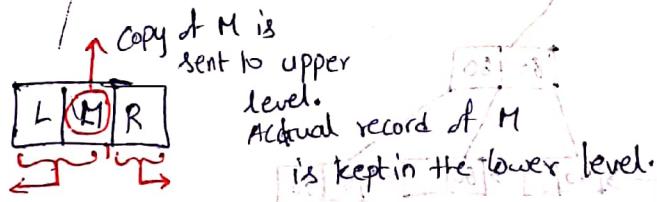
$$n^*(7+9) + 6 \leq 512$$

$$16n \leq 506 \Rightarrow n = 31$$

$$\Rightarrow \text{no of nodes} = (34)^3 \times 31$$

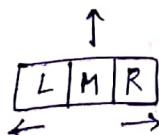
Insertion:

leaf split:



key less than or equal to middle are moved to left subtree. (or) alternate way can also be followed.

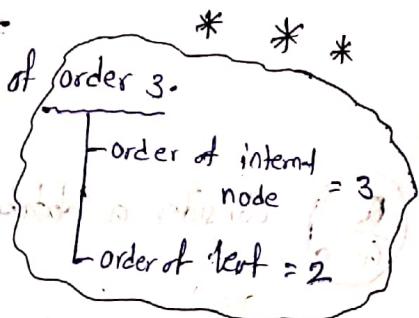
Internal split:



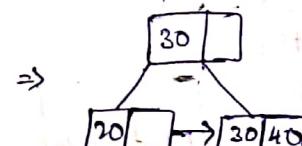
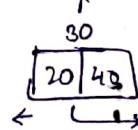
This is same as B-Tree's split.

E: Insert the following keys in a B⁺-Tree of order 3.

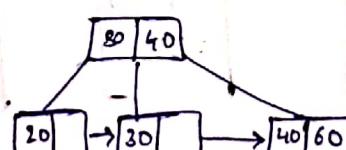
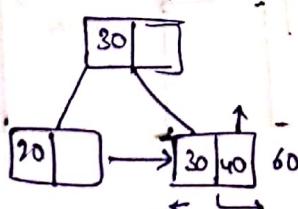
keys: 20, 40, 30, 60, 10, 80, 90

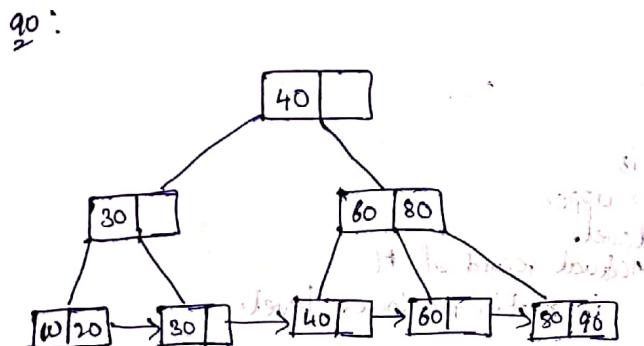
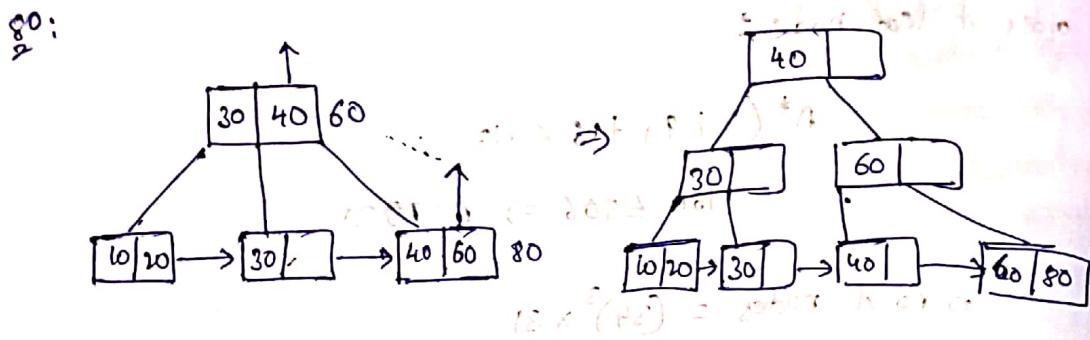
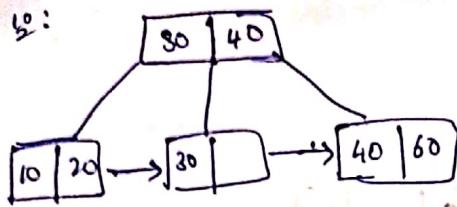


20, 40, 30:



60:





Note: If order of internal node is n ,

$$\text{Min no of ptrs} = \left\lceil \frac{n}{2} \right\rceil$$

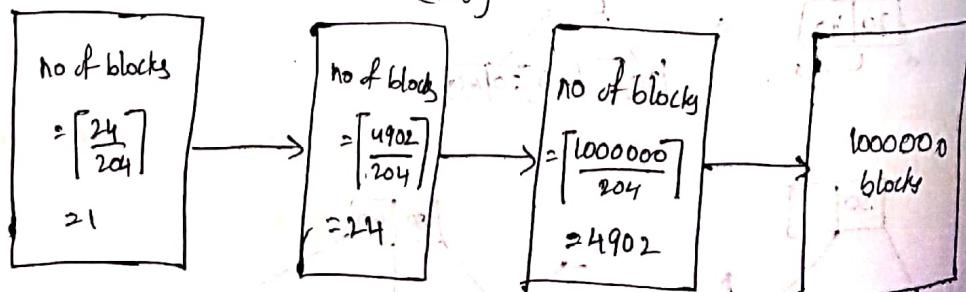
$$\Rightarrow \text{Min no of keys} = \left\lceil \frac{n}{2} \right\rceil - 1$$



Consider a database ...

Q1:

$$B_f = \frac{4096}{(2048)} = 204$$



Method 2:

order of internal node = 205

" " leaf node = 204

| level | block | ptrs |
|-------|-------|------|
| 1 | 205 | 205 |

2 205 $(205)^2$

3 $(205)^2$ $(205)^2 + 204$

no of blocks in leaf level = $(205)^2$

\Rightarrow no of ~~key, record~~ ^{ptr} in leaf level

$$[(205)^2 + 204] \approx 106$$

\therefore 4 block accesses

Deletion:

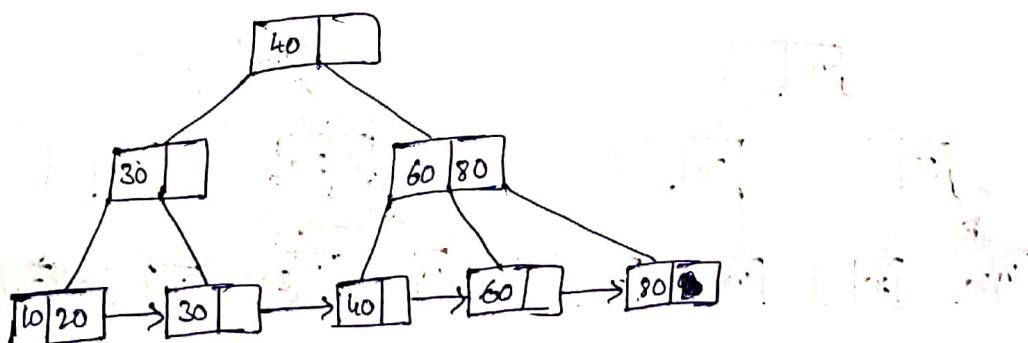
while deletion below points must be checked

i) height balanced

ii) Minimum Order

iii) Properly (i.e., \leq left node, \geq right node)

Consider below BT-Tree

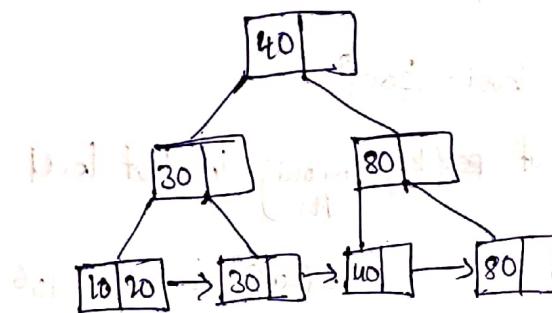


Delete 60:

60 has to be deleted from both leaf & internal node.

Deleting 60 from internal node, leaves the node empty.

so we replace 60 with 80 (we should not replace it with 40)

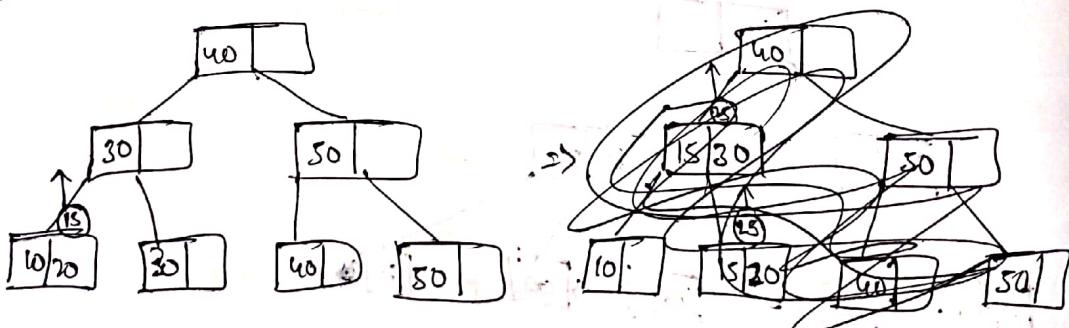


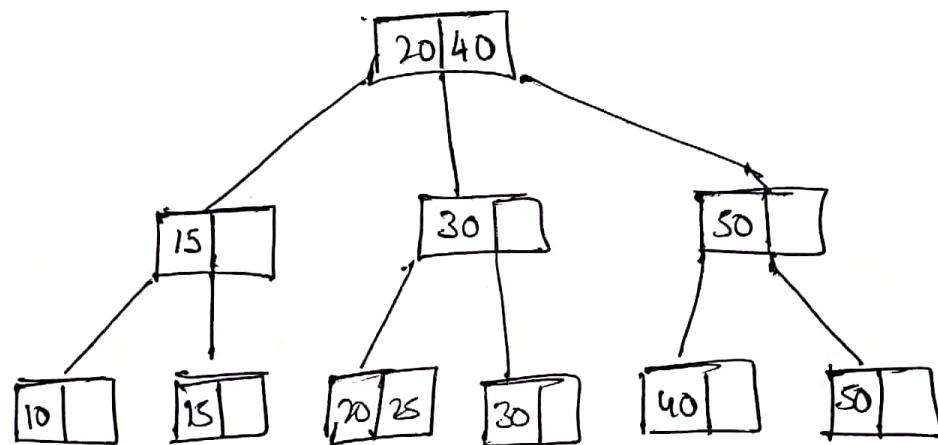
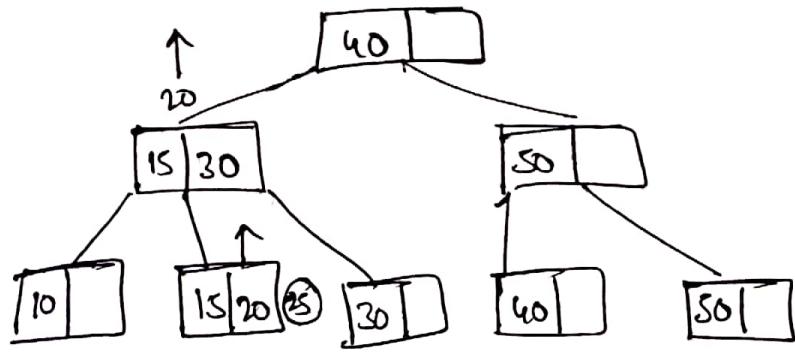
Delete 80:

Deleting 80 causes violation of minimum order in both leaf node and internal node. Also there is no possibility of ~~merging~~ borrowing.



(P/15) 15:





$\therefore 1$
 $\therefore \text{opt}(a)$

P/16

