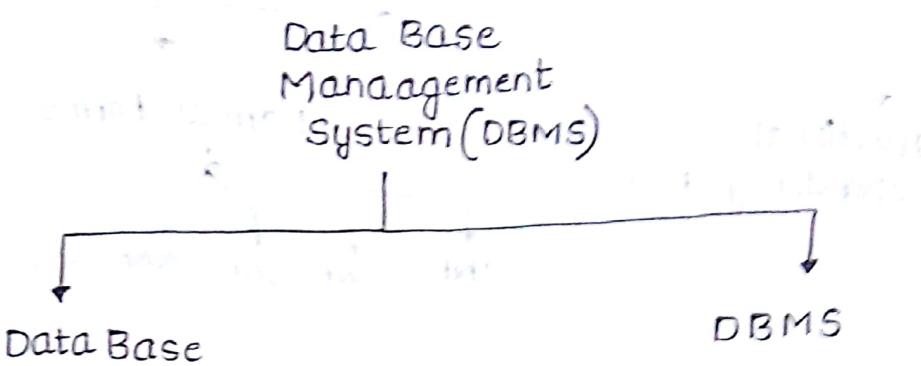


Date: 26/12/2017,
Tuesday.

Introduction



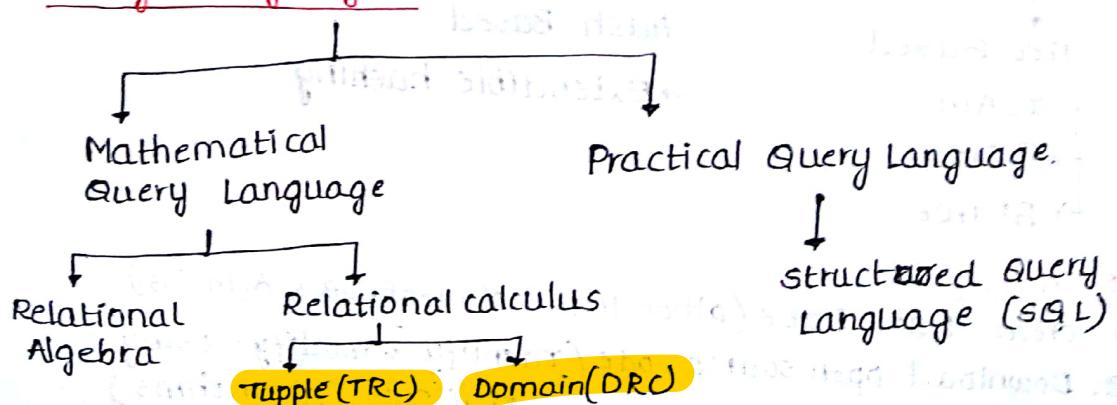
- * A data base is collection of interrelated data

* DBMS is collection of interrelated data and set of application programs which are used to perform some operations on the data
Eg: oracle, DB2, mysql etc

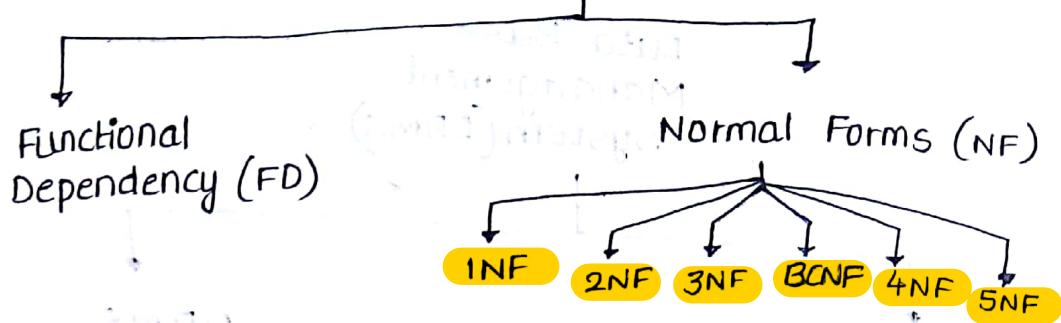
Data Modelling

- * → Entity Relationship Model (pictorial representation)
- * → Relational Model (All DBMS are RDBMS. commercial.)
- Network Model
- Hierarchical Model
- Object oriented Model.

Query Languages



⇒ After designing DBMS, we conduct a testing process called "Normalization" (It's a DBMS theory)



⇒ Normalization is used to design efficient Data Base.

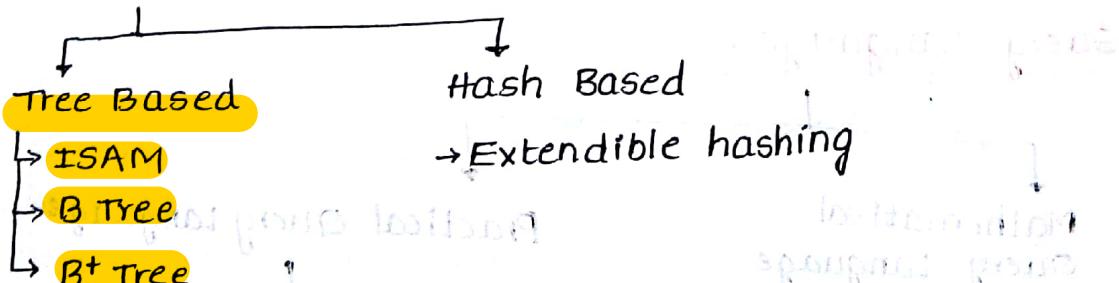
Transaction Management (TM)

- Transaction processing (TP)
- Transaction Management (TM)
- Transaction Failure (TF)
- Concurrent transactions (CT)

Query Processor

- Query processing (Q.P)
- Query Evaluating (QE)
- Query Optimisation (QO)

Indexing



* Projects

1. Create Data Base (other than student DB & ATM DB)
2. Download open source code (Identify & modify using efficient Algorithms)

File systems

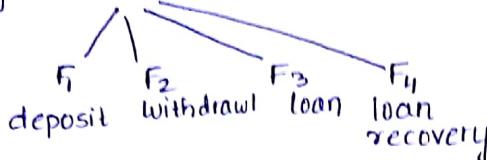
vs

DBMS

- ⇒ Different files are meant for different purposes

⇒ **No Data Redundancy***

Eg:- Bank



$F_i \rightarrow$ data + Applications

Drawback: **customer information is repeating.**

⇒ **No Data Inconsistency**

There is "**Data Redundancy**"

→ space complexity

→ data inconsistency

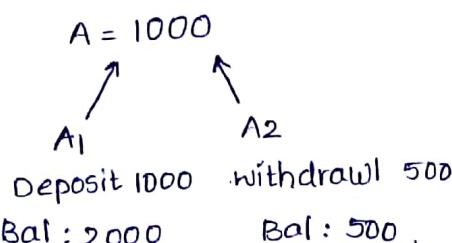
Eg: data must be updated in all file systems.

⇒ Ex:- if we want, the customers info whose balance is above 1,00,000. This is a new application. It is difficult to develop a new application with the existing applications and integrating them.

⇒ No problem of Integrating new application with existing application.

⇒ **Easy Data Access.**

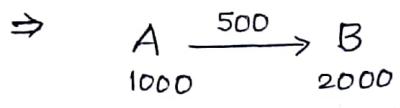
⇒



But none is true.

⇒ **No concurrent execution mismatch problem.**

∴ concurrent execution problems.



- $F_1 \left\{ \begin{array}{l} 1) \text{Read } A \\ 2) A.\text{bal} = A.\text{bal} - 500; \\ 3) \text{write } A \end{array} \right.$
- $F_2 \left\{ \begin{array}{l} 4) \text{Read } B \\ 5) B.\text{bal} = B.\text{bal} + 500 \\ 6) \text{write } B \end{array} \right.$

Due to some bugs, OS problems
if the transaction is
failed after step 3, there
is no recovery mechanism
in file systems.

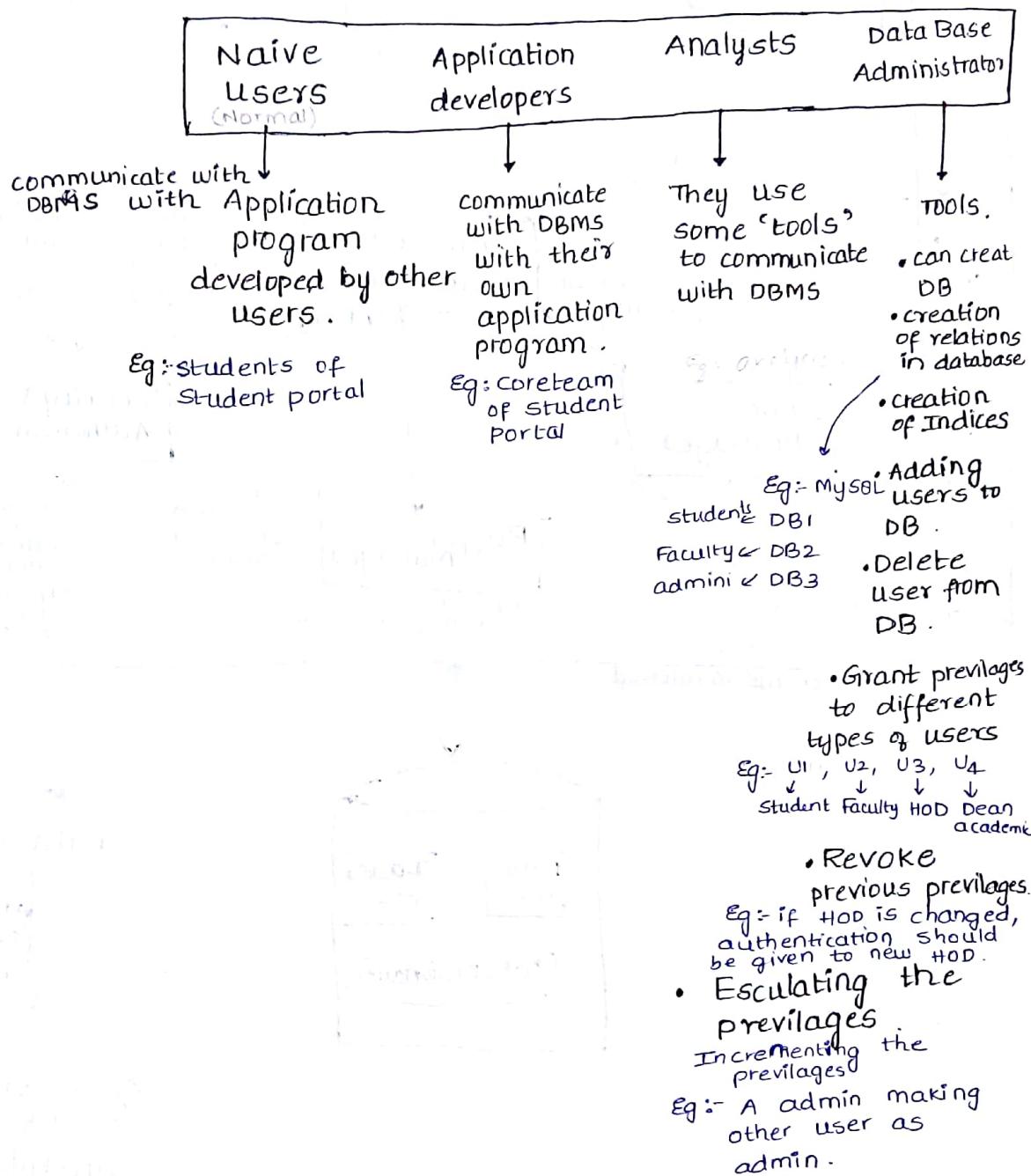
\Rightarrow There is a data
recovery mechanism
in DBMS, during
failures.

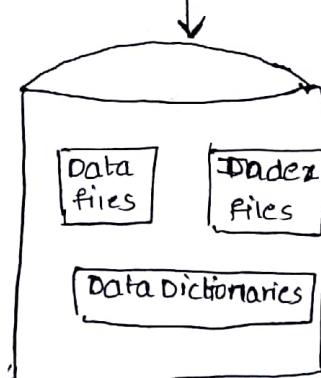
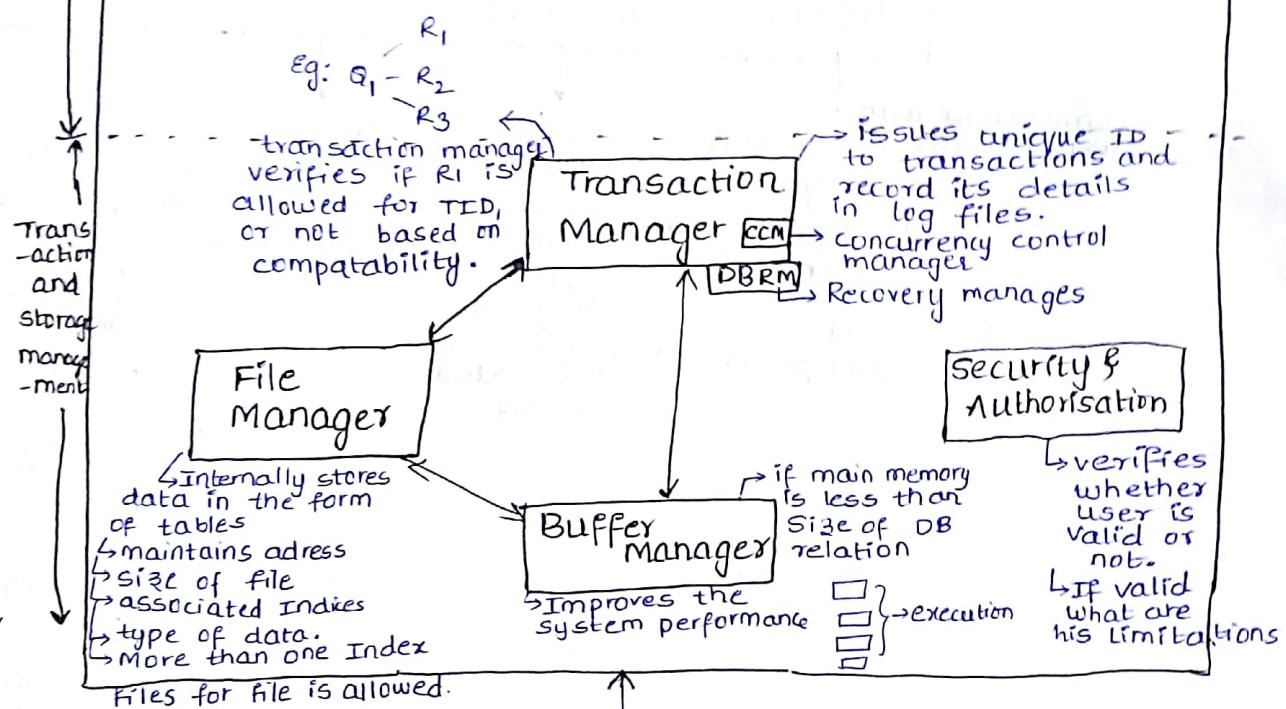
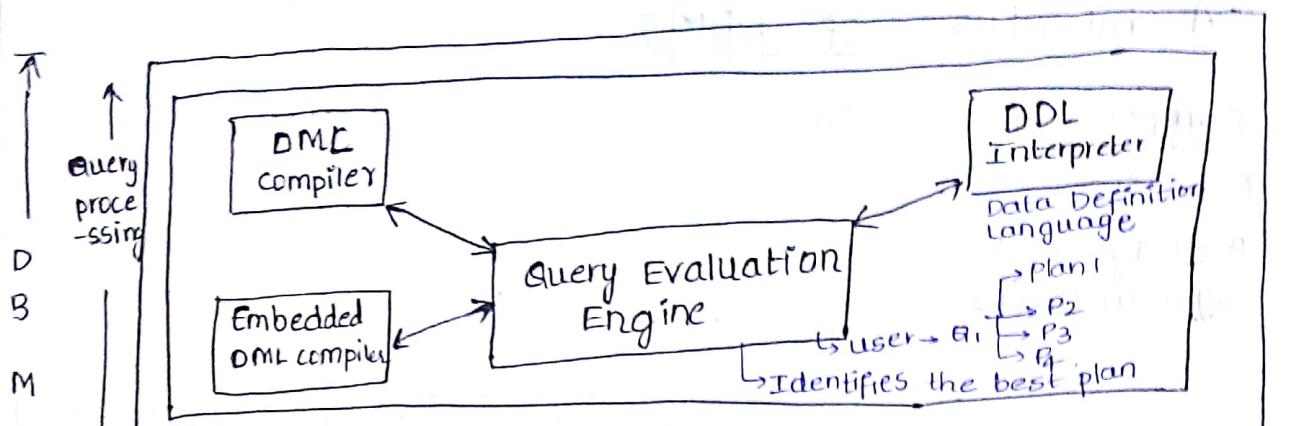
Architecture of DBMS

Components of DBMS

Data Base

Management system users





Data dictionary

Provides statistical information

- metaData
- Data about Data
- Index on each data.

Eg: of statistical information

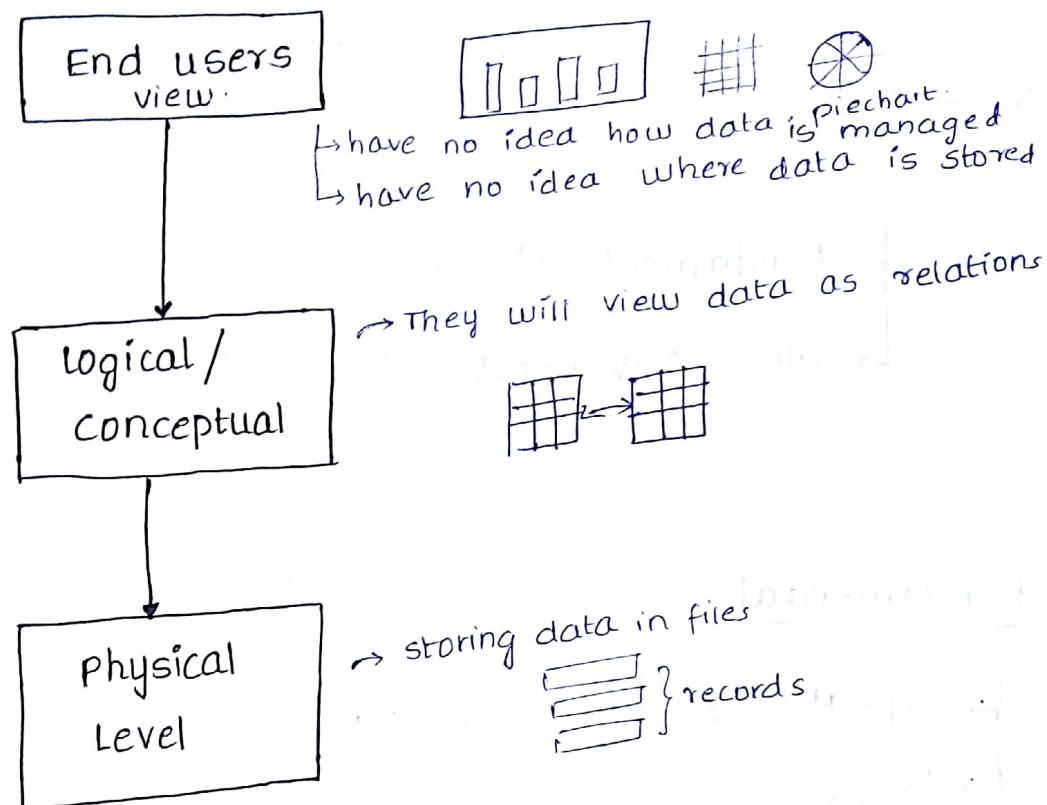
attribute, domain, size etc.

Statistical data present in data dictionaries

→ which helps the Query engine to prepare plans.

Data Abstraction

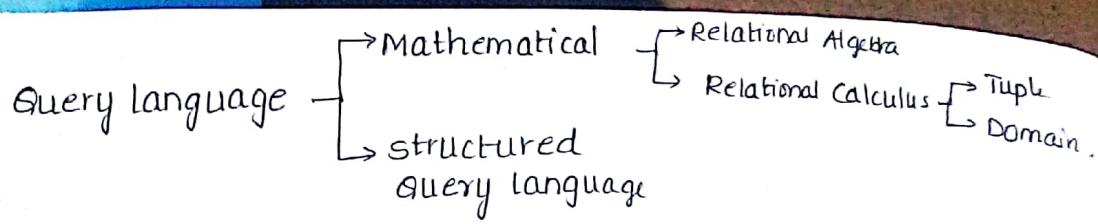
3-Level schema:



Advantage of 3 level schema

⇒ Change of one level does not effect other levels.
This is called data Independency?

⇒ Hiding the complexities at different Level



Relational Algebra

- Fundamental / Basic
- other / Advanced. (we can use fundamental operations to develop new operations)

Fundamental

- select (σ)
- Project (π)
- cartesian product (\times)
- set union (U)
- set Difference (-)
- Rename (ρ)

Advanced

- Join
 - NJ (Natural)
 - TJ (Theta)
 - outer join
 - Left
 - Right
 - full
- set instruction.
- Database Modification
 - Insert
 - delete
 - update
- Aggregate operations

Tuple:- A horizontal row in a relation.

Select (σ) :-

It divides horizontally, the given relation.

✓
✗
✓
✓
✗
✓

→ only from main memory not from DB
it deletes those tuple which doesn't satisfy given condition and outputs the remaining

$\sigma_{\theta}(\text{Rel})$

↓
original relation in D/B
(or) Result of a R.A expression.

Simple
Compound.

(condition) $\theta = \theta_1 \wedge \theta_2 \wedge \theta_3$

$\theta = \theta_1 \vee \theta_2$

Eg:- $\sigma_{m > 50 \wedge d = 'CSE'}(\text{student})$

its equivalent ^{structured} query

```
SELECT *
FROM student
WHERE m > 50
AND d = 'CSE';
```

Project (Π) :-

Attribute: A column in a relation.

✓	✗	✓	✓
A ₁	A ₂	A ₃	A ₄

Π_{A_1, A_3, A_4} → It projects only those attributes mentioned.

$\Pi_{SNO}(\text{student})$

Selection followed by Projection :-

$$\Pi_{SNO} (\sigma_{m > 50} (\text{student}))$$

→ outputs only student roll no whose marks > 50

from 'student' relation

SELECT SNO FROM student WHERE m > 50;

$$t_1 \in r_1$$

$$r_1 = \sigma_\theta (R)$$

$$r_2 = \Pi_{A_1 A_2 A_3} (R)$$

no. of tuples in r_1

min = 0
max

} depends on number of condition satisfied tuples

no. of tuples in $r_2 = \text{max}$

no. of attributes in $r_1 = \text{max}$

no. of attributes in $r_2 = 3$

Projection followed by selection.

$$\Pi_{SNO, SNAME} (\sigma_{m > 50} (\text{student})) \equiv \Pi_{SNO, SNAME} \left(\sigma_{m > 50} \left(\Pi_{SNO, SNAME} (\text{student}) \right) \right)$$

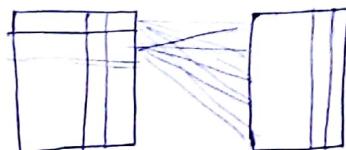
* Note:-

$$\Pi_{A_1, A_2} (\sigma_\theta (R)) \equiv \Pi_{A_1, A_2} \left(\sigma_\theta \left(\Pi_{A_1, A_2 \cup \text{Attr}} (R) \right) \right)$$

Cartesian product :-

$$DB = \{R_1, R_2, \dots, R_n\}$$

When we require some portion from R_1 & some from R_2 first combine both R_1, R_2 using cartesian product.



$$\gamma = R_1 \times R_2 \\ (m) \quad (n)$$

- no. of tuples in $\gamma = m * n$

- $t_1 \in \gamma$
- size of t_1 is $= \underbrace{\text{size } T_1 \in R_1}_{\text{Attributes}} + \text{size of } T_2 \in R_2$

Eg:1 $R_1(ABC)$

$R_2(DE)$

$$R_1 \times R_2 = \begin{array}{|c|c|c|c|c|} \hline A & B & C & D & E \\ \hline & & & & \\ \hline \end{array}$$

Eg:2 $R_1(ABC)$

$R_2(CDE)$

$$R_1 \times R_2 = \begin{array}{|c|c|c|c|c|} \hline A & B & C & C & D & E \\ \hline & & & & & \\ \hline \end{array}$$

Set operations:-

Union :

$$S_1 = \{1, 2, 3, 4\}$$

$$S_2 = \{1, 3, 5\}$$

$$S_1 \cup S_2 = \{1, 2, 3, 4, 5\}$$

→ eliminates duplicates
(Internal Algorithm)
→ if we need duplicates
perform union all.

Pre-conditions to perform set operations

- 1) Arity (no. of attributes) should match.
- 2) Domain value of i th attribute in R_1 should match with domain value of i th attribute in R_2 .

$$R_1 = \{ A_1, A_2, A_3 \}$$

$$R_2 = \{ X_1, X_2, X_3 \}$$

s.no s.no ✓
s.no Id ✗
s.name s.name ✓

Set difference (-) :-

$$S_1 = \{ 1, 2, 3, 4 \}$$

$$S_2 = \{ 1, 3, 5 \}$$

$$S_1 - S_2 = \{ 2, 4 \} \quad \text{elements in } S_1 \text{ but not in } S_2$$

Rename (P) :

$$P \underset{\text{newname}}{\underset{\text{oldname}}{\wedge}}$$

$$P \underset{(R)}{\wedge} (A_1, A_2, A_3)$$

$R \rightarrow$ original relation in DB
 \rightarrow result of relational algebra expression.

$$\text{Eg:- } P (\pi_{\text{sno, grad}} (\text{student}))$$

Intersection (\cap) :

$$S_1 = \{ 1, 2, 3, 4 \}$$

$$S_2 = \{ 1, 3, 5 \}$$

$$S_1 \cap S_2 = \{ 1, 3 \}$$

$$S_1 \cap S_2 = S_1 - (S_1 - S_2)$$

join (\bowtie):-

$$R_1 \bowtie_{JC} R_2 = \sigma_{\theta}(R_1 \times R_2)$$

Eg:- $R_1 = (A, B, C)$ +
 $R_2 = (C, D, E)$

$$R_1 \bowtie_{R_1.C=R_2.C} R_2 = \sigma_{R_1.C=R_2.C}(R_1 \times R_2)$$

$$R_1 = (SNO, Sname, Add)$$

$$R_2 = (SNO, marks, grade)$$

$$\prod_{Sname, marks} \left(R_1 \bowtie_{R_1.SNO=R_2.SNO} R_2 \right) = \prod_{Sname, marks} \left(\sigma_{R_1.SNO=R_2.SNO}(R_1 \times R_2) \right)$$

$$\tau = R_1 \bowtie_{JC} R_2$$

No. of attributes in τ = No. of attributes in R_1 + No. of attributes in R_2 .

No. of tuples in τ \rightarrow min = 0.
 \rightarrow max = $\min(R_1, R_2)$

\bowtie \rightarrow Natural join (includes only one attribute of, duplicate attributes)
 \rightarrow Theta join

Equivalent query in SQL

SELECT Sname, marks

FROM R_1, R_2

WHERE $R_1.SNO = R_2.SNO$;

(Or)

R_1 natural join R_2

Aggregate Operations

- Min
- Max
- Avg
- Sum
- Count

* group by Σ (Rel)
Attribute Aggregate operation

dept Σ (emp)
 $\max(sal)$

* Generalised Projection.

$\Pi_{A_1, F_1 A_2 F_2} (R)$

eg:- $\Pi_{SNO, m+5} \left(\sigma_{m < 30} (\text{student}) \right)$

Adding 5 marks to students having
marks less than 30.

101

Database Modification

- ① Insert
 - ② Delete
 - ③ Update.
- } advanced
↓
using basic operations.

Aggregate Operations

- Min
- Max
- Avg
- Sum
- Count

* group by Σ (Rel)
Attribute Aggregate operation

dept Σ (emp)
 $\max(sal)$

* Generalised Projection.

$\Pi_{A_1 F_1 A_2 F_2} (R)$

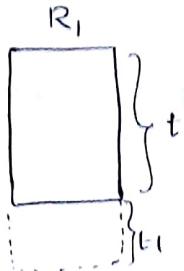
Eq: $\Pi_{SNO, m+5} \left(\sigma_{m < 30} (\text{student}) \right)$

Adding 5 marks to students having
marks less than 30.

Database Modification

- ① Insert
 - ② Delete
 - ③ Update.
- } advanced
↓
using basic operations.

Insert :-



(union)

$$= t + b_1$$

$$R_1 \leftarrow R_1 \cup t_1$$

Delete :-

- (1) Search (set difference)
- (2) Delete

$$\text{initially } R = t.$$

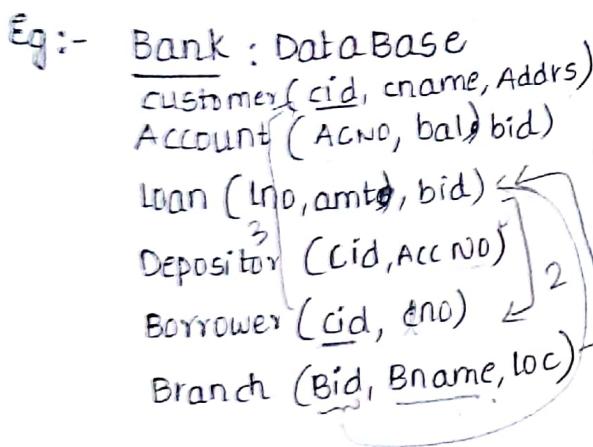
$$\text{after search } t_2$$

$$\text{after delete } t - t_2$$

$$\begin{aligned} & t_1 \leftarrow \sigma_{\theta}(R_1) \\ & R_1 \leftarrow R_1 - t_1 \end{aligned} \quad \left. \begin{array}{l} \text{deletion.} \\ \text{insertion.} \end{array} \right\}$$

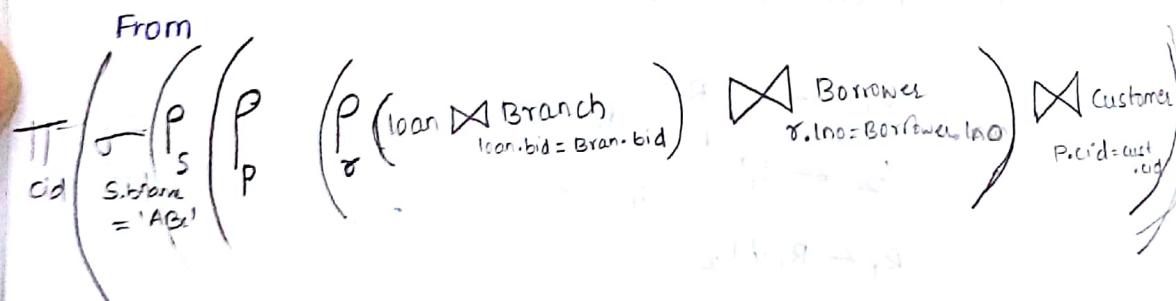
Update : delete & insert

$$\begin{aligned} & t_1 \leftarrow \sigma_{\theta}(R_1) \\ & R_1 \leftarrow R_1 - t_1 \\ & t_2 \leftarrow t_1' \\ & R_1 \leftarrow R_1 \cup t_2 \end{aligned} \quad \left. \begin{array}{l} \text{delete} \\ \text{insert} \end{array} \right\}$$



1. Find all customer id's those who are having an account at 'ABC' branch.
2. Find all customer names those who are having an account at 'ABC' branch.
3. Find all customer name those who are having an account and a loan in the bank.

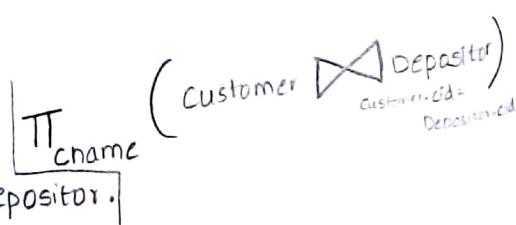
1. Select cid



4. Find all customers names those who are having an account in the bank.
5. similarly a loan in the bank.

4) Select distinct cname

From customer natural join Depositor.



5) $\Pi_{cname} \left(\left(\text{customer} \bowtie_{customer.cid = Borrower.cid} \text{Borrower} \right) \right)$

Select distinct cname
from customer natural join
Borrower;

$$③ \left(\prod_{cname} (\text{customer} \bowtie \text{Depositor}) \text{ on } \text{cust_id} = \text{depos_id} \right) \cup \prod_{cname} (\text{customer} \bowtie \text{Borrower}) \text{ on } \text{co_cid} = \text{b_cid}$$

- Natural join (⊗) Theta join retrieves the data that is common in both the relations and loss of the data. So, to overcome this we use outer join.

Eg:	R_1	R_2	$R_1 \text{ join } R_2$ (Theta)
	A B C	D E	
	1 b ₁ c ₁	1 e ₁	
	2 b ₂ c ₂	3 e ₃	
	3 b ₃ c ₃	5 e ₅	
	4 b ₄ c ₄		
	5 b ₅ c ₅		

Outer join :- It is an extension of natural join.

It join all tuples from R_1 OR R_2 OR BOTH without missing any tuples.

outerjoin

↳ left outer join :  $R_1 \Join R_2$
all tuples from R_1

	R_1			R_2	
	A	B	C	D	E
1	b ₁	c ₁		1	e ₁
2	b ₂	c ₂		3	e ₃
3	b ₃	c ₃		5	e ₅
4	b ₄	c ₄			

$$\delta = R_1 \Join R_2$$

	A	B	C	D	E
1	b ₁	c ₁		1	e ₁
3	b ₃	c ₃		3	e ₃
2	b ₂	c ₂		null	null
4	b ₄	c ₄		null	null

Right outer join : 

$$\delta = R_1 \Join R_2$$

	A	B	C	D	E
1	b ₁	c ₁		1	e ₁
3	b ₃	c ₃		3	e ₃
	null	null	null	5	e ₅

Full outer join : \bowtie

$$f = R_1 \bowtie R_2$$

	A	B	C	D	E	F
1		b ₁	c ₁		e ₁	
null		• null	• null	5	e ₅	
3	b ₃	c ₃		3	e ₃	
2	b ₂	c ₂		• null	• null	
4	b ₄	c ₄		• null	• null	

$t_1 \in \alpha$ \rightarrow no. of tuples = no. of tuples in left side relation

$t_2 \in \gamma$ \rightarrow no. of tuples = no. of tuples in right side relation

$t_3 \in f$ \rightarrow no. of tuples = no. of tuples in left side + no. of tuples in right side

- no. of tuples ($R_1 \bowtie R_2$)
natural join

VIEWS

↓
SQL

create view view-name as <RA expression>

Eg: Create view V₁ as $\pi_{A_1, A_2, A_3} (\sigma_0 (R_1 \Delta R_2))$

Database Design life cycle

SDLC
(subway design life cycle)

- Problem definition
- Analysis of Problem
type complexity possible solutions

- Design
 - Algorithm
 - Flow Diagram

DBDLC
(database development for life cycle)

- Define the problems for DB

- Problem Analysts
 - Identifying the entities (obj) → (relation)
 - Identifying the properties of an entities (attributes)
 - Identify the relationship exists b/w entities
 - type of data, user, facilities, system (app)

- Design → pictorial representation of DB

Modelling
 ER model
 NO DBMS to implement it. so convert relational model (tabular)
 All, current DBMS's are RDBMS

(R_1, R_2, \dots, R_n)
 in acceptable forms

-1NF
 -2NF
 -3NF
 -BCNF
 Normalization

Testing these anomalies

→ Size of relations (no data redundancy)
 → Insert, Delete, Update } problems

- Implementation

- Testing

- Implementation
(i.e., creation of tables
using any DBMS)

- Testing
(insertion,
deletion,
retrieval).

22/11/18

RELATIONAL CALCULUS

⇒ Non procedural query language whereas
Relational algebra and SQL are procedural query languages

Procedural QL

↳ To perform operation on data.
we write procedure.

Relational calculus

- ↳ non procedural QL
- ↳ TRC → based on tuples.
- ↳ DRC → based on domain value.

TRC :-

$\{ t / P(t) \}$
 ↓ ↳ predicate of relation.
 set of output tuples

e.g:- 1. $\{ t / t \in \text{student} \} \equiv \text{select * from student}$
 $\equiv \prod_{\text{All attributes}} (\text{student})$

2. $\{ t / t \in \text{student} (t[m] > 70) \} \equiv \text{select * from student}$
 where ~~m~~ $> 70 ;$
 $\equiv \prod_{\text{All attributes}} (\sigma_{m > 70} \text{student})$

** no. of tuples $\begin{cases} 0 \\ m \text{ (size of relation)} \end{cases}$

To output only specific attributes.

3. $\{ t / \exists s \in \text{student} (s[\text{sno}] = t[\text{sno}] \wedge s[m] > 70) \}$

O/P:- s.no.

||
||

4) $\{ t \mid \exists s, t \in \text{student} (s[\text{SNO}] = t[\text{SNO}] \wedge s[\text{Name}] = t[\text{Name}] \wedge s[m] > 70) \}$

Q/P

SNO	Name

5) combining two relations

student (SNO, Name, m)

Reg (SNO, dept).

$$\begin{aligned} \{ t \mid \exists s \in \text{student} (s[\text{SNO}] = t[\text{SNO}] \wedge s[m] > 70 \wedge s[\text{Name}] = t[\text{Name}]) \\ \wedge \exists r \in \text{Reg} (r[\text{SNO}] = t[\text{SNO}] \wedge r[\text{dept}] = 'CSE' \wedge r[\text{dept}] = t[\text{dept}]) \\ \wedge s[\text{SNO}] = r[\text{SNO}] \} \end{aligned}$$

Q/P

SNO	Name	dept

DRC :-

$$\{ \langle A_1, A_2, \dots, A_n \rangle \mid Q(\langle A_1, A_2, \dots, A_n \rangle) \}$$

↗ output attributes ↗ Attributes participating in query condition. ↗ output attributes.

Eg 1.

student (SNO, Name, Marks)

$$\{ \langle SNO, Name, Marks \rangle \mid \langle SNO, Name, Marks \rangle \in \text{student} \wedge \text{marks} > 70 \}$$

\equiv select * from student where marks > 70;

2. $\{ \langle SNO \rangle \mid \langle SNO, \text{marks} \rangle \in \text{student} \wedge \text{marks} > 70 \}$.

\equiv select SNO from student where marks > 70

3. combining two relations.

$\text{Reg}(s, d)$

$\text{student}(s, N, M)$

$$\left\{ \begin{array}{l} \langle s, N, d \rangle / \exists t_1 \langle s, N, M \rangle \in \text{student} \wedge M > 70 \\ \wedge \exists t_2 \langle s, d \rangle \in \text{Reg} \wedge d = 'CSE' \\ \wedge t_1[s] = t_2[s] \end{array} \right\}$$

— x — x —

In both TRC & DRC
select customers ids who have an account and
loan in the bank.

Given

$\text{dep}(Acno, cid)$

$\text{Bor}(Lno, cid)$.

1) TRC

$$\left\{ \begin{array}{l} t / \exists r \in \text{dep} (t[cid] = r[cid]) \\ \wedge \exists s \in \text{Bor} (t[cid] = s[cid]) \\ \wedge r[cid] = s[cid] \end{array} \right\}$$

2) DRC

$$\left\{ \begin{array}{l} \langle cid \rangle / \exists r \langle cid \rangle \in \text{dep} \\ \wedge \exists s \langle cid \rangle \in \text{Bor} \\ \wedge r[cid] = s[cid] \end{array} \right\}$$

DATABASE DESIGN

- * ER model is a pictorial representation of database.

Database Design

→ Database Modelling.

→ Data Models

→ object based → object-oriented
(class & objects)
* Entity-Relationship Model

→ Record-based. (works on set of records)

* → Relational
→ Network
→ Hierarchical

⇒ Object is a real life entity which is distinguishable from other with a set of properties.

⇒ ER Model.

- Works on entities (objects)
- Attributes (properties)

ER Model.

Entity

→ e.g. S1 101 XYZ

→ an object

→ Attributes (properties) e.g. {S1, S2, S3}

Entity set : Set of similar entities.

Types of Attributes

1. Simple Attribute :- value of an attribute is constructed from a single component.
No subparts in it.

Eg :- Gender ↗ male
 ↘ female

2. Composite Attribute :- value of an attribute is constructed from more than one components. subparts will be there.

5)

Eg:- Name \rightarrow First name + Middle name + Last name

6)

Address \rightarrow door no + street name + country

3. Single Value Attribute :-

only one value for each entity.

Eg:- Rollno of student.

4. Multivalue Attribute :-

More than one value for the same entity.

Eg:- Subjects registered by a student.

phonenumbers

5. Derived Attribute :-

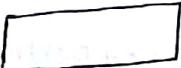
If the value of an attribute is derived from other attribute.

Eg:- Age derived from D.O.Birth

Experience " " D.O.join

Total marks of a student in S1 + S2 + S3

Symbols for an ER model

1) Entity set -  (rectangle)

2) Attribute -  (ellipse)

(simple,
composite,
Single value)

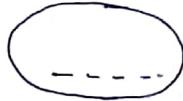
 (double ellipse)

3) Multivalue Attribute -

 (dashed ellipse)

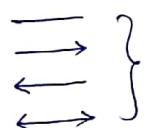
4) Derived Attribute -

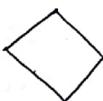
5) primary key Attribute - 

6) Descriminator - 
(an attribute
part of
weak
entity set)

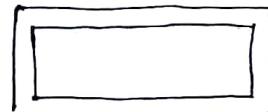
Lines:

— between Entity and its attributes.

 between two entities with
a Relationship.

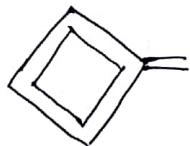
Relationship:  (diamond)

weak-Entity set:



(double
rectangle)

Total participation.



Relationship:

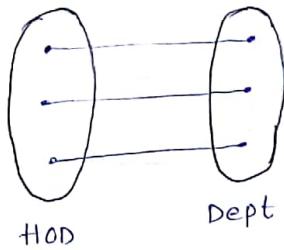
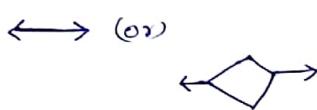
It is an association among entities.

Relationship set : set of similar relations

$S_1 \xrightarrow{R_1} D_1$
 $S_2 \xrightarrow{R_2} D_2$

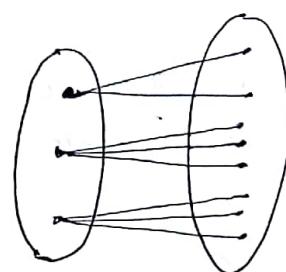
Types of Relation set.

1) one - to - one



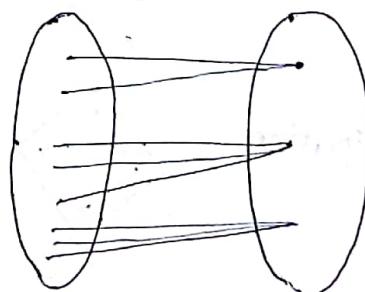
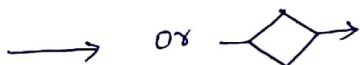
only one faculty as HOD for one dept, and
only one HOD for one dept.

2) One - to - many



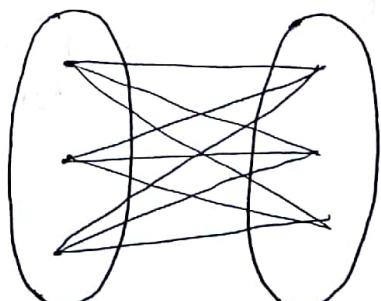
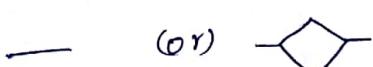
department students

3) Many to one



students department

④ Many to many

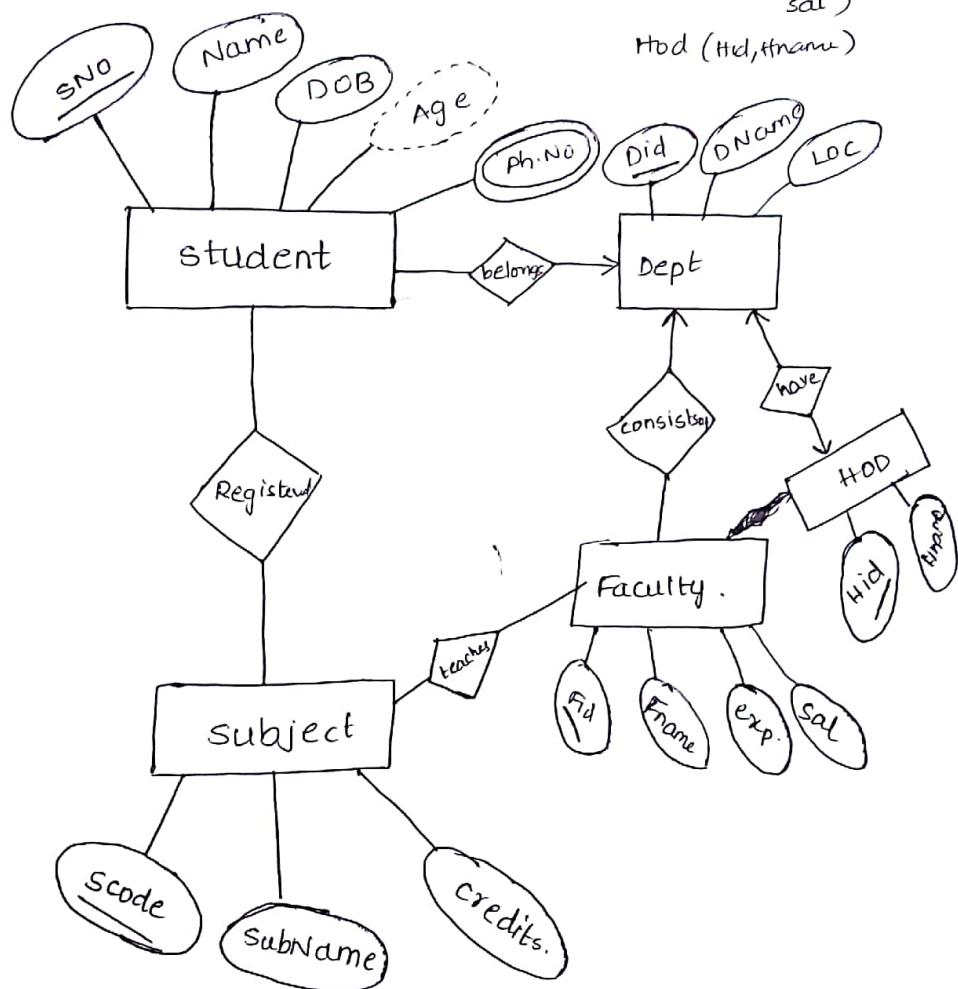


products

Customers

Example ER model for student database.

Dept (Did, DName, loc)
 Lab (Lid, LName, LL)
 faculty (Fid, FName, exp, sal)
 HOD (Hid, HName)



Assumptions:-

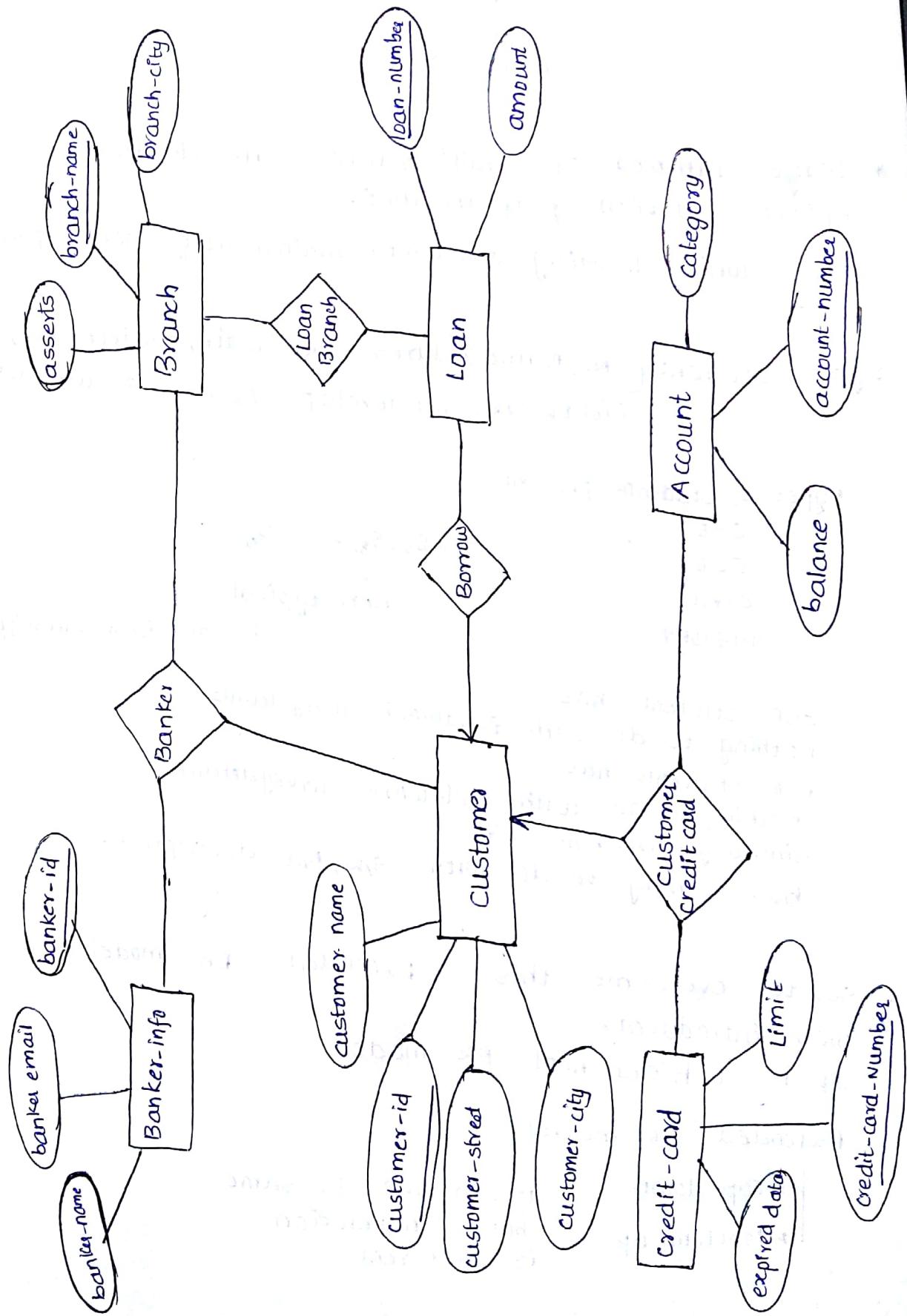
- many students may belong to each department, whereas each student cannot belong to more than one department
- Each subject can be registered by many students, each student can register to more than one subject.
- Each subject can be taught by more than one faculty and each faculty can teach more than one subject.
- Each faculty belongs to only one department, one department can have more than one faculty.

5) Each department will have only one HOD,
and each HOD will be there for each department



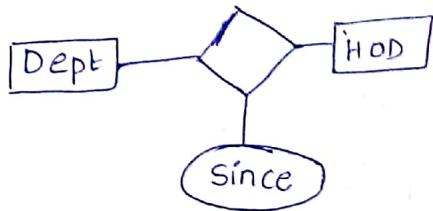
Advantages of centralization
1. Better coordination and planning of work
2. Better control over all activities
3. Better utilization of available resources
4. Better division of work
5. Better utilization of available funds
6. Better utilization of available labour
7. Better utilization of available equipment
8. Better utilization of available energy

ER Model. Bank database



Note:- joining attribute to a related entity.

* This attribute, ~~isn't~~ separately with each relation doesn't ~~seem~~ meaningful.



* Large number of 'null' values in database effect system performance.
like while loading data into main memory lack of space etc.

Eg:- Student (SNO, name, Address, practicals, projectwork, fieldwork, s/w develop, h/w develop, surveydetails)

types of students present

CSE

$s_1, s_2, s_3, \dots, s_n$

ECE

non typical

CIVIL

HISTORY

↳ so first classify

CSE students has nothing to do with fieldwork, survey details

ECE students has nothing to do with fieldwork, survey details

similarly CIVIL & HISTORY

has nothing to do with s/w, h/w development

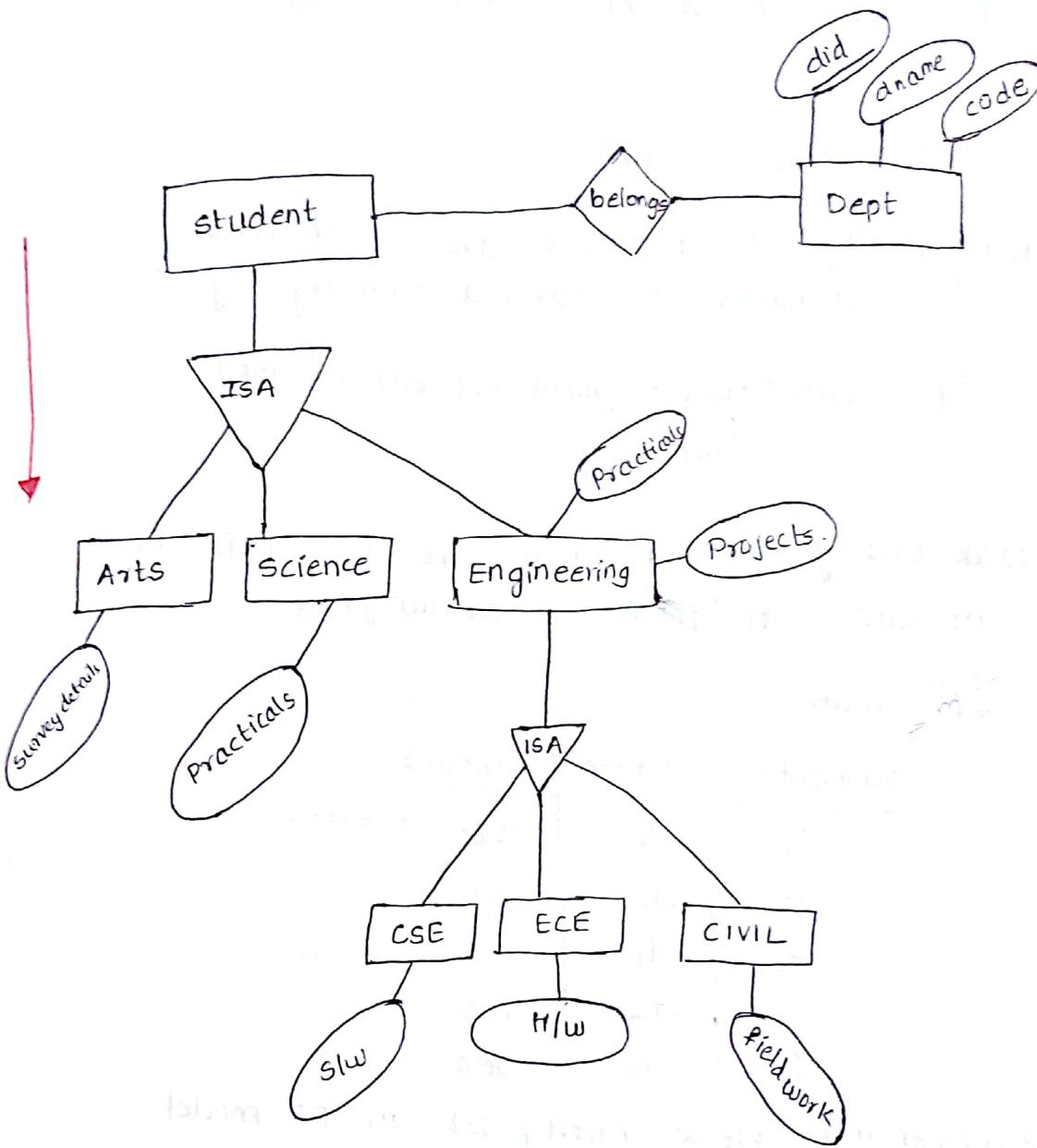
* so, to overcome this Extended ER model was introduced.

It is a hierarchical ER model.

Extended ER model

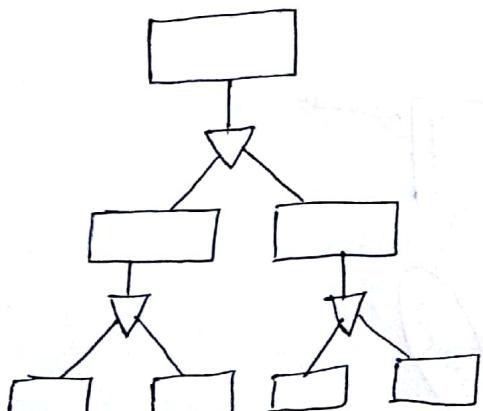
→ Top-down
→ Bottom up } appearance is same
but construction is different

Top-down method
⇒ "specialization"



Bottom-up method

⇒ "Generalization"



Key attribute: one or more attributes which
(Primary Key) uniquely identifies a tuple.

D

* Primary key should not be NULL & should be unique.

$E_1 (a_1, a_2, a_3, \dots, a_n)$

* Strong Entity set: If there are sufficient attributes to form a primary key.

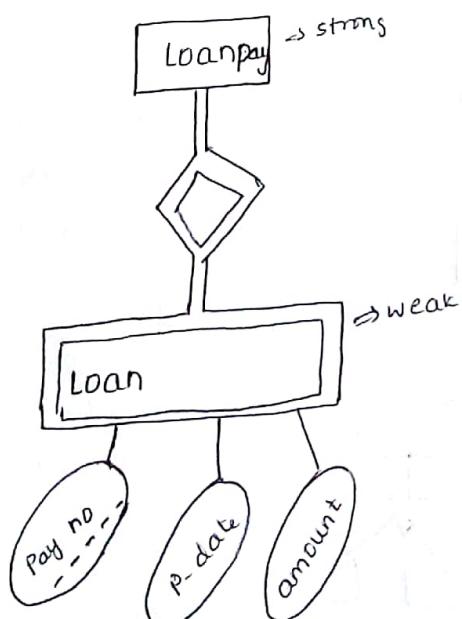
Eg:- Loan (Lno, paymentno, p-date, amount)
↓
primary key.

* Weak Entity set: If there are no sufficient attributes to form a primary key.

Eg:- Loan

paymentno	pdate	amount
P ₁	d ₁	100 ↗ person's
P ₂	d ₁	200
P ₁	d ₁	300 ↗ person
P ₁	d ₂	400
P ₂	d ₂	500

Representing weak Entity set in ER model



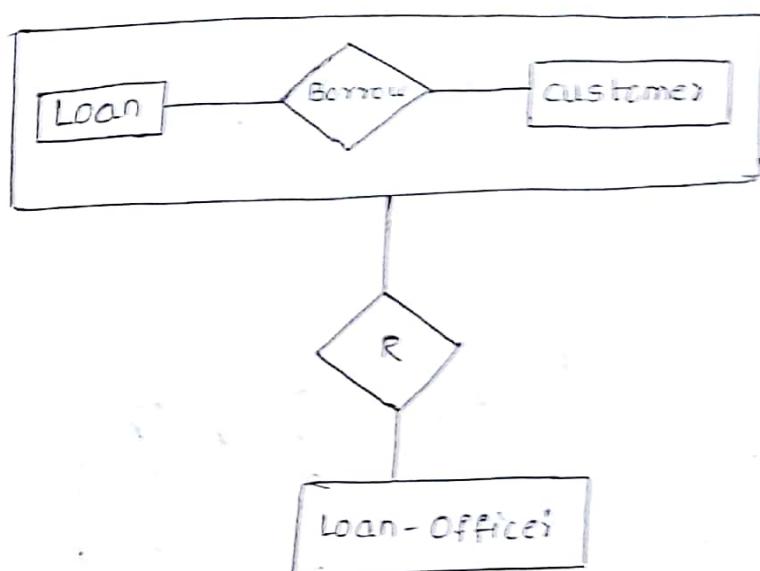
Discriminator:

It is an attribute, addition of it to a weak Entity set changes it to a strong Entity set.

A E.g:-

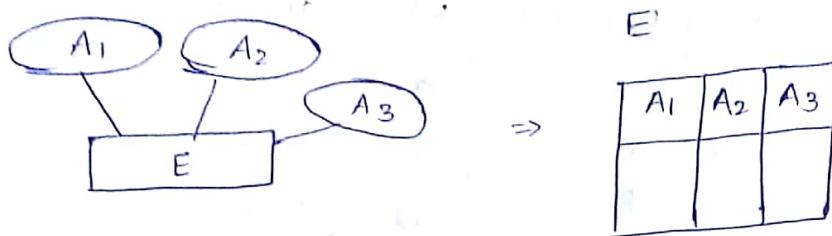
LNO	Payment no	pdate	amount
L ₁	P ₁	d ₁	100
L ₃	P ₂	d ₁	200
L ₂	P ₁	d ₂	300
L ₄	P ₁	d ₂	400
L ₁	P ₂	d ₂	500
L ₂	P ₂	d ₂	600

Aggregation



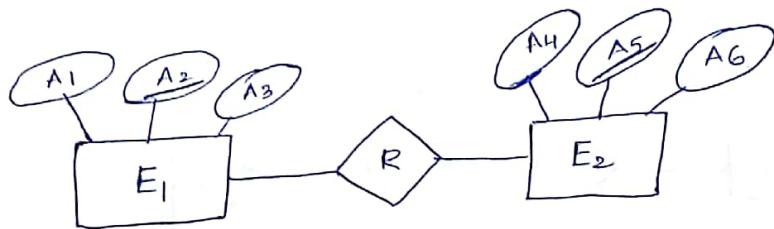
Converting an ER model to Relational Model

- ① Converting an Entity into relation (Table)



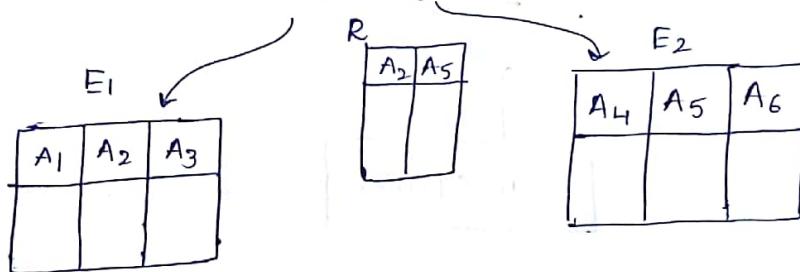
* Relation name is same as Entity name

- ② Converting a Relationship into relation (Table)₂



$$A\ell\ell(R) = PK(E_1) \cup PK(E_2) \cup A\ell\ell(R)$$

$$= A_2 \cup A_5 \cup \text{null.} = A_2 \cup A_5$$



- ③ Converting a multivalued attribute to Relation (Table)₃

student

S NO	S NAME	PHNO
101	x	P ₁
101	x	P ₂
101	x	P ₃

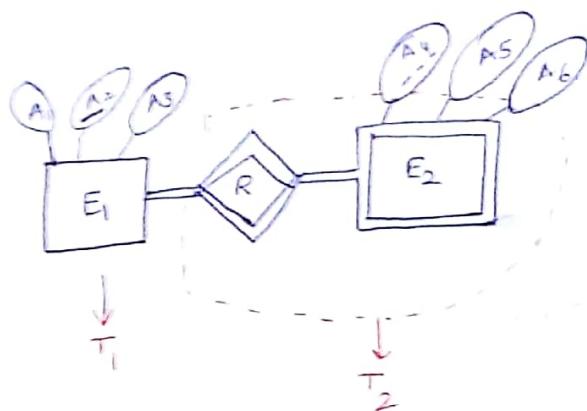
data redundancy

$$\text{Attr}(T_3) = \text{mv Attr} \cup \text{PK}(E)$$

Eq:- $E(A_1, A_2, \underline{A_3}, A_4, \underline{A_5})$
 ↓
 multi valued

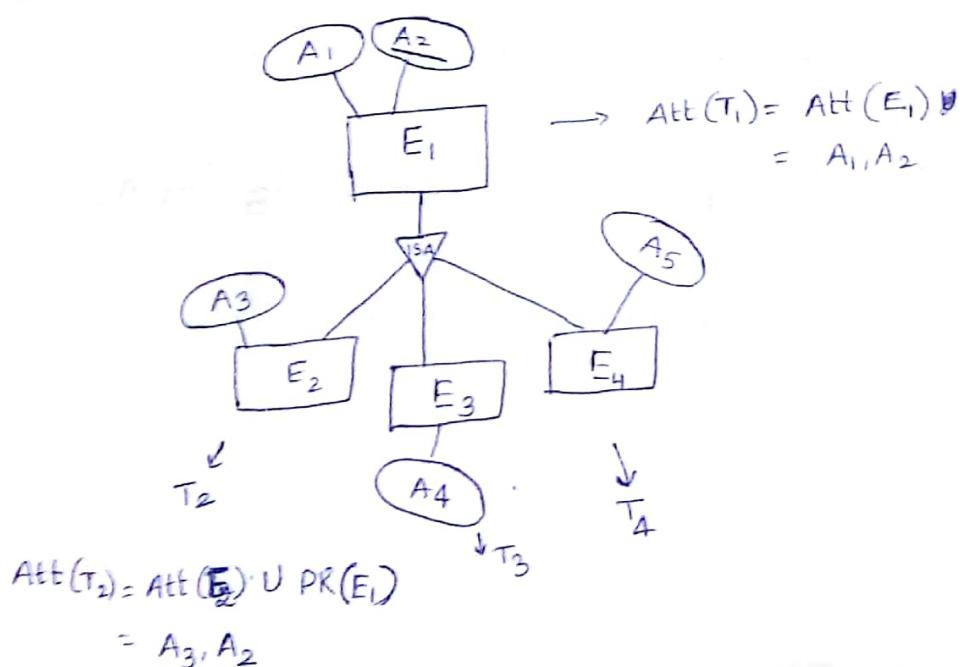
$$\begin{aligned}\text{Attr}(T_3) &= A_2 \cup A_3, A_5 \\ &= A_2, A_3, A_5\end{aligned}$$

4) Weak Entity set in to Table.



$$\begin{aligned}\text{Attr}(T_2) &= \text{Attr}(\text{weak Entityset}) \cup \text{PK}(\text{strong Entity}) \\ &= (A_4, A_5, A_6) \cup A_2 \\ &= A_2, A_4, A_5, A_6\end{aligned}$$

5) Conversion of specialization / Generalization into Relation (Table)



Relational Model for student database,
ER model is discussed before (10 pages)

Attr(student) = sno, Name, Dob, Age, phno

Attr(phno) = phno U (sno)
= phno, sno.

Attr(Dept) = Did, Dname, loc

Attr(HOD) = Hid, Hname

Attr(faculty) = fid, fname, exp, sal.

Attr(subject) = scode, subName, credit

Attr(belongs) = sno, did.

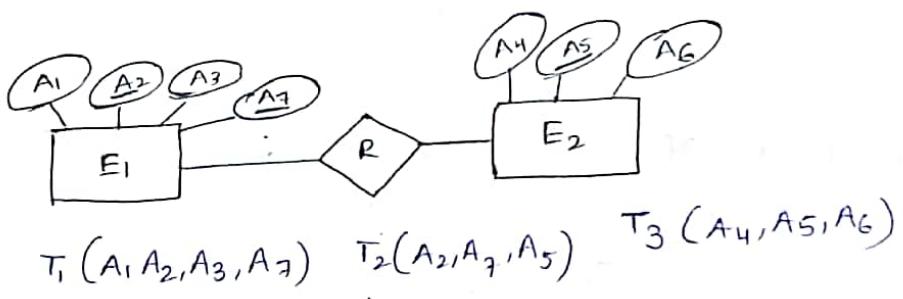
Attr(have) = Did, Hid

Attr(consists of) = Fid, Did

Attr(Registered) = sno, scode

Attr(teaches) = scode, fid.

Step: 2



$T_1(A_1, A_2, A_3, A_5, A_7)$, $T_3(A_4, A_5, A_6)$

minimum number of tables = 2

Need For Database design :

Student (SNO, Sname, Address, DNO, DName, loc)

1	A	A ₁	d ₁	Y	L ₁
2	B	A ₂	d ₁	Y	L ₁
3	C	A ₁	d ₁	Y	L ₁
4	D	A ₂	d ₂	Y	L ₂
5	E	A ₃	d ₂	Y	L ₂
6	F	A ₄	d ₃	Z	L ₃
7	G	A ₅	d ₃	Z	L ₃
8	H	A ₅	d ₄	K	L ₄
9	I	A ₆	d ₁	K	L ₄
10	J	A ₇	d ₅	P	L ₅

Problems with above relation

1) Data Redundancy.

2) Insertion anomaly.

- Eg:- if we insert a dept d₆ in which no student registered.
then studentcount gives 11 but there are only 10 students

3) Deletion problem

Eg:- if last tuple is deleted

then dept count = 4

but physically there are 5 departments

4) Updation anomaly

Eg:- if loc, Dname of tuples 1,2 are done

but tuple 3 is not done \Rightarrow Inconsistent.

Database Design



Normalization.

It is a process or theory of decomposing the larger relations with anomalies into smaller relations without anomalies based on Functional Dependency theory.

It is a process to eliminate data redundancy.

Functional Dependancy :- Dependencies exists b/w the Attribute of the relation

$$DB = (R_1, R_2, \dots, R_n)$$

\downarrow

$$R_i(A_1, A_2, \dots, A_n)$$

Let 'R' be a given relation

$x, y \in R$. x, y are the set of Attributes

$x \rightarrow y$
 $\swarrow \quad \downarrow$
Determinant Dependent.

x is functionally determines y

y is functionally depending on x

Let $\gamma_1(R)$ is an instance of R

\downarrow_{100}

$t_1, t_2 \in \gamma_1$

if $t_1[x] = t_2[x]$ then

$t_1[y] = t_2[y]$

Eg:

	S.NO	Name.
t_1	1	A
	2	B
	3	C
	4	D
t_2	1	A
	2	B
	3	C
	4	D
	5	All
	6	B,

$\Rightarrow F = \{ f_1, f_2, f_3 \}$ Let R be a given set of attributes
 ↓ set of FD's. $\alpha, \beta, r, S \in R$ \hookrightarrow set of Attributes

If one FD is known then to know hidden FD's.

Closure properties of FD's

- 1) Reflexivity: if $B \subseteq \alpha$ then $\alpha \rightarrow B$
- 2) Transitivity: if $\alpha \rightarrow B$ and $B \rightarrow f$ then $\alpha \rightarrow f$
- 3) Augmentation: if $\alpha \rightarrow B$ and f then $\alpha f \rightarrow Bf$
- 4) Union: if $\alpha \rightarrow B$ and $\alpha \rightarrow f$ exists then $\alpha \rightarrow Bf$
- 5) Decomposition: if $\alpha \rightarrow Bf$ exists then $\alpha \rightarrow B$ and $\alpha \rightarrow f$ exists
- 6) pseudo transitivity: if $\alpha \rightarrow B$ and $Bf \rightarrow S$ then $\alpha \rightarrow S$

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

$$\begin{array}{l} \text{RNO} \rightarrow \{\text{name}\} \\ A \rightarrow B \\ \text{RNO} \rightarrow \{\text{Marks}\} \\ A \rightarrow C \end{array} \quad \left\{ \begin{array}{l} A \rightarrow BC \\ \text{RNO} \rightarrow \{\text{name}, \text{Marks}\} \end{array} \right.$$

Consider Bank Database

Customer (cid, cname, Address, Age)

Depositor (AccNo, bal, bid, bname, loc)

Borrower (LNo, amt, bid, bname, loc)

$$\begin{array}{l} \text{AccNo} \rightarrow \text{bal} \\ \text{AccNo} \rightarrow \text{bid} \\ \text{AccNo} \rightarrow \text{bname} \\ \text{AccNo} \rightarrow \text{loc} \\ \text{bid} \rightarrow \text{bname} \\ \text{bid} \rightarrow \text{loc} \end{array} \quad \left\{ \begin{array}{l} \text{AccNo} \rightarrow \{\text{bal}, \text{bid}, \text{bname}, \text{loc}\} \\ \text{bid} \rightarrow \{\text{bname}, \text{loc}\} \end{array} \right.$$

Finding a key

input: R is a given relation, set of F

output: set of candidate keys

Step 1: α is a set of attributes $\alpha \subseteq R$

Step 2: if $\alpha^+ \rightarrow R$ and posses the following two properties.

Step 3: Uniqueness: Uniquely identifies each tuple in R ($\alpha^+ \rightarrow R$)

Irreducibility: β is any subset of α

if $\beta \neq \emptyset$ then is uniquely identifies each tuple in R then α is a superkey and β is candidate key.

Closer of attributes

i/p : R, F, α .

o/p : α^+

Step 1: result = α

Step 2: For each FD $B \rightarrow f$ in F

if $f \subseteq \text{result}$ then

result = result $\cup f$

Step 3: if any changes in result then
goto step 2

Step 4: $\alpha^+ = \text{result}$

Ex: ① R(A, B, C, D)

F: $A \rightarrow B$

$A \rightarrow C$

$B \rightarrow D$

$\underline{\underline{B^+}}: \quad \alpha = B \quad \text{result} = B.$

$B \rightarrow f$

$A \rightarrow B \quad A \subseteq B \quad \text{No}$

$A \rightarrow C \quad A \subseteq B \quad \text{No}$

$B \rightarrow D \quad B \subseteq B \quad \text{Yes.}$

result = $B \cup D = BD$

change in result, so

$A \rightarrow B \quad A \subseteq BD \quad \text{No}$

$A \rightarrow C \quad A \subseteq BD \quad \text{No.}$

$B \rightarrow D \quad B \subseteq BD \quad \text{Yes.}$

result = $BD \cup B = BD$.

no change in result

so, $\alpha^+ = BD$

$B^+ = BD$

Eg ②: $R(A B C \quad G H I)$

$$F = \{ \begin{array}{l} A \rightarrow B, \\ A \rightarrow C, \\ CG \rightarrow H \\ CG \rightarrow I \\ B \rightarrow H \end{array} \}$$

find \check{AB}^+ , \check{AC}^+ , \check{AG}^+ , \check{AH}^+ , \check{AI}^+

$$\check{AB}^+ := \alpha = AB, \text{ result} = AB$$

$$A \rightarrow B, \quad A \subseteq AB, \text{ yes} \\ \text{result} = B \cup AB = AB.$$

$$A \rightarrow C, \quad A \subseteq AB, \text{ yes} \\ \text{result} = C \cup AB = ABC$$

$$CG \rightarrow H, \quad AC \subseteq ABC, \text{ NO}$$

$$CG \rightarrow I, \quad CG \subseteq ABC, \text{ NO.}$$

$$B \rightarrow H, \quad B \subseteq ABC, \text{ yes.}$$

$$\text{result} = H \cup ABC = ABCH$$

Now, change in result, so loop once again

$$\text{so, } \alpha^+ = ABCH$$

$$\boxed{\check{AB}^+ = ABCH}$$

$$\check{AC}^+ := \alpha = AC, \text{ result} = AC$$

$$A \rightarrow B, \quad A \subseteq AC, \text{ yes} \\ \text{result} = B \cup AC = ABC$$

$$A \rightarrow C, \quad A \subseteq ABC, \text{ yes} \\ \text{result} = C \cup ABC = ABC$$

$$CG \rightarrow H, \quad CG \subseteq ABC, \text{ NO.}$$

$$CG \rightarrow I, \quad CG \subseteq ABC, \text{ NO.}$$

$$B \rightarrow H, \quad B \subseteq ABC, \text{ yes}$$

$$\text{result} = H \cup ABC = ABCH$$

change in result, so

$$A \rightarrow B, \quad A \subseteq ABCH, \text{ yes}$$

$$\text{result} = B \cup ABCH = ABCH$$

$A \rightarrow C$, $A \subseteq ABCH$, yes
 result = $C \cup ABCH = ABCH$

$CG \rightarrow H$, $CG \subseteq ABCH$, no

$CG \rightarrow I$, $CG \subseteq ABCH$, no.

$B \rightarrow H$, $B \subseteq ABCH$, yes

result = $H \cup ABCH = ABCH$.

no change in result

$\therefore A^+ = ABCH$

$\boxed{A^+ = ABCH}$

— * —

Eg: for finding superkey and candidate key

$R(ABCD)$

$F:$
 $A \rightarrow B$
 $A \rightarrow C$
 $B \rightarrow D$
 $D \rightarrow A$.

$X^+ \rightarrow R$
 Key \downarrow
 superkey $(X-A)^+ \rightarrow R$
 $A; \in X-A \rightarrow$ superkey
 $((X-A)-A_1)^+ \rightarrow R$.

$\{A, B, AB, BC, AC, AD, BD, CD$
 ~~$C, D, ABC, ABD, ACD, BCD, ABCD\}$~~

$(A)^+$; $d = A$ result = A.

$A \rightarrow B$.

$A \subseteq A$, yes result = $A \cup B = AB$

$A \rightarrow C$, $A \subseteq AB$, result = ABC

$B \rightarrow D$, $B \subseteq ABC$, result = $ABCD$

$D \rightarrow A$, $D \subseteq ABCD$, result = $ABCD$.

change in result, so repeat

we get $A^+ = ABCD$

A is candidate key,

similarly check all subsets of R.

$(B)^+$; $d = B$, result = B

$A \rightarrow B$

$A \not\subseteq B$, no.

$A \rightarrow C$

$A \not\subseteq B$ no

$B \rightarrow D$

$B \subseteq B$, yes, result = BD

$D \rightarrow A$,

$D \subseteq BD$, yes result = BD

change in result, so repeat

$B^+ = ABCD$

$\rightarrow A \rightarrow B$
 $A \subseteq ABD$, yes, result = $B \cup ABD = ABD$

$A \rightarrow C$
 $A \subseteq ABD$, yes, result = $C \cup ABD = ABCD$

$B \rightarrow D$. " " = $ABCD$

$D \rightarrow A$ " " $\rightarrow ABCD$

change in result, repeat

$(B)^+ = ABCD \therefore B$ is candidate key.

2. $R(ABCDEF GH) \xrightarrow{\text{g independent}}$

$$F = \{ CH \rightarrow G, A \rightarrow BC, \\ B \rightarrow CFH, E \rightarrow A, \\ F \rightarrow EG \}$$

List all candidate keys

{ A, B, C, E, F, G, H, AB, AC, AE, AF, AG, AH, BC, BE, BF, BG, BH, CE, CF, CG, CH, EF, EG, EH, FG, FH, GH, ABC, ABE, ABF, ABG, ABH, ... }

Closure of FD's (F^+)

F is a given FD's on R .

F^+ is logically derived from F using closure properties of FD's.

Ex1. $R(ABCD)$

$$A \rightarrow B$$

$$A \rightarrow C$$

$$B \rightarrow D$$

i/p: R, F

$$O/P: F^+ A \rightarrow B, A \rightarrow C \Rightarrow A \rightarrow BC$$

$$A \rightarrow B, B \rightarrow D \Rightarrow A \rightarrow D$$

Step

① for all possible attribute set X , calculate X^+

Step

② For each attribute $a \in X^+$ Add $X \rightarrow a$ to F^+

Step

3: return F^+

end

For Ex: 1

$$A^+ = ABCD$$

$$\begin{array}{llll} A \in A^+ & B \in A^+ & C \in A^+ & D \in A^+ \\ A \rightarrow A, & A \rightarrow B & A \rightarrow C & A \rightarrow D \end{array}$$

$$F_1^+ = \{ A \rightarrow A, A \rightarrow B, A \rightarrow C, A \rightarrow D \}$$

$$B^+ = BD$$

$$B \in B^+ \quad D \in B^+$$

$$B \rightarrow B, \quad B \rightarrow D$$

$$F_2^+ = \{ B \rightarrow B, B \rightarrow D \}$$

$$C^+ \neq \alpha = C, \text{ result} = C$$

$$A \rightarrow B, \quad A \not\subseteq C$$

$$A \rightarrow C, \quad A \not\subseteq C$$

$$B \rightarrow D, \quad B \not\subseteq C$$

no change in result

Note: $A \rightarrow A$
 $y \subseteq X$ then
 $x \rightarrow y$ is trivial

$C^t = C$

$$F_3^t = \{C \rightarrow C\}$$

(D)^t: $\alpha = D$, result = D.

$A \rightarrow B$, $A \subseteq D$, NO

$A \rightarrow C$, $A \subseteq D$, NO

$B \rightarrow D$, $B \subseteq D$, NO

No change in result

$$D^t = D$$

$$F_4^t = \{D \rightarrow D\}$$

similarly calculate for, $AB^t, AC^t, AD^t, BC^t, BD^t, CD^t, ABC^t$ - - - -

$$F^t = F_1^t \cup F_2^t \cup F_3^t \cup F_4^t \cup \dots \cup F_n^t$$

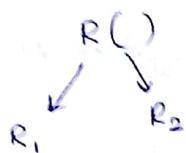
$$= \{ A \rightarrow A, A \rightarrow B, A \rightarrow C, A \rightarrow D, \\ B \rightarrow B, B \rightarrow D, \\ C \rightarrow C, \\ D \rightarrow D \} - - - - -$$

Decomposition of Relation :-

To overcome, drawbacks of large relation like redundancy, insertion, updation & deletion anomalies
We need to Decompose,

$$R(A_1, A_2, \dots, A_p, A_n)$$

Valid decomposition: if it is possible to retrieve data from the decomposed relations.



$$\begin{aligned} R_1 \cup R_2 &= R \\ \text{Att}(R_1) \cup \text{Att}(R_2) &= \text{Att}(R) \\ R_1 \cap R_2 &\neq \emptyset \end{aligned}$$

lossless join: if data is not lost after joining the decomposed relations.

Lossy decomposition :- if information is lost during decomposition.

Eg:-

$$R(A B C D)$$

$$\begin{array}{l} A \rightarrow B \\ A \rightarrow C \\ B \rightarrow D \end{array} \quad \left. \begin{array}{l} A \rightarrow BC \\ A \rightarrow C \end{array} \right\} A \rightarrow BC \quad R_1(A B C)$$

$$R_2(B, D)$$

$$R_1 \cap R_2 = A B C$$

$$R_1 \cap R_2 \rightarrow R_2$$

$$B \rightarrow R_2 \quad \checkmark$$

\therefore It is a lossless decomposition or valid decomposition.

Normalization

It is process of successive decomposition of relation into number of relations with any anomalies.

1NF

2NF

3NF

BCNF

4NF

5NF

1NF :- A relation R is in 1NF, if all values are atomic.

. NO repeated columns & NO set of values

Not in 1NF			→ convert into 1NF		
SNO	Name	phno.	SNO	Name	phno.
101	x	{111, 222}	101	x	111
102	y	{333, 444, 555}	101	x	222
			102	y	333
			102	y	444
			102	y	555

* But it creates data redundancy.

2NF :- A relation R is in 2NF, iff it is in 1NF and all non key attributes are fully dependent on **key** Attributes.

↳ candidate key
not superkey

Eg:- Company (cno, pno, aty, city, status)

cno	pno	aty	city	status
101	P1	100	Q1	A
101	P2	200	C1	A
102	P1	200	C1	B
102	P2	300	C1	A
103	P1	400	C2	B
103	P2	500	C2	B

FD = { cno → city
cno → status.
city → status. }

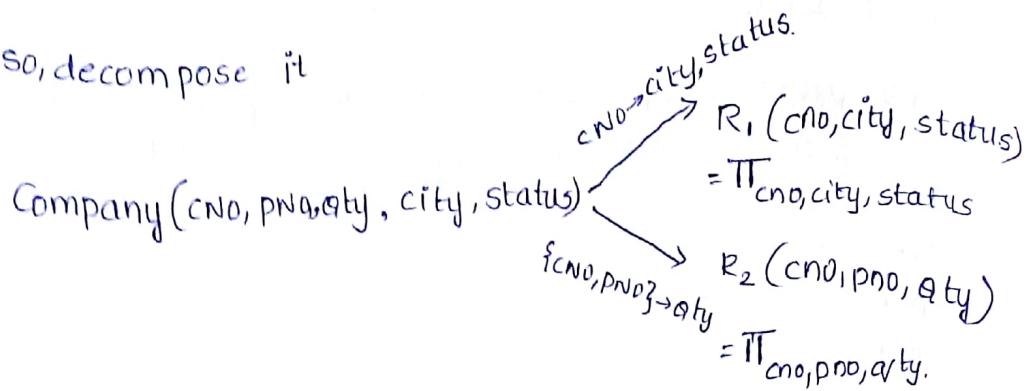
{ cno, pno } → aty

}

city is partially dependant on cno
so not

(cno) + → city, status.
it is violating the property of 2NF

so, decompose it



Now, check if the decomposition is valid or not

$$R_1 \cap R_2 = \text{cno}$$

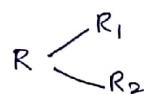
$$R_1 \cap R_2 \rightarrow R_1 \quad [\because \text{cno} \rightarrow \{\text{city, status}\}]$$

\therefore It is a valid lossless join
decomposition

Dependancy Preservation:

Let R be a given relation

F is a set of F.D's



$$(F_{R_1} \cup F_{R_2})^+ = F^+$$

$$\begin{aligned} \text{cno} \rightarrow \text{city, status} &= \left\{ \begin{array}{l} \text{cno} \rightarrow \text{city, state} \\ \vdots \\ \{ \text{cno, pno} \} \rightarrow \text{qty} \end{array} \right. \\ &\cup \\ &\underline{\left\{ \begin{array}{l} \{ \text{cno, pno} \} \rightarrow \text{qty} \\ \text{cno} \rightarrow \text{city, pno} \end{array} \right\}} \quad \left. \begin{array}{l} \vdots \\ \} \end{array} \right\} \end{aligned}$$

— x — x —

$R_2 \rightarrow$ is in 2NF (\because all non key attributes are
fully depending on key attribute)

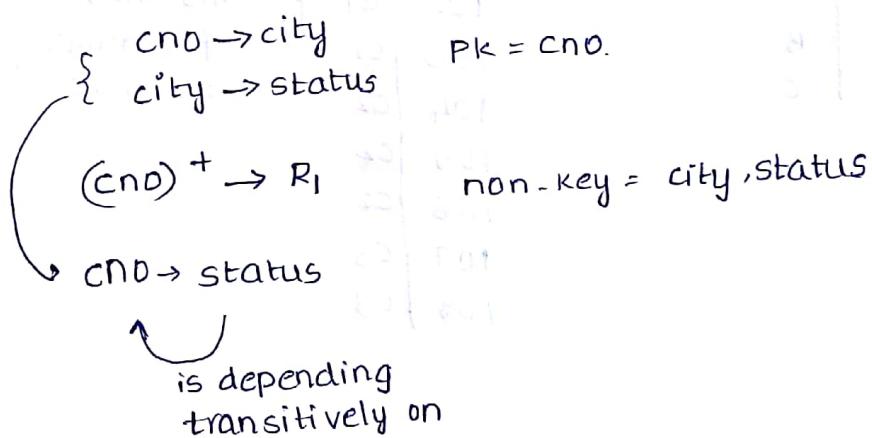
$R_1 \rightarrow$ is in 2NF

cno	city	status
101	C1	A
102	C1	A
103	C1	A
104	C2	B
105	C2	B
106	C2	B
107	C3	C
108	C3	C
-	C4	D

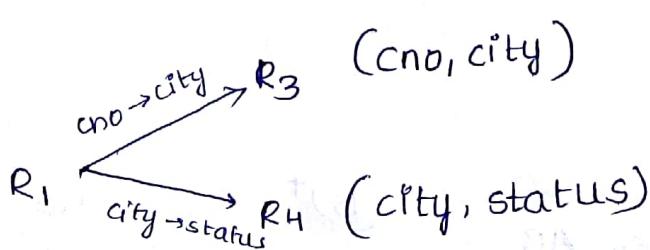
Insertion anomaly

3NF :- If a relation R is in 3NF
iff it is in 2NF and non-key attributes
are non-transitively dependent on PK attribute.

Eg:- From the above table.



\therefore city \rightarrow status is violating the property of 3NF
so, decompose it.



$$R_3 \cap R_4 \rightarrow R_1$$

\therefore It's a valid decomposition.

Testing dependency preservation.

~~R₃~~

$$\left(F_{R_3} \cup F_{R_4} \right)^+ = F^+$$

$$\left(\text{cno} \rightarrow \text{city} \cup \text{city} \rightarrow \text{status} \right) = \left(\text{cno} \rightarrow \text{city} \right)^+ \cup \left(\text{city} \rightarrow \text{status} \right)^+$$

Check if individual relations are in 3NF or not

$R_3 \rightarrow$ in 3NF

$R_4 \rightarrow$ is also in 3NF.

city	status
C_1	A
C_2	B
C_3	C

cno	city
101	C_1
102	C_1
103	C_1
104	C_2
105	C_2
106	C_2
107	C_3
108	C_3

this much percentage
of redundancy is
acceptable.

Question ①: $R(A, B, C, D)$, Test if it is in 3NF

$A \rightarrow B$

$A \rightarrow C$

$B \rightarrow D$

$R(A B C D E) \rightarrow$ keep the

relation in

$A \rightarrow B$

$B \rightarrow C$

$D \rightarrow E$

it highest NF

$$(A)^+ = \begin{array}{l} A \rightarrow B \text{ yes} = A \\ A \subseteq A \text{ res} = AB \\ A \rightarrow C \text{ yes} = ABC \\ A \subseteq AB \text{ res} = ABCD \\ B \rightarrow D \text{ yes} = ABCD \\ B \subseteq ABC \text{ res} = ABCD \end{array}$$

$$(B)^+ = \begin{array}{l} A \subseteq B \text{ no} \\ A \subseteq B \text{ no} \\ B \subseteq B \text{ yes, res} = BD \end{array}$$

$$(B)^+ = BD$$

$$(C)^+ = C$$

$$(D)^+ = D$$

$$R \xrightarrow{ } R_1(A, B, C) \\ R \xrightarrow{ } R_2(B, D)$$

$$R_1 \cap R_2 = B$$

$$B \rightarrow R_2 \quad \checkmark$$

\therefore It is a valid decomposition

$$(F_{R_1} \cup F_{R_2})^+ = (F)^+ \quad \text{dependancy preservation}$$

R_1 is in 2NF

R_2 is in 2NF

$R_1(A, B, C)$

$$A \rightarrow B$$

$$A \rightarrow C$$

non key attributes B, C are depending non transitively
on ' A '

$\therefore R_1$ is in 3NF

$R_2(B, D)$

$$B \rightarrow D$$

R_2 is also in 3NF.

Question 2: $R(A, B, C, D, E)$

$$A \rightarrow B$$

$$B \rightarrow C$$

$$D \rightarrow E$$

keep the relation in its highest normal form.

7/9/18

Q) $R = (ABCDEF GH)$ keep the relation in highest Normal form.

$$F = \{ \begin{array}{l} CH \rightarrow G \\ A \rightarrow BC \\ B \rightarrow CFH \\ E \rightarrow A \\ F \rightarrow EG \end{array} \}$$

Sol:- Assume that R is 1NF

'D' is an independent key

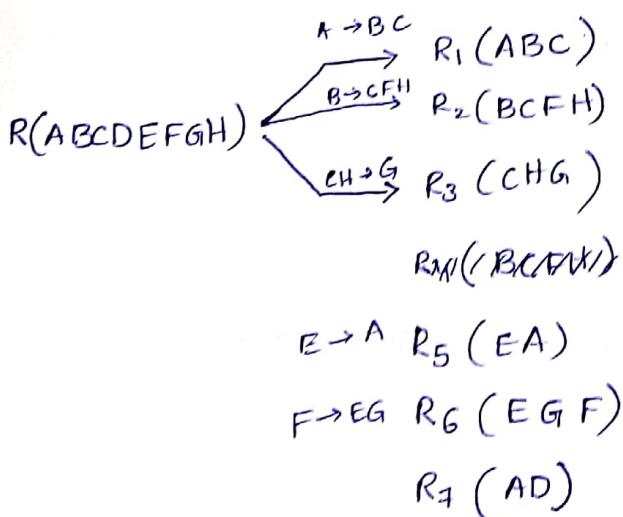
Candidate keys:- AD, ED, FD, BD

$$\left[\begin{array}{l} A \rightarrow R \\ E \rightarrow A \\ F \rightarrow E \\ B \rightarrow F \end{array} \right]$$

Non key attributes:- G, C, H, A, B, D, E, F

∴ non key attributes are not fully depending on key attributes

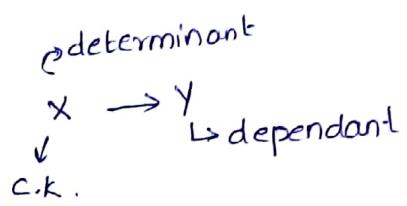
∴ R is not in 2NF then decompose it



BCNF :-

⇒ If a relation is in BCNF.

⇒ All the determinants are candidate key's



⇒ Eg:-

Company (cno, pno, qty, city, status)

FD's:

- $cno \rightarrow \{city, status\}$
- $city \rightarrow status$
- $\{cno, pno\} \rightarrow qty$

CK $\{cno, pno\}$

$cno \rightarrow city$
 $\because city \rightarrow status$ } is violating BCNF rule.

so, decompose

R → R₁ (cno, pno, qty) in BCNF ✓

R → R₂ (cno, city, status)
 $cno \rightarrow city$
 $\boxed{city \rightarrow status}$ not in BCNF
 violating

decompose

$cno \rightarrow city$

R₃ (cno, city)

in BCNF ✓

$city \rightarrow status$

R₄ (city, status)

in BCNF ✓

no
Find R(A, B C D) in which NF?

FDS:
 $AB \rightarrow D$
 $BC \rightarrow D$
 $A \rightarrow B$
 $B \rightarrow A$

Sol:



prime attribute: A, B, D

Non-prime attribute: C

Primary key: A, B

Foreign key: D

Relationship: A-B-C-D

Step 1: Find FDs in R(A, B, C, D)

Step 2: Find primary key

Step 3: Find non-prime attribute

୪୧୮

In Which NF?

S.NO	Subject	Language
1	S1	L1
1	S1	L2
1	S2	L1
1	S2	L2
2	S1	L1
2	S1	L2
2	S2	L1
2	S2	L2
3	S1	L3
3	S1	L2
3	S1	L1

Here, there is no functional dependency.

It is 2NF \because no partial dependancies

IE is 3NF :: no transitive dependancies

It is in BCNF also.

But still there is data redundancy.

4NF:-

multi valued Dependancy (MVD) :-
 Let R be a given Relation $x, y \in R$ are set of attributes

butes

$x \rightarrow$
↓
Multiple
determinant

multivalue
dependant

i.e., multiple values

of y are depending

on single value of x

single values of x

• 100 •

\Rightarrow is a multi valued
 \Leftrightarrow dependency

value
dant

multiple values

of y are depending

on singular values of \mathbf{x}

the single values of λ .

$$t_1, t_2, t_3, t_4 \in R$$

$$t_1[x] = t_2[x] = t_3[x] = t_4[x]$$

$$t_1[y] = t_3[y]$$

$$t_2[y] = t_4[y]$$

$$t_1[R-Y] = t_4[R-Y]$$

$$t_2[R-Y] = t_3[R-Y]$$

so, for the previous example
MVD's are $SNO \rightarrow\!\!\! \rightarrow$ Language.

4NF

\Rightarrow A Relation R is in 4NF, if it is in BCNF
and there are ^{no} multivalued independent
parts in the primary key

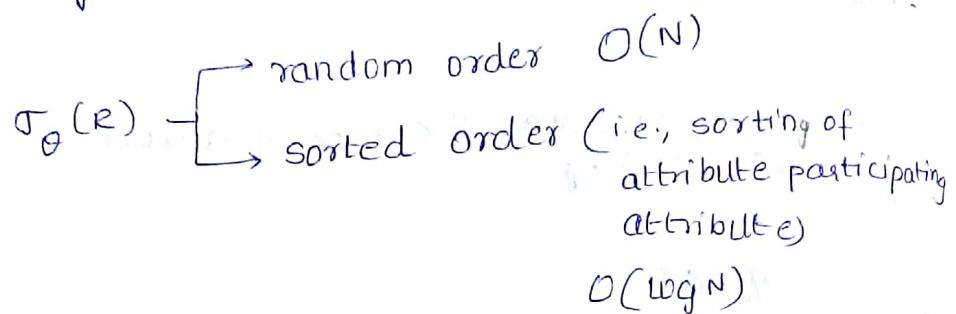
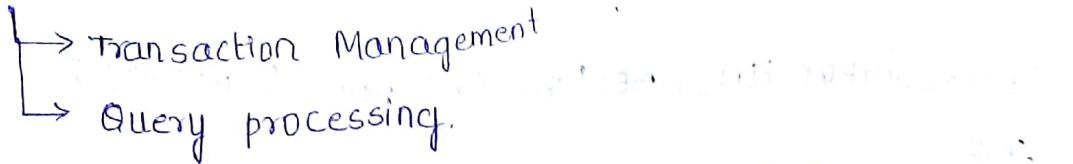
The previous example is not in 4NF,
so, decompose it

MVD $SNO \rightarrow\!\!\! \rightarrow$ subject $R_1(SNO, \text{subject})$
 $SNO \rightarrow\!\!\! \rightarrow$ Language $R_2(SNO, \text{language})$

$R_1(SNO, \text{subject})$	
1	S1
1	S2
2	S1
2	S2
3	S1

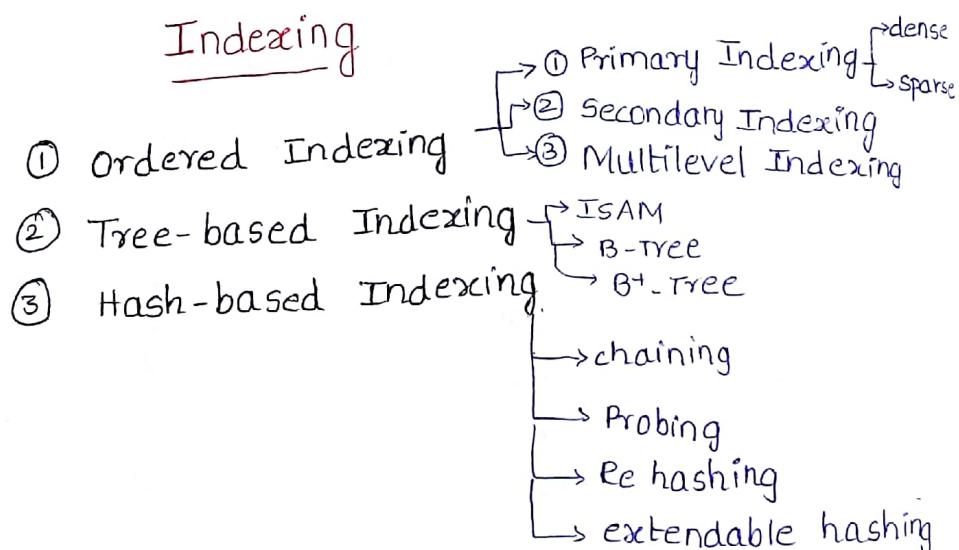
$R_2(SNO, \text{language})$	
1	L1
1	L2
2	L1
2	L2
3	L1
3	L2
3	L3

DBMS



Indexing :- Indexing is used for improving system performance i.e., speed up access to desired data.

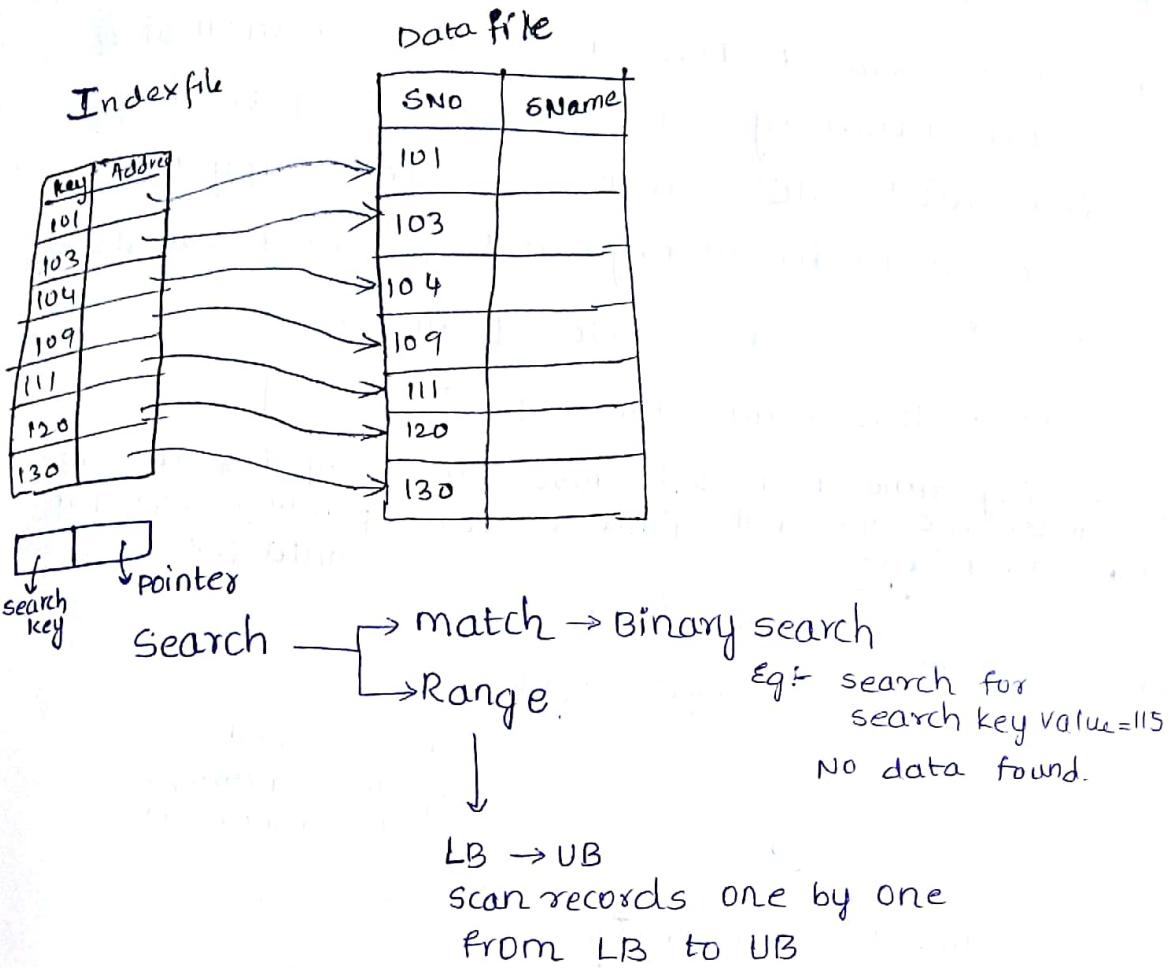
Indexing



Ordered Indexing

1) **Primary Indexing** : Assumes that data file is in sorted order.

Dense :- For each search key value we find an entry in index file corresponding to data file.



if a record is found it will point to that and does sequential scan upto upper boundary.

Insert →^{Eg} insert a new record with search key value = 105.

First it searches for the record with Search key value just less than 105 i.e., 104

After inserting, update the Index file, with new record 105.

Deletion → search the data record with the given search key value.

Update Index file

because [no. of records in Index file = Data file]
in "dense" indexing

Advantages of Primary dense Indexing

- 1) When size of D.B. is greater than that of main memory, we need not partition the data file, just we can copy index file to main memory and perform operations on it. [Because size of index file is much less than size of data file]

- 2) Searching time is less, because for every key value in index file we will find a corresponding record in data file.

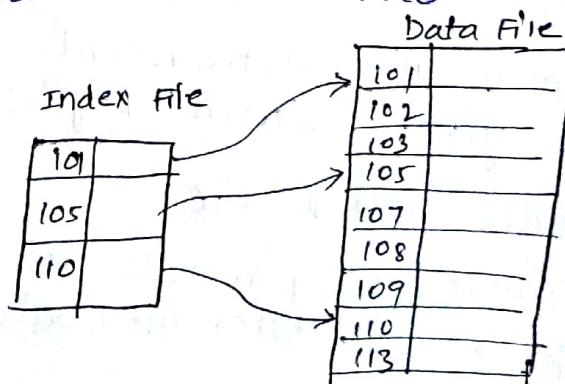
Disadvantage:-

- 1) search
 - match $\rightarrow t_1 + C$
 - range $\rightarrow t_1 + t_2$
i.e., we may find only two matching records in 100 records.
- Insert $\rightarrow T = t_1 + \underset{\text{updating}}{t_2} + C$
- Delete $\rightarrow T = t_1 + \underset{\text{deleting}}{t_2} + C$

- 2) When no. of data file records are increasing, no. of index file records increases and it may lead to partitioning of Index File.

Sparse primary Indexing

Here only few records of data file are present in Index File



search → match → search for the search key value less than the required search key value.

If found, from there do sequential search till we find the required record.

$$T = t_1 + c$$

↳ Range: similar to match query but start searching from search key value less than lower boundary

$$T = t_1 + t_2$$

Insert: $T = t_1 + t_2 + c$

↓
search for search key value less than or equal to inserting record key value

↓
linear search to find exact position

** no need to update.

Delete: $T = t_1 + t_2 + c$

↓
search for search key value nearest to the required key value

↓
linear search till we ~~not~~ find exact record

↓
delete.

** no need to update.

Note: Don't decide by scanning Index file that a record is found.

Advantages of sparse Indexing

- ⇒ It reduces the size of the index file saving space
- ⇒ we need not update the index file all the time

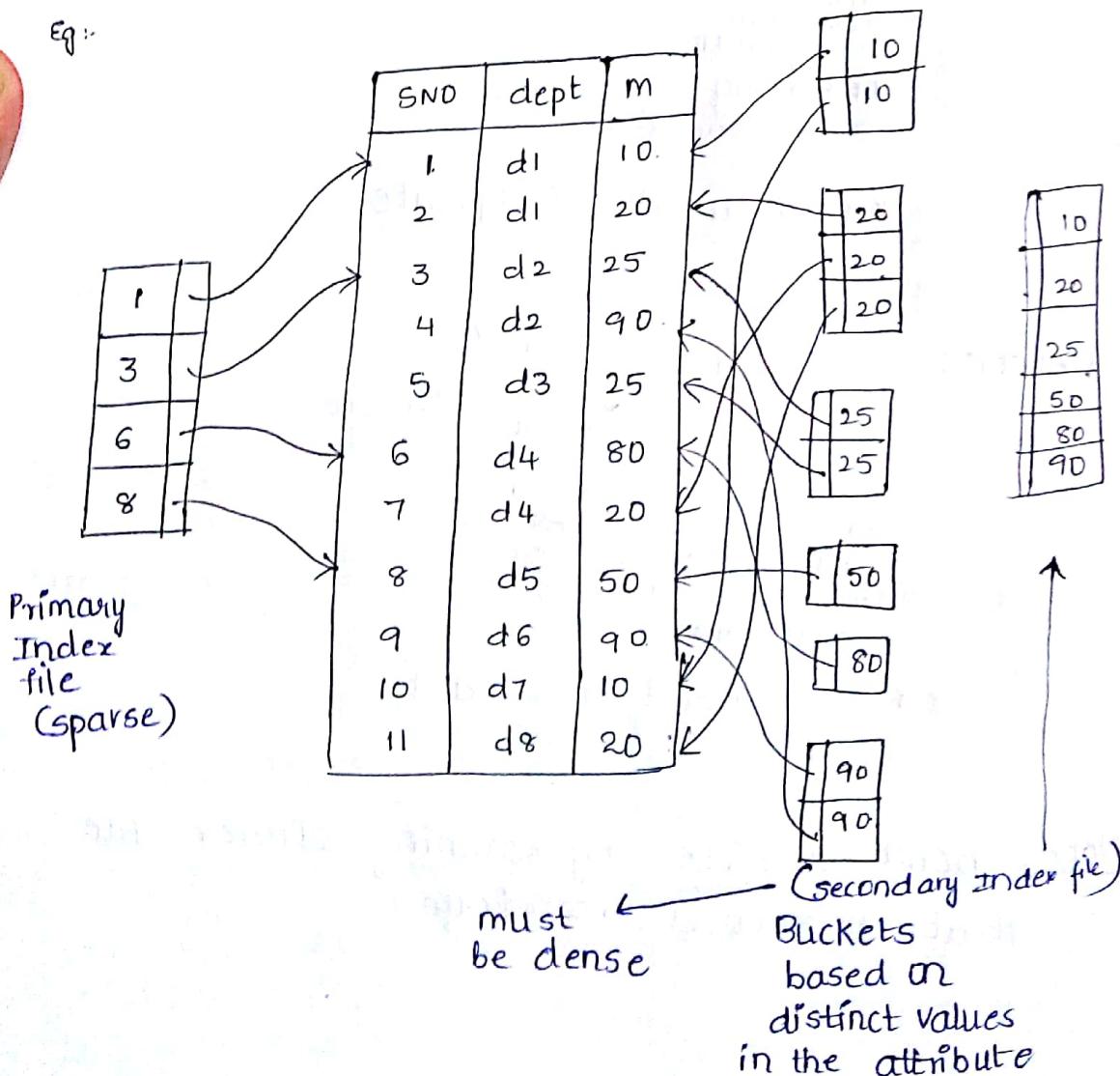
Disadvantages of sparse Indexing

- ⇒ Generally slower than dense indexing

Secondary Indexing :-

- ⇒ Index File is in sorted order.
- ⇒ If queries are based on other attributes other than primary key, create secondary index based on other attributes.

Eg:-



size of bucket = no. of records with that value.

search → M → search in secondary index file using Binary search.
→ R → start from LB

start :

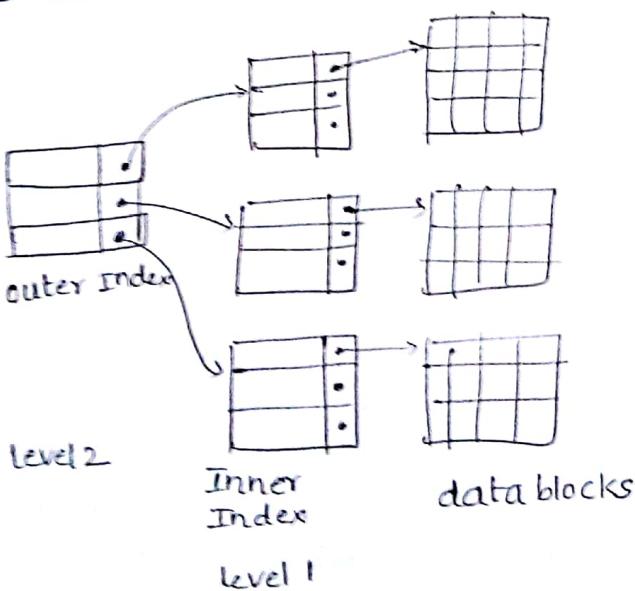
Insert :- From primary index file

Delete :- start from secondary index file.

if after deletion if bucket becomes empty
delete record from secondary index file.

Multilevel Indexing :-

- As the size of the database grows, size of index file also grows.
- Multi-level indexing helps in breaking down the index into several smaller index files in order to make the outermost level so small that it can be saved on prim main memory

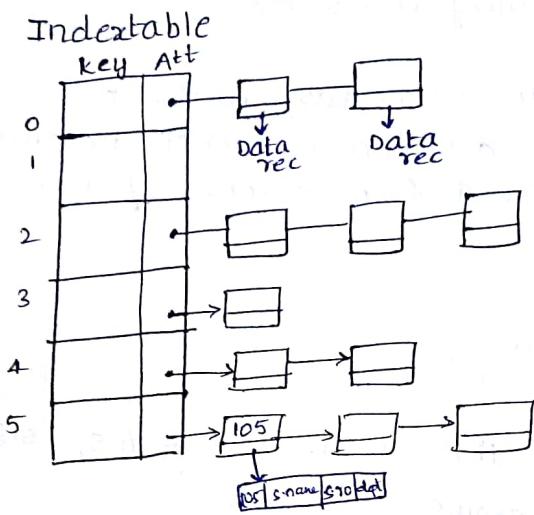


Search → M → $T = t_1 + t_2 + t_3$ → in data file
→ R → in level 1 in level 2
} start from LB

Hash based Indexing

- Separate chaining.
- Extendable Hashing.

Separate chaining



hash Function = $h(x)$

hash key $h(x) = x \bmod \text{TableSize}$

Insertion

→ Exact (Match query)

$$\text{Eg: } 0_{\text{RND}} = 105$$

$$h_k(105) = 5$$

Find the hash key and insert a node in the corresponding key in the Index table.

Delete

→ search using hash key.

If the node is not found in the corresponding list then decide it is not present.

Else delete the record.

update

Insert

Delete

Upda

Range

Ex

⇒ It
so,
rar

Ext

Head
tab

Bucket

Data
record
[7]

update:- Go and search in the hash key value list of Index table.

If found update it
else decide not exist.

$$\text{Insert } T = C + T_1 + C \quad \begin{array}{l} \xrightarrow{\text{hashkey generation}} \\ \xrightarrow{\text{linear search}} \end{array}$$

$$\text{Deletion } T = C + T_1 + C$$

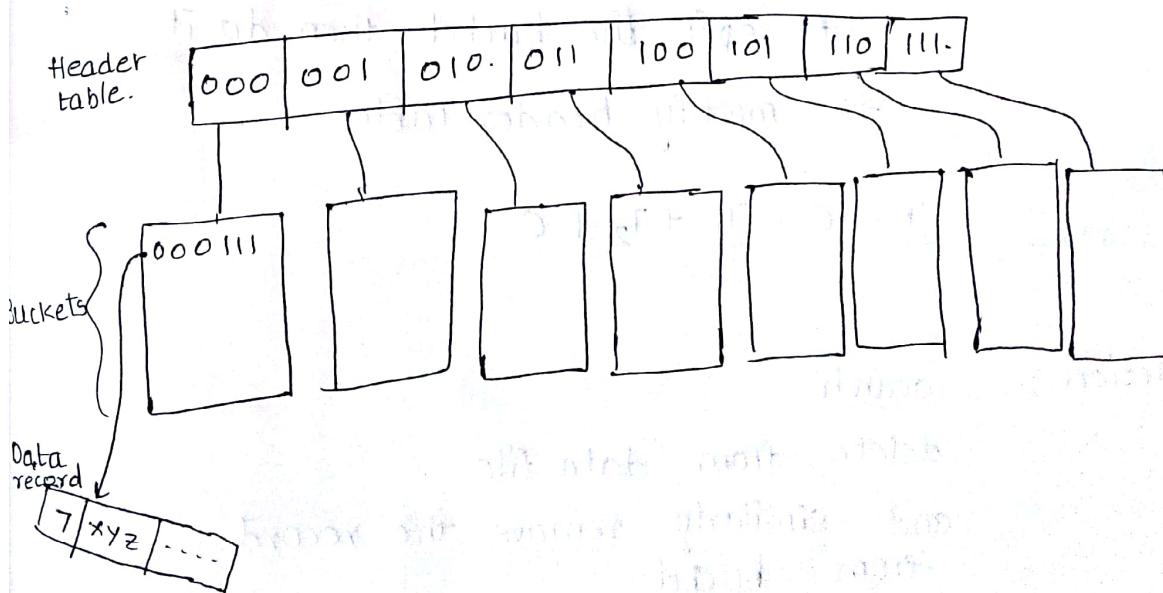
$$\text{Update } T = C + T_1 + C$$

Range Query

Eg:- RNO 110 - 120.

It has to scan Index table almost totally so, hash based indices are not good for range queries.

Extendable Hashing



Search

Exact

Eg; $\overline{0} \text{ (s)}$
 $RNO=7$

→ $0000\ 0111$

↓
Search in header file

Then search in corresponding
bucket

$$T = C + T_1 + C$$

↓
Conversion
to binary ↓
depends
on
bucket size
& search
key values
in bucket

Insertion :-

Eg:-

$19 |$

insert this record in datafile.

and based on first three bits ~~and~~
search in header file.

if bucket is full and is possible
to split the bucket then do it.

else modify header table.

$$T = C + T_1 + T_2 + C$$

Deletion :- search

delete from data file

and similarly remove the record
from bucket.

If, after removal if bucket is empty
then merge the adjacent buckets

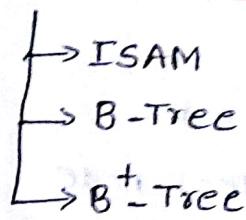
→ merging time

updation: similar to search.

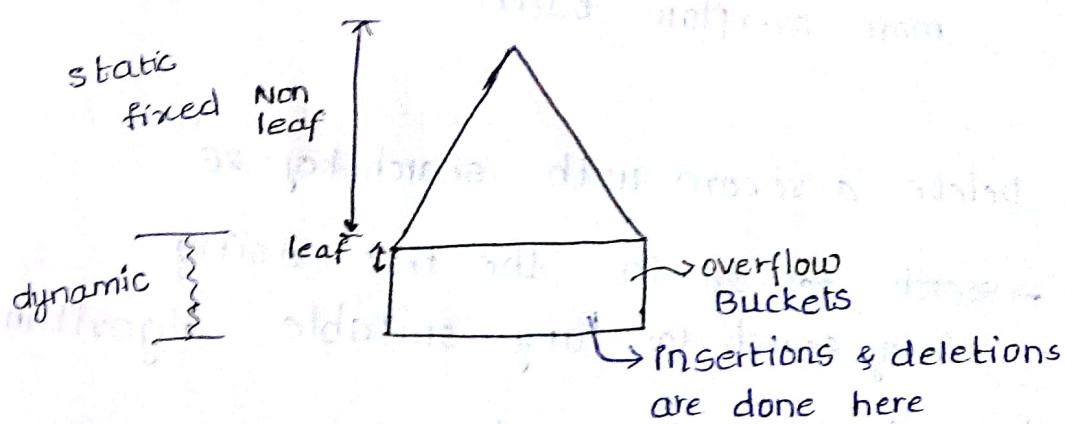
range search:

It is not necessary that all the search key values are found in the same bucket.
we need to scan all buckets almost.
 \therefore It is not a good practice.

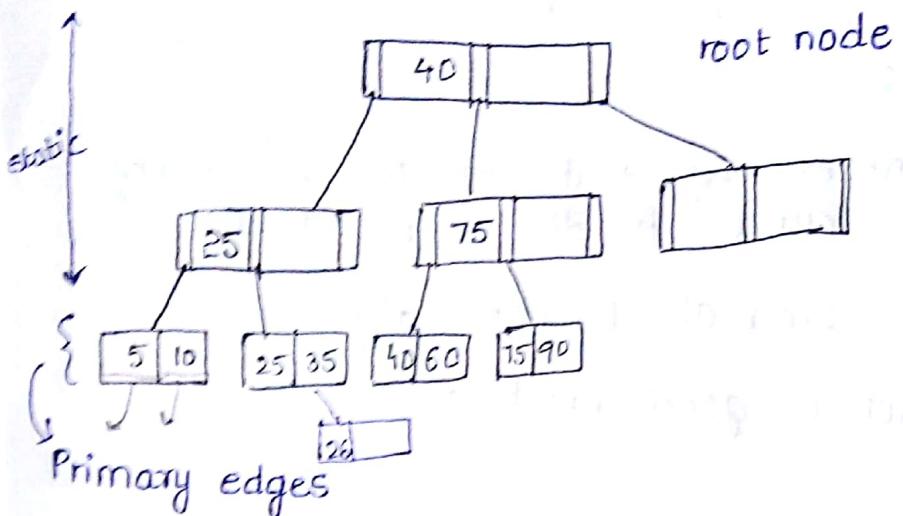
TREE BASED INDEXING TECHNIQUES



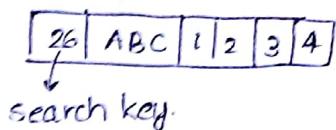
ISAM: node structure is same as B⁺ Tree



Indexed Sequential Access Method (ISAM)



Eg:- 1. Insert



- perform Binary search on root node.
- If space is not found at leaf node.
So, create overflow bucket.
- Even, if space is not found in
Overflow bucket , then create one
more overflow bucket

2. Delete a record with search key = 26.

- search for 26 in the tree using
binary search (or) any suitable algorithm
- If found remove 26
and don't update
(or) merge nodes
- Nodes are merged only when node
is empty.

Answers

If it doesn't exist in leaf node then
decide data record doesn't exist

TRANSACTION MANAGEMENT

Transaction

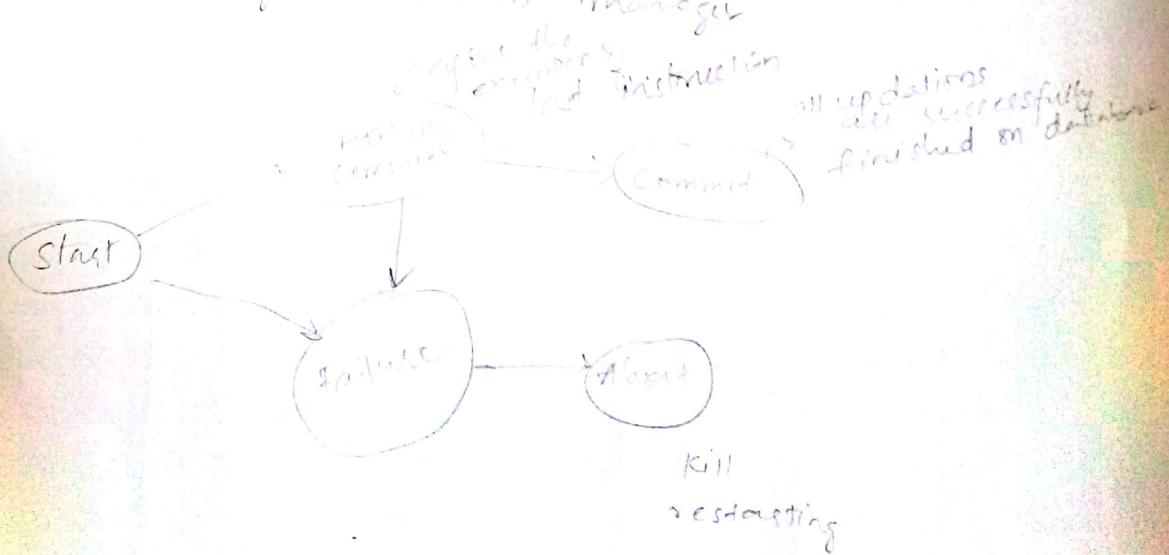
→ set of operations

Eg: Bank database

T₁: A → B

T₁: Begin
 = Read A → selection
 A-balance → update
 write A → modification
 Read B → selection
 B-balance → deposit
 write B → update
 end.

Each transaction is uniquely identified by unique TID assigned by transaction manager



Transaction Properties (ACID)

- **Atomicity**: All the updations are done
 - or
 - None of the updations are done

It avoids inconsistency of Database.



If transaction is

Successful then discards old one and
points to new copy
else points to old copy only.

- **Consistency**: Before and after transaction
DB must ^{be} consistent.

Eg:-

$$\begin{array}{ccc} A & \xrightarrow{100} & B \\ 1000 & & 500 \end{array} \quad A + B = 1500.$$

$$\begin{array}{ccc} A & \xrightarrow{ } & B \\ 900 & & 510 \end{array} \quad A + B = 1500.$$

- **Isolation**: Let t_1, t_2 be two concurrent executing transaction

Isolation property states that

at a time only one transaction takes place in serial order

$$t_1 \rightarrow t_2$$

(or) $t_2 \rightarrow t_1$

► Durability: Once the transaction executes 'COMMIT' all the updations are permanent.

Serial scheduling:-

S₁

T ₁	T ₂
R(A)	
W(A)	
R(B)	
W(B)	
commit	
R(A)	
W(A)	
R(B)	
W(B)	
commit	

S₂

T ₁	T ₂
R(A)	
W(A)	
R(B)	
W(B)	
commit	
R(A)	
W(A)	
R(B)	
W(B)	
commit	

$$\text{waiting time}(T_2) = \text{Execution time}(T_1)$$

$$WT(T_1) = ET(T_2)$$

S₃

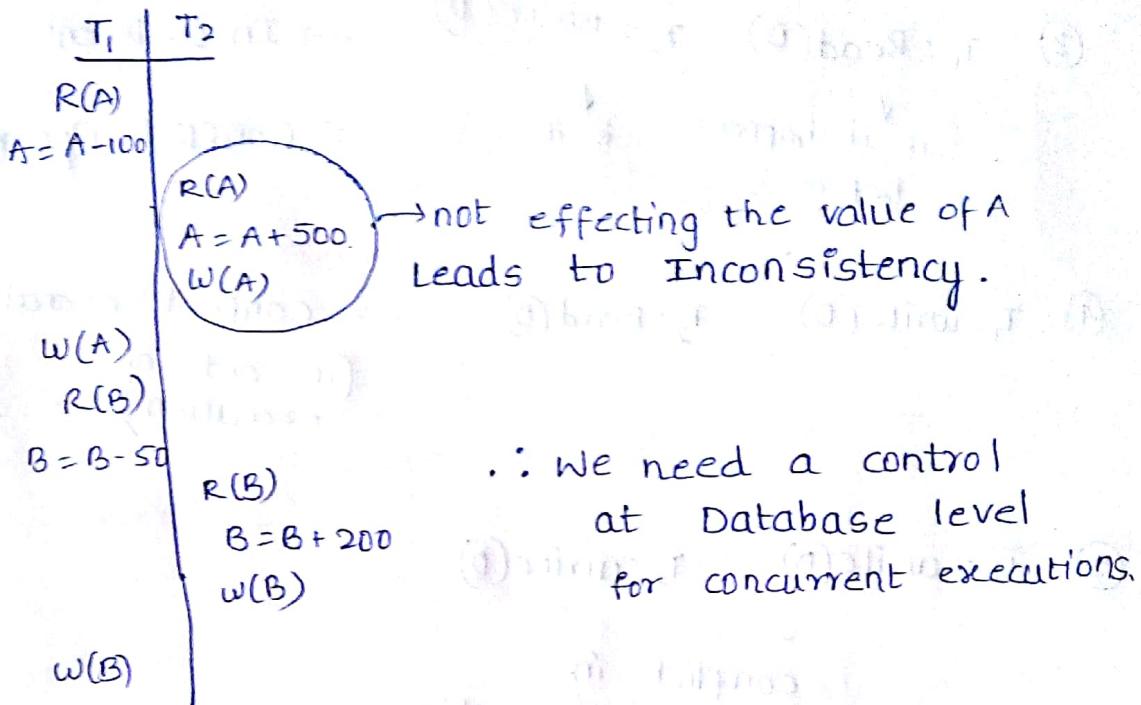
T ₁	T ₂
R(A)	
W(A)	
	R(A)
	W(A)
R(B)	
W(B)	
commit	
R(B)	
W(B)	
commit	

S₄

T ₁	T ₂
R(A)	
W(A)	
R(A)	
W(A)	
R(B)	
W(B)	
commit	
R(B)	
W(B)	
commit	

waiting time decreased.

Example for S₃



S₃ → is a serializable.

→ Allowed to execute.

→ guarantees for Consistency.

Eg:- S₁, S₂

But waiting time is more.

* S is called a serializable schedule, if it is equivalent to serial schedule.

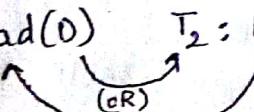
→ Conflict serializable schedule

→ View serializable schedule.

Conflict Serializable Schedule

↳ conflict operation.

Eg:- T₁ : Read(0) T₂ : Read(0)



→ consistent

∴ Non conflict
T₁ followed by T₂

(2)

T_1 : Read(A)

T_2 : Read(B)

→ Non conflict

(3) T_1 : Read(D) T_2 : Write(D) → Inconsistent

if it happens before T_2

∴ Conflict Operations

(4) T_1 : Write(D) T_2 : Read(D) → conflict operations
(in order of execution)

(5) T_1 : write(D) T_2 : write(D)

conflict in
order of executions

* S, S' are conflict equivalent schedules if one schedule is obtained by swapping some non conflict operations of other schedule

Eg:

	T_1	T_2		T_1	T_2	
	R(A)			R(A)		
	W(A)			W(A)		
		R(A)			R(B)	
		W(A)	swap.		W(B)	
non conflicts	R(B)			R(A)		
	W(B)			W(A)		

* if S' is serial schedule then S is conflict serializable schedule

Let T_i and T_j be two transactions.

$G(V, E)$ G is directed graph.
 V is vertex set.
 E is Edge set

$$E = \{e_1, e_2, \dots, e_n\}$$

$V \rightarrow$ Transactions in S

$$e = T_i \rightarrow T_j, \text{ data } 'D'$$

1. If T_i executes $\text{read}(D)$ before the T_j executes $\text{write}(D)$
2. If T_i executes $\text{write}(D)$ before the T_j executes $\text{Read}(D)$
3. If T_i executes $\text{write}(D)$ before the T_j executes $\text{write}(D)$

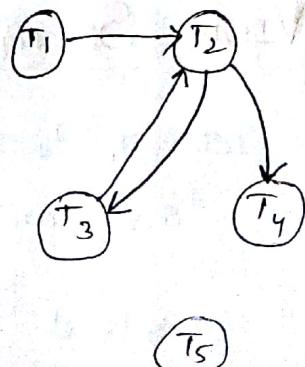
Eg :

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

$$V = \{T_1, T_2\}$$

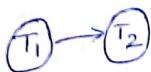


T_1	T_2	T_3	T_4	T_5
$R(A)$				
$W(A)$				
	$R(A)$			
	$W(A)$			
		$R(B)$		
		$W(B)$		
			$R(C)$	
			$W(C)$	
				$R(D)$
				$W(D)$
				$R(E)$
				$W(E)$



* A schedule is said to be conflict serializable if their graphs are equivalent.

Eg:-	T ₁	T ₂
	R(A)	
	W(A)	
	E(B)	
	W(B)	
		R(A)
		W(A)
		E(B)
		W(B)



	T ₁	T ₂
	R(A)	
	W(A)	
	R(A), W(A)	
	R(B)	
	W(B)	
	R(B)	
	W(B)	



*

VIEW SERIALIZATION

* Let T_i, T_j, D dataitem and S, S' are two schedules.

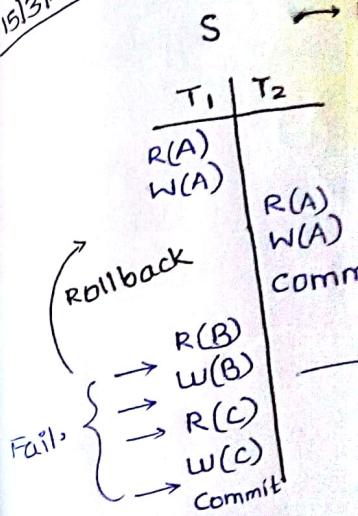
S & S' are view equivalent if it satisfies the following conditions

- 1) If T_i reading the initial value of D in S, the T_i also read the initial value of D in S'

2) If T_i reads other T_j T_i read(S)

3) If T_i performs in S' also

15/3/8



* All non-recovered allowed to only recover

* In a shed roll-backs are called cascaded

* Non-cascading

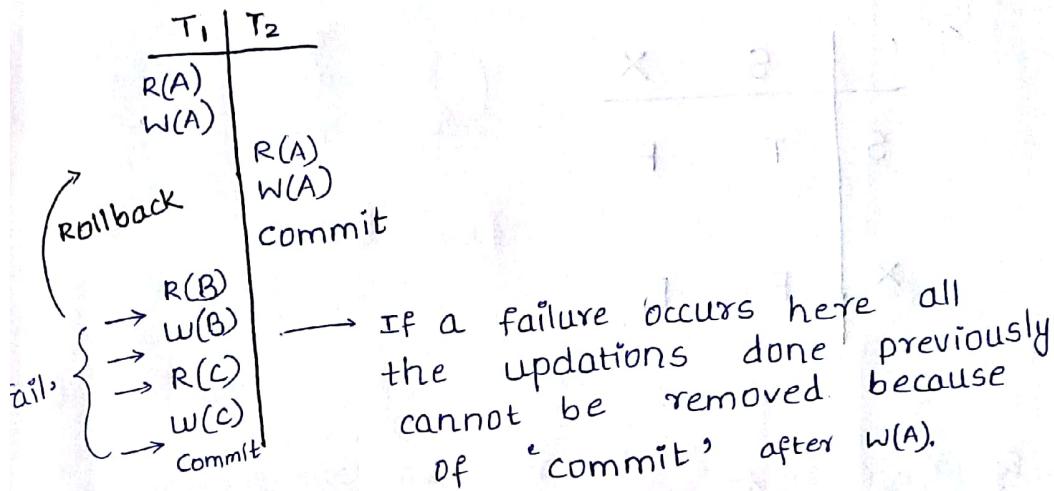
T₁
R(A)
W(A)

2) If T_i reads (D) when written by the other T_j in S , then there exists T_i read(D) where T_j in S'

3) If T_i performs final write on (D) in S , in S' also T_i only performs the final write (D)

13/8

$S \rightarrow$ Non recoverable schedule.

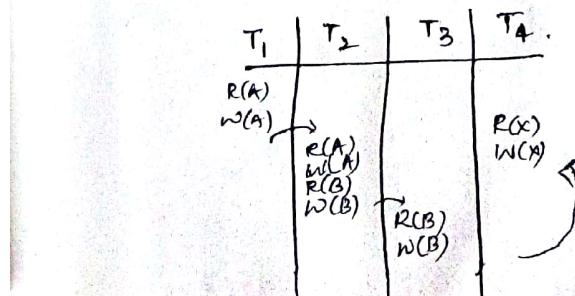


* All non-recoverable schedules are not allowed to execute.

only recoverable schedules are allowed.

* In a schedule when there are series of roll-backs are there, then those schedules are called cascading schedules.

* Non-cascading schedules are allowed but not cascading schedules.



Concurrency Control Management

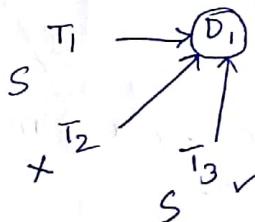
LOCK :

- shared lock S-lock(D) → Read(D)
- Exclusive lock X-lock(D) → Read and write(D)

Lock Compatable Table

	S	X
S	T	F
X	F	F

Eg:-



T_2 waits in waiting queue

; it is not compatible

T_1	Concurrency Control Manager (CCM)
Begin $LOCK_X(A)$ $R(A)$ $W(A)$ $unlock(A)$ $LOCK_X(B)$ $R(B)$ $W(B)$ $unlock(B)$ $commit$	T_1 requesting $LOCK_X(A)$ grant $LOCK_X(B)$ T_1 permit read on A

T_1	T_2
lock_X(A)	
R(A)	
W(A)	
unlock(A)	
lock_X(B)	lock_X(A)
	R(A)
	W(A)
R(B)	
W(B)	
unlock(B)	lock_X(B)
	R(B)
	W(B)
	unlock(B)

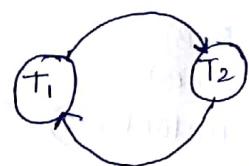
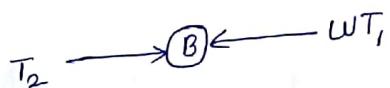
Eg:- Bank

waiting queue (wq): T_2 withdraw } \rightarrow This has to wait till the completion of all other transaction.

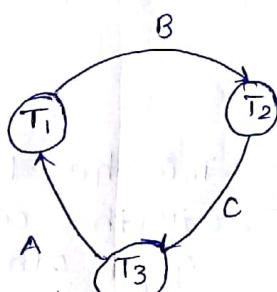
This is known as Starvation.

* But OLY schedule must be of starvation free.

T_1	T_2	
LOCK_X(A)		granted x-lock on A to T_1
R(A)	LOCK_X(B)	granted x-lock on B to T_2
W(A)	R(B)	
	W(B)	
LOCK_X(B)	LOCK_X(A)	
R(B)		$T_1 \leftarrow A$
W(B)		$T_2 \rightarrow B \leftarrow WT_1$
unlock(A)	R(A)	
unlock(B)	W(A)	
	unlock(A)	
	unlock(B)	



Dead lock

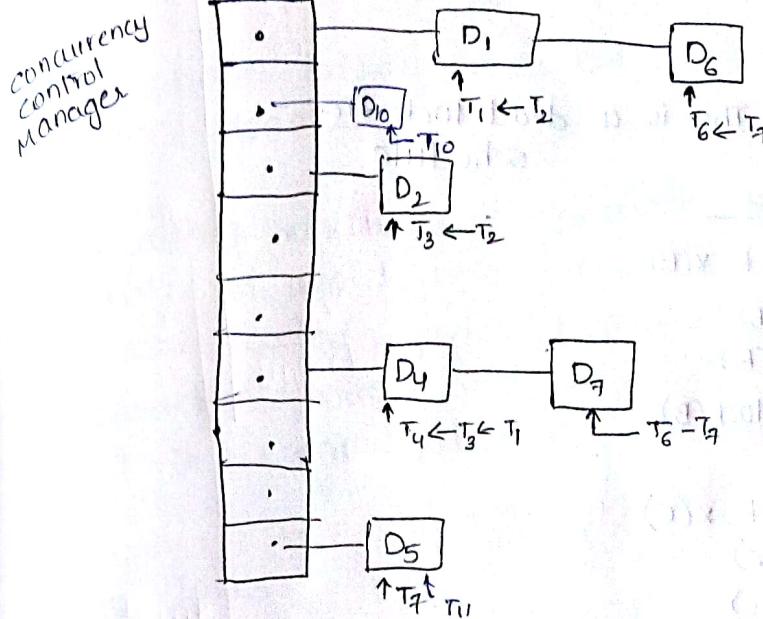


T_1	T_2	T_3
LOCK_X(A)		
R(A)	LOCK_X(B)	
W(A)	R(B)	
	W(B)	
LOCK_X(B)	LOCK_X(C)	LOCK_X(C)
R(B)	R(C)	R(C)
W(B)	W(C)	W(C)
unlock(A)		LOCK_X(A)
unlock(B)	unlock(B)	R(A)
	unlock(C)	W(A)
unlock(A)		unlock(C)
unlock(B)		unlock(A)

→ our schedule must be dead lock free

LOCK Table

Data structure which maintains waiting queue



Eg :- $T_{10} (D_1, D_4, D_{10})$

if D_{10} is not there in the Lock Table
then D_{10} will be placed according to hashing.

⇒ If T_2 is rolled back then T_2 will be removed from lock Table.

→ Implementation:
using hashtable + linked list.

How to identify starvation situation occurred or not?

→ count

How to identify DeadLock situation?

→ Graph, check cycles.

2-PHASE LOCKING Mechanism

→ To avoid non-recoverability

- ① Growing phase → allowed to lock
→ not allowed to unlock
- ② Shrinking phase → allowed to unlock
→ NOT allowed to lock the data item.

Ex:-

This is a dead lock free schedule.

T_1	T_2
Lock_x(A)	LOCK_x(B)
R(A)	R(B)
W(A)	W(B)
unlock(A)	unlock(B)
Lock_x(B)	lock_x(C)
R(B)	R(C)
W(B)	W(C)
unlock(B)	unlock(C)
	lock_x(A)
	R(A)
	W(A)
	unlock(A)

↓
2 phase locking schedule

T_1	T_2
lock_x(A)	lock_x(B)
lock_x(B)	lock_x(C)
	lock_x(A)
R(A)	R(B)
W(A)	W(B)
R(B)	R(C)
W(B)	W(C)
unlock(A)	R(A)
unlock(B)	W(A)
→ commit	unlock(A)
	unlock(C)
	unlock(B)
	commit

There is possibility for cascading.

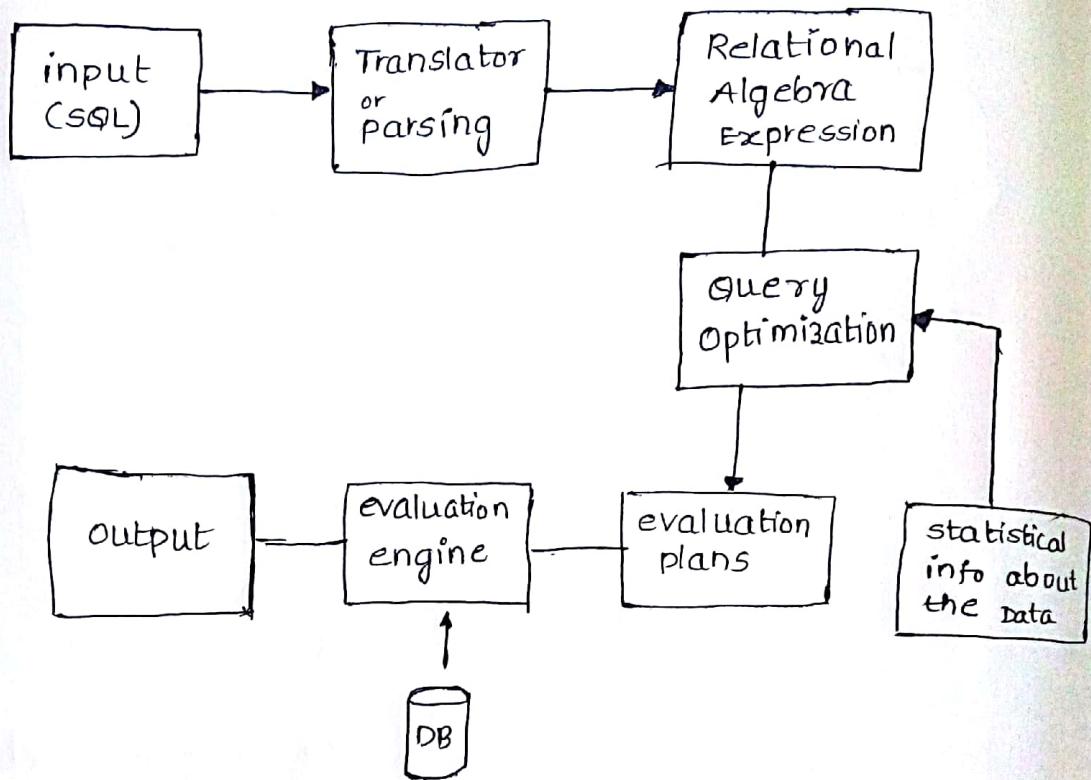
∴ It is helpful for recovery.

- * This mechanism/protocol assures serializability. It can be proved that the transactions can be serialized in the order of their lock points (i.e., the point where a transaction acquired its final lock)
- * Two-phase locking does not ensure freedom from deadlocks.
- * Cascading roll-back is possible under two-phase locking. To avoid this, follow a modified protocol called strict two-phase locking. Here a transaction must hold all its exclusive locks till it commits/aborts.
- * Rigorous two phase locking is even stricter: here all locks

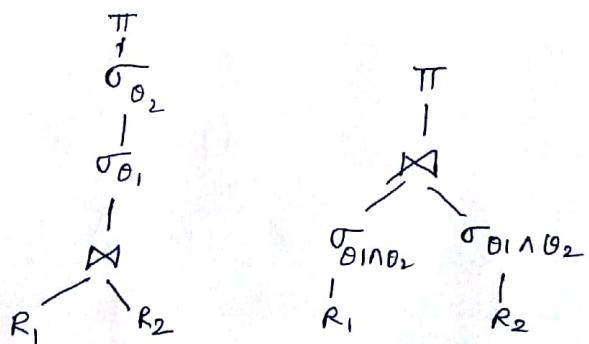
Lock-conversions

- * Two phase locking with lock conversions.
 - First phase :
 - can acquire a lock-s on item.
 - can acquire a lock-x on item.
 - can convert a lock-s to a lock-x (upgrade)
 - Second phase
 - can release a lock-s
 - can release a lock-x
 - can convert a lock-x to a lock-s (downgrade).
- * This protocol assures serializability - But still relies on the programmer to

QUERY PROCESSING

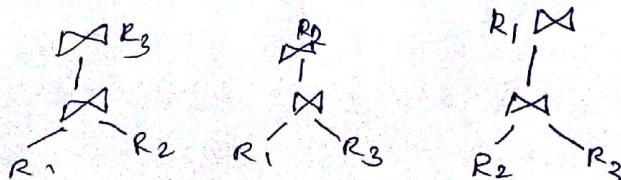


evaluation plans for $\Pi(\sigma_{\theta_1 \wedge \theta_2}(R_1 \bowtie R_2))$



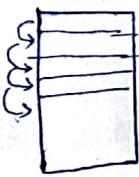
choose a cheaper plan.

$R_1 \bowtie R_2 \bowtie R_3$



Algorithms for Selection

Simple selection $\sigma_{\theta_1}(R_1)$



① Linear scan
(Linear search)

'n' number of comparisons

② θ_1 based on A_1 attribute

Binary search (if input data is sorted
based on A_1)

③ Index based selection. (If index is available
use hash function)

④ compound selection.

$\sigma_{\theta_1 \wedge \theta_2}(R)$

If there is no index or it is unsorted then

use linear scan

If there is index (or) if it is sorted

then use $\sigma_{\theta_2}(\sigma_{\theta_1}(R))$

For compound selection on $\theta_1, \theta_2, \dots, \theta_n$
 $\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_n$

$\sigma_{\theta_n}(\dots (\sigma_{\theta_2}(\sigma_{\theta_1}(R))) \dots)$

here use
Linear scan
(or) Binary search
(or) Index based search

For $\theta_1 \vee \theta_2 \vee \theta_3$
 t_1 tuples \downarrow t_2 tuples $\rightarrow t_3$ tuples
 $t_1 \cup t_2 \cup t_3$

For Range Queries

$$x_1 - x_2$$

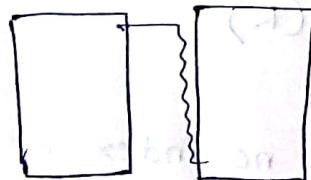
$$\theta = A.\text{val} \geq x_1 \wedge A.\text{val} \leq x_2$$

Algorithms for join

① Nested loop join

$$\sigma_{\theta}(R_1 \times R_2)$$

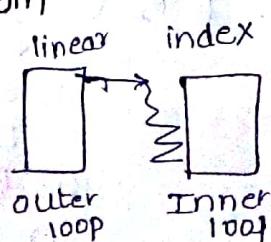
$$\approx 2(m * n)$$



Algorithm

- 1) For each tuple t_i in R_1
 - For each tuple t_j in R_2
 - if (t_i, t_j) satisfies the condition
 - then result = combine(t_i, t_j)

2) Index based loop join



For each tuple t_i in R_1 ,

$x = \text{Index}(t_i, A)$

For each tuple t_j in R_2

if ($\theta(t_j, x)$)

then combine (t_i, t_j)

③ Block Nested loop join.

If the size of relation is greater than size of memory, then we go for Block nested loop join.

For each block b_1 of R_1 do begin

For each block b_2 of R_2 do begin

For each tuple t_i in B_1 of R_1 do begin

For each tuple t_j in B_2 of R_2 do begin

if (t_i, t_j) satisfies the condition

then combine (t_i, t_j)

end

end

end

end

④ Merge-Join.

1. Sort both the relations on their join attribute.

2. Merge the sorted relations to join them.

→ join step is similar to the merge stage of sort-merge algorithm

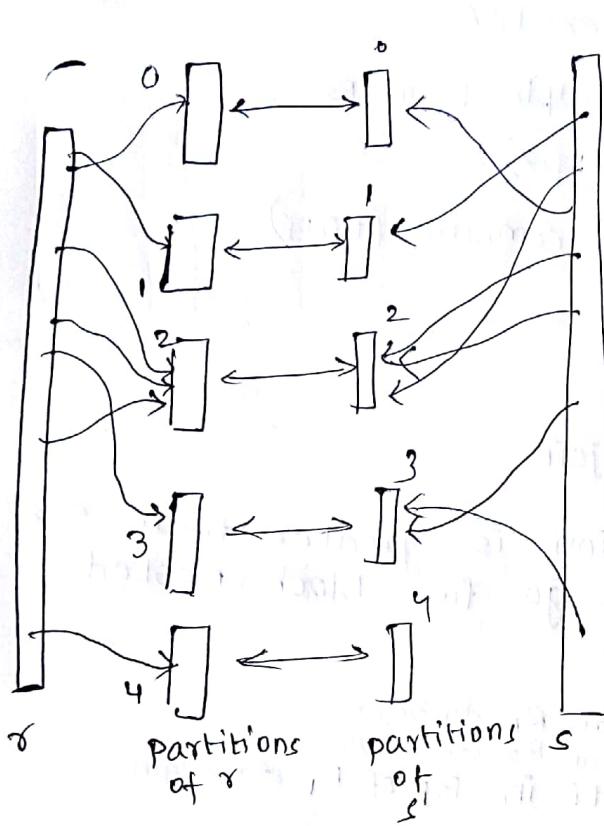
→ Main difference is handling of duplicate values in join attribute - every pair with same value on join attribute must be matched

	a_1	a_2
p1	a	3
	b	1
	c	8
	d	13
	e	7
	f	5

\rightarrow

	a_1	a_3
p2	a	A
	b	G
	c	L
	d	N

5) Hash join.



Complex joins

$\tau \bowtie_{\theta_1, \theta_2, \dots, \theta_n} s$

⇒ either use nested loops/block nested loops

σ  S
 $\theta_1, \theta_2, \dots, \theta_n$

* Duplicate elimination:

can be implemented via hashing or sorting.

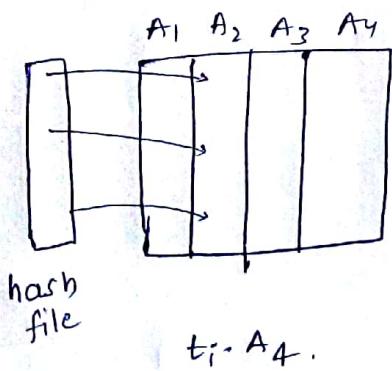
- on sorting duplicates will come adjacent to each other, and all but one set of duplicates can be deleted.
- Optimization: duplicates can be deleted during run generation as well as at intermediate merge steps in external sort merge.
- Hashing is similar - duplicates will come into the same bucket.

* Projection.

- perform projection on each tuple
- followed by duplicate elimination.



Projection



$t_i \cdot A_4$.

$$T = t_1 + t_2$$

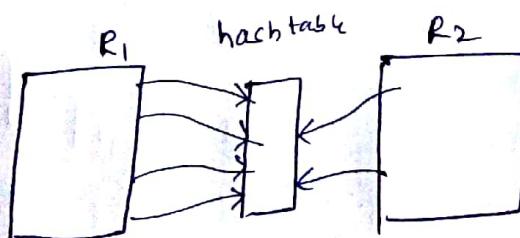
↓
linear scanning of all records to verify if it is a duplicate value or not by hashing.

Set Operation

$\cup, \cap, -$

Union (\cup) :

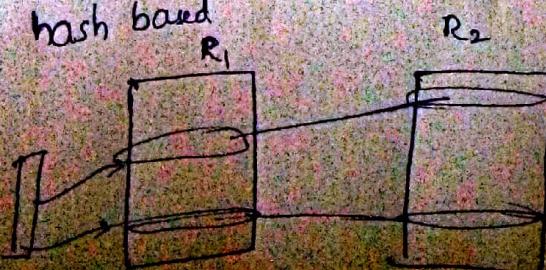
hash based



sort based

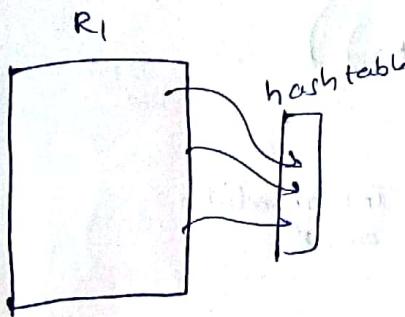
intersection (\cap) :

hash based



set difference (-):

$$R_1 - R_2$$



NOW scan R_2 , for each tuple of R_1 i
if its corresponding hash value
is already, delete its hash value.

result = hashtable.

Aggregate operations

- SUM
- COUNT
- MIN
- MAX

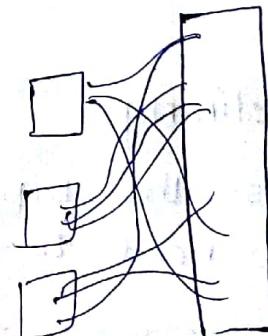
- (group by)

hash based.

sort based

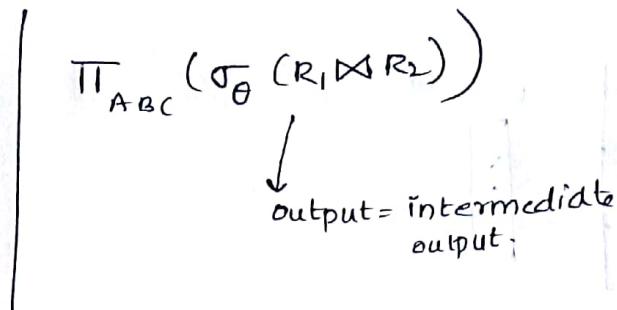
→ calculations
like min, max,
sum are done
on the resultant
tuples present
consecutively

→ buckets relation.



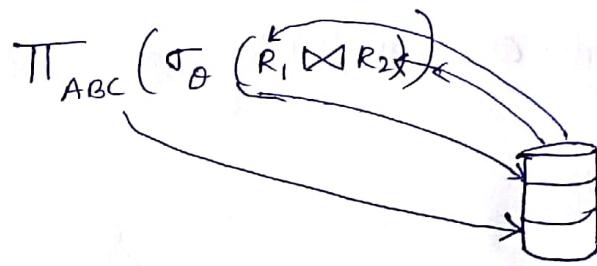
perform calculation
like min, max, count
on each bucket

Evaluation of Expressions



Materialization

Generates results of an expression whose inputs are relations or are already computed, materialize (store) it on disk. Repeat.

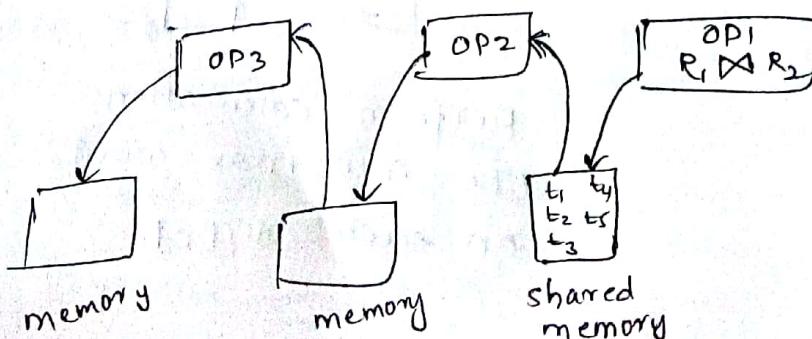


Drawbacks

More number of I/O operations are present which are costlier.

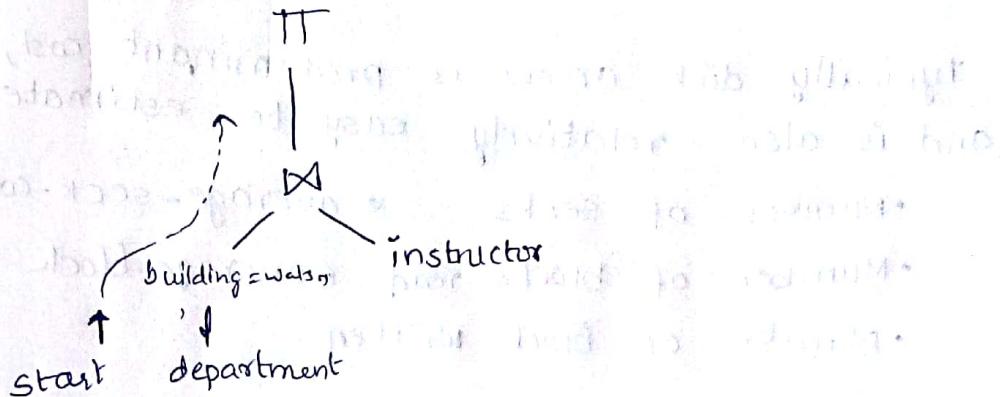
↳ Pipelining:

Pass the tuples to next parent operation.
No need to store just pass on.



Eg:- for Materialized

$$\Pi_{\text{non}}^r (\sigma_{\text{building} = \text{"watson"}}, \text{department} \bowtie \text{instructor})$$



⇒ pipeline may not always be possible
eg., sort, hash-join.

⇒ For pipeline to be effective, use evaluation algorithms that generate output tuples even as tuples are received for inputs to the operation.

⇒ pipeline can be executed in two ways:

Producer driver

demand driven

⇒ Producer starts producing only on the demand of the consumer.

⇒ if there is no overflow, just keep on producing (or) writing.

Query cost :- Total elapsed time for answering a query.

Factors contributing to time cost:-

→ disk accesses, CPU or network communication.

Typically disk access is predominant cost, and is also relatively easy to estimate.

- Number of seeks \times average-seek-cost
- Number of blocks read \times average-block-read-cost
- Number of blocks written

\Rightarrow cost of write is greater than cost of read

(Data is read back after being written to ensure that the write was successful.)

$t_T \rightarrow$ time to transfer one block.

$t_S \rightarrow$ time for one seek

\therefore cost of b block transfers + S seeks

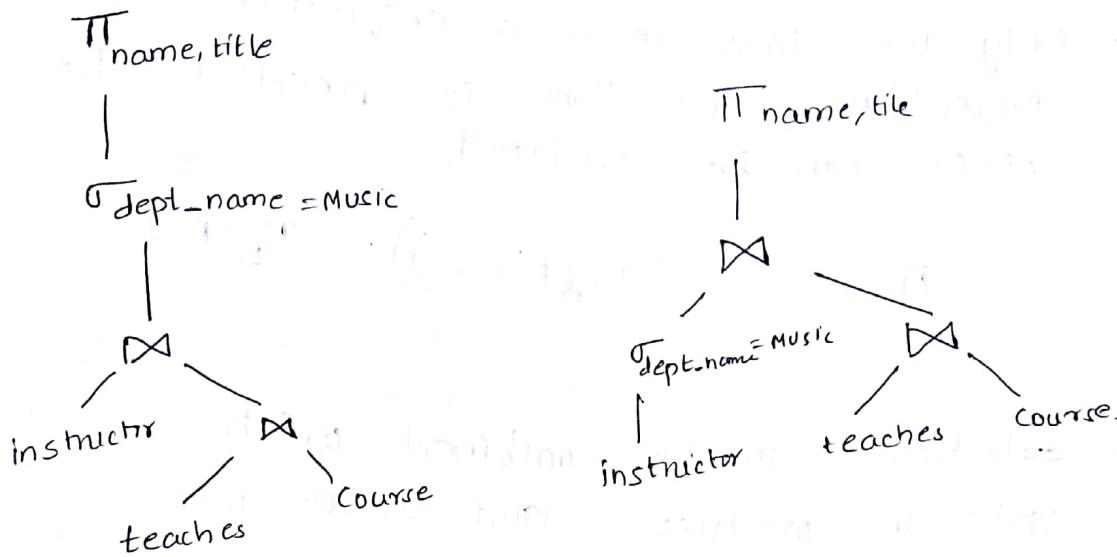
$$\boxed{\text{cost} = b * t_T + S * t_S}$$

Note:- we ignore CPU costs for simplicity.

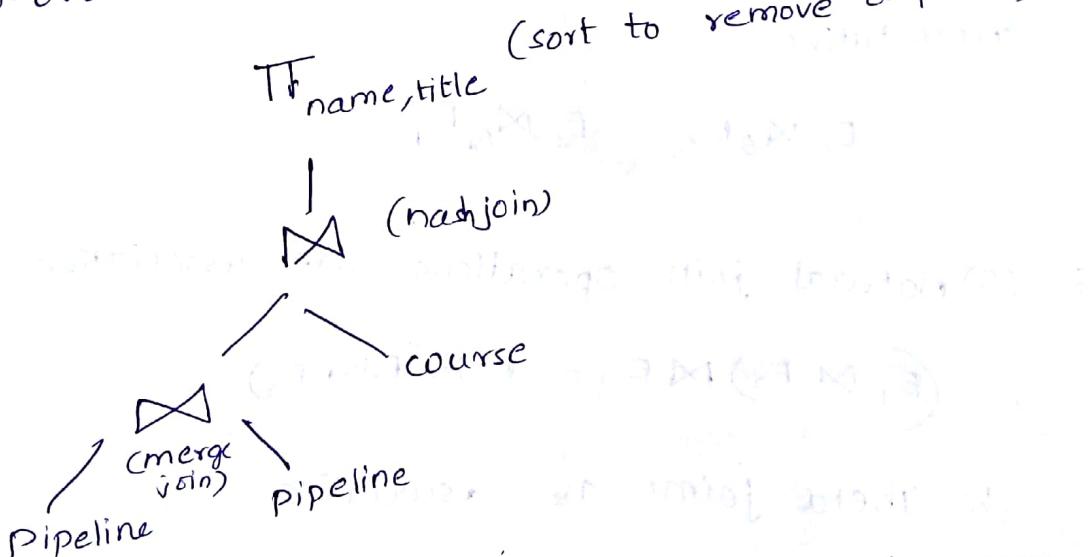
===== (Real systems do take CPU cost into account).

$$\begin{aligned} \text{Cost of operation} &= \text{Access time (cost)} + \text{Execution cost} + \text{Communication cost}. \end{aligned}$$

QUERY OPTIMIZATION



- * An evaluation plan defines exactly what algorithm is used for each operation and how the execution of the operations co-ordinate.



- If two queries are accepting same input and giving same output then they are Equivalent.

* Equivalence Rules Operations can be deconstructed.

$$\sigma_{\theta_1 \wedge \theta_2}(E) = \sigma_{\theta_1}(\sigma_{\theta_2}(E))$$

2. Selection operations are commutative

$$\sigma_{\theta_1}(\sigma_{\theta_2}(E)) = \sigma_{\theta_2}(\sigma_{\theta_1}(E))$$

3. Only the last ~~is~~ in a sequence of projection operations is needed, the others can be omitted.

$$\pi_{L_1}(\pi_{L_2} \dots (\pi_{L_n}(E)) \dots) = \pi_{L_1}(E)$$

4. Selections can be combined with cartesian products and theta joins.

a) $\sigma_{\theta}(E_1 \times E_2) = E_1 \bowtie_{\theta} E_2$

b) $\sigma_{\theta_1}(E_1 \bowtie_{\theta_2} E_2) = E_1 \bowtie_{\theta_1 \wedge \theta_2} E_2$

5. Theta-join operations (and natural joins) are commutative

$$E_1 \bowtie_{\theta} E_2 = E_2 \bowtie_{\theta} E_1$$

6. (a) Natural join operations are associative.

$$(E_1 \bowtie E_2) \bowtie E_3 = E_1 \bowtie (E_2 \bowtie E_3)$$

(b) Theta joints are associative

$$(E_1 \bowtie_{\theta_1} E_2) \bowtie_{\theta_2 \wedge \theta_3} E_3 = E_1 \bowtie_{\theta_1 \wedge \theta_3} (E_2 \bowtie_{\theta_2} E_3)$$

7. selection operation distributes over the theta join only when
- when all the attributes in θ involve only the attributes of one of the being joined

$$\sigma_{\theta_0}(E_1 \bowtie_{\theta} E_2) = (\sigma_{\theta_0}(E_1)) \bowtie_{\theta} E_2$$

- when θ_1 involves only the attributes of E_1 and θ_2 involves only the attributes of E_2

$$\sigma_{\theta_1 \wedge \theta_2}(E_1 \bowtie_{\theta} E_2) = (\sigma_{\theta_1}(E_1)) \bowtie_{\theta} (\sigma_{\theta_2}(E_2))$$

- 8) projection operation distributes over the theta join only when

- If θ involves only attributes from $L_1 \cup L_2$

$$\pi_{L_1 \cup L_2}(E_1 \bowtie_{\theta} E_2) = (\pi_{L_1}(E_1)) \bowtie_{\theta} (\pi_{L_2}(E_2))$$

- b) consider a join $E_1 \bowtie_{\theta} E_2$
- Let L_1 and L_2 be sets of attributes from E_1 and E_2 respectively.
 - Let L_3 be attributes of E_1 that are involved in join condition θ , but are not in $L_1 \cup L_2$

and

- Let L_4 be attributes of E_2 that are involved in join condition θ , but are not in $L_1 \cup L_2$

$$\pi_{L_1 \cup L_2}(E_1 \bowtie_{\theta} E_2) = \pi_{L_1 \cup L_2}((\pi_{L_3 \cup L_4}(E_1)) \bowtie_{\theta} (\pi_{L_3 \cup L_4}(E_2)))$$

9) set operations Union and intersection
are commutative

$$E_1 \cup E_2 = E_2 \cup E_1$$

$$E_1 \cap E_2 = E_2 \cap E_1$$

(set difference is not commutative)

10) Associative
 $(E_1 \cup E_2) \cup E_3 = E_1 \cup (E_2 \cup E_3)$

$$(E_1 \cap E_2) \cap E_3 = E_1 \cap (E_2 \cap E_3)$$

11) Selection Operation distributes over U, ∩ & -

$$\sigma_\theta(E_1 - E_2) = \sigma_\theta(E_1) - \sigma_\theta(E_2)$$

and similarly for ∪ and ∩ in place of -

Also $\sigma_\theta(E_1 - E_2) = \sigma_\theta(E_1) - E_2$

and similarly for ∩ in place of '-'
but not for ∪

12) Projection Operation distributes over union

$$\pi_L(E_1 \cup E_2) = \pi_L(E_1) \cup \pi_L(E_2)$$

** Practical query optimizers incorporate the following approaches.

- Search all the plans and choose the best plan in a cost-based fashion.
- Uses heuristics to choose a plan.

Cost-Based Optimization with Equivalence Rules

- 1) Physical Equivalence rules allow logical query plan to be converted to physical query plan specifying what algorithms are used for each step.
- 2) Efficient optimizer based on equivalent rules depends on
 - A space efficient representation of expressions which avoids making multiple copies of subexpression.
 - Techniques for detecting duplicate derivation of expressions.
 - A form of dynamic programming based on memoization, which stores the best plan for a subexpression the first time it is optimized, and reuses it on repeated optimization calls on same sub-expression.



Heuristic Optimization

- Cost-based optimization is expensive, even with dynamic programming.
- Systems ~~may~~ use heuristics to reduce the number of choices that must be made in a cost-based fashion.

* Heuristic Optimization transforms the query-tree by using a set of rules that typically (but not in all cases) improves execution performance.

- (i) Perform Selection early (reduces the number of tuples)
- ii) Perform Projection early (reduces the number of attributes)
- iii) Perform most restrictive selection and join operations (i.e., with smallest result size) before other similar operations.
- iv) Some systems use only heuristics, others combine heuristics with partial cost-based optimization

RECOVERY Systems

- Failure classification
- storage structure
- Log-Based Recovery
- ARIES

Failure classification

1) Transaction Failure.

→ logical errors:

• inconsistency of data
• inconsistency of log
• inconsistency of transaction
• inconsistency of state

→ system errors:

• power failure

2) System Crash:

3) Disk Failure:

Storage structure:

1) volatile storage:

does not survive system crashes.

Ex: main memory, cache memory

2) Non volatile storage:

survives system crashes.

Ex: disk, tape, flash memory,

non-volatile (battery backed up) RAM

but may still fail losing data.

3) Stable storage

- Maintain multiple copies of each block on separate disks.
- Failure during data transfer can still result in inconsistent copies.

Log-Based Recovery :-

Recovery algorithms.

During normal transaction processing to ensure enough information exists to recover from failure

Actions taken after a failure to recover the database contents to a state that ensures atomicity, consistency and durability

Log :

- History of actions executed by the DBMS.
- Sequence of log records , and maintains a record of update activities on the database.
- stored in stable storage.

LSN

- Every log record given a unique ID called the log sequence number.
- LSN are given in increasing order.

PageLSN

Every page in DB contains the LSN of the most recent log record that describes a change to this page

① when transaction by writing a T_i starts, it registers itself $\langle T_i \text{ start} \rangle$ log record.

② Before T_i executes write(x), a log record $\langle T_i, x, v_1, v_2 \rangle$ is written,

where v_1 is the value of x before the write (the old values) and v_2 is the value to be written to x (the new value)

③ When transaction T_i finishes its last statement the log record

$\langle T_i \text{ commits} \rangle$ is written

④ When transaction T_i has aborted $\langle T_i \text{ abort} \rangle$

Undo & Redo

undo of a log record $\langle T_i, x, v_1, v_2 \rangle$ writes the old value v_1 to x .

Redo of a log record $\langle T_i, x, v_1, v_2 \rangle$ writes the new value v_2 to x .

check point

- streamline recovery procedure by periodically performing checkpointing.
- output all log records currently residing in main memory onto stable storage.
 - output all modified buffer blocks to the disk.
 - write a log record $\langle \text{checkpoint } L \rangle$ onto stable storage where L is list of all transactions active at the time of checkpoint.
 - All updates are stopped while doing checkpointing.

SECURITY AND AUTHORIZATION

* **Secrecy** :- Users should not be able to see things they are not supposed to be.

Eg:- A student can't see other students' grade.

* **Integrity** :- Users should not be able to modify things they are not supposed to.

Eg:- Only instructors can assign grades.

* **Availability** :- Users should be able to view and modify data.

** A Security policy specifies who is authorized to do what.

** A security mechanism allows us to enforce a chosen security policy.

* Two main mechanisms at the DBMS level.

- i) Discretionary access control.
- ii) Mandatory access control.

Discretionary access control.

* Based on the concept of access rights or privileges for objects (tables and views) and mechanisms for giving users privileges (and revoking privileges).

* Creator of a table or a view automatically gets all privileges on it.

• DBMS keeps track of who and loses privileges, and subsequently gains ensures that only

requests from users who have the necessary privileges (at the time the request is issued) are allowed.

GRANT command

GRANT privileges ON object TO users [WITH GRANT OPTION]

privileges

1. SELECT : can read all columns
2. INSERT : can insert tuples.
3. DELETE : can delete tuples.
4. REFERENCES : can define foreign keys. that refer to this column

* If a user has a privilege with the GRANT OPTION can pass privilege on to other users.

* only owner can execute CREATE, ALTER, and DROP

Eg:- 1. GRANT INSERT, SELECT ON Sailors TO Horatio.

2. GRANT DELETE ON SAILORS TO GUPPY WITH GRANT OPTION.

3. GRANT UPDATE(rating) ON Sailors TO Dustin,

4. GRANT SELECT ON ActiveSailors TO GUPPY, YUPPY.

5. REVOKE: When a privilege is revoked
from X

THERE IS NO SUBSTITUTE

FOR HARD WORK. NEVER GIVE UP.
NEVER STOP BELIEVING."

- Manam Ramu